**Handbook of**

# THEORETICAL and COMPUTATIONAL
# NANOTECHNOLOGY

### Nanodevice Modeling and Nanoelectronics

**10**

$$g_1(r) = \exp\left[-\frac{\upsilon_1(r)}{k_B T}\right]\gamma_1(r)$$

$$\ln\left[\frac{g(r)}{g_1(r)}\right]$$

$$\langle k+q|w_0|k\rangle = \frac{4\pi}{\Omega_0}\int_0^\infty w_0(R)\frac{\sin qR}{qR}R^2\,dR$$

*Edited by*

**Michael Rieth and Wolfram Schommers**

Fo

Pierre-Gilles de Gennes, Nobel Prize Laureate

**ASP**

**AMERICAN
SCIENTIFIC
PUBLISHERS**

*Handbook of*

# THEORETICAL and COMPUTATIONAL NANOTECHNOLOGY

# Titles in Nanotechnology Book Series

*Founding Editor*

## Dr. Hari Singh Nalwa

1. **Encyclopedia of Nanoscience and Nanotechnology, 10-Volume Set**
   *Edited by* Hari Singh Nalwa

2. **Handbook of Theoretical and Computational Nanotechnology, 10-Volume Set**
   *Edited by* Michael Rieth and Wolfram Schommers

3. **Bottom-up Nanofabrication: Supramolecules, Self-Assemblies, and Organized Films, 6-Volume Set**
   *Edited by* Katsuhiko Ariga and Hari Singh Nalwa

4. **Handbook of Semiconductor Nanostructures and Nanodevices, 5-Volume Set**
   *Edited by* A. A. Balandin and K. L. Wang

5. **Handbook of Organic-Inorganic Hybrid Materials and Nanocomposites, 2-Volume Set**
   *Edited by* Hari Singh Nalwa

6. **Handbook of Nanostructured Biomaterials and Their Applications in Nanobiotechnology, 2-Volume ɛ ⅛ he Set**
   *Edited by* Hari Singh Nalwa

7. **Handbook of Electrochemical Nanotechnology, 2-Volume Set**
   *Edited by* Yuehe Lin and Hari Singh Nalwa

8. **Polymeric Nanostructures and Their Applications, 2-Volume Set**
   *Edited by* Hari Singh Nalwa

9. **Soft Nanomaterials, 2-Volume Set**
   *Edited by* Hari Singh Nalwa

10. **Functional Nanomaterials**
    *Edited by* Kurt E. Geckeler and Edward Rosenberg

11. **Synthesis, Functionalization and Surface Treatment of Nanoparticles**
    *Edited by* I. M. Baraton

12. **Quantum Dots and Nanowires**
    *Edited by* S. Bandyopadhyay and Hari Singh Nalwa

13. **Nanoclusters and Nanocrystals**
    *Edited by* Hari Singh Nalwa

14. **Molecular Nanoelectronics**
    *Edited by* Mark A. Reed and T. Lee

15. **Magnetic Nanostructures**
    *Edited by* Hari Singh Nalwa

16. **Nanoparticles for Pharmaceutical Applications**
    *Edited by* J. Domb, Y. Tabata, M. N. V. Ravi Kumar, and S. Farber

17. **Cancer Nanotechnology**
    *Edited by* Hari Singh Nalwa and Thomas Webster

18. **Biochips Nanotechnology**
    *Edited by* Nongyue He and Hari Singh Nalwa

19. **Nanotoxicology**
    *Edited by* Yuliang Zhao and Hari Singh Nalwa

20. **Polymer Nanocomposites and Their Applications**
    *Written by* Suprakas Sinha Ray and Mosto Bousmina

21. **Nanoscale Science and Engineering Education**
    *Edited by* Aldrin E. Sweeney and Sudipta Seal

22. **Hard Nanomaterials**
    *Edited by* Hari Singh Nalwa

**Additional Volumes in Preparation**

*Visit: www.aspbs.com*

# *Handbook of*
# THEORETICAL and
# COMPUTATIONAL
# NANOTECHNOLOGY

Volume 10

**Nanodevice Modeling and
Nanoelectronics**

*Edited by*

**Michael Rieth** and **Wolfram Schommers**

Forschungszentrum Karlsruhe, Karlsruhe, Germany

# AMERICAN SCIENTIFIC PUBLISHERS

Handbook of Theoretical and Computational Nanotechnology
edited by Michael Rieth and Wolfram Schommers.

The image on the cover of this handbook was provided by Professor Jeong Won Kang, Chung-Ang University, Seoul, Korea. See Jeong Won Kang, Won Young Choi, and Ho Jung Hwang. *Journal of Computational and Theoretical Nanoscience*, Vol. 1(2), pp. 199–203 (2004). Copyright © 2004, American Scientific Publishers.

This book is printed on acid-free paper. ⊗

The information provided in this handbook is compiled from reliable sources but the contributing authors, editors, and the publisher cannot assume any responsibility whatsoever for the validity of all statements, illustrations, data, procedures, and other related materials contained herein or for the consequences of their use.

**Dr. Michael Rieth**
Forschungszentrum Karlsruhe
Institute of Materials Research I
D-76021 Karlsruhe, GERMANY

**Prof. Dr. Wolfram Schommers**
Forschungszentrum Karlsruhe
Institute for Scientific Computing
D-76021 Karlsruhe, GERMANY

# Foreword

Nanoscience is fashionable. All administrations in the Western world have stressed their interest in nanoobjects and nanotechnologies. As usual, this type of large scientific movement has its pluses and minuses. Many scientists join the crowd without necessarily changing anything in their actual work. Most chemists, for instance, build new molecules that may be called nanoobjects; but again, as usual, the movement does generate significant new content.

Let us, for instance, follow the role of nanostructures in *chemistry*. On one side, nature has provided us with beautiful, robust objects such as fullerenes and carbon tubes, which have some admirable properties. The current challenge is to obtain them in large amounts and at a reasonable price. Here is the real problem.

A completely different sector is obtained from *chemical nanomachines*, for which a molecular unit of nanometric size moves with respect to another one through a change in redox potential or pH. Some of these machines have been built. At the moment, I feel rather skeptical about them because they are extremely costly, extremely fragile (sensitive to poisons), and not easy to protect with a suitable coating—or by a local "antipoison" center. But, here again, there is a challenge.

Let us now turn to *biology*. Here we find an immense group of working nanomachines, enzymes, ionic channels, sensor proteins, adhesion molecules, and so on. They are extremely impressive, but of course they represent progressive construction by trial and error over more than a billion years. Should we try to mimic these machines or, rather, use them for technological purposes, *as they are*, for instance, to grow plants or create proteins at an industrial level according to the techniques of molecular genetics? This is a major question.

A third, open side is *quantum physics* and the (remote) possibility of quantum computers. In my youth, I had hopes for digital storage via quantized flux quanta: The corresponding technology, based on Josephson functions, was patiently built by IBM, but they ultimately dropped out. This shows the hardship of nanotechnologies even when they are handled by a large, competent group. But the cause is not lost, and it may well be that our children use some unexpected form of quantum computers.

Thus, we are facing real challenges, not just the vague recommendations of some anonymous boards. And, we need the tools. We need to know the behavior of materials at the nanolevel, the clever tricks of physical chemistry required to produce nanoparticles or nanopores, the special properties of small cooperative systems (nanomagnets, nanosuperconductors, nanoferroelectrics, etc.), the ability for assembling functional units, and so on.

The aim of the present handbook is to help us with the tools by suitable modelizations. It is written by leading experts, starting from general theoretical principles and progressing to detailed recipes.

In the second half of the 18th century, all the knowledge (fundamental and practical) of the Western world was condensed into an outstanding encyclopedia constructed energetically by Denis Diderot just after the industrial revolution started. Here, at a more modest level, we can hope for something similar. Soon after the first wave, including this handbook, a certain form of nanoindustry may be born.

The discussions started in this handbook will continue in a journal (*Journal of Computational and Theoretical Nanoscience*) launched by the present editors. I wish them the best.

**Professor Pierre-Gilles de Gennes**

Nobel Prize Laureate, Physics
Collège de France
Paris, France

v

# Preface

This is the first handbook that deals with theoretical and computational developments in nanotechnology. The 10-volume compendium is an unprecedented single reference source that provides ideal introduction and overview of the most recent advances and emerging new aspects of nanotechnology spanning from science and engineering to neurogenetics. Many works in the field of theoretical and computational nanotechnology have been published to date, but no book or handbook has focused on all aspects in this field that deal with nano-machines, electronics, devices, quantum computing, nanostructured materials, nanorobotics, medicine, biology, biotechnology, and more.

There is no doubt that nanoscience will be the dominant direction for technology in this new century, and this science will influence our lives to an extent impossible in years past: Specific manipulations of matter at its ultimate level will open completely new perspectives on all scientific and technological disciplines. To be able to produce optimal nanosys-tems with tailor-made properties, it is necessary to analyze and construct such systems in advance by adequate theoretical and computational methods. The handbook gives a com-plete overview of the essential methods, models, and basic pictures.

But, as is well known, there are also threats connected with nanotechnology, specifically with respect to biological systems: Self-assembly can be an uncontrolled process, and the final state of a developing system is in general uncertain in such cases. To avoid undesir-able developments, the theoretical (computational) analysis of such processes is not only desirable but also absolutely necessary. Thus, the computational and theoretical methods of nanoscience are essential for the prediction of new and custom nanosystems and can help keep nanoscience under control. There is basically no alternative. Therefore, one possible answer to the question, "Why a book on theoretical and computational nanotechnology?" is *to give nanotechnology a direction!*

In the design of macroscopic and microscopic systems, engineering is essentially based on *intuitive concepts*, which are tailored to observations in everyday life. Classical mechanics is also based on these macroscopic observations, and its notions have been chosen with respect to our intuitive demands for *visualizability*. However, when we approach the nanolevel, the tools used for the design of macroscopic and microscopic systems become more and more useless. At the nanolevel, *quantum phenomena* are dominant, and the main features in con-nection with quantum effects are not accessible to our intuitive concepts, which are merely useful at the macroscopic level; the framework of quantum theory is in striking conflict with our intuitive demands for visualizability, and we are forced to use abstract physical laws expressed by mathematical equations. In other words, effects at the nanolevel are (almost) not accessible to our usual engineering concepts. Therefore, here we rely on the abstract mathematical relations of theoretical physics. In nanotechnology functional systems, machines and the like cannot be adequately designed without the use of these abstract theoretical laws and the application of suitable computational methods. Therefore, in nano-technology, theoretical and computational methods are centrally important: This makes the present handbook an indispensable compendium.

Nanometer-scale units are by definition very small atomic structures and functional sys-tems; it is the smallest level at which functional matter can exist. We already learned to manipulate matter at this ultimate level: Atoms can be moved experimentally in a controlled manner from one position to another. This is astonishing because one nanometer only cor-responds to one millionth of a millimeter. For example, an electrical nanogenerator could be designed consisting of various parts that included a very fast revolving rotator. One mil-lion of these generators could be arranged side by side on a length of two centimeters; it is remarkable that not only *static* nanostructures could in principle be produced and sig-nificantly manipulated but also artificial *dynamical* nanosystems. But, the downscaling of functional structures from the macroscopic to the nanometer scale is only one of the essential

points in connection with nanotechnology. In addition—and maybe much more important—nanosystems provide unique properties in comparison to those we observe at the macroscopic level. For example, a metal nanocluster shows a melting temperature that strongly deviates from that of a macroscopic piece of metal; its melting point is significantly lower. A decrease down to a fraction of only 20% is typical, depending. however, on the material and particle number.

A professional treatment of the various problems in nanoscience and nanotechnology makes the application and development of theoretical and computational methods in this field absolutely necessary. In other words, the discipline of theoretical and computational nanotechnology has to be considered as a key topic to be able to treat nanotechnology adequately and to reach optimal solutions for certain tasks. It is therefore desirable to get a timely overview about the specific topics presently relevant in this field. In this respect, the handbook gives a complete overview of the specific topics so far established in nanotechnology. Each chapter gives a certain overview of actual activities of the envisaged topic and in most cases an adequate description of the basics, so advanced students also can benefit from the handbook. It was our strategy to provide consistent and complete representations so the reader would be able to study each chapter without consulting other works. This of course leads to certain overlaps, which was also part of our strategy to enable an approach to the same topic from various points of view.

The handbook reflects the spectrum of questions and facts that are and could be relevant in the field of nanotechnology. Not only formal developments and methods are outlined, but also descriptions of a broad variety of applications in particular are typical for the handbook. All relevant topics have been taken into account, from functional structures—like an electrical nanogenerator—or quantum computing to questions that deal directly with basic physics. Almost all fields related to theoretical and computational nanotechnology could be covered, including *multiscale modeling*, which is important for the transition from microscale to nanoscale and vice versa.

All theoretical and computational methods used in connection with the various topics in nanoscience are directly based on the *same* theoretical physical laws. At the nanolevel, all properties of our world emerge at the level of the *basic* theoretical laws. In traditional technologies, engineers do not work at the ultimate level. They use more or less phenomenological descriptions that generally cannot be deduced from the basic physical theoretical laws. We have as many phenomenological descriptions as there are technological disciplines, and each is tailor-made to a specific topic. An exchange of concepts is either not possible or rather difficult. In contrast, at the ultimate nanolevel the world is based on only one theory for all disciplines, and this is expressed by basic theoretical physics. This situation opens the possibility for interconnections between the various topics in nanotechnology to bring about new effects and chances for further applications. In other words, nanotechnology and nanoscience can be considered interdisciplinary. Clearly, the handbook reflects the interdisciplinary character of this new science and technology.

The *Handbook of Theoretical and Computational Nanotechnology* includes 138 chapters written by hundreds of the world's leading scientists. Topics cover mainly the following areas:

(i) Computational biology: DNA, enzymes, proteins, biomechanisms, neurogenetic information processing, and nanomedicine

(ii) Computational chemistry: quantum chemistry, molecular design, chemical reactions, drugs, and design

(iii) Computational methods and simulation techniques from *ab initio* to multiscale modeling

(iv) Materials behavior at the nanolevel, such as mechanics, defects, diffusion, and dynamics

(v) Nanoscale processes: membranes, pores, diffusion, growth, friction, wear, catalysis

(vi) Nanostructured materials: metals, composites, polymers, liquid crystals, photonic crystals, colloids, and nanotubes

(vii) Nanostructures: fullerenes, nanotubes, clusters, layers, quantum dots, thin films, surfaces, and interfaces

(viii) Nanoengineering and nanodesign: nanomachines, nano-CAD, nanodevices, and logic circuits

(ix) Nanoelectronics: molecular electronics, nanodevices, electronic states, and nanowires
(x) Nanomagnetism: magnetic properties of nanostructures and nanostructured materials
(xi) Nanooptics: optical response theory, quantum dots, luminescence, and photonic crystals
(xii) Quantum computers: theoretical aspects, devices, and computational methods for simulating quantum computers and algorithms

The handbook provides broad information on all basic and applied aspects of theoretical and computational nanotechnology by considering more than two decades of pioneering research. It is the only scientific work of its kind since the beginning of nanotechnology, bringing together core knowledge and the very latest advances. The handbook is written for audiences of various levels while providing the latest up-to-date information to active scientists and experts in the field. This handbook is an indispensable source for research professionals and developers seeking the most up-to-date information on theoretical and computational nanotechnology among a wide range of disciplines, from science and engineering to medicine.

This handbook was written by leading experts, and we are highly grateful to all contributing authors for their tremendous efforts in writing these outstanding state-of-the-art chapters that altogether form a unified whole. K. Eric Drexler (designer of nanomachines, founder of the Foresight Institute, coiner of the term *nanotechnology*) gives an excellent introductory chapter about possible trends of future nanotechnology. We especially express our sincere gratitude to Dr. Drexler for his instructive and basic representation.

We cordially extend our special thanks to Professor Pierre-Gilles de Gennes for his valuable and insightful Foreword.

The editors are particularly thankful to Dr. Hari Singh Nalwa, President and CEO of American Scientific Publishers, for his continuous support of the project and the enthusiastic cooperation in connection with all questions concerning the development of the handbook. Furthermore, we are grateful to the entire team at Bytheway Publishing and especially to Kate Brown for copyediting.

**Dr. Michael Rieth**
**Prof. Dr. Wolfram Schommers**
Karlsruhe, Germany

# Contents

## CHAPTER 1. Computational Nanoelectronics
### Dragica Vasileska, David K. Ferry, Stephen M. Goodnick

# CHAPTER 2.  Process Simulation for Silicon Nanoelectronic Devices

*Wolfgang Windl*

# CHAPTER 3.  Electron Transport in Nanostructured Systems—*Ab Initio* Study

*Yoshiyuki Kawazoe, Hiroshi Mizuseki, Rodion Belosludov, Amir Farajian*

# CHAPTER 4.  Single-Electron Functional Devices and Circuits

*Takashi Morie, Yoshihito Amemiya*

# CHAPTER 5.  Modeling of Single-Electron Transistors for Efficient Circuit Simulation and Design

*YunSeop Yu, SungWoo Hwang, Doyeol Ahn*

# CHAPTER 6.   Electric Properties of Nanostructures

*K. Palotás, B. Lazarovits, P. Weinberger, L. Szunyogh*

# CHAPTER 7.   Transport Theory for Interacting Electrons Connected to Reservoirs

*Akira Oguri*

# CHAPTER 8.   Computational Nanotechnology: Computational Design and Analysis of Nanosize Electronic Components and Circuits

*Jerry A. Darsey, Dan A. Buzatu*

# CHAPTER 9.   Tunneling Models for Semiconductor Device Simulation

*Andreas Gehring, Siegfried Selberherr*

# CHAPTER 10.   Electronic Structure of Quantum Dots

*J. B. Wang, C. Hines, R. D. Muhandiramge*

# CHAPTER 11.   Spatiotemporal Dynamics of Quantum-Dot Lasers

*Edeltraud Gehrig, Ortwin Hess*

## CHAPTER 12.   Theoretical Investigations of Silicon Quantum Dots
*Lin-Wang Wang*

## CHAPTER 13.   Nanoscale Device Modeling
*Massimo Macucci, Luca Bonci*

# CHAPTER 14.  Wigner Function–Based Device Modeling

*Hans Kosina, Mihail Nedjalkov*

# CHAPTER 15.  Logic Design of Nanodevices

*Svetlana N. Yanushkevich*

## CHAPTER 16.   Nanoelectromechanical Systems and Modeling

*Changhong Ke, Horacio D. Espinosa*

# About the Editors

**Dr. Michael Rieth** has been a research scientist at the Institute of Materials Research I (IMF-I) in the Forschungszentrum Karlsruhe, Germany, since 2002. He has been head of the consulting company AIFT, Karlsruhe, since 1987. He worked as a researcher at the Institute of Materials Research II (IMF-II), Forschungszentrum Karlsruhe, from 1995 to 1999 and at the Engineering Science Department of the University of Patras (Greece) from 1999 to 2000. He was product manager at AMA Systems, Pforzheim, Germany, from 2000 to 2001. He received his master of science (German Dipl. Ing.) degree in electrical engineering from the University of Karlsruhe in 1991 and his doctoral degree in physics from the University of Patras (Greece) in 2001. Dr. Rieth published 23 research articles in refereed journals, 2 book chapters, and four patents. He is the author of *Nano-Engineering in Science and Technology* (World Scientific, Singapore, 2003) and was the editor-in-chief of the *Journal of Computational and Theoretical Nanoscience* (2004–2005). His main scientific interests are in atomistic modeling of metallic nanosystems and materials development for advanced fusion reactor applications.

**Prof. Dr. Wolfram Schommers** is a theoretical physicist and is presently at the Research Center of Karlsruhe in Germany. He is also professor of theoretical physics, professor of physics and materials sciences, and distinguished professor in Europe, China, and the United States. He began his studies of theoretical physics at the Technical University of Munich and continued his course work at the University of Münster, receiving a diploma in physics. After a brief intermezzo in the industry, Professor Schommers joined the Research Center of Karlsruhe. He received his doctoral degree (Dr. rer. nat.) in theoretical physics from the University of Karlsruhe.

Professor Schommers concentrates his scientific activities on computational and theoretical physics. His main fields of interest include foundations of physics, liquids, solids, and gases; superionic conductors; surface science; and nanophysics as the basis for the investigation of properties of nanometer-scale atomic devices, junctions, quantum dots, and nanomachines. He has published the results of his research and thoughts in various scientific journals (214 articles and book chapters).

Some topics concerning liquids, solids, and gases concern interaction potentials, single-particle motion, diffusion, generalized phonon density of states, collective motion, and liquid-solid phase transition. Selected articles include "The Effect of van der Waals-Type Interactions in Metals: A Pseudopotential Model" (*Zeitschrift für Physik B* 121, 1976); "Liquid Argon: The Influence of Three-Body Interactions on Atomic Correlations" (*Physical Review A* 16, 327, 1977); "Theoretical Investigation of the Liquid Solid Transition. A Study for Gallium" (*Solid State Communications* 21, 65, 1977); "Pair Potentials in Disordered Many-Particle Systems: A Study for Liquid Gallium" (*Physical Review A* 28, 3599, 1983); "Many-Body Polarization and Overlap Effects in the Dynamic Structure Factor of Dense Krypton Gas" (with P. A. Egelstaff, J. J. Salacuse, and J. Ram; *Physical Review A* 34, 1516, 1986); and "Comment on 'Pair Interaction from Structural Data for Dense Classical Liquids'" (*Physical Review Letters* 58, 427, 1987).

Topics in connection with superionic conductors involve structure and dynamics, correlated motions, and collective behavior. Selected articles include "Correlations in the Motions of Particles in AgI: A Molecular-Dynamics Study" (*Physical Review Letters* 38, 1536, 1977);

"Current–Current Correlations in AgI" (*Physical Review B* 16, 327, 1977); "Structure and Dynamics of Superionic Conductors" (*Physical Review B* 21, 847, 1979); "Triplet Correlations in Solid Electrolytes" (*Solid State Ionics* 1, 473, 1980).

Topics concerning surface physics touch on temperature effects, structure, dynamics, and interaction. Selected works are as follows: "Structural and Dynamical Behaviour of Noble-Gas Surfaces" (*Physical Review A* 32, 6845, 1985); "Statistical Mechanics of the Liquid Surface and the Effect of Premelting," in *Structure and Dynamics of Surfaces II* (Springer-Verlag, Heidelberg, 1987); "The Effect of Non-Linear Interactions at the Surface of Solids" (*Surface Science* 269/270, 180, 1992); and "Steps, Point Defects and Thermal Expansion at the Au(100) Surface" (with H. Zimmermann, M. Nold, U. Romahn, and P. von Blanckenhagen; *Surface Science* 287/288, 76, 1993).

Some details regarding the work of Professor Schommers on nanophysics include study of nanoclusters, nanostructures, and nanomachines; temperature effects; and electronic states. Selected works include "Phonons and Structure in Nano-Clusters: A Molecular Dynamics Study for Al" (*Nanostructured Materials* 9, 621, 1997); "Excited Nano-Clusters" (*Applied Physics A* 68, 187, 1999); "Thermal Stability and Specific Properties of Nanosystems" (with S. Baskoutas and M. Rieth; *Modern Physics Letters B* 14, 621, 2000); "Computational Atomic Nanodesign," in *Encyclopedia of Nanoscience and Nanotechnology* (with M. Rieth; American Scientific Publishers, Stevenson Ranch, CA, 2004); "Computational Engineering of Metallic Nanostructures and Nanomachines" (with M. Rieth; *Journal of Nanoscience and Nanotechnology* 2, 679, 2002); and "Electron in an Interaction Potential of General Shape" (with M. Rieth; *Journal of Computational and Theoretical Nanoscience* 2, 362, 2005).

Concerning the foundations of physics, Professor Schommers has discussed new aspects in connection with reality, and his basic ideas can be summarized as follows: Information about reality outside flows via sense organs into the body of the observer, and the brain forms a picture of reality. On the basis of many facts, Schommers concluded that the symbols in this picture of reality should have in general no similarity with the objects in the outside world; that is, the reality outside is transformed. On the one hand, we have the reality; on the other hand, we have a picture of reality. The reality is projected on space and time, and we obtain a picture of reality; the structures in the pictures are different from those in the reality outside. This conception is discussed mathematically by Professor Schommers in connection with quantum phenomena leading to new aspects in connection with relevant basic topics. Like both Whitehead and Bergson, Schommers argues for the primacy of processes and shows that space and time are closely tied to real processes. Selected work in this regard are "Inertial Frames of Reference: Mass Coupling to Space and Time" (*International Journal of Theoretical Physics* 20, 411, 1981); "Raum-Zeit, Quantentheorie und Bilder von der Wirklichkeit" (*Philosophia Naturalis* 23, 238, 1986); "Being and Becoming at the Microscopic Level" (*International Journal of Modern Physics B* 3, 1, 1989); "Space-Time and Quantum Phenomena," in *Quantum Theory and Pictures of Reality* (Springer-Verlag, Heidelberg, 1989); and "Truth and Knowledge," in *What Is Life?* (World Scientific, Singapore, 2002).

Professor Schommers is author and editor of the following books: *Fundamentals of Nanometer Structuring; Structure and Dynamics of Surfaces I and II; Quantum Theory and Pictures of Reality; The Visible and the Invisible; Das Sichtbare und das Unsichtbare; Elemente des Lebens; What is Life?; Formen des Kosmos; Space and Time, Matter and Mind; Symbols, Pictures and Quantum Reality.*

Professor Schommers is the editor-in-chief of the *Journal of Computational and Theoretical Nanoscience*. He is also an editorial board member of various scientific journals, and he is principal editor-in-charge of the book series *Foundations of Natural Science and Technology*. He is an invited member of the Humboldt Academy, an invited member of the Academic Board of the Humboldt Society, and an invited member of the Advisory Board of Medical Ethics of the 21st Century. Professor Schommers is also deputy governor of the American Biographical Institute (inauguration 2000).

Professor Schommers has been honored by various awards, medals, and appointments. He has been cited in *Who's Who in the World, Who's Who in Science and Technology, Living Science, The Europe 500, The Barons 500, 2000 Outstanding Intellectuals of the 21st Century, Leading Intellectuals of the World, 500 Leaders of Influence,* and *International Register of Profiles* (no. 123 of 200), and elsewhere.

# List of Contributors

Number in parentheses indicates the page on which the author's contribution begins.

**Doyeol Ahn** (319)
Institute of Quantum Information Processing and Systems, University of Seoul, Jeonnong, Dongdaemun, Seoul, Republic of Korea

**Yoshihito Amemiya** (239)
Department of Electrical Engineering, Hokkaido University, Sapporo, Japan

**Rodion Belosludov** (211)
Institute for Materials Research, Tohoku University, Sendai, Japan

**Luca Bonci** (687)
Dipartimento di Ingegneria dell'Informazione, Università degli studi di Pisa, Pisa, Italy

**Dan A. Buzatu** (437)
Division of Chemistry, National Center for Toxicological Research, Jefferson, Arkansas, USA

**Jerry A. Darsey** (437)
Department of Chemistry, University of Arkansas at Little Rock, Little Rock, Arkansas, USA

**Horacio D. Espinosa** (817)
Department of Mechanical Engineering, Northwestern University, Illinois, USA

**Amir Farajian** (211)
Institute for Materials Research, Tohoku University, Sendai, Japan

**David K. Ferry** (1)
Department of Electrical Engineering, Arizona State University, Tempe, Arizona, USA

**Edeltraud Gehrig** (605)
Advanced Technology Institute, School of Electronics and Physical Sciences, University of Surrey, Guildford, Surrey, United Kingdom

**Andreas Gehring** (469)
Institute for Microelectronics, Technical University Vienna, Vienna, Austria

**Stephen M. Goodnick** (1)
Department of Electrical Engineering, Arizona State University, Tempe, Arizona, USA

**Ortwin Hess** (605)
Advanced Technology Institute, School of Electronics and Physical Sciences, University of Surrey, Guildford, Surrey, United Kingdom

**C. Hines** (545)
School of Physics, The University of Western Australia, Crawley, Australia

**SungWoo Hwang** (319)

Department of Computer and Electronics Engineering, Korea University, Sungbuk, Seoul, Republic of Korea, and Institute of Quantum Information Processing and Systems, University of Seoul, Jeonnong, Dongdaemun, Seoul, Republic of Korea

**Yoshiyuki Kawazoe** (211)

Institute for Materials Research, Tohoku University, Sendai, Japan

**Changhong Ke** (817)

Department of Mechanical Engineering, Northwestern University, Illinois, USA

**Hans Kosina** (731)

Institute for Microelectronics, TU Vienna, Vienna, Austria

**B. Lazarovits** (363)

Center for Computational Materials Science, Technical University of Vienna, Vienna, Austria

**Massimo Macucci** (687)

Dipartimento di Ingegneria dell'Informazione, Università degli studi di Pisa, Pisa, Italy

**Hiroshi Mizuseki** (211)

Institute for Materials Research, Tohoku University, Sendai, Japan

**Takashi Morie** (239)

Graduate School of Life Science and Systems Engineering, Kyushu Institute of Technology, Kitakyushu, Japan

**R. D. Muhandiramge** (545)

School of Physics, The University of Western Australia, Crawley, Australia

**Mihail Nedjalkov** (731)

Institute for Microelectronics, TU Vienna, Vienna, Austria

**Akira Oguri** (409)

Department of Material Science, Osaka City University, Sumiyoshi-ku, Osaka, Japan

**K. Palotás** (363)

Center for Computational Materials Science, Technical University of Vienna, Vienna, Austria

**Siegfried Selberherr** (469)

Institute for Microelectronics, Technical University Vienna, Vienna, Austria

**L. Szunyogh** (363)

Center for Computational Materials Science, Technical University of Vienna, Vienna, Austria; and Department of Theoretical Physics, Center for Applied Mathematics and Computational Physics Budapest University of Technology and Economics, Budapest, Hungary

**Dragica Vasileska** (1)

Department of Electrical Engineering, Arizona State University, Tempe, Arizona, USA

**J. B. Wang** (545)

School of Physics, The University of Western Australia, Crawley, Australia

**Lin-Wang Wang** (635)

Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, California, USA

**P. Weinberger** (363)

Center for Computational Materials Science, Technical University of Vienna, Vienna, Austria

**Wolfgang Windl** (137)

Department of Materials Science and Engineering, The Ohio State University, Columbus, Ohio, USA

**Svetlana N. Yanushkevich** (765)

Department of Electrical and Computer Engineering, University of Calgary, Calgary, Alberta, Canada

**YunSeop Yu** (319)

Department of Information and Control Engineering, Hankyong National University, Anseong, Kyconggi-do, Republic of Korea

# Handbook of Theoretical and Computational Nanotechnology

## Edited by
## Michael Rieth and Wolfram Schommers

## Volume 1. BASIC CONCEPTS, NANOMACHINES, AND MEDICAL NANODEVICES

## Volume 2. ATOMISTIC SIMULATIONS—ALGORITHMS AND METHODS

## Volume 3.  QUANTUM AND MOLECULAR COMPUTING, QUANTUM SIMULATIONS

# Volume 4.   NANOMECHANICS AND MULTISCALE MODELING

# Volume 5.   TRANSPORT PHENOMENA AND NANOSCALE PROCESSES

## Volume 6.   BIOINFORMATICS, NANOMEDICINE, AND DRUG DESIGN

# Volume 7.   MAGNETIC NANOSTRUCTURES AND NANO-OPTICS

# Volume 9. NANOCOMPOSITES, NANO-ASSEMBLIES, AND NANOSURFACES

# Volume 10. NANODEVICE MODELING AND NANOELECTRONICS

# CHAPTER 1

# Computational Nanoelectronics

## Dragica Vasileska, David K. Ferry, Stephen M. Goodnick

*Department of Electrical Engineering, Arizona State University, Tempe, Arizona, USA*

## CONTENTS

# 1. NANOELECTRONICS DEFINED

Nanotechnology, by its name, implies technology at small scales, literally at nanometer scale dimensions ($10^{-9}$ m). Certainly, nanotechnology in the broadest sense is not new. Nanocrystalline materials have been used since antiquity to enhance the chemical and material properties of man-made objects. Also, the transition from the "nano" to "micro" to "macro" worlds is not abrupt, but occurs smoothly over multiple length scales. As a result, there is often confusion, and a certainly ambiguity, in what is truly "nanotechnology" as opposed to microelectronics, micromachining, and cellular biology. Somewhat arbitrarily, we define *nanometer scale* to characteristic feature sizes on the order of 100 nm, or less in terms of the separation of the micro and nanoworlds. The fact that almost all such structures contain nanoscale features in one form or another has led to nanotechnology being regarded as a somewhat broad umbrella encompassing a host of scientific and engineering disciplines.

*Nanoelectronics* generally refers to nanometer scale devices, circuits, and architectures impacting continued scaling of information processing systems, including communication and sensor systems, as well as providing an interface between the electronic and biological worlds. The present attention on nanotechnology and nanoelectronics has been driven from the top down by the continued scaling of semiconductor device dimensions into the nanometer scale regime, as discussed in more detail below. It is predicted that the scaling down of dimensions in present semiconductor technologies will continue for the next 10–12 years, until a hard limit of Moore's Law is finally reached due to manufacturability, or finally due to reaching atomic dimensions themselves. By the end of that time it will be necessary for radical new technologies to be introduced if continued progress in reducing device dimensions and increasing chip density is to be maintained. This "end of the roadmap" implies that industry faces an enormous challenge of developing commercially viable nanoscale chip technologies within the next 10 years. Fundamental advances are needed in new switching mechanisms, new computing paradigms realized from locally connected architectures such as cellular nonlinear networks (CNN), new ways to design for fault tolerance, new methods to achieve low power circuit design, and new methods for testing very dense and highly integrated nanoscale systems-on-a-chip.

From the molecular scale side or 'bottom up', the nanotechnology 'revolution' has been enabled by remarkable advances in atomic scale probes and nanofabrication tools. Structures and images at the atomic scale have been made possible by the invention of the scanning tunneling microscope (STM), and the associated atomic force microscope (AFM) [1]. Such scanning probe microscopy (SPM) techniques allow atomic scale resolution imaging of atomic positions, spectroscopic features, and positioning of atoms on a surface. Concurrently, there have been significant advances in the synthesis and control of self-assembled systems, semiconductor nanowires, molecular wires, novel states of carbon such as fullerenes and carbon nanotubes, etc. These advances have led to an explosion of scientific breakthroughs in studying the properties of individual molecular structures with potential application as components of molecular electronic (moletronic) devices and circuits. As discussed below, such bottom up technology for novel materials growth and potential device fabrication is more closely akin to the self-assembly and complex templated structure formation found in biological systems, that is, biomimetic structures.

## 1.1. Issues in Semiconductor Device Scaling

As the density of integrated circuits continues to increase, there is a resulting need to shrink the dimensions of the individual devices of which they are comprised. Smaller circuit dimensions would reduce the overall die area, thus allowing for more transistors on a single die without negatively impacting the cost of manufacturing. However, getting more functions into each circuit generally leads to larger die size, and this requires larger wafers. As semiconductor feature sizes shrink into the nanometer scale regime, device behavior becomes increasingly complicated as new physical phenomena at short dimensions occur, and limitations in material properties are reached.

For conventional silicon MOSFET device scaling, the device size is scaled in all dimensions, resulting in smaller oxide thickness, junction depth, channel length, channel width, and isolation spacing. Advances in lithography have driven device dimensions to the deep-submicrometer range, where gate lengths are currently 65 nm and below. The Semiconductor Industry Association (SIA) projects that by the end of 2009, leading edge production devices will employ 25 nm gate lengths and have oxide thickness of 1.5 nm, or less [2]. In fact, laboratory MOSFET devices with gate lengths down to 15 nm have been reported at the time of the present review, which exhibit excellent I–V characteristics [3]. Beyond that, there has been extensive work over the past decade related to nanoelectronic or quantum scale devices which operate on very different principles from conventional MOSFET devices, but may allow the continued scaling beyond the end of the current scaling roadmap [4]. This trend has been motivated by the fact that the performance of the scaled devices in the 25-nm regime and below is itself problematic [5], as discussed below.

For example, to enhance device performance, the gate oxide thickness has to be aggressively scaled. However, as the gate oxide thickness approaches 1 nm through scaling, tunneling through the gate oxide results in unacceptably large off-state currents, dramatically increasing quiescent power consumption [6], and rendering the device impractical for analog applications due to unacceptable noise levels. Another consequence of scaling is that the stack of layered materials that comprise electronic devices is becoming more like a continuum of interfaces rather than a stack of bulk thin films. Therefore, topology effects arising from surface (interface)-to-surface (interface) interactions now dominate the formation of potential barriers at interfaces. Interface inhomogeneity effects include morphological and compositional inhomogeneities. Morphological inhomogeneities, typically manifested as atomic-scale roughness, are often responsible for increased leakage currents in MOSFET gates and degraded transport properties. Fluctuations in the elemental distribution are expressions of compositional inhomogeneities. For finite dimensions and number of atoms, interface domains cannot be represented as superpositions of a few homogeneous thin film regions. Instead, the challenge of characterizing this complex system requires accurate atomic level information about the three-dimensional structure, geometry and composition of atomic-scale interfaces.

The anticipated replacement of silicon dioxide by new gate oxide materials may ensure a combination of superior dielectric properties (small, effective oxide thickness) with very low leakage current density. However, the diffusion of oxygen through the growing metal oxide films and the consequent oxidation of silicon lead to the formation of a multilayer stack of materials with different electrical properties. Current CAD tools for device design are not able to describe the complex electronic structure and transport behavior needed in these structures to provide engineers a reasonable estimate of the tunneling current at the different gate voltages. This leaves expensive experiments as the only option to derive these highly important electric characteristics of new devices.

Yet another issue that will pose serious problems for the operation of future ultra-small devices is related to substrate doping, particularly the high values required in bulk Si devices for control of the channel in short gate-length devices. The distribution of dopants is traditionally treated as a continuum in semiconductor physics, which implies: (1) the number of impurity atoms is small as compared to the total number of atoms in the semiconductor matrix, and (2) the impurity atom distribution is statistically uniform, while the position of an individual atom in the lattice is not defined, for example, is random. The assumption of statistical uniformity requires large number of atoms, which is not

the case in, for example, a 25-nm MOSFET device in which one has less than one hundred dopant atoms in the junction region. In these future ultrasmall devices, the number and location of each dopant atom will play an important role in determining the overall device behavior. The challenge of precisely placing small number of dopants may represent an insurmountable barrier, which could end conventional MOSFET scaling. Even in undoped alternative device concepts, the presence of a single impurity may significantly influence the device threshold voltage, and its on-state current, based on the location of the unintentional dopant atom. Hence, a full statistical analysis must be performed on the role of the randomly placed impurity atoms within the channel region of the device on the device subthreshold slope, threshold voltage, and transconductance degradation, as these parameters are exploited in circuit design for fault-tolerant nano-system integrated circuits.

Quantum-mechanical effects, due to, for example, spatial quantization in the device channel region, also play a significant role in the operation of nanoscale field-effect transistors. In MOSFETs, such quantization influences two basic quantities affecting device performance: (1) the channel charge induced by the gate at the surface of the substrate and (2) the carrier transport from source to drain along the channel. In the second case, because of the two-dimensional (2D) confinement of carriers in the channel, the mobility (or microscopically speaking, the carrier scattering) will be different from the three-dimensional (3D) case. Theoretically speaking, the 2D mobility should be larger than its 3D counterpart due to the reduced density of states function, that is, reduced number of final states into which the carriers can scatter. It is important to note, however, that in the smallest size devices, carriers experience very little or no scattering at all (ballistic limit), which makes this second issue less critical when modeling nano-scale devices.

Regarding channel charge, quantum effects in the surface potential will have a significant impact on both the amount of charge which can be induced by the gate electrode through the gate oxide and the profile of the channel charge in the direction perpendicular to the surface. The critical parameter in this direction is the gate-oxide thickness, which for a 25-nm MOSFET device, as noted earlier, is on the order of 1 nm. One of the most important consequences of quantization effects in the MOSFET channel is the displacement of the charge carriers away from the interface proper due to the vanishing of the envelope function there. This displacement gives rise to an effective increase in oxide thickness, and hence a reduction of the total gate capacitance. Additional degradation of the total gate capacitance arises from polydepletion effects (if a poly-Si gate is used), thus leading to an additional capacitance component in series with the oxide and the inversion layer capacitances, which reduces drive current. The total gate capacitance degradation, on the other hand, when combined with the quantum-mechanical band-gap widening effect and reduced density of states for a quasi–two-dimensional (Q2D) system, gives rise to a reduction of the sheet electron density. This effect, in turn, increases the threshold voltage and, at the same time, degrades the device transconductance. Hence, to properly describe the operation of future ultra-small devices, it is clearly necessary to incorporate quantum-mechanical effects into device simulators as described in more detail in Sections 5 and 6 of this review article.

## 1.2. Nonclassical and Quantum Effect Devices

To fabricate devices beyond current scaling limits, CMOS technology is rapidly moving toward quasi-3D structures such as dual-gate, tri-gate, and Fin-FET structures [7], in which the active channel is increasingly a nanowire or nanotube rather than bulk region. Such 3D gate structures are needed to maintain charge control in the channel, as channel lengths scale toward nanometer dimensions. The heavily doped Si substrate is increasingly being replaced by SiGe and Si on insulator technology as well.

Beyond field effect transistors, there have been numerous studies over the past two decades of alternatives to classical CMOS at the nanoscale. As dimensions become shorter than the phase-coherence length of electrons, the quantum mechanical wave nature of electrons becomes increasingly apparent, leading to phenomena such as interference, tunneling,

and quantization of energy and momentum as discussed earlier. Indeed, for a one-dimensional wire, the system may be considered a waveguide with "modes," each with a conductance less than or equal to a fundamental constant $2e^2/h$. Such quantization of conductance was first measured in split-gate field-effect transistors at low temperatures [8, 9], but manifestations of quantized conductance appear in many transport phenomena such as universal conductance fluctuations [10] and the quantum Hall effect [11]. While various early schemes were proposed for quantum interference devices based on analogies to passive microwave structures (see, for example, [12–14]), most suffer from difficulty in control of the desired waveguide behavior in the presence of unintentional disorder. This disorder can arise from the discrete impurity effects discussed earlier, as well as the necessity for process control at true nanometer scale dimensions. More recently, promising results have been obtained on ballistic $Y$-branch structures [15], where nonlinear switching behavior has been demonstrated even at room temperature [16].

In the previous section, the role of discrete impurities as an undesirable element in the performance of nanoscale FETs was detailed. However, the discrete nature of charge in individual electrons, and control of charge motion of single electrons has in fact been the basis of a great deal of research in single electron devices and circuits (see, for example, Ref. [17]). The understanding of single electron behavior is most easily provided in terms of the capacitance, $C$, of a small tunnel junction, and the corresponding change in electrostatic energy, $E = e^2/2C$, when an electron tunnels from one side to the other. When physical dimensions are sufficiently small, the corresponding capacitance (which is a geometrical quantity in general) is correspondingly small, so that the change in energy is greater than the thermal energy, resulting in the possibility of a "Coulomb blockade," or suppression of tunnel conductance due to the necessity to overcome this electrostatic energy. This Coulomb blockade effect allows the experimental control of electrons to tunnel one by one across a junction in response to a control gate bias (see, for example, Refs. [4, 18]). Single-electron transistors [19], turnstiles [20, 21], and pumps [22] have been demonstrated, even at room temperature [23]. Computer-aided modeling tools have even been developed based on Monte Carlo simulation of charge tunneling across arrays of junctions to facilitate the design of single-electron circuits [24]. As in the case of quantum interference devices, the present-day difficulties arise from fluctuations due to random charges and other inhomogenieties, as well as the difficulty in realizing lithographically defined structures with sufficiently small dimensions to have charging energies approaching $kT$ and above.

There has been rapid progress in realizing functional nanoscale electronic devices based on self-assembled structures such as semiconductor nanowires (NWs) [25] and carbon nanotubes (CNTs) [26]. Semiconductor nanowires have been studied over the past decade in terms of their transport properties [4], and for nanodevice applications such as resonant tunneling diodes [27], single electron transistors [28, 29], and field effect structures [25]. Recently, there has been a dramatic increase in interest in NWs due to the demonstration of directed self-assembly of NWs via in situ epitaxial growth [30, 31]. Such semiconductor NWs can be elemental (Si,Ge) or III–V semiconductors, where it has been demonstrated that such wires may be controllably doped during growth [32], and abrupt compositional changes forming high quality 1D heterojunctions can be achieved [33, 34]. Nanowire FETs, bipolar devices and complementary inverters have been synthesized using such techniques [35, 36]. The ability to controllably fabricate heterostructure nanowires has led to demonstration of nanoelectronic devices such as resonant tunneling diodes [37] and single electron transistors [38]. The scalability of arrays of such nanowires to circuits and architectures has also begun to be addressed [39], although the primary difficulty is in the ability to grow and orient NWs with desired location and direction.

Likewise, CNTs have received considerable attention due to the ability to synthesize NTs with metallic, semiconducting and insulating behavior, depending primarily on the chirality (i.e., how the graphite sheets forming the structure of the CNT wrap around and join themselves) [40]. Complementary $n$- and $p$-channel transistors have been fabricated from CNTs, and basic logic functions demonstrated [41]. The primary difficulty faced today is the directed growth of CNTs with the desired chirality, and positioning on a semiconductor surface, suitable for large-scale production.

## 1.3. Computational Nanoelectronics

In the previous sections, a brief synopsis was given of developments in nanoelectronics ranging from the rapid scaling of present-day field-effect transistors to nanoscale dimensions, through alternative devices and potential architectures based on quantum interference and single-electron behavior, to self-assembly of potential components for nanoelectronic devices based on nanowires and nanotubes.

In addition to the problems related to the actual operation of such ultrasmall devices, the reduced feature sizes require more complicated and time-consuming manufacturing processes. This fact signifies that a pure trial-and-error approach to device optimization will become impossible since it is both too time consuming and too expensive. Because computers are considerably cheaper resources, simulation is becoming an indispensable tool for the device engineer. Besides offering the possibility to test hypothetical devices that have not (or could not) yet been manufactured, simulation offers unique insight into device behavior by allowing the observation of phenomena that can not be measured on real devices. *Computational nanoelectronics* in this context refers to the physical simulation of nanoscale devices in terms of charge and thermal transport and the corresponding electrical behavior of devices at the nanoscale.

In general, transport itself is a phenomena described theoretically at the microscopic level by methods of nonequilibrium statistical mechanics. As discussed in more detail in Section 7, many nanoelectronic devices have a generic three-terminal representation as illustrated in Fig. 1, in which an active device region (whose properties are controlled by a separate gate contact), is coupled to left (source) and right (drain) contacts that source and sink carriers from the active region. The whole nonequilibrium device structure is coupled to the "environment" with which energy and particles are exchanged. A complete microscopic description that includes the full details of the coupled active region, contacts and environment is computationally intractable; hence, idealized simplifications of this model must be employed. The environment itself most often is treated as an infinite heat bath which transfers energy and particles into and out of the system represented by the device itself. The contacts are most often idealized as equilibrium conductors characterized by quasi-equilibrium chemical potentials, $\mu_s$ and $\mu_D$, as shown in Fig. 1, from which particles are injected and absorbed from the active region, the net flux of which represents the current flowing through the device. Hence, the most detailed level of transport modeling occurs in the active device regions. This approximate decomposition of the device into an active region coupled to ideal contacts and environment becomes increasingly inaccurate as device dimensions approach the nanoscale, where device, contacts and the environment are all strongly coupled.

In Section 1.1, we have detailed how the physics of device behavior becomes increasingly complicated as nanoscale dimensions are approached. The goal of *Computational Nanoelectronics* is to provide simulation tools with the necessary level of sophistication to



Figure 1. Conceptual representation of a generic three-terminal device composed of left and right contacts, gate, and active region.

capture the essential physics, while at the same time minimizing the computational burden, so that results may be obtained within a reasonable time frame. Figure 2 illustrates the main components of device simulation at any level. There are two main kernels, which must be solved self-consistently with one another, the transport equations governing charge flow, and the fields driving charge flow. Both are coupled strongly to one another and hence must be solved simultaneously. The fields arise from external sources, as well as the charge and current densities which act as sources for the time varying electric and magnetic fields obtained from the solution of Maxwell's equations. Under appropriate conditions, only the quasi-static electric fields arising from the solution of Poisson's equation are necessary.

The fields in turn are driving forces for charge transport. As shown in the upper part of Fig. 2, a fundamental description of transport requires knowledge of the electronic states in a device structure or material (which in general is modified as the system is driven out of equilibrium), as well as the coupling of the charge carriers with the host material through energy exchange and scattering with the vibrational modes of the atoms in the environment. These vibrational modes are the basis for thermal transport of heat, which is often as important a problem as charge transport itself in device behavior.

Figure 3 illustrates the various levels of approximation for treating transport in devices within a hierarchical structure ranging from compact modeling at the top to an exact quantum mechanical description at the bottom. At the beginning of semiconductor technology development, the device electrical characteristics could be estimated using simple one-dimensional (1D) analytical models (the gradual channel approximation for MOSFETs), based on the so-called drift-diffusion (DD) formalism. Various approximations had to be made to obtain closed-form solutions, but the resulting nonlinear models captured the basic device features. These approximations include simplified doping profiles and device geometries. Such nonlinear models form the basis of circuit analysis and design using compact models in an equivalent circuit representation.

With ongoing refinements and improvements in technology in terms due to scaling down of device dimensions, these simplified 1D approximations lost their basis as an accurate predictor of device behavior, and a more accurate description was required. This goal was realized by solving the DD equations numerically. Numerical simulation of carrier transport in semiconductor devices dates back to the famous work of Scharfetter and Gummel [42], who proposed a robust discretization of the DD equations which is still in use today.

However, as semiconductor device dimensions scaled into the submicrometer regime, the assumptions underlying the DD model lost their validity. Therefore, the transport models have been continuously refined and extended to more accurately capture transport phenomena occurring in these devices. The need for refinement and extension is primarily caused by the ongoing feature size reduction in state-of-the-art technology. As the supply voltages cannot be scaled accordingly without jeopardizing the circuit performance, the electric field inside the devices has increased. A large electric field which rapidly changes over small length scales gives rise to nonlocal and hot-carrier effects that begin to dominate device



Figure 2. Schematic description of the device simulation sequence. Reprinted with permission from D. Vasileska et al., *Mater. Sci. Eng. R* 38, 181 (2002). © 2002, Elsevier.

| | Model | Improvements |
|---|---|---|
| | Compact models | Appropriate for circuit design |
| | Drift-diffusion equations | Good for devices down to 0.5 μm, include μ(E) |
| | Hydrodynamic equations | Velocity overshoot effect can be treated properly |
| | Boltzmann transport equation monte Carlo/CA methods | Accurate up to the classical limits |
| | Quantum hydrodynamics | Keep all classical hydrodynamic features + quantum corrections |
| | Quantum monte Carlo/CA methods | Keep all classical features + quantum corrections |
| | Quantum kinetic equation Liouville, Wigner Boltzmann | Accurate up to single particle description |
| | Green's functions method | Includes correlations in both space and time domain |
| | Direct solution of the *n*-body Schrödinger equation | Can be solved only for small number of particles |

*Approximate* → *Exact* (Semiclassical approaches / Quantum approaches)
*Easy, fast* → *Difficult*

**Figure 3.** Illustration of the hierarchy of transport models. Reprinted with permission from D. Vasileska and S. Goodnick, "Encyclopedia of Materials: Science and Technology," Elsevier, 2001, p. 1. © 2001, Elsevier.

performance. An accurate description of these phenomena is required and is becoming a primary concern for industrial applications.

To overcome some of the limitations of the DD model, extensions have been proposed which basically add an additional balance equation for the average carrier energy [4]. Furthermore, an additional driving term is added to the current relation, which is proportional to the gradient of the carrier temperature. However, a vast number of these models exist, and there is a considerable amount of confusion as to their relation to each other. It is now a common practice in industry to use standard hydrodynamic models in trying to understand the operation of as-fabricated devices, by adjusting any number of phenomenological parameters (e.g., mobility, impact ionization coefficient, etc.). However, such tools do not have predictive capability for ultrasmall structures, for which it is necessary to relax some of the approximations in the Boltzmann transport equation. Therefore, one needs to move downward to the quantum transport area in the hierarchical map of transport models shown in Fig. 3, where, at the very bottom, we have the Green's function approach. The latter is the most exact, but at the same time the most difficult of all. In contrast to, for example, the Wigner function approach (which is Markovian in time), the Green's functions method allows one to consider simultaneously correlations in space and time, both of which are expected to be important in nanoscale devices. However, the difficulties in understanding the various terms in the resultant equations and the enormous computational burden needed for its actual implementation make the usefulness in understanding quantum effects in actual devices of limited value. For example, the only successful utilization of the Green's function approach commercially is the NEMO (nanoelectronics modeling) simulator [43], which is primarily 1D.

From the discussion above it follows that, contrary to the recent technological advances discussed in Section 1.1, the present state of the art in device simulation is currently lacking in the ability to treat these new challenges in scaling of device dimensions from conventional down to quantum scale devices. For silicon devices with active regions below 0.2 microns in diameter, macroscopic transport descriptions based on drift-diffusion models (see Fig. 3) are clearly inadequate. As already noted, even standard hydrodynamic models do not usually provide a sufficiently accurate description since they neglect important contributions from the tail of the phase space distribution function in the channel regions [44, 45]. Within the

requirement of self-consistently solving the coupled transport-field problem in this emerging domain of device physics, there are several computational challenges, which limit this ability. One is the necessity to solve both the transport and the Poisson's equations over the full 3D domain of the device (and beyond if one includes radiation effects). As a result, highly efficient algorithms targeted to high-end computational platforms (most likely in a multi-processor environment) are required to fully solve even the appropriate field problems. The appropriate level of approximation necessary to capture the proper nonequilibrium transport physics relevant to a future device model is an even more challenging problem both computationally and from a fundamental physics framework.

In this review chapter, we give an overview of the basic techniques used in the field of computational electronics and nanoelectronics related to nanoscale device simulation. Since electronic structure is the basic input to transport as illustrated in Fig. 2, we begin with a review of the electronic bandstructure of semiconductors, and the associated dynamics of carriers under external fields (Section 2). This allows one to calculate relevant material parameters, such as effective masses and effective density-of-states function. Afterward, we present a discussion of the basic equations governing transport in semiconductors, leading to the description of the Monte Carlo (MC) method for the solution of the semiclassical Boltzmann transport equation (BTE) (Section 3.1), and the hydrodynamic and drift-diffusion models for device simulation, that follow from moments of the BTE (Section 3.3). In Sections 4.1 and 4.2, we give an overview of field solvers for both high-frequency (solution of the Maxwell equations) and low-frequency (solution of quasi-static Poisson equation) applications, respectively. Some key elements of particle-based simulation, such as grid-size and time-step criteria, charge-assignment scheme, inclusion of the short-range electron-electron and electron-ion interactions, are described in Section 5.1. In Section 5.2, we give an overview of commercially available drift-diffusion/hydrodynamics device simulators. The simulation of the optoelectronic and high-frequency devices via the solution of the full set of Maxwell's equations coupled with a Monte Carlo transport kernel is discussed in Section 5.3. The inclusion of quantum corrections into particle-based simulators, using the effective potential approach, is discussed in Section 6.1. A brief description of the quantum hydrodynamic model (QHD) for device simulation and its application to modulation-doped high-electron mobility transistors (HEMTs) is given in Section 6.2. An overview of some of the major quantum transport approaches listed in the hierarchy of Fig. 3, as well as issues in quantum transport, are then discussed in Section 7.

## 2. ELECTRONIC STRUCTURE CALCULATION

The basis for discussing transport in semiconductors and other crystalline solids is the under-lying electronic *band* structure of the material arising from the solution of the many-body Schrödinger equation in the presence of the periodic potential of the lattice, as discussed in a host of solid-state physics textbooks. The solutions in the presence of the periodic potential of the lattice are in the form of Bloch functions.

$$\psi_{n,\mathbf{k}} = u_n(\mathbf{k})e^{i\mathbf{k}\cdot\mathbf{r}} \tag{1}$$

where $\mathbf{k}$ is the wavevector, and $n$ labels the band index corresponding to different solutions for a given wave vector. The cell-periodic function, $u_n(\mathbf{k})$, has the periodicity of the lattice and modulates the traveling wave solution associated with free electrons.

A brief look at the symmetry properties of the eigenfunctions greatly enhances the under-standing of the evolution of the band structure (Fig. 4). First, consider the energy eigenvalues of the individual atoms that constitute the semiconductor crystal. Almost all semiconductors have tetrahedral bonds corresponding to $sp^3$ hybridization. However, the individual atoms are comprised of outermost (valence) electrons in $s$- and $p$-type orbitals. The symmetry (or geometric) properties of these orbitals are made most clear by considering the angular part of the wave function of each:

$$s = 1$$

$$p_x = \frac{x}{r} = \sqrt{3}\sin\theta\cos\varphi$$

Figure 4. The typical semiconductor bandstructure. For direct-gap semiconductors, the conduction band state at k = 0 is s-like. The valence band states are linear combinations of p-like orbitals. For indirect-gap semiconductors, however, even the conduction band minima states are an admixture of s- and p-like states.

$$p_x = \frac{x}{r} = \sqrt{3}\sin\theta\sin\varphi \tag{2}$$

$$p_z = \frac{z}{r} = \sqrt{3}\cos\theta$$

We denote these states by $|S\rangle$, $|X\rangle$, $|Y\rangle$, and $|Z\rangle$. When the individual atoms form in a crystal, the valence electrons hybridize into $sp^3$ orbitals that lead to tetrahedral bonding. The crystal develops its own band structure with gaps and allowed bands. For semiconductors, one is typically worried about the band structure of the conduction and the valence bands only. Because of symmetry, the states near the center of the Brillouin zone ($k = 0$) are very similar to the $|S\rangle$ and three $p$-type states of the individual atoms, as shown in Fig. 4.

Electronic band structure calculation methods can be grouped into two general categories [46]. The first category consists of ab initio methods, such as Hartree-Fock or density functional theory (DFT), which calculate the electronic structure from first principles, that is, without the need for empirical fitting parameters. In general, these methods utilize a variational approach to calculate the ground state energy of a many-body system, where the system is defined at the atomic level. The original calculations were performed on systems containing a few atoms. Today, calculations are performed using approximately 1000 atoms but are computationally expensive, sometimes requiring massively parallel computers.

In contrast to ab initio approaches, the second category consists of empirical methods, such as the orthogonalized plane wave (OPW) [47], tight-binding [48] (also known as the linear combination of atomic orbitals (LCAO) method), the k · p method [49], and the local [50] or the nonlocal [51] empirical pseudopotential method (EPM). These methods involve empirical parameters to fit experimental data such as the band-to-band transitions at specific high-symmetry points derived from optical absorption experiments. The appeal of these methods is that the electronic structure can be calculated by solving a one-electron Schrödinger wave equation (SWE). Thus, empirical methods are computationally less expensive than ab initio calculations, and provide a relatively easy means of generating the electronic band structure necessary for transport calculations. Because of their wide spread usage in this context, in the rest of this section we will review some of the most commonly used techniques, namely, the empirical pseudopotential method, the tight-binding and the k · p method. The empirical pseudopotential method is described in Section 2.1, the tight-binding is discussed in Section 2.2, and the k · p method is described in Section 2.3. Applications of the k · p method are given in Section 2.4, which is followed by solutions of the effective mass Schrödinger equation for metal-oxide-semiconductor devices and for heterostructures (Section 2.5). We finish this chapter by a brief description of the carrier dynamics that is given in Section 2.6.

*Spin-Orbit Coupling.* Before proceeding with the description of the various empirical band structure methods, it is useful to introduce the spin-orbit interaction Hamiltonian.

The effects of spin-orbit coupling are most easily considered by regarding the spin-orbit interaction energy $H_{so}$ as a perturbation. In its most general form, $H_{so}$ operating on the wave functions $\psi_k$ is then given by

$$H_{so} = \frac{\hbar}{4m^2c^2}[\nabla V \times \mathbf{p}] \cdot \sigma \tag{3}$$

where $V$ is the potential energy term of the Hamiltonian, and $\sigma$ is the Pauli spin tensor. It can also be written in the following form as an operator on the cell-periodic function:

$$H_{so} = \frac{\hbar}{4m^2c^2}[\nabla V \times \mathbf{p}] \cdot \sigma + \frac{\hbar^2}{4m^2c^2}[\nabla V \times \mathbf{k}] \cdot \sigma \tag{4}$$

The first term is $k$-independent and is analogous to the atomic spin-orbit splitting term. The second term is proportional to $\mathbf{k}$ and is the additional spin-orbit energy arising from the crystal momentum. Rough estimates indicate that the effect of the second term on the energy bands is less than 1% of the effect of the first term. The relatively greater importance of the first term comes from the fact that the velocity of the electron in its atomic orbit is very much greater than the velocity of a wave packet made up of wave vectors in the neighborhood of $\mathbf{k}$.

The spin-orbit splitting occurs in semiconductors in the valence band because the lower-lying valence electrons are more tightly bound to the nucleus, just like electrons around the proton in the hydrogen atom. Furthermore, we can make some predictions about the magnitude of the splitting—in general, the splitting should be greater for crystals whose constituent atoms have a higher atomic number—since the field in the vicinity of the nuclei is much greater due to the greater number of protons. In fact, the spin-orbit splitting energy, $\Delta$, of semiconductors increases as the fourth power of the atomic number (i.e., number of protons) of the constituent elements. In Fig. 5, the spin-orbit splitting energy $\Delta$ is plotted against an average atomic number and a rough fit using power law is used.

**Rashba and Dresselhaus Spin Splitting.** The manipulation of the spin of charge carriers in semiconductors is one of the key problems in the field of *spintronics* [52]. In the paradigmatic spin transistor proposed by Datta and Das [53], the electron spins, injected from a ferromagnetic contact into a two-dimensional electron system, are controllably rotated during their passage from source to drain by means of the Rashba spin-orbit coupling [54]. The coefficient, $\alpha$, which describes the strength of the Rashba spin-orbit coupling, and hence the degree of rotation, is dependent on the average electric field which can be tuned externally by a gate voltage. This coupling stems from the inversion asymmetry of the confining potential of two-dimensional electron (or hole) systems. In addition to the Rashba coupling, caused by structural inversion asymmetry (SIA), a Dresselhaus type of



Figure 5. The spin-orbit splitting energy, $\Delta$, for different semiconductors plotted against the average atomic number $Z_{av}$.

coupling also contributes to the spin-orbit interaction [55]. The latter is due to bulk inversion asymmetry (BIA), and the interface inversion asymmetry (IIA). The BIA and the IIA contributions are phenomenologically inseparable and described below by the generalized Dresselhaus parameter $\beta$. Both Rashba and Dresselhaus couplings result in spin splitting of the bands and give rise to a variety of spin-dependent phenomena that allow one to evaluate the magnitude of the total spin splitting of electron subbands.

However, it is not usually possible to extract the relative contributions of Rashba and Dresselhaus terms to the spin-orbit coupling. To obtain the Rashba coefficient, $\alpha$, the Dresselhaus contribution is normally neglected. At the same time, the Dresselhaus and Rashba terms can interfere in such a way that macroscopic effects vanish though the individual terms are large. For example, both terms can cancel each other, resulting in a vanishing spin splitting in certain $k$-space directions. This cancellation leads to the disappearance of antilocalization, the absence of spin relaxation in specific crystallographic directions, and the lack of SdH beating. In Ref. [56], the importance of both Rashba and Dresselhaus terms was pointed out: turning $\alpha$ such that $\alpha = \beta$ holds, allows one to build a nonballistic spin-effect transistor.

The consequences of the Rashba and Dresselhaus terms on the electron dispersion and on the spin orientation of the electronic states of the two-dimensional electron gas are summarized below. If we consider QWs of the zinc-blende structure grown in the [001] direction, the spin-orbit part of the total Hamiltonian contains the Rashba as well as the Dresselhaus term that are calculated according to

$$\alpha(\sigma_x k_y - \sigma_y k_x) + \beta(\sigma_x k_x - \sigma_y k_y) \tag{5}$$

where $\mathbf{k}$ is the electron wave vector, and $\sigma$ is the vector of the Pauli matrices. Here, the $x$-axis is aligned along the [100] direction, $y$-axis is aligned along the [010] direction and $z$-axis is the growth direction. Note that this Hamiltonian contribution contains only terms linear in $\mathbf{k}$. As confirmed experimentally [57], terms cubic in $\mathbf{k}$ change only the strength of $\beta$ leaving the Hamiltonian unchanged.

## 2.1. The Empirical Pseudopotential Method

The concept of pseudopotentials was introduced by Fermi [58] to study high-lying atomic states. Afterward, Hellman proposed that pseudopotentials be used for calculating the energy levels of the alkali metals [59]. The widespread usage of pseudopotentials did not occur until the late 1950s, when activity in the area of condensed matter physics began to accelerate. The main advantage of using pseudopotentials is that only valence electrons have to be considered. The core electrons are treated as if they are frozen in an atomic-like configuration. As a result, the valence electrons are treated as moving in a weak one-electron potential.

The pseudopotential method is based on the orthogonalized plane wave (OPW) method due to Herring [47]. In this method, the crystal wave function $\psi_{\mathbf{k}}$ is constructed to be orthogonal to the core states. This is accomplished by expanding $\psi_{\mathbf{k}}$ as a smooth part of symmetrized combinations of Bloch functions $\varphi_{\mathbf{k}}$, augmented with a linear combination of core states. This expansion is expressed as

$$\psi_{\mathbf{k}} = \varphi_{\mathbf{k}} + \sum_i b_{\mathbf{k},i} \Phi_{\mathbf{k},i} \tag{6}$$

where $b_{\mathbf{k},i}$ are orthogonalization coefficients and $\Phi_{\mathbf{k},i}$ are core wave functions. For Si-14, the summation over $i$ in Eq. (6) is a sum over the core states, $1s^2 2s^2 2p^6$. Since the crystal wave function is constructed to be orthogonal to the core wave functions, the orthogonalization coefficients can be calculated, thus yielding the final expression

$$\psi_{\mathbf{k}} = \varphi_{\mathbf{k}} - \sum_i \langle \Phi_{\mathbf{k},i} \mid \varphi_{\mathbf{k}} \rangle \Phi_{\mathbf{k},i} \tag{7}$$

To obtain a wave equation for $\varphi_{\mathbf{k}}$, the Hamiltonian operator,

$$H = \frac{p^2}{2m} + V_c \tag{8}$$

is applied to Eq. (6), where $V_c$ is the attractive core potential, and the following wave equation results:

$$\left(\frac{p^2}{2m} + V_c + V_R\right)\varphi_k = E\varphi_k \qquad (9)$$

where $V_R$ represents a short-range, non-Hermitian repulsion potential, of the form

$$V_R = \sum_t \frac{(E - E_t)\langle\Phi_{k,t} \mid \varphi_k\rangle\Phi_{k,t}}{\varphi_k} \qquad (10)$$

$E_t$ in Eq. (10) represents the atomic energy eigenvalue, and the summation over $t$ represents a summation over the core states. The result given in Eq. (9) can be thought of as the wave equation for the pseudowave function, $\varphi_k$, but the energy eigenvalue $E$ corresponds to the true energy of the crystal wave function, $\psi_k$. Furthermore, as a result of the orthogonalization procedure, the repulsive potential, $V_R$, which serves to cancel the attractive potential, $V_c$, is introduced into the pseudo–wave function Hamiltonian. The result is a smoothly varying pseudopotential $V_P = V_C + V_R$. This result is known as the Phillips-Kleinman cancellation theorem [60], which provides justification why the electronic structure of strongly bound valence electrons can be described using a nearly free electron model and weak potentials.

To simplify the problem further, model pseudopotenials are used in place of the actual pseudopotential. Figure 6 summarizes the various models employed. Note that the 3D Fourier transforms (for bulk systems) of each of the above-described model potentials are of the following general form

$$V(q) \sim \frac{Ze^2}{\varepsilon_0 q^2}\cos(qr_c) \qquad (11)$$



(a) Constant effective potential in the core region:

$$V(r) = \begin{cases} \dfrac{-Ze^2}{4\pi\varepsilon_0 r}; & r > r_c \\[2mm] \dfrac{-Ze^2}{4\pi\varepsilon_0 r_c}; & r \leq r_c \end{cases}$$

(b) Empty core model:

$$V(r) = \begin{cases} \dfrac{-Ze^2}{4\pi\varepsilon_0 r}; & r > r_c \\[2mm] 0; & r \leq r_c \end{cases}$$

(c) Model potential due to Heine and Abarenkov:

$$V(r) = \begin{cases} \dfrac{-Ze^2}{4\pi\varepsilon_0 r}; & r > r_c \\[2mm] A; & r \leq r_c \end{cases}$$

(d) Lin and Kleinman model potentials:

$$V(r) = \begin{cases} \dfrac{-Ze^2}{4\pi\varepsilon_0 r}\{1\text{-exp}\,[-\beta(r-r_c)]\}; & r > r_c \\[2mm] 0; & r \leq r_c \end{cases}$$

**Figure 6.** Various model potentials. Reprinted with permission from S. Gonzalez, M.S. Thesis, Arizona State University. 2002. © 2002.

This $q$-dependent pseudopotential is then used to calculate the energy band structure along different crystallographic directions, using the procedure outlined in the following section.

## 2.1.1. Description of the Empirical Pseudopotential Method

Recall from the previous section that the Phillips-Kleinman cancellation theorem provides a means for the energy band problem to be simplified into a one-electron-like problem. For this purpose, Eq. (9) can be rewritten as

$$\left(\frac{p^2}{2m} + V_P\right)\varphi_k = E\varphi_k \tag{12}$$

where $V_P$ is the smoothly varying crystal pseudopotential. In general, $V_P$ is a linear combination of atomic potentials, $V_a$, which can be expressed as a summation over the lattice translation vectors, $\mathbf{R}$, and atomic basis vectors $\tau$, to arrive at the following expression:

$$V_P(\mathbf{r}) = \sum_{\mathbf{R}}\sum_{\tau} V_a(\mathbf{r} - \mathbf{R} - \tau) \tag{13}$$

To simplify further, the inner summation over $\tau$ can be expressed as the total potential, $V_0$, in the unit cell located at $\mathbf{R}$. Eq. (13) then becomes

$$V_P(\mathbf{r}) = \sum_{\mathbf{R}} V_0(\mathbf{r} - \mathbf{R}) \tag{14}$$

Because the crystal potential is periodic, the pseudopotential is also a periodic function and can be expanded into a Fourier series over the reciprocal lattice to obtain

$$V_P(\mathbf{r}) = \sum_{G} V_0(G)e^{iG \cdot r} \tag{15}$$

where the expansion coefficient is given by

$$V_0(G) = \frac{1}{\Omega} \int d^3r V_0(\mathbf{r})e^{-iG \cdot r} \tag{16}$$

and $\Omega$ is the volume of the unit cell.

To apply this formalism to the zincblende lattice, it is convenient to choose a two-atom basis centered at the origin ($\mathbf{R} = 0$). If the atomic basis vectors are given by $\tau_1 = \tau = -\tau_2$, where $\tau$, the atomic basis vector, is defined in terms of the lattice constant $a_0$ as $\tau = a_0(1/8, 1/8, 1/8)$, $V_0(\mathbf{r})$ can be expressed as

$$V_0(\mathbf{r}) = V_1(\mathbf{r} - \tau) + V_2(\mathbf{r} + \tau) \tag{17}$$

where $V_1$ and $V_2$ are the atomic potentials of the cation and anion. Substituting Eq. (17) into Eq. (16), and using the displacement property of Fourier transforms, $V_0(\mathbf{r})$ can be recast as

$$V_0(G) = e^{iG \cdot \tau}V_1(G) + e^{-iG \cdot \tau}V_2(G) \tag{18}$$

Writing the Fourier coefficients of the atomic potentials in terms of symmetric ($V_S(G) = V_1 + V_2$)) and antisymmetric ($V_A(G) = V_1 - V_2$)) *form factors*, $V_0(G)$ is given by

$$V_0(G) = \cos(G \cdot \tau)V_S(G) + i\sin(G \cdot \tau)V_A(G) \tag{19}$$

where the prefactors are referred to as the symmetric and antisymmetric structure factors. The form factors above are treated as adjustable parameters that can be fit to experimental data, hence the name empirical pseudopotential method. For diamond-lattice materials, with two identical atoms per unit cell, the $V_A = 0$ and the structure factor is simply $\cos(G \cdot \tau)$. For zinc-blende lattice, like the one in GaAs material system, $V_A \neq 0$ and the structure factor is more complicated.

With the potential energy term specified, the next task is to recast the Schrödinger equation in a matrix form. Recall that the solution to the Schrödinger wave equation in a periodic

lattice is a Bloch function, which is composed of a plane wave component and a cell periodic part that has the periodicity of the lattice, that is,

$$\varphi_k(r) = e^{ik \cdot r} u_k(r) = e^{ik \cdot r} \sum_G U(G') e^{iG' \cdot r} \qquad (20)$$

By expanding the cell periodic part $u_k(r)$ of the Bloch function appearing in Eq. (20) into Fourier components, and substituting the pseudo–wave function $\varphi_k$ and potential $V_0$ into the Schrödinger wave equation, the following matrix equation results

$$\sum_G \left\{ \left[ \frac{\hbar^2(k+G)^2}{2m} - E \right] U(G) + \sum_{G'} V_0(|G - G'|) U(G') \right\} = 0 \qquad (21)$$

The expression given in Eq. (21) is zero when each term in the sum is identically zero, which implies the following condition

$$\left[ \frac{\hbar^2(k+G)^2}{2m} - E \right] U(G) + \sum_{G'} V_0(|G - G'|) U(G') = 0 \qquad (22)$$

In this way, the band structure calculation is reduced to solving the eigenvalue problem specified by Eq. (22) for the energy $E$. As obvious from Eq. (20), $U(G')$ is the Fourier component of the cell periodic part of the Bloch function. The number of reciprocal lattice vectors used determines both the matrix size and calculation accuracy.

The eigenvalue problem of Eq. (22) can be written in the more familiar form $\mathbf{HU} = E\mathbf{U}$, where $\mathbf{H}$ is a matrix, $\mathbf{U}$ is a column vector representing the eigenvectors, and $E$ is the energy eigenvalue corresponding to its respective eigenvector. For the diamond lattice, the diagonal matrix elements of $\mathbf{H}$ are then given by

$$H_{i,j} = \frac{\hbar^2}{2m} |k + G_i|^2 \qquad (23)$$

for $i = j$, and the off-diagonal matrix elements of $\mathbf{H}$ are given by

$$H_{i,j} = V_S(|G_i - G_j|) \cos[(G_i - G_j) \cdot \tau] \qquad (24)$$

for $i \neq j$. Note that the term $V_S(0)$ is neglected in arriving at Eq. (23), because it will only give a rigid shift in energy to the bands. The solution to the energy eigenvalues and corresponding eigenvectors can then be found by diagonalizing the matrix $\mathbf{H}$.

### 2.1.2. Implementation of the Empirical Pseudopotential Method for Si and Ge

For a typical semiconductor system, 137 plane waves are sufficient, each corresponding to vectors in the reciprocal lattice, to expand the pseudopotential. The reciprocal lattice of a face-centered cubic (FCC), that is diamond or zinc-blende structure, is a body-centered cubic (BCC) structure. Reciprocal lattice vectors up to and including the 10th-nearest neighbor from the origin are usually considered which results in 137 plane waves for the zinc-blende structure. The square of the distance from the origin to each equivalent set of reciprocal lattice sites is an integer in the set $|G^2| = 0, 3, 4, 8, 11, 12, \ldots$, where $|G^2|$ is expressed in units of $(2\pi/a_0)^2$. Note that the argument of the pseudopotential term $V_S$ in Eq. (24) is the difference between reciprocal lattice vectors. It can be shown that the square of the difference between reciprocal lattice vectors will also form the set of integers previously described. This means that $V_S$ is only needed at discrete points corresponding to nearest-neighbor sites. The pseudopotential, on the other hand, is a continuous quantity. Therefore, its Fourier transform $V_S(q)$ is also a continuous function that is shown in Fig. 7. The points corresponding to the first three nearest neighbors are also indicated on this figure.

Recall that the pseudopotential is only needed at a few discrete points along the $V(q)$ curve. The discrete points correspond to the $q^2$-values that match the integer set described previously. There is some controversy, however, regarding the value of $V_S$ as $q$ vanishes.

**Figure 7.** Fourier transform of the pseudopotential. (Note that $q = |G - G'|$). Reprinted with permission from S. Gonzalez, M.S. Thesis, Arizona State University, 2002. © 2002.

There are two common values seen in the literature: $V_1(0) = -3/2E_F$ and $V_1(0) = 0$. In most cases, the term $V_S(0)$ is ignored because it only gives a rigid shift in energy to the bands. The remaining form factors needed to compute the band structure for nonpolar materials correspond to $q^2 = 3, 8$, and $11$. For $q^2 = 4$, the cosine term in Eq. (24) will always vanish. Furthermore, for values of $q^2$ greater than $11$, $V(q)$ quickly approaches zero. This convergence comes from the fact that the pseudopotential is a smoothly varying function, and only few plane waves are needed to represent it, that is, if a function is rapidly varying in space, then many more plane waves would be required. Another advantage of the empirical pseudopotential method is that only three parameters are needed to describe the band structure of nonpolar materials.

Using the form factors listed in Table 1, where the Si form factors are taken from Ref. [61] and the Ge form factors are taken from Ref. [62], the band structures for Si and Ge are plotted in Fig. 8 [63]. Note that spin-orbit interaction is not included in these simulations. The lattice constants specified for Si and Ge are 5.43 Å and 5.65 Å, respectively. Si is an indirect band gap semiconductor. Its primary gap, that is, minimum gap, is calculated from the valence band maximum at the Γ-point to the conduction band minimum along the Δ direction, 85% of the distance from Γ to X. The band gap of Si is calculated to be $E_g^{Si} = 1.08$ eV, in agreement with experimental findings. Ge is also an indirect band-gap semiconductor. Its band gap is defined from the top of the valence band at Γ to the conduction band minimum at L. The band gap of Ge is calculated to be $E_g^{Ge} = 0.73$ eV. The direct gap, which is defined from the valence band maximum at Γ to the conduction band minimum at Γ, is calculated to be 3.27 eV and 0.82 eV for Si and Ge, respectively. Note that the curvature of the top valence band of Ge is larger than that of Si. This corresponds to the fact that the effective hole mass of Si is larger than that of Ge. Note that the inclusion of the spin-orbit interaction will lift the triple degeneracy of the bands at the Γ point, leaving doubly degenerate heavy and light-hole bands and a split-off band moved downward in energy by few tens of microelectronvolts (depending on the material under consideration).

In summary, the local empirical pseudopotential method (EPM) described in this section is rather good for an accurate description of the optical gaps. However, as noted by Chelikowsky and Cohen [64], when these local calculations are extended to yield the

**Table 1.** Local pseudopotential form factors.

| Form Factor (Ry) | Si | Ge |
|---|---|---|
| $V_3$ | −0.2241 | −0.2768 |
| $V_8$ | 0.0551 | 0.0582 |
| $V_{11}$ | 0.0724 | 0.0152 |

Figure 8. (Left panel) EPM band structure of silicon. (Right panel) EPM band structure of germanium. Reprinted with permission from S. Gonzalez, M.S. Thesis, Arizona State University, 2002. © 2002.

valence-band electronic density of states, the results obtained are far from satisfactory. The reason for this discrepancy arises from the omission of the low cores in the derivation of the pseudopotential in the previous section. This, as previously noted, allowed the usage of a simple plane wave basis. To correct for the errors introduced, an energy-dependent non-local correction term is added to the local atomic potential. This increases the number of parameters needed but leads to better convergence and more exact band-structure results [65, 66].

### 2.1.3. Empirical Pseudopotential Method for Hexagonal GaN

In the previous section, the EPM implementation for the diamond and zinc-blende material systems was explained in detail, and representative bandstructure calculations for Si and Ge were presented. The method is not limited to cubic structures, it is used as well for state-of-the-art materials such as GaN, AlN, InN, and their alloys, that crystallize in the hexagonal wurtzite structure ($\alpha$-nitrides). The current interest in the group-III nitride material system is due to their wide-bandgaps, which has application in short-wavelength optoelectronic devices (LEDs, lasers), and high-power electronic devices. The main computational difference compared to zincblende and diamond is in the set of reciprocal lattice vectors, G, defining the hexagonal lattice, as well as the more complicated atomic basis vectors, $\tau$, corresponding to four atoms per unit cell rather than two for diamond and zincblende. Generalizing the EPM method for this less symmetric material system, a limited set of pseudopotential form factors can again be introduced, resulting in the bandstructure shown in Fig. 9 [67]. In this particular example, since transport was the major goal of the calculation, the form factors were adjusted to optimize the effective mass compared, at the cost of a reduced bandgap, which has been artificially adjusted to match the experimental value. Necessary ingredients in these calculations are the use of the continuous ionic model potentials, which are screened by the model dielectric function derived for semiconductors by Levine and Louie [68]. Such an approach allows for a continuous description in the reciprocal space, the explicit inclusion of bond charges, and the exploitation of the ionic model potential transferability to other crystal structures, namely, the wurtzite crystal. In Ref. [69], it was shown by way of an example of wurtzite phase nitrides, that crystal-specific anisotropies can be taken into account via proper choice of the screening function.

### 2.2. The Tight-Binding Method

Tight binding (TB) is another semiempirical method for electronic structure calculations. While it retains the underlying quantum mechanical description of electrons in a solid, the Hamiltonian is parameterized and simplified before the calculation, rather than constructing it from first principles. The method is detailed by Slater and Koster [70], who laid the initial ground work. Conceptually, tight binding works by postulating a basis set which consists of atomic-like orbitals (i.e., they share the angular momentum components of the atomic

797

2007

n

**Figure 9.** GaN band structure calculated from via EPM (after Yamakawa et al. [67]). The direct gap has been shifted to obtain better fit of the effective masses. Reprinted with permission from [67], S. Yamakawa et al., *Comput. Electron.* 2, 481 (2003). © 2003, Springer.

orbitals and are easily split into radial and angular parts) for each atom in the system, and the Hamiltonian is then parameterized in terms of various high symmetry interactions between these orbitals. For tetrahedral semiconductors, as already noted, a conceptual basis set of one $s$-like orbital and three $p$-like orbitals has been used. In the most common form of tight binding (nearest neighbor, orthogonal TB), the orbitals are assumed to be orthogonal and interactions between different orbitals are only allowed to be non-zero within a certain distance, which is placed somewhere between the first and second nearest neighbors in the crystal structure. A further simplification made, is to neglect three-center integrals (i.e., an interaction between orbitals on atoms A and B which is mediated by the potential on atom C), meaning that each interaction is a function of the distance between the atoms only.

The quantitative description of the method presented below is due to Chadi and Cohen [48]. First denote the position of the atom in the primitive cell as

$$r_{jl} = R_j + r_l \qquad (25)$$

where $R_j$ is the position of the $j$th primitive cell and $r_l$ is the position of the atom within the primitive cell. Let $h_l(r)$ be the Hamiltonian of the $l$th isolated atom, such that

$$h_l \phi_{ml}(r - r_{jl}) = E_{ml}\phi_{ml}(r - r_{jl}) \qquad (26)$$

where $E_{ml}$ and $\phi_{ml}$ are the eigenvalues and the eigenfunctions of the state indexed by $m$. The atomic orbitals. $\phi_{ml}$, are called Löwdin orbitals [71], and they are different from the usual atomic wave functions in that they have been constructed in such a way that wave functions centered at different atomic sites are orthogonal to each other. The total Hamiltonian of the system is then

$$H_0 = \sum_{l,j} h_l(r - r_{jl}) \qquad (27)$$

Note that the sum over $l$ refers to a sum within the different atoms in the basis, therefore, $l = 1, 2$ for diamond and zinc-blende crystals. The unperturbed Bloch functions, that have the proper translational symmetry, are constructed to be of the following form:

$$\Phi_{mlk} = \frac{1}{\sqrt{N}} \sum_j e^{ir_{jl}\cdot k}\phi_{ml}(r - r_{jl}) \qquad (28)$$

The eigenvalues of the total Hamiltonian $H = H_0 + H_{int}$ (where $H_{int}$ is the interaction part of the Hamiltonian) are then represented as a linear combination of the Bloch functions

$$\Psi_L = \sum_{ml} c_{ml} \Phi_{mlk} \tag{29}$$

Operating with the total Hamiltonian of the system $H$ on $\Psi_k$, and using the orthogonality of the atomic wave functions, one arrives at the following matrix equation

$$\sum_{ml} [H_{m'l',ml} - E_k \delta_{mm'} \delta_{ll'}] c_{ml} = 0 \tag{30}$$

where the matrix element appearing in the above expression is given by

$$H_{m'l',ml}(\mathbf{k}) = \sum_i e^{i(\mathbf{R}_i + \mathbf{r}_l - \mathbf{r}_{l'}) \cdot \mathbf{k}} \langle \phi_{mlk}(\mathbf{r} - \mathbf{r}_{il}) | H | \phi_{mlk}(\mathbf{r} - \mathbf{r}_{il'}) \rangle \tag{31}$$

Note that in the simplest implementation of this method, instead of summing over all the atoms, one sums over the nearest-neighbor atoms only. Also note that the index $m$ represents the $s$- and $p$-states of the outermost electrons ($|s\rangle$, $|X\rangle$, $|Y\rangle$ and $|Z\rangle$), and $l$ is the number of distinct electrons in the basis. For the case of tetrahedrally coordinated semiconductors, the number of nearest neighbors is four and are located at

$$\begin{cases} d_1 = (1, 1, 1)\dfrac{a_0}{4} \\[2mm] d_2 = (1, -1, -1)\dfrac{a_0}{4} \\[2mm] d_3 = (-1, 1, -1)\dfrac{a_0}{4} \\[2mm] d_4 = (-1, -1, 1)\dfrac{a_0}{4} \end{cases} \tag{32}$$

For a diamond lattice, one also defines the following matrix elements:

$$\begin{cases} V_{ss} = 4V_{ss\sigma} \\[2mm] V_{sp} = -\dfrac{4V_{sp\sigma}}{\sqrt{3}} \\[2mm] V_{xx} = 4\left[\dfrac{V_{pp\sigma}}{3} + \dfrac{2V_{pp\pi}}{3}\right] \\[2mm] V_{xy} = 4\left[\dfrac{V_{pp\sigma}}{3} - \dfrac{V_{pp\pi}}{3}\right] \end{cases} \tag{33}$$

As an example, consider the matrix element between two $s$-states

$$H_{s_1, s_2} = [e^{i\mathbf{k} \cdot \mathbf{d}_1} + e^{i\mathbf{k} \cdot \mathbf{d}_2} + e^{i\mathbf{k} \cdot \mathbf{d}_3} + e^{i\mathbf{k} \cdot \mathbf{d}_4}] \langle s_1 | H_{int} | s_2 \rangle = g_1(\mathbf{k}) V_{ss} \tag{34}$$

Notice the appearance of the Bloch sum $g_1(\mathbf{k})$ in Eq. (34). This observation suggests that for different basis states, there will be four different Bloch sums, $g_1$ through $g_4$, of the form

$$\begin{cases} g_1(\mathbf{k}) = [e^{i\mathbf{k} \cdot \mathbf{d}_1} + e^{i\mathbf{k} \cdot \mathbf{d}_2} + e^{i\mathbf{k} \cdot \mathbf{d}_3} + e^{i\mathbf{k} \cdot \mathbf{d}_4}] \\[2mm] g_2(\mathbf{k}) = [e^{i\mathbf{k} \cdot \mathbf{d}_1} + e^{i\mathbf{k} \cdot \mathbf{d}_2} - e^{i\mathbf{k} \cdot \mathbf{d}_3} - e^{i\mathbf{k} \cdot \mathbf{d}_4}] \\[2mm] g_3(\mathbf{k}) = [e^{i\mathbf{k} \cdot \mathbf{d}_1} - e^{i\mathbf{k} \cdot \mathbf{d}_2} + e^{i\mathbf{k} \cdot \mathbf{d}_3} - e^{i\mathbf{k} \cdot \mathbf{d}_4}] \\[2mm] g_4(\mathbf{k}) = [e^{i\mathbf{k} \cdot \mathbf{d}_1} - e^{i\mathbf{k} \cdot \mathbf{d}_2} - e^{i\mathbf{k} \cdot \mathbf{d}_3} + e^{i\mathbf{k} \cdot \mathbf{d}_4}] \end{cases} \tag{35}$$

It is also important to note that the Hamiltonian matrix elements between $s$- and $p$-states on the same atom, or two different $p$-states on the same atom, are zero because of symmetry

in diamond and zincblende crystals. The $8 \times 8$ secular determinant representing al possible nearest-neighbor interactions between the tight-binding $s$- and $p$-orbitals centered on each atom in the crystal is

$$
\begin{array}{c|cccccccc}
 & S1 & X1 & Y1 & Z1 & S2 & X2 & Y2 & Z2 \\
\hline
S1 & E_s - E_k & 0 & 0 & 0 & V_{ss}g_1 & V_{sp}g_2 & V_{sp}g_3 & V_{sp}g_4 \\
X1 & 0 & E_p - E_k & 0 & 0 & -V_{sp}g_2 & V_{xx}g_1 & V_{xy}g_4 & V_{xy}g_3 \\
Y1 & 0 & 0 & E_p - E_k & 0 & -V_{sp}g_3 & V_{xy}g_4 & V_{xx}g_1 & V_{xy}g_2 \\
Z1 & 0 & 0 & 0 & E_p - E_k & -V_{sp}g_4 & V_{xy}g_3 & V_{xy}g_2 & V_{xx}g_1 \\
S2 & V_{ss}g_1^* & -V'_{sp}g_2^* & -V_{sp}g_3^* & -V_{sp}g_4^* & E_s - E_k & 0 & 0 & 0 \\
X2 & V_{sp}g_2^* & V_{xx}g_1^* & V_{xy}g_4^* & V_{xy}g_3^* & 0 & E_p - E_k & 0 & 0 \\
Y2 & V_{sp}g_3^* & V_{xy}g_4^* & V_{xx}g_1^* & V_{xy}g_2^* & 0 & 0 & E_p - E_k & 0 \\
Z2 & V_{sp}g_4^* & V_{xy}g_3^* & V_{xy}g_2^* & V_{xx}g_1^* & 0 & 0 & 0 & E_p - E_k
\end{array}
$$

$$\tag{36}$$

The tight-binding parameters appearing in Eqs. (33) and (36) are often obtained by comparison with empirical pseudopotential calculations, as reported in Ref. [48] (Table 2).

Using the method described one can quite accurately describe the valence bands, whereas the conduction bands are not reproduced that well due to the omission of the interaction with the higher-lying bands. The accuracy of the conduction bands can be improved with the addition of more overlap parameters. However, there are only four conduction bands and the addition of more orbitals destroys the simplicity of the method.

## 2.3. The k · p Method

In contrast to the previously described empirical pseudopotential and the tight-binding methods, the **k · p** method is based upon perturbation theory [72, 73]. In this method, the energy is calculated near a band maximum or minimum by considering the wave number (measured from the extremum) as a perturbation.

### 2.3.1. k · p General Description

For a better understanding of the method, first assume that the Schrödinger equation is one-dimensional and stationary. To further elaborate the problem, also assume that the particle sees a potential, $V = V_- + V_0$, where $V_-$ is the periodic potential that has the periodicity of the 1D lattice, and $V_0$ is the confinement potential. The one-dimensional Schrödinger wave equation is written

$$
H_0\psi(x) = \left[\frac{p^2}{2m} + V(x)\right]\psi(x) = \lambda\psi(x) \tag{37}
$$

and $V_0 = 0$ if $x \notin [-x_0, x_0]$; and $V_0 = -V_0$ otherwise. Here, $V_0$ and $x_0$ are some arbitrary positive constants. If $V_0$ is small, then the solutions to the one-dimensional Schrödinger equation are of the Bloch form (as discussed in the introduction part of this section), repeated here for completeness for a 1D case:

$$
\psi_k(x) = e^{ikx}u_k(x) \tag{38}
$$

Table 2. Chadi and Cohen tight-binding parameters [48].

| | $E_p - E_s$ | $V_{ss}$ | $V_{sp}$ | $V_{xx}$ | $V_{xy}$ |
|---|---|---|---|---|---|
| C | 7.40 | -15.2 | 10.25 | 3.0 | 8.3 |
| Si | 7.20 | -8.13 | 5.88 | 1.71 | 7.51 |
| Ge | 8.41 | -6.78 | 5.31 | 1.62 | 6.82 |

where $u_k(x)$ is cell periodic part of the Bloch function. The Schrödinger equation can then be written as

$$H_k u_k(x) = \left[\frac{p^2}{2m} + V(x) + \frac{\hbar}{m}k \cdot p\right]u_k(x) = \left[E_k - \frac{\hbar^2 k^2}{2m}\right]u_k(x) \tag{39}$$

The term $(\hbar/m)\mathbf{k} \cdot \mathbf{p}$ is treated as a perturbation to $H_0$ for determining $u_k(x)$ and $E_k$ in the vicinity of $k = 0$ in terms of the complete set of cell-periodic wave functions and energy eigenvalues at $k = 0$, which are assumed known. To simplify the form of Eq. (39), it is convenient to define

$$E'_k = E_k - \frac{\hbar^2 k^2}{2m} \tag{40}$$

To deal with this problem, we now assume that we have an orthonormal basis $\{\zeta_i\}_{i=1}^{n}$ of eigenvectors (associated to their eigenvalues, $\lambda_i$) of the operator $p^2/2m + V$, that are of a fixed parity (the orbitals may be of $s$- or $p$-type). We then project operator $H_k$ on the finite dimensional space generated by the $\zeta_i$'s, to obtain

$$\langle \zeta_i|H_k|\zeta_j\rangle = \lambda_j \delta_{ij} + \frac{\hbar}{m}k\langle \zeta_i|p|\zeta_j\rangle + \langle \zeta_i|V_J|\zeta_j\rangle$$

$$= \lambda_j \delta_{ij} + kP_{ij} + Q_{ij} \tag{41}$$

that is, we arrive at the symmetric eigenvalue matrix

$$H(k) = Q + \begin{bmatrix} \lambda_1 & \cdots & (kP_{ij}) \\ \vdots & \ddots & \vdots \\ (k\overline{P}_{ij}) & \cdots & \lambda_n \end{bmatrix} \tag{42}$$

the solutions of which provide us the eigenvalues and the corresponding eigenvectors.

## 2.3.2. k · p Theory near the Γ Point for Bulk Materials

In general, one either has a bulklike system or lower-dimensional systems such as 2D and 1D electron gases, in which there is a confinement in one and two directions, respectively. Such lower dimensional systems are frequently encountered in nanoscale devices, which makes this general discussion of the $\mathbf{k} \cdot \mathbf{p}$ method very useful. For a general system, with spin-orbit interaction included in the model, and using the result for the 1D case given in the previous section, the Schrödinger equation is of the following general form

$$\left[\frac{p^2}{2m_0} + \frac{\hbar}{4m_0^2 c^2}(\sigma \times \nabla V) \cdot \mathbf{p} + \frac{\hbar^2 k^2}{2m_0} + \frac{\hbar k}{m_0} \cdot \left(\mathbf{p} + \frac{\hbar}{4m_0 c^2}(\sigma \times \nabla V)\right)\right]u_{nk}(r) = E_{nk}u_{nk}(r) \tag{43}$$

The Hamiltonian in Eq. (43) can be divided into two terms

$$[H(\mathbf{k} = 0) + W(\mathbf{k})]u_{nk} = E_{nk}u_{nk} \tag{44}$$

where the only $\mathbf{k}$-dependence is preserved in $W(\mathbf{k})$. Next, as in the 1D case, we assume that the local, single particle solutions of the Hamiltonian $H(\mathbf{k} = 0)$ has a complete set of eigenfunctions $u_{n0}$, that is,

$$H(\mathbf{k} = 0)u_{n0} = E_{n0}u_{n0} \tag{45}$$

An arbitrary ("well behaving") lattice periodic function can be written as a series expansion using the eigenfunctions $u_{n0}$. We then insert an expansion

$$u_{nk} = \sum_m c_m^n(\mathbf{k})u_{m0} \tag{46}$$

in Eq. (44), and find matrix equation for determining the unknown coefficients, $c_m^n(\mathbf{k})$. We multiply from left by $u_{n0}^*$, integrate, and use the orthogonality of the basis functions, obtain

$$\sum_m \left[ \left( E_{n0} - E_{nk} + \frac{\hbar^2 k^2}{2m_0} \right) \delta_{nm} + \frac{\hbar \mathbf{k}}{m_0} \cdot \langle u_{n0}| \left( \mathbf{p} + \frac{\hbar}{4m_0 c^2}(\sigma \times \nabla V) \right) |u_{m0}\rangle \right] c_m^n \mathbf{k}) = 0 \quad (47)$$

Solving the above matrix equation then gives the exact eigenstates of Eq. (43). However, this looks good only in principle, the reason being that the calculation becomes increasingly complicated as $\mathbf{k}$ increases. One has to increase the number of states in the expansion given in Eq. (47), and the calculations become numerically unfeasible. Therefore, this approach is practical only for small wave vector values.

When $\mathbf{k}$ is small, the nondiagonal terms are small, and the lowest-order solution for eigenstate $u_{nk} = u_{n0}$ is $c_m^n(\mathbf{k}) = \delta_{nm}$, and the corresponding eigenvalue is given by

$$E_{nk} = E_{n0} + \frac{\hbar^2 k^2}{2m_0} \quad (48)$$

If the nondiagonal terms are small, one can improve the above result by using the second-order perturbation theory

$$E_{nk} = E_{n0} + \left\langle u_{n0} \left| \frac{\hbar^2 k^2}{2m_0} \right| u_{n0} \right\rangle + \sum_{m \neq n} \frac{\langle u_{n0}|H_I|u_{m0}\rangle \langle u_{m0}|H_I|u_{n0}\rangle}{E_{n0} - E_{m0}} \quad (49)$$

where

$$H_I = \frac{\hbar \mathbf{k}}{m_0} \cdot \left( \mathbf{p} + \frac{\hbar}{4m_0 c^2}(\sigma \times \nabla V) \right) \quad (50)$$

Since the kinetic energy operator is a scalar, the second-order eigen energies can be written as

$$E_{nk} = E_{n0} + \frac{\hbar^2 k^2}{2m_0} + \frac{\hbar^2}{m_0^2} \sum_{m \neq n} \frac{|\pi_{nm} \cdot \mathbf{k}|^2}{E_{n0} - E_{m0}} \quad (51)$$

where

$$\left\{ \begin{array}{l} \pi = \mathbf{p} + \frac{\hbar}{4m_0 c^2}(\sigma \times \nabla V) \\[2mm] \pi_{nm} = \langle u_{n0}|\pi|u_{m0}\rangle \end{array} \right. \quad (52)$$

The vector $\mathbf{k}$ can be taken outside the integral in Eq. (51). It is seen that the eigenvalue depends quadratically on the wave vector in the vicinity of the $\Gamma$ point. Then, Eq. (51) is often written as

$$E_{nk} = E_{n0} + \frac{\hbar^2}{2} \sum_{\alpha\beta} k_\alpha \frac{1}{\mu_{\alpha\beta}} k_\beta, \quad \alpha, \beta: x, y, z \quad (53)$$

where

$$\frac{1}{\mu_n^{\alpha\beta}} = \frac{1}{m_0}\delta_{\alpha\beta} + \frac{2}{m_0^2} \sum_{m \neq n} \frac{\pi_{nm}^\alpha \pi_{nm}^\beta}{E_{n0} - E_{m0}} \quad (54)$$

is an effective mass tensor.

### 2.3.3. Kane's Theory

$\mathbf{k} \cdot \mathbf{p}$ theory, as discussed in Section 2.3.2, is essentially based on perturbation theory. A more exact approach, capable of including strong band to band interactions, is provided by Eq. (47). Note that the inclusion of a complete set of basis states in Eq. (47) is not feasible numerically. However, one can improve the $\mathbf{k} \cdot \mathbf{p}$ theory drastically if it includes in Eq. (47) those bands that are strongly coupled, and correct this approximation by treating the influence of distant (energetically) bands perturbatively. This procedure may be made consistent if the electron bands can be divided into two groups. In the first group of bands, there is a strong interband coupling—the number of bands in this group is very limited (up to 8, say). The second group of bands is only weakly interacting with the first set. This interaction is

treated by perturbation theory. This approach is called Kane's [72, 73] model, and has been shown to be very predictive for the III-V compound semiconductors.

Within Kane's theory, one constructs a new basis of $p$-symmetric atomic Bloch states as a linear combination of the "directed orbital" atomic Bloch states $|u_{nn}\rangle$ discussed in the previous section. The new basis set will consist of eigenfunctions of operators $J$ and its component in the $z$-direction $J_z$. The new basis set is then denoted by $|jm_j\rangle$, where $j = 1/2$, $3/2$ and $m_j = j, j - 1, \ldots, -j$, giving six subbands. These can be considered together with the s-symmetric conduction band. The resulting eight-band model gives a good description of the electronic structure of III-V semiconductors near the $\Gamma$ point. The new basis set is given in terms of the directed orbitals shown in Table 3. In the literature, there are other sets of basis functions that differ from the set of functions given here by a unitary transformation $|jm_j\rangle' = U|jm_j\rangle$, $UU^\dagger = I$.

The atomic Bloch states in Table 3 are eigenstates of the Hamiltonian $H(\mathbf{k} = 0)$, and include spin-orbit interaction. $\Gamma_6$ corresponds to the conduction band, $\Gamma_8$ denotes the heavy-hole ($m_j = \pm 3/2$) and $\Gamma_8$ ($m_j = \pm 1/2$) the light-hole band. $\Gamma_7$ is known as split-off band. If we neglect in Eq. (4) the spin-orbit term that depends on the wave vector; that is the term

$$\frac{\hbar \mathbf{k}}{m_0} \cdot (\pi - \mathbf{p}) = \frac{\hbar \mathbf{k}}{m_0} \cdot \left[\frac{\hbar}{4m_0 c^2}(\sigma \times \nabla V)\right] \tag{55}$$

then the matrix representation of the Hamiltonian

$$H(\mathbf{k}) = H(\mathbf{k} = 0) + \frac{\hbar^2 k^2}{2m_0} + \frac{\hbar \mathbf{k} \cdot \mathbf{p}}{m_0} \tag{56}$$

is given in Table 4 using the basis set of Table 3. Note that this Hamiltonian does not yet include the influence of distant bands, which makes the effective mass of the valence band to differ from the electron rest mass. Some of the notations used in Table 4 are given below

$$k_\pm = \frac{1}{\sqrt{2}}(k_x \pm ik_y)$$

$$E_0 = E_{\Gamma_6} - E_{\Gamma_8}$$

$$\Delta = E_{\Gamma_8} - E_{\Gamma_7} \tag{57}$$

$$P = \frac{-i}{m_0}\langle S|p_x|X\rangle = \frac{-i}{m_0}\langle S|p_y|Y\rangle = \frac{-i}{m_0}\langle S|p_z|Z\rangle$$

**Table 3.** The atomic basis states at $\Gamma$ point. The eigenvalues in the fourth column correspond to Eq. (44). The zero point of energy has been set to the bottom of the conduction band.

| $u_i$ | $|j, m_j\rangle$ | $\psi_{j, m_j}$ | $E_i(k = 0)$ | |
|---|---|---|---|---|
| $u_1$ | $\left|\frac{1}{2}, \frac{1}{2}\right\rangle$ | $i|S\uparrow\rangle$ | $0$ | $\Gamma_6$ |
| $u_2$ | $\left|\frac{3}{2}, \frac{1}{2}\right\rangle$ | $-\sqrt{\frac{2}{3}}|Z\uparrow\rangle + \frac{1}{\sqrt{6}}|X + iY\rangle\downarrow$ | $-E_0$ | $\Gamma_8$ |
| $u_3$ | $\left|\frac{3}{2}, \frac{3}{2}\right\rangle$ | $\frac{i}{\sqrt{2}}|X + iY\rangle\uparrow$ | $-E_0$ | $\Gamma_8$ |
| $u_-$ | $\left|\frac{1}{2}, \frac{1}{2}\right\rangle$ | $\frac{1}{\sqrt{3}}|X + iY\rangle\downarrow + \frac{1}{\sqrt{3}}|Z\rangle\uparrow$ | $-E_0 - \Delta$ | $\Gamma_7$ |
| $u_2$ | $\left|\frac{1}{2}, -\frac{1}{2}\right\rangle$ | $i|S\downarrow\rangle$ | $0$ | $\Gamma_6$ |
| $u_4$ | $\left|\frac{3}{2}, -\frac{1}{2}\right\rangle$ | $\sqrt{\frac{2}{3}}|Z\downarrow\rangle - \frac{1}{\sqrt{6}}|X - iY\rangle\uparrow$ | $-E_0$ | $\Gamma_8$ |
| $u_6$ | $\left|\frac{3}{2}, -\frac{3}{2}\right\rangle$ | $\frac{1}{\sqrt{2}}|X - iY\rangle\downarrow$ | $-E_0$ | $\Gamma_8$ |
| $u_5$ | $\left|\frac{1}{2}, -\frac{1}{2}\right\rangle$ | $-\frac{1}{\sqrt{3}}|X - iY\rangle\uparrow + \frac{1}{\sqrt{3}}|Z\rangle\downarrow$ | $-E_0 - \Delta$ | $\Gamma_7$ |

**Table 4.** Eight band Hamiltonian $H(\mathbf{k})$.

$$
\begin{bmatrix}
 & |iS\uparrow\rangle & |\frac{3}{2},\frac{3}{2}\rangle & |\frac{3}{2},\frac{3}{2}\rangle & |\frac{1}{2},\frac{1}{2}\rangle & |iS\downarrow\rangle & |\frac{3}{2},-\frac{1}{2}\rangle & |\frac{3}{2},-\frac{3}{2}\rangle & |\frac{1}{2},-\frac{1}{2}\rangle \\
|iS\uparrow\rangle & \frac{\hbar^2 k^2}{2m_0} & \sqrt{\frac{2}{3}}P\hbar k_+ & P\hbar k_z & \sqrt{\frac{1}{3}}P\hbar k_+ & 0 & \sqrt{\frac{1}{3}}P\hbar k_- & 0 & \sqrt{\frac{2}{3}}P\hbar k_- \\
|\frac{3}{2},\frac{3}{2}\rangle & \sqrt{\frac{2}{3}}P\hbar k_- & -E_0+\frac{\hbar^2 k^2}{2m_0} & 0 & 0 & \sqrt{\frac{1}{3}}P\hbar k_- & 0 & 0 & 0 \\
|\frac{3}{2},\frac{3}{2}\rangle & P\hbar k_- & 0 & -E_0+\frac{\hbar^2 k^2}{2m_0} & 0 & 0 & 0 & 0 & 0 \\
|\frac{1}{2},\frac{1}{2}\rangle & \sqrt{\frac{1}{3}}P\hbar k_- & 0 & 0 & -E_0+\frac{\hbar^2 k^2}{2m_0} & \sqrt{\frac{2}{3}}P\hbar k_- & 0 & 0 & 0 \\
|iS\downarrow\rangle & 0 & \sqrt{\frac{1}{3}}P\hbar k_+ & 0 & \sqrt{\frac{2}{3}}P\hbar k_+ & \frac{\hbar^2 k^2}{2m_0} & \sqrt{\frac{2}{3}}P\hbar k_- & P\hbar k_- & \sqrt{\frac{1}{3}}P\hbar k_- \\
|\frac{3}{2},-\frac{1}{2}\rangle & \sqrt{\frac{1}{3}}P\hbar k_+ & 0 & 0 & 0 & \sqrt{\frac{2}{3}}P\hbar k_+ & -E_0+\frac{\hbar^2 k^2}{2m_0} & 0 & 0 \\
|\frac{3}{2},-\frac{3}{2}\rangle & 0 & 0 & 0 & 0 & P\hbar k_+ & 0 & -E_0-\frac{\hbar^2 k^2}{2m_0} & 0 \\
|\frac{1}{2},-\frac{1}{2}\rangle & \sqrt{\frac{2}{3}}P\hbar k_+ & 0 & 0 & 0 & \sqrt{\frac{1}{3}}P\hbar k_+ & 0 & 0 & -E_0+\Delta-\frac{\hbar^2 k^2}{2m_0}
\end{bmatrix}
$$

We now calculate the dispersion as a function of $\mathbf{k}$ for the eight-band $\mathbf{k} \cdot \mathbf{p}$ theory by diagonalizing the Hamiltonian in Table 4. For bulk systems, or heterostructures in which we have confinement only in one direction, the Hamiltonian in Table 4 is easy to diagonalize if the $z$-axis of the coordinate system is taken in the direction of the wavevector. In this case, $k_z = k$, and accordingly $k_\perp = 0$. This choice is possible since it can be shown that the Hamiltonian is isotropic and, therefore, the eigenvalues and eigenvectors depend on the magnitude of $\mathbf{k}$-only. In this coordinate system, the Hamiltonian is brought into a block form

$$
H = \begin{bmatrix} H_{4 \times 4} & 0 \\ 0 & H_{4 \times 4} \end{bmatrix}
\tag{58}
$$

where the $4 \times 4$ matrix is given by

$$
H_{4 \times 4} = \begin{bmatrix}
\dfrac{\hbar^2 k^2}{2m_0} & -\sqrt{\dfrac{2}{3}}P\hbar k_z & 0 & \sqrt{\dfrac{1}{3}}P\hbar k_z \\[2ex]
-\sqrt{\dfrac{2}{3}}P\hbar k_z & -E_0+\dfrac{\hbar^2 k^2}{2m_0} & 0 & 0 \\[2ex]
0 & 0 & -E_0+\dfrac{\hbar^2 k^2}{2m_0} & 0 \\[2ex]
\sqrt{\dfrac{1}{3}}P\hbar k_z & 0 & 0 & -E_0-\Delta+\dfrac{\hbar^2 k^2}{2m_0}
\end{bmatrix}
\tag{59}
$$

and where for bulk case, one can assume $k_z = k$. The eigenvalues (doubly degenerate) are obtained by finding the roots of the determinant equation

$$
|H - E(\mathbf{k})\mathbf{1}| = 0
\tag{60}
$$

Denoting $\lambda(\mathbf{k}) = E(\mathbf{k}) - \frac{\hbar^2 k^2}{2m_0}$, the eigenvalues, that is, the roots of Eq. (60) are

$$
\lambda(\mathbf{k}) = -E_0
$$

$$
\lambda(\mathbf{k})[\lambda(\mathbf{k}) + E_0][\lambda(\mathbf{k}) + E_0 + \Delta] = \hbar^2 k^2 P^2 \left[\lambda(\mathbf{k}) + E_0 + \frac{2\Delta}{3}\right]
\tag{61}
$$

This last equation corresponds to the original formulation of Kane [72, 73]. In his derivation, he uses a different basis set corresponding to eigenfunctions of the operators $L^2$, $L_z$, $S^2$, $S_z$,

but the eigenvalues and dispersion relations are equal since the basis sets are related by a unitary transformation. From the first equation, one obtains the dispersion for the $HH$-band:

$$E_{hh} = -E_0 - \frac{\hbar^2 k^2}{2m_{hh}}; \quad \frac{1}{m_{hh}} = \frac{1}{m_0}$$

(62)

Note that the effective hole mass is still equal to the bare electron mass. From the second equation, we obtain the other three dispersion relations as follows. We assume that the coefficient $\hbar^2 k^2 P^2$ is small. We then obtain the lowest-order solution by setting this term equal to zero and obtain (as expected) the original band-edge positions:

$$\lambda_{\Gamma_6}^0 = 0, \quad \lambda_{\Gamma_8,lh}^0 = -E_0, \quad \lambda_{\Gamma_7}^0 = -E_0 - \Delta$$

(63)

The zero of the energy scale is taken to be at the conduction band edge. Now, the first order solution is obtained for each band by inserting the zero order solution on the rhs of Eq. (61) and also in the left-hand side in all other terms except for the one becoming zero if the substitution is made. The first-order eigenvalues are then obtained analytically. For example, for the conduction band one obtains

$$\lambda_{\Gamma_6}^1(\mathbf{k})[\lambda_{\Gamma_8}^0(\mathbf{k}) + E_0][\lambda_{\Gamma_6}^0(\mathbf{k}) + E_0 + \Delta] = \hbar^2 k^2 P^2\left[\lambda_{\Gamma_6}^0(\mathbf{k}) + E_0 + \frac{2\Delta}{3}\right] \Leftrightarrow$$

$$\lambda_{\Gamma_6}^1(\mathbf{k})[E_0][E_0 + \Delta] = \hbar^2 k^2 P^2\left[E_0 + \frac{2\Delta}{3}\right] \Rightarrow$$

(64)

$$\lambda_{\Gamma_6}^1(\mathbf{k}) = \hbar^2 k^2 P^2 \frac{[E_0 + \frac{2\Delta}{3}]}{[E_0][E_0 + \Delta]}$$

that is,

$$E_{\Gamma_6}(\mathbf{k}) = \hbar^2 k^2 P^2 \frac{[E_0 + \frac{2\Delta}{3}]}{[E_0][E_0 + \Delta]} + \frac{\hbar^2 k^2}{2m_0} = \frac{\hbar^2 k^2}{2}\left(\frac{1}{m_0} + \frac{4P^2}{3E_0} + \frac{2P^2}{3(E_0 + \Delta)}\right)$$

(65)

which means that the effective mass of the electrons in the vicinity of the conduction band edge is

$$\frac{1}{m_c} = \frac{1}{m_0} + \frac{4P^2}{3E_0} + \frac{2P^2}{3(E_0 + \Delta)}$$

(66)

For the light hole and the split-off bands, one obtains by similar procedure

$$E_{lh} = -E - \frac{\hbar^2 k^2}{2m_{lh}}; \quad \frac{1}{m_{lh}} = \frac{1}{m_0} - \frac{4P^2}{3E_0}$$

$$E_{so} = -E - \Delta - \frac{\hbar^2 k^2}{2m_{lh}}; \quad \frac{1}{m_{so}} = \frac{1}{m_0} - \frac{2P^2}{3(E_0 + \Delta)}$$

(67)

Note that because of the relative magnitudes of the matrix elements of the dipole operator, the conduction band has a positive effective mass, whereas the light-hole and split-off bands have negative effective electron masses. Kane [72, 73] used this method to describe the energy band structure in a $p$-type germanium and silicon, and indium antimonide.

### 2.3.4. Coupling with Distant Bands

To describe the coupling with distant bands, we consider the wave equation

$$(H_0 + W)\psi = E\psi$$

(68)

where

$$W = \frac{\hbar}{m_0}\mathbf{k} \cdot \mathbf{p} + \frac{\hbar^2 k^2}{2m_0}$$

(69)

We assume that the eigenvalues $E_l$ corresponding to eigenstates $|l\rangle l = 1, 2, \ldots, 8$ of the Hamiltonian $H_0$ are close to each others on the energy scale. These are the eight bands considered in Kane's model. These eigenstates are strongly coupled by the operator $W$.

We assume that there is another set of eigenstates $|\nu\rangle$ of $H_0$ only weakly coupled by $W$. We now calculate the correction to the eight lowest eigenvalues of $H_0$ caused by the distant bands. Let $\psi$ be solution to Eq. (68) including this correction, that is,

$$\psi = \sum_l c_l |l\rangle + \sum_\nu c_\nu |\nu\rangle \tag{70}$$

Inserting Eq. (70) into Eq. (68), we obtain

$$\sum_l c_l [(E_l - E)\delta_{lm} + \langle m|W|l\rangle] + \sum_\nu c_\nu \langle m|W|\nu\rangle = 0$$

$$\sum_\nu c_\nu [(E_\nu - E)\delta_{\mu\nu} + \langle \mu|W|\nu\rangle] + \sum_l c_l \langle \mu|W|l\rangle = 0 \tag{71}$$

Since the coupling to the distant bands $|\nu\rangle$ is weak, one can conclude that if $\psi$ is one of the lowest eight eigenvalues, the relative magnitudes of the expansion coefficients are: $|c_l| \cong 1$ and $|c_\nu| \ll 1$. The second of the above equations then gives

$$c_\nu \cong \frac{1}{E - E_\nu} \sum_l c_l \langle \nu|W|l\rangle \tag{72}$$

Inserting this into the first of Eq. (71), we obtain

$$\sum_l c_l \left[ (E_l - E)\delta_{lm} + \langle m|W|l\rangle + \langle m|W \sum_\nu \frac{|\nu\rangle\langle\nu|}{E - E_\nu} W|l\rangle \right] = 0 \tag{73}$$

It is obvious that the influence of the distant bands can be taken into account by the replacement

$$W \rightarrow \tilde{W} = W + W \sum_\nu \frac{|\nu\rangle\langle\nu|}{E - E_\nu} W \tag{74}$$

It can be shown that the Hamiltonian of the distant band interaction, $\tilde{W} - W$, is given by Table 5, where the following notation has been used

$$F(\mathbf{k}) = Ak^2 + \frac{B}{2}(k^2 - 3k_z^2)$$

$$G(\mathbf{k}) = Ak^2 - \frac{B}{2}(k^2 - 3k_z^2)$$

$$H(\mathbf{k}) = -iDk_z(k_x - ik_y) \tag{75}$$

Table 5. Hamiltonian of the distant band interaction.

$$I(\mathbf{k}) = \frac{\sqrt{3}}{2} B(k_x^2 - k_y^2) - iDk_x k_y$$

$$A = \frac{L + 2M}{3}, \quad B = \frac{L - M}{3}, \quad C^2 = D^2 - 3B^2, \quad D = \frac{N}{\sqrt{3}}$$

where

$$L = \frac{\hbar^2}{2m_0} + \frac{\hbar^2}{m_0^2} \sum_r \frac{|\langle X|p_x|\nu\rangle|^2}{E - E_r}$$

$$M = \frac{\hbar^2}{2m_0} + \frac{\hbar^2}{m_0^2} \sum_r \frac{|\langle X|p_y|\nu\rangle|^2}{E - E_r}$$

(76)

$$N = \frac{\hbar^2}{m_0^2} \sum_r \frac{\langle X|p_x|\nu\rangle\langle\nu|p_y|Y\rangle + \langle X|p_y|\nu\rangle\langle\nu|p_x|Y\rangle}{E - E_r}$$

$$\frac{1}{m_e^*} = \frac{2}{m_0^2} \sum_r \frac{|\langle X|p_x|\nu\rangle|^2}{E - E_r} + \frac{1}{m_0}$$

## 2.3.5. The Luttinger-Kohn Hamiltonian

In the Luttinger-Kohn approximation [74] it is assumed that $E_0$ and $\Delta$ are large enough, so that coupling of the $\Gamma_8$ bands with the split-off band is weak. This allows one to derive a $4 \times 4$ Hamiltonian submatrix. The derivation consists of unitary transformation from the basis set $|X\rangle, |Y\rangle, |Z\rangle$ multiplied by the spin functions $|\uparrow\downarrow\rangle$ parts to the $jj$—coupled subspace $|3/2, 3/2\rangle, |3/2, -3/2\rangle$ ($HH$-states) and $|3/2, 1/2\rangle, |3/2, -1/2\rangle$ ($LH$-states). The Hamiltonian is

$$H = \begin{array}{c|cccc}
 & \left|\frac{3}{2},\frac{3}{2}\right\rangle & \left|\frac{3}{2},\frac{1}{2}\right\rangle & \left|\frac{3}{2},-\frac{1}{2}\right\rangle & \left|\frac{3}{2},-\frac{3}{2}\right\rangle \\
\hline
\left|\frac{3}{2},\frac{3}{2}\right\rangle & F - E_0 & H & -I & 0 \\
\left|\frac{3}{2},\frac{1}{2}\right\rangle & H^* & G - E_0 & 0 & I \\
\left|\frac{3}{2},-\frac{1}{2}\right\rangle & -I^* & 0 & G - E_0 & H \\
\left|\frac{3}{2},-\frac{3}{2}\right\rangle & 0 & I^* & H^* & F - E_0 \\
\end{array}$$

(77)

The eigenvalues of the above Hamiltonian are obtained from the determinant equation $|\mathbf{H} - E\mathbf{1}| = 0$, which gives

$$E_{\Gamma_8}(k) = -E_0 + \frac{1}{2}(F + G) \pm \sqrt{\left(\frac{F - G}{2}\right)^2 + |I|^2 + |H|^2}$$

(78)

$$E_{\Gamma_8}(k) = -E_0 + Ak^2 \pm \sqrt{B^2 k^4 + C^2(k_x^2 k_y^2 + k_y^2 k_z^2 + k_z^2 k_x^2)}$$

If one assumes that $\mathbf{k}||z$-axis, one obtains

$$E_{\Gamma_8}(k) = -E_0 + (A \pm B)k^2$$

(79)

Inserting for the matrix elements A and B, the $HH$-$LH$ dispersion is given by

$$E = -E_0 - \frac{\hbar^2 k^2}{2m^*} : \quad \frac{1}{m^*} = -\frac{1}{m_0} + \frac{2}{m_0^2} \sum_{m\neq\Gamma_8} \frac{\langle\frac{3}{2},\pm\frac{3}{2}|p_z|u_{m0}\rangle}{E_m + E_0}$$

(80)

The last term is a correction coming from the coupling with the distant bands. This gives the $HH$-band a negative electron effective mass (and a positive hole effective mass).

The Hamiltonian given in Eq. (77) can be written in a more familiar form in terms of Luttinger parameters. First, the Luttinger parameters $\gamma_1$, $\gamma_2$, $\gamma_3$ are written in terms of $L$, $M$, $N$, $R$, $S$, and $T$ (see the discussion in Ref. [75]). The matrix elements of the $HH$-$LH$ Hamiltonian are then given by

$$F(\mathbf{k}) = -\frac{\hbar^2}{2m_0}\left[(\gamma_1 + \gamma_2)(k_x^2 + k_y^2) + (\gamma_1 - 2\gamma_2)k_z^2\right]$$

$$G(\mathbf{k}) = -\frac{\hbar^2}{2m_0}\left[(\gamma_1 - \gamma_2)(k_x^2 + k_y^2) + (\gamma_1 + 2\gamma_2)k_z^2\right]$$

$$H(\mathbf{k}) = \frac{\hbar^2}{2m_0}\sqrt{3}\gamma_3 k_z (k_x - ik_y)$$                              (81)

$$I(\mathbf{k}) = -\frac{\hbar^2}{2m_0}\sqrt{3}\left[\gamma_2(k_x^2 - k_y^2) - 2i\gamma_3 k_x k_y\right]$$

## 2.4. Applications of k · p to Quasi-2D Electron and Hole Systems

### 2.4.1. Heterostructure Devices

Development of molecular beam epitaxy (MBE) [76, 77] has been pushed by device technology to achieve structures with nanoscale and atomic scale dimensions, and this has led to an entirely new area of condensed matter physics and investigation of structures exhibiting strong quantum size effects. MBE has played a key role in the discovery of phenomena like the two dimensional electron and hole gases, quantum Hall effect [78], and new structures like quantum wires [79] and quantum dots [329]. The continued miniaturization of solid-state devices is leading to the point where quantization-induced phenomena become more and more important. These phenomena have shown that the role of material purity, native defects, and interface quality are very critical to the device performance.

Modulation doping is a technique employed to achieve adequate carrier densities in one region of the device which is physically separated from the source of the carriers, the ionized impurities. The low-temperature mobility of modulation doped GaAs/AlGaAs structures is a good measure of the GaAs/AlGaAs material quality. This depends very strongly on the epitaxial structure, particularly the placement and quantity of dopant impurities. The two-dimensional electron gas (2DEG) that exists at the interface between GaAs and the wider band gap AlGaAs exhibits a very high mobility at low temperatures. Even at room temperatures, the mobility is larger than that of bulk GaAs. Two factors contribute to this higher mobility, both arising from the selective doping of AlGaAs buffer layers rather than the GaAs layers in which the carriers reside. The first is the natural separation between the donor atoms in the AlGaAs and the electrons in the GaAs. The second is the inclusion of an undoped AlGaAs spacer layer in the structure. Such structures are quite complicated but can be easily fabricated using MBE techniques. A typical heterostructure begins with the bulk GaAs wafer upon which a GaAs buffer layer or superlattice is grown. The latter is used to act as a barrier to the out-diffusion of impurities and defects from the substrate. It also consists of a GaAs cap layer and alternating layers of AlGaAs and GaAs. The common practice is to use a doping for the AlGaAs layers in the active region, but nowadays undoped AlGaAs layers are used and a delta doped layer is included. This delta doped layer along with the growth of superlattices restricts the formation of defects, known as D-X centers [81], to a minimum. There are two important AlGaAs layers on either side of the δ-doped layer and they are called buffer and spacer layer, respectively. The spacer layer is closer to the GaAs quantum well and is of high purity to prevent scattering of the channel carriers by the ionized impurities. A usual practice is to use undoped AlGaAs layers to have very good confinement of the charge carriers in the well.

A prototypical GaAs/AlGaAs heterostructure used, for example, for quantum wires and dots formation that utilizes the split-gate technique, is shown in Fig. 10. The calculated

Figure 10. Modulation (delta) doped GaAs/AlGaAs quantum well structure. The left is a schematic of the layer structure, the right side the calculated conduction band profile. Reprinted with permission from [83], A. Ashwin, Master's Thesis, Arizona State University, 2005. © 2005.

zero-temperature conduction band profile along the growth direction is shown on the right panel of Fig. 10. One can see relative to the position of the Fermi energy, that the electrons are localized in the quantum-well region, forming a 2DEG. The top surface has a patterned Schottky gate structure, called a quantum point contact (QPC), which depletes the electrons underneath with negative bias, forming a quasi–one-dimensional waveguide through which electrons are injected [see, for example, Ref. [9]]. The particular structure shown in Fig. 10 has been used in the investigation of spin filtering of electrons through the quantum point contact structure [82, 83].

### 2.4.2. Spin-Resolved Band Structure in GaAs Quantum Wells

Today, research in *spintronics*—the field of semiconductor electronics based upon exploiting the quantum mechanical property of electron spin as an information carrying entity—owes its existence, in general, to two primary factors. The first is theoretical: the spin component of the electron's wave function can retain its form (i.e., it's orientation or, if referring to an ensemble of electrons, the coherence) in semiconductor transport for much longer times than can the spatial components. In fact, the spin relaxation time, $\tau_s$, can be on the order of nanoseconds, as opposed to picoseconds for the spatial momentum relaxation time [84, 85]. The second factor is technological: improvements in lithography and the continuing progress in developing efficient spin filter injection/detection mechanisms that polarize (or detect polarized) electrons by various means other than by applying a cumbersome external magnetic field [86], have prompted the development of novel spintronic devices such as the spin transistor (SPINFET) [53]. The starting point for understanding spin injection/detection, spin-flip scattering mechanisms and their relative impact on transport is, of course, the determination of an accurate spin-resolved band structure.

This section is organized as follows. We first present an application of the multiband k · p to calculate the full band-structure in quantum confined systems in the absence of stress, strain, or magnetic fields. The effect of bulk inversion asymmetry, a source of spin-splitting in bands, is then introduced into the model. The k · p model discussed in Section 2.3.1 is further extended to account for spin splitting in the conduction band due to any structural inversion asymmetry (SIA) that may exist.

### 2.4.2.1. General Characteristics of the k · p Solver.
The Fermi energy of the electrons involved in the transport through the QPCs is on the order of 10 meV. This energy is low enough such that a band structure calculation accurate only around a point of high symmetry in the Brillouin zone (Γ point) is sufficient. For this situation, the k · p method of calculating band structure, explained in details in Section 2.3, is applicable.

Eppenga et al. [87], using the Kane basis set and including spin-orbit coupling and remote band effects via Löwdin perturbation theory [71], arrive at the following Hamiltonian

$$
H = \begin{bmatrix}
E_{EL} & -\sqrt{2}P_z & P_z & \sqrt{3}P_+ & 0 & -P_- & -\sqrt{2}P_- & 0 \\
& E_{LH} & G_1 & \sqrt{2}G_z & -P_+^* & 0 & -\sqrt{3}G_- & G_2 \\
& & E_{SO} & -G_- & -\sqrt{2}P_+^* & \sqrt{3}G_- & 0 & \sqrt{2}G_2 \\
& & & E_{HH} & 0 & -G_2 & -\sqrt{2}G_2 & 0 \\
& & & & E_{EL} & \sqrt{2}P_z & -P_z & -\sqrt{3}P_- \\
& & & & & E_{LH} & G_1 & \sqrt{2}G_- \\
& & & & & & E_{SO} & -G_- \\
& & & & & & & E_{HH}
\end{bmatrix}
\tag{82}
$$

The "empty spaces" in Eq. (82) denote Hermitian conjugates. The diagonal matrix elements are

$$
\begin{aligned}
E_{EL} &= E_g + s\tilde{e} \\
E_{LH} &= -\gamma_1\tilde{e} + \gamma_2\tilde{e}_1 \\
E_{SO} &= -\Delta - \gamma_1\tilde{e} \\
E_{HH} &= -\gamma_1\tilde{e} + \gamma_2\tilde{e}_1
\end{aligned}
\tag{83}
$$

where $E_g$ is the band gap and $\Delta$ is the spin-orbit split-off energy. The $\gamma_1$, $\gamma_2$, and $s$ parameters are effective mass parameters that modify the free-electron term. Also,

$$
\tilde{e} = \frac{\hbar^2}{2m}(k_x^2 + k_y^2 + k_z^2), \quad \tilde{e}_1 = \frac{\hbar^2}{2m}(2k_z^2 - k_x^2 - k_y^2), \quad e_2 = \frac{\hbar^2}{2m}(k_y^2 - k_x^2)
\tag{84}
$$

The off-diagonal terms are given as

$$
\begin{aligned}
P_z &= \sqrt{\frac{1}{3}}(ip\tilde{k}_z + \beta k_x k_y) \\
P_+ &= \sqrt{\frac{1}{6}}[ip(k_x \pm ik_y) + \beta\tilde{k}_z(k_x \pm ik_y)] \\
G_1 &= \sqrt{2}\gamma_2\tilde{e}_1 \\
G_2 &= -\sqrt{3}\gamma_2 e_2 + i2\sqrt{3}\gamma_3 k_x k_y \\
G_- &= \sqrt{6}\gamma_3\tilde{k}_z(k_x \pm ik_y)
\end{aligned}
\tag{85}
$$

In these terms, $p = -\frac{\hbar^2}{m}\int_{unit\,cell} dr\phi_s^* \frac{d}{dz}\phi_z$, and accounts for the coupling of the conduction band $s$-states of the Kane basis with the valence band $z$-state. Terms containing $\gamma_3$ result in an anisotropic band structure near the $\Gamma$ point if $\gamma_2 \neq \gamma_3$.

The parameter $\beta$ is due to the bulk inversion asymmetry (BIA) and causes spin-splitting of the bands. Also known as Dresselhaus splitting [55], BIA induced spin-splitting occurs in zinc-blende semiconductors because two different kinds of atoms (e.g., Ga and As, Ga and Sb, etc...) exist, resulting in asymmetrical wave functions about an axis of symmetry (e.g., [100] axis). This means that while Kramer's theorem, $E(k \uparrow) = E(-k \downarrow)$ is satisfied for all values of $k$, the situation away from $k = 0$ is such that $E(k \uparrow) \neq E(-k \downarrow)$. In fact, spin-splitting of the conduction band ($\Gamma_6$ band) due to BIA is proportional to $k^3$ for small values of $k$ in bulk zinc-blende semiconductors. However, in 2DEG systems such as heterostructures, a linear dependence on $k$ occurs too. The valence bands exhibit linear BIA splitting in both the bulk and the 2DEG cases. Thorough discussions of the $k$-dependence of BIA spin-splitting can be found in Zawadzki and Pfeffer [88], Silsbee [89], and Winkler [90].

Finally, note that the matrix operator of Eq. (82) does not include the effects of stress and strain or the influence of a magnetic field. Of course, these effects undoubtedly will modify the $k$-dependence on spin-splitting. However, it must be mentioned that the Kane model inclusive of these effects has been derived by several researchers, notably Trebin et al. [91].

### 2.4.2.2. Quantum Well Structures.

To test the applicability of the Kane approach, the Hamiltonian given in Eq. (82) is applied to the *symmetrical* quantum well shown in Fig. 11, Ref. [92].

To start, one applies the operator $\hat{p} \to -i\hbar\nabla$ to all $k_z$ terms in the eigenvalue equation with matrix operator given in Eq. (82). This operation is done to the $k_z$ terms only, since quantization is assumed to be along the $z$-direction. In doing this, it is convenient to note that the resultant matrix, with the matrix operator given explicitly included, then takes the form,

$$\begin{pmatrix} F_{11}\nabla_z^2 + B_{11}\nabla_z + C_{11} + D_{11} & \cdots & F_{18}\nabla_z^2 + B_{18}\nabla_z + C_{18} + D_{18} \\ \vdots & \ddots & \vdots \\ F_{81}\nabla_z^2 + B_{81}\nabla_z + C_{81} + D_{81} & \cdots & F_{88}\nabla_z^2 + B_{88}\nabla_z + C_{88} + D_{88} \end{pmatrix} \begin{pmatrix} \chi_{1,z} \\ \chi_{2,z} \\ \vdots \\ \chi_{7,z} \\ \chi_{8,z} \end{pmatrix} = \varepsilon \begin{pmatrix} \chi_{1,z} \\ \chi_{2,z} \\ \vdots \\ \chi_{7,z} \\ \chi_{8,z} \end{pmatrix} \quad (86)$$

where the $F$ and $B$ terms denote coefficients to second- and first-order partial derivatives w.r.t. $z$, respectively. The $C$ terms denote potential offsets, and the $D$ terms indicate all other terms not operated on by the momentum operator. All terms $B$, $C$, $D$, and $F$ are actually $z$-dependent functions, though not denoted as such, to ease the notation.

To ensure Hermiticity of the resulting matrix at the heterostructure interfaces, the discretization scheme as noted in Eppenga et al. [17] is used, where

$$B(z)\frac{\partial}{\partial z} \to \frac{1}{2}\left(B(z)\frac{\partial}{\partial z} + \frac{\partial}{\partial z}B(z)\right) \quad (87)$$

$$F(z)\frac{\partial^2}{\partial z^2} \to \frac{\partial}{\partial z}F(z)\frac{\partial}{\partial z} \quad (88)$$

Equations (87) and (88) yield

$$B(z)\frac{\partial\chi}{\partial z} = \frac{1}{4\Delta z}[2B(z)\chi_{z+1} - 2B(z)\chi_{z-1} + (B(z+1) - B(z-1))\chi_z]$$

$$\approx B(z)\frac{\chi_{z+1} - \chi_{z-1}}{2\Delta z} \quad (89)$$



Figure 11. Symmetrical quantum well. Though not shown here, the valence split off band also exhibits a band offset. The dimensions of the well are chosen so as to be able to compare to the tight-binding model results of Chang and Schulman [93].

and

$$\frac{\partial \chi}{\partial z} F(z) \frac{\partial \chi}{\partial z} = \left( \frac{F(z)}{\Delta z^2} + \frac{F(z+1) - F(z-1)}{4\Delta z^2} \right) \chi_{z+1}$$

$$+ \left( \frac{F(z)}{\Delta z^2} + \frac{F(z+1) - F(z-1)}{4\Delta z^2} \right) \chi_{z-1} - \left( \frac{2F(z)}{\Delta z^2} \right) \chi_z$$

$$\approx \frac{F(z)\chi_{z+1} + F(z)\chi_{z-1} - 2F(z)\chi_z}{\Delta z^2} \tag{90}$$

Equations (89) and (90) are then used to transform the matrix from the bulk to the $z$-quantized situation. Thus,

$$F(z)\frac{\partial^2 \chi}{\partial z^2} \rightarrow \frac{F(z)\chi_{z+1} + F(z)\chi_{z-1} - 2F(z)\chi_z}{\Delta z^2} \tag{91}$$

$$B(z)\frac{\partial \chi}{\partial z} \rightarrow B(z)\frac{\chi_{z+1} - \chi_{z-1}}{2\Delta z} \tag{92}$$

Applying Eqs. (91) and (92) to the eigenvalue equation results in a matrix of the form shown in Table 6.

The results, depicted in Fig. 12 for the valence bands, are in good agreement with the Chang and Schulman tight binding calculation [93], which requires several more adjustable fitting parameters other than $\gamma_1$, $\gamma_2$, $\gamma_3$, $s$, and $\rho$ used here.

Although Fig. 12 indicates that the 8-band **k · p** model is useful for accurately describing band structure in the vicinity of a point of symmetry in the Brillouin zone, the model does of course have severe limitations. Indeed, it was found that in calculating the bulk band structure of GaAs, the conduction band suddenly curves down and goes negative at about 10% of the path length along 111 ($\Gamma \rightarrow$ L). Thus, at large real values of $k$, the bandgap can "disappear" and "spurious" solutions to the eigenvalue problem exist [94]. Therefore, at large values of $k$, the 8-band **k · p** approach is not appropriate. Other methods such as tight-binding methods must be used. Higher order **k · p** models, such as the 24-band **k · p** model recently reported by Radhia et al. [95], could also possibly be used, though they would be computationally inefficient and computer-memory intensive.

### 2.4.2.3. Inclusion of SIA Effects into the Eight-Band Model.

As previously mentioned, SIA is of great interest in spintronics, since the promise of many proposed spintronic devices relies on the principle of being able to modulate the SIA via application of an external

**Table 6.** General form of the expanded eight-band Kane Matrix. Each point in real-space, along the quantized $z$-axis corresponds to an eight-row block in this matrix.

GaAs/AlGaAs QW dispersion, 19.2 nm QW width, 20nm AlGaAs sides.

Figure 12. Valence band dispersion in 19.2-nm quantum well, as calculated with eight-band k·p method.

electric field. The "Rashba effect" [54] predicts that an electric field applied perpendicular to the plane of a 2DEG will cause SIA. In turn, this electric field will then relativistically induce an effective magnetic field in the plane of the 2DEG, effectively lifting the spin-degeneracy of the charge carriers.

Aside from external electric fields, SIA can in principle be caused by anything that results in an asymmetrical quantum well. Effectively, this means that the penetration of the wave function in the cladding layers is not identical to both sides of the well. As reported by Silsbee [89], this can occur during MBE growth. For example, if one attempts to grow a heterostructure formed by a cation and anion of one kind of atom (denoted by C1 and A1, respectively) and a cation and anion of another type of atom (C2 and A2), one would desire the following growth pattern:

$$\ldots\text{A1-C1-A1-C1-A1-C1-A1} = \text{C2-A2-C2-A2-C2-A2-C2} = \text{A1-C1-A1-C1-A1-C1}\ldots,$$

where the "=" sign indicates the heterostructure interface. Typically, though, MBE growth yields

$$\ldots\text{A1-C1-A1-C1-A1-C1-A1} = \text{C2-A2-C2-A2-C2-A2-C2-A2} = \text{A1-C1-A1-C1-A1-C1}\ldots.$$

Note that the second structure is asymmetrical.

The general form of the SIA term to be added to Eq. (82) to account for SIA splitting of the conduction band is given in Eq. (5), and repeated here for convenience

$$H_{so} = \alpha_1 \sigma \cdot (k \times \hat{z}) \tag{93}$$

The term $\alpha_1$ is a coefficient, typically called the "Rashba coefficient." In fact, the precise nature of this term has been controversial. Until recently many people believed that this term is proportional to the electric field of the conduction band. However, Zawadzki and Pfeffer [88] point out that the average electric field in a bound state of a quantum well is zero. Interestingly, they report that $\alpha_1$ has a dependence on the valence band offsets at the interfaces. Furthermore, Winkler [96] has stated that the SIA splitting of the $HH$ and $LH$ valence bands ($\Gamma_8$ bands) should take a form,

$$H_{8v} = \alpha_2 J \cdot (k \times \hat{z}) + \alpha_3 J' \cdot (k \times \hat{z}) \tag{94}$$

Since the basis set considered is the same as the Kane basis, $J$ is defined as the angular momentum matrix operator for particles of momentum $j = 3/2$, and $J'$ is defined as $J^3$.

Most models of SIA induced spin-splitting include only a $2 \times 2$ term. However, both Eqs. (93) and (94) may be added to the Hamiltonian matrix operator given in Eq. (82).

Using an eight-band solver, Cartioxa et al. [97] have recently investigated the relative contributions of the BIA and SIA terms to conduction band splitting in an asymmetrical AlSb/InAs/GaSb/AlSb quantum well, as depicted in Fig. 13. Note that the band offset of the InAs and GaSb form a so-called broken-gap Type II heterostructure, where the VB of GaSb overlaps the CB of InAs. In this structure, each layer is $\sim$3.0 nm thick. In Fig. 13(b), they show the BIA and SIA splitting for a particular magnitude of $k$ ($k = 0.01(2\pi/a)$, $a$ is the lattice spacing) as this vector is rotated through the $k_x k_y$ plane.

### 2.4.3. k · p for NanoScale p-Channel MOSFET Devices

Electron transport in Si inversion layers has been a primary subject of research for many years now; however, hole transport has been relegated to the background mainly due to the complex valence band-structure in Si. Hole transport is affected by the warping and anisotropy of the valence bands, and the band structure cannot be approximated with a simple effective mass picture. The advent of alternate device structures [98–100] aimed at boosting the speed and density of VLSI circuits, however, seems to have revived interest in hole transport. The important alternate device technologies are buried channel strained SiGe $p$-channel MOSFETs and surface channel strained Si.

As an example of application of the **k · p** method in simulation of hole transport in a nanoscale Si p-channel MOSFETs, we develop here a method for incorporating band structure and quantum effects. This is achieved by coupling a 2D Poisson–1D discretized 6 × 6 **k · p** Hamiltonian solver (discussed in Section 2.4.2, for the case when the conduction band contribution is neglected) self-consistently to a Monte Carlo transport kernel. Monte Carlo methods for semiclassical transport are described in more detain in Section 3.1. This method is generic and can easily be extended to model strained layer MOSFETs by incorporating an additional strain Hamiltonian into the band-structure kernel.

The band structure is calculated using the **k · p** method, where the Hamiltonian is written

$$H = \begin{bmatrix} \mathbf{H}_{k \cdot p} & 0 \\ 0 & \mathbf{H}_{k \cdot p} \end{bmatrix} + \mathbf{H}_{so} + \mathbf{1}V(z) \tag{95}$$

where $\mathbf{H}_{k \cdot p}$ and $\mathbf{H}_{so}$ are the 6 × 6 **k · p** and the spin-orbit Hamiltonians respectively, **1** is a 6 × 6 identity matrix, and $V(z)$ is the confining potential along the device depth. Replacing the vector $k_z$ with its operator notation as $k_z = -i\partial/\partial z$, and using a finite difference discretization, Eq. (96) can be recast into an eigenvalue equation for the eigenenergies in the $xy$-plane, for different values of the in-plane **K**-vector, $\mathbf{K}_{||}(k_x, k_y)$. The solution of the



Figure 13. (a) Asymmetrical QW of Cartioxa et al. (b) Spin splitting for well of (a) [91].

eigenvalue problem involves the diagonalization of a tridiagonal block matrix whose rank is given by 6 × $N_z$, where $N_z$ is the number of mesh points along the depth direction. For the 3D (bulk) carriers in the source and drain, we only have the first two terms of Eq. (95). This 6 × 6 Hamiltonian can easily be diagonalized to give the eigenvalues of 3D carriers at $(k_x, k_y, k_z)$.

To include carrier scattering within the transport kernel, the density of states of the system (2D and 3D) are required. For the 2D case, we tabulate the in-plane K- vector, $K_{||}$, as a function of carrier energy $(\varepsilon_{2d})$, band $(\nu)$ and subband $(n)$ indices, and the in-plane azimuth angle $(\phi)$. For the 3D case, the K-vector, $K_{3D}$ is tabulated as a function of carrier energy $(\varepsilon_{3d})$, band index $(\nu)$, and the azimuthal $(\phi)$ and polar $(\theta)$ angles. In order to set up the inverse problem, the discretized eigenvalue Eq. (95) for the 2D system can be recast into an eigenvalue equation for $|K_{||}|$ [101] as shown by the following equation, where $D_n$ operates on $|K_{||}|^n$.

$$\begin{bmatrix} \mathbf{0} & \mathbf{I} \\ -\mathbf{D}_2^{-1}\cdot[\mathbf{D}_0 - I E] & -\mathbf{D}_2^{-1}\cdot\mathbf{D}_1 \end{bmatrix}\begin{bmatrix} \psi_K \\ \psi_K^{(1)} \end{bmatrix} = K\begin{bmatrix} \psi_K \\ \psi_K^{(1)} \end{bmatrix} \tag{96}$$

Since $\varepsilon_\nu^n(k_x, k_y)$ is quadratic in $|K_{||}|$, the problem involves diagonalizing a matrix whose rank is twice as large as that of the discretized $k\cdot p$ Hamiltonian (i.e., 12 × $N_z$). In the 3D case, using a similar technique, one can show that the problem involves diagonalizing a matrix whose rank is twice that of the $k\cdot p$ Hamiltonian (i.e., 12). Thus, for the 3D case, one can tabulate the values initially and these can be used throughout the simulation. The computational complexity for the 2D case led us to make the following simplifying assumptions.

Using a sufficiently high vertical electric field ~5MV/cm, a triangular test potential was generated and used to tabulate the dispersions and density of states (DOS) of the ground state subbands in each band (heavy-hole HH, light-hole LH, and split-off SO). It was then assumed that for the case of a real confining potential in the device, the dispersion in each subband for a particular band would be given by the tabulated (triangular-well) dispersion of the ground state subband of the corresponding band, thus allowing the capture of basic features of the subband anisotropy, warping and nonparabolicity. The only effect of the "real" confining potential in the device would be the translation of the dispersion on the energy axis by the subband energies at the $\Gamma$ point.

$$\varepsilon_\nu^n(k_x, k_y) \approx [\varepsilon_\nu^0(k_x, k_y) - \varepsilon_\nu^0(0,0)] + \varepsilon_\nu^n(0,0) \tag{97}$$

For the inverse problem, a similar approach is used. The triangular test potential is used in the inverse solver, in order to tabulate the in-plane K-vectors $K_{||}^{n,\nu}(\varepsilon_{2d}, \phi)$ for a set of chosen $(\varepsilon_{2D}, \phi)$. Having tabulated the in-plane K-vectors for the lowest subband in each band, we assume that the same dispersion holds also when employing the actual device potential for all the subbands of the given band, that is,

$$K_{||}^{n,\nu}(\varepsilon_{2D}, \phi) \approx K_{||}^{0,\nu}(\varepsilon_{2D}, \phi) \tag{98}$$

The Monte Carlo particle-based simulator handles the transport of holes through the device and is described in much more details in Section 3.1.3. Having calculated the hole band structure in the contacts and the active device region, that is, under the gate, the quantum mechanical hole density in the channel, is calculated self consistently with the Poisson equation and the 2D band-structure code. Holes are then initialized in real space based on the local carrier density and their energy is initialized by assuming a thermal distribution. As the carriers drift under the influence of the electric field due to the applied bias, the confining potential changes and this in turn changes the eigen energies and the eigenfunctions. As a result, the scattering rates must be updated frequently during the simulation. Within the scope of the current model, it is assumed that holes are quasi-3D particles in the source and drain regions, and have used appropriate models to treat these boundary conditions effectively.

The isosurfaces of the lowest heavy, light, and split-off subbands, for the case of the triangular test potential are shown in Fig. 14. Note the strong warping of the heavy hole band when compared with the fairly regular shapes for the light hole and the split-off bands, which makes it extremely difficult for analytical band models to describe the valence band-structure accurately. The hole density of states is determined by performing a surface integral over these isosurfaces, and these are then used to determine the carrier scattering rates in the channel.

The density of states for the confined carriers is shown in Fig. 15. The deviation of the density of states obtained by a full band calculation from a regular steplike profile expected for an ideal two-dimensional systems described in the effective-mass approximation is clearly seen in the case of the light-hole and split-off bands, whereas the heavy hole density of states looks more like a step function.

The output characteristics of a 25-nm $p$-channel conventional Si MOSFET are shown in Fig. 16. Significant drain-induced barrier lowering (DIBL) is seen in the output characteristics in this case, as evidenced by the increasing drain current with drain-source voltage after saturation is reached. Only phonon scattering (acoustic and optical phonons) are included in this calculation, as discussed later. An equivalent effective mass two band (heavy and light hole bands) model with similar scattering mechanisms included, underestimates the current by about 14%. Thus, it is clear that the effective mass approximations is not reliable and, therefore, band-structure calculations are required to accurately predict the output current under high field transport conditions in nanoscale MOSFETs, in particular $p$-channel ones.

## 2.5. Solving the Effective Mass Schrödinger Equation in State-of-the Art Devices

It has been known for many years that carriers in the inversion layer of a Si MOSFET are confined by the barrier between the semiconductor-oxide interface on one side and the band bending of the conduction band on the other side. Since the average thickness of the inversion layer is comparable to the de Broglie wavelength of the electrons, this confinement is sufficient to produce quantization in the direction normal to the oxide-semiconductor interface. The space quantization effect is very important in determining the number of carriers in the inversion layer for devices with very high substrate doping, representative of current state-of-the art technology. In principle, the solution of even the equilibrium problem in these structures requires self-consistent solution of the Poisson and the Schrödinger equations using any of the methods described in Sections 2.1–2.3. This is a difficult and time-consuming problem, however. In many practical situations, it is sufficient to utilize the band-structure solvers to get the proper band-edge effective masses, which are then used in the time-independent Schrödinger equation for stationary potentials. One such tool that has been successfully utilized in the calculation of the energy-level structure in simple MOS or dual-gate capacitor structures is SCHRED that has been developed at Arizona State University and is currently residing on the Purdue NanoHUB [102].



Figure 14. Isosurfaces of the lowest lying *HH*, *LH*, and *SO* subbands on a (001) oriented substrate. Reprinted with permission from S. Krishnan et al., *Microelectronics* (in press). © 2005, Elsevier.

**Figure 15.** Density of states for channel (triangular test potential) and bulk (3D) carriers, respectively, from **k · p** calculations. Reprinted with permission from S. Krishnan et al., *Microelectronics* (2005). © 2005, Elsevier.

## 2.5.1. Description of SCHRED

SCHRED is a one-dimensional solver that solves self-consistently the 1D Poisson equation

$$\frac{\partial}{\partial z}\left[\varepsilon(z)\frac{\partial \varphi}{\partial z}\right] = -e\left[N_D^+(z) - N_A^-(z) + p(z) - n(z)\right] \tag{99}$$

and the 1D one-electron effective-mass Schrödinger equation

$$\left[-\frac{\hbar^2}{2m_i^\perp}\frac{\partial^2}{\partial z^2} + V_{\text{eff}}(z)\right]\psi_{ij}(z) = E_{ij}\psi_{ij}(z) \tag{100}$$

In Eqs. (99) and (100), $\varphi(z)$ is the electrostatic potential, $\varepsilon(z)$ is the spatially dependent dielectric constant, $N_D^+(z)$ and $N_A^-(z)$ are the ionized donor and acceptor concentrations, $n(z)$ and $p(z)$ are the electron and hole densities, $V_{\text{eff}}(z)$ is the effective potential energy term, $m_i^\perp$ is the effective mass normal to the semiconductor-oxide interface of the $i$th valley, and $E_{ij}$ and $\psi_{ij}(z)$ are the energy level and the corresponding wave function of the electrons residing in the $j$th subband from the $i$th valley. The effective potential energy term, $V_{\text{eff}}(z)$, in the 1D Schrödinger equation equals the sum of the Hartree $[V_H = -e\varphi(z)]$, image $[V_{\text{im}}(z)]$,



**Figure 16.** Output characteristics of a 25-nm *p*-channel Si MOSFET calculated using a full band and an effective mass model. Reprinted with permission from S. Krishnan et al., *Microelectronics* (2005). © 2005, Elsevier.

and exchange-correlation $[V_{exc}(z)]$ terms. The Hartree term represents the solution of the 1D Poisson equation, including the average charge density of the free electrons.

If Poisson's equation is solved fully as a boundary value problem, then appropriate electrostatic boundary conditions are imposed at the dielectric interface between regions of different permittivity. If, rather, the electrostatic potential in 1D is derived from integration of the charge density in the Hartree approximation, then one can account for the dielectric boundary conditions by introducing a fictitious "image" charge, which gives the proper boundary conditions at the dielectric interface between, for example, Si and $SiO_2$. The image potential felt by and electron at a position $z$ from the interface, is given by

$$V_{im}(z) = \frac{e^2}{16\pi\varepsilon_{sc}z} \frac{\varepsilon_{sc} - \varepsilon_{ox}}{\varepsilon_{sc} + \varepsilon_{ox}}$$                (101)

where $\varepsilon_{sc}$ and $\varepsilon_{ox}$ are the semiconductor and the oxide dielectric constants, respectively.

*2.5.1.1. Inclusion of Exchange-Correlation Corrections.* In silicon inversion layers, due to the large effective mass, many-body effects such as exchange and correlation can play an important role. For example, Stern [103] has calculated that the exchange energy is comparable to or even larger than the energy separation between subbands calculated in the Hartree approximation. In general, the exchange energy is the contribution to the overall energy of the electron gas that arises from the correlation between two electrons whose positions are reversed, or exchanged [104]. In other words, as a consequence of the Pauli exclusion principle, electrons with equal spin tend to avoid each other (exchange repulsion) so that each electron is surrounded by an *exchange hole*. The presence of the exchange hole indicates that the mean separation between electrons with equal spin is larger than it would be without the Pauli principle. The existence of the exchange hole reduces the overall Coulomb repulsion, which explains the reduction in the ground-state energy of the system.

According to Hartree-Fock theory, electrons with different spin do not avoid each other, since the states are chosen to satisfy the exchange principle, but they do not include Coulomb correlations [218]. In reality, there exists an additional correlation, which leads to the so-called Coulomb hole. To treat these effects, one has to go beyond the Hartree-Fock theory. Therefore, if one writes the exact ground state energy of the system as

$$E = E^{HF} + E_{corr} = E_{kin}^{HF} + E_{exc}^{HF} + E_{corr}$$                (102)

it is obvious that the correlation energy represents the correction to the ground-state energy of the system beyond the Hartree-Fock approximation. Therefore, the correlation energy is not a quantity with physical significance; it merely represents the error incurred in making a fairly crude first-order approximation. Since an exact calculation of $E_{corr}$ is generally not possible, one of the main tasks of the many-body theory is to obtain good estimates for $E_{corr}$.

In a series of three papers, Hohenberg and Kohn [106], Kohn and Sham [107], and Sham and Kohn [108] laid the foundations for a "new" theory of electronic structure. The theory represents a systematic extension of the Thomas-Fermi ideas, and is capable in principle of providing exact answers. The theory is based on two theorems which center on the particle density as a fundamental variable for the description of any many-body system. The first theorem states that the total ground-state energy $E$ of any many-body system is a functional of the one-particle density, $n(r)$. In this context, different many-body systems differ only by the local potential felt by the electrons. Furthermore, splitting off from the total energy the explicit interaction with the external potential, $V_{ext}(r)$, the theorem also states that the rest is a universal functional of $n(r)$, that is, independent of the external potential. Thus, if

$$E[n] = F[n] + \int d^3r V_{ext}(r)n(r)$$                (103)

then the functional, $F$, depends only on $n$ and not on $V_{ext}(r)$. The second theorem states that for any system (any external potential), the functional, $E[n]$, for the total energy has a minimum equal to the ground-state energy at the physical ground-state density of the

system. These theorems, although rather abstract in nature, were of immense importance to the rapid development of density-functional theory. It is customary to extract from $F[n]$ the classical Coulomb energy and write

$$G[n] = F[n] - \frac{1}{8\pi\varepsilon_0} \int d^3r \int d^3r' \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \tag{104}$$

In this notation, the energy functional becomes

$$E[n] = G[n] + \int d^3r V_{ext}(\mathbf{r})n(\mathbf{r}) + \frac{1}{8\pi\varepsilon_0} \int d^3r \int d^3r' \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \tag{105}$$

The stationary functional, $E[n]$, allows in principle a much simpler determination of the ground-state energy, $E$, and density, $n(\mathbf{r})$, than the conventional Rayleigh-Ritz method [109]. The functional, $G[n]$, is further divided into two parts

$$G[n] = T_s[n] + E_{xc}[n] \tag{106}$$

where $T_s[n]$ is the kinetic energy of a non-interacting electron gas of density, $n(\mathbf{r})$, in its ground state and $E_{xc}[n]$ represents the exchange and correlation energy. With these new quantities, we can write

$$E[n] = T_s[n] + \int d^3r V_{ext}(\mathbf{r})n(\mathbf{r}) + \frac{1}{8\pi\varepsilon_0} \int d^3r \int d^3r' \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} + E_{xc}[n] \tag{107}$$

The energy functional given in Eq. (107) has to be minimized with respect to the electron density, $n(\mathbf{r})$, subject only to the normalization condition

$$N = \int d^3r\, n(\mathbf{r}) \tag{108}$$

where $N$ is the total number of the electrons in the system under consideration. The standard method for taking care of the constraint given in Eq. (108) is to write the variational principle as

$$\delta(E - \mu N) = 0 \tag{109}$$

where $\mu$ is a Lagrange multiplier. Carrying out the variation, one obtains the Euler condition:

$$\mu = \frac{\delta T_s[n]}{\delta n(\mathbf{r})} + V_{ext}(\mathbf{r}) + \frac{1}{4\pi\varepsilon_0} \int d^3r' \frac{n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} + V_{xc}(\mathbf{r}) \tag{110}$$

where

$$V_{xc}(\mathbf{r}) = \frac{\delta E_{xc}[n]}{\delta n(\mathbf{r})} \tag{111}$$

is the exchange-correlation potential. The variational derivative of the kinetic energy, $\delta T_s[n]/\delta n(\mathbf{r})$, is then replaced with the kinetic operator, $-\hbar^2 \nabla^2/2m^*$. At this point, Kohn and Sham [107, 108] make a crucial observation, that the Euler Eq. (110) is the Euler equation of noninteracting particles subject to the effective external potential, $V_{eff}(\mathbf{r})$, given by

$$V_{eff}(\mathbf{r}) = V_{ext}(\mathbf{r}) + \frac{1}{4\pi\varepsilon_0} \int d^3r' \frac{n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} + V_{xc}(\mathbf{r}) = V_H(\mathbf{r}) + V_{xc}(\mathbf{r}) \tag{112}$$

where the Hartree potential is obtained from the solution of the corresponding Poisson equation. This scheme allows one to construct an equivalent one-particle formulation of the complicated many-body problem at hand.

The exchange-correlation energy, $E_{xc}[n]$, is in general an unknown functional of the electron density. However, for slowly varying density, one can make the local-density approximation (LDA):

$$E_{xc}[n] \approx \int d^3r\, \varepsilon_{xc}(n(\mathbf{r}))n(\mathbf{r}) \tag{113}$$

where $\varepsilon_{xc}(n(r))$ is the exchange and correlation energy per electron of a uniform electron gas with density $n_0 = n(r)$. The original idea for this approximation comes from Slaer [110]. Using LDA, we find

$$V_{xc}(r) \approx \frac{d(n\varepsilon_{xc}(n))}{dn} = \varepsilon_{xc}(n) + n\frac{d\varepsilon_{xc}(n)}{dn} \tag{114}$$

which, according to the Seitz theorem [111], is equivalent to the definition of the chemical potential. Therefore, $V_{xc}(r)$ can be interpreted as the exchange-correlation contribution to the chemical potential of a homogeneous electron gas of density $n_0$ equal to the local electron density $n(r)$ of the inhomogeneous system. As pointed out by Kohn and Batita [112], the LDA works surprisingly well in calculating the electronic structure of confined electronic systems where the electron density

$$n(r) = n(z) = \sum_i N_i |\psi_i^2(z)| \tag{115}$$

is not slowly varying in space. The exchange-correlation potential, $V_{xc}$, for LDA has been parameterized by many authors. A standard parameterized form due to Hedin and Lundqvist [113–116], used also in SCHRED, is

$$V_{xc}(z) = -\frac{e^2}{4\pi^2\varepsilon_0}[3\pi^2 n(z)]^{1/3}\left[1 + 0.7734x\ln\left(1 + \frac{1}{x}\right)\right] \tag{116}$$

where $x = x(z) = r_s/21$, $r_s = r_s(z) = [4\pi b^3 n(z)/3]^{-1/3}$ and $b = 4\pi\varepsilon_{sc}\hbar^2/m^*e^2$. The first term on the rhs of Eq. (116) is the exchange energy correction due to the attractive interaction between other electrons and the Fermi hole resulting from the displaced charge. The second term represents the correlation energy correction to the chemical potential, $\mu$. Using this parameterized expression for $V_{xc}(z)$, one calculates the electronic subband wave functions and the corresponding subband energies by solving the so-called Hohenberg-Kohn-Sham (HKS) equation, which is formally the same as the Schrödinger equation in which, as already noted, one takes $V_{eff}(z) = V_H(z) + V_{xc}(z) + V_{im}(z)$. In such a calculation, one obtains not only the total energy and the electron density, but also the eigenvalues of the LS equations [117]. For silicon inversion layers, by analogy to the spin-density formalism [118], the exchange-correlation correction to the chemical potential is different for the unprimed and primed valleys, and it depends only on the volume density of electrons in the unprimed and primed subbands.

The extension of this formalism for nonzero temperatures was formally set up by Mermin [119], and the finite temperature exchange-correlation functions that enter the Kohn-Sham-Mermin formulation were calculated by Gupta et al. [120, 121]. The finite-temperature exchange correction to the chemical potential calculated by Gupta et al. is

$$\frac{V_{exc}(r, T)}{V_{exc}(r, T = 0)} \approx \frac{2}{3}\left(\frac{T_F}{T}\right) \tag{117}$$

where the Fermi temperature, $T_F$, is defined in terms of the zero-temperature variables $k_F = [3\pi^2 n(r)]^{1/3}$, $E_F = \hbar^2 k_F^2/2m^*$ and $k_B T_F = E_F$. The result given in Eq. (117) is valid for Boltzmann statistics. In the Debye limit, the corresponding correlation energy correction is given by

$$V_{corr}(r, T) = -\frac{e^2}{8\pi\varepsilon_0}\sqrt{\frac{e^2 n(r)}{\varepsilon_0 k_B T}} \tag{118}$$

From the results given in Eqs. (117) and (118), it is obvious that the correlation contributions can still be important for temperatures where the exchange contribution has become vanishingly small.

A comparison of the calculated self-consistent potentials for (100) p-type Si with $N_a = 2.8 \times 10^{15}$ cm$^{-3}$ (corresponding to $N_{depl} = 2.02 \times 10^{11}$ cm$^{-2}$), $N_s = 4 \times 10^{12}$ cm$^{-2}$, and $T = 0$ K, with (thick lines) and without (thin lines) the inclusion of the exchange-correlation correction to $V_{eff}(z)$, is given in Ref. [172]. Here it was found that the subband energies are

lowered considerably by the exchange-correlation effect. The energy of the ground subband is lowered by 35 meV, whereas the energy of the first excited subband is lowered by 20 meV, which is in agreement with results obtained by Vinter [123] (Fig. 17). Since the inclusion of the exchange-correlation effects increases the subband separation, this many-body correction leads to an increase of the carrier concentration at which the occupation of the second subband begins. We also find that, in contrast to the image term which tends to increase the spatial extent of the wave functions, the exchange-correlation term tends to compress the wave functions.

In Fig. 18, comparison is made of the simulation results for the energy spacing between the lowest two subbands (subbands $\varepsilon_0$ and $\varepsilon_1$ from the unprimed ladder) with the infrared absorption measurements of Kneschaurek et al. [124] on a $p$-type Si(100) at $T = 4.2$ K. The doping concentration is $N_a - N_d = (2 \pm 0.2) \times 10^{15}$ cm$^{-3}$. The experimental data shown in the figure represent the *dark sweep* spectroscopy results. For these experimental conditions, the depletion layer length and the depletion charge density do not reach their thermal equilibrium values and the measured value of the experimentally relevant effective depletion charge density is $N^*_{depl} = (1.0 \pm 0.1) \times 10^{11}$ cm$^{-2}$. To be in agreement with the experimental conditions, a value of $N_{depl} = 1 \times 10^{11}$ cm$^{-2}$ is assumed. The experimental data shows a faster increase in the level splitting than the Hartree theory (with and without the image term). The inclusion of the exchange-correlation correction to the chemical potential in the Kohn-Sham equation significantly improves the situation, especially at higher inversion charge densities (In the space-charge layer, the condition of slowly varying potential translates into $N_s \gg N^*_{depl}$). It is believed that the so-called exciton-like and depolarization corrections nearly cancel each other except at very high electron concentrations [125, 126]. (The exciton shift is the interaction of the excited electron with the hole in the ground state, analogous to the exciton associated with the valence-to-conduction band transition. The depolarization shift is a plasmon shift of the transition caused by the screening response of the electron gas.) The simulation results shown in Fig. 18 for the energy spacing are in close agreement with those obtained by Ando [127].

In principle, the finite-temperature extension of the density-functional theory presented in this section is obtained by using the finite-temperature expressions for the exchange and correlation corrections to the chemical potential and through the change of the occupancies of various subbands. However, Das Sarma and Vinter [128, 129] have shown that neglecting any temperature dependence of the exchange-correlation potential, but retaining its implicit temperature dependence through the electron density $n(z)$, which is calculated at finite temperature, leads to results which are in very good agreement with the measured



Distance from the interface [A]

Figure 17. Calculated potential energy profile, subband structure and normalized wave functions for (100) $p$-type Si with $N_a = 2.8 \times 10^{15}$ cm$^{-3}$, $N_s = 4 \times 10^{12}$ cm$^{-2}$ and interface-trap density $N_{it} = 9.5 \times 10^{10}$ cm$^{-2}$. The thick (thin) lines correspond to the case when the exchange-correlation effect is included (omitted) in the calculation. Reprinted with permission from [122], D. Vasileska et al., *J. Vac. Sci. Technol. B* 13, 1841 (1995). © 1995, American Institute of Physics.

Figure 18. Density dependence of the separation of the $\varepsilon_1$ and $\varepsilon_0$ subbands in a (100) p-type silicon inversion layer. The filled triangles represent the infrared absorption measurements. Reprinted with permission from D. Vasileska, Ph.D. Thesis, Arizona State University, 1995. © 1995.

subband separations, especially the ones for the unprimed ladder. To check this argument, we calculated the subband separation of a p-type Si with effective depletion charge density $N_{depl}^* = 6 \times 10^{10}$ cm$^{-2}$ and (100) orientation of the surface at $T = 300$ K using first the parameterized expression given in Eq. (116), and then the finite-temperature results for the exchange-correlation corrections to the chemical potential given in Eqs. (117) and (118). The difference in the calculated subband energies for various inversion charge densities $N_s$ was found to be always less then 5%, even for the worst case. The simulation results for the subband separations $\varepsilon_{10}$ and $\varepsilon_{1'0}$ and various inversion charge densities, for the same sample, are shown in Figs. 19–20, respectively. The filled triangles in both figures represent the room-temperature infrared resonant absorption measurements due to Schäffler et al. [130]. It is believed that the net correction to the subband separation due to depolarization and exciton-like shifts is less that 4%. A total of 10(5 + 5) and 5(3 + 2) subbands were used in these simulations. It is observe that the use of 10 instead of 5 subbands leads to the increase in the subband separations in both cases throughout the whole range of $N_s$. However, this increase is more pronounced for the primed ladder of subbands. For comparison, in both figures the Hartree results are given for the subband separation. It is seen that the Hartree approximation becomes a better approximation for the subband energy difference at elevated temperatures due to the decrease of the exchange energy correction to the chemical



Figure 19. Subband energy difference $\varepsilon_{10}$ versus inversion charge density at $T = 300$ K. Reprinted with permission from D. Vasileska, Ph.D. Thesis, Arizona State University, 1995. © 1995.

Figure 20. Subband energy difference $\varepsilon_{1'0'}$ versus inversion charge density at $T = 300$ K. Reprinted with permission from D. Vasileska, Ph.D. Thesis, Arizona State University, 1995. © 1995.

potential. These simulation results for the subband energy difference for the unprimed ladder of subbands are in agreement with Refs. [128, 129]. However, for the primed ladder of subbands, they are in better agreement with the experimental data compared to the results of Das Sarma and Vinter. The major difference comes from the fact that they use 5 instead of 10 subbands as well as the conductivity instead of the density-of-states mass.

**2.5.1.2. Other Simulator Details.** In all the calculations presented here, it is assumed that the SiO$_2$/Si interface is parallel to the [100] plane. For this particular case, the six equivalent minima of the bulk silicon conduction band split into two sets of subbands (Fig. 21).

The first set ($\Delta_2$-band) consists of the two equivalent valleys with in-plane effective mass $m_{||} = 0.19\,m_0$ and perpendicular effective mass $m_\perp = 0.91\,m_0$. The second set ($\Delta_4$-band) consists of the four equivalent valleys with $m_{||} = 0.42\,m_0$ and $m_\perp = 0.19\,m_0$. The energy levels associated with the $\Delta_2$-band comprise the so-called unprimed ladder of subbands, whereas those associated with the $\Delta_4$-band comprise the primed ladder of subbands.



Figure 21. (a) Constant energy surfaces in Si together with the description of the $\Delta_2$- and $\Delta_4$-bands. Also shown are the appropriate transverse and in-plane masses for the two equivalent bands. (b) Schematics of the band bending in MOS capacitors. Also shown are the energy levels belonging to the unprimed and primed ladder of subbands corresponding to the $\Delta_2$- and $\Delta_4$-bands, respectively. The first index describes the band (=1 for the $\Delta_2$-band, and =4 for the $\Delta_4$-band), whereas the second one refers to the appropriate energy level within the band. Reprinted with permission from D. Vasileska, Ph.D. Thesis, Arizona State University, 1995. © 1995.

The self-consistent solution of the 1D Schrödinger-Poisson problem is obtained in the following way: One starts with some initial guess for the electrostatic potential and uses it to solve the 1D Schrödinger equation numerically [109]. After the eigenfunctions are determined, and the eigenvalues that characterize the electrons in the inversion layer, the inversion layer electron density appearing in the 1D Poisson equation is obtained by summing over all subbands to yield

$$n(z) = \sum_{i,j} N_{ij}\psi_{ij}^2(z) = \sum_{i,j} g_i \frac{m_i k_B T}{\pi \hbar^2} \ln\left[1 + \exp\left(\frac{E_F - E_{ij}}{k_B T}\right)\right]\psi_{ij}^2(z) \qquad (119)$$

In Eq. (119), $E_F$ is the Fermi level, $k_B$ is the Boltzmann constant, $T$ is the temperature, $m_i^1$ is the in-plane effective mass of the $i$th band, $N_{ij}$ is the sheet-charge density corresponding to the $j$th subband from the $i$th band, and $g_i$ ($g_1 = 2$ for the $\Delta_2$-band, and $g_2 = 4$ for the $\Delta_4$-band) is the band degeneracy. It is important to note that the inversion layer electrons are treated quantum mechanically only when confined by the surface field. If the electrons are not confined, or if one relied on the validity of the classical description of the inversion layer electron density, then the solution of the 1D Schrödinger equation is skipped, and one uses instead

$$n(z) = N_c F_{1/2}\left[\frac{E_F - E_c(z)}{k_B T}\right] \qquad (120)$$

where $N_c$ is the effective density of states of the conduction band. For holes, which are always treated classically for $p$-type substrates, the classical density is correspondingly written

$$p(z) = N_v F_{1/2}\left[\frac{E_c(z) - E_G - E_F}{k_B T}\right] \qquad (121)$$

where $N_v$ is the effective density of states of the valence band, and $E_G$ is the semiconductor bandgap. For the evaluation of the Fermi-Dirac integrals, which appear in Eqs. (120) and (121), we use the analytical approximation due to Bednarczyk and Bednarczyk [131]

$$F_{1/2}(x) = \left[e^{-x} + \frac{3\sqrt{\pi}}{4\nu^{3/8}}\right]^{-1} \qquad (122)$$

where

$$\nu(x) = x^4 + 50 + 33.6x\{1 - 0.68\exp[-0.17(x+1)^2]\} \qquad (123)$$

The polysilicon gates are modeled as heavily-doped single-crystal silicon. Both the electrons and holes are treated classically and assuming general Fermi-Dirac statistics, valid for degenerate semiconductors.

After the electron and hole concentrations are updated in the semiconductor and/or the polysilicon gates, the 1D Poisson equation is solved numerically for the electrostatic potential using finite-difference discretization scheme and LU decomposition method. The 1D Schrödinger equation is then solved to find the updated values for the electron density at each mesh point, and the previously described procedure is repeated until a self-consistent solution is found. It is important to note that the potential energy profile for the next iteration is obtained by using fixed-convergence factor scheme for the first 10 iterations, and the extrapolated convergence-factor scheme thereafter. The error criterion for the convergence of the self-consistent field iterations is that the absolute value of the difference between the input and output potentials at each mesh point is less than 0.01 mV.

At self-consistency, that is, once the self-consistent results are determined for the variation of the charge distribution on the semiconductor side of the MOS capacitor as a function of the gate voltage, $V_G$, the calculation is performed of the total gate capacitance $C_{tot}$. $C_{tot}$ is determined by differentiating the total induced charge density in the channel with respect to $V_G$. In contrast to some previous studies [132], where $C_{inv}$ was approximated with $\varepsilon_{sc}/\langle z\rangle_{av}$,

where $\langle z \rangle_{av}$ is the centroid of the electron density distribution, here the inversion layer capacitance is calculated by differentiating the total sheet charge density,

$$N_s = \sum_{i,j} N_{ij}$$                                                                (124)

with respect to the surface potential [133]. The depletion-layer capacitance $C_{depl}$ and the polygate capacitance $C_{poly}$ are evaluated in an analogous manner.

## 2.5.2. Some Sample Simulation Results Obtained with SCHRED

To demonstrate the existence of the two physical origins of the inversion layer capacitance, $C_{inv}$, discussed in Ref. [134], in Fig. 22 the variation of $C_{inv}$ is shown with inversion charge density $N_s$ in the channel of a MOS capacitor with substrate doping $N_A = 5 \times 10^{17}$ cm$^{-3}$, oxide thickness $t_{ox} = 4$ nm and metal gates. Exchange-correlation and image contributions to the effective potential energy term $V_{eff}(z)$, appearing in the 1D Schrödinger equation, have not been included in these simulations.

The pronounced double-slope behavior of the quantum mechanically calculated $C_{inv}$ comes from the fact that the total inversion layer capacitance can be represented as a series capacitance of two contributions. The first contribution is classical and comes from the finite density of states, that is, because a finite change in the surface potential is always necessary to increase $N_s$ (inset of Fig. 22), which, in turn, leads to finite value for $C_{inv}$. This term dominates at low values of $N_s$ (low gate voltages). The second contribution to $C_{inv}$ is due to the finite inversion layer thickness, which effectively increases the oxide thickness in terms of the total gate capacitance, thus providing an additional capacitance component. This term dominates at large gate voltages, where the inversion charge density $N_s$ significantly influences the band bending and leads to a steeper rise of the conduction band near the SiO$_2$/Si interface.

In Fig. 23, simulated $C_{tot}$ to oxide capacitance $C_{ox}$ is shown for metal/$p$-substrate and $n+$-poly/$p$-substrate MOS capacitors, as a function of the physical oxide thickness $t_{ox}$ and the doping of the polysilicon gates $N_p$, assuming $V_G = 3$ V. The high value for $V_G$, used here, may overestimate the severity of the bias dependent attenuation for thinner oxides, but a consistent value for $V_G$ is useful for the purpose of tabulating the simulated results. The results shown clearly demonstrate that classical charge model and Maxwell-Boltzmann (nondegenerate) statistics are clearly inadequate for oxide thickness below 10 nm. Even use of Fermi-Dirac statistics in the classical charge description can lead to significant errors in the estimate of the total gate capacitance for devices with metal gates and oxide thickness



Figure 22. Variation of the inversion layer capacitance with inversion charge density at $T = 300$ K. In the inset the self-consistent results are shown for the variation of the surface potential with $N_s$ when using both SC and QM descriptions of the electron density in the inversion layer. The surface potential is calculated using $\varphi_s = (E_i - E_i)_{surf} - (E_F - E_i)_{bulk}$, where $E_i$ is the intrinsic energy level. Reprinted with permission from D. Vasileska and D. K. Ferry, Proceeding of the 1st Conference MSM, Santa Clara, CA, 1998, p. 408. © 1998, D. Vasileska.

**Figure 23.** Simulated $C_{tot}$ to oxide capacitance $C_{ox}$ for metal/$p$-substrate and $n+$-poly/$p$-substrate MOS capacitors, as a function of the physical oxide thickness $t_{ox}$ and the doping of the polysilicon gates $N_p$. We use $V_G = 3$ V.

less than 5 nm because of the higher surface fields and, therefore, pronounced quantum-mechanical, size-quantization effects in the channel. For example, the classical model that uses Maxwell-Boltzmann (Fermi-Dirac) statistics predicts that, for the device with $t_{ox} = 1$ nm, $C_{tot}/C_{ox} = 0.983$ (0.882). However, the quantum-mechanical model predicts that $C_{tot}/C_{ox} = 0.795$, which leads to relative error of 23.65 (10.94)%. As previously noted, the depletion of the poly-silicon gates will further degrade the total gate capacitance.

The linear region threshold voltage shift between the QM and SC predictions for a device with $N_A = 5 \times 10^{17}$ cm$^{-3}$ and $t_{ox} = 4$ nm as a function of the doping of the polysilicon gate is shown in Fig. 24. The threshold voltage $V_{th}$ equals the gate voltage for which $Q_{inv} = 10^{-3}Q_{depl}$. As expected, the QM description of the charge in the channel increases $V_{th}$ and the shift in the threshold voltage is about 74 mV. This is due to the fact that the QM picture differs from the SC one in two ways: First, the energy spectrum is not continuous, but consists of discrete energy levels which, in turn, reduces the DOS function. Second, the energy of the ground subband from the unprimed ladder of subbands does not coincide with the bottom of the conduction band and the energy difference $\Delta E = E_{11} - E_C$ increases with increasing substrate doping. The depletion of the polysilicon gate, due to insufficient doping, further increases the threshold voltage. The additional shift in the threshold voltage due to the inclusion of the polygate depletion can be as large as 68 mV for $N_D = 10^{19}$ cm$^{-3}$, and drops down to about 18 mV for $N_D = 2 \times 10^{20}$ cm$^{-3}$.



**Figure 24.** Linear region threshold voltage shift between the QM and the SC predictions versus $N_p$. We use $N_A = 5 \times 10^{17}$ cm$^{-3}$ and $t_{ox} = 4$ nm.

The linear region threshold voltage shift for the device with $t_{ox} = 4$ nm, $N_D = 10^{20}$ cm$^{-3}$, and different substrate doping is shown in Fig. 25. Also shown in this figure are the van Dort et al. [135] experimental data for a device with metal gates and oxide thickness $t_{ox} = 14$ nm. Very close agreement between the experimentally derived threshold voltage shifts and the SCHRED simulation results for the device with 14-nm-thick oxide can be observed. A major difference from the results shown in Fig. 24 is that the inclusion of both the QM effects in the channel and polygate depletion leads to strong dependence of the threshold voltage shift upon the substrate doping $N_A$. For example, for a device with $t_{ox} = 4$ nm, $N_A = 10^{18}$ cm$^{-3}$ and $N_D = 10^{20}$ cm$^{-3}$, the inclusion of the quantum-mechanical, space-quantization effect leads to a threshold voltage shift of about 106 mV. The addition of polygate depletion leads to a further shift in the threshold voltage of about 34 mV. This observation, together with the results shown in Fig. 25, suggests that both a QM description of the charge density distribution in the channel and polygate depletion must be accounted for if accurate results for the threshold voltage are desired.

### 2.5.3. Modification of the Effective Mass Schrödinger Equation for Heterostructures

Note that in a solid in general, the momentum space is periodic and the true wave function is approximately the product of a periodic Bloch function and an envelope function. The Schrödinger equation can be used to study the evolution of the envelope wave function for an electron in the conduction band, provided that the effective mass $m^*$ is used in the Hamiltonian. When the Schrödinger equation is applied to semiconductors in the effective mass approximation, the potential $V(r)$ is assumed to be only the electrostatic potential, since the effect of the periodic crystal potential is accounted for by the effective mass itself. Such models can be used for relatively low energies close to the bottom of the conduction band, where a parabolic dispersion relation is a good approximation. In semiconductors, some of the most interesting applications of the Schrödinger equation involve spatially varying material compositions and heterojunctions. The effective mass approximation can still be used with some caution. Since the effective mass is a property of a bulk, it is not well defined in the neighborhood of a sharp material transition. In the hypothesis of slow material composition variations in space, one can adopt the Schrödinger equation with a spatially varying effective mass, taken to be the mass of a bulk with the local material properties. However, it can be shown that the Hamiltonian operator is no longer Hermitian for varying mass. A widely used Hermitian form brings the effective mass inside the differential operator as

$$-\frac{\hbar^2}{2}\nabla\cdot\left(\frac{1}{m^*}\nabla\psi\right)$$

(125)



Figure 25. Linear region threshold voltage shift between the QM and the SC predictions versus $N_A$.

This approach is extended to abrupt heterojunctions, as long as the materials on the wo sides have similar properties and bandstructure, as in the case of the GaAs/AlGaAs sysem in a certain range of the Al concentration. One has to keep in mind that very close to the heterojunctions the effective mass Schrödinger equation provides a reasonable mathematcal connection between the two regions, but the physical quantities are not necessarily vell defined. For instance, in the case of a narrow potential barrier obtained by using a thin layer of AlGaAs surrounded by GaAs, it is not clear at all what effective mass should be used for the AlGaAs, since such a region cannot be certainly approximated by a bulk. Even more difficult is to treat the case when there is a transition between direct and indirect bandgap materials (for example, GaAs and AlGaAs with large Al concentration).

Assuming a uniform mesh size $\Delta x$, the Hamiltonian given in Eq. (125) can be discretized in 1D by introducing midpoints in the mesh intervals on the two sides of the generic grid point $i$. First, the outer derivative at point $i$ is evaluated with centered finite differences, using quantities defined at points $(i - 1/2)$ and $(i + 1/2)$,

$$-\frac{\hbar^2}{2}\frac{\partial}{\partial x}\left[\frac{1}{m^*}\frac{\partial\psi}{\partial x}\right] \approx -\frac{\hbar^2}{2\Delta x}\left[\left(\frac{1}{m^*}\frac{\partial\psi}{\partial x}\right)_{i+1/2} - \left(\frac{1}{m^*}\frac{\partial\psi}{\partial x}\right)_{i-1/2}\right]$$
(126)

and then the derivatives defined on the midpoints are also evaluated with centered differences using quantities on the grid points:

$$-\frac{\hbar^2}{2\Delta x^2}\left[\frac{\psi(i+1)-\psi(i)}{m^*(i+1/2)} - \frac{\psi(i)-\psi(i-1)}{m^*(i-1/2)}\right]$$
(127)

The effective mass is the only quantity which must be known at the midpoints. If an abrupt heterojunction is located at point $i$, the abrupt change in effective mass is treated without ambiguity. It can be shown that the box integration procedure yields the same result.

Another hermitian Hamiltonian operator proposed for variable mass has the form

$$-\frac{\hbar^2}{4}\left[\frac{1}{m^*}\nabla^2\psi + \nabla^2\left(\frac{1}{m^*}\psi\right)\right]$$
(128)

which is the linear combination of two non-Hermitian operators. It is instructive to compare the two formulations. In 1D, the operators can be rewritten as follows

$$-\frac{\hbar^2}{2}\frac{\partial}{\partial x}\left[\frac{1}{m^*}\frac{\partial\psi}{\partial x}\right] = -\frac{\hbar^2}{2}\left[\frac{1}{m^*}\frac{\partial^2\psi}{\partial x^2} + \frac{\partial\psi}{\partial x}\frac{\partial}{\partial x}\left(\frac{1}{m^*}\right)\right]$$
(129)

$$-\frac{\hbar^2}{4}\left[\frac{1}{m^*}\frac{\partial^2\psi}{\partial x^2} + \frac{\partial^2}{\partial x^2}\left(\frac{1}{m^*}\psi\right)\right] = -\frac{\hbar^2}{2}\left[\frac{1}{m^*}\frac{\partial^2\psi}{\partial x^2} + \frac{\partial\psi}{\partial x}\frac{\partial}{\partial x}\left(\frac{1}{m^*}\right) + \frac{1}{2}\psi\frac{\partial^2}{\partial x^2}\left(\frac{1}{m^*}\right)\right]$$
(130)

The second operator has an additional term involving the second order derivative of the effective mass. For smoothly varying mass, the two approaches are approximately equivalent If one were to use the form on the right-hand side of Eq. (130) for discretization of the operator, it is easy to see that a direct application of finite differences is awkward. The proper procedure is to apply box integration to the interval $[i - 1/2; i + 1/2]$

$$-\frac{\hbar^2}{2}\int_{i-1/2}^{i+1/2} dx\left[\frac{1}{m^*}\frac{\partial^2\psi}{\partial x^2} + \frac{\partial\psi}{\partial x}\frac{\partial}{\partial x}\left(\frac{1}{m^*}\right)\right]$$
(131)

Integration by parts of the first term yields

$$-\frac{\hbar^2}{2}\left[\left[\frac{1}{m^*}\frac{\partial\psi}{\partial x}\right]_{i-1/2}^{i+1/2} - \int_{i-1/2}^{i+1/2} dx\frac{\partial\psi}{\partial x}\frac{\partial}{\partial x}\left(\frac{1}{m^*}\right) + \int_{i-1/2}^{i+1/2} dx\frac{\partial\psi}{\partial x}\frac{\partial}{\partial x}\left(\frac{1}{m^*}\right)\right]$$
(132)

The two integrals cancel, and if the result is divided by the integration length $\Delta x$, Eq. (127) is recovered.

## 2.6. Carrier Dynamics

Under the influence of an external field, Bloch electrons in a crystal change their wavevector according to the acceleration theorem

$$\hbar \frac{d\mathbf{k}(t)}{dt} = \mathbf{F} \tag{133}$$

where $\mathbf{F}$ is the external force acting on the particle. The effect on the actual velocity or momentum of the particle is, however, not straightforward as the velocity is related to the group velocity of the wave packet associated with the particle, and is given by

$$\mathbf{v} = \frac{1}{\hbar} \nabla_k E(\mathbf{k}) \tag{134}$$

where $E(\mathbf{k})$ is one of the dispersion relations from Fig. 8. As the particle moves through k-space under the influence of an electric field, for example, its velocity can be positive or negative, eventually leading to *Bloch oscillations* if scattering did not limit the motion. Only near extremum of the bands, for example, at the $\Gamma$ point in Fig. 8 for the valence band, or close to the L point in the conduction band, does the dispersion relation resemble that of the free electrons, $E(\mathbf{k}) = \hbar^2 k^2 / 2m^*$, where $m^*$ is the *effective mass*, which is different from the free electron mass. There, the electron velocity is simply given by $\mathbf{v} = \hbar \mathbf{k}/m^*$, and the momentum is $\mathbf{p} = \hbar \mathbf{k}$.

In the case of the valence band, the states are nearly full, and current can only be carried by the absence of electrons in a particular state, leading to the concept of *holes*, whose dynamics are identical to that of electrons except their motion is in the opposite direction of electrons, hence they behave as positively charged particles. In relation to transport and device behavior, these holes are then treated as positively charged particles in the presence of external fields, and one has to simulate the motion of both electrons and holes.

For device modeling and simulation, different approximate band models are employed. As long as carriers (electrons and holes) have relatively low energies, they may be treated using the so called *parabolic band approximation*, where they simply behave as free particles having an effective mass. If more accuracy is desired, corrections due to deviation of the dispersion relation from a quadratic dependence on $\mathbf{k}$ may be incorporated in the *nonparabolic band model*. If more than one conduction band minimum is important, this model may be extended to a *multivalley model*, where the term valley refers to different conduction minima. Finally, if the entire energy dispersion is used, one usually referes to the model as *full band*, and some of the previously described methods is usually employed.

## 3. SEMICLASSICAL TRANSPORT MODELING

Figure 3 illustrated various levels of approximation in describing charge transport within a hierarchical structure ranging from the exact quantum mechanical solution of the $n$-particle problem at the bottom, to analytical 1D phenomenological modeling used in circuit simulation at the top. The exact quantum mechanical solution of even a few particle system is a challenging computational task, and clearly impossible for a semiconductor device with typical free-carrier electron densities that are on the order of $10^{17}$ cm$^{-3}$ or more. Hence, simplifying approximations are necessary.

For conventional semiconductor devices, such as bipolar junction transistors (BJTs) and field effect transistors (FETs), the device behavior has been adequately described within the semiclassical model of charge transport, since the characteristic dimensions are typically at length scales much larger than those over which quantum mechanical phase coherence is maintained. Hence a particle-based description is adequate as described within the Boltzmann equation framework, and approximations thereof. As device dimensions continue to shrink, the channel lengths are now approaching the characteristic wavelength of particles (the de Broglie wavelength at the Fermi energy, for example), and quantum effects are expected to be increasingly important. It has in fact been well known for 30 years that quantum confinement effects occur for electrons in the inversion layers of Si metal-oxide-semiconductor field effect transistor (MOSFET) devices as discussed in Section 2.5.

However, at room temperature and under strong driving fields, such quantum effects have usually been found to be second order at best in terms of the overall device behavior. However, as discussed in Section 1.1, it is not clear that this situation will persist as all spatial dimensions are reduced, and consideration of quantum effects, such as tunneling and interference, may in fact dominate.

As mentioned earlier, the classical description of charge transport is given by the BTE in the hierarchy of Fig. 3. The BTE is an integral-differential kinetic equation of motion for the probability distribution function for particles in the 6-dimensional phase space of position and (crystal) momentum

$$\frac{\partial f(\mathbf{r},\mathbf{k},t)}{\partial t} + \frac{1}{\hbar}\nabla_{\mathbf{k}} E(\mathbf{k}) \cdot \nabla_{\mathbf{r}} f(\mathbf{r},\mathbf{k},t) + \frac{\mathbf{F}}{\hbar} \cdot \nabla_{\mathbf{k}} f(\mathbf{r},\mathbf{k},t) = \left.\frac{\partial f(\mathbf{r},\mathbf{k},t)}{\partial t}\right|_{\text{Coll}} \tag{135}$$

where $f(\mathbf{r},\mathbf{k},t)$ is the one-particle distribution function. The right-hand side is the rate of change of the distribution function due to randomizing collisions, and is an integral over the in-scattering and the out-scattering terms in momentum (wavevector) space. Once $f(\mathbf{r},\mathbf{k},t)$ is known, physical observables, such as average velocity or current, are found from averages of $f$. Equation (135) is semiclassical in the sense that particles are treated as having distinct position and momentum in violation of the quantum uncertainty relations; yet, their dynamics and scattering processes are treated quantum mechanically through the electronic band structure (discussed in Section 2) and the use of time-dependent perturbation theory. Through moment expansion of the BTE, a set of approximate partial differential equations in position space, similar to those arising in the field of fluid dynamics, are obtained leading to the so-called hydrodynamic model for charge transport, which will be discussed in Section 3.3. The simplification of the hydrodynamic model to include just the continuity equation and the current density written in terms of the local electric field and concentration gradients leads to the so-called drift-diffusion model, also discussed in Section 3.3. Finally, the reduction of the drift-diffusion model to one dimensional nonlinear analytical expressions allows for the development of lumped parameter behavioral models suitable for circuit level simulation of many individual devices as well as passive elements.

The BTE itself is an approximation to the underlying many body classical Liouville equation, and quantum mechanically by the Liouville–von Neumann equation of motion for the density matrix. The main approximations inherent in the BTE are the assumption of *instantaneous* scattering processes in space and time, the Markov nature of scattering processes (i.e., they are uncorrelated with the prior scattering events), and the neglect of multiparticle correlations (i.e., the system may be characterized by a single particle distribution function). In semiclassical simulation, some of these assumptions are relaxed through the use of molecular dynamics techniques discussed in Sections 3.1.5 and 3.1.2 (in the context of device simulations). However, the inclusion of quantum effects such as particle interference, tunneling, and so on, which take one further down the hierarchy of Fig. 3, is more problematic in the semiclassical *Ansatz* and is an active area of research today as device dimensions approach the quantum regime.

## 3.1. Direct Solution of Boltzmann Transport Equation: Monte Carlo Method

The ensemble Monte Carlo technique has been used now for more than 30 years as a numerical method to simulate nonequilibrium transport in semiconductor materials and devices and has been the subject of numerous books and reviews [136–138]. In application to transport problems, a random walk is generated to simulate the stochastic motion of particles subject to collision processes in some medium. This process of random walk generation may be used to evaluate integral equations and is connected to the general random sampling technique used in the evaluation of multidimensional integrals [139].

The basic technique is to simulate the free particle motion (referred to as the free flight) terminated by instantaneous random scattering events. The Monte Carlo algorithm consists of generating random free flight times for each particle, choosing the type of scattering occurring at the end of the free flight, changing the final energy and momentum of the

particle after scattering, and then repeating the procedure for the next free flight. Sampling the particle motion at various times throughout the simulation allows for the statistical estimation of physically interesting quantities such as the single particle distribution function, the average drift velocity in the presence of an applied electric field, the average energy of the particles, and so on. By simulating an *ensemble* of particles, representative of the physical system of interest, the nonstationary time-dependent evolution of the electron and hole distributions under the influence of a time-dependent driving force may be simulated.

The particle-based picture, in which the particle motion is decomposed into free flights terminated by instantaneous collisions, is basically the same picture underlying the derivation of the semi-classical BTE. In fact, it may be shown that the one-particle distribution function obtained from the random walk Monte Carlo technique satisfies the BTE for a homogeneous system in the longtime limit [140].

### 3.1.1. Free Flight Generation

In the Monte Carlo method, the dynamics of particle motion is assumed to consist of free flights terminated by instantaneous scattering events, which change the momentum and energy of the particle. To simulate this process, the probability density $P(t)$ is required, in which $P(t)dt$ is the joint probability that a particle will arrive at time $t$ without scattering after the previous collision at $t = 0$, and then suffer a collision in a time interval $dt$ around time $t$. The probability of scattering in the time interval $dt$ around $t$ may be written as $\Gamma[k(t)]dt$, where $\Gamma[k(t)]$ is the scattering rate of an electron or hole of wavevector $k$. The scattering rate, $\Gamma[k(t)]$, represents the sum of the contributions from each individual scattering mechanism, which are usually calculated using perturbation theory, as described later. The implicit dependence of $\Gamma[k(t)]$ on time reflects the change in $k$ due to acceleration by internal and external fields. For electrons subject to time independent electric and magnetic fields, Eq. (133) may be integrated to give the time evolution of $k$ between collisions as

$$k(t) = k(0) - \frac{e(E + v \times B)t}{\hbar} \tag{136}$$

where $E$ is the electric field, $v$ is the electron velocity [given by Eq. (136)], and $B$ is the magnetic flux density. In terms of the scattering rate, $\Gamma[k(t)]$, the probability that a particle has not suffered a collision after a time $t$ is given by $\exp(-\int_0^t \Gamma[k(t')]dt')$. Thus, the probability of scattering in the time interval $dt$ after a free flight of time $t$ may be written as the joint probability

$$P(t)dt = \Gamma[k(t)]\exp\left[-\int_0^t \Gamma[k(t')]dt'\right]dt \tag{137}$$

Random flight times may be generated according to the probability density $P(t)$ above using, for example, the pseudorandom number generator implicit on most modern computers, which generate uniformly distributed random numbers in the range [0,1]. Using a direct method (see, for example, Ref. [136]), random flight times sampled from $P(t)$ may be generated according to

$$r = \int_0^{t_r} P(t)\,dt \tag{138}$$

where $r$ is a uniformly distributed random number and $t_r$ is the desired free flight time. Integrating Eq. (138) with $P(t)$ given by Eq. (137) above yields

$$r = 1 - \exp\left[-\int_0^{t_r} \Gamma[k(t')]dt'\right] \tag{139}$$

Since $1 - r$ is statistically the same as $r$, Eq. (139) may be simplified to

$$-\ln r = \int_0^{t_r} \Gamma[k(t')]dt' \tag{140}$$

Equation (140) is the fundamental equation used to generate the random free flight time after each scattering event, resulting in a random walk process related to the underlying

particle distribution function. If there is no external driving field leading to a change of k between scattering events (for example, in ultrafast photoexcitation experiments with no applied bias), the time dependence vanishes, and the integral is trivially evaluated. In the general case where this simplification is not possible, it is expedient to introduce the so called self-scattering method [141], in which we introduce a fictitious scattering mechanism whose scattering rate always adjusts itself in such a way that the total (self-scattering plus real scattering) rate is a constant in time

$$\Gamma = \Gamma[\mathbf{k}(t')] + \Gamma_{self}[\mathbf{k}(t')] \tag{141}$$

where $\Gamma_{self}[\mathbf{k}(t')]$ is the self-scattering rate. The self-scattering mechanism itself is defined such that the final state before and after scattering is identical. Hence, it has no effect on the free flight trajectory of a particle when selected as the terminating scattering mechanism, yet results in the simplification of Eq. (140) such that the free flight is given by

$$t_r = -\frac{1}{\Gamma} \ln r \tag{142}$$

The constant total rate (including self-scattering) $\Gamma$ is chosen a priori so that it is larger than the maximum scattering encountered during the simulation interval. In the simplest case, a single value is chosen at the beginning of the entire simulation (constant gamma method), checking to ensure that the real rate never exceeds this value during the simulation. Other schemes may be chosen that are more computationally efficient, and which modify the choice of $\Gamma$ at fixed time increments [142].

### 3.1.2. Final State after Scattering

The algorithm just described determines the random free flight times during which the particle dynamics is treated semiclassically according to Eq. (136). For the scattering process itself, we need the type of scattering (i.e., impurity, acoustic phonon, photon emission, etc.), which terminates the free flight, and the final energy and momentum of the particle(s) after scattering. The type of scattering which terminates the free flight is chosen using a uniform random number between 0 and $\Gamma$, and using this pointer to select among the relative total scattering rates of all processes including self-scattering at the final energy and momentum of the particle:

$$\Gamma = \Gamma_{self}[n, \mathbf{k}] + \Gamma_1[n, \mathbf{k}] + \Gamma_2[n, \mathbf{k}] + \cdots + \Gamma_N[n, \mathbf{k}] \tag{143}$$

with $n$ the band index of the particle (or subband in the case of reduced-dimensionality systems), and $\mathbf{k}$ the wave vector at the end of the free-flight.

Once the type of scattering terminating the free flight is selected, the final energy and momentum (as well as band or subband) of the particle due to this type of scattering must be selected. For this selection, the scattering rate, $\Gamma_j[n, \mathbf{k}; m, \mathbf{k}']$, of the $j$th scattering mechanism is needed, where $n$ and $m$ are the initial and final band (subband) indices, and $\mathbf{k}$ and $\mathbf{k}'$ are the particle wavevectors before and after scattering. Defining a spherical coordinate system around the initial wavevector $\mathbf{k}$, the final wavevector $\mathbf{k}'$ is specified by $|\mathbf{k}'|$ (which depends on conservation of energy) as well as the azimuthal and polar angles, $\varphi$ and $\theta$ around $\mathbf{k}$. Typically, the scattering rate, $\Gamma_j[n, \mathbf{k}; m, \mathbf{k}']$, only depends on the angle $\theta$ between $\mathbf{k}$ and $\mathbf{k}'$. Therefore, $\varphi$ may be chosen using a uniform random number between 0 and $2\pi$ (i.e., $2\pi r$), while $\theta$ is chosen according to the cross section for scattering arising from $\Gamma_j[n, \mathbf{k}; m, \mathbf{k}']$. If the probability for scattering into a certain angle $P(\theta)d\theta$ is integrable, then random angles satisfying this probability density may be generated from a uniform distribution between 0 and 1 through inversion of Eq. (140). Otherwise, a rejection technique (see, for example, Ref. [136, 137]) may be used to select random angles according to $P(\theta)$.

### 3.1.3. Ensemble Monte Carlo Simulation

The algorithm above may be used to track a single particle over many scattering events in order to simulate the steady-state behavior of a system. Transient simulation requires the use of a synchronous ensemble of particles in which the algorithm above is repeated for each particle in the ensemble representing the system of interest until the simulation is completed.

Figure 26 illustrates an ensemble Monte Carlo simulation in which a fixed time step, $\Delta t$, is introduced to which the motion of all the carriers in the system is synchronized. The symbols illustrate random, instantaneous, scattering events, which may or may not occur during one time step. Basically, each carrier is simulated only up to the end of the time step, and then the next particle in the ensemble is treated. Over each time step, the motion of each particle in the ensemble is simulated independent of the other particles. Nonlinear effects such as carrier-carrier interactions or the Pauli exclusion principle are then updated at each times step, as discussed in more detail below.

The nonstationary one-particle distribution function and related quantities such as drift velocity, valley or subband population, and so on, are then taken as averages over the ensemble at fixed time steps throughout the simulation. For example, the drift velocity in the presence of the field is given by the ensemble average of the component of the velocity at the $n$th time step as

$$\bar{v}_z(n\Delta t) \cong \frac{1}{N} \sum_{j=1}^{N} v_z^j(n\Delta t) \tag{144}$$

where $N$ is the number of simulated particles and $j$ labels the particles in the ensemble. This equation represents an estimator of the true velocity, which has a standard error given by

$$s = \frac{\sigma}{\sqrt{N}} \tag{145}$$

where $\sigma^2$ is the variance which may be estimated from [139]

$$\sigma^2 \cong \frac{N}{N-1} \left\{ \frac{1}{N} \sum_{j=1}^{N} (v_z^j)^2 - \bar{v}_z^2 \right\} \tag{146}$$

Similarly, the distribution functions for electrons and holes may be tabulated by counting the number of electrons in cells of $k$-space. From Eq. (145), we see that the error in estimated average quantities decreases as the square root of the number of particles in the ensemble, which necessitates the simulation of many particles. Typical ensemble sizes for good statistics are in the range of $10^4$–$10^5$ particles. Variance reduction techniques to decrease the standard error given by Eq. (146) may be applied to enhance statistically rare events such as impact ionization or electron-hole recombination [137].

## 3.1.4. Scattering Processes

Free carriers (electrons and holes) interact with the crystal and with each other through a variety of scattering processes which relax the energy and momentum of the particle. On the basis of first-order, time-dependent perturbation theory, the transition rate from an initial state $\mathbf{k}$ in band $n$ to a final state $\mathbf{k}'$ in band $m$ for the $j$th scattering mechanism is given by



Figure 26. Ensemble Monte Carlo simulation in which a time step, $\Delta t$, is introduced over which the motion of particles is synchronized. The crosses (×) represent scattering events. Reprinted with permission from D. Vasileska and S. Goodnick, "Encyclopedia of Materials, Science and Technology." Elsevier, 2001, p. 1. © 2001, Elsevier.

Fermi's golden rule [143]:

$$\Gamma_j[n, \mathbf{k}; m, \mathbf{k}'] = \frac{2\pi}{\hbar} |\langle m, \mathbf{k}'|V_j(\mathbf{r})|n, \mathbf{k}\rangle|^2 \delta(E_{\mathbf{k}'} - E_{\mathbf{k}} \mp \hbar\omega) \tag{147}$$

where $V_j(\mathbf{r})$ is the scattering potential of this process, and $E_k$ and $E_{k'}$ are the initial and final state energies of the particle. The delta function results in conservation of energy for long times after the collision is over, with $\hbar\omega$ the energy absorbed (upper sign) or emitted (lower sign) during the process. Scattering rates calculated by Fermi's golden rule above are typically used in Monte Carlo device simulation as well as simulation of ultrafast processes. The total rate used to generate the free flight in Eq. (142), discussed in the previous section, is then given by

$$\Gamma_j[n, \mathbf{k}] = \frac{2\pi}{\hbar} \sum_{m, \mathbf{k}'} |\langle m, \mathbf{k}'|V_j(\mathbf{r})|n, \mathbf{k}\rangle|^2 \delta(E_{\mathbf{k}'} - E_{\mathbf{k}} \mp \hbar\omega) \tag{148}$$

There are major limitations to the use of the golden rule due to effects such as *collision broadening* and *finite collision duration time* [140]. The energy conserving delta function is only valid asymptotically for times long after the collision is complete. The broadening in the final state energy is given roughly by $\Delta E \approx \hbar/\tau$, where $\tau$ is the time after the collision, which implies that the normal $E(\mathbf{k})$ relation is only recovered at long times. Attempts to account for such *collision broadening* in Monte Carlo simulation have been reported in the literature [144, 145], although this is still an open subject of debate. Inclusion of the effects of *finite collision duration* in Monte Carlo simulation have also been proposed [146, 147]. Beyond this, there is still the problem of dealing with the quantum mechanical phase coherence of carriers, which is neglected in the scatter free-flight algorithm of the Monte Carlo algorithm. This topic is discussed later in Section 6.

Figure 27 lists the scattering mechanisms one should in principle consider in a typical Monte Carlo simulation. They are roughly divided into scattering due to crystal defects, which is primarily elastic in nature, lattice scattering between electrons (holes) and lattice vibrations or phonons, which is inelastic, and finally scattering between the particles themselves, including both single particle and collective type excitations. Phonon scattering involves different modes of vibration, either acoustic or optical, as well as both transverse and longitudinal modes. Carriers may either emit or absorb quanta of energy from the lattice, in the form of phonons, in individual scattering events. The designation of inter- versus intravalley scattering comes from the multivalley band-structure model mentioned in Section 2 and refers to whether the initial and final states are in the same valley or in different valleys. The scattering rates $\Gamma_j[n, \mathbf{k}; m, \mathbf{k}']$ and $\Gamma_j[n, \mathbf{k}]$ are calculated using time-dependent perturbation theory using Fermi's rule, Eqs. (147) and (148), and the calculated rates are then tabulated in a scattering table to select the type of scattering and final state after scattering as discussed earlier.



Figure 27. Scattering mechanisms in a typical semiconductor. Reprinted with permission from D. Vasileska and S. Goodnick, "Encyclopedia of Materials. Science and Technology." Elsevier. 2001. p. 1. © 2001, Elsevier.

### 3.1.5. Multicarrier Effects

Multiparticle effects relate to the interaction between particles in the system, which is a nonlinear effect when viewed in the context of the BTE, due to the dependence of such effects on the single particle distribution function itself. Most algorithms developed to deal with such effects essentially linearize the BTE by using the previous value of the distribution function to determine the time evolution of a particle over the successive time-step. Multicarrier effects may range from simple consideration of the Pauli exclusion principle (which depends on the exact occupancy of states in the system), to single particle and collective excitations in the system. Inclusion of carrier-carrier interactions in Monte Carlo simulation has been an active area of research for quite some time and is briefly discussed below. Another carrier-carrier effect that is of considerable importance when estimating leakage currents in MOSFET devices, is impact ionization, which is a pure generation process involving three particles (two electrons and a hole or two holes and an electron). The latter is also discussed below.

#### 3.1.5.1. Pauli Exclusion Principle.
The Pauli exclusion principle requires that the bare scattering rate given by Eq. (147) be modified by a factor $1 - f_m(\mathbf{k}')$ in the collision integral of the BTE, where $f_m(\mathbf{k}')$ is the one-particle distribution function for the state $\mathbf{k}'$ in band (subband) $m$ after scattering. Since the net scattering rate including the Pauli exclusion principle is always less than the bare scattering rate, a self-scattering rejection technique may be used in the Monte Carlo simulation as proposed by Bosi and Jacoboni [148] for one particle simulation and extended by Lugli and Ferry [149] for EMC. In the self-scattering rejection algorithm, an additional random number $r$ is generated (between 0 and 1), and this number is compared with $f_m(\mathbf{k}')$, the occupancy of the final state (which is also between 0 and 1 when properly normalized for the numerical $\mathbf{k}$-space discretization). If $r$ is greater than $f_m(\mathbf{k}')$, the scattering is accepted and the particle's momentum and energy are changed. If this condition is not satisfied, the scattering is rejected, and the process is treated as a self-scattering event with no change of energy or momentum after scattering. Through this algorithm, no scattering to this state can occur if the state is completely full.

#### 3.1.5.2. Carrier-Carrier Interactions.
Carrier-carrier interactions, apart from degeneracy effects, may be treated as a scattering process within the Monte Carlo algorithm on the same footing as other mechanisms. In the simplest case of bulk electrons in a single parabolic conduction band, the process may be treated as a binary collision where the scattering rate for a particle of wave vector $\mathbf{k}_0$ due to all the other particles in the ensemble is given by [150]:

$$\Gamma_{cc}(\mathbf{k}_0) = \frac{n m_e e^4}{4\pi \hbar^3 \varepsilon^2 \beta^2} \int d\mathbf{k} f(\mathbf{k}) \frac{|\mathbf{k} - \mathbf{k}_0|}{(|\mathbf{k} - \mathbf{k}_0|^2 + \beta^2)} \tag{149}$$

where $f(\mathbf{k})$ is the one-particle distribution function (normalized to unity), $\varepsilon$ is the permittivity, $n$ is the electron density, and $\beta$ is the screening constant. In deriving Eq. (149), one assumes that the two particles interact through a statically screened Coulomb interaction, which ignores the energy exchange between particles in the screening which in itself is a dynamic, frequency-dependent effect. Similar forms have been derived for electrons in 2D [151, 152] and 1D [153], where carrier-carrier scattering leads to inter-subband as well as intra-subband transitions. Since the scattering rate in Eq. (149) depends on the distribution function of all the other particles in the system, this process represents a nonlinear term as discussed earlier. One method is to tabulate $f(\mathbf{k})$ on a discrete grid, as is done for the Pauli principle, and then numerically integrate Eq. (149) at each time step. An alternate method is to use a self-scattering rejection technique [154], where the integrand excluding $f(\mathbf{k})$ is replaced by its maximum value and taken outside the integral over $\mathbf{k}$. The integral over $f(\mathbf{k})$ is just unity, giving an analytic form used to generate the free flight. Then, the self-scattering rejection technique is used when the final state is chosen to correct for the exact scattering rate compared to this artificial maximum rate, similar to the algorithm used for the Pauli principle.

The treatment of intercarrier interactions as binary collisions above neglects scattering by collective excitations such as plasmons or coupled plasmon-phonon modes. These effects

may have a strong influence on carrier relaxation, particularly at high carrier density. One approach is to make a separation of the collective and single-particle spectrum of the interacting many-body Hamiltonian, and treat them separately (i.e., as binary collisions for the single particle excitations), and as electron-plasmon scattering for the collective modes [155]. Another approach is to calculate the dielectric response within the random phase approximation, and associate the damping given by the imaginary part of the inverse dielectric function with the electron lifetime [156].

A semiclassical approach to carrier-carrier interaction, which is fully compatible with the Monte Carlo algorithm, is the use of Molecular Dynamics [157], in which carrier-carrier interaction is treated continuously in real space during the free-flight phase through the Coulomb force of all the particles. A very small time step is required when using molecular dynamics to account for the dynamic distribution of the system. A time step on the order of 0.5 $fs$ is often sufficiently small for this purpose. The small time step assures that the forces acting on the particles during the time of flight are essentially constant, that is $f(t) \cong f(t + \Delta t)$, where $f(t)$ is the single particle distribution function.

Using Newtonian kinematics, we can write the real space trajectories of each particle as

$$\mathbf{r}(t + \Delta t) = \mathbf{r}(t) + \mathbf{v}\Delta t + \frac{1}{2}\frac{\mathbf{F}(t)}{m}\Delta t^2 \tag{150}$$

and

$$\mathbf{v}(t + \Delta t) = \mathbf{v}(t) + \frac{\mathbf{F}(t)}{m}\Delta t \tag{151}$$

Here, $\mathbf{F}(t)$ is the force arising from the applied field as well as that of the Coulomb interactions. We can write $\mathbf{F}(t)$ as

$$\mathbf{F}(t) = q\left[\mathbf{E} - \sum_i \nabla\varphi(\mathbf{r}_i(t))\right] \tag{152}$$

where $q\mathbf{E}$ is the force due to the applied field and the summation is the interactive force due to all particles separated by distance $\mathbf{r}_i$, with $\varphi(\mathbf{r}_i)$ the electrostatic potential. As in Monte Carlo simulation, one has to simulate a finite number of particles due to practical computational limitations on execution time. In real space, this finite number of particles corresponds to a particular simulation volume given a certain density of carriers, $V = N/n$, where $n$ is the density. Since the carriers can move in and out of this volume, and since the Coulomb interaction is a long-range force, one must account for the region outside $V$ by periodically replicating the simulated system. The contributions due to the periodic replication of the particles inside $V$ in cells outside has a closed form solution in the form of an Ewald sum [158], which gives a linear as well as $1/r^2$ contribution to the force. The equation for the total force in the molecular dynamics technique then becomes

$$\mathbf{F} = \frac{-e^2}{4\pi\varepsilon}\sum_i^N\left(\frac{1}{r_i^3}\mathbf{a}_i + \frac{2\pi}{3V}\mathbf{r}_i\right) \tag{153}$$

The above equation is easily incorporated in the standard Monte Carlo simulation discussed up to this point. At every time step the forces on each particle due to all the other particles in the system are calculated from Eq. (152). From the forces, an interactive electric field is obtained which is added to the external electric field of the system to couple the molecular dynamics to the Monte Carlo.

The inclusion of the carrier-carrier interactions in the context of particle-based device simulations is discussed in Section 3.1.2. The main difficulty in treating this interaction term in device simulations arises from the fact that the long-range portion of the carrier-carrier interaction is included via the numerical solution of the quasi-static Poisson equation (see Section 4.2). Under these circumstances, special care has to be taken when incorporating the short-range portion of this interaction term to prevent double counting of the force.

**3.1.5.3. Band-to-Band Impact Ionization.**    Another carrier-carrier scattering process is that of impact ionization, in which an energetic electron (or hole) has sufficient kinetic energy to create an electron-hole pair. Impact ionization therefore leads to the process of

carrier multiplication. This process is critical for example in the avalanche breakdown of semiconductor junctions, and is a detrimental effect in short channel MOS devices in terms of excess substrate current and decreased reliability.

The ionization rate of valence electrons by energetic conduction band electrons is usually described by Fermi's rule Eq. (147), in which a screened Coulomb interaction is assumed between the two particles, as discussed earlier in this section, where screening is described by an appropriate dielectric function such as that proposed by Levine and Louie [159]. In general, the impact ionization rate should be a function of the wavevector of the incident electron, hence of the direction of an electric field in the crystal, although there is still some debate as to the experimental and theoretical evidence. More simply, the energy dependent rate (averaged over all wavevectors on a constant energy shell) may be expressed analytically in the power law form

$$\Gamma_{ii}(E) = P[E - E_{th}]^a \tag{154}$$

where $E_{th}$ is the threshold energy for the process to occur, which is determined by momentum and energy conservation considerations, but minimally is the bandgap of the material itself. $P$ and $a$ are parameters that may be fit to more sophisticated models. The Keldysh formula [160] is derived by expanding the matrix element for scattering close to threshold, which gives $a = 2$, and the constant $P = C/E_{th}^2$, with $C = 1.19 \times 10^{14}$/s and assuming a parabolic band approximation,

$$E_{th} = \frac{3 - 2m_v/m_c}{1 - m_v/m_c} E_g \tag{155}$$

where $m_v$ and $m_c$ are the effective masses of the valence and conduction band respectively, and $E_g$ is the bandgap. More complete full-bandstructure calculations of the impact ionization rate have been reported for Si [161, 162], GaAs [162, 163] and wide bandgap materials [164], which are fairly well fit using the power law model given in Eq. (159).

Within the ensemble Monte Carlo method, the maximum scattering rate is used to generate the free flight time. The state after scattering of the initial electron plus the additional electron and hole must satisfy both energy and momentum conservation within the Fermi rule model, which is somewhat complicated unless simple parabolic band approximations are made.

### 3.1.6. Full-Band Particle-Based Simulation

The Monte Carlo algorithm discussed in this section initially evolved during the 1970s and early 1980s using simplified representations of the electronic bandstructure in terms of a multivalley parabolic or nonparabolic approximation close to band minima and maxima. This simplifies the particle tracking in terms of the $E$-$k$ relationship and particle motion in real space, and greatly simplifies the calculated scattering rates such that analytical forms may be used. It soon became apparent that for devices where high field effects are important, or for the correct simulation of high energy processes like impact ionization, the full band structure of the material is required. Particle based simulation which incorporates part or all of the band structure directly into the particle dynamics and scattering is commonly referred to as *full-band* Monte Carlo simulation [138].

Typically, the empirical pseudopotential method (EPM) discussed in Section 2.1 has been utilized in full band Monte Carlo codes due to the relative simplicity of the calculation, and the plane wave basis which facilitates calculation of some scattering processes. Early full band codes developed at the University of Illinois utilized the full bandstructure for the particle dynamics, but assumed isotropic energy dependent scattering rates using the full band density of states [138]. This is due to the computational difficulty and memory requirements to store the full $k$-dependent scattering rates throughout the whole Brillouin zone. Later simulators relaxed this restriction, although often assuming quasi-isotropic rates. Probably the most completely developed full-band code for full-band Monte Carlo device simulation is the DAMOCLES code developed at IBM by Fischetti and Laux [261], which has been used extensively for simulation of a variety of device technologies [166].

The full-band codes above are based on essentially the same algorithm as was discussed in Section 3.1, in which a particle scatters based on the total scattering rate, then the type of scattering and the final state after scattering are selected using the full $k$-dependent rates for each mechanism. An alternative approach, referred to as Cellular Monte Carlo [167], stores the entire transition table for the total scattering rate for all mechanisms from every initial state $k$ to every final state $k'$. Particle scattering is accomplished in a single step, at the expense of large memory consumption (on the order of 2 Gbytes of RAM) necessary to store the necessary scattering tables.

Figure 28 shows the calculated steady state drift velocity and average energy for Si as a function of electric field for the CMC method and the earlier results from DAMOCLES which are essentially the same. In such simulations, steady state is typically reached after 2 ps of simulation time, and then averages are calculated over the ensemble and in time for several picoseconds thereafter.

## 3.2. Semiclassical Transport for Low-Dimensional Systems

As mentioned earlier, confinement of electrons in a heterojunction or oxide-semiconductor interface result in quantum confinement of the motion, and reduction of the dimensionality of carriers to form a quasi two-dimensional electron gas (2DEG). The properties of two-dimensional electrons in the inversion layer of a Si-SiO$_2$ MOS structure were extensively studied in the 1970s and 1980s, where a thorough review is in Ref. [168]. The band edge profile across an MOS structure is shown in Fig. 29. Typically MOS structures are fabricated on (100) Si substrates, where the SiO$_2$ layer is an insulator due to its wide bandgap. A 2DEG is induced electrostatically by application a positive voltage $V_G$, forming an inversion layer if the substrate is $p$-type. The sheet density of 2DEG can be described as

$$N_s = \frac{\varepsilon_{ox}}{ed_{ox}}(V_G - V_T)$$                                   (156)

where $V_T$ is the threshold voltage for the formation of an inversion layer.

Another important 2DEG system is the modulation-doped GaAs-AlGaAs heterostructures, shown in Fig. 30. The bandgap of AlGaAs is wider than in GaAs, hence confining electrons into the GaAs side. As Fig. 30 illustrates, when chemical equilibrium is established after formation of a heterojunction, an inversion layer is formed at the interface.

The 2DEG created by modulation doping can be further confined into narrow one-dimensional (1D) channels by selective depletion in spatially separated regions. The simplest lateral confinement technique is to create split metallic gates [8, 9], or quantum point contacts, as shown in Fig. 31.

### 3.2.1. Wave Functions

In Section 2.5, a 1D self-consistent Schrödinger-Poisson solver (SCHRED) for calculating the energy eigenvalues and the corresponding eigenfunctions for a MOS capacitor structure



Figure 28. Comparison of full-band Monte Carlo simulation results using DAMOCLES [165] (triangles) to those using the CMC approach [167]. The upper plot is the steady state drift velocity and the lower plot the average energy versus electric field. Reprinted with permission from [167], M. Saraniti and S. M. Goodnick, *IEEE Trans. Electron Devices* 47, 1909 (2000). © 2000, IEEE.

**Figure 29.** Band diagram showing conductance band $E_c$, valence band $E_v$ and quasi-Fermi level $E_f$. A 2DEG is formed at the interface between the oxide (SiO)$_2$ and $p$-type silicon substrate as a consequence of the gate voltage $V_g$.

was described. As discussed earlier in connection with SCHRED, $V_H$ is the Hartree contribution to the total potential energy, which is obtained as a solution of the 1D Poisson equation. Under certain conditions, one can approximate the potential energy term and obtain approximate analytical values for the eigenvalues and the eigenvectors.

For example, if the inversion layer contribution to the Hartree potential energy is negligible compared to the depletion layer contribution $w$, for small values of $z(z \ll w)$ one can approximate the potential energy profile given in a MOS capacitor by

$$V(z) \approx \frac{e^2 N_A w}{\varepsilon_{sc}} = eE_s z \tag{157}$$

where $E_s$ is the electric field at the semiconductor-insulator interface ($z = 0$) and $w$ is the width of the depletion region. This is called the triangular-potential approximation. It leads to the Airy equation with solutions [169, 170]:

$$\psi_i(z) = Ai\left[\frac{2m_z^* eE_s}{\hbar^2}\left(z - \frac{\varepsilon_i}{eE_s}\right)\right] \tag{158}$$

and

$$\varepsilon_i \approx \left(\frac{\hbar^2}{2m_z^*}\right)^{1/3}\left[\frac{3\pi eE_s}{2}\left(i + \frac{3}{4}\right)\right]^{2/3} \tag{159}$$



**Figure 30.** Band structure of the interface between $n$-AlGaAs and intrinsic GaAs. (a) before and (b) after the charge transfer.

**Figure 31.** On the formation of a 1D channel by a split gate superimposed over a 2DEG structure.

The eigenvalues $\varepsilon_i$ are asymptotic values for large $i$, but they are quite close even for the ground state $i = 0$. The exact eigenvalues for the three lowest states have $i + \frac{3}{4}$ in Eq. (159) replaced by 0.7587, 1.7540, and 2.7525, respectively.

The triangular-potential approximation is a reasonable approximation when there is little or no charge in the inversion layer, but fails when the inversion charge density is comparable to or greater than the depletion charge density $N_{\text{depl}} = N_I w$. At low temperatures and moderately high inversion charge densities and (100) orientation of the surface, only the lowest subband of the two equivalent valleys with longitudinal mass perpendicular to the interface is usually occupied. In this case, a variational approach gives a good estimate for the energy of the lowest subband. One can approximate the wave function of the lowest subband with some trial function. The trial function proposed by Fang and Howard [171] is

$$\psi_0(z) = \left(\frac{b^3}{2}\right)^{1/2} z e^{-bz/2} \tag{160}$$

where $b$ is a parameter which is determined by minimizing the total energy of the system for given values of the inversion and depletion charges. The total energy per electron, the quantity that needs to be minimized, is then given by

$$\frac{E}{N} = T + V_d + V_i = \left(\frac{\hbar^2 b^2}{8m_z^*}\right) + \left(\frac{3e^2 N_{\text{depl}}}{\varepsilon_{sc} b} - \frac{6e^2 N_I}{\varepsilon_{sc} b^2}\right) + \frac{1}{2}\left(\frac{33e^2 N_s}{16\varepsilon_{sc} b}\right) \tag{161}$$

The first term in Eq. (161) represents the expectation value of the kinetic energy of the electron, while the second and third terms correspond to the average potential energies of the electron interacting with the depletion and inversion charge. The factor 1/2 in Eq. (161) prevents double counting of the electrons. After some algebra, one finds that the value of $b$ that minimizes the average energy per electron is

$$b = \left(\frac{12m_z^* e^2 N^*}{\hbar^2 \varepsilon_{sc}}\right)^{1/3} \tag{162}$$

where

$$N^* = N_{\text{depl}} + \frac{11}{32}N_s \tag{163}$$

The energy of the lowest state is found after a straightforward calculation to be

$$\varepsilon_0 = T + V_d + V_i \tag{164}$$

and the inversion charge contribution to $V_{ii}(z)$ reduces to

$$V_i(z) = \frac{e^2}{\varepsilon_{sc}} \sum_i N_i \left[\frac{3}{h}(1 - e^{-bz}) - ze^{-bz}\left(2 + \frac{1}{2}bz\right)\right] \tag{165}$$

### 3.2.2. Density of States

The density of states $g(E)$ is defined as the number of states per energy interval $E$, $E + dE$. It is clear that

$$g(E) = \sum_{\alpha} \delta(E - E_{\alpha}) \tag{166}$$

where $\alpha$ is the set of quantum numbers characterizing the states. In the present case, it includes the subband quantum number $n$, spin quantum number $\sigma$, valley quantum number $v$ and the in-plane quasi-momentum $\mathbf{k}$. If the spectrum is degenerate with respect to spin and valleys, one can define the spin degeneracy $\nu_s$ and the valley degeneracy $\nu_v$ to get

$$g(E) = \frac{\nu_s \nu_v}{(2\pi)^d} \sum_n \int d^d k \, \delta(E - E_n) \tag{167}$$

Here we calculate the number of states per unit volume, $d$ being the dimension of the space. For 2D case, we obtain easily

$$g(E) = \frac{\nu_s \nu_v m}{2\pi \hbar^2} \sum_n \Theta(E - E_n) \tag{168}$$

Within a given subband the 2D density of states function is energy independent. Since there can exist several subbands in the confining potential, the total density of states can be represented as a set of steps, as shown in Fig. 32. At low temperature ($k_B T \ll E_F$) all the states are filled up to the Fermi level. Because of energy-independent density of states, the sheet electron density is linear in the Fermi energy, namely

$$N_s = N \frac{\nu_s \nu_v m E_F}{2\pi \hbar^2} \tag{169}$$

The Fermi momentum in each subband can be determined as

$$k_{Fn} = \frac{1}{\hbar} \sqrt{2m(E_F - E_n)} \tag{170}$$

In Eq. (169), $N$ is the number of transverse modes having the edges $E_n$ below the Fermi energy.

The situation is more complicated if the gas is confined in a narrow, 1D channel, say, along $y$-axis. In a similar way, the in-plane wave function can be decoupled as a product

$$\psi(\mathbf{r}) \propto \eta(y) e^{ik \cdot x} \tag{171}$$

the corresponding energy being

$$E_{n,s,k} = E + E_s(k_x) + \frac{\hbar^2 k_x^2}{2m} \tag{172}$$



Figure 32. Density of states for a quasi-2D system.

In Eq. (172), $E_{n_s} = E_n + E_s$ characterizes the energy level in the potential confined in both ($z$ and $y$) directions. For square-box confinement, the terms are

$$E_s = \frac{(s\pi\hbar)^2}{2mW^2} \tag{173}$$

where $W$ is the channel width, while for the parabolic confinement $U(y) = (1/2)m\omega_0^2 y^2$ (typical for split-gate structures)

$$E_s = \left(s - \frac{1}{2}\right)\hbar\omega_0 \tag{174}$$

For such system, the total density of states is

$$g(E) = \frac{V_x V_y \sqrt{m}}{2^{3/2}\pi\hbar} \sum_{n_s} \frac{\Theta(E - E_{n_s})}{\sqrt{E - E_{n_s}}} \tag{175}$$

which is singular at the 1D subband edge. The energy dependence of the density of states for the case of parabolic confinement is shown in Fig. 33.

### 3.2.3. Scattering Rate Calculation for Low-Dimensional Systems

Calculation of the scattering rates for confined carriers proceeds in a similar manner as in the 3D case (see Section 3.1), although with initial and final states defined by not only the initial final wavevector of the electron, but the confined states corresponding to initial and final subbands as well. Before going into the details of the matrix elements of some of the major scattering mechanisms listed in Fig. 27, some general expressions are first derived. Suppose we want to calculate the scattering out of some initial state $\mathbf{k}$ in subband $n$, Fermi's golden rule can then be employed, which gives the transition rate from a state $\mathbf{k}$ in subband $n$ into a state $\mathbf{k}'$ belonging to a subband $m$:

$$S_{nm}(\mathbf{k}, \mathbf{k}') = \frac{2\pi}{\hbar}|M(\mathbf{k}, \mathbf{k}')|^2_{nm}\delta(E - E' \pm \hbar\omega) \tag{176}$$

Assuming a plane-wave basis for the wave functions in the unconfined direction ($xy$-plane), the total wave functions of the initial and the final states are of the following general form (for a quasi-2DEG):

$$\psi_n(\mathbf{k}, z) = \frac{1}{\sqrt{A}}e^{i\mathbf{k}\cdot\mathbf{r}}\varphi_n(z), \quad \psi_m(\mathbf{k}', z') = \frac{1}{\sqrt{A}}e^{i\mathbf{k}'\cdot\mathbf{r}}\varphi_m(z') \tag{177}$$



Figure 33. Density of states for a quasi-1D system (solid line) and the number of states (dashed lines).

where $A$ is the area of the sample, $\mathbf{r}$ is the wavevector in the $xy$-plane and $\mathbf{R} = (\mathbf{r}, z)$ is a 3D wavevector. The matrix element for scattering between states $\mathbf{k}$ and $\mathbf{k}'$ in subbands $n$ and $m$ is then given by

$$M(\mathbf{k}, \mathbf{k}')_{nm} = \frac{1}{A} \int e^{i(\mathbf{k} - \mathbf{k}') \cdot \mathbf{r}} d^2 r \int dz \varphi_m^*(z) H_{qr}(\mathbf{R}) \varphi_n(z) \tag{178}$$

where $H_{qr}$ is the interaction potential and the form of the integral with respect to $z$ depends upon the type of the scattering mechanics considered [172]. It is also important to note that in low-dimensional systems, since the momentum is quantized in one or two directions to form subbands, there are additional *intrasubband* and *intersubband* transitions, which complicates the generation of scattering tables and choice of the final state after scattering. Below, the matrix elements for some of the most important scattering mechanisms are discussed in detail for quasi-2D systems.

### 3.2.3.1. Coulomb Scattering.
Scattering associated with the Coulomb centers near the plane of the 2DEG in either a heterostructure system or MOS devices can be separated in contributions from the depletion layer, the interface charge and the barrier/oxide charge. An extensive discussion of the role of multiple-scattering contributions to the electron mobility of a doped semiconductor and the apparent difficulties with the impurity averaging is presented in the papers by Moore [173], and Kohn and Luttinger [174, 175]. The expressions for the potential due to a single charge located in the region of interest in which the image term is properly included, are given by Stern and Howard [176], which is generalized below for the many-subband case.

Using the usual method of images [177], one easily finds that, in the presence of a discontinuity at $z = 0$ in the dielectric constant between two materials, the potential due to a charged center located at $\mathbf{R}_i = (\mathbf{r}_i, z_i)$ equals

$$U_i(\mathbf{r}, z) = \frac{e^2}{4\pi k} \frac{1}{\sqrt{(\mathbf{r} - \mathbf{r}_i)^2 + (z - z_i)^2}} \tag{179}$$

for $z_i < 0$, and

$$U_i(\mathbf{r}, z) = \frac{e^2}{4\pi k} \left[ \frac{1}{2}\left(1 + \frac{\varepsilon_{ox}}{\varepsilon_{sc}}\right) \frac{1}{\sqrt{(\mathbf{r} - \mathbf{r}_i)^2 + (z - z_i)^2}} + \frac{1}{2}\left(1 - \frac{\varepsilon_{ox}}{\varepsilon_{sc}}\right) \frac{1}{\sqrt{(\mathbf{r} - \mathbf{r}_i)^2 + (z - z_i)^2}} \right] \tag{180}$$

for $z_i > 0$. In Eqs. (179) and (180), $\kappa = 0.5(\varepsilon_{sc} + \varepsilon_{ox})$ is the average dielectric constant at the interface.

Depletion charge scattering is due to the ionized charges in the depletion layer, which is relevant for both the MOS case, and for a simple AlGaAs/GaAs heterojunction. Using Eq. (180), the matrix element squared for scattering between subbands $n$ and $m$ due to the depletion charge is equal to

$$|\langle n|U^{\text{depl}}(q)|m\rangle|^2 = |U_{nm}^{\text{depl}}(\mathbf{q})|^2 = N_{\text{depl}} \left(\frac{e^2}{2\kappa q}\right)^2 A_{nm}^2(q) \int_0^\infty dz_i O_{nm}^2(q, z_i) \tag{181}$$

where $N_{\text{depl}}$ is the depletion charge density, $\mathbf{R}_i = (\mathbf{r}_i, z_i)$ is the location of an arbitrary charge center in the depletion region and $A_{nm}(q)$ and $O_{nm}(q, z_i)$ are the form factors due to the finite extension of the electron gas in the quantization direction of the form

$$A_{nm}(q) = \int_0^\infty dz \psi_n(z) e^{-qz} \psi_m(z) \tag{182}$$

and

$$O_{nm}(q, z_i) = 0.5\left(1 + \frac{\varepsilon_{ox}}{\varepsilon_{sc}}\right) e^{qz_i} + 0.5\left(1 - \frac{\varepsilon_{ox}}{\varepsilon_{sc}}\right) e^{-qz_i}$$

$$+ 0.5\left(1 + \frac{\varepsilon_{ox}}{\varepsilon_{sc}}\right)\left[e^{-qz_i}\frac{a_{nm}^{l+1}(q, z_i)}{A_{nm}(q)} - e^{qz_i}\frac{a_{nm}^{l}(q, z_i)}{A_{nm}(q)}\right] \tag{183}$$

respectively, where

$$a_{nm}^{(\pm)}(q, z_i) = \int_0^{z_i} dz\,\psi_n(z)e^{\pm q z}\psi_m(z) \tag{184}$$

In the above expressions, $\mathbf{q}$ is a wavevector in the plane parallel to the interface ($xy$-plane in our case).

At the $Si/SiO_2$ interface, there are always a large number of Coulomb centers near the interface, due to the disorder and defects in the crystalline structure in the neighborhood of the interface. They are associated with the dangling bonds and can lead to charge-trapping centers which scatter the free carriers through the Coulomb interaction. In MBE grown heterolayers, this interface state density is usually much lower, as there are relative few dangling bonds in such lattice matched structures. Using Eq. (179), the matrix element squared for the scattering from a sheet charge with charge density $N_{it}$ located in the oxide, at a distance $z_t (z_t < 0)$ from the interface is

$$|\langle n|U''(\mathbf{q})|m\rangle|^2 = |U_{nm}''(\mathbf{q})|^2 = N_{it}\left(\frac{e^2}{2\kappa}\right)^2\left[\frac{A_{nm}(q)}{q}\right]^2 e^{2qz_t} \tag{185}$$

For interface-trap scattering, $z_t = 0$. By similar arguments, one finds that the matrix element for scattering from the charges in the oxide or AlGaAs barrier region, with charge density $N_{ox}$, is given by

$$|\langle n|U^{ox}(\mathbf{q})|m\rangle|^2 = |U_{nm}^{ox}(\mathbf{q})|^2 = N_{ox}\left(\frac{e^2}{2\kappa}\right)^2\left[\frac{A_{nm}(q)}{q}\right]^2\frac{1 - e^{2qd_{ox}}}{2q} \tag{186}$$

where $d_{ox}$ is the oxide thickness. For a modulation doped heterostructure system, where an undoped spacer layer of thickness $d_{sp}$ exists, the lower limit of the integration over the impurity configuration is no longer zero, resulting in multiplication of Eq. (186) by a factor $e^{-qd_{sp}}$, which attenuates the scattering rate, resulting in higher mobility generally.

### 3.2.3.2. Surface-Roughness Scattering.
This scattering mechanism is associated with interfacial disorder due to the random variation of the position of the interface, again in either oxide-semiconductor or heterostructure systems. For a MOS system, the degree of interface roughness depends upon the oxidation temperature and ambient as well as postoxidation anneal and removal of the wafer from the furnace. In heterostructure systems, the degree of roughness is related to issues in epitaxial growth, such as the relative diffusivity of column V elements on the surface during growth. Early theories on surface-roughness were based on the Boltzmann equation, in which the surface is incorporated via boundary conditions on the electron distribution function [178–180]. The first quantum-mechanical treatment of the problem was given by Prange and Nee [181]. Subsequently, the theory followed two different paths.

The basic idea of the first approach is to incorporate the variations in the confining potential of the rough surface as a boundary condition on the Hamiltonian of the system. Since there is no simple perturbation theory to treat arbitrary changes in the boundary conditions, the problem of a free-electron Hamiltonian with complicated boundary conditions is then transformed by an appropriate coordinate transformation into a problem with simpler boundary conditions (i.e., into a problem where we have flat surfaces). This coordinate transformation technique has been proposed by Tesanovic et al. [182] and was later used by Trivedi and Ashcroft [183]. As a consequence of this transformation, the Hamiltonian of the system now has additional terms that play the role of potential interaction terms. These additional terms are treated by perturbative techniques, which are valid when the roughness of the surface is small compared to the thickness of the well.

In the second approach [104], the effect of the surface roughness is taken into account through a random local potential term

$$V_0\theta[-z + \Delta(\mathbf{r})] - V_0\theta(-z) \cong V_0\delta(z)\Delta(\mathbf{r}) \tag{187}$$

which is then treated perturbatively. The random function $\Delta(\mathbf{r})$ is a measure of the roughness and is most conveniently expressed in terms of the auto covariance function of $\Delta(\mathbf{r})$. The power spectrum $S(q)$ is the two-dimensional Fourier transform of the auto covariance function of $\Delta(\mathbf{r})$. For the Gaussian correlated roughness that is usually assumed [184–188], the power spectrum is given by

$$S_G(q) = \pi\Delta^2\zeta^2 \exp\left(-\frac{q^2\zeta^2}{4}\right) \qquad (188)$$

Parameters $\Delta$ and $\zeta$ characterize the r.m.s. height of the bumps on the surface and the roughness correlation length, respectively. Goodnick et al. [189] made extensive analysis of high-resolution transmission electron microscopy (HRTEM) measurements to test the assumption of Gaussian correlation. They found that exponential correlation describes roughness much better than the Gaussian correlation irrespective of growth conditions. Roughly speaking, it means that the interface may be regarded as consisting of terraces of a few nanometers in size separated by atomic steps of a few tenths of nanometers, as shown in Fig. 34. This result has recently been confirmed by Atomic Force Microscope (AFM) measurements [190]. The power spectrum for the exponential correlation is given by

$$S_E(q) = \frac{\pi\Delta^2\zeta^2}{(1 + q^2\zeta^2/2)^{3/2}} \qquad (189)$$

A generalization of the result given in Eq. (189) is a self-affine roughness correlation function, which in two-dimensions is of the form

$$S_{SA}(q) = \frac{\pi\Delta^2\zeta^2}{(1 + q^2\zeta^2/4n)^{n+1}} \qquad (190)$$

where, $n > 0$ is an exponent describing the high-$q$ fall-off of the distribution. It reduces to exponential correlation for $n = 0.5$.

For identical roughness parameters, the Gaussian spectrum decays slower for small wavevectors and then falls to zero rapidly for large wavevectors. The exponential model also leads to a rougher interface due to the tails in the spectrum, which allows for short-range fluctuations to be considered as well. For $n < 0.5$ and small values of q, the power spectrum of the self-affine model decays faster compared with the previous two models, but then falls slowly for large wavevectors. This essentially means that, in this regime, it also allows for short-range fluctuations to be considered. For large exponents, the power spectral density of the self-affine model approaches the one for the Gaussian model.

In general, the matrix element for scattering between subbands $n$ and $m$ for this scattering mechanism is of the form

$$|\langle n|U^{sr}(\mathbf{q})|m\rangle|^2 = S(q)\Gamma_{nm}^2(q) \qquad (191)$$



Figure 34. (a) High-resolution transmission electron micrograph of the interface between Si and $SiO_2$. The oxide is in the top half of the picture, while the rows of Si atoms can be observed in the bottom half. The image is a lattice plane image lying in the (111) plane, while the interface is a (100) plane. (b) Relevant dimensions for the steps occurring at the interface.

For large $V_0$, the matrix element $\Gamma_{nm}$ reduces to

$$\Gamma_{nm}^{(0)} = \frac{\hbar^2}{2m_z} \frac{d\psi_n}{dz} \frac{d\psi_m}{dz}\bigg|_{z=0}$$

$$= \int_0^\infty dz \left\{ \psi_n(z) \frac{\partial V(z)}{\partial z} \psi_m(z) - \varepsilon_m \frac{d\psi_n}{dz} \psi_m(z) + \varepsilon_n \psi_n(z) \frac{d\psi_m}{dz} \right\} \tag{192}$$

The last expression is the result obtained by Prange and Nee [181]. Matsumoto and Uemura [191] calculated that in the electronic quantum limit, $\Gamma_{nm} = eE_{av}$, where $E_{av} \propto (\frac{1}{2}N_s + N_{depl})$.

The change in the potential energy of the system due to surface-roughness was corrected by Ando [192], by considering the change in the electron density distribution and the effective dipole moment of the deformed Si-SiO$_2$ surface. The later scattering rate becomes

$$\Gamma_{nm}(q) = \Gamma_{nm}^{(0)} + \frac{e^2}{\varepsilon_{sc}} \frac{\varepsilon_{sc} - \varepsilon_{ox}}{\varepsilon_{sc} + \varepsilon_{ox}} A_{nm}(q) \left\{ (N_{depl} + N_s) - \frac{1}{2} \sum_l N_l A_{il}(q) \right\} \tag{193}$$

where $\Gamma_{nm}^{(0)}$ is given by Eq. (192). Because of the presence of the dielectric medium, one needs to correct the expression given in Eq. (193) with the contribution of the image term

$$\Gamma_{nm}^{image}(q) = \frac{e^2}{\varepsilon_{sc}} \frac{\varepsilon_{sc} - \varepsilon_{ox}}{\varepsilon_{sc} + \varepsilon_{ox}} \frac{q^2}{16\pi} \int_0^\infty dz \psi_n(z) \left[ \frac{K_1(qz)}{qz} - \frac{1}{2} \frac{\varepsilon_{sc} - \varepsilon_{ox}}{\varepsilon_{sc} + \varepsilon_{ox}} K_0(qz) \right] \psi_m(z) \tag{194}$$

In Eq. (194), $K_0$ and $K_1$ are the modified Bessel functions. An additional complications associated with the finite oxide thickness, which may further reduce the mobility via scattering with remote roughness, will be ignored in the present treatment.

**3.2.3.3. Electron-Phonon Interaction.** Phonon scattering can cause three different types of electronic transitions in the Si-inversion layer: transitions between states within a single valley via acoustic phonons (called intravalley acoustic-phonon scattering) and nonpolar optical phonons (called intravalley optical phonon scattering), and transitions between different valleys via nonpolar optical phonons (called intervalley scattering) [193–200]. The intravalley acoustic-phonon scattering involves phonons with low energies and is almost an elastic process. The intravalley optical-phonon scattering is induced by optical phonons of low momentum and high energy. The intervalley scattering can be induced by the emission and absorption of high-momentum, high-energy phonons, which can be of either acoustic- or optical-mode nature. Intervalley scattering can therefore be important only for temperatures high enough that an appreciable number of suitable phonons is excited or for hot electrons which can emit high energy phonons [201].

In order to evaluate the scattering potential that describes the electron-phonon interaction, a Hamiltonian is needed that describes the coupled electron-phonon system. The total Hamiltonian of the system may be written [202, 203]

$$\hat{H} = \hat{H}_e + \hat{H}_a + \hat{H}_{ea} \tag{195}$$

where $\hat{H}_e$ is the electronic part, $\hat{H}_a$ is the atomic part of the Hamiltonian that describes the normal modes of vibration of the solid and $\hat{H}_{ea}$ is the electron-ion interaction term of the form

$$\hat{H}_{ea} = \sum_{i,j} V_{ea}(\mathbf{r}_i - \mathbf{R}_j) \tag{196}$$

In general, each ion is at a position $\mathbf{R}_j = \mathbf{R}_j^{(0)} + \mathbf{Q}_j$, which is the sum of the equilibrium position $\mathbf{R}_j^{(0)}$ and the displacement $\mathbf{Q}_j$. Under the assumption of small displacements, one can expand $V_{ea}$ in a Taylor series

$$V_{ea}(\mathbf{r} - \mathbf{R}_j) = V_{ea}(\mathbf{r} - \mathbf{R}_j^{(0)}) - \mathbf{Q}_j \cdot \nabla V_{ea}(\mathbf{r} - \mathbf{R}_j^{(0)}) + O(Q^2) \tag{197}$$

The zero-order term is the potential function for the electrons when the atoms are in their equilibrium positions, which forms a periodic potential in the crystal. The solution of the

Hamiltonian for electron motion in this periodic potential gives the Bloch states of the solids. Since the first-order term is much smaller than the zero-order term, the electron-phonon interaction can be treated perturbatively. Therefore, the lowest order term for the electron-phonon interaction is of the form

$$H_{e-ph}(\mathbf{r}) = -\sum_{j} \mathbf{Q}_j \cdot \nabla_i V_{ea}(\mathbf{r} - \mathbf{R}_j^{(0)}) \tag{198}$$

It is obvious that this interaction Hamiltonian does not act on the spin variables in this approximation. Under the assumption that the electron-atom potential possesses a Fourier transform, one can write

$$V_{ea}(\mathbf{r}) = \frac{1}{N} \sum_{\mathbf{q}} V_{ea}(\mathbf{q}) e^{i\mathbf{q}\cdot\mathbf{r}} \tag{199}$$

where $N$ is a number of primitive cells, and wave vector $\mathbf{q}$ spans the whole $q$-space. Ionic displacement may be decomposed into normal-mode representation and it is customary to write

$$\mathbf{Q}_j = i \sum_{\mathbf{k}} \left( \frac{\hbar}{2MN\omega_{\mathbf{k}\lambda}} \right)^{1/2} \left( \hat{a}_{\mathbf{k}\lambda} \mathbf{e}_{\mathbf{k}\lambda} e^{i\mathbf{k}\cdot\mathbf{R}_j^{(0)}} + \hat{a}_{\mathbf{k}\lambda}^{+} \mathbf{e}_{\mathbf{k}\lambda}^{+} e^{i\mathbf{k}\cdot\mathbf{R}_j^{(0)}} \right) \tag{200}$$

where $\mathbf{e}_{\mathbf{k}\lambda}$ is the unit polarization vector that obeys the standard orthonormality and completeness relations. $\omega_{\mathbf{k}\lambda}$ is the phonon frequency for wavevector $\mathbf{k}$ running over the whole Brillouin zone of the phonon branch $\lambda$. $\hat{a}_{\mathbf{k}\lambda}(\hat{a}_{\mathbf{k}\lambda}^{+})$ are the phonon annihilation (creation) operators. In acoustic waves, $\mathbf{Q}_j$ refers to the relative displacement of the unit cell; in optical waves it refers to the relative displacement of the two atoms in the unit cell. Thus

$$H_{e-ph}(\mathbf{r}) = \sum_{\mathbf{q},\mathbf{G}} e^{i(\mathbf{q}+\mathbf{G})\cdot\mathbf{r}} V_{ea}(\mathbf{q}+\mathbf{G})(\mathbf{q}+\mathbf{G}) \cdot \mathbf{e}_{\mathbf{q}\lambda} \left( \frac{\hbar}{2\rho V \omega_{\mathbf{q}\lambda}} \right) (\hat{a}_{\mathbf{q}\lambda} + \hat{a}_{\mathbf{q}\lambda}^{+}) \tag{201}$$

where $MN = \rho V$, and $\rho$ is the density of the solid. The summation over $\mathbf{G}$ represents summation over all reciprocal lattice vectors of the solid. If one defines a function

$$M_{\mathbf{q}\lambda} = \left( \frac{\hbar}{2\rho V \omega_{\mathbf{q}\lambda}} \right) \sum_{\mathbf{G}} e^{i\mathbf{G}\cdot\mathbf{r}} (\mathbf{q}+\mathbf{G}) \cdot \mathbf{e}_{\mathbf{q}\lambda} V_{ea}(\mathbf{q}+\mathbf{G}) \tag{202}$$

then the Hamiltonian for the electron-phonon interaction becomes

$$H_{e-ph}(\mathbf{r}) = \sum_{\mathbf{q}} M_{\mathbf{q}\lambda} e^{i\mathbf{q}\cdot\mathbf{r}} (\hat{a}_{\mathbf{q}\lambda} + \hat{a}_{\mathbf{q}\lambda}^{+}) \tag{203}$$

The exact form of the matrix elements for acoustic and nonpolar-optical phonon scattering (zero- and first-order terms) are given below. Since distinction between 3D and 2D vectors needs to remain clear, in the following, we use the following notation: capital bold letters will refer to three-dimensional vectors, whereas small bold letters will be used for two-dimensional wavevectors that lie in the $xy$-plane.

*Deformation potential scattering.* In general, the application of mechanical stress alters the band structure by shifting energies, and, where it destroys symmetry, by removing degeneracies. It is usually assumed that the mechanical stress does not change the band curvature, and therefore does not change the effective masses, but introduces shift in the energy states that are close to the band extremum [204–207].

For isotropic elastic continuum, matrix element for deformation potential scattering (acoustic phonons) can be obtained by taking the long-wavelength limit of the Eq. (202) [208]. For small values of $\mathbf{q}$, the summation over reciprocal lattice vectors can be neglected, except for the term $\mathbf{G} = 0$. The screened electron-ion interaction becomes a constant which is usually denoted as $\Xi$ (it gives the shift of the band edge per unit elastic strain). Under these assumptions, $M_{\mathbf{Q}\lambda}$ simplifies to

$$M_{\mathbf{Q}\lambda} = i\mathbf{Q} \cdot \mathbf{e}_{\mathbf{Q}\lambda} \left( \frac{\hbar \Xi^2}{2\rho V \omega_{\mathbf{Q}\lambda}} \right)^{1/2} \tag{204}$$

where $\omega_{Q\lambda} = v_{s\lambda}Q$, where $v_{s\lambda}$ is the sound velocity, $\omega_{Q\lambda}$ is the phonon frequency. Long wavelength acoustic phonons (LA mode) have $Q||e_{Q\lambda}$ which makes the matrix element nonzero. TA phonons have $Q\perp e_{Q\lambda}$ which makes the matrix element vanish. Therefore, the deformation potential mainly couples electrons to LA phonons.

For anisotropic elastic continuum such as silicon, the deformation potential constant $\Xi$ becomes a tensor. The anisotropy of the intravalley deformation potential in the ellipsoidal valleys in silicon has been extensively studied by Herring and Vogt [207]. Expanding the electron-phonon matrix element over spherical harmonics and retaining only the leading terms, they have expressed the anisotropy of the interaction in terms of the angle $\theta_Q$ between the wavevector $Q$ of the emitted (absorbed) phonon and the longitudinal axis of the valley. They have shown that the matrix element is proportional to $Q$ via the deformation potential $\Delta_\lambda(\theta_Q)$ ($\lambda$ = LA or TA) given by [209–211]

$$\Delta_{LA}(\theta_Q) \approx \Xi_d + \Xi_u \cos^2(\theta_Q) \tag{205}$$

and

$$\Delta_{TA}(\theta_Q) \approx \Xi_u \cos(\theta_Q)\sin(\theta_Q) \tag{206}$$

Equation (206) accounts for the contribution of both TA branches. Therefore, the acoustic mode scattering is characterized by two constants: $\Xi_u$ (uniaxial shear potential) and $\Xi_d$ (dilatation potential) that is believed to have values of approximately 9.0 eV and −11.7 eV, respectively (See Fig. 35).

In bulk silicon, this anisotropy is usually ignored by using an effective deformation potential constant $\Xi_{LA}^{eff}$ for the interaction with longitudinal modes, and ignoring the role of the lower-energy TA modes. This approximation can be justified due to the following reasons. The acoustic modes are most effective at low energy. In this regime and in the usual elastic and equipartition approximation that will be described later, due to the linear dependence on $Q$, scattering of electrons at some energy $\varepsilon$ samples almost uniformly the equal energy ellipsoid, so that one can take the average values of $\Delta_\lambda$ over the ellipsoid. Since there is nothing to fix the energy scale in the problem, this averaging procedure is independent of the electron energy. Moving to the two-dimensional situation, one cannot follow a parallel path to arrive at an isotropic, energy independent effective deformation potential, which complicates the treatment of this scattering process.

Since the wave functions of the initial and final states are usually expressed as a product of a one-electron wave functions (Bloch functions) and harmonic oscillator wave functions, after the averaging over the phonon states is performed, the terms inside the brackets of Eq. (203) that represent phonon absorption (term $\hat{a}_{q\lambda}$) and phonon emission (term $\hat{a}_{q\lambda}^+$) processes reduce to $\sqrt{N_{q\lambda}}$ and $\sqrt{N_{q\lambda} + 1}$, respectively. In thermal equilibrium with a lattice at temperature $T$, the phonon occupation number $N_{q\lambda}$ is given by Bose-Einstein statistics

$$N_{q\lambda} = \frac{1}{e^{\hbar\omega_{q\lambda}/k_B T} - 1} \tag{207}$$



Figure 35. Angular dependence of the deformation potential for longitudinal modes.

where $k_B$ is the Boltzmann constant. At high enough temperatures, the acoustic phonon energies are much smaller than the thermal energies of the electrons. Therefore, one can expand the exponent in the denominator of Eq. (207) into a series and, in the *equipartition approximation*, approximately obtain

$$N_{q\lambda} + 1 \approx N_{q\lambda} \approx \frac{k_B T}{\hbar \omega_{q\lambda}} \gg 1 \tag{208}$$

Incorporating these terms as well as the exponential term $e^{iq_z z}$ into the definition of $M_{q\lambda}$, after a straightforward calculation one finds that the matrix element squared for scattering between subbands $n$ and $m$ due to acoustic phonons for both, absorption plus emission process, after the averaging over $q_z$ is performed, reduces to

$$|\langle n|U_\lambda^{ao}(\mathbf{q})|m\rangle|^2 = \frac{k_B T}{\rho V v_{s\lambda}^2} [\Delta_{\lambda, nm}^{\text{eff}}(\mathbf{q})]^2 F_{nm} \tag{209}$$

where

$$F_{nm} = \int_0^\infty dz \psi_n^2(z) \psi_m^2(z) \tag{210}$$

The effective deformation potential constant is calculated from

$$[\Delta_{\lambda, nm}^{\text{eff}}(\mathbf{q})]^2 = \frac{1}{F_{nm}} \int_0^\infty dq_z \Delta_\lambda^2(\theta_Q) |\mathcal{F}_{nm}(q_z)|^2 \tag{211}$$

where

$$\mathcal{F}_{nm}(q_z) = \int_0^\infty dz \psi_n(z) e^{iq_z z} \psi_m(z) \tag{212}$$

The main result that can be deduced from the expression above is that the form-factor, $\mathcal{F}_{nm}(q_z)$, introduces an energy scale in the problem by fixing the *fuzzy* component, the wavevector $q_z$. This result is an expected one and follows immediately from the uncertainty principle $\Delta_z \Delta_{p_z} \geq \hbar$. Since the electrons are *frozen* into their wave functions, and cannot oscillate in the quantized direction, the uncertainty in the particles location along the $z$-axis has been reduced. Therefore, there must be a corresponding increase in the uncertainty in the particles $z$-directed momentum [212].

*Nonpolar optical phonon scattering.* The scattering of electrons by zone-center optical and intervalley phonons in semiconductor crystals has been treated rather extensively by Ferry [201, 213, 214]. The nonpolar optical interaction is important for intra-subband scattering as well as for scattering of electrons (and holes) between different minima of the conduction (or valence) band. This later interaction is important for scattering of carriers in semiconductors with many-valley band structure, as it is usually the case in Si and Ge, and in the Gunn effect, where scattering occurs between different sets of equivalent minima. Harrison [206] pointed out that the nonpolar optical matrix element may be either of zero or higher order in the phonon wavevector. In subsequent treatments of electron transport in which the nonpolar interaction is important, only the zero-order term was considered, generally owing to the impression that the higher order terms are much smaller. Although this is usually the case, there arise many cases in which the zero-order term is forbidden by the symmetry of the state involved. In these cases, the first order term becomes the leading term, and can become significant in many instances. For example, the first-order intervalley scattering plays an important role in hot-electron transport in the $n$-type inversion layer in Si. Ignoring this scattering process means that there will be no saturation of the drift velocity at high electric fields, because the zero-order intervalley scattering rate is weakly dependent on the electron energy of high-energy electrons, while the first-order intervalley scattering rate increases as the electron energy increases.

The matrix element for nonpolar optical phonon scattering is generally found from a deformable ion model explained in the introduction part of this section. If one thinks of an optical phonon as occurring at finite $\mathbf{G}$, then the $\mathbf{q}$ dependence is unimportant, so that the

entire matrix element becomes constant, resulting in

$$M_{0\lambda} = M_{n\lambda} = \left( \frac{\hbar D_\lambda^2}{2\rho V \omega_{o\lambda}} \right)^{1/2} \tag{213}$$

where $D_\lambda$ is the deformation field (usually given in eV/cm) and $\omega_{o\lambda}$ is the frequency of the relevant phonon mode which is usually taken to be independent of the phonon wave vector for optical and intervalley processes. Fourier transforming back to real space, a constant in $q$-space produces a delta-function in real space. Therefore, this zero-order term represents a short-ranged interaction. A local dilatation or compression of the lattice produces a local fluctuation in the energy of the electron or hole. Incorporating the exponential term $e^{iq_z z}$ into the definition of $M_{n\lambda}$, after averaging over $q_z$, the following result is obtained

$$|\langle n|U_\lambda^{opt\,0}|m\rangle|^2 = \frac{\hbar D_\lambda^2}{2\rho V \omega_{o\lambda}} F_{nm} \tag{214}$$

for the squared matrix element for scattering between subbands $n$ and $m$ that belong to the $\alpha$ and $\beta$ valley, respectively.

When the zero-order matrix element for the optical or intervalley interaction vanishes, then $D_\lambda$ is identically zero. In this case, one has to consider the first-order term of the interaction whose matrix element is

$$M_{Q\lambda} = i\mathbf{Q} \cdot \mathbf{e}_{Q\lambda} \left( \frac{\hbar D_{1\lambda}^2}{2\rho V \omega_{n\lambda}} \right)^{1/2} \tag{215}$$

In this context, a first-order process means a process similar to acoustic phonon scattering. Following the previously explained procedure, we find that the matrix element squared for scattering between subbands $n$ ($\alpha$-valley) and $m$ ($\beta$-valley) is given by

$$|\langle n|U_\lambda^{opt\,1}|m\rangle|^2 = \frac{\hbar D_{1\lambda}^2}{2\rho V \omega_{o\lambda}} (q^2 F_{nm} + c_{nm}) \tag{216}$$

where

$$c_{nm} = \int_{-0}^{\infty} dz \left\{ \frac{d}{dz} [\psi_n(z)\psi_m(z)] \right\}^2 \tag{217}$$

The constant term $c_{nm}$ is a small correction term that we have found [172].

In the scattering among the equivalent valleys, there are two types of phonons that might be involved in the process (see Fig. 36). The first type, so-called g-phonon couples the two



Figure 36. Diagrammatic representation of intervalley transitions due to g- and f-phonons.

valleys along opposite ends of the same axis (i.e., $\langle 100 \rangle$ to $\langle \bar{1}00 \rangle$). This is an Umklapp process and has a net phonon wave vector $0.3\pi/a$. The $f$-phonons couple the $\langle 100 \rangle$ valley with $\langle 010 \rangle$, $\langle 001 \rangle$, and so on. The reciprocal lattice vector involved in the $g$-process is $G_{100}$ and that for an $f$-process is $G_{111}$. Degeneracy factors ($g_r$) for transition between unprimed ($\alpha = 1$) and primed ($\alpha = 2$) set of subbands, for both $g$- ($r = 1$) and $f$-phonons ($r = 2$) are summarized in Table 7.

Within a three-subband approximation, scattering between the two valleys in the $\varepsilon_0$ and $\varepsilon_1$ subbands involves only $g$-type phonons. The scattering between these two minima is usually treated by using a high-energy phonon of 750-K activation temperature (treated as zero-order interaction) and 134-K phonon treated via first order interaction. Scattering between the two $\varepsilon_0$, $\varepsilon_1$ and the four $\varepsilon'_0$ subbands involves $f$-phonons with activation temperatures of 630 K and 230 K treated via zero-order and first-order interaction, respectively. Scattering between the subbands $\varepsilon'_0$ involves both $g$- and $f$-phonons with activation temperatures of 630-K (zero-order interaction) and 190-K (first-order interaction). All of the high-energy phonons are assumed to be coupled with a value of $D_\lambda = 9 \times 10^8$ [eV/cm] and all of the first-order coupled phonons are assumed to be coupled with $D_{1\lambda} = 5.6$ [eV] (This value is consistent with the results given in Ref. [215].) The first Born approximation result for the total electron-bulk phonon scattering rate for $p$-type silicon with $N_a = 1 \times 10^{15}$ cm$^{-3}$, $N_s = 1 \times 10^{12}$ cm$^{-2}$ and $T = 300$ K, with (thick line) and without (thin line) the inclusion of the correction term for the first order process, for the lowest subband of the unprimed ladder of subbands, is given in Fig. 37. We see that, throughout the whole energy range, there is an increase of approximately 10% of the total electron-bulk-phonon scattering rate due to the correction term introduced previously that could lead to mobility reduction. The same trend was also observed for the higher-lying subbands.

### 3.2.4. Screening of Coulomb and Surface-Roughness Scattering

It is well known that Coulomb and surface-roughness scattering significantly affect the electron mobility in Si inversion layers, particularly at low- and high-inversion charge densities. As mentioned earlier, the scattering potentials for these two processes are strongly affected by screening of the mobile charges in the inversion layer. Therefore, any theory that tries to explain the density and temperature dependence of the electron mobility must account for these screening corrections. Since the calculation of the exact dielectric function of homogeneous electron gas is a formidable problem, various approximate solutions for the dielectric function exist in the literature [216–218]. Some of these have been very successful, perhaps because they are simple (Thomas-Fermi method) or perhaps because they are accurate (mean-field approximation, also known as random-phase approximation [RPA]). The Thomas-Fermi method is basically the semiclassical limit of the Hartree calculation. On the other side, RPA is an exact Hartree calculation of the charge density in the presence of the self-consistent field of the external charge plus electron gas. More precisely, in the mean-field approximation, one includes only the long-range Coulomb interaction in the dielectric response, leaving out all exchange-correlation corrections. It leads to the so-called Lindhard dielectric function that is extensively employed in the literature [219–226].

The Thomas-Fermi theory of screening is described first. The quantum theory of plasma screening is discussed afterwards. Application of this theory for the calculation of the screened matrix elements for Coulomb and surface-roughness scattering is given at the end of this section.

*Thomas-Fermi theory.* Suppose that a positive test-charge $\rho_{ext}(\mathbf{r}) = en_{ext}(\mathbf{r})$ is placed at a given position $\mathbf{r}$ in a 3D-electron gas. The test-charge attracts the electrons, inducing a disturbance $\delta\rho(\mathbf{r}) = -e\delta n(\mathbf{r})$ in the charge-density distribution within the plasma. The net

**Table 7.** Degeneracy factors for transition between unprimed and primed subbands, for both $g$- and $f$-phonons.

|  | $\alpha = 1$ | $\alpha = 2$ |
|---|---|---|
| $\alpha = 1$ | $g_1 = 1; g_2 = 0$ | $g_1 = 0; g_2 = 4$ |
| $\alpha = 2$ | $g_1 = 0; g_2 = 2$ | $g_1 = 1; g_2 = 2$ |

**Figure 37.** Total electron-bulk phonon scattering rate for the electrons in the lowest unprimed subband calculated within the first Born approximation. Reprinted with permission from D. Vasileska, Ph.D. Thesis, Arizona State University, 1995. © 1995.

result of the movement of the entire ensemble of charges is a potential that does not behave as a simple Coulomb potential due to the initial charge. In essence, this is a self-consistent process, in which the charges produce a potential, which modifies the total charge, which modifies the total potential, and so on. The modification of the strength of the Coulomb potential of a single charge in the presence of electron plasma is called *screening effect*.

In treating screening, it is convenient to define two electrostatic potentials. The first one, $\varphi_{ext}(\mathbf{r})$, arises solely from the test-charge itself, and satisfies the following Poisson equation

$$-\nabla^2 \varphi_{ext}(\mathbf{r}) = \frac{\rho_{ext}(\mathbf{r})}{\varepsilon_a} = \frac{e}{\varepsilon_a} n_{ext}(\mathbf{r}) \qquad (218)$$

The second one, the so-called effective electrostatic potential, $\varphi_{eff}(\mathbf{r})$, arises from both the positively charged test-charge and the induced charge-density. It therefore satisfies

$$-\nabla^2 \varphi_{eff}(\mathbf{r}) = \frac{\rho_{ext}(\mathbf{r}) + \delta\rho(\mathbf{r})}{\varepsilon_a} = \frac{e}{\varepsilon_a}[n_{ext}(\mathbf{r}) - \delta n(\mathbf{r})] \qquad (219)$$

Introducing $V_{eff/ext}(\mathbf{r}) = e\varphi_{eff/ext}(\mathbf{r})$ and Fourier transforming Eqs. (218) and (219) yields

$$q^2 V_{ext}(\mathbf{q}) = \frac{e^2}{\varepsilon_a} n_{ext}(\mathbf{q}) \qquad (220)$$

and

$$q^2 V_{ext}(\mathbf{q}) = \frac{e^2}{\varepsilon_a}[n_{ext}(\mathbf{q}) - \delta n(\mathbf{q})] \qquad (221)$$

respectively, where $\mathbf{q}$ is a 3D-wave vector and $\delta n(\mathbf{q})$ is the Fourier transform of the induced electron density. Using the results given in Eqs. (220) and (221), the expression for $V_{eff}(\mathbf{q})$ simplifies to

$$V_{eff}(\mathbf{q}) = \frac{V_{ext}(\mathbf{q})}{1 + (e^2/\varepsilon_o q^2)[\delta n(\mathbf{q})/V_{eff}(\mathbf{q})]} = \frac{V_{ext}(\mathbf{q})}{1 + V_c(\mathbf{q})[\delta n(\mathbf{q})/V_{eff}(\mathbf{q})]} \qquad (222)$$

where $V_c(\mathbf{q}) = e^2/\varepsilon_o q^2$ is a Fourier transform of the Coulomb interaction for a 3D-system. From the theory of dielectric media, it follows that the effective electrostatic potential and the external one are linearly related to each other

$$\varphi_{ext}(\mathbf{r}) = \int d^3r' \chi(\mathbf{r} = \mathbf{r}')\varphi_{eff}(\mathbf{r}') \qquad (223)$$

which implies that the corresponding Fourier transforms satisfy

$$\varphi_{\text{ext}}(\mathbf{q}) = \chi(\mathbf{q})\varphi_{\text{eff}}(\mathbf{q}) \tag{224}$$

where $\chi(\mathbf{q})$ is the wave vector–dependent relative dielectric constant of the medium. From the above results, one has that

$$\chi(\mathbf{q}) = 1 + V_c(\mathbf{q})\left[\frac{\delta n(\mathbf{q})}{V_{\text{eff}}(\mathbf{q})}\right] \tag{225}$$

If the effective potential energy, $V_{\text{eff}}(\mathbf{r})$, varies slowly on the length scale of the Fermi wavelength then, within the Hartree approximation, the electron energy is modified from its free-electron value by the total local potential, i.e.,

$$\varepsilon(\mathbf{k}) = \frac{\hbar^2 k^2}{2m^*} - V_{\text{eff}}(\mathbf{r}) = \varepsilon_k - V_{\text{eff}}(\mathbf{r}) \tag{226}$$

For an electron plasma in thermal equilibrium, the electron-density distribution is represented by the Fermi-Dirac distribution function, so that the plasma density is calculated from [227]

$$n(\mathbf{r}) = s \iint \frac{d^3k}{(2\pi)^3} \frac{1}{1 + \exp[(\varepsilon_k - V_{\text{eff}}(\mathbf{r}) - \mu)/k_B T]} \tag{227}$$

If the change induced by the potential, $V_{\text{eff}}(\mathbf{r})$, is small, the induced electron density can be approximated as

$$\delta n(\mathbf{r}) \approx \frac{\partial n}{\partial \mu} V_{\text{eff}}(\mathbf{r}) \tag{228}$$

Inserting the Fourier transformed result yields

$$\chi(\mathbf{q}) = 1 + V_c(q)\frac{\partial n}{\partial \mu} = 1 + \left(\frac{q_s}{q}\right)^2 \tag{229}$$

where

$$q_s = \sqrt{\frac{e^2}{\varepsilon_o}\frac{\partial n}{\partial \mu}} \tag{230}$$

is the so-called Thomas-Fermi screening wave vector. For a nondegenerate Boltzmann distribution, $\partial n/\partial \mu = n/k_B T$. Inserting this result into Eq. (230) leads to the Debye-Hückel screening wavevector for a 3D system.

From the above results, one immediately finds that for a Q2D-system, the relative dielectric function is given by

$$\chi(\mathbf{q}) = 1 + \frac{1}{q}\sum_i q_{si} \tag{231}$$

where $\mathbf{q}$ is a 2D-wavevector in the plane parallel to the interface and

$$q_{si} = g_i \frac{m^* e^2}{2\pi \varepsilon_o \hbar^2}\left[1 - \exp\left(-\frac{\pi \hbar^2 N_i}{g_i m^* k_B T}\right)\right] \tag{232}$$

is the screening wave vector of the $i$th subband. In Eq. (232), $g_i$ is the valley degeneracy factor. It is interesting to note that, in contrast to the 3D case (Eq. [230]), the screening wavevector for a Q2D-system becomes density-independent at very low temperatures and/or high-enough inversion charge densities.

*Mean-field approximation.* In this section we present the derivation of the Lindhard dielectric function for a 3D system from the equations of motion for the density operator

$$\hat{n}(\mathbf{q}, t) = \sum_k \hat{c}^+_{k-q}\hat{c}_k \tag{233}$$

which is a Fourier transform of the real-space operator $\hat{n}(\mathbf{r}, t) = \hat{\psi}^+(\mathbf{r}, t)\hat{\psi}(\mathbf{r}, t)$. In the self-consistent Hartree approximation, the effective single particle Hamiltonian is

$$\hat{H} = \int d^3r\hat{\psi}^+(\mathbf{r}, t)\left(-\frac{\hbar^2}{2m^*}\nabla^2\right)\hat{\psi}(\mathbf{r}, t) + \int d^3r V_{\text{eff}}(\mathbf{r}, t)\hat{\psi}^+(\mathbf{r}, t)\hat{\psi}(\mathbf{r}, t)$$

$$= \sum_k \varepsilon_k \hat{c}_k^+ \hat{c}_k + \sum_{k, q} V_{\text{eff}}(\mathbf{q}, t)\hat{c}_{k+q}^+\hat{c}_k \tag{234}$$

where the effective potential $V_{\text{eff}}(\mathbf{r}, t)$ equals to the sum of the Coulomb potential $V_c(\mathbf{r})$ of a test-charge and the induced potential $V_s(\mathbf{r}, t)$ of the screening particles. The induced potential satisfies the Poisson equation

$$\nabla^2 V_s(\mathbf{r}, t) = -\frac{e^2}{\varepsilon_0}\delta n(\mathbf{r}, t) \tag{235}$$

where $\delta n(\mathbf{r}, t)$ is the deviation of the electron density from its equilibrium value $n_0$. To calculate the induced potential, it is necessary to derive the expression for the screening particle density $\delta n(\mathbf{r}, t)$. This is obtained by writing an equation of motion for the density operator given in Eq. (233), and then solving it approximately. In the homogeneous electron gas, the particle-density operator has an expectation value zero, unless there is a perturbation in the system. A perturbation in the system polarizes the plasma, so that the average value $\delta n(\mathbf{q}, t) = \langle \hat{n}(\mathbf{q}, t)\rangle$ becomes finite. In the linear screening model, one assumes that this average is proportional to the effective potential, that is, $\langle \hat{n}(\mathbf{q}, \omega)\rangle \propto V_{\text{eff}}(\mathbf{q}, \omega)$, where $\omega$ is the frequency. Therefore, the goal of the derivation presented below is to determine the form of this constant of proportionality.

The density operator satisfies the Liouville equation

$$i\hbar\frac{\partial\hat{n}(\mathbf{q}, t)}{\partial t} = [\hat{H}, \hat{n}(\mathbf{q}, t)] \tag{236}$$

In deriving the Lindhard dielectric function, it is more convenient to evaluate the equation of motion of one component of the density operator, i.e.,

$$i\hbar\frac{\partial}{\partial t}\hat{c}_{k-q}^+\hat{c}_k = [\hat{H}, \hat{c}_{k-q}^+\hat{c}_k]$$

$$= (\varepsilon_{k-q} - \varepsilon_k)\hat{c}_{k-q}^+\hat{c}_k + \sum_p V_{\text{eff}}(\mathbf{p}, t)(\hat{c}_{p+k-q}^+\hat{c}_k^+ - \hat{c}_k^+\hat{c}_{k-p}) \tag{237}$$

If the Coulomb potential, $V_c(\mathbf{r})$, of a test-charge is assumed to oscillate at a single frequency $\exp[-i(\omega + i\delta)]$, the time derivative on the left-hand side of Eq. (237) is $(\hbar\omega + i\delta)\hat{c}_{k-q}^+\hat{c}_k$. The term $\omega + i\delta(\delta \to 0)$ establishes an adiabatic switch-on of the test-charge potential. With this assumption, Eq. (237) becomes

$$(\hbar\omega + \varepsilon_k - \varepsilon_{k-q} + i\delta)\hat{c}_{k-q}^+\hat{c}_k = \sum_p V_{\text{eff}}(\mathbf{p}, \omega)(\hat{c}_{p+k-q}^+\hat{c}_k - \hat{c}_k^+\hat{c}_{k-p}) \tag{238}$$

Within the mean-field approximation, the summation on the right-hand side of Eq. (238) is approximated by keeping only the term which has $\mathbf{p} = \mathbf{q}$ and neglecting all other terms. (It is assumed that the terms with other values of $\mathbf{p}$ average out to zero.) In other words, one keeps only those terms on the right-hand side of Eq. (238) that represent the local density. These terms are then replaced with their expectation values $f_k = \langle \hat{c}_k^+\hat{c}_k\rangle$ and $f_{k-q} = \langle \hat{c}_{k-q}^+\hat{c}_{k-q}\rangle$, where $f_k$ is the Fermi-Dirac distribution. Under these circumstances, the approximate equation can be easily solved, to give

$$\langle \hat{c}_{k-q}^+\hat{c}_k\rangle \cong V_{\text{eff}}(\mathbf{q}, \omega)\frac{f_k - f_{k-q}}{\hbar\omega + \varepsilon_k - \varepsilon_{k-q} + i\delta} \tag{239}$$

Performing the summation over $\mathbf{k}$, one finally arrives at

$$\delta n(\mathbf{q}, \omega) = \sum_k \langle \hat{c}_{k-q}^+\hat{c}_k\rangle = V_{\text{eff}}(\mathbf{q}, \omega)P^{(0)}(\mathbf{q}, \omega) \tag{240}$$

where

$$P^{(0)}(\mathbf{q}, \omega) = \sum_k \frac{f_k - f_{k-q}}{\hbar\omega + \varepsilon_k - \varepsilon_{k-q} + i\delta} \tag{241}$$

is the so-called bare polarizability function. Equations (235) and (240) may now be solved to obtain the dielectric function. In the Fourier transformed equation (235), one substitutes first the result given in Eq. (240) and then solves for the effective potential. This gives

$$V_{\text{eff}}(\mathbf{q}, \omega) = \frac{V_c(\mathbf{q})}{1 - V_c(\mathbf{q})P^{(0)}(\mathbf{q}, \omega)} \tag{242}$$

Unlike the bare interaction, $V_c(\mathbf{q})$, the effective interaction, $V_{\text{eff}}(\mathbf{q}, \omega)$, is frequency dependent. If $V_{\text{eff}}(\mathbf{q}, \omega)$ is Fourier transformed to $(\mathbf{q}, t)$-space, it will thus be a time-dependent interaction. This is due to the inertia of the polarization charge. In other words, it takes a finite amount of time for the electrons to come to their screening positions, which makes the whole system behave dynamically.

The ratio of the Coulomb potential of the test charge and the effective potential in Eq. (242) is just the RPA dielectric function

$$\chi(\mathbf{q}, \omega) = 1 - V_c(\mathbf{q})P^{(0)}(\mathbf{q}, \omega) \tag{243}$$

which was *the* dielectric function in the early days of the electron gas theory, since it is rather easy to derive and it also predicts correctly a number of properties of the electron gas such as *plasmons*. The results given in Eqs. (241) and (243) complete the derivation of the Lindhard dielectric function from the method of self-consistent fields.

At this point, it is interesting to show that in the static ($\omega = 0$) and long-wavelength ($\mathbf{q} \to 0$) limit, the Lindhard dielectric function reduces to the classical Thomas-Fermi result. In this limit, the denominator in Eq. (241) is

$$\varepsilon_k - \varepsilon_{k-q} \approx \mathbf{q} \cdot \frac{\partial \varepsilon_k}{\partial \mathbf{k}} \tag{244}$$

and the numerator can be approximated with

$$f_k - f_{k-q} \approx \mathbf{q} \cdot \frac{\partial f_k}{\partial \mathbf{k}} = \left( \mathbf{q} \cdot \frac{\partial \varepsilon_k}{\partial \mathbf{k}} \right) \frac{\partial f_k}{\partial \varepsilon_k} \tag{245}$$

In this case, the bare-polarizability function simplifies to

$$P^{(0)}(\mathbf{q}, 0) \approx \sum_k \frac{\partial f_k}{\partial \varepsilon_k} \tag{246}$$

Since $\partial f_k / \partial \varepsilon_k = -\partial f_k / \partial \mu$ and

$$n = \sum_k f_k \tag{247}$$

where $n$ is the electron density, we find

$$\chi(\mathbf{q}, 0) = 1 + V_c(\mathbf{q}) \frac{\partial n}{\partial \mu} \tag{248}$$

which is identical to the classical result given in Eq. (231).

*Calculation of screened matrix elements.* For external potentials which are local in space (such as those that represent impurity and surface-roughness scattering), screened matrix elements are calculated using

$$W_{ij}^{\text{eff}}(\mathbf{q}, \omega) = V_{ij}^{\text{bare}}(\mathbf{q}) + \frac{1}{q} \sum_{m,n} F_{ij,nm}(\mathbf{q}) q_{nm}'(\mathbf{q}, \omega) W_{nm}^{\text{eff}}(\mathbf{q}, \omega) \tag{249}$$

where $V_{ij}^{bare}(\mathbf{q})$ is the matrix element of the unscreened scattering potential for scattering between subbands $i$ and $j$. and $W_{ij}^{eff}(\mathbf{q}, \omega)$ represents the matrix element of the effective or screened scattering potential and $F_{ij,nm}$ is the form factor. For small values of $\mathbf{q}$, the intersubband form-factors (terms with $i \neq j$ and/or $n \neq m$) are at least two-orders of magnitude smaller than the intrasubband form-factors (terms with $i = j$ and $n = m$) due to the orthonormality of the subband basis. Therefore, it is reasonable to assume

$$F_{ii,nm}(\mathbf{q}) = F_{ij,nm}(\mathbf{q})\delta_{ij}\delta_{nm} \tag{250}$$

so that the diagonal terms of the screened matrix elements can be calculated from

$$W_{ii}^{eff}(\mathbf{q}, \omega) = V_{ii}^{bare}(\mathbf{q}) + \frac{1}{q}\sum_{n} F_{ii,nm}(\mathbf{q})q_{nm}^{s}(\mathbf{q}, \omega)W_{nm}^{eff}(\mathbf{q}, \omega) \tag{251}$$

It is obvious that, adopting the diagonal approximation given in Eq. (250), we have reduced considerably the size of the matrix that needs to be inverted in order to calculate the screened matrix elements. The off-diagonal elements are then calculated from

$$W_{ij}^{eff}(\mathbf{q}, \omega) = V_{ij}^{bare}(\mathbf{q}) + \frac{1}{q}\sum_{n} F_{ij,nm}(\mathbf{q})q_{nm}^{s}(\mathbf{q}, \omega)W_{nm}^{eff}(\mathbf{q}, \omega) \tag{252}$$

It is interesting to point out that if we replace the screened matrix elements on the RHS of Eqs. (251) and (252) with the bare ones, in the static and long-wavelength limit we arrive at the classical result (Hartree approximation).

Note that screening significantly affects the transport properties in a 2DEG layer, so that an appropriate screening approximation needs to be employed. It is straightforward to employ static screening for Coulomb and surface-roughness scattering, but its applicability for phonon scattering is questionable. For example, if we consider the long-wavelength acoustic phonons, they always have a nonzero component of the wavevector in the direction perpendicular to the interface. Hence, for $\mathbf{q} \to 0$, the frequency of the phonons is finite, whereas the plasma frequency,

$$\omega_{pl}(q) = \sqrt{\frac{e^2 N_s q}{2\kappa m_{xy}^*}} \tag{253}$$

of the 2DEG (in the electronic quantum limit case) approaches zero as $\mathbf{q} \to 0$. Since the phonon frequency is much larger than the plasma frequency, the screening is ineffective and the use of the bare matrix element, rather than using a statically screened one is *less wrong* [static screening is valid in the limit when $\omega_{pl}(q) \gg \omega^{ex}(q)$, where $\omega^{ex}(q)$ is the frequency of the external perturbation]. Similar arguments are valid for intervalley scattering with short-wavelength phonons. Because the large wave vectors enter the transition and also because the frequencies of these transitions are usually much larger than the frequencies at which the plasma can respond, the intervalley scattering is weakly affected by screening, so that static screening is not justified and one has to perform dynamical calculations or leave the bare matrix elements in order to treat screening properly.

## 3.3. Hydrodynamic and Drift-Diffusion Model

In a number of practical applications, it is not necessary to know the exact distribution function obtained by solving the Boltzmann transport equation (BTE). Instead, it suffices to know only the lowest moments, like the mean and the variance, for example. For this purpose, semi-classical transport equations are derived based on either the first three or four moments of the distribution function that describe the carrier concentration, current, average carrier energy and energy flux variation (if four moments are retained). The various coefficients that appear in these equations may be assumed to be a function of the average carrier energy. The relationship between these coefficients and the energy is usually determined from steady-state Monte Carlo calculations and experimental data for homogeneous samples.

The first three moments of the BTE for electrons, which describe conservation of particles, momentum (mass flow) and energy, are expressed by the following set of equations [228]

$$\frac{\partial n}{\partial t} = -\nabla \cdot (n\mathbf{v}) \tag{254}$$

$$\frac{\partial P_i}{\partial t} = -\sum_j 2\frac{\partial W_{ji}}{\partial x_j} - neE_i - P_i\left\langle\frac{1}{\tau_m}\right\rangle \tag{255}$$

$$\frac{\partial W}{\partial t} = -\sum_i \frac{\partial J_{w_i}}{\partial x_i} + \mathbf{J} \cdot \mathbf{E} - (W - W_0)\left\langle\frac{1}{\tau_e}\right\rangle \tag{256}$$

where $n$ is the electron density, $\mathbf{v}$ is the average electron velocity, $P_i$ and $E_i$ ($i = 1, 2, 3$ for $x$-, $y$-, and $z$-coordinate) are the $i$th components of the total momentum and the electric field, $W_{ij}$ is a component of the total kinetic energy density tensor, $W$ is the total kinetic energy density ($W_0$ being the equilibrium electron energy corresponding to the lattice temperature $T_L$), $\mathbf{J} = -e\mathbf{P}/m^*$ is the current density, and $\mathbf{J}_w$ is the kinetic energy flux. The momentum and energy relaxation rates that appear in Eqs. (255) and (256) are defined as

$$\left\langle\frac{1}{\tau_m}\right\rangle = \frac{\sum_k Pf(\mathbf{k})/\tau_m(\mathbf{k})}{P}, \quad \text{and} \quad \left\langle\frac{1}{\tau_e}\right\rangle = \frac{\sum_k \varepsilon(\mathbf{k})f(\mathbf{k})/\tau_e(\mathbf{k})}{W - W_0} \tag{257}$$

Assuming that the carrier velocity $\mathbf{v}$ equals the sum of a drift ($\mathbf{v}_d$) and a thermal component ($\mathbf{c}$), that is, $\mathbf{v} = \mathbf{v}_d + \mathbf{c}$, one can express the total kinetic energy as being equal the sum of a drift and a thermal energy component due to the random thermal motion of the carriers; that is,

$$W = \frac{1}{2}nm^*v_d^2 + \frac{3}{2}k_BT_c = K + \frac{3}{2}k_BT_c \tag{258}$$

where $k_B$ is the Boltzmann constant and $K$ is the drift component of the kinetic energy density. The kinetic energy flux term appearing in Eq. (258) then reduces to

$$\mathbf{J}_w = W\mathbf{v}_d + \mathbf{v}_d \cdot (nk_B\overline{T}_c) + \mathbf{Q} \tag{259}$$

where $\overline{T}_c$ is the temperature tensor and $\mathbf{Q}$ is the heat flux vector.

Further simplifications to the momentum and energy balance equations are usually made assuming a displaced Maxwellian form for the distribution function, which leads to a diagonal temperature tensor. This approximation is valid for systems in which the electron-electron interactions play a significant role. The use of a displaced Maxwellian distribution function leads to the following set of balance (also known as hydrodynamic) equations

$$\frac{\partial n}{\partial t} = -\nabla \cdot (n\mathbf{v}) \tag{260}$$

$$\frac{\partial \mathbf{v}_d}{\partial t} + \mathbf{v}_d \cdot \nabla\mathbf{v}_d + \frac{1}{nm^*}\nabla(nk_BT) = -\mathbf{v}_d\left\langle\frac{1}{\tau_m}\right\rangle \tag{261}$$

$$\frac{\partial W}{\partial t} = -\nabla \cdot (W\mathbf{v}_d + nk_BT\mathbf{v}_d - \kappa\nabla T) + \mathbf{J} \cdot \mathbf{E} - (W - W_0)\left\langle\frac{1}{\tau_e}\right\rangle \tag{262}$$

Note that the displaced Maxwellian distribution, which is symmetric in momentum space, will lead to zero heat flux. since it involves the third moment of the distribution function. However, Bløtekjær [229] has pointed out that this term may be significant for non-Maxwellian distributions, so that a phenomenological description for the heat flux $\mathbf{Q} = -\kappa\nabla T$ has been used in Eq. (262), where $\kappa$ is the thermal conductivity. As already mentioned, the ensemble averaged energy dependent momentum and energy relaxation rates that appear in Eqs. (261) and (262), are determined by steady-state Monte Carlo simulation for bulk material under uniform electric fields.

For simulations where steady state solutions are required, or transient events with relatively large time scales are being investigated, it is possible to neglect the terms $\partial\mathbf{v}_d/\partial t$,

$\partial W / \partial t$ and $\nabla \cdot (n \mathbf{v}_d)$ in the momentum and energy balance equations. Furthermore, if carrier-heating effects are negligible, then the drift component of the kinetic energy density can be ignored. With these simplifications, and assuming that there are no temperature gradients in the system, the steady-state momentum balance equation leads to the following expression for the current density

$$\mathbf{J} = -en\mathbf{v}_d = en\mu\mathbf{E} + eD\nabla n \qquad (263)$$

where the mobility and the diffusion coefficient are calculated using [230]

$$\mu_n = \frac{e}{m^* \langle 1/\tau_m \rangle} \quad \text{and} \quad D_n = \frac{1}{e} k_B T \mu_n \qquad (264)$$

The result given in Eq. (264), together with the Eq. (260), constitutes the drift-diffusion model for electrons. A similar set of equations can be written for holes. Since a low-field limit has been assumed in order to arrive at the result given in Eq. (263), the mobility and the diffusion coefficients are energy independent quantities. To extend the validity of the drift-diffusion model to the high-field regime, *ad hoc* inclusion of field dependent mobilities and diffusion coefficients is usually used in standard device simulators, such as Silvaco's ATLAS. However, the applicability of such an approach becomes questionable in nanoscale devices in which non-stationary and ballistic transport effects play significant role.

## 4. FIELD EQUATIONS

In the previous sections, we have discussed transport models within the context of the semi-classical BTE. Of equal or greater importance in terms of the behavior of electronic devices are the self-consistent fields inside the device associated with the external bias and internal charge and current distributions. In particular, the carriers within a semiconductor are accelerated by the electric and magnetic fields according to the Lorentz force equation

$$\mathbf{F} = \hbar \frac{d\mathbf{k}}{dt} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}) \qquad (265)$$

where $\mathbf{B}$ and $\mathbf{E}$ are the magnetic flux density and electric field intensity respectively. In general, these fields correspond to the solution of Maxwell's equations originating from the microscopic charges and currents in the device. For high frequency device modeling, in which the device dimensions are comparable to the wavelength, wave propagation effects may be important, and full wave solutions of Maxwell's equations are necessary. In considering opto-electronic devices, direct solutions of either Maxwell's or the associated Helmholtz equation are necessary to represent optical field within the device. Such approaches has been taken, for example, in modeling microwave transistors under the context of "global modeling" [231], as well as for the analysis of semiconductor device, in which optical cavity modes are coupled to semiconductor device simulation [232]. For most device modeling applications, the magnetic contribution to the Lorentz force is much smaller than the electric field contribution, and wave propagation effects are negligible, so that quasi-static representation of the fields in terms of the solution of Poisson's equation are sufficient.

In Section 4.1, we first discuss direct solution of Maxwell's equations using the finite difference time domain method (FDTD), which has been used extensively in the electromagnetics community, and is employed in simulation of high frequency devices and circuits. The subsequent section (Section 4.2) is devoted to the description of efficient numerical solution methods for the Poisson's equation, as applied to semiconductor device simulation. The coupling of the various field solvers in semiconductor device simulation is then discussed in Section 5.

### 4.1. Finite Difference Time Domain Techniques

For electromagnetic solvers in general, there exist a number of general commercial packages for solving Maxwell's equations, such as ANSOFT's HFSS program [233]. For semiconductor device simulation, the time domain FDTD method mentioned above is convenient since time

domain methods are also used in the solution of the transport equations as discussed earlier. Commercial codes are available for FDTD electromagnetic simulation (see for example the XFDTD [234] and Fidelity [235] codes). ISE has recently released a commercial simulation tool combining their device simulation tool DESSIS with an FDTD based simulator (EMLAB) for high-frequency simulation of semiconductor devices [236]. In the following, we give a brief description of the FDTD method itself, and its coupling to particle based simulators.

Maxwell's equations in SI units are written as

$$\Gamma \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad \nabla \cdot \mathbf{D} = \rho$$

$$\nabla \times \mathbf{H} = \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t}, \quad \nabla \cdot \mathbf{B} = 0$$

(266)

Here $\rho$ is the free charge density

$$\rho(\mathbf{r}) = q(N_D - N_A + p - n)$$

(267)

where $N_D$ and $N_A$ are the ionized donor and acceptor concentrations, while $p$ and $n$ are the hole and electron concentrations which are functions of position. For linear isotropic media, the relations for the various fields above is simplified by the constitutive relationships

$$D = \varepsilon E, \quad \text{and} \quad B = \mu H$$

(268)

where $\mu$ is the permeability ($\mu_0 = 4\pi \times 10^{-7}$ for nonmagnetic semiconductors) and $\varepsilon$ is the permittivity.

In Cartesian coordinates, the curl equations are expanded as

$$\frac{\partial H_x}{\partial t} = \frac{1}{\mu}\left(\frac{\partial E_y}{\partial z} - \frac{\partial E_z}{\partial y}\right)$$

$$\frac{\partial H_y}{\partial t} = \frac{1}{\mu}\left(\frac{\partial E_z}{\partial x} - \frac{\partial E_x}{\partial z}\right)$$

$$\frac{\partial H_z}{\partial t} = \frac{1}{\mu}\left(\frac{\partial E_x}{\partial y} - \frac{\partial E_y}{\partial x}\right)$$

$$\frac{\partial E_x}{\partial t} = \frac{1}{\varepsilon}\left(\frac{\partial H_z}{\partial y} - \frac{\partial H_y}{\partial z}\right) + J_x$$

$$\frac{\partial E_y}{\partial t} = \frac{1}{\varepsilon}\left(\frac{\partial H_x}{\partial z} - \frac{\partial H_z}{\partial x}\right) + J_y$$

$$\frac{\partial E_z}{\partial t} = \frac{1}{\varepsilon}\left(\frac{\partial H_y}{\partial x} - \frac{\partial H_x}{\partial y}\right) + J_z$$

(269)

These may then be discretized [237] using the so-called Yee cell [238]. The Yee cell consists of a set of interpenetrating finite difference grids, one representing the electric the other the magnetic fields, over which the derivatives in space are expanded. The electric fields are assumed to be updated at time step $n$ while the magnetic fields are updated at time step $n + 1/2$. Using central differences, the discretized Maxwell's equations are written as

$$E_z^{n-1}(i, j, k + 0.5) = E_z^n(i, j, k + 0.5)$$

$$+ \frac{\Delta t}{\varepsilon}\left[\frac{H_y^{n+0.5}(i + 0.5, j, k + 0.5) - H_y^{n+0.5}(i - 0.5, j, k + 0.5)}{\Delta x}\right.$$

$$- \frac{H_x^{n+0.5}(i, j + 0.5, k + 0.5) - H_x^{n+0.5}(i, j - 0.5, k + 0.5)}{\Delta y}$$

$$\left. + J_z^{n+0.5}(i, j, k + 0.5)\right]$$

(270)

for the electric field, and

$$H_x^{n+0.5}(i, j + 0.5, k + 0.5) = H_x^{n-0.5}(i, j + 0.5, k + 0.5)$$

$$- \frac{\Delta t}{\mu} \left[ \frac{E_z^n(i, j + 1, k + 0.5) - E_z^n(i, j, k + 0.5)}{\Delta y} \right.$$

$$\left. - \frac{E_y^n(i, j + 0.5, k + 1) - E_y^n(i, j + 0.5, k)}{\Delta z} \right] \quad (271)$$

for the magnetic field strength. Similar equations hold for the other components of **E** and **H**. As can be seen in Eqs. (270) and (271), the electric field at time step $n + 1$ is determined explicitly by the electric field at time step $n$, and the magnetic fields in the adjoining Yee cell mesh points at the previous half time step, $n + 1/2$. Likewise, the magnetic field at time step $n + 1/2$ is calculated explicitly by the magnetic field at time step $n - 1/2$ and the electric field in adjacent mesh points at time step $n$. Hence the time evolution of the electric and magnetic fields is calculated noniteratively in a time marching fashion in half time-step intervals. The stability of this technique naturally depends on the time step and grid spacing.

The grid cell size is typically chosen to minimize the effects of numerical dispersion. If the highest frequency component in the simulation is characterized by a wavelength $\lambda = c/\nu$, then empirically, the cell size should be smaller than approximately $\lambda/10$ to avoid artificial dispersion effects. Other considerations in the grid size depend on the geometrical considerations of the structure being simulated. Once the spatial grid has been determined, the time step, $\Delta t$, used to propagate the electric and magnetic fields forward in time, has an upper bound determined by the furthest distance a signal can propagate over this time interval [239]. Setting this distance equal to the minimum grid spacing, this constraint yields

$$\Delta t \leq \frac{dr\sqrt{\mu\varepsilon}}{\sqrt{3}} = \frac{dr}{\sqrt{3}v_p} \quad (272)$$

where $dr$ is the minimum space increment in any direction, and $v_p$ is the phase velocity. More generally, the *Courant stability condition* [240] is of the form

$$\Delta t \leq \frac{1}{v_p\sqrt{\frac{1}{(\Delta x)^2} + \frac{1}{(\Delta y)^2} + \frac{1}{(\Delta z)^2}}} \quad (273)$$

where $\Delta x$, $\Delta y$, and $\Delta z$ are the minimum grid spacings in the three Cartesian coordinate directions.

In the FDTD method, the current density, **J**, appearing above is the primary coupling between charge transport and the coupled electromagnetic fields. During the half time interval between the calculation of the electric or magnetic field components, charge transport is simulated over this time interval using the frozen field components of the previous time step. The updated current density at each grid point is then used in Eq. (270) to update the electric field in the FDTD algorithm. Within the drift diffusion or hydrodynamic models, the current density is calculated directly from the continuity equation. Within a particle-based scheme, one has to map the continuous particle motion onto the discrete grid points. In a nearest grid point (NGP) scheme, the weighted velocities of the particle in the ensemble that lie within a unit cell volume around a given grid point, are summed according to

$$J(i, j, k, t) = \frac{1}{\Delta x \Delta y \Delta z} \sum_{n=1}^{N(i,j,k)} S_n v_n \quad (274)$$

where $S_n$ and $v_n$ refer respectively to the charge and velocity of the $n$th particle associated with the grid point, and $N(i, j, k)$ is the total number of particles within a unit cell around the grid point $(i, j, k)$.

The solution of Maxwell's equations using the FDTD technique requires the imposition of boundary conditions. The usual conditions for the continuity of the tangential and perpendicular components of the electric and magnetic fields are applied at dielectric and metallic boundaries. The simulation of open systems (e.g., an infinite domain) requires special care in that boundary conditions on the simulation domain must be specified to minimize the

artificial reflection of the outgoing wave. Various types of absorbing boundary conditions have been proposed in the literature [241, 242]. Currently, the most popular technique that minimizes artificial reflection on open boundaries is the so-called perfectly matched layer (PML) method proposed by Berenger [243]. This method utilizes a fictitious magnetic loss for impedance matching of the outgoing wave to a highly lossy material. The PML has been found effective in attenuating outgoing waves for a wide range of frequencies and incident angles on the boundary surface. and is currently the state of the art in the FDTD method.

## 4.2. Poisson's Equation

For most nanoscale semiconductor device modeling, the device dimensions themselves are much smaller than the characteristic wavelengths associated with the maximum frequency of operation, so that quasi-static solutions of Maxwell's equations are sufficient. In this case, the electric field may be expressed as the gradient of a scalar potential. The divergence equation for the electric displacement in Eq. (266) then yields Poisson's equation, which in 2D or 3D is written

$$\nabla^2 V(\mathbf{r}) = -\frac{q(N_D - N_A + p - n)}{\varepsilon_{sc}} = -\frac{\rho(\mathbf{r})}{\varepsilon_{sc}} \tag{275}$$

where, as previously noted, $N_D$ and $N_A$ are the ionized donor and acceptor concentrations, while $p$ and $n$ are the hole and electron concentrations, which are functions of position. Poisson's equation assumes that the field may be described as the gradient of a scalar potential $V$, valid in the quasi-static limit for the associated temporal variation.

In the numerical solution of the 2D or 3D Poisson equation, the application of a conventional finite-difference or finite-elements scheme leads to algebraic equations having a well-defined structure, defined over a finite mesh or grid. For example, using central differences in 2D to write the Laplacian appearing in Eq. (275), the discretized form becomes a system of linear equations

$$a_N V_{i,j+1} + a_S V_{i,j-1} + a_E V_{i+1,j} + a_W V_{i-1,j} + a_C V_{i,j} = -\frac{\rho_{i,j}}{\varepsilon_{sc}} \tag{276}$$

where $i$ and $j$ label the two-dimensional coordinates of a particular grid point, and the coefficients representing the grid cells to the north ($N$), south ($S$), and so on, are determined from the grid spacing in the usual way.

In general, the resulting system of equations can be represented by the matrix equation $\mathbf{Ax} = \mathbf{b}$ [244]. The most suitable methods for the solution of this matrix equation are direct methods, but the computational cost becomes prohibitive as the number of equations increases, which is normally the case in 2D and 3D device simulations. This has led to the development of iterative procedures that utilize the well-defined structure of the coefficient matrix. The simplest and most commonly used iterative procedures are the successive overrelaxation (SOR) and the Alternating Direction Implicit (ADI) methods [245]. Both methods lose their effectiveness when complex problems are encountered and when the equation set becomes large, as it is usually the case in 3D problems. One alternative to avoiding this problem is the Incomplete Lower Upper (ILU) decomposition method described in Section 4.2.1, which provides a significant increase in the power of iterative methods and is, therefore, more suitable for solving 3D problems. Other alternatives for solving large-scale 3D problems are the multi-grid and the conjugate gradient methods that are also discussed below.

### 4.2.1. Incomplete Lower-Upper Decomposition (ILU) Method

Within incomplete factorization schemes [246] for 2D problems, the matrix $\mathbf{A}$ is decomposed into a product of lower ($\mathbf{L}$) and upper ($\mathbf{U}$) triangular matrices, each of which has four nonzero diagonals in the same locations as the ones of the original matrix $\mathbf{A}$. The unknown elements of the $\mathbf{L}$ and $\mathbf{U}$ matrices are selected in such a way that the five diagonals common to both $\mathbf{A}$ and $\mathbf{A}' = \mathbf{LU}$ are identical and the four superfluous diagonals represent the matrix $\mathbf{N}$; that is, $\mathbf{A}' = \mathbf{A} + \mathbf{N}$. Thus, rather than solving the original system of equations

$Ax = b$, one solves the modified system $LUx = b + Nx$, by solving successively the matrix equations $LV = b + Nx$ and $V = Ux$, where $V$ is an auxiliary vector. It is important to note that the four superfluous terms of $N$ affect the rate of convergence of the ILU method. Stone [247] suggested the introduction of partial cancellation, which minimizes the influence of these additional terms and accelerates the rate of convergence of the ILU method. By using a Taylor series expansion, the superfluous terms appearing in $A'$ are partially balanced by subtracting approximately equal terms.

### 4.2.2. Multigrid Methods

The multi-grid method represents an improvement over the SOR and ILU methods in terms of iterative techniques available for solving large systems of equations [248]. The basic principle behind the multigrid method is to reduce different Fourier components of the error on grids with different mesh sizes. Most iterative techniques work by quickly eliminating the high-frequency Fourier components, while the low-frequency ones are left virtually unchanged. The result is a convergence rate that is initially fast, but slows down dramatically as the high-frequency components disappear. The multigrid method utilizes several grids, each with consecutively coarser mesh sizes. Each of these grids acts to reduce a different Fourier component of the error, therefore increasing the rate of convergence with respect to single grid based methods, such as an SOR.

The initial setup for the multigrid solver is to create a sequence of grids. The finest grid is generated according to the device structure, and each consecutive grid is obtained by doubling the spacing of the previous one. This is repeated until the final, coarsest grid contains $3 \times 3$ points. The coarsening process must ensure propagation of the boundary conditions to all grids in order to obtain a unique solution. The Poisson equation is solved on the finest grid, and the residual, which is a computational measure of the error, is passed down or restricted to the next, coarser grid. The next grid solves the error equation. This results in a reduction of the relatively lower-frequency error components, as compared with the initial error on the previous grid. This process is continued through all subsequent grids. At the coarsest grid, the error equation is solved exactly and the error is prolonged up through the finer grids, adding its correction at each grid level. At the finest grid, the correction is used to update the final solution. In this way, the multiple error components are reduced simultaneously and the procedure is repeated until convergence of a solution is obtained.

### 4.2.3. Conjugate Gradient Methods with Preconditioning

The basic conjugate gradient (CG) algorithm is one of the best known iterative techniques for solving sparse symmetric positive definite (SPD) systems, but it loses its applicability when the resulting system of equations is not SPD. In such circumstances, the best alternative are the Lanczos-type algorithms, which solve not only the original system $Ax = b$ but also solve the dual linear system $A^T x^* = b^*$. In recent years, the conjugate gradient squared (CGS) method due to Sonneveld [249] has been recognized as an attractive transpose-free variant of the biconjugate gradient (Bi-CG) iterative method [250]. This method works quite well in many cases, but the very high variations in the residual vectors often cause the residual norms to become inaccurate, which can lead to substantial buildup of rounding errors and overflow. The Bi-CGSTAB method due to Van der Vorst [251] is a variant of the CGS algorithm, which avoids squaring of the residual polynomial. It has been demonstrated that the convergence behavior of this method is smoother because it produces more accurate residual vectors and, therefore, more accurate solutions. In conjunction with the Bi-CGSTAB method, a successful preconditioning matrix can be obtained by using ILU factorization [252]. If $L$ and $U$ are the strictly lower and the strictly upper triangular parts of $A$, then the preconditioning matrix is

$$K_{ILU} = (L + \tilde{D})\tilde{D}^{-1}(U + \tilde{D}) \tag{277}$$

where $diag(K_{ILU}) = diag(A)$. The case $k = 0$ is used here, and means no fill-ins are allowed. Once the diagonal $\tilde{D}$ is computed, scaling of the original matrix $A$ is performed using

$$\tilde{A} = \tilde{D}^{-1/2}A\tilde{D}^{-1/2} = diag(\tilde{A}) + \tilde{L} + \tilde{U} \tag{278}$$

The preconditioning matrix for this symmetrically scaled matrix is of the form

$$\tilde{K} = (\tilde{L} + I)(I + \tilde{U}) \tag{279}$$

where I is the identity matrix. The Bi-CGSTAB method is now applied to the preconditioned system:

$$(\tilde{L} + I)^{-1}\tilde{A}(I + \tilde{U})^{-1}\tilde{x} = (\tilde{L} + I)^{-1}\tilde{b} \tag{280}$$

where $\tilde{b} = \tilde{D}^{-1/2}b$, and the solution of the original system of equations is obtained via $x = \tilde{D}^{-1/2}(I + \tilde{U})^{-1}\tilde{x}$. In the calculation of the product $(\tilde{L} + I)^{-1}\tilde{A}(I + \tilde{U})^{-1}\tilde{x}$, extra work is avoided by using the Eisenstat's trick [253]

$$(\tilde{L} + I)^{-1}\tilde{A}(I + \tilde{U})^{-1}\tilde{p} = \tilde{t} + (\tilde{L} + I)^{-1}\{\tilde{p} + [diag(\tilde{A}) - 2I]\tilde{t}\} \tag{281}$$

where $\tilde{t} = (I + \tilde{U})^{-1}\tilde{p}$.

## 5. SEMICLASSICAL DEVICE SIMULATION

In previous sections, we introduced the numerical solution of the BTE using Monte Carlo methods (Section 3.1), the approximate solutions of the BTE using either hydrodynamic or drift-diffusion model (Section 3.3), and the solution of Maxwell's equations (Section 4.1) and Poisson's equation (Section 4.2) over a finite mesh. Within a device, both the transport kernel and the field solver are coupled to each other. The field associated with the potential coming from Poisson's equation is the driving force accelerating particles in the Monte Carlo phase, for example, while the distribution of mobile (both electrons and holes) and fixed charges (e.g., donors and acceptors) provides the source of the electric field in Poisson's equation corresponding to the right-hand side of Eq. (275). Below we give an extensive description of the Monte Carlo particle-based device simulators with emphasis on the particle-mesh coupling and the inclusion of the short-range Coulomb interaction (Section 5.1). This discussion is followed by a brief summary of hydrodynamic/drift-diffusion device simulators with reference to commercially available simulation software (Section 5.2). We finish this section with an application of the FDTD methods coupled with a MC transport kernel on the example of a coplanar strip on a GaAs substrate (Section 5.3).

### 5.1. Particle-Based Device Simulations

Within the particle-based EMC method with its time-marching algorithm, Poisson's equation may be decoupled from the BTE over a suitably small time step (typically less than the inverse plasma frequency corresponding to the highest carrier density in the device). Over this time interval, carriers accelerate according to the frozen field profile from the previous time-step solution of Poisson's equation, and then Poisson's equation is solved at the end of the time interval with the frozen configuration of charges arising from the Monte Carlo phase (see discussion in Ref. [157]). Note that Poisson's equation is solved on a mesh, whereas the solution of charge motion using EMC occurs over a continuous range of coordinate space in terms of the particle position. Therefore, a particle-mesh (PM) coupling is needed for both the charge assignment and the force interpolation. The PM coupling is broken into four steps: (1) assign particle charge to the mesh; (2) solve the Poisson equation on the mesh; (3) calculate the mesh-defined forces; and (4) interpolate to find forces on the particle. There are a variety of schemes that can be used for the PM coupling and these are discussed in the next section.

Another issue that has to be addressed in particle-based simulations is the real space boundary conditions for the particle part of the simulation. Reflecting boundary conditions are usually imposed at the artificial boundaries. As far as the Ohmic contacts are concerned, they require more careful consideration because electrons crossing the source and drain contact regions contribute to the corresponding terminal current. In order to conserve charge in the device, the electrons exiting the contact regions must be re-injected. Commonly employed models for the contacts include [254]:

- Electrons are injected at the opposite contact with the same energy and wavevector $k$. If the source and drain contacts are in the same plane, as in the case of MOSFET

simulations, the sign of k, normal to the contact will change. This is an unphysical model, however [255].

- Electrons are injected at the opposite contact with a wavevector randomly selected based upon a thermal distribution. This is also an unphysical model.

- Contact regions are considered to be in thermal equilibrium. The total number of electrons in a small region near the contact are kept constant, with the number of electrons equal to the number of dopant ions in the region. This is a very good model most commonly employed in actual device simulations.

- Another method uses 'reservoirs' of electrons adjacent to the contacts. Electrons naturally diffuse into the contacts from the reservoirs, which are not treated as part of the device during the solution of Poisson's equation. This approach gives results similar to the velocity weighted Maxwellian [254], but at the expense of increased computational time due to the extra electrons simulated. It is an excellent model employed in few most sophisticated particle-based simulators.

There are also several possibilities for the choice of the distribution function—Maxwellian, displaced Maxwellian, and velocity-weighted Maxwellian [256].

To simulate the steady-state behavior of a device, the system is started in some initial condition, with the desired potentials applied to the contacts, and then the simulation proceeds in a time stepping manner until steady state is reached. This typically takes several picoseconds of simulation time, and consequently several thousand time-steps based on the usual time increments required for stability. Figure 38 shows the particle distribution in a 3D metal semiconductor field effect transistor (MESFET) structure, where the dots indicate the individual simulated particles [257]. In these simulations, the charge-neutral method, discussed earlier, is used.

After sufficient time has elapsed, so that the system is driven into a steady-state regime, one can calculate the steady-state current through a specified terminal. The device current can be determined via two different, but consistent methods. First, by keeping track of the charges entering and exiting each terminal, the net number of charges over a period of the simulation can be used to calculate the terminal current. The method is quite noisy due to the discrete nature of the carriers. In a second method, the sum of the carrier velocities in a portion of the device are used to calculate the current. For this purpose, the device is divided into several sections along, for example, the x-axis (from source to drain for the case of a MOSFET or MESFET simulation). The number of carriers and their corresponding velocity is added for each section after each free-flight time step. The total x-velocity in each section is then averaged over several time steps to determine the current for that section. The total device current can be determined from the average of several sections, which gives a much smoother result compared to counting the terminal charges. By breaking the device into sections, individual section currents can be compared to verify that the currents are uniform. In addition, sections near the source and drain regions of a MOSFET or a MESFET may have a high y-component in their velocity and should be excluded from the current calculations. Finally, by using several sections in the channel, the average energy and velocity of electrons along the channel is checked to ensure proper physical characteristics.

As in the case of solving the full Maxwell's equations, for a stable Monte Carlo device simulation, one has to choose the appropriate time step, $\Delta t$, and the spatial mesh size ($\Delta x$, $\Delta y$, and/or $\Delta z$). The time step and the mesh size may correlate to each other in connection



Figure 38. Example of the particle distribution in a MESFET structure simulated in 3D using an EMC approach. Reprinted with permission from [286], S. M. Goodnick et al., Int. J. Num. Model. 8, 205 (1995). © 1995, Wiley.

with the numerical stability. For example, the time step $\Delta t$ must be related to the plasma frequency

$$\omega_p = \sqrt{\frac{e^2 n}{\varepsilon_s m^*}}$$ (282)

where $n$ is the carrier density. From the viewpoint of the stability criterion, $\Delta t$ must be much smaller than the inverse plasma frequency. The highest carrier density specified in the device model is used to estimate $\Delta t$. If the material is a multi-valley semiconductor, the smallest effective mass to be experienced by the carriers must be used in Eq. (282) as well. In the case of GaAs, with the doping of $5 \times 10^{17}$ cm$^{-3}$, $\omega_p \cong 5 \times 10^{13}$; hence, $\Delta t$ must be smaller than 0.02 ps.

The mesh size for the spatial resolution of the potential is dictated by the charge variations. Hence, one has to choose the mesh size to be smaller than the smallest wavelength of the charge variations. The smallest wavelength is approximately equal to the Debye length (for degenerate semiconductors the relevant length is the Thomas-Fermi wavelength), given as

$$\lambda_D = \sqrt{\frac{\varepsilon_s k_B T}{e^2 n}}$$ (283)

The highest carrier density specified in the model should be used to estimate $\lambda_D$ from the stability criterion. The mesh size must be chosen to be smaller than the value given by Eq. (283). In the case of GaAs, with the doping density of $5 \times 10^{17}$ cm$^{-3}$, $\lambda_D \cong 6$ nm.

On the basis of previous discussion, the time step ($\Delta t$), and the mesh size ($\Delta x$, $\Delta y$, an/or $\Delta z$) are specified independently based on physical arguments. However, there are numerical constraints as well. This means that $\Delta t$ chosen must be checked again by calculating the distance $l_{max}$, defined as

$$l_{max} = v_{max} \times \Delta t$$ (284)

where $v_{max}$ is the maximum carrier velocity that can be approximated by the maximum group velocity of the electrons in the semiconductor (on the order of $10^8$ cm/s). The distance $l_{max}$ is the maximum distance the carriers can propagate during $\Delta t$. The time step is therefore chosen to be small enough so that $l_{max}$ is smaller than the spatial mesh size chosen using Eq. (283). This constraint arises because for too large of a time step, $\Delta t$, there may be substantial change in the charge distribution, while the field distribution in the simulation is only updated every $\Delta t$, leading to unacceptable errors in the carrier force.

## 5.1.1. Particle-Mesh (PM) Coupling

The charge assignment and force interpolation schemes usually employed in self-consistent Monte Carlo device simulations are the nearest-grid-point (NGP) and the cloud-in-cell (CIC) schemes [258]. In the NGP scheme, the particle position is mapped into the charge density at the closest grid point to a given particle. This has the advantage of simplicity, but leads to a noisy charge distribution, which may exacerbate numerical instability. Alternately, within the CIC scheme a finite volume is associated with each particle spanning several cells in the mesh, and a fractional portion of the charge per particle is assigned to grid points according to the relative volume of the 'cloud' occupying the cell corresponding to the grid point. This method has the advantage of smoothing the charge distribution due to the discrete charges of the particle based method, but may result in an artificial 'self-force' acting on the particle, particularly if an inhomogeneous mesh is used.

To better understand the NGP and the CIC scheme, consider a tensor product mesh with mesh lines $x_i$, $i = 1, \ldots, N_x$ and $y_j$, $j = 1, \ldots, N_y$. If the mesh is uniformly spaced in each axis direction, then $(x_{i-1} - x_i) = (x_{i+2} - x_{i+1})$. The permittivities are considered constant within each mesh element and are denoted by $\varepsilon_{kl}$, $k = 1, \ldots, N_x - 1$ and $l = 1, \ldots, N_y - 1$. Define centered finite-differences of the potential $\psi$ in the $x$- and $y$-axis at the midpoints of

element edges as follows:

$$
\begin{cases}
\Delta^x_{k+\frac{1}{2},l} = -\dfrac{\psi_{k+1,l} - \psi_{k,l}}{x_{k+1} - x_k} \\[2ex]
\Delta^y_{k,l+\frac{1}{2}} = -\dfrac{\psi_{k,l+1} - \psi_{k,l}}{y_{l+1} - y_l}
\end{cases}
\tag{285}
$$

where the minus sign is included for convenience because the electric field is negative of the gradient of the potential. Consider now a point charge in 2-D located at $(x, y)$ within an element $\langle i, j \rangle$. If the restrictions for the permittivity (P) and the tensor-product meshes with uniform spacing in each direction (M) apply, the standard NGP/CIC schemes in two dimensions can be summarized by the following four steps:

1. *Charge assignment to the mesh.* The portion of the charge $\rho_L$ assigned to the element nodes $(k, l)$ is $w_{kl}\rho_L$, $k = i, i+1$ and $l = j, j+1$, where $w_{kl}$ are the four charge weights which sum to unity by charge conservation. For the NGP scheme, the node closest to $(x, y)$ receives a weight $w_{kl} = 1$, with the remaining three weights set to zero. For the CIC scheme, the weights are $w_{ij} = w_x w_y$, $w_{i+1,j} = (1 - w_x)w_y$, $w_{i,j+1} = w_x(1 - w_y)$, and $w_{i+1,j+1} = (1 - w_x)(1 - w_y)$. $w_x = (x_{i+1} - x)/(x_{i+1} - x_i)$ and $w_y = (y_{j+1} - y)/(y_{j+1} - y_j)$.

2. *Solve the Poisson equation.* The Poisson equation is solved by some of the numerical techniques discussed in Section 4.2.

3. *Compute forces on the mesh.* The electric field at mesh nodes $(k, l)$ is computed as $E^x_{kl} = (\Delta^x_{k-\frac{1}{2},l} + \Delta^x_{k+\frac{1}{2},l})/2$ and $E^y_{kl} = (\Delta^y_{k,l-\frac{1}{2}} + \Delta^y_{k,l+\frac{1}{2}})/2$, for $k = i, i+1$ and $l = j, j+1$.

4. *Interpolate to find forces on the charge.* Interpolate the field to position $(x, y)$ according to $E^x = \sum_{kl} w_{kl} E^x_{kl}$ and $E^y = \sum_{kl} w_{kl} E^y_{kl}$, where $k = i, i+1$, $l = j, j+1$ and the $w_{ij}$ are the NGP or CIC weights from step 1.

The requirements (P) and (M) severely limit the scope of devices that may be considered in device simulations using the NGP and the CIC schemes. Laux [259] proposed a new particle-mesh coupling scheme, namely, the nearest-element-center (NEC) scheme, which relaxes the restrictions (P) and (M). The NEC charge assignment/force interpolation scheme attempts to reduce the self-forces and increase the spatial accuracy in the presence of nonuniformly spaced tensor-product meshes and/or spatially dependent permittivity. In addition, the NEC scheme can be utilized in one axis direction (where local mesh spacing is nonuniform) and the CIC scheme can be utilized in the other (where local mesh spacing is uniform). Such hybrid schemes offer smoother assignment/interpolation on the mesh compared to the pure NEC. The new steps of the pure NEC PM scheme are as follows:

1'. *Charge assignment to the mesh.* Divide the line charge $\rho_L$ equally to the four mesh points of the element $\langle i, j \rangle$.

3'. *Compute forces on the mesh.* Calculate the fields $\Delta^x_{i+(1/2),l}$, $l = j, j+1$, and $\Delta^y_{k,j+(1/2)}$, $k = i, i+1$.

4'. *Interpolate to find force on the charge.* Interpolate the field according to the following $E^x = (\Delta^x_{i+(1/2),j} + \Delta^x_{i+(1/2),j+1})/2$ and $E^y = (\Delta^y_{i,j+(1/2)} + \Delta^y_{i+1,j+(1/2)})/2$.

The NEC designation derives from the appearance, in step (1') of moving the charge to the center of its element and applying a CIC-like assignment scheme. The NEC scheme involves only one mesh element and its four nodal values of potential. This locality makes the method well-suited to nonuniform mesh spacing and spatially varying permittivity. The interpolation and error properties of the NEC scheme are similar to the NGP scheme.

### 5.1.2. The Short-Range Force

In modern deep-submicrometer devices, for achieving optimum device performance and eliminating the so-called punch-through effect, the doping densities must be quite high. This necessitates a careful treatment of the electron-electron ($e$-$e$) and electron-impurity ($e$-$i$) interactions, an issue that has been a major problem for quite some time. Many of the approaches used in the past have included the short-range portions of the $e$-$e$ and $e$-$i$

interactions in the k-space portion of the Monte Carlo transport kernel, as discussed in Section 3.1.5, thus neglecting the important inelastic properties of these two interaction terms [260, 261]. An additional problem with this screened scattering approach is that, unlike the other scattering processes, $e$-$e$ and $e$-$i$ scattering rates need to be re-evaluated frequently during the simulation process to take into account the changes in the distribution function and the screening length. The calculation of the distribution function is highly CPU intensive, and it cannot account for local variations of the electron density in real space. Furthermore, ionized impurity scattering is usually treated as a simple two-body event, thus ignoring the multi-ion contributions to the overall scattering potential. A simple screening model is usually used that ignores the dynamical perturbations to the Coulomb fields caused by the movement of the free carriers. To overcome the above difficulties, several authors have advocated the use of the coupled ensemble Monte-Carlo-molecular dynamics approach [262–264], that gives simulation mobility results in excellent agreement with the experimental data for high substrate doping levels [264]. However, it is proven to be quite difficult to incorporate this coupled ensemble Monte-Carlo-molecular dynamics approach when inhomogeneous charge densities, characteristic of semiconductor devices, are encountered [261, 265]. An additional problem with this approach in a typical particle-based device simulation arises from the fact that both the $e$-$e$ and $e$-$i$ interactions are already included, at least within the Hartree approximation (long-range carrier-carrier interaction), through the self-consistent solution of the three-dimensional (3D) Poisson equation via the PM coupling discussed in Section 5.1.1c. The magnitude of the resulting mesh force that arises from the force interpolation scheme, depends upon the volume of the cell, and, for commonly employed mesh sizes in device simulations, usually leads to double counting of the force.

To overcome the above-described difficulties of incorporation of the short-range $e$-$e$ and $e$-$i$ force into the problem, one can follow two different paths. One way is to use the $P^3M$ scheme introduced by Hockney and Eastwood [258]. An alternative to this scheme is to use the corrected-Coulomb approach due to Gross et al. [266–269].

### 5.1.2.1. The $P^3M$ Method.

The particle-particle-particle-mesh ($P^3M$) algorithms are a class of hybrid algorithms developed by Hockney and Eastwood [258]. These algorithms enable correlated systems with long-range forces to be simulated for a large ensemble of particles. The essence of the method is to express the interparticle forces as a sum of two component parts; the short range part $F_{sr}$, which is nonzero only for particle separations less than some cutoff radius $r_c$, and the smoothly varying part $F$, which has a transform that is approximately band-limited. The total short-range force on a particle $F_{sr}$ is computed by direct particle-particle (PP) pair force summation, and the smoothly varying part is approximated by the particle-mesh (PM) force calculation.

Two meshes are employed in the $P^3M$ algorithms: the charge-potential mesh and a coarser mesh, the so-called chaining mesh. The charge potential mesh is used at different stages of the PM calculation to store, in turn, charge density values, charge harmonics, potential harmonics and potential values. The chaining mesh is a regular array of cells whose sides have lengths greater than or equal to the cutoff radius $r_c$ of the short-range force. Associated with each cell of this mesh is an entry in the head-of-chain array: This addressing array is used in conjunction with an extra particle coordinate, the linked-list coordinate, to locate pairs of neighboring particles in the short-range calculation.

The particle orbits are integrated forward in time using the leapfrog scheme:

$$x_i^{n+1} = x_i^n + \frac{p_i^{n+1/2}}{m}\Delta t \tag{286}$$

$$p_i^{n+1/2} = p_i^{n-1/2} + (F_i + F_i^{sr})\Delta t \tag{287}$$

The positions $\{x_i\}$ are defined at integral time-levels and momenta $\{p_i\}$ are defined at half-integral time levels. Momenta $\{p_i\}$ are used rather than velocities for reasons of computational economy.

To summarize, the change in momentum of particle $i$ at each time step is determined by the total force on that particle. Thus, one is free to choose how to partition the total

force between the short range and the smoothly varying part. The reference force **F** is the interparticle force that the mesh calculation represents. For reasons of optimization, the cutoff radius of $F_{sr}$ has to be as small as possible, and therefore **F** to be equal to the total interparticle force down to as small a particle separation as possible. However, this is not possible due to the limited memory storage and the required CPU time even in the state-of-the-art computers.

The harmonic content of the reference force is reduced by smoothing. A suitable form of reference force for a Coulombic long-range force is one which follows the point particle force law beyond the cutoff radius $r_c$, and goes smoothly to zero within that radius. The smoother the decay of $F(x)$ and the large $r_c$ becomes, the more rapidly the harmonics $R(k)$ decay with increasing **k**. Such smoothing procedure is equivalent to ascribing a finite size to the charged particle. As a result, a straightforward method of including smoothing is to ascribe some simple density profile $S(x)$ to the reference interparticle force. Examples of shapes, which are used in practice, and give comparable total force accuracy are the uniformly charged sphere, the sphere with uniformly decreasing density, of the form

$$S(r) = \begin{cases} \dfrac{48}{\pi a^4}\left(\dfrac{a}{2} - r\right), & r < \dfrac{a}{2} \\ 0, & \text{otherwise} \end{cases} \tag{288}$$

and the Gaussian distribution of density. The second scheme gives marginally better accuracies in 3D simulations. Note that the cutoff radius of the short-range force implied by Eq. (288) is $a$ rather than $r_c$. In practice, one can make $r$ significantly smaller than $a$, because continuity of derivatives at $r = a$ causes the reference force to closely follow the point particle force for radii somewhat less than $a$. It has been found empirically that a good measure of the lower bound of $r_c$ is given by the cube root of the autocorrelation volume of the charge shapes, which for the case of uniformly decreasing density gives

$$r_c \geq \left(\frac{5\pi}{48}\right)^{1.3} a \approx 0.7a \tag{289}$$

Once the reference interparticle force **F** for the PM part of the calculation is chosen, the short-range part $F_{sr}$ is found by subtracting **F** from the total interparticle force, i.e.,

$$F_{sr} = F^{tot} - F \tag{290}$$

### 5.1.2.2. The Corrected Coulomb Approach.

This second approach is a purely numerical scheme that generates a corrected Coulomb force look-up table for the individual $e\text{-}e$ and $e\text{-}i$ interaction terms. To calculate the proper short-range force, one has to define a 3D box with uniform mesh spacing in each direction. A single (fixed) electron is then placed at a known position within a 3D domain, while a second (target) electron is swept along the "device" in, for example, 0.2 nm increments so that it passes through the fixed electron. The 3D box is usually made sufficiently large so that the boundary conditions do not influence the potential solution. The electron charges are assigned to the nodes using one of the charge-assignment schemes discussed previously [259]. A 3D Poisson equation solver is then used to solve for the node or mesh potentials. At self-consistency, the force on the swept electron $F = F_{mesh}$ is interpolated from the mesh or node potential. In a separate experiment, the Coulomb force $F_{tot} = F_{coul}$ is calculated using standard Coulomb law. For each electron separation, one then tabulates $F_{mesh}$, $F_{coul}$ and the difference between the two $F' = F_{coul} - F_{mesh} = F_{sr}$, which is called the corrected Coulomb force or a short-range force. The later is stored in a separate look-up table.

As an example, the corresponding fields to these three forces for a simulation experiment with mesh spacing of 10 nm in each direction are shown in Fig. 39. It is clear that the mesh force and the Coulomb force are identical when the two electrons are separated several mesh points (30–50 nm apart). Therefore, adding the two forces in this region would result in double-counting of the force. Within three to five mesh points, $F_{mesh}$ starts to deviate from $F_{coul}$. When the electrons are within the same mesh cell, the mesh force approaches

Figure 39. Mesh, Coulomb, and corrected Coulomb field versus the distance between the two electrons. Note: F = eE. Reprinted with permission from [266], W. S. Gross et al., *IEEE Electron Devices Lett.* 20, 463 (1999). © 1999, IEEE.

zero, due to the smoothing of the electron charge when divided amongst the nearest node points. The generated look-up table for F' also provides important information concerning the determination of the minimum cutoff range based upon the point where $F_{coul}$ and $F_{mesh}$ begin to intersect, i.e., F' goes to zero.

Figure 40 shows the simulated doping dependence of the low-field mobility, derived from 3D resistor simulations, which is a clear example demonstrating the importance of the proper inclusion of the short-range electron-ion interactions. For comparison, also shown in this figure are the simulated mobility results reported in [270], calculated with a bulk EMC technique using the Brooks-Herring approach [271] for the e-i interaction, and finally the measured data [272] for the case when the applied electric field is parallel to the ⟨100⟩ crystallographic direction. From the results shown, it is obvious that adding the corrected Coulomb force to the mesh force leads to mobility values that are in very good agreement with the experimental data. It is also important to note that, if only the mesh force is used in the free-flight portion of the simulator, the simulation mobility data points are significantly higher than the experimental ones due to the omission of the short-range portion of the force.

The short-range e-e and e-i interactions also play significant role in the operation of semiconductor devices. For example, carrier thermalization at the drain end of the MOSFET channel is significantly affected by the short-range e-e and e-i interactions. This is illustrated in Fig. 41 on the example of a 80-nm channel-length n-MOSFET. Carrier thermalization



Figure 40. Low-field electron mobility derived from 3D resistor simulations versus doping. Also shown on this figure are the Ensemble Monte Carlo results and the appropriate experimental data. Reprinted with permission from [266], W. S. Gross et al., *IEEE Electron Devices Lett.* 20, 463 (1999). © 1999, IEEE.

Figure 41. Average energy of the electrons coming to the drain from the channel. Filled (open) circles correspond to the case when the short-range $e$-$e$ and $e$-$i$ interactions are included (omitted) in the simulations. The channel length extends from 50 to 130 nm. Reprinted with permission from [267]. W. J. Gross et al., *VLSI Design* 10, 437 (2000). © 2000, Taylor and Francis.

occurs over distances that are on the order of few nm when the $e$-$e$ and $e$-$i$ interactions are included in the problem. Using the mesh force alone does not lead to complete thermalization of the carriers along the whole length of the drain extension, and this can lead to inaccuracies when estimating the device on-state current.

## 5.2. Hydrodynamic/Drift-Diffusion Device Simulations

As already discussed in Section 3.3, the balance equations are a set of coupled conservation laws in the form of differential equations that can be easily derived from the BTE. In fact, BTE can be fully represented by an infinite set of such conservation laws, starting with particle density conservation, followed by momentum conservation, energy conservation, etc. This set can be regarded as equivalent to a series expansion of the BTE. However, in order to be of any use, the expansion has to be truncated after a suitable number of terms. Hence, in practice, only a limited number of the most important conservation laws are retained, which may suffice for a satisfactory analysis of most devices. In fact, the much used drift-diffusion formalism, discussed at the end of Section 3.3, is based on the first two conservation laws only—the particle density and the momentum balance equations. But in order to describe important effects in modern-day devices related to nonstationary electron transport and heating of the carrier gas, the third conservation law—the energy balance equation—is also needed, which leads to what is known as the hydrodynamic formulation, described in details in Section 3.3. With this model, phenomena such as velocity overshoot and thermalization of energetic carriers via collisions are described. When the corresponding hydrodynamic or drift-diffusion equations are coupled with a field solver, one has a hydrodynamic or a drift-diffusion device simulator.

Commercial 2D and even 3D simulators based on either drift-diffusion or the hydrodynamic formalism, such as PISCES [273], MEDICI [274], MINIMOS [275], DESSIS from ISE [276], SILVACO [277], etc., are quite popular, especially for analyzing silicon devices, because the effects of electron heating are not as pronounced in silicon as in compound semiconductors. But these simulators are usually not accurate enough for a quantitative description of short-channel compound semiconductor devices. Still, the drift-diffusion formalism may be quite useful for numerically challenging tasks, such as three-dimensional FET modeling. Another example may be the analysis of the field distribution near the pinch-off for estimates of the breakdown voltage in MESFETs and HFETs. They also have been very successfully used in investigating fluctuations in the threshold voltage and the off-state power dissipation in nano-scale MOSFETs, in which there are very few impurity atoms in the device active region, and the position of each impurity will have significant influence on the actual device performance [278].

An example for the role of the atomistic description of the impurity atoms on the potential profile and current stream lines is given in Figure 42 for a device with gate-length equal to 0.1 $\mu$m and gate-width equal to 0.05 $\mu$m. The effect of the randomly sited impurities can

**Figure 42.** Conduction band edge (left panel) and current stream lines in n-channel ultrasmall MOSFET devices in which atomistic description of the impurity atoms is used in the device active region.

be seen on the potential plots on the left for two devices with different number and different impurity distribution. The potential fluctuations force the current to divert around the potential peak of the random impurity. This may be seen on the figures on the right where the current flow vectors avoid several regions, which represent the role of the impurities. Such nonuniform current flow leads to threshold voltage and off-state current fluctuations amongst devices fabricated on the same chip. The effect becomes more prominent as channel lengths are scaled into the nm range, thus giving rise to device reliability concerns.

## 5.3. Electromagnetic Device Simulation

Simulation of optoelectronic and high frequency devices requires the solution of the full set of Maxwell's equations rather than just the Poisson's equation. We previously discussed the FDTD method for solving Maxwell's equations in Section 4.1. Numerous other solution techniques are of course available such as finite elements methods (FEM), moment methods, and frequency domain techniques [279]. For semiconductor lasers, drift diffusion and hydrodynamic models have been coupled to solutions of, for example, the Helmholtz wave equation in the optical cavity to simulate laser performance in MINILASE-II [280, 281].

The modeling of optoelectronic devices such as semiconductor lasers and light-emitting diodes using particle based device simulation is computationally difficult due to the characteristic time scale for spontaneous emission, which is on the order of nanoseconds. In contrast, the time-step in a Monte Carlo simulation is typically a femtosecond, which requires an enormous number of time-steps just to characterize a few radiative transitions. Monte Carlo is still used to calibrate the moment method models used in MINILASE. For example, Rota et al. [282] used Monte Carlo simulation to investigate the capture process for electrons into a quantum well laser for calibrating the rate models used in MINILASE. However, the issue of dealing with vastly different time scale phenomena within a time-domain algorithm, such as the EMC method, is challenging.

For high frequency devices and circuits, the characteristic time scales of scattering events and the inverse frequency are somewhat commensurate, and FDTD methods, as discussed in Section 4.1, have been successfully applied. For time domain techniques, such as the FDTD method, the same coupled algorithm for device simulation is used as already discussed in

this section, in which the field equations are decoupled from the transport equations over a small time interval, and the solution of one is used as the input for the other. For coupled Monte Carlo/FDTD simulation, the carriers are accelerated by the full Lorentz force given by Eq. (136), while the source for Maxwell's equations is provide by the current density on the FDTD mesh, calculated by projecting the carrier velocities onto the nearest mesh point.

As an example of the application of the coupled Monte Carlo/FDTD simulation method, a structure which has been utilized to study carrier dynamics under high field conditions is the biased co-planar strip configuration shown in Fig. 43 [283, 284]. In this structure, an ultra-short optical pulse (switching beam) is used to excite electron-hole pairs between the co-planar strips. Due to the dc bias on the strips, the electrons and holes are excited in opposite directions giving rise to a time-dependent photocurrent induced in the waveguide structure. The pulse propagates down the waveguide, where it is detected by a time-delayed pulse (sampling beam). The change in electric field due to the propagating pulse is detected (for example) by the shift in an excitonic resonance due to the Stark effect [284]. Measurement of the time-dependent photocurrent detects the time-dependent velocity of the electrons and holes as they accelerate in the electric field from an initial state of essentially zero velocity.

Monte Carlo simulation has been employed in the interpretation of these results using FDTD solutions of Maxwell's equations, which are solved self-consistently with the particle dynamics [285, 286]. In these simulations, Maxwell's equations for the electric and magnetic fields are discretized onto a 3D Yee cell grid using the FDTD method described in Section 4.1, and solved at each time step using as a source term the current density calculated from the previous Monte Carlo phase of the simulation during the previous time step. Using the solutions for **E** and **B** from this step, the particles then accelerate under the influence of the Lorentz force during the next Monte Carlo phase. Assuming the time step is properly chosen, this method allows the evolution of the system to be modeled during and after photoexcitation by the switching beam. Figure 44 shows a typical result of the calculated particle current induced by the switching beam for an average field of 40 kV/cm in a GaAs co-planar strip structure with an excited carrier population of $1 \times 10^{17}/cm^3$ in a 1 $\mu$m spot diameter between the strip-lines. The particle current shows an overshoot behavior, which is expected for the short-time dynamics of carriers accelerated in an electric field [137]. However, the decay of the current back to zero, and the undershoot is not expected from simple carrier dynamics in a constant field, and arises due to the self-consistent field of the electrons and holes themselves which collapses the dc field existing in the gap. Whether one uses full solutions to Maxwell's equations, or simply Poisson's equation (quasi-static solutions), the result is fairly similar, and this is clearly seen from the results shown in Fig. 44.

The time-dependent separation of electrons and holes after photoexcitation in the above experiments resembles a Hertzian dipole, which has a characteristic frequency in the



Figure 43. Illustration of the experimental configuration for electro-optic sampling of a biased co-planar strip on a GaAs substrate. Reprinted with permission from [286], S. M. Goodnick et al., *Int. J. Num. Model.* 8, 205 (1995). © 1995, Wiley.

Figure 44. Calculated particle current versus time for the solution based on Poisson's equation only, and solutions considering FDTD solutions to Maxwell's equations for the structure shown in Fig. 43 and an average field of 40 kV/cm. Reprinted with permission from [286], S. M. Goodnick et al., *Int. J. Num. Model.* 8, 205 (1995). © 1995, Wiley.

terahertz range (based upon the time scale shown in Fig. 44). The emission of tera-hertz radiation in such structures has long been recognized and has potential applications in sources and detectors in this frequency range [287]. Son et al. [288] have used the measured tera-hertz radiation in a similar experimental structure to Fig. 43 in order to study the particle dynamics after photoexcitation. Modeling of such tera-hertz radiation using the coupled Monte Carlo/FDTD simulation described above is accomplished by modeling a much larger domain than the co-planar strip structure of Fig. 43 to include the free space outside. Attention must be paid to the proper boundary conditions on the larger domain, which should be purely absorbing to first order. Recent improvements in the FDTD method to approximate absorbing boundary conditions have been developed which result in very little reflected radiation [237]. Figure 45 illustrates the calculated FDTD results for the normal and tangential electric fields in the near-field regime at an observation point directly above the co-planar strips [289]. The results are shown for a time corresponding to the delay time for the electromagnetic pulse to arrive at the observation point. The results bear a close resemblance to the expected wave forms due to radiation from an ideal Hertzian dipole.

# 6. QUANTUM CORRECTIONS TO SEMICLASSICAL APPROACHES

As devices scale toward nanometer dimensions, and new nanoscale devices emerge, the semiclassical techniques in computational electronics discussed in the previous few sections become inaccurate, and new phenomena need to be accounted for, such as quantum mechanical tunneling. In the past, quantum effects have been known to dominate the operation



Figure 45. Calculated near-field electric field outside of the microstrip. (a) The tangential component. (b) The radial component [289]. The *x* and *y* scales are in microns. Reprinted with permission from K. A. Remley et al., *IEEE Trans. Microwave Theory Tech.* 46, 2476 (1998). © 1998, IEEE.

of resonant tunneling diodes [290], quantum cascade lasers [291], etc. As discussed in the introduction (Section 1.1), tunneling through the gate oxide [292], source to drain tunneling and space-quantization effects are expected to be important in nanoscale MOSFETs and will require solution of the one-dimensional (1D) Schrödinger-Poisson problem. Solutions of the two-dimensional (2D) Schrödinger-Poisson problem are needed, for example, for describing the channel charge in narrow-width MOSFETs. With regard to gate-oxide tunneling, the one-electron effective-mass approximation may not be sufficiently accurate and *ab initio* calculations will most probably be needed [293, 294].

It is also relevant to recall from the discussion given in the Section 1.1 that discrete impurity effects in nanoscale MOSFETs will lead to potential fluctuations which, in turn, will affect the magnitude of the device terminal characteristics (threshold voltage, off-state current, off-state power dissipation, etc.). In nanoscale MOSFETs, these potential fluctuations will eventually lead to further electron confinement into small boxes containing only a few electrons as device dimensions scale even further. This implies that understanding transport in future ultra-small devices will also require understanding transport in single and coupled quantum dots. The quantum dots by themselves have been the focus of numerous studies (see, e.g., Ref. [295]). For example, controllable loading of these dots with few electrons has already been achieved, thus allowing one to speak of artificial quantum-dot hydrogen atoms and quantum-dot helium atoms [296]. Computing architectures for quantum devices, so called quantum cellular automata, which consist of cells of coupled quantum dots occupied by only a few electrons, have also been proposed [297, 298] and realized experimentally [299].

In these nanoscale devices, the time scale for the carrier transport is relatively short, as they leave the source with a memory of its distribution, traverse the channel under high fields, and enter the drain. Questions that arise in this context are, for example, what are the requirements for proper transport equations, and how can these be incorporated into existing simulators? Understanding transport in quantum dot structures is yet another challenging problem that needs further consideration. For example, one of the main difficulties in explaining transport in open quantum dots is the determination of the exact energy spectrum and how the dot states couple with the leads (the quasi-two-dimensional electron gas) via the quantum point contacts (QPC). Fluctuations in the confining potential, due to the atomistic nature of the impurity atoms and how they affect the energy level spectrum in the dot, is yet another issue that prevents one in establishing a one to one correspondence between experiments and device simulations.

Because of the complexity of dealing with quantum transport at the lowest level of the hierarchy of Fig. 3 (Green's function method or direct solution of the *n*-body Schrödinger equation), and due to the desire to have device simulation tools which are able to deal with multiple levels of length scales and complexity, from the quantum regime to the classical regime, increasing interest is being focused at present on the use of quantum mechanically derived potentials that may be added as "corrections" to the semi-classical simulation tools. A way to include quantum effects into classical simulation tools is to add such quantum potentials to the mean field potential computed from Poisson's equation. In the past, such potential corrections have been employed mostly in the context of fluid approximations leading to the so-called quantum-hydrodynamic (QHD) equations [300], where the corresponding equations are usually derived under the assumption that the electron gas is near thermal equilibrium. Even so, they were expected to be more generally valid and allow one to simulate quantum effects in ultra-small scale semiconductor and nano-electronic devices. More recently, these quantum corrections are introduced as modifications of the Hartree potential obtained from solving the Poisson's equation. Also note that this concept of multiscale simulation, currently in the focus of the scientific research, is particularly critical in semiconductor devices where much of the device domain is in quasi-equilibrium and behaves classically (e.g., substrate, source-drain contacts, gate, etc.), whereas the critical regions governing the current are spatially small, subject to high fields and high degrees of potential confinement leading to quantum effects. It is, therefore, expected to be quite successful in overcoming some of the limitations of the semi-classical transport approaches discussed in Sections 3 and 5 of this review article. An in-depth description of the effective potential approach, utilized in particle-based simulations, is given in Section 6.1. In Section 6.2, we

give brief description of the quantum hydrodynamic model and its use in device simulations on the example of a high electron mobility modulation doped SiGe device structure.

## 6.1. The Effective Potential Approach

From a circuit modeling point of view, even the 1D solution of the Schrödinger-Poisson problem is a burdensome approach in terms of both complexity and computational cost. Because of this, it is common practice in industry to use analytical and macroscopic (in the sense of retaining the classical transport framework by adding correction terms to account for the quantum-mechanical effects) models that have provided some practical solutions. However, there are a number of problems associated with these approaches and all of them are directly related to the nonstationary nature of carrier transport (velocity overshoot) in deep submicrometer devices. Hence, more sophisticated models are needed that are able to capture the appropriate transport physics of the processes occurring in the smallest device sizes.

The idea of quantum potentials originates from the hydrodynamic formulation of quantum mechanics, first introduced by de Broglie and Madelung [301–303], and later developed by Bohm [304, 305]. In this picture, the wave function is written in complex form in terms of its amplitude $R(\mathbf{r}, t)$ and phase $\psi(\mathbf{r}, t) = R(\mathbf{r}, t)\exp[iS(\mathbf{r}, t)/\hbar]$. These are then substituted back into the Schrödinger equation to obtain the following coupled equations of motion for the density and phase

$$\frac{\partial \rho(\mathbf{r}, t)}{\partial t} + \nabla \cdot \left( \rho(\mathbf{r}, t)\frac{1}{m}\nabla S(\mathbf{r}, t) \right) = 0 \tag{291}$$

$$-\frac{\partial S(\mathbf{r}, t)}{\partial t} = \frac{1}{2m}[\nabla S(\mathbf{r}, t)]^2 + V(\mathbf{r}, t) + Q(\rho, \mathbf{r}, t) \tag{292}$$

where $\rho(\mathbf{r}, t) = R^2(\mathbf{r}, t)$ is the probability density. By identifying the velocity as $\mathbf{v} = \nabla S/m$, and the flux as $\mathbf{j} = \rho\mathbf{v}$, Eq. (291) becomes the continuity equation. Hence, Eqs. (291) and (292) arising from this so-called Madelung transformation to the Schrödinger equation have the form of classical hydrodynamic equations with the addition of an extra potential, often referred to as the *quantum*, or *Bohm potential*, written as

$$V_Q = -\frac{\hbar^2}{2mR}\nabla^2 R \rightarrow -\frac{\hbar^2}{2m\sqrt{n}}\nabla^2 \sqrt{n} \tag{293}$$

where the density, $n$, is related to the probability density as $n(\mathbf{r}, t) = N\rho(\mathbf{r}, t) = NR^2(\mathbf{r}, t)$, where $N$ is the total number of particles. The Bohm potential essentially represents a field through which the particle interacts with itself. It has been used, for example, in the study of wave packet tunneling through barriers [306], where the effect of the quantum potential is shown to lower or smoothen barriers, and hence *allow* for the particles to leak through.

An alternate form of the quantum potential was proposed by Iafrate, Grubin and Ferry [307], who derived a form of the quantum potential based on moments of the *Wigner-Boltzmann equation*, the kinetic equation describing the time evolution of the Wigner distribution function [308]. Their form, based on moments of the Wigner function in the pure state, and involving an expansion of order $O(\hbar^2)$, is given by

$$V_Q = -\frac{\hbar^2}{8m}\nabla^2(\ln n) \tag{294}$$

which is sometimes referred to as the *Wigner potential*, or as the density gradient correction. Such quantum potentials have been extensively used in *density-gradient* and *quantum-hydrodynamic* methods. Their use in particle-based simulation schemes becomes questionable due to the presence of statistical noise in the representation of the electron density and the considerable difficulty to calculate the second derivative of the density on a completely unstructured mesh given by the particle discretization.

To avoid this problem, Ferry and Zhou derived a form for a smooth quantum potential [309], based on the effective classical partition function of Feynman and Kleinert [310].

More recently, Gardner and Ringhofer [311] derived a smooth quantum potential for hydrodynamic modeling, valid to all orders of $\hbar^2$, which involves a smoothing integration of the classical potential over space and temperature. There, it was shown that, close to the equilibrium regime, the influence of the potential on the ensemble can be replaced by the classical influence of a smoothed non-local barrier potential. While this effective potential depends nonlocally on the density, it does not directly depend on its derivatives. Through this effective quantum potential, the influence of the barriers on an electron is felt at quite some distance from the barrier. The smoothed effective quantum potential has been used successfully in quantum-hydrodynamic simulations of resonant tunneling effects in one dimensional double-barrier structures [312].

In analogy to the smoothed potential representations discussed above for the quantum hydrodynamic models, it is desirable to define a smooth quantum potential for use in quantum particle-based simulations. Ferry [313] has suggested an *effective potential scheme* that emerges from a wave packet description of the particle motion. where the extent of the wave packet spread is obtained from the range of wavevectors in the thermal distribution function (characterized by an electron temperature). The effective potential, $V_{eff}$, is related to the self-consistent Hartree potential, $V$, obtained from the Poisson equation, through an integral smoothing relation

$$V_{eff}(\mathbf{x}) = \int V(\mathbf{x} + \mathbf{y})G(\mathbf{y}. a_0)d\mathbf{y} \tag{295}$$

where $G$ is a Gaussian with standard deviation $a_0$. The effective potential, $V_{eff}$, is then used to calculate the electric field that accelerates the carriers in the transport kernel of the Monte Carlo particle-based device simulator as discussed in Ref. [314]. The calculation of $V_{eff}$ has a fairly low computational cost. However, the explicit treatment of particles moving in the smoothed electric field at for example in the vicinity of the Si/SiO$_2$ interface (where the fields are the strongest) requires the use of a small time step less than 0.01 fs to eliminate artificial heating of the carriers, which adds to the computational cost. Note, also, that within this approach the parameter, $a_0$, has to be adjusted in the initial stages of the simulation via comparisons of the sheet/line density of the Q2D/Q1D structure being investigated using the effective potential approach and the 1D/2D Schrödinger-Poisson simulations.

The effective potential approach due to Ferry [314] is based on an Ansatz given by Eq. (295) of a Gaussian smoothing associated with the spatial extent of the wave packet characterized by an effective radius, $a_0$. While some guide to the choice of this parameter is given by effective classical partition function of Feynman and Kleinert [310], it is more typically chosen as a fit parameter in comparing to the self-consistent solution of the 1D Schrödinger-Poisson equation corresponding to the direction of confinement experience by carriers.

In order to put this on a more rigorous basis, an effective quantum potential for use in Monte Carlo device simulators is described below based on the Wigner-Boltzmann equation. This approach is based on perturbation theory around thermodynamic equilibrium and leads to an effective potential which depends on the energy and wavevector of each individual electron, thus effectively lowering step-function barriers for high-energy carriers [315]. The quantum potential is derived from the idea that the Wigner and the Boltzmann equation with the quantum corrected potential should possess the same steady state. The resultant quantum potential is, in general, two-degrees smoother than the original Coulomb and barrier potentials, i.e., possesses two more classical derivatives, which essentially eliminates the problem of statistical noise. The computation of the quantum potential involves only the evaluation of pseudodifferential operators and can, therefore, be effectively facilitated using Fast Fourier Transform (FFT) algorithms. The approach is quite general and can easily be modified to modeling of, for example. triangular quantum wells. The previously described approach has been used in simulation of 25-nm MOSFET device with oxide thickness of 1.2 nm and dual-gate device structures, as discussed in Section 6.1.2.

## 6.1.1. Thermodynamic Effective Potential

The basic idea of the thermodynamic approach to effective quantum potentials is that the resulting semiclassical transport picture should yield the correct thermalized equilibrium

quantum state. Using quantum potentials, one generally replaces the quantum Liouville equation

$$\partial_t \rho + \frac{i}{\hbar}[H, \rho] = 0 \tag{296}$$

for the density matrix, $\rho(x, y)$, by the classical Liouville equation

$$\partial_t f + \frac{\hbar}{2m^*}k \cdot \nabla_x f - \frac{1}{\hbar}\nabla_x V \cdot \nabla_k f = 0 \tag{297}$$

for the classical density function $f(x, k)$. Here, the relation between the density matrix and the density function, $f$, is given by the Weyl quantization

$$f(x, k) = W[\rho] = \int \rho\left(x + \frac{y}{2}, x - \frac{y}{2}\right)\exp(ik \cdot y)dy \tag{298}$$

The thermal equilibrium density matrix in the quantum mechanical setting is given by $\rho^{eq} = e^{-\beta H}$, where $\beta = 1/k_B T$ is the inverse energy, and the exponential is understood as a matrix exponential; that is, $\rho^{eq}(x, y) = \sum_\lambda \psi_\lambda(x)\exp(-\beta\lambda)\psi_\lambda(y)^*$ holds, with $\{\psi_\lambda\}$ the orthonormal eigen-system of the Hamiltonian, $H$. In the semiclassical transport picture, on the other hand, the thermodynamic equilibrium density function, $f_{eq}$, is given by the Maxwellian $f_{eq}(x, k) = \exp(-\frac{\beta\hbar^2|k|^2}{2m^*} - \beta V)$. Consequently, to obtain the quantum mechanically correct equilibrium states in the semiclassical Liouville equation with the effective quantum potential, $V^Q$, we set

$$f_{eq}(x, k) = \exp\left(-\frac{\beta\hbar^2|k|^2}{2m^*} - \beta V^Q\right) = W[\rho^{eq}] = \int e^{-\beta H}\rho\left(x + \frac{y}{2}, x - \frac{y}{2}\right)\exp(ik \cdot y)dy \tag{299}$$

This basic concept was originally introduced by Feynman and Kleinert [310]. Different forms of the effective quantum potential arise from different approaches to approximate the matrix exponential $e^{-\beta H}$.

In the approach presented in this paper, we represent $e^{-\beta H}$ as the Green's function of the semigroup generated by the exponential. Introducing an artificial dimensionless parameter, $\gamma$, and defining $\rho(x, y, \gamma) = \sum_\lambda \psi_\lambda(x)\exp(-\gamma\beta\lambda)\psi_\lambda(y)^*$, we obtain a heat equation for $\rho$ by differentiating $\rho$ w.r.t. $\gamma$ and using the eigenfunction property of the wave functions $\psi_\lambda$. This heat equation is referred to as the Bloch equation

$$\partial_\gamma \rho = -\frac{\beta}{2}(H \cdot \rho + \rho \cdot H), \quad \rho(x, y, \gamma = 0) = \delta(x - y) \tag{300}$$

and $\rho^{eq}(x, y)$ is given by $\rho(x, y, \gamma = 1)$. Under the Weyl quantization, this equation becomes, with the usual Hamiltonian $H = -\frac{\hbar^2}{2m^*}\Delta_x + V$ and defining the effective energy $E$ by $f = W[\rho] = e^{-\beta E}$

$$\partial_\gamma E = \frac{\beta\hbar^2}{8m^*}(\Delta_x E - \beta|\nabla_x E|^2)$$

$$+ \frac{\hbar^2|k|^2}{2m^*} + \frac{1}{2(2\pi)^3}\sum_{r=\pm 1}\int V\left(x + \frac{ry}{2}\right)\exp[\beta E(x, k, \gamma) - \beta E(x, q, \gamma) + iy(k - q)]dqdy,$$

$$E(x, k, \gamma = 0) = 0 \tag{301}$$

The effective quantum potential is in this formulation given by $E(x, k, \gamma = 1) = V^Q + \frac{\hbar^2|k|^2}{2m^*}$. The logarithmic Bloch equation is now solved 'asymptotically', using the *Born approximation*, i.e., by iteratively inverting the highest order differential operator (the Laplacian). This involves successive solution of a heat equation for which the Green's function is well known, giving (see Ref. [316] for the details)

$$V^Q(x, k) = \frac{1}{(2\pi)^3}\int \frac{2m^*}{\beta\hbar^2 k \cdot \xi}\sinh\left(\frac{\beta\hbar^2 k \cdot \xi}{2m^*}\right)\exp\left(-\frac{\beta\hbar^2}{8m^*}|\xi|^2\right)V(y)e^{i\xi(x-y)}dyd\xi \tag{302}$$

Note that the effective quantum potential, $V^Q$, now depends on the wave vector $k$. For electrons at rest; that is, for $k = 0$, the effective potential, $V^Q$, reduces to the Gaussian smoothing given in Ref. [313]. Also note that there are no fitting parameters in this approach; that is, the size of the wavepacket is determined by the particle's energy.

The potential $V(y)$ that appears in the integral of Eq. (302) can be represented as a sum of two potentials: the barrier potential, $V_B(x)$, which takes into account the discontinuity at the Si/SiO$_2$ interface due to the difference in the semiconductor and the oxide affinities, and the Hartree potential, $V_H(x)$, that results from the solution of the Poisson equation. Note that the barrier potential is 1D, independent of time, and needs to be computed only once in the initialization stage of the code. On the other hand, the Hartree potential is 2D and time dependent, as it describes the evolution of charge from quasi-equilibrium to a nonequilibrium state. Since the evaluation of the effective Hartree potential, as given by Eq. (302), is CPU intensive, approximate solution methods have been pursued to resolve this term within a certain level of error tolerance.

We recall from the above discussion that the barrier potential is just a step function. Under these circumstances, $e\nabla_x V_B(x) = B(1,0,0)^t \delta(x_1)$, where $B$ is the barrier height (on the order of 3.2 eV for the Si-SiO$_2$ interfaced) and $x_1$ is a vector perpendicular to the interface. We actually need only the gradient of the potential, so that, using the pseudodifferential operators, one obtains

$$\nabla_x V_B^Q(x, p) = \exp\left[\frac{\beta\hbar^2 |\nabla_x|^2}{8m^*}\right] \frac{2m^* \sin(\beta\hbar p \cdot \nabla_x / 2m^*)}{\beta\hbar p \cdot \nabla_x} \nabla_x V_B(x) \tag{303}$$

This equation gives

$$e\nabla_x V_B^Q(x, p) = \frac{B}{2\pi}(1,0,0)^t \int \exp\left[-\beta\frac{\hbar^2 |\xi_1|^2}{8m^*}\right] \frac{2m^* \sinh(\beta\hbar p_1 \cdot \xi_1 / 2m^*)}{\beta\hbar p_1 \cdot \xi_1} e^{i x_1 \cdot} d\xi_1 \tag{304}$$

Note that $V_B^Q$ is only a function of $(x_1, p_1)$, that is, it remains to be strictly one-dimensional, where $x_1$ and $p_1$ are the position and the momentum vector perpendicular to the interface. When combined with the fact that one has to calculate this integral only once, this fact motivates the tabulation of the result given by Eq. (304) on a mesh.

The Hartree potential, as computed by solving the $d$-dimensional Poisson equation, depends in general upon $d$ particle coordinates. For example, on a rectangular mesh, the 2D Hartree potential is given by $V_H(x_1, x_2, t)$, and one has to evaluate $V_H^Q(x_1, x_2, p_1, p_2, t)$ using Eq. (302) $N$ times each time step for all particles position and momenta: $x^n$, $p^n$, $n = 1, \ldots, N$ (where $N$ is the number of electrons, which is large). This task is, of course, impossible to accomplish in finite time on present state-of-the-art computers. The following scheme is therefore suggested. According to (302), one evaluates the quantum potential by multiplying the Hartree potential by a function of $\hbar\nabla_x$, or by multiplying the Fourier transform of the Hartree potential by a function of $\hbar\xi$. The expression in Eq. (302) may be factored into

$$V_H^Q(x, k) = \frac{2im^*}{\beta\hbar^2 k \cdot \nabla_x} \sinh\left(\frac{\beta\hbar^2 k \cdot \nabla_x}{2im^*}\right) \exp\left(\frac{\beta\hbar^2}{8m^*} |\nabla_x|^2\right) V_H(x)$$

$$= \frac{2im^*}{\beta\hbar^2 k \cdot \nabla_x} \sinh\left(\frac{\beta\hbar^2 k \cdot \nabla_x}{2im^*}\right) V_H^1(x) \tag{305}$$

with

$$V_H^1(x) = \exp\left(\frac{\beta\hbar^2}{8m^*} |\nabla_x|^2\right) V_H(x) \tag{306}$$

The evaluation of the potential, $V_H^1(x)$, which is a version of the Gaussian smoothed potential due to Ferry [313], is computationally inexpensive since it does not depend on the wave vector $k$. However, because of the Gaussian smoothing, $V_H^1(x)$ will be a smooth function of position, even if the Hartree potential $V_H(x)$ is computed via the Poisson equation where the electron density is given by a particle discretization. Therefore, the Fourier transform

of the potential, $V_H^0(x)$, will decay rapidly as a function of $\xi$, and it is admissible to use a Taylor expansion for small values of $\hbar\xi$ in the rest of the operator. This approximation gives

$$\frac{2im^*}{\beta\hbar^2 k \cdot \Gamma_1} \sinh\left(\frac{\beta\hbar^2 k \cdot \Gamma_1}{2im^*}\right) \approx 1 - \frac{\beta^2\hbar^4(k \cdot \Gamma_1)^2}{24(m^*)^2} \tag{307}$$

or

$$\partial_{x_i} V_H^Q(x'', p'') = \partial_{x_i} V_H^0(x'') - \frac{\beta^2\hbar^2}{24m^{*2}} \sum_{i,k=1}^{3} p_i'' p_k'' \partial_{x_i}\partial_{x_k}\partial_{x_i} V_H^0(x''), \quad n = 1, \ldots, N \tag{308}$$

for all particles. This computation is done by simple numerical differentiation of a sufficiently smooth grid function, $V_H^0$, and interpolation. The evaluation of (308) is the price we have to pay when comparing the computational cost of this approach as opposed to the original approach of Ferry [313], which uses simple forward, backward or centered difference scheme for the calculation of the electric field. The advantage is that this generalized effective potential approach avoids the use of adjustable parameters.

## 6.1.2. Quantum Effects in a Conventional Nanoscale MOSFET

As an application of the generalized effective potential approach detailed above, examples are given of particle based simulation of nanoscale gate-length MOSFETs including quantum effects through this method. The parameters of the device structure simulated are as follows: the average channel/substrate doping equals $10^{19}$ cm$^{-3}$, the doping of the source and drain regions is $10^{19}$ cm$^{-3}$, the junction depth is 30 nm, the oxide thickness is 1.2 nm and the gates are assumed to be metal gates with work function equal to the semiconductor electron affinity. The gate/channel length is 25 nm. Figure 46 shows the carrier confinement within the triangular potential well, with and without the inclusion of the quantum-mechanical size-quantization effects. These results are shown for the bias conditions $V_G = V_D = 1$ V. From these results, it is evident that the low-energy electrons are displaced little more than the high-energy electrons; the reason being the fact that the high-energy electrons tend to behave as classical particles and hence are displaced relatively less. Also note that there is practically no carrier heating for the case when the effective potential is used in calculating the driving electric field. The carrier displacement from the interface proper is also seen from the results presented in Figure 47. Notice that there is approximately 2-nm-average shift of the electron density distribution near the source end of the channel when quantization effects are included in the model. Also note that carriers behave more like bulk carriers at the drain end of the channel and are displaced in the same manner when using both the classical and the quantum-mechanical model.



Figure 46. Electron localization within the triangular potential barrier for the case when quantization effects are not included in the model (left panel) and for the case when we include quantum-mechanical space-quantization effects by using the effective potential approach presented in this paper (right panel). Reprinted with permission from S. Ahmed et al., *IEEE Trans. Nanotechnol.* 4, 465 (2005). © 2005, IEEE.

**Figure 47.** Electron distribution in the device without (left panel) and with (right panel) the incorporation of quantum-mechanical, size-quantization effects. Reprinted with permission from S. Ahmed et al., *IEEE Trans. Nanotechnol.* 4, 465 (2005). © 2005, IEEE.

The device transfer characteristics are shown in the left panel of Fig. 48. Again, it is clear that the full quantum potential and the barrier potential give similar values for the current. Looking more in detail at the device transfer characteristics, one finds that quantization effects lead to a threshold voltage increase of about 220 mV. When properly adjusted for the oxide thickness difference, this result is consistent with previously published data. Evidently, as deduced from the output characteristics shown in the right panel of Fig. 48, the shift in the threshold voltage leads to a decrease in the on-state current by 30%. The later observation confirms earlier findings that one must include quantum effects into the theoretical model to be able to properly predict the device threshold voltage and its on-state current.

Next, the simulation results of a 15 nm conventional $n$-channel MOSFET device are discussed. Similar devices have been fabricated by Intel Corporation. The physical gate length of the device used is 15 nm. The source/drain length equals 15 nm and the junction depth is also 15 nm. The bulk substrate thickness used for simulations is 45 nm. The height of the fabricated polysilicon gate electrode for this device is 25 nm. The gate oxide used was $SiO_2$ with physical thickness of only 0.8 nm. The source/drain doping density is $2 \times 10^{19}$ cm$^{-3}$ and the channel doping is $1.5 \times 10^{19}$ cm$^{-3}$. The substrate doping used is $1 \times 10^{18}$ cm$^{-3}$. The simulated device output characteristics are shown in Fig. 49.

There are again several noteworthy features in these results:

1. Quantum-mechanical size quantization increases the threshold voltage as observed from the decrease in the slope in the linear region, and hence degrades the device transconductance.

2. Drain current degradation due to quantum effects is not uniform, rather decreases with the increase in drain bias. The reason may be attributed again to the fact that the



**Figure 48.** Device transfer characteristic for $V_D = 0.1$ V (left panel) Device output characteristics VG = 1.0 V (right panel). Reprinted with permission from S. Ahmed et al., *IEEE Trans. Nanotechnol.* 4, 465 (2005). © 2005, IEEE.

Figure 49. (Left panel) Conventional 15-nm MOSFET device output characteristics. (Right panel) Average electron velocity along the channel. Reprinted with permission from S. S. Ahmed, Ph.D. Thesis, Arizona State University, 2005. © 2005.

electrons tend to behave as classical particles as average carrier energy increases with the increase in drain bias.

3. There is a considerable difference between the barrier-correction and the barrier-Hartree (full) correction which is mainly due to the use of higher doping density ($1.5 \times 10^{19}$ cm$^{-3}$) in the channel region than was used in the 25 nm MOSFET ($1 \times 10^{19}$ cm$^{-3}$) case.

The higher doping density has a direct impact on the Hartree potential, making the triangular channel potential steeper and hence introducing a pronounced quantum effect. But the overall degradation of the drain current as compared to the 25-nm MOSFET device structure is reduced in the 15-nm device because of the ballistic nature of the carrier motion in the latter case. This fact becomes clear if one observes the velocity profile of the device as depicted in the right panel of Fig. 49.

What is important in this figure is that the carriers attain a velocity, which is comparable to that in the 25-nm device structure even with a lesser biases applied; i.e., $V_G = V_D = 0.8V$. Also, the gate oxide thickness is less in the 10-nm device, which means that the gate oxide capacitance constitutes the major portion of the total effective gate capacitance thereby reducing the impact of the quantum capacitance.

4. The discrepancy between the experimental and the simulated results is attributed mainly to two reasons: (a) the series resistance coming from the finite width of the actual device structure and the contact resistances, and (b) the gate polysilicon depletion effects which as previously mentioned can introduce further degradation of the drain current on the order of 10–30% depending on the doping density and the height of the polysilicon gate used. The limited data as supplied by the Intel Corporation shows that the polysilicon gate is 25 nm in height, which can indeed contribute to a significant degradation of the drain current.

5. The use of a commercial simulator like the drift-diffusion based SILVACO Atlas fails to predict the device behavior mainly because of the ballistic and quantized nature of the carriers in these nanoscale device structures.

### 6.1.3. Size Quantization in Nanoscale SOI Devices

Size-quantization effects in nanoscale fully depleted silicon on insulator (SOI) devices arise due to the physical nature of the confined region, which is sandwiched between the two oxide layers. In order to verify the applicability of the quantum potential approach discussed above for this technology, a single gated SOI device structure is simulated. The SOI device used here has the following specifications: gate length is 40 nm, the source/drain length is 50 nm each, the gate oxide thickness is 7 nm with a 2 nm source/drain overlap, the buried oxide (BOX) layer thickness is 200 nm, the channel doping is uniform at $1 \times 10^{17}$ cm$^{-3}$, the doping of the source/drain regions equals $2 \times 10^{19}$ cm$^{-3}$, and the gate is assumed to be a metal gate with work function equal to the semiconductor electron affinity. There is a 10-nm

spacer region between the gate and the source/drain contacts. The silicon (SOI) film thickness is varied over a range of 1–10 nm for the different simulations that were performed to capture the trend in the variations of the device threshold voltage. Similar simulations were performed in Refs. [317, 318] using the Schrödinger-Poisson solver and Ferry's effective potential approaches, respectively. For comparison purposes, the threshold voltage is extracted from the channel inversion density versus gate bias profile and extrapolating the linear region of the characteristics to a zero value. This method also corresponds well to the linear extrapolation technique using the drain current-gate voltage characteristics.

The results showing the trend in the threshold voltage variation with respect to the SOI film thickness are depicted in Fig. 50. One can see that Ferry's effective potential approach overestimates the threshold voltage for a SOI thickness of 3 nm due to the use of rather approximate value for the standard deviation of the Gaussian wave packet which results in a reduced sheet electron density. As the silicon film thickness decreases, the resulting confining potential becomes more rectangular due to the combined effects of both the inversion layer quantization and the SOI film (physical) quantization, which also emphasizes the need for using a more realistic quantum-mechanical wavepacket description for the confined electrons. The favorable comparison of the generalized quantum potential to the exact Schrödinger-Poisson results indicate that this generalized approach can be applied to the simulation of SOI devices with a greater accuracy and predictive capability as it will be seen from the results presented in the next section.

## 6.1.4. Size Quantization in Nanoscale DG SOI Devices

As a final example of application of the effective potential method, a dual gate SOI (DG SOI) MOSFET device is considered. Such dual gate devices were originally designed to achieve the ITRS performance specifications for the year 2016. Figure 51 shows the simulated DG SOI device structure used in this section, which is similar to the devices reported in Ref. [319]. The dotted portion of the device, referred to as the *intrinsic* device, is the portion of the device in which quantum effects are taken into consideration.

The effective intrinsic device consists of two gate stacks (the gate contact and SiO$_2$ gate dielectric) above and below a thin silicon film. For the intrinsic device, the thickness of the silicon film is 3 nm. Use of a thicker body reduces the series resistance, and the effect of process variation but it also degrades the short channel effects (SCE). From the SCE point of view, a thinner body is preferable, but it is harder to fabricate very thin films of uniform thickness, and the same amount of process variation (±10%) may give intolerable fluctuations in the device characteristics. A thickness of 3 nm seems to be a reasonable compromise, but other body thicknesses have also been examined. The top and bottom gate insulator thickness is 1 nm, which is near the scaling limit for SiO$_2$. As for the gate contact, a metal gate with tunable work function, $\Phi_G$, is assumed, where $\Phi_G$ is adjusted to 4.188 to provide a specified off-current value of 4 $\mu A/\mu m$. The background doping of the silicon film is taken to be intrinsic; however, because of diffusion of the dopant ions, the



Figure 50. Threshold voltage variation with SOI film thickness. Reprinted with permission from S. S. Ahmed, Ph.D. Thesis, Arizona State University, 2005. © 2005.

$T_{ox} = 1$ nm          $T_{si} = 3$ nm
$L_G = 9$ nm             $L_T = 17$ nm
$L_{sd} = 10$ nm         $N_{sd} = 2 \times 10^{20}$ cm$^{-3}$
$N_b = 0$                $g = 1$ nm/decade
$\Phi_G = 4.188$         $V_G = 0.4$ V

**Figure 51.** DG device structure. Reprinted with permission from S. S. Ahmed, Ph.D. Thesis, Arizona State University, 2005. © 2005.

doping profile from the heavily doped S/D extensions to the intrinsic channel is graded with a coefficient of $g$ which equals to 1 nm/dec. For convenience, the doping scheme is also shown in Fig. 51. According to the ITRS roadmap, high performance devices should have a gate length of $L_G = 9$ nm by the year 2016. At this scale, both 2D electrostatic and quantum mechanical effects play an important role, and traditional device simulators may not provide reliable projections. The length, $L_T$, is an important design parameter in determining the on-current, while the gate metal work function, $\Phi_G$, directly controls the off-current. The doping gradient, $g$, affects both on-current and off-current.

The intrinsic device is simulated using the generalized effective potential approach in order to gauge the impact of size-quantization effects on DG SOI performance. The results are then compared to that from a full quantum approach based on the nonequilibrium Green's function (NEGF) formalism (NanoMOS-2.5) developed at Purdue University [320]. The NEGF method is discussed later in Section 7.7. In this method, scattering inside the intrinsic device is treated by a simple Buttiker probe model, which gives a phenomenological description of scattering, and is easy to implement under the Greens' function formalism. The simulated output characteristics are shown in Fig. 52. Devices with both 3 nm and 1 nm



**Figure 52.** Generic DG SOI device output characteristics. Reprinted with permission from S. S. Ahmed, Ph.D. Thesis, Arizona State University, 2005. © 2005.

channel thickness are simulated, with an applied gate bias of 0.4 V. The salient features of this figure are as follows:

1. Even with an undoped channel region, the devices achieve a significant improvement with respect to the short channel effects (SCEs) as depicted in flatness of the saturation region. This improvement is due to the use of the two gate electrodes and an ultrathin SOI film, which allow the gates more control on the channel charge.

2. Reducing the channel SOI film thickness to 1 nm further reduces SCEs, and improves device performance. However, the reduction in the drive current at higher drain bias is due to series resistance effects in the ultrathin body, which are naturally more pronounced for large drain currents.

3. Regarding quantum effects, one can see that quantum-mechanical size quantization does not play a major role in degrading the device drive current, primarily because of using an undoped channel region. Also, looking at the 3 nm (or 1 nm) case alone, one can see that the impact of quantization effects reduces as the drain voltage increases because of the growing bulk nature of the channel electrons.

4. The percentage reduction in the drain current is more pronounced in the 1 nm case throughout the range of applied drain bias because of the stronger physical confinement arising from the two SiO$_2$ layers sandwiching the silicon film.

5. Finally, the comparison between the quantum potential formalism and the NEGF approach for the device with 3-nm SOI film thickness shows reasonable agreement which further establishes the applicability of this method in the simulations of different technologically viable nanoscale classical and nonclassical MOSFET device structures.

## 6.2. Quantum Hydrodynamic Model

As mentioned earlier, the hydrodynamic model described in Section 5.2 arising from moments of the semiclassical Boltzmann transport equation can itself be corrected for quantum mechanical effects based on introduction of a quantum potential added to the electrostatic potential. The hydrodynamic equations which explicitly include quantum corrections and describe the particle conservation, momentum conservation and energy conservation (see Ref. [321] for a detailed discussion) are the following:

$$\frac{\partial n}{\partial t} + \nabla \cdot (n\mathbf{v}) = 0 \tag{309}$$

$$\frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} = -\frac{q\mathbf{E}}{m^*} - \frac{1}{nm^*}\nabla(nk_B T_q) - \frac{\mathbf{v}}{\tau_m} \tag{310}$$

$$\frac{\partial T}{\partial t} + \frac{1}{3\gamma}\mathbf{v} \cdot \nabla T_q = -\frac{2}{3\gamma}\nabla \cdot (\mathbf{v}T_q) + \frac{m^*\mathbf{v}^2}{3\gamma k_B}\left(\frac{2}{\tau_m} - \frac{1}{\tau_w}\right) - \frac{T - T_0}{\tau_w} \tag{311}$$

where $n$ is the average electron density, $\mathbf{v}$ is the average electron velocity, $T$ is the effective electron temperature, $m^*$ is the effective electron mass, $\mathbf{E}$ is the electric field, $\tau_m$ is the momentum relaxation time, $\tau_w$ is the energy relaxation time, and $T_q$ is given by

$$T_q = \gamma T + \frac{2}{3k_B}U_q \tag{312}$$

with

$$U_q = -\frac{\hbar^2}{8m^*}\nabla^2 \ln(n) \tag{313}$$

where $U_q$ is the quantum correction. The explicit quantum correction, as already discussed in Section 6.1, involves the second order space derivative of the log of the density. Hence, it tends to smooth the electron distribution, especially where the electron density has sharp changes. The factor, $\gamma$, is the degeneracy factor [322], given by

$$\gamma = \frac{F_{3/2}(\mu_f/k_B T)}{F_{1/2}(\mu_f/k_B T)} \tag{314}$$

Figure 53. (Left panel) Schematic description of the device structure under investigation. (Right panel) Simulated *I-V* characteristics. Reprinted with permission from J. R. Zhou et al., in " Proceedings of the International Workshop on Computational Electronics." OSU, 1993. © 1993, D. K. Ferry.

where $\mu_f$ is the Fermi energy measured from the conduction band edge, and is introduced as a correction to the total average electron kinetic energy

$$w = \frac{1}{2}m^*v^2 + \frac{3}{2}\gamma k_B T + U_q \qquad (315)$$

The relaxation times, $\tau_m$ and $\tau_w$, are functions of energy, and, as discussed in Section 3.3, are determined by fitting the homogeneous hydrodynamic equations to the velocity-field and energy-field relations from Monte Carlo simulations described in Section 3.1.

This quantum hydrodynamic model has been applied to a variety of device structures realized in different material systems. As an example, the investigation of transport in a 0.18 μm gate-length, modulation-doped structure, is discussed here as shown schematically in the right panel of Fig. 53. The doping of the top $Si_{0.7}Ge_{0.3}$ layer is 3.5 × 10¹⁸ cm⁻³, and a doping of 1 × 10¹⁴ cm⁻³ is used in the $Si_{0.7}Ge_{0.3}$ substrate. The lattice temperature in the simulation is taken to be 300 K. The typical simulation domain is 1 μm × 0.09 μm. The thickness of the top $Si_{0.7}Ge_{0.3}$ layer is 19 nm, and the strained-Si channel is 18 nm.

The simulated *I-V* characteristics for gate biases 0.7, 0.5, 0.2, and 0 V, respectively, are shown on the right panel of Fig. 53. The small thickness of the top $Si_{0.7}Ge_{0.3}$ layer provides a normally off device, since the Schottky barrier height of 0.9 eV leads to an estimated depletion width of 18.4 nm. The peak transconductance is about 300 mS/mm, and a good saturation with a drain conductance of 4.6 mS/mm is obtained for 0.5 V on the gate. Approximately the same current level and transconductance is found in a 0.25 μm device. These simulation results are comparable with corresponding experimental measurements. The relatively larger current level (0.3 μA/mm) and transconductance (330 mS/mm) found in the experiment is thought to be because of a higher sheet-charge density (2.5 × 10¹² cm⁻² compared with 1 × 10¹² cm⁻² in this simulation) in the quantum well for their particular modulation-doped structure. It is interesting to note that the transconductance of this device approaches the same order of magnitude as that of an AlGaAs/GaAs device with the same geometry, although the transconductance of the SiGe device is about three times smaller. The inclusion of quantum corrections leads to about 15% current increase for gate voltage of 0.5 V. By inspecting the electron density distribution along the channel region of the device (not shown here), one can see that this is due to the rapid change in the electron density at the gate end close to the drain contact within a region that is much shorter than the gate-length. The inclusion of the quantum potential also leads to increase of the electron density in the channel.

## 7. FULLY QUANTUM APPROACHES

As noted above, channel quantization in the direction normal to the oxide-semiconductor interface of MOSFETs has been a fact of life for many years. This leads to important modifications which are readily seen in smaller devices. Two such effects are a shift in the

threshold voltage, due to the rise of the lowest occupied subband above the conduction minimum, and a reduction in the gate capacitance, due to the setback of the maximum in the inversion density away from the interface. This latter produces a so-called quantum capacitance, which is effectively in series with the normal gate capacitance [323]. If these are the major effects produced by the quantization, then they can be readily handled in a normal semiclassical theory by the introduction of an effective potential, as was discussed in the previous section. However, if the individual quantum levels in the inversion layer become resolved, or if the lateral quantization (in either width or thickness of an SOI layer) becomes important, then a full quantum mechanical model is required to handle the device. In the following, we turn to the description of a full quantum mechanical simulation for ultrasmall devices, although we concentrate mostly upon the MOSFET.

There have been many suggestions for different quantum methods to model ultrasmall semiconductor devices [324–326]. However, in each of these approaches, the length and the depth were modeled rigorously as a two-dimensional simulation, while the third dimension (width) is usually included through the assumption that there is no interesting physics in this dimension (lateral homogeneity). Others have made the assumption that only a single mode (or a few modes) is important (discussed later) and that the mode does not change shape as it propagates from the source of the device to the drain of the device [327, 328]. These are not likely to be valid assumptions, especially as we approach devices whose width is comparable to the channel length, both of which may be less than 10 nm.

It is important to consider all the modes that may be excited in the source (or drain) region, as this is known to be responsible for some of the interesting physics that we wish to capture. In the source, the modes that are excited are three dimensional (3D) in nature, even in a thin SOI device. These modes are then propagated from the source to the channel, and the coupling among the various modes will be dependent upon the details of the total confining potential at each point along the channel. Moreover, as the doping and the Fermi level in short-channel MOSFETs increases, we can no longer assume that there is only one occupied subband. Indeed, it is well known that the quantization of the channel in the MOSFET is strongest at the source end, and is generally much weaker, or non-existent, at the drain end (the latter is true in the case of a pinched-off channel). Thus, the number of modes varies down the channel, and the coupling of the modes is not subject to an orthogonality rule, as the dimension over which the mode is defined varies with the position in the channel. This latter breaks the normal orthogonality relationship.

As interest in nanodevices extends beyond just the MOSFET, but to more general types of "device" structures which involve quantum transport, a general development of quantum transport theory will be given, and its application to such device-like structures. In the next section, we will deal with the general idea of transport between two contact regions, and how it may be viewed in terms of a modal transport. In the subsequent section, we will address the problems of the contacts themselves as well as some important constraints upon the solutions, and, in Section 7.3, we turn to a general treatment of a "device" embedded in its environment. This latter is quite important as quantum mechanics normally deals with *closed* systems. On the other hand, devices are *open* systems in which current and energy flow through the active region, but are measured in an external region, known as the environment. Thus, one must understand how quantum mechanics is applied to such open systems.

In Section 7.4, we develop the Schrödinger equation solutions for a self-consistent study of devices, primarily through the iterative scattering matrix approach. This will be applied to a simple two-dimensional structure and to a three-dimensional simulation of a small MOSFET. We then turn to a short discussion of density matrices and Wigner functions as they are applied to device simulation, before finally turning to the nonequilibrium Green's function.

In each of these approaches, we deal with quantization of a single electron in the region of interest. That is, we deal with the single particle wave function, rather than the full many-body wave function. Simulations still are described with the normal electron densities and Fermi functions, but the many-body interactions within the electron (or hole) gas are treated at best through a mean-field approximation, such as the local density approximation (LDA). Here, the details of the exchange and correlation within the electron (or hole) gas are approximated by a functional of the local density, which provides a correction to the

solutions obtained from the Poisson equation. The fact is that exact many-body wave functions are extremely complicated and require $O(N^5)$ partial wave functions in a Fock space just to diagonalize the Hamiltonian. Even with a simple quantum dot, only total numbers of electrons of the order of 10 have been solved [329, 330]. However, it is known that LDA approximations can still give relatively good results for the energy levels of these structures [331], although there are some limitations to this statement, which will be examined in Section 7.8.

## 7.1. General Conductivity and the Landauer Formula

The general treatment of transport through quantum systems is best described by a treatment of the conductivity, or conductance, of the overall device structure. In particular, for a structure with very small transverse dimensions, the conductance can be described in terms of a few transverse modes. In general, the conductivity in a homogeneous semiconductor can be written as

$$\sigma_d = \frac{n_d e^2 \tau_m}{m^*} \tag{316}$$

where $d (= 1, 2, 3)$ is the dimensionality of the semiconductor, $\tau_m$ is the momentum relaxation time, and the other symbols have their usual meaning. In a degenerate semiconductor, which is the normal case in MOSFETs and other quantum systems, the density is related to the Fermi energy or, more important, to the Fermi wave vector, as

$$n = \begin{cases} \dfrac{2k_F}{\pi}, & d = 1 \\[2mm] \dfrac{k_F^2}{2\pi}, & d = 2 \\[2mm] \dfrac{k_F^3}{3\pi^2}, & d = 3 \end{cases} \tag{317}$$

Hence, the conductivity is written as

$$\sigma = \frac{e^2}{\pi h} \cdot \begin{cases} 2v_F \tau_m, & d = 1 \\[2mm] v_F \tau_m \dfrac{k_F}{2}, & d = 2 \\[2mm] v_F \tau_m \dfrac{k_F^2}{3\pi}, & d = 3 \end{cases} \tag{318}$$

where $v_F$ is the Fermi velocity $(= \hbar k_F / m^*)$. If the conductance is now computed, using the transverse area as $W$ (the width) in two dimensions and $\pi W^2/4$ in three dimensions, we find

$$G = \frac{2e^2}{h} \cdot \begin{cases} \dfrac{2v_F \tau_m}{L}, & d = 1 \\[2mm] \dfrac{v_F \tau_m}{2L} \dfrac{k_F W}{2}, & d = 2 \\[2mm] \dfrac{v_F \tau_m}{3L}\left(\dfrac{k_F W}{2}\right)^2, & d = 3 \end{cases} \tag{319}$$

The factor involving the Fermi velocity arises only when the transit time is large compared to the relaxation time. That is, for ballistic transport, the term in $v_F \tau_m / L$ will disappear from (319), leaving just the universal unit of conductance

$$G_u = \frac{2e^2}{h} \tag{320}$$

and the number of transverse modes, given by the term in $k_F W$. Hence, this relationship, in the ballistic limit, reduces to the famous Landauer formula [332, 333]

$$G = \frac{2e^2}{h} N \tag{321}$$

where $N$ is the number of transverse modes. To be sure, the actual Landauer formula also includes the probability (or transmission) for each of those modes.

To see how the transmission enters into the picture, we begin not with the conductivity for a homogeneous semiconductor, but with the general form of the current for a "tunneling" structure. If an arbitrary barrier between a "left" contact and a "right" contact is considered, the current in terms of a left-originated current and a right-originated current can be developed. The left contact may be thought of as the source and the right contact thought of as the drain, as in a MOSFET. The general gradual channel approximation for the current can be rewritten as

$$
\begin{aligned}
I_S &= \frac{WC_0\mu}{L}\left(V_G - V_T - \frac{V_D}{2}\right)V_D \\
&= \frac{WC_0\mu}{2L}[(V_G - V_T)^2 - (V_G - V_T - V_D)^2] \\
&= I_{S,0} - I_{D,0}
\end{aligned}
\tag{322}
$$

That is, the channel current is a difference between a source-derived current and a drain-derived current, and saturation occurs when the latter vanishes (at $V_D = V_G - V_T$). In a tunneling structure, the left-derived current is written as [334]

$$
I_l = A(2e)\int \frac{d^3k}{(2\pi)^3}v_z(k_z)T(k_z)f_{FD}(E_L)
\tag{323}
$$

In a similar manner, the right-derived current is written as

$$
I_R = -A(2e)\int \frac{d^3k'}{(2\pi)^3}v'_z(k'_z)T(k'_z)f_{FD}(E_R)
\tag{324}
$$

In general, the two energy scales are separated by the applied voltage,

$$
E_L = E, \qquad E_R = E + eV_a
\tag{325}
$$

which leads to the conservation of the velocity through

$$
v_z dk_z = v'_z dk'_z
\tag{326}
$$

Hence, the total current may be written as

$$
I = A\frac{e}{\pi\hbar}\int \frac{d^2k_t}{(2\pi)^2}\int dE_z T(E_z)[f_{FD}(E) - f_{FD}(E + eV_a)]
\tag{327}
$$

Now, it is clear that the integration over the transverse momentum provides the number of transverse modes in a small, laterally quantized structure.

In the quantum limit, the transverse integration in (327) becomes a summation over the discrete states. It is also important to note that the energy in the Fermi-Dirac functions is the total energy, the sum of the $z$-component plus the transverse energy of the quantized mode. Thus, Eq. (327) may be written in the quantized limit as

$$
I = \frac{e}{\pi\hbar}\sum_{i,j}\int dE_z T_{ij}(E_z)[f_{FD}(E_z + E_i) - f_{FD}(E_z + E_j + eV_a)]
\tag{328}
$$

The sum is over the various transmissions from mode $i$ in the left contact to mode $j$ in the right contact. If is assumed that the applied voltage is vanishingly small, the Fermi function can be expanded, approximating the derivative as a delta function to obtain

$$
I = \frac{2e^2}{h}\sum_{i,j} T_{ij}
\tag{329}
$$

under the constraint that modes $i$, $j$ lie below the respective Fermi energies. If $i = j = N$, then Eq. (329) reduces to Eq. (321). Equation (329) is the more general form of the Landauer formula, but *it must be emphasized that the left contact and the right contact are*

*regions which are in equilibrium (or near equilibrium) with a well defined Fermi-Dirac distribution function.* The transmission must be computed from the far left, where this constraint is valid, to the far right, where this constraint is also valid. No approximations, which compute the transmission from only small regions within the solution space of the Poisson equation, can be made—they will yield incorrect results. This means that the transmission must be recomputed for each iteration of the Poisson equation in any self-consistent calculation! In fact, as will be seen later, the role of the shift in the Fermi energy in the contacts, which produces current-carrying contacts, must often be part of the self-consistent loop—the current itself is self-consistently varied during the solution iteration [335, 336].

## 7.2. Modes, Contacts, and Constraints on Solutions

In the foregoing, we have developed the relationship between conditions in left- and right-contact regions and the current that results from transmission through the arbitrary structure which lies between these two contact regions. The general result (328) connects the current to the quantum transmission between modes in the left-contact and modes in the right-contact. This assumes that the left and right contacts are of finite extent so that one may develop a mode structure within these contacts. This, in itself, is an approximation, as real contacts are usually metallic and contain a very high density of electrons. Thus, the assumption of modes in the contacts needs to be supported via a relationship that shows a connection to the more likely metallic behavior. In particular, it is a well known fact that a potential barrier arises between the source contact and the channel in a MOSFET. This barrier is not unlike the potential barrier that arises near any electron-emitting cathode. The nature of the barrier and how it responds to applied potentials largely determines the analytical behavior of the current in such a device. Hence, the conditions that are placed upon a "contact" in any device simulation may well determine the results of the simulation, and often can result in nonphysical output. The nature of the contact has long been known to be an experimental problem in measurements of transport; it has not been recognized that it is also a significant theoretical problem in computations of the transport.

Contacts have been discussed relative to their effect on transport for quite some time. One of the earliest important discussions dealt with the role that the contact played in non-equilibrium transport in semiconductors [337, 338]. If the region is equilibrated at an effective temperature, and carrying a current, the most general form of the "drifted" Fermi-Dirac function has been given by an expansion of the density matrix in generalized integrals of the motion developed by Zubarev [339] and by Fano [340]. In general, then, one expects to have the Fermi-Dirac show a shift due to the dynamical momentum arising from the current. In the case of the right contact, one may also expect to have an elevated temperature, similar to the "electron temperature" discussed previously. These integral invariants must be evaluated within the overall self-consistent loop of the device simulation.

But, there is more to the role of the "contact." In general, we expect the need for a quantum treatment of the transport to arise from coherence in the transport itself. That is, quantum interference effects can be expected to occur in the device, and these will be exhibited in the observable transport. Thus, one needs to fully incorporate this coherence. Yet, the connection of the device to the outside world is through the contacts. A device within a circuit must not allow for coherent effects to impact other devices, if the circuit is to be evaluated as a collection of discrete devices [341]. For this to occur, correlation functions, which exist within the active device, must be terminated by relaxation processes within the contact regions [342, 343]. The relaxation processes must decay faster than the buildup of correlations due to the current flow in the device, if the system is to reach a steady state. Thus, the contact must not only provide a smooth transition into the quantum system, but it must also provide the equilibration necessary for the coherence and the nonequilibrium nature of the carriers. In a modal situation, this will require that the number of modes in the contact be much larger than the number of modes in the active region of the device. This situation will also allow for excitation via evanescent modes of the transition into the active region, an effect well known to be important in modal transmission mismatch situations [4].

Finally, we need to consider again the role of the evanescent states at interfaces such as that between the contact and the active region of the device. People familiar with microwave

waveguides are aware that discontinuities give rise to the excitation of evanescent waves. The number and type of these determine inductive and capacitive "storage" effects at the discontinuity (or interface). It is similar in quantum "waveguides." Evanescent waves are excited at interfaces, and these correspond to capacitive and inductive effects or, more exactly, produce phase shifts in the quantum wave that are important in resonance effects. Tunneling itself is via the propagation of evanescent waves through a barrier. In the resonant tunneling diode (RTD), where two barriers are separated by a quantum well, the phase shifts introduced by the barriers are important in determining the exact energy level of the bound state in the quantum well. Since these phase shifts change as bias is applied to the RTD, the position of the energy level also changes with bias [344]. In the MOSFET, carriers get into the channel by excitation over, or tunneling through, a potential barrier. Hence, proper treatment of the evanescent waves in this region is quite important in ballistic transport and quantum resonances within the device. Even in the absence of the tunneling barrier, the excitation of evanescent waves at a discontinuity is very important in quantum transport as it changes the coupling of the various modes on either side of the discontinuity [4, 345]. Thus, it is important to consider an adequate number of evanescent modes at each discontinuity or interface if accurate results are to be obtained.

## 7.3. Separating the Device from Its Environment— Open Quantum Systems

The central feature of a small quantum device in which quantum transport is to be evaluated is the coupling of the device to its environment. As mentioned, one normally considers a quantum system as closed and isolated from any environment. One then solves the Schrödinger equation for the eigenstates of that system and no transport is considered. When transport is of interest, then the system (device) must be opened to its environment, and the interactions between the system and the environment determine a large part of the transport problem. Consider Fig. 54(a) for example. Here we consider a device embedded within its environment. The red "wall" totally isolates the system from its environment. The Hamiltonian for this system may be written as

$$H = H_s + H_e \tag{330}$$

where the first term describes the Hamiltonian of the system and the second term that of the environment. The density matrix may be written as a product of that for the system and that for the environment

$$\rho(\mathbf{r}, \mathbf{r}') = \langle \Psi^*(\mathbf{r}')\Psi(\mathbf{r}) \rangle = \rho_s \otimes \rho_e \tag{331}$$

where the symbol "$\otimes$" represents a tensor product. Hence, by tracing over the environmental variables, the system density matrix is immediately retrieved, and the Liouville equation for the system is merely that of the isolated entity

$$i\hbar\frac{\partial \rho_s}{\partial t} = [H_s, \rho_s] = \bar{H}_s\rho_s \tag{332}$$

The last form of Eq. (332) introduces the commutator-generating superoperator notation.



Figure 54. (a) A system embedded in an environment, but isolated from that environment. (b) An open system which interacts with its environment.

When the system is opened to the environment, as indicated in Fig. 54(b), there is an interaction between the system and its environment. The Hamiltonian is no longer the simple version of Eq. (330), but now includes this interaction as

$$H = H_s + H_e + H_{sd}$$  (333)

While one can still write the density matrix as a tensor product of the system and environment states, tracing out the environment states no longer yields a simple expression for the system evolution. Rather, we now have the more complicated situation of

$$i\hbar\frac{\partial\rho_s}{\partial t} = \hat{H}_s\rho_s + Tr_v[H_{sd}, \rho_s \otimes \rho_c]$$  (334)

Now, the evolution of the system can be completely governed by the last term, which represents the effect of the environment on the system (device). This results is both good and bad, as one wants external potentials to yield controllable currents, but, for quantum transport, it is necessary to know the nature of the quantum states within the system. *This is no trivial task.*

Open quantum systems interact with their environment, which usually contains the measurement system itself. The output of a measurement is presumed to be classical in nature. The manner in which the quantum properties of the system are revealed in the classical results of the measurement, as well as the manner in which these quantum properties evolve into intrinsic classical properties, have been the focus of investigation since the formulation of quantum theory. One interpretation, which explicitly includes the coupled systems, is that of *decoherence* [346]. Decoherence is thought to be an important part of the measurement process, especially in selecting the classical results; that is, in passing from the quantum states to the measured classical states of a system [347]. However, the description (and interpretation) of the decoherence process has varied widely. A key point is that the *interaction of the system upon the environment*, as well as the interaction of the environment upon the system, is important. Zurek has proposed that the interaction of the system on the environment leads to a preferred, *discrete* set of quantum states, known as *pointer* states, which remain robust, as their superposition with other states, and among themselves, is reduced by the decoherence process. In general, any quantum system interacts with external degrees of freedom, which are representative of the environment in which the quantum system is embedded. This interaction is important if the properties of the quantum system are to be observed in the environment. This interaction, though, causes decoherence, which results in the loss of "purity" of the states in the quantum system. Not all the states in the quantum system are equally susceptible to this, however, and there remains a smaller set of initial states that is relatively robust with respect to the interaction with the environment. These are the pointer states, and their existence is a universal property of the quantum system [348]. This decoherence-induced selection of the preferred pointer states was termed *einselection* by Zurek [347].

In addition to the pointer states, there exists a sea of states that are heavily damped by the decoherence process. It has been recently shown that. in an open quantum dot, most of the eigenstates (of the closed dot) are heavily coupled to the environment and are generally washed out. There remain, however, a set of states which are not washed out, as they are only very weakly coupled to the environment [349]. These pointer states are responsible for the quasi-periodic oscillatory behavior (with gate voltage or magnetic field) in these dots [350]. Classically, these dots are found to have a mixed phase space, which has a set of Kolmogorov-Arnold-Moser (KAM) islands surrounded by a sea of chaos [351]. Full self-consistent Poisson solutions coupled to quantum transport simulations have shown that the behavior of the states on the KAM islands (with voltage) is exactly that of the quantum pointer states [352]. Moreover, it has also been demonstrated that these pointer states have classical Poissonian statistics, whereas the totality of the states in the dots exhibit Gaussian statistics [353]. Consequently, it is clear that the pointer states, which remain observable in the open quantum system and are connected to the environment via phase-space tunneling, provide the transition to the classical regular orbits found in the device.

While the pointer states in the open quantum dot have been identified. we may easily conjecture that the quantum modes which remain in the channel of the MOSFET are the

equivalent pointer states in the latter system. An important point is that the pointer states compose a set of orthogonal basis vectors for the remaining quantization within the system (device). That is, the density matrix for the pointer states (which would be a portion of the system density matrix $\rho_s$) is diagonal with no connection to the washed out states, and only a weak connection to the environment. It has been shown in a generalized density matrix formulation that one can develop a projection superoperator for which the basis states are a diagonal set, just as for the pointer states [354]. In this formulation, the interaction between the environment and the system is characterized by a generalized memory term [355], as has been previously found by a variety of authors [356–358]. This may be described schematically as shown in Fig. 55. While the figure is expressed in terms of the less-than (correlation function) Green's function, which will be described later, the terms are directly interpretable in terms of Eq. (334). There remains a term arising from the first term on the right-hand side of this equation, which is the evolution of the system density matrix without any effect of the environment, although these states may be renormalized due to the interaction. Then, there are three terms which arise directly from the last term in Eq. (334). The first of these is a second-order process in which the system operates on the environment followed by an interaction of the environment on the system, as shown schematically in Fig. 55(b). Here, this provides a continuing entanglement of the system states with the environment. The next two are direct interactions of the environment upon the system, the first at the initial time (and therefore representing memory of the environments initial state) and the second at a later time, which includes evolution of the environment itself by these interactions. These terms yield an effective interaction of the environment on the system, which is a modification of the real environment. With the effective interaction, one can solve for an approximate solution for the system evolution [359].

It is clear that, while we can deal with open quantum systems, it is quite complicated due to the interaction with the environment. In classical transport, the environment might be the phonon bath which provides the scattering processes for the carriers, and this is included in a relatively simple and straightforward manner, as discussed in previous sections. Here, however, the environment plays a much larger effect here, as discussed above in terms of the contacts. This makes it much more difficult to properly include all of these effects, and is a general limitation on most treatments of quantum transport that have been presented to date. We will try to elucidate some of this in the following.

## 7.4. The Usuki Solution to the Schrödinger Equation

We now turn to numerical solutions of the Schrödinger equation for a device connected to two "reservoirs," which represent the contacts and the environment. Within this approach,



Figure 55. The four terms in the less-than Green's function for the system evolution. (a) The closed system-like evolution, which is independent of the environment. (b) The "entanglement" term, which is an interaction of the system on the environment followed by a reverse interaction upon the system. (c) and (d) Memory terms representing initial an delayed interactions of the environment on the system. Reprinted with permission from [355]. J. Knezevic and D. K. Ferry, Phys. Rev. E 67, 066122 (2003). © 2003, American Physical Society.

the environmental effects are treated only through the initial set of modes in the contacts and the Fermi-Dirac distribution for the occupancy of these modes, using Eq. (328) to describe the total current through the structure. As displayed in Fig. 56, the quantum mechanical problem can be solved either directly or via an iterative technique [360]. Usually, the number of discretized nodes needed to represent the Schrödinger equation in the solution domain is such that the latter approach is to be preferred, and this is the one which we will describe in this section.

One may write the Schrödinger equation in terms of a wave function $\psi_j$ for the slice $j$. Then the discretized version of the Schrödinger equation is obtained, keeping only terms up to first order in the approximation for the derivatives, as

$$(E_F \mathbf{I} - \mathbf{H}_j)\psi_j + \mathbf{H}_{j,j-1}\psi_{j-1} + \mathbf{H}_{j,j+1}\psi_{j+1} = 0 \tag{335}$$

Here, $\psi_j$ is an $M$-dimensional vector whose elements are described by the index $i$. The problem is solved on a square lattice whose grid spacing is $a$, with the wires extending $M$ lattice sites across in the $x$-direction and $N$ lattice sites in the $y$-direction. In this equation, the $\mathbf{H}_j$ matrices represent Hamiltonians for each slice, and contain the local potential at each grid point. The matrices $\mathbf{H}_{j,j\pm1}$ give the interslice coupling. By approximating the derivative, the kinetic energy terms are mapped onto a tight-binding model with the nearest-neighbor hopping energy given by

$$t = -\frac{\hbar^2}{2m^*a^2}, \quad \mathbf{H}_{j,j\pm1} = t\mathbf{I} \tag{336}$$

This equation can be used to derive a transfer matrix that allows us to translate across the system, evaluate the transmission coefficients, and evaluate the conductance using the Landauer formula. Transfer matrices, however, are notoriously unstable due to the exponentially growing evanescent waves, which must be included. This difficulty can be overcome by performing some clever matrix manipulations and calculating via the scattering matrices, commonly used in microwave waveguide approaches [361, 362].

### 7.4.1. Two-Dimensional Approaches

The wave function for each slice may now be written. That is, one can write the wave function as a vector

$$\psi_j = \begin{bmatrix} \varphi_j^M \\ \cdots \\ \varphi_j^2 \\ \varphi_j^1 \end{bmatrix} \tag{337}$$

where the superscript refers to the $x$-axis and the index $i$. There are $M$ values for the matrix, and this represents a wire of width $(M + 1)a$. From the nature of the discretization,



Figure 56. Geometry of a quantum dot embedded in a quantum wire which serves as contacts to the structure. The grid represents the underlying mesh on which the calculations are performed, though, in practice, the mesh is much finer. Reprinted with permission from [360], R. Akis et al., *Phys. Rev. B* 54, 17705 (1996). © 1996, American Physical Society.

it is automatically assumed that the wave function vanishes at $i = 0$ and $i = M + 1$. There are other important points that relate to the hopping energy. The discretization of the Schrödinger equation introduces an artificial band structure, due to the periodicity that this discretization introduces. As a result, the band structure in any one direction has a cosinusoidal variation with momentum eigenvalue (or mode index), and the total width of this band is $4t$. Hence, if one wants to properly simulate the real band behavior, which is quadratic in momentum, one needs to keep the energies of interest below a value where the cosinusoidal variation deviates significantly from the parabolic behavior desired. For practical purposes, this means that $E_{\text{max}} < t$.

With the discrete form of the Schrödinger equation defined, the scattering matrices relating adjacent slices in our solution space need to be obtained. The Hamiltonian matrix is defined by

$$H_i = \begin{bmatrix} V(N,j) - 4t & \cdots & t & 0 \\ \cdots & \cdots & \cdots & \cdots \\ t & \cdots & V(2,j) - 4t & t \\ 0 & \cdots & t & V(1,j) - 4t \end{bmatrix} \tag{338}$$

With this formulation of the matrices, the general procedure follows that of Usuki et al. [361]. One first solves the eigenvalue problem on slice 0 at the end of the source (away from the channel), which determines the propagating and evanescent modes for a given Fermi energy in this region. The wave function is thus written in a mode basis, but this is immediately transformed to the site basis, and one propagates from the drain end, using the scattering matrix iteration:

$$\begin{bmatrix} C_1(j+1) & C_2(j+1) \\ 0 & I \end{bmatrix} = \begin{bmatrix} 0 & I \\ -I & (t I)^{-1}(EI - H_i) \end{bmatrix} \times \begin{bmatrix} C_1(j) & C_2(j) \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ P_1(j) & P_2(j) \end{bmatrix} \tag{339}$$

The dimension of these matrices is $2M \times 2M$, but the effective propagation is handled by submatrix computations, through the fact that the second row of this equation sets the iteration conditions

$$C_2(j+1) = P_2(j) = [-C_2(j) + (t I)^{-1}(EI - H_i)]^{-1} \tag{340}$$
$$C_1(j+1) = P_1(j) = P_2(j)C_1(j)$$

At the source end, $C_1(0) = 1$, and $C_2(0) = 0$ are used as the initial conditions. These are now propagated to slice $N$, which is the end of the active region, and then onto the $N - 1$ slice. At this point, the inverse of the mode-to-site transformation matrix is applied to bring the solution back to the mode representation, so that the transmission coefficients of each mode can be computed. These are then summed to give the total transmission, and this result is used in a version of equation (328) to compute the current through the device (there is no integration over the transverse modes, only over the longitudinal density of states and energy).

If it is desired to incorporate a self-consistent potential within the device, Poisson's equation must additionally be solved. Here, the density at each point in the device is determined from the wave function squared magnitude at that point, and this result is used as the charge source in Poisson's equation. The solution for $C_1(N + 2)$ is the wave function at this point, and this solution is back-propagated using the recursion algorithm

$$\Phi_\xi^{N-2-i}(i,j) = P_1(j) + P_2\Phi_\xi^{N-i+1}(i,j) \tag{341}$$

Here, Eq. (341) is in the mode representation, and $\xi$ is the mode index. The density at any site $(i,j)$ is found by taking the sum over $\xi$ of the occupied modes at that site, as

$$n(i,j) = \sum_\xi |\Phi_\xi^{N-2-i}(i,j)|^2 \tag{342}$$

Within the group at Arizona State University (ASU) and many others, this iterative scattering matrix has been applied to simulate transport in a variety of systems, including the quantum dots discussed above, quantum wires with a quantum point contact, and so on. In many of these applications, a magnetic field has been present, and this modifies the hopping energy by a Peierl's phase factor, which is why the inverse matrices have been kept in the form shown in Eqs. (339) and (340), Ref. [360].

## 7.4.2. The Three Dimensional Case

We now present the extension of the Usuki recursive scattering matrix approach to a full 3D quantum simulation, which is being used at ASU to simulate short-channel, fully depleted SOI MOSFET devices [363, 364]. One major change in the notation is that the $x$-axis is taken as the direction along the channel, and in the formulation it is assumed that this direction is parallel to the (100) direction. This notation makes the conduction band of Si, with which we are primarily interested in the present section, to be composed of three sets of valleys. In practice, the actual direction down the channel is the [110] direction, but the present notation allows the most general situation. The $y$-direction is taken in the plane of the semiconductor-oxide interface, normal to the channel direction, which makes the $z$-direction normal to this interface. This choice of axes is most useful, as the resulting Hamiltonian matrix will be diagonal. In contrast, if we had chosen the [110] direction to lie along the channel, the six ellipsoids would have split into a twofold pair (those normal to the [100] plane) and a fourfold pair, but the Hamiltonian would not be diagonal since the current axis makes an angle with each ellipsoid of the fourfold pair. Using the present orientation complicates the wave function, as will be seen, but allows for simplicity in terms of the amount of memory needed to store the Hamiltonian and to construct the various scattering matrices (as well as the amount of computational time that is required). It is relatively simple to change to the [110] case, as has been done in a number of simulations [365].

We can now write a total wave function which is composed of three major parts, one for each of the three sets of valleys. That is, the wave function can be written as a vector

$$\Psi_j = \begin{bmatrix} \Psi^{(x)} \\ \Psi^{(y)} \\ \Psi^{(z)} \end{bmatrix} \tag{343}$$

where the superscript refers to the coordinate axis along which the principal axis of the ellipsoid lies (the longitudinal mass direction). Thus, $\Psi^{(x)}$ refers to the two ellipsoids oriented along the $x$ axis (the $\langle 100 \rangle$ ellipsoids). Each of these three component wave functions is a complicated wave function on its own. Consider the Schrödinger equation for one of these sets of valleys ($i$ corresponds to $x$, $y$, or $z$ valleys):

$$\frac{-\hbar^2}{2}\left(\frac{1}{m_x}\frac{d^2}{dx^2} + \frac{1}{m_y}\frac{d^2}{dy^2} + \frac{1}{m_z}\frac{d^2}{dz^2}\right)\Psi^{(i)} + V(x,y,z)\Psi^{(i)} = E\Psi^{(i)} \tag{344}$$

Here, it is assumed that the mass is constant, in order to simplify the equations (for non-parabolic bands, the reciprocal mass enters between the partial derivatives). The mass is labeled corresponding to the principle coordinate axes, and these take on the values of $m_l$ and $m_t$ as appropriate. One then implements this on a finite difference grid with uniform spacing $a$. The derivatives appearing in the discrete Schrödinger equation are replaced with finite difference representations of the derivatives. The Schrödinger equation then reads

$$-t_x(\psi_{i-1,j,k} + \psi_{i-1,j,k}) - t_y(\psi_{i,j+1,k} + \psi_{i,j-1,k}) - t_z(\psi_{i,j,k+1} + \psi_{i,j,k-1})$$
$$+ (V_{i,j,k} + 2t_x + 2t_y + 2t_z)\psi_{i,j,k} = E\psi_{i,j,k} \tag{345}$$

where $t_x$, $t_y$, and $t_z$ are the hopping energies

$$t_x = \frac{\hbar^2}{2m_x a^2}$$

$$l_x = \frac{\hbar^2}{2m_x a^2} \qquad (346)$$

$$l_z = \frac{\hbar^2}{2m_z a^2}$$

Each hopping energy corresponds with a specific direction in the silicon crystal. The fact that there are now three sets of hopping energies is quite important. The smallest value of $t$ corresponds to the longitudinal mass, and if energies are desired of the order of the source-drain bias ($\sim 1$ V), then one must have $a < 0.2$ nm. That is, one must take the grid size to be comparable to the Si lattice spacing!

With the discrete form of the Schrödinger equation defined. one can obtain the transfer matrices relating adjacent slices in the solution space. For this, the method is developed in terms of *planar* slices, and follows a procedure first put forward by Usuki et al. [361, 366]. This approach is modified here by the two dimensions in the transverse plane. The transverse plane has $N_y \times N_z$ grid points. Normally, this would produce a second-rank tensor (matrix) for the wave function. and it would propagate via a fourth-rank tensor. However, the coefficients can be re-ordered into a $N_y N_z \times 1$ first-rank tensor (vector), so that the propagation is handled by a simpler matrix multiplication. Since the smaller dimension is the $z$ direction, $N_z$ is used for the expansion. and the vector wave function is written as

$$\Psi^{(i)} = \begin{bmatrix} \psi^{(i)}_{1,N_y} \\ \psi^{(i)}_{2,N_y} \\ \cdots \\ \psi^{(i)}_{N_z,N_y} \end{bmatrix} \qquad (347)$$

Now. Eq. (345) can be rewritten as a matrix equation as, with $s$ an index of the distance along the $x$ direction.

$$H^{(i)}\Psi^{(i)}(s) - T_x^{(i)}\Psi^{(i)}(s-1) - T_x^{(i)}\Psi^{(i)}(s+1) = EI\Psi^{(i)}(s) \qquad (348)$$

Here. $I$ is the unit matrix, $E$ is the energy to be found from the eigenvalue equation, and

$$H^{(i)} = \begin{bmatrix} H_0^{(i)}(\mathbf{r}) & t_z^{(i)} & \cdots & 0 \\ t_z^{(i)} & H_0^{(i)}(\mathbf{r}) & \cdots & \cdots \\ \cdots & \cdots & \cdots & t_z^{(i)} \\ 0 & \cdots & t_z^{(i)} & H_0^{(i)}(\mathbf{r}) \end{bmatrix} \qquad (349)$$

and

$$T_x^{(i)} = \begin{bmatrix} t_x^{(i)} & 0 & \cdots & 0 \\ 0 & t_x^{(i)} & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & t_x^{(i)} \end{bmatrix} \qquad (350)$$

The dimension of these two super-matrices is $N_z \times N_z$, while the basic Hamiltonian terms of Eq. (349) have dimension $N_y \times N_y$. so that the total dimension of the above two matrices is $N_y N_z \times N_y N_z$. In general. if one takes $k$ and $j$ as indices along $y$. and $\eta$ and $\nu$ as indices along $z$. then

$$(t_x^{(i)})_{\eta\nu} = t_x^{(i)}\delta_{\eta\nu}. \qquad (t_x^{(i)})_{kj} = t_x^{(i)}\delta_{kj}. \qquad (t_z^{(i)})_{\nu\nu} = t_z^{(i)}\delta_{\nu\nu} \qquad (351)$$

and

$$
H_n^{(i)}(\mathbf{r}) =
\begin{bmatrix}
V(s,1,\eta) + W & t_x^{(i)} & \cdots & 0 \\
t_x^{(i)} & V(s,2,\eta) + W & \cdots & 0 \\
\cdots & \cdots & \cdots & t_x^{(i)} \\
0 & 0 & t_x^{(i)} & V(s,N_y,\eta) + W
\end{bmatrix}
\tag{352}
$$

The quantity $W$ is $2(t_x^{(i)} + t_y^{(i)} + t_z^{(i)})$, and is therefore independent of the valley index.

With this setup of the matrices, the general procedure follows that laid out in the two-dimensional case above. One first solves the eigenvalue problem on slice 0 at the end of the source (away from the channel), which determines the propagating and evanescent modes for a given Fermi energy in this region. The wave function is thus written in a mode basis, but this is immediately transformed to the site basis, and one propagates from the drain end, using the scattering matrix iteration

$$
\begin{bmatrix}
C_1^{(i)}(s+1) & C_2^{(i)}(s+1) \\
0 & 1
\end{bmatrix}
=
\begin{bmatrix}
0 & 1 \\
-1 & (T_x^{(i)})^{-1}(EI - H^{(i)})
\end{bmatrix}
$$

$$
\times
\begin{bmatrix}
C_1^{(i)}(s) & C_2^{(i)}(s) \\
0 & 1
\end{bmatrix}
\begin{bmatrix}
1 & 0 \\
P_1^{(i)}(s) & P_2^{(i)}(s)
\end{bmatrix}
\tag{353}
$$

The dimension of these matrices is $2N_y N_z \times 2N_y N_z$, but the effective propagation is handled by submatrix computations, through the fact that the second row of this equation sets the iteration conditions

$$
C_2^{(i)}(s+1) = P_2^{(i)}(s) = [-C_2^{(i)}(s) + (T_x^{(i)})^{-1}(EI - H^{(i)})]^{-1}
$$

$$
C_1^{(i)}(s+1) = P_1^{(i)}(s) = P_2^{(i)}(s)C_1^{(i)}(s)
\tag{354}
$$

At the source end, $C_1(0) = 1$, and $C_2(0) = 0$ are used as the initial conditions. These are now propagated to the $N_x$ slice, which is the end of the active region, and then onto the $N_z + 1$ slice. At this point, the inverse of the mode-to-site transformation matrix is applied to bring the solution back to the mode representation, so that the transmission coefficients of each mode can be computed. These are then summed to give the total transmission and this is used in a version of Eq. (328) to compute the current through the device (there is no integration over the transverse modes, only over the longitudinal density of states and energy).

As mentioned above, if one is to incorporate a self-consistent potential within the device, then Poisson's equation must be solved. Here, the density at each point in the device is determined from the wave function squared magnitude at that point, and this is used to drive Poisson's equation. The solution for $C_1(N_x + 2)$ is again the wave function at this point, and this is back-propagated using the recursion algorithm:

$$
\Phi_\xi^{(N_x+2,s,i)}(j,\eta) = P_1^{(i)}(s) + P_2^{(i)}\Phi_\xi^{(N_x+2,s+1,i)}(j,\eta)
\tag{355}
$$

Here, as before, the superscript "$i$" denotes the valley, while $j$ and $\eta$ denote the transverse position. Again, we are in the mode representation, and $\xi$ is the mode index. The density at any site $(s,j,\eta)$ is found by taking the sum over $\xi$ of the occupied modes at that site, as

$$
n(s,j,\eta) = \sum_\xi |\Phi_\xi^{(N_x+2,s,i)}(j,\eta)|^2
\tag{356}
$$

The devices considered here are trigate quantum wire SOI MOSFETs. In general, the source and drain regions are doped $6 \times 10^{19}$ cm$^{-3}$ $n$-type. The dimensions of the source and the drain are 18.47 nm wide, 10.32 nm long, and 6.51 nm high (integer multiples of the Si atom spacing) corresponding to the thickness of the silicon layer. The channel of the device is a $p$-type region which is left undoped. The channel region has dimensions identical to that

of the source and drain dimensions. A uniform 1-nm oxide layer covers the top and sides of the device to isolate the gates from the semiconductor.

Because there are only a finite number of actual dopants in the entire device, they are placed discretely on the lattice sites. To achieve this, the silicon lattice is scanned to distribute dopant atoms to the various regions of the device using the method presented in Ref. [367]. A typical potential profile is shown in Fig. 57, Ref. [365]. Note that the channel is undoped, although this does not eliminate fluctuations in threshold voltage either from device to device, or as the potentials are varied within a single device. The reason for this can be seen from the figure—the separation between the source and the channel is a random variable of lateral (and depth) position.

We now examine results obtained from the simulation of the SOI MOSFET, but with a lightly doped channel region. The $I_d-V_g$ curves obtained from six different dopant distributions are shown in Fig. 58(a), Ref. [368]. The spikes present in the plots give an excellent example of the quantum interference effects that occur in this system. When the discreteness of the doping is taken into account, the landscape of the potential is drastically altered [367]. The potential peaks now present in the channel set up additional reflections in the device, as well as forming resonant levels. These effects lead to the observance of the spikes in the curves shown in Fig. 58(a). Further, device to device variation in the threshold voltage is seen due to the differing dopant distributions. In Fig. 58(b), the $I_d-V_d$ curves are plotted. Additional peaks in the current give further evidence of the formation of resonant levels in the channel based on the position of the dopants in the channel. It can be concluded that this resonant behavior persists even at greatly elevated drain voltages. These results outline the impact of quantum interference and the existence of discrete dopant induced effects in SOI quantum wire MOSFETs.

## 7.5. The Density Matrix and the Wigner Function

The simulations discussed in Section 7.4 were based on calculation of the wave function itself, which arises from the discretized solution of the Schrödinger equation. There are other functions which are useful in treating quantum transport, and were already mentioned earlier—the density matrix and the less-than Green's function. In this section, and the next two, we will define and discuss these functions in more detail.

In general, any arbitrary time-dependent wave function may be expanded in a basis set, as

$$\Psi(\mathbf{r}, t) = \sum_n c_n \varphi_n(\mathbf{r}) e^{-iE_n t/\hbar} = \sum_n c_n \varphi_n(\mathbf{r}, t) \tag{357}$$



Figure 57. A typical self-consistent potential profile for a quantum wire trigate transistor. The localized donors in the source and drain form preferential sites in the potential. Reprinted with permission from [365], M. J. Gilbert and D. K. Ferry, IEEE Trans. Nanotechnol. 4, 355 (2005). © 2005, IEEE.

**Figure 58.** (a) $I_d$–$V_g$ curves for six different dopant distributions for $V_d = 10$ mV, exhibiting the effects of the quantum interference. The electron density interacting with the acceptors causes the observed spikes. (b) $I_d$–$V_d$ curves for the SOI MOSFET devices. From bottom to top the gate voltages are: 1.5, 2, 2.5, and 3 V. The peaks are evidence of the formation of resonant levels in the channel. Reprinted with permission from [365]. M. J. Gilbert and D. K. Ferry, *IEEE Trans. Nanotechnol.* 4, 355 (2005). © 2005, IEEE.

where a time-dependent basis set has been introduced in the last term. At time $t = 0$ (or an arbitrary initial time $t_0$), the coefficients can be expressed as

$$c_n = \int \varphi_n^*(\mathbf{r})\Psi(\mathbf{r}, 0)d^3\mathbf{r} \tag{358}$$

Then, the general solution can then be written as

$$\Psi(\mathbf{r}, t) = \sum_n \int \varphi_n^*(\mathbf{r})\Psi(\mathbf{r}, 0)d^3\mathbf{r}\varphi_n(\mathbf{r})e^{-iE_n t/\hbar}$$

$$= \sum_n \int \varphi_n^*(\mathbf{r})\varphi_n(\mathbf{r})e^{-iE_n t/\hbar}\Psi(\mathbf{r}, 0)d^3\mathbf{r} \tag{359}$$

In general, the initial time could be any initial time, as mentioned above, and the argument of the exponential would take on the time value $t - t_0$. We could then sum over all values of this initial time so that we could write Eq. (359) as

$$\Psi(\mathbf{r}, t) = \int_{t_0}^{t} K(\mathbf{r}, \mathbf{r}'; t, t_0)\Psi(\mathbf{r}, t_0)dt_0 \tag{360}$$

with

$$K(\mathbf{r}, \mathbf{r}'; t, t') = \sum_n \varphi_n^*(\mathbf{r}')\varphi_n(\mathbf{r})e^{iE_n(t-t')/\hbar}$$

$$= \sum_n \varphi_n^*(\mathbf{r}', t')\varphi_n(\mathbf{r}, t) \tag{361}$$

In the last line, the time-dependent basis functions defined in (357) have been reintroduced. The quantity $K$ is termed the *propagator kernel*, which is related to the Green's functions discussed in the next section. A more general form of Eq. (361) allows for the fact that the basis functions at different times may not be orthogonal to one another, so that the more general expression may be written as

$$K(\mathbf{r}, \mathbf{r}'; t, t') = \sum_{n,m} c_{nm}\varphi_n^*(\mathbf{r}')\varphi_m(\mathbf{r})e^{i(E_m t - E_n t')/\hbar}$$

$$= \sum_{n,m} c_{nm}\varphi_n^*(\mathbf{r}', t')\varphi_m(\mathbf{r}, t) \tag{362}$$

The equal time version of the propagator kernel is termed the density matrix [369], and is written as

$$\rho(\mathbf{r}, \mathbf{r}'; t) = \sum_{n, m} c_{nm} \varphi_n^*(\mathbf{r}', t) \varphi_m(\mathbf{r}, t) \tag{363}$$

This quantity is written in this reduced manner for convenience, but the density matrix is properly a square matrix whose dimension is defined by the size of the basis set (which is often infinite in many systems). Normalization requires that the trace of the matrix be unity:

$$Tr(\rho) = \sum_n c_{nn} = 1 \tag{364}$$

The density matrix is a particularly useful form in discussing entanglement between two subspaces (or subsets of wave functions). The density matrix has been used to study the changes arising from the quantum distribution in a high uniform, high electric field [370–373], and some ensemble Monte Carlo techniques have been developed to study the density matrix in this regime of transport [374, 375]. It has been used in numerical studies of quantum transport in resonant tunneling diodes [376, 377], and a general review of this has appeared earlier [378]. The density matrix has also been coupled to a Poisson solver to study small dual-gated MOSFETs [379].

The problem with the earlier approach using the wave functions, and with approaches that utilize the density matrix, is that one is limited to use of the position representation. Normally, in classical transport, the distribution function is also a function of carrier momentum (or energy). This has not been the case so far. However, one can introduce the center-of-mass and difference coordinates through

$$\mathbf{R} = \frac{\mathbf{r} + \mathbf{r}'}{2}, \quad \mathbf{s} = \mathbf{r} - \mathbf{r}' \tag{365}$$

and then Fourier transform on the difference coordinate to obtain the Wigner distribution function

$$f_w(\mathbf{R}, \mathbf{p}, t) = \frac{1}{h^3} \int d^3 s \rho(\mathbf{R}, \mathbf{s}, t) e^{i \mathbf{p} \cdot \mathbf{s} / \hbar} \tag{366}$$

When this is done, the equation of motion for the Wigner distribution function has a form reminiscent of the Boltzmann transport equation (in the absence of scattering):

$$\frac{\partial f_w}{\partial t} + \frac{1}{m} \mathbf{p} \cdot \nabla_R f_w - \frac{1}{h^3} \int d^3 P W(\mathbf{R}, \mathbf{P}) f_w(\mathbf{R}, \mathbf{p} + \mathbf{P}, t) = 0 \tag{367}$$

where

$$W(\mathbf{R}, \mathbf{P}) = \int d^3 q \sin\left(\frac{\mathbf{P} \cdot \mathbf{q}}{\hbar}\right) \left[ V\left(\mathbf{R} + \frac{\mathbf{q}}{2}\right) - V\left(\mathbf{R} - \frac{\mathbf{q}}{2}\right) \right] \tag{368}$$

is the nonlocal potential. The use of the Wigner function is particularly important in scattering problems [380], and it clearly shows the transition to the semiclassical world [381]. As with the density matrix, it has been used to study the resonant tunneling diode, and its use in a variety of transport venues has been reviewed recently [378].

More recent approaches have used variants of Monte Carlo techniques to study the scattering process with the Wigner function [382], and to study the resonant tunneling diode via particle techniques [383, 384]. In these simulations, it is necessary to discretize both real space and momentum space, and then to develop a technique to handle the phase interference that is not well treated by particle techniques. In work at ASU, an additional property of the particle, known as the *affinity* has been introduced, which relates to its effectiveness [385]. In Fig. 59, we plot the current-voltage curve for a resonant tunneling diode at 300 K that is obtained by this technique [383]. Here, the contact regions were doped to $10^{18}$ cm$^{-3}$, and the barriers were 0.3 eV, 3-nm barriers surrounding a 5-nm quantum well. A lightly

**Figure 59.** The current-voltage characteristics for an upsweep and a down sweep of the voltage. An additional sweep is shown to demonstrate the level of Monte Carlo noise in the computation. Reprinted with permission from [383], L. Shifren et al., *IEEE Trans. Electron Devices* 50, 769 (2003). © 2003, IEEE.

doped ($10^{16}$ cm$^{-3}$) region, of 30 nm length was placed on either side of the barrier structure, between the latter and the contacts. It is clear that the negative differential conductance is nicely reproduced with a peak-to-valley ratio that is comparable to that observed experimentally at this temperature.

In the above simulations, scattering by all the normal phonon processes present in GaAs were included. Normally, one does this with the scattering rates computed from the Fermi golden rule, just as done for semiclassical transport. Here, however, a technique has been included for evaluating the intracollisions field effect [387, 388]. In this approach, the similarities between the path integral form of the Wigner equation of motion and that for the retarded density matrix [389] has been used to evolve a numerical technique for inclusion of the intracollisional field effect [390]. This intracollisional field effect is a distinctly quantum effect that treats directly the temporal duration of the scattering process itself, the so-called collisional retardation.

## 7.6. The Recursive Green's Function

As with the direct solution of the wave function, one can approach the use of Green's functions in two ways. One can either write the entire Green's function for the system or one can develop a recursive approach which propagates from one end to the other. It is this latter approach which is discussed in this section, primarily as it is utilized for low temperature calculations of degenerate systems, where only transport at the Fermi surface is of interest. Here, we will limit ourselves to a two-dimensional system (the single Green's function for the entire system is reviewed in [140]). This recursive Green's seems to date from a one-dimensional discussion of Czycholl and Kramer [391], but it was popularized by Thouless and Kirkpatrick [392] and Lee and Fisher [393, 394]. Incorporation of a local potential, needed for self-consistency, has been reviewed by Kramer and Masek [395, 396]. This approach has been utilized to study fluctuations in semiconductor quantum wires by Baranger et al. [397].

In the previous section, we developed the propagator kernel for the wave function. In a similar manner, we can develop the *retarded* Green's function as [22]

$$G_r(\mathbf{r}, \mathbf{r}'; t, t') = -i\theta(t - t')\langle\{\psi(\mathbf{r}, t), \psi^+(\mathbf{r}', t')\}\rangle \qquad (369)$$

where the angle brackets have been added to symbolize an ensemble average at non-zero temperatures or a summation over the basis set and the curly brackets indicate the anti-commutator of the two field operators (wave functions in some sense). The *advanced* Green's function is given by

$$G_a(\mathbf{r}, \mathbf{r}'; t, t') = i\theta(t' - t)\langle\{\psi^+(\mathbf{r}', t'), \psi(\mathbf{r}, t)\}\rangle \tag{370}$$

so that the kernel can be rewritten in terms of these functions. The field operators themselves satisfy the anticommutation relationship

$$\{\psi(\mathbf{r}, t), \psi^+(\mathbf{r}', t')\} = \psi(\mathbf{r}, t)\psi^+(\mathbf{r}', t') + \psi^+(\mathbf{r}', t')\psi(\mathbf{r}, t)$$
$$= \delta(\mathbf{r} - \mathbf{r}')\delta(t - t') \tag{371}$$

In equilibrium (or very near to equilibrium), the Green's functions can be related to a basic Green's function $G_0$, which is a function only of the differences in the two positions and the two times, as expected for a homogeneous system. Then, this latter Green's function can be written as, in the Fourier transform representation,

$$G_0(\mathbf{k}, \omega) = \frac{\hbar}{\hbar\omega - E(\mathbf{k})} \tag{372}$$

which leads to the form for the retarded and advanced equilibrium functions

$$G_0^{r,a}(\mathbf{k}, \omega) = \frac{\hbar}{\hbar\omega - E(\mathbf{k}) \pm i\eta} \tag{373}$$

where $\eta$ is a small convergence factor that insures the correct time ordering for (361) and (370). Here, the upper sign corresponds to the first superscript, and this convention will be followed in the remainder of the section.

When there are interactions in the system, either between the electrons, or due to impurities or phonons, then the retarded and advanced functions can be found from Dyson's equation

$$G^{r,a}(\mathbf{k}, \omega) = G_0^{r,a}(\mathbf{k}, \omega) + G_0^{r,a}(\mathbf{k}, \omega)\frac{1}{\hbar}\Sigma^{r,a}(\mathbf{k}, \omega)G^{r,a}(\mathbf{k}, \omega) \tag{374}$$

*This equation is written in the Hartree approximation as a self-consistent Born approximation.* That is, the last two parts of the second term on the right is more properly written as a two-particle Green's function. Here, it has been expanded using perturbation theory with the perturbation terms resummed into the self energy $\Sigma^{r,a}$. One must examine each and every circumstance to be assured that this perturbation and re-summation is valid, especially in devices with strongly inhomogeneous charge and potential distributions. The more common form of (374) is

$$G^{r,a}(\mathbf{k}, \omega) = \frac{1}{[G_0^{r,a}(\mathbf{k}, \omega)]^{-1} - \frac{1}{\hbar}\Sigma^{r,a}(\mathbf{k}, \omega)} \tag{375}$$

It is Dyson's equation that allows one to develop the recursive Green's function for transport in a two-dimensional quantum "wire." Here, as in the recursive scattering matrix approach above, the Green's function is developed in terms of slices transverse to the current flow direction. The Green's function may be written in terms of evaluations at two longitudinal slices $j$ and $j'$ as [397] (We try to keep the same notation as for the recursive scattering matrix approach above.):

$$\langle j|G(ij, i'j'; E)|j'\rangle = G_{jj'}(i, i'; E) \tag{376}$$

In this equation, we have Fourier transformed on the time difference and rewritten the frequency as $E = \hbar\omega$, which is not the same as the kinetic energy, but the latter does not appear in the position representation. We are interested at an interior slice $q$, and the sections to the left and right of this point are represented by Green's functions $G^L$ and $G^R$. The Green's function of the coupled system is denoted $G^{L+R}$. In the wave function approach above, the coupling between sites was described in terms of a hopping term $t$.

Here, the interaction an energy, $V$, couples two adjacent slices. Hence, Dyson's equation may be written as

$$G_{pq}^{L,R} = G_{p,r-1}^{K} V_{r-1,r} G_{rq}^{L+R} \tag{377}$$

One can use this to build up sets of equations for the on-slice Green's functions and the slice-to-slice coupling Green's functions, which eventually leads to the two equations [140]

$$G_{rp}^{L-R} = G_{r,q+1}^{R} V_{q+1,q}(1 - G_{qq}^{L} \Sigma_{q}^{R}) G_{qp}^{L}$$
$$G_{pp}^{L+R} = G_{pp}^{L} + G_{p,q}^{L} \Sigma_{q}^{R}(1 - G_{qq}^{L} \Sigma_{q}^{R}) G_{qp}^{L} \tag{378}$$

with

$$\Sigma_{q}^{R,L} = V_{q,q-1} G_{q-1,q-1}^{R,L} V_{q-1,q} \tag{379}$$

In practice, the recursive Green's function approach treats a region that is coupled to two perfect quantum wires. The transverse modes of these wires are critical to the excitation of the active region and are used to compute the transmission from one side to the other. The coupling potential from one slice to the next is the hopping energy discussed before. The perfect wire is supposed to have a width $(M + 1)a$ as previously. The value of the wave function at site $i$ in the wire for mode $\xi$ is (which defines the mode to site transformation matrix)

$$U(i, \xi) \sim \sin\left(\frac{i\xi\pi}{M+1}\right) \tag{380}$$

which is unnormalized. The contact wire is assumed to have a node of the wave function at the point $j = -1$, the initial slice. If the mode is propagating, then the wave function has a sinusoidal variation away from this point (to the left). If the mode is evanescent, then it decays away from this point, which gives it an imaginary velocity. Now, the slice Hamiltonian matrix is set up as (338), including the Peierl's phase for the off diagonal terms if there is a magnetic field. Similarly, the hopping potential term is $tI$, with the Peierl's phase for the off diagonal terms if there is a magnetic field. The Green's function at slice zero is then

$$G_{00} = \frac{1}{H_0 + G_{ww} V_{w0}'} \tag{381}$$

Here, $G_{ww}$ is the adjusted site representation for the contact wire at the slice $j = -1$, and has the velocity of the mode included and $V_{w0} = V$ is the hopping potential coupling. Then, the recursion is built up with iterating on the slice index $j$ as

$$G_{jj} = [H_j - VG_{j-1,j-1}V^*]^{-1}$$
$$G_{0j} = G_{0,j-1}V^*G_{j,j}$$
$$G_{j0} = G_{jj}VG_{j-1,0}$$
$$G_{00} = G_{00} + G_{0j}VG_{j0} \tag{382}$$

This is then connected to the end boundary contact and the conversion back to modes made. The conductivity then comes from the mode to mode transmission, including the various velocities of each mode [398].

As examples of this approach, Takagaki has considered a number of structures, including crossed wires with a box resonator [399, 400], tunneling from a quantum point contact into an subsidiary channel [401], the role of disorder on quantum point contacts [402], and edge state interactions with a repulsive potential [403, 404]. In Fig. 60, we plot the case of a single antidote embedded in a quantum wire whose width is three times the size of the antidot. The potential height is twice the Fermi energy, and the Fermi wavelength is $6a$, where $a$ is the grid spacing (this allows fully normalized units).

**Figure 60.** Magnetoconductance of an antidot embedded in a quantum waveguide with a magnetic field applied. Reprinted with permission from [404], Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 48, 8152 (1993). © 1993, American Physical Society.

## 7.7. Nonequilibrium Green's Functions

In the discussion of the previous two sections, the system (device) was assumed to be in equilibrium with its environment, so that the distribution function was given by the Fermi-Dirac distribution. In the non-equilibrium case, this is no longer the case and we must now find the distribution function, exactly as was the situation for the semi-classical transport where we needed to solve the Boltzmann transport equation. In quantum transport, this entails two new Green's functions, the so-called correlation functions [405, 406]:

$$G^<(\mathbf{r}, t; \mathbf{r}', t') = i\langle \psi^+(\mathbf{r}', t')\psi(\mathbf{r}, t)\rangle$$
$$G^>(\mathbf{r}, t; \mathbf{r}', t') = -i\langle \psi(\mathbf{r}, t)\psi^+(\mathbf{r}', t')\rangle \tag{383}$$

These are called the *less-than* function and the *greater-than* function, respectively. Because of the dependence upon the wave functions, there are of course relationships between these Green's functions and the previous ones. These four form a complete set, although other combinations are possible [111]. As with the other Green's functions, there are equations of motion for these two correlation functions. The entire set is usually written together in a matrix form as [407]

$$G_K = \begin{bmatrix} G^r & G^< + G^> \\ 0 & G^a \end{bmatrix} \tag{384}$$

Under the same situation as the previous section, when the two-particle Green's function can be separated via a perturbation expansion and resumed into a self-energy, the equations of motion can be written as [408]

$$\left(i\hbar\frac{\partial}{\partial t} - H_0(\mathbf{r}) - V(\mathbf{r})\right)G_K = \hbar I + \Sigma_K G_K$$
$$\left(-i\hbar\frac{\partial}{\partial t'} - H_0(\mathbf{r}') - V(\mathbf{r}')\right)G_K = \hbar I + G_K \Sigma_K \tag{385}$$

Here, the self energy $\Sigma_K$ is a matrix having the same form as Eq. (384), and the product of the self-energy and the Green's function implies a convolution integration over an internal set of space and time variables. It is important to note, however, that this form neglects important correlations that may exist in the initial state of the system. Here, the initial state at which external potentials are applied may, in fact, be nonequilibrium situations such as those which exist in semiconductor devices. The carrier densities are likely to be very inhomogeneous and therefore significant built-in potentials will exist [409, 410]. In one form, the set of equations (385) is replaced with an equivalent set in which the potential is the

built-in potential, and the self-energy terms are replaced by the Hartree-Fock self-energy terms [409]

$$\Sigma_{HF}(\mathbf{r}, t; \mathbf{r}', t') = iV(\mathbf{r} - \mathbf{r}')G\ (\mathbf{r}, t; \mathbf{r}', t')\big|_{t=t'} \tag{386}$$

Even in this form, however, we are still ignoring important initial correlations (for when the external potentials are switched "on." These correlations can only be included via an additional set of Green's functions, which are not often expressible in terms of the set defined so far [411–415]. Many of these initial state problems arise from the failure of the Wick decomposition, which is crucial to the perturbation expansion and the formulation of the resummed self-energy. Nevertheless, in spite of these concerns, the non-equilibrium Green's functions have been applied to transport and devices.

Transport in homogeneous semiconductors subjected to high electric fields has been treated by Barker [416] and others [417, 418]. Datta has formulated a general treatment for device simulation, and has applied it to the resonant tunneling diode [419], even in the presence of optical phonon scattering [420], although the latter was treated in the local approximation rather than considering the temporal retardation effects (mentioned above). In fact, a relatively successful simulation package was developed for treating the resonant tunneling diode, including a more detailed energy band structure [421]. This has been extended to two-dimensional studies of MOSFETs [422], and to molecular structures connected between two metal or semiconductor contacts [423].

Other work has used the Green's function method to couple a quantum dot to two contact leads [424, 425]. In this approach, the dot is coupled to the leads via a parameter $\Gamma$. As was discussed above, in connection with tunneling, this parameter must be a complex quantity, since it must describe both the tunneling probability (a real part) and the phase shift (an imaginary part) which the transparent barrier produces at the edge of the dot. Without the latter, one cannot determine the proper energy levels within the dot, since this phase affects the latter just as in the resonant tunneling diode.

As remarked earlier, one usually utilizes a recursive algorithm, whether solving for the wave function or for the Green's functions. The reason for this is that the matrices for direct solvers (at all grid points simultaneously) are usually too large. A novel and efficient method has recently been introduced [426–430], which allows one to calculate the ballistic transport properties of a two- or three-dimensional device of arbitrary shape, potential profile, and number of leads. In this method, which is termed the contact block reduction (CBR) method, such quantities as the transmission function and the charge density of the open system can be obtained from the eigenstates of a corresponding closed system, that need to be calculated only once as the solution of a very small linear algebraic system. Importantly, the calculation of relatively few eigenstates of the closed system is believed to be sufficient to obtain very accurate results.

The first step of the CBR method consists of dividing the device coordinate space into two regions: the boundary region corresponding to the contact with the leads, $C$, and the region corresponding to the rest the device, $D$. In the ballistic case, the self-energy matrix, representing the coupling of the device to the leads, is non-zero only in the region $C$ and has the following structure

$$\Sigma = \begin{bmatrix} \Sigma_C & 0 \\ 0 & 0 \end{bmatrix} \tag{387}$$

The three submatrices are identically zero in this approach, so that only that of the contact-coupled region is required. In a similar manner, one can subdivide all the Green's function matrices as

$$G = \begin{bmatrix} G_C & G_{CD} \\ G_{DC} & G_D \end{bmatrix} \tag{388}$$

Importantly, one can rewrite the Dyson's equation (375) as

$$G = \frac{1}{G_0^{-1} - \Sigma} = [I - \Sigma]^{-1} G_0 \equiv A^{-1} G_0 \tag{389}$$

in order to introduce the inverse operator A. Then, the retarded Green's function matrix can be found to be

$$G^r = \begin{bmatrix} A_C^{-1}G_{0,C} & A_C^{-1}G_{0,C} \\ (G_{0,DC} - A_{DC}A_C^{-1}G_{0,C}) & (G_{0,D} - A_{DC}A_C^{-1}G_{0,CD}) \end{bmatrix} \tag{390}$$

The transmission function is determined by the elements of the following submatrix (of the small size) of the contacts

$$G_t^r = A_C^{-1}G_{0,C} \tag{391}$$

Any other quantity of interest, such as the density matrix and particle (charge) density, can be found with similar computational effort. The calculation of the Green's function of the decoupled device is made through its spectral representation. The eigenfunctions are obtained by solving the Schrödinger equation of the decoupled device only once, employing generalized von Neumann boundary conditions at the contacts. The use of these boundary conditions significantly reduces the required range of the eigenvalue spectrum of $H_0$ in the calculation of the transmission function, as only a few percent of the eigenvalues are usually sufficient to define the interaction with the contact. The remaining task of solving Dyson's equation for each energy value is also reduced to a manageable level.

As flexible as it is, the CBR method provides a very efficient way to solve the ballistic quantum transport problem in an open system. The CBR computational cost is mainly determined by the partial solution of the eigenvalue problem of a closed system. This method is capable of handling several connected leads precisely (i.e., not neglecting any off-diagonal elements of the self-energy as in the recursive Green's function method). In Fig. 61, the characteristic curves for a small double-gate MOSFET are shown, in which the dimensions are taken from the experimental device [431]. Here, the gates and oxide tunneling barriers are fully included in the self-consistent scheme. The band-structure is modeled by assuming parabolic bands with anisotropic effective masses for silicon and an isotropic effective mass for the oxide. Figure 62 depicts the charge density in the entire device region. Using the self-consistent CBR simulator, the calculation of a single bias point for the entire 2D device with more than 7000 grid points requires 2 to 3 hours on a regular 3 GHz Pentium 4 PC. The potential profile and density are shown in Fig. 62 for a case when the device is turned on.

## 7.8. Problems in Quantum Transport, Especially for the Future

In the discussion of quantum transport so far, it has been assumed that the energy band structure of the semiconductor remains that of the bulk material. In fact, in future ultrasmall devices, this may no longer be the case. Even when one couples, for example, a molecule to the a pair of leads, be they semiconductors or metals, usually the energy levels need to be determined of the molecule (or quantum wire or carbon nanotube), especially as the local configuration of the atoms may well have changed in this situation. One normally



Figure 61. Calculated transfer characteristics of ballistic DGFET showing the drain and (tunneling-induced) gate current for two source-drain voltages.

**Figure 62.** Potential profile (a) and electron density (b) for the double-gate MOSFET when the channel is fully formed at a gate voltage of $V_g = 0.4$ V.

approaches this via some form of a density-functional calculation (DFT). Normally, these are first-principles calculations, but with the exchange and correlation energy expressed in some form of a local density approximation (LDA). DFT calculations have some well-known problems. They are techniques for computing the ground state of the system. However, electrons in the conduction band represent excited states of the semiconductor, and the DFT approaches, in general, do not give good results for even the band gap of most semiconductors. To be sure, there have been some new approaches which have improved upon this situation, such as the GW approach [432] or the "exact exchange" approach [433], but the problem still remains.

One might think that metals would be easier, but there one has to deal with the surface potential. The fact that electrons do not fall out of the metal means that there is a barrier to this process, and this barrier is commonly called the work function. The presence of the work function is due to a charge dipole at the surface of the metal, where the electron cloud moves closer to the surface than the bulk atomic charge. DFT approaches, in general, do not model this potential correctly. The importance of this lies in the situation where one chemisorbs a molecule, such as a dithiol, onto the metal (such as gold). The chemical bonding of the sulpher to the gold results in a modification of the surface dipole in the gold and a different barrier between the bulk gold and the sulpher states. Without a good knowledge of the details of this barrier, one cannot use the above approaches to determine the tunneling through this interfacial barrier and the consequent current flow through the molecule. Consequently, it is not unexpected that simulations give orders of magnitude more expected current than is seen experimentally.

Approaches to the study of molecules have often just solved for the DFT properties of the molecule itself, ignoring the metal, and then applied the transport calculation treating the metal as a boundary condition [434]. While the results are interesting, leaving the interfacial potential out of the self-consistent loop is not the best approach. Nevertheless, these are important first steps to obtaining transport information, and will not be improved upon until adequate DFT calculations can be done to yield the interface potential to acceptable accuracy.

While this review has dealt with semiconductors principally, the idea of transport through small regions, such as molecules or quantum wires, is similar to problems in other arenas. One which has drawn interest of late is that of ion channels, where we have to deal with transport through a small channel that connects two reservoirs. Here, the interface potential is the lipid layer potential that exists on either side of the layer through which the ion channel passes. As in the above situation, no good approaches exist to calculate the value of this layer from first principles, and it is usually introduced via an assumed dipole charge layer [435]. Nevertheless, semiclassical approaches have made some inroad in studying the transport through these channels, but no quantum transport has been made to date.

## 8. CONCLUSIONS

To summarize this review, we have attempted to overview the field of Computational Nano-electronics with particular emphasis on the main numerical methods used in semiconductor nanodevice simulation, with examples taken primarily from our own research. The review is by no means inclusive of the extensive research in this field since the mid-1970s, which is continuing at an uninterrupted pace. The interested reader is referred to several of the books and review articles referenced in the present review.

In this review, we have tried to emphasize contemporary issues and techniques used in the simulation of both scaled semiconductor devices, and new devices such as quantum dots, res-onant tunneling diodes and quantum wave guides. The challenges of simulating increasingly smaller devices with more complicated geometries include the necessity of full 3D modeling, inclusion of atomistic effects in terms of discrete dopant profiles and other device inho-mogeneities, nonstationary/ballistic transport with proper treatment of both the long-range and the short-range particle-particle interactions, quantum mechanical interference effects, tunneling, quantization of motion, and nonlocal effects. The inclusion of all these effects comes at the cost of vastly increased computational burden, as one would expect. Fortu-nately, there has been a concurrent improvement in the performance in terms of speed and memory of the computational platforms available, based on the same technologies this field is attempting to ameliorate. The desktop computer system today has essentially the same computational performance of the worlds fastest supercomputer of only a decade ago. There is increasingly improvement in the performance of supercomputer clusters composed of net-worked arrays of microprocessors. Over the same time-frame, there has also been significant algorithmic speedup of many of the techniques discussed herein, which has greatly aided the development of this field. The prospect of having a fully "first-principles" approach to nanoelectronic modeling in which the input is ultimately the position of atoms themselves, seems less of a fantasy today with this continued exponential growth in computation.

## ACKNOWLEDGMENTS

## REFERENCES

1. G. Binnig and H. Rohrer, *Appl. Phys. Lett.* 40, 178 (1982).
2. International Technology Roadmap of Semiconductors, 2004. http://public.itrs.net/.
3. R. Chau, D. Doyle, M. Doczy, S. Datta, S. Hareland, B. Jin. J. Kavalieros, and M. R. Metz, in "Proceedings of the 61st Device Research Conference." Salt Lake City, 2003, p. 123.
4. D. K. Ferry and S. M. Goodnick, "Transport in Nanostructures." Cambridge University Press, Cambridge, 1997.
5. T. H. Ning, *IEEE 2000 Custom Circuits Conf.* 49 (2000).
6. M. Depas, B. Vermeire, P. W. Mertens, R. L. Van Meirhaegne, and M. M. Heyns, *Solid-State Electron* 38, 1465 (1995).
7. M. Leong, H.-S. Wong, E. Nowak, J. Kedzierski, and E. Jones, *ISQED* 492 (2002).
8. B. J. van Wees, H. van Houten, C. W. J. Beenakker, J. G. Williamson, L. P. Kouwenhoven, D. van der Marel, and C. T. Foxon, *Phys. Rev. Lett.* 60, 848 (1988).
9. D. A. Wharam, T. J. Thornton, R. Newbury, M. Pepper, H. Ahmed, J. E. F. Frost, D. G. Hasko, and D. C. Peacock, *J. Phys. C* 21, L209 (1988).
10. S. Washburn, in "Mesoscopic Phenomena in Solids" (B. L. Altshuler, P. A. Lee, and R. A. Webb, Eds.), pp. 1–36. North-Holland, Amsterdam, 1991.
11. R. E. Prange and S. M. Girvin, (Eds.), "The Quantum Hall Effect," 2nd Ed. Springer-Verlag, New York, 1990.
12. F. Sols, M. Macucci, U. Ravaioli, and K. Hess, *J. Appl. Phys.* 66, 3892 (1989).
13. S. Datta, *Superlatt. Microstruct.* 6, 83 (1989).
14. A. Weisshaar, J. Lary, S. M. Goodnick, and V. K. Tripathi, *IEEE Electron Device Lett.* 12, 2 (1991).
15. L. Worschech, B. Weidner, S. Reitzenstein, and A. Forchel, *Appl. Phys. Lett.* 78, 3325 (2001).
16. K. Hieke and M. Ulfward, *Phys. Rev. B* 62, 16727 (2000).
17. K. K. Likharev, *Proc. IEEE* 87, 606 (1999).

18. H. Grabert and M. H. Devoret. (Eds.), "Single Charge Tunneling, Coulomb Blockade Phenomena in Nanostructures." NATO ASI Series B 294. Plenum Press, New York, 1992.
19. K. Likharev, *IBM J. Res. Dev.* 32, 144 (1988).
20. L. J. Geerligs, V. F. Anderegg, P. A. M. Holweg, J. E. Mooij, H. Pothier, D. Esteve, C. Urbina, and M. H. Devoret, *Phys. Rev. Lett.* 54, 2691 (1990).
21. L. P. Kouwenhouven, A. T. Johnson, N. C. van der Vaart, C. J. P. M. Harmans, and C. T. Foxon, *Phys. Rev. Lett.* 67, 1626 (1991).
22. H. Pothier, P. Lafarge, C. Urbina, D. Esteve, and M. H. Devoret, *Europhys. Lett.* 17, 249 (1992).
23. D. H. Kim, S.-K. Sung, K. R. Kim, J. D. Lee, B.-G. Park, B. Ho Choi, S. W. Hwang, and D. Ahn, *IEEE Trans. Electron Devices* 49, 627 (2002).
24. C. Wasshuber, H. Kosina, and S. Selberherr, *IEEE Trans. CAD* 16, 937 (1997).
25. Y. Cui and C. M. Lieber, *Science* 291, 851 (2001).
26. R. Martel, V. Derycke, C. Lavoie, J. Appenzeller, K. K. Chan, J. Tersoff, and P. Avouris, *Phys. Rev. Lett.* 87, 256805 (2001).
27. M. A. Reed, J. N. Randall, R. J. Aggarwal, R. J. Matyi, T. M. Moore, and A. E. Wetsel, *Phys. Rev. Lett.* 60, 535 (1988).
28. L. Zhuang, L. Guo, and S. Y. Chou, *Appl. Phys. Lett.* 72, 1205 (1998).
29. D. H. Kim, S.-K. Sung, K. R. Kim, J. D. Lee, B.-G. Park, B. Ho Choi, S. W. Hwang, and D. Ahn, *IEEE Trans. Electron Devices* 49, 627 (2002).
30. K. Hiruma, M. Yazawa, T. Katsuyama, K. Haraguchi, M. Koguchi, and H. Kakibayashi, *J. Appl. Phys.* 77, 447 (1995).
31. H. Dai, E. W. Wong, Y. Z. Lu, S. Fan, and C. M. Lieber, *Nature (London)* 375, 769 (1995).
32. Y. Cui, X. Duan, J. Hu, and C. M. Lieber, *J. Phys. Chem. B* 104, 5213 (2000).
33. M. T. Björk, B. J. Ohlsoon, T. Sass, A. I. Persson, C. Thelander, M. H. Magnusson, K. Deppert, L. R. Wallenbeg, and L. Samuelson, *Appl. Phys. Lett.* 80, 1058 (2002).
34. M. T. Björk, B. J. Ohlsoon, T. Sass, A. I. Persson, C. Thelander, M. H. Magnusson, K. Deppert, L. R. Wallenbeg, and L. Samuelson, *Nano Lett.* 2, 87 (2002).
35. Y. Cui, Z. Zhong, D. Wang, W. U. Wang, and C. M. Lieber, *Nano Lett.* 3, 149 (2003).
36. X. Duan, C. Niu, V. Sahl, J. Chen, J. W. Parce, S. Empedocies, and J. L. Goldman, *Nature (London)* 425, 274 (2003).
37. M. T. Björk, B. J. Ohlsoon, C. Thelander, A. I. Persson, K. Deppert, L. R. Wallenbeg, and L. Samuelson, *Appl. Phys. Lett.* 81, 4458 (2002).
38. C. Thelander, T. Martensson, M. T. Björk, B. J. Ohlsson, M. W. Larsson, L. R. Wallenberg, and L. Samuelson, *Appl. Phys. Lett.* 83, 2052 (2003).
39. Z. Zhong, D. Wang, Y. Cui, M. W. Bockrath, and C. M. Lieber, *Science* 302, 1377 (2003).
40. M. S. Dresselhaus, G. Dresselhaus, and P. C. Eklund, "Science of Fullerenes and Carbon Nanotubes." Academic Press, Inc, New York. 1996.
41. See, for example P. L. McEuen, M. S. Fuhrer, and H. Park, *IEEE Trans. Nanotechnol.* 1, 78 (2002).
42. D. Scharfetter and H. Gummel, *IEEE Trans. Electron Devices* 16, 64 (1969).
43. R. Lake, G. Klimeck, R. C. Bowen, and D. Jovanovic, *J. Appl. Phys.* 81, 7845 (1997).
44. G. Baccarani and M. Wordeman, *Solid-State Electron.* 28, 407 (1985).
45. S. Cordier, *Math. Mod. Meth. Appl. Sci.* 4, 625 (1994).
46. P. Y. Yu and M. Cardona, "Fundamentals of Semiconductors." Springer, Berlin, 1999.
47. C. Herring, *Phys. Rev.* 57, 1169 (1940).
48. D. J. Chadi and M. L. Cohen, *Phys. Status Solidi B* 68, 405 (1975).
49. J. Luttinger and W. Kohn, *Phys. Rev.* 97, 869 (1955).
50. M. L. Cohen and T. K. Bergstresser, *Phys. Rev.* 141, 789 (1966).
51. J. R. Chelikowsky and M. L. Cohen, *Phys. Rev. B* 14, 556 (1976).
52. D. D. Awshalom, D. Loss, and N. Samarth, "Semiconductor Spintronics and Quantum Computation." Springer, Berlin, 2002.
53. S. Datta and B. Das, *Appl. Phys. Lett.* 56, 665 (1990).
54. Y. A. Bychkov and E. I. Rashba, *J. Phys. C* 17, 6039 (1984).
55. G. Dresselhaus, *Phys. Rev.* 100, 580 (1955).
56. J. Schliemann, J. Carlos Egues, and D. Loss, *Phys. Rev. Lett.* 90, 146801 (2003).
57. S. D. Ganichev, V. V. Bel'kov, L. E. Golub, E. L. Ivchenko, P. Schneider, S. Giglberger, J. Eroms, J. De Boeck, G. Borghs, W. Wegscheider, D. Weiss, and W. Prettl, *Phys. Rev. Lett.* 92, 256601 (2004).
58. E. Fermi, *Nuovo Cimento* 11, 157 (1934).
59. H. J. Hellman, *J. Chem. Phys.* 3, 61 (1935).
60. J. C. Phillips and L. Kleinman, *Phys. Rev.* 116, 287 (1959).
61. J. R. Chelikowsky and M. L. Cohen, *Phys Rev. B* 10, 12 (1974).
62. L. R. Saravia and D. Brust, *Phys. Rev.* 176, 915 (1968).
63. S. Gonzalez, Master's thesis, Arizona State University, 2001.
64. J. R. Chelikowsky and M. L. Cohen, *Phys. Rev. B* 10, 5059 (1974).
65. K. C. Padney and J. C. Phillips, *Phys. Rev. B* 9, 1552 (1974).
66. D. Brust, *Phys. Rev. B* 4, 3497 (1971).
67. S. Yamakawa, S. Aboud, M. Saraniti and S. M. Goodnick, *Comput. Electron.* 2, 481 (2003).
68. Z. H. Levine and St. G. Louie, *Phys. Rev. B* 25, 6310 (1982).

69. D. Fritsch, H. Schmidt, and M. Grundmann, *Phys. Rev. B* 67, 235205 (2003).
70. J. C. Slater and G. F. Coster, *Phys. Rev.* 94, 1498 (1954).
71. P.-O. Löwdin, *J. Chem. Phys.* 19, 1396 (1951).
72. E. O. Kane, *J. Phys. Chem. Solids* 1, 82 (1956).
73. E. O. Kane, *J. Phys. Chem. Solids* 1, 249 (1957).
74. J. M. Luttinger and W. Kohn, *Phys. Rev.* 97, 869 (955).
75. See www.wsi.tu-muenchen.de/nextnano3.
76. A. Cho, *J. Vac. Sci. Tech.* 8, S31 (1971).
77. A. Cho and J. Arthur, *Prog. Solid-State Chem.* 10, 157 (1975).
78. K. Von Klitzing, "In Lex Prixes Nobel." Almquist & Wiksell International, Stockholm, 1985. p. 56.
79. B. J. van Wees, H. van Houten, C. W. J. Beenaker, J. G. Williamson, L. P. Kouwenhoven, D. van de Marel, and C. T. Foxon, *Phys. Rev. Lett.* 60, 848 (1988).
80. D. Pfannkuche, R. R. Gerhardts, P. A. Maksym, and V. Gudmundsson, *Physica B* 189, 6 (1993).
81. E. Calleja, P. M. Mooney, S. L. Wright, and M. Heiblum, *Appl. Phys. Lett.* 49, 657 (1986).
82. M. J. Gilbert and J. P. Bird, *Appl. Phys. Lett.* 77, 1050 (2000).
83. A. Ashwin, Master's thesis, Arizona State University, 2005.
84. S. Das Sarma, J. Fabian, J. Xuedong Hu, and I. Zutic, *IEEE Trans. Magn.* 36, 2821 (2000).
85. I. Zutic, J. Fabian, and S. Das, *Rev. Mod. Phys.* 76, 323 (2004).
86. S. A. Wolf, D. D. Awschalom, R. A. Buhrman, J. M. Daughton, S. von Molnar, M. L. Roukes, A. Y. Chtchelkanova, and D. M. Treger, *Science* 294, 1488 (2001).
87. R. Eppenga, M. F. Schuurmans, and S. Colak, *Phys. Rev. B* 36, 1554 (1987).
88. W. Zawadzki and P. Pfeffer, *Semicond. Sci. Technol.* 19, R1 (2004).
89. R. Silsbee, *J. Phys.* 16, R179 (2004).
90. R. Winkler, *Phys. Rev. B* 62, 4245 (2000).
91. H. R. Trebin, U. Rossler, and R. Ranvaud, *Phys. Rev. B* 20, 686 (1979).
92. B. Tierney, unpublished.
93. Y. Chang and J. N. Schulman, *Phys. Rev. B* 31, 2069 (1985).
94. A. T. Meney, B. Gonul, and E. P. O'Reilly, *Phys. Rev. B* 50, 893 (1994).
95. S. Ben Radhia, K. Boujdaria, S. Ridene, H. Bouchriha, and G. Fishman, *J. Appl. Phys.* 94, 5726 (2003).
96. R. Winkler, *Phys. Rev. B* 62, 4245 (2000).
97. X. Cartoixa, D. Ting, and T. C. McGill, *Nanotechnology* 14, 308 (2003).
98. P. M. Garone, V. Venkataraman, and J. C. Sturm, *IEEE Electron Device Lett.* 13, 56 (1992).
99. N. Collaert, P. Verheyen, K. De Meyer, R. Loo, and M. Caymax, *IEEE Trans. Nanotechnol.* 1, 190 (2002).
100. R. Oberhuber, G. Zandler, and P. Vogl, *Phys. Rev. B* 58, 9941 (1998).
101. M. V. Fischetti, Z. Ren, P. M. Solomon, M. Yang, and K. Rim, *J. Appl. Phys.* 94, 1079 (2003).
102. http://nanohub.org/Online_Simulation_on_the_nanoHUB.
103. F. Stern, *Phys. Rev. Lett.* 30, 278 (1973).
104. D. K. Ferry, "Semiconductors." MacMillan, New York, 1991.
105. N. W. Ashcroft and N. D. Mermin, "Solid State Physics." W. B. Saunders, New York, 1976.
106. P. Hohenberg and W. Kohn, *Phys. Rev. B* 136, 864 (1964).
107. W. Kohn and L. J. Sham, *Phys. Rev. A* 140, 1133 (1965).
108. L. J. Sham and W. Kohn, *Phys. Rev.* 145, 561 (1966).
109. S. E. Koonin and D. C. Meredith, "Computational Physics." Addison-Wesley, New York, 1990.
110. J. C. Slater, *Phys. Rev.* 81, 385 (1951).
111. G. D. Mahan, "Many-Particle Physics." Plenum, New York, 1981.
112. W. Kohn and P. Batista, in "Theory of the Inhomogeneous Electron Gas" (S. Lundqvist and N. H. March, Eds.). Plenum Press, New York, 1983.
113. L. Hedin and B. I. Lundqvist, *J. Phys. C* 4, 2064 (1971).
114. O. Gunnarsson and B. I. Lundqvist, *Phys. Rev. B* 13, 4274 (1976).
115. I. K. Marmorkos and S. Das Sarma, *Phys. Rev. B* 48, 1544 (1993).
116. F. Stern and S. Das Sarma, *Phys. Rev. B* 30, 840 (1984).
117. E. Gawlinski, T. Dzurak, and R. A. Tahir-Kheli, *J. Appl. Phys.* 72, 3562 (1992).
118. U. von Barth, in "Electron Correlations in Solids, Molecules and Atoms" (J. T. Devreese and F. Brosens, Eds.). Plenum Press, New York, 1983.
119. N. D. Mermin, *Phys. Rev. A* 137, 1441 (1965).
120. U. Gupta and A. K. Rajagopal, *Phys. Rev. A* 21, 2064 (1980).
121. U. Gupta and A. K. Rajagopal, *Phys. Rev. A* 22, 2792 (1980).
122. D. Vasileska, P. Bordone, T. Eldridge, and D. K. Ferry, *J. Vac. Sci. Technol. B* 13, 1841 (1995).
123. B. Vinter, *Phys. Rev. B* 15, 3947 (1977).
124. P. Kneschaurek, A. Kamgar, and J. F. Koch, *Phys. Rev. B* 14, 1610 (1976).
125. T. Ando, *Z. Physik B* 26, 263 (1977).
126. W. A. Bloss, *J. Appl. Phys.* 66, 3639 (1989).
127. T. Ando, *Phys. Rev. B* 13, 3468 (1976).
128. S. Das Sarma and B. Vinter, *Phys. Rev. B* 23, 6832 (1981).
129. S. Das Sarma and B. Vinter, *Phys. Rev. B* 26, 960 (1982).
130. F. Schäffler and F. Koch, *Solid State Comm.* 37, 365 (1981).
131. J. S. Blakemore, *Solid-State Electron.* 25, 1067 (1982).

*132.* S. Luryi, *Appl. Phys. Lett.* 52, 501 (1988).

*133.* D. Vasileska, D. K. Schroder, and D. K. Ferry, *IEEE Trans. Electron Devices* 44, 584 (1997).

*134.* S. Takagi and A. Toriumi, *IEEE Trans. Electron Devices* 42, 2125 (1995).

*135.* M. J. van Dort, P. H. Woerlee, A. J. Walker, C. A. H. Juffermans, and H. Litka, *IEEE Trans. Electron Devices* 39, 932 (1992).

*136.* C. Jacoboni and L. Reggiani, *Rev. Mod. Phys.* 55, 645 (1983).

*137.* C. Jacoboni and P. Lugli, "The Monte Carlo Method for Semiconductor Device Simulation." Springer, Vienna, 1989.

*138.* K. Hess, "Monte Carlo Device Simulation: Full Band and Beyond." Kluwer Academic, Boston, 1991.

*139.* M. H. Kalos and P. A. Whitlock, "Monte Carlo Methods." Wiley, New York, 1986.

*140.* D. K. Ferry, "Semiconductors." Macmillan, New York, 1991.

*141.* H. D. Rees, *J. Phys. Chem. Solids* 30, 643 (1969).

*142.* R. M. Yorston, *J. Comp. Phys.* 64, 177 (1986).

*143.* L. I. Schiff, "Quantum Mechanics." McGraw-Hill, New York, 1955.

*144.* Y.-C. Chang, D. Z.-Y. Ting, J. Y. Tang, and K. Hess, *Appl. Phys. Lett.* 42, 76 (1983).

*145.* L. Reggiani, P. Lugli, and A. P. Jauho, *Phys. Rev. B* 36, 6602 (1987).

*146.* D. K. Ferry, A. M. Kriman, H. Hida, and S. Yamaguchi, *Phys. Rev. Lett.* 67, 633 (1991).

*147.* P. Bordone, D. Vasileska, and D. K. Ferry, *Phys. Rev. B* 53, 3846 (1996).

*148.* S. Bosi S and C. Jacoboni, *J. Phys. C* 9, 315 (1976).

*149.* P. Lugli and D. K. Ferry, *IEEE Trans. Electron Devices* 32, 2431 (1985).

*150.* N. Takenaka, M. Inoue, and Y. Inuishi, *J. Phys. Soc. Jpn.* 47, 861 (1979).

*151.* S. M. Goodnick and P. Lugli, *Phys. Rev. B* 37, 2578 (1988).

*152.* M. Moško, A. Mošková, and V. Cambel, *Phys. Rev. B* 51, 16860 (1995).

*153.* L. Rota, F. Rossi, S. M. Goodnick, P. Lugli, E. Molinari, and W. Porod, *Phys. Rev. B* 47, 1632 (1993).

*154.* R. Brunetti, C. Jacoboni, A. Matulionis, and V. Dienys, *Physica B/C* 134, 369 (1985).

*155.* P. Lugli and D. K. Ferry, *Phys. Rev. Lett.* 56, 1295 (1986).

*156.* J. F. Young and P. J. Kelly, *Phys. Rev. B* 47, 6316 (1993).

*157.* R. W. Hockney and J. W. Eastwood, "Computer Simulation Using Particles." Institute of Physics, Bristol, 1988.

*158.* D. J. Adams and G. S. Dubey, *J. Comput. Phys.* 72, 156 (1987).

*159.* Z. H. Levine and S. G. Louie, *Phys. Rev. B* 25, 6310 (1982).

*160.* L. V. Keldysh, *Z. Eksp. Teor. Fiz.* 37, 713 (1959).

*161.* N. Sano and A. Yoshii, *Phys. Rev. B* 45, 4171 (1992).

*162.* M. Stobbe, R. Redmer, and W. Schattke, *Phys. Rev. B* 47, 4494 (1994).

*163.* Y. Wang and K. Brennan, *J. Appl. Phys.* 71, 2736 (1992).

*164.* M. Reigrotzki, R. Redmer, N. Fitzer, S. M. Goodnick, M. Dür, and W. Schattke, *J. Appl. Phys.* 86, 4458 (1999).

*165.* M. V. Fischetti and S. E. Laux, *Phys. Rev. B* 38, 9721 (1988).

*166.* For a complete overview, see www.research.ibm.com/DAMOCLES.

*167.* M. Saraniti and S. M. Goodnick, *IEEE Trans. Electron Devices* 47, 1909 (2000).

*168.* T. Ando, A. B. Fowler, and F. Stern, *Rev. Mod. Phys.* 54, 437 (1982).

*169.* M. Abramovitz and I. A. Stegun, "Handbook of Mathematical Functions." U.S. Government Printing Office, Washington, D.C., 1964.

*170.* I. S. Gradshteyn and I. M. Ryzhik, "Table of Integrals, Series, and Products." Academic, San Diego, 1994.

*171.* F. F. Fang and W. E. Howard, *Phys. Rev. Lett.* 16, 797 (1966).

*172.* D. Vasileska, Ph.D. dissertation, Arizona State University, 1995.

*173.* E. J. Moore, *Phys. Rev.* 160, 607 (1967).

*174.* W. Kohn and J. M. Luttinger, *Phys. Rev.* 108, 590 (1957).

*175.* J. M. Luttinger, in "Mathematical Methods in Solid State and Superfluid Theory" (R. C. Clark and G. H. Derrick, Eds.), p. 157. Plenum Press, New York, 1967.

*176.* F. Stern and W. E. Howard, *Phys. Rev.* 163, 816 (1967).

*177.* J. D. Jackson, "Classical Electrodynamics." John Willey, New York, 1975.

*178.* J. J. Thomson, *Proc. Cambridge Philos. Soc.* 11, 1120 (1901).

*179.* K. Fuchs, *Proc. Cambridge Philos. Soc.* 34, 100 (1938).

*180.* E. H. Sondheimer, *Adv. Phys.* 1, 1 (1952).

*181.* E. Prange and T. Nee, *Phys. Rev.* 168, 779 (1968).

*182.* Z. Tesanovic, M. V. Jaric, and S. Maekawa, *Phys. Rev. Lett.* 57, 2760 (1986).

*183.* N. Trivedi and N. W. Ashcroft, *Phys. Rev. B* 38, 12298 (1988).

*184.* S. M. Goodnick, R. G. Gann, D. K. Ferry, and C. W. Wilmsen, *Surf. Sci.* 113, 233 (1982).

*185.* A. Gold, *Solid State Commun.* 60, 531 (1986).

*186.* A. Gold, *Phys. Rev. B* 35, 723 (1987).

*187.* G. Fishman and D. Calecki, *Phys. Rev. Lett.* 62, 1302 (1989).

*188.* H. Sakaki, T. Noda, K. Hirakawa, M. Tanaka, and T. Matsusue, *Appl. Phys. Lett.* 51, 1934 (1987).

*189.* S. M. Goodnick, D. K. Ferry, C. W. Wilmsen, Z. Liliental, D. Fathy, and L. Krivanek, *Phys. Rev. B* 32, 8171 (1985).

*190.* R. M. Feenstra, *Phys. Rev. Lett.* 72, 2749 (1994).

*191.* Y. Matsumoto and Y. Uemura, *Jpn. J. Appl. Phys. Suppl.* 2, 367 (1974).

*192.* T. Ando, *J. Phys. Soc. Jpn.* 43, 1616 (1977).

*193.* P. J. Price, *Ann. Phys.* 133, 217 (1981).

194  B. K. Ridley, *Rep. Prog. Phys.* 54, 169 (1991).
195  D. Roychoudhury and P. K. Basu, *Phys. Rev. B* 22, 6325 (1980).
196  C. Hao, J. Zimmermann, M. Charef, R. Fauquembergue, and E. Constant, *Solid-State Electron* 28, 733 (1985).
197  S. M. Goodnick and P. Lugli, *Phys. Rev. B* 37, 2578 (1988).
198  S. Imanaga and Y. Hayafuji, *J. Appl. Phys.* 70, 1522 (1991).
199  W. Magnus, *Solid-State Electron* 36, 843 (1993).
200  T. Yamada, Jing-Rong Zhou, H. Miyata, and D. K. Ferry, *IEEE Trans. Electron Devices* 41, 1513 (1994).
201  D. K. Ferry, *Phys. Rev. B* 14, 5364 (1976).
202  G. D. Mahan, "Many-Particle Physics." Plenum, New York, 1981.
203  J. M. Ziman, "Electrons and Phonons." Oxford University Press, London, 1960.
204  B. K. Ridley, "Quantum Processes in Semiconductors." Oxford University Press, New York, 1993.
205  F. Seitz, *Phys. Rev.* 73, 549 (1948).
206  W. A. Harrison, *Phys. Rev.* 104, 1281 (1956).
207  C. Herring and E. Vogt, *Phys. Rev.* 101, 944 (1956).
208  D. K. Ferry and C. Jacoboni, (Eds.), "Quantum Transport in Semiconductors." Plenum, New York, 1992.
209  H. J. G. Meyer, *Phys. Rev.* 112, 298 (1958).
210  H. A. Gómez de Cerdeira, *Solid State Commun.* 12, 511 (1973).
211  L. Pintschovius, J. A. Vergéz, and M. Cardona, *Phys. Rev. B* 26, 5658 (1982).
212  M. Lundstrom, "Fundamentals of Carrier Transport." Addison-Wesley, New York, 1992.
213  D. K. Ferry, *Surf. Sci.* 57, 218 (1976).
214  D. K. Ferry, *Phys. Rev. B* 14, 1605 (1976).
215  S. Zollner, S. Gopalan, and M. Cardona, *J. Appl. Phys.* 68, 1682 (1990).
216  D. Vasileska, P. Bordone, and D. K. Ferry, in "Proceedings of the Third International Workshop on Computational Electronics." Portland, Oregon, 1994.
217  D. Vasileska, P. Bordone, T. Eldridge, and D. K. Ferry, in "Quantum Transport in Ultrasmall Devices" (D. K. Ferry, H. L. Grubin, A.-P. Jauho, and C. Jacoboni, Eds.), p. 525. Plenum, New York, NY, 1995.
218  N. W. Ashcroft and N. D. Mermin, "Solid State Physics." W. B. Saunders, New York, 1976.
219  F. Stern, *Phys. Rev. Lett.* 18, 546 (1967).
220  E. D. Siggia and P. C. Kwok, *Phys. Rev. B* 2, 1024 (1970).
221  N. D. Mermin, *Phys. Rev. B* 1, 2362 (1970).
222  J. P. Walter and M. L. Cohen, *Phys. Rev. B* 5, 3101 (1972).
223  L. Wendler and R. Pechstedt, *Phys Status Solidi B* 138, 197 (1986).
224  Wei-Ye Chung and D. K. Ferry, *Solid-State Electron* 31, 1369 (1988).
225  H. Haug and C. Ell, *Phys. Rev. B* 46, 2126 (1992).
226  K. S. Yi, A. M. Kriman, and D. K. Ferry, *Semicond. Sci. Technol.* 7, B316 (1992).
227  H. Haug and S. W. Koch, "Quantum Theory of the Optical and Electronic Properties of Semiconductors." World Scientific, Singapore, 1990.
228  C. M. Snowden, "Introduction to Semiconductor Device Modeling." World Scientific, Singapore, 1986.
229  K. Bløtekjær, *IEEE Trans. Electron Devices* 17, 38 (1970).
230  K. Tomizawa, "Numerical Simulation of Submicron Semiconductor Devices." Artech House, Boston, 1993.
231  R. O. Grondin, S. M. El-Ghazaly, and S. M. Goodnick, *IEEE Trans. Microwave Theory Tech* 47, 817 (1999).
232  M. Grupen and K. Hess, *IEEE J. Quantum Electron.* 34, 120 (1998).
233  See www.ansoft.com/products/hf/hfss for details.
234  www.remcom.com
235  www.reland.com
236  See www.isc.ch/news/release_7.html.
237  A. Taflove, "Computational Electrodynamics: The Finite-Difference Time-Domain Method." Artech House, Boston, 1995.
238  Y. S. Yee, *IEEE Trans. Antennas Propag.* 14, 302 (1966).
239  W. L. Stutzman and G. A. Thiele, "Antenna Theory and Design." Wiley, New York, 1998, p. 493.
240  A. Taflove and M. E. Brodwin, *IEEE Trans. Microwave Theory Tech.* 23, 623 (1975).
241  G. Mur, *IEEE Trans. Electromagn. Compat.* 23, 1073 (1981).
242  T. G. Moore, F. G. Blaschak, A Taflove, and G. A. Kriegsmann, *IEEE Trans. Antennas Propag* 36, 1797 (1988).
243  J.-P. Berenger, *J. Comp. Phys.* 114, 185 (1994).
244  S. Selberherr, "Analysis and Simulation of Semiconductor Devices." Springer, New York, 1984.
245  G. Dahlquist and Å. Björck, "Numerical Methods." Prentice-Hall, Englewood Cliffs, New Jersey, 1974.
246  G. V. Gadiyak and M. S. Obrecht, "Simulation of Semiconductor Devices and Processes." *Proceedings Second International Conference* (1986), p. 147.
247  H. L. Stone, *SIAM J. Numer. Anal.* 5, 536 (1968).
248  W. Hackbush, "Multi-Grid Methods and Applications." Springer, Berlin, 1985.
249  P. Sonneveld, *SIAM J. Sci. Stat. Comput.* 10, 36 (1989).
250  Y. Saad, "Iterative Methods for Sparse Linear Systems." PWS Publishing, Boston, 1996.
251  H. A. Van der Vorst, *SIAM J. Sci. Stat. Comput.* 13, 631 (1992).
252  H. A. Van der Vorst, *SIAM J. Sci. Stat. Comput.* 10, 1174 (1989).
253  S. C. Eisenstat, *SIAM J. Sci. Stat. Comput.* 2, 1 (1981).
254  T. Gonzalez and D. Pardo, *Solid-State Electron.* 39, 555 (1996).
255  P. A. Blakey, S. S. Cherensky, and P. Sumer, "Physics of Submicron Structures." Plenum, New York, 1984.

256. T. Gonzalez and D. Pardo, *Solid-State Electron.* 39, 555 (1996).
257. S. S. Pennathur and S. M. Goodnick, *Inst. Phys. Conf. Ser.* 141 793 (1995).
258. R. W. Hockney and J. W. Eastwood, "Computer Simulation Using Particles," Institute of Physics, Bristol, 1988.
259. S. E. Laux, *IEEE Trans. Comp.-Aided Des. Int. Circ. Sys.* 15, 1266 (1996).
260. M. E. Kim, A. Das, and S. D. Senturia, *Phys. Rev. B* 18, 6890 (1978).
261. M. V. Fischetti and S. E. Laux, *Phys. Rev. B* 38, 9721 (1988).
262. P. Lugli and D. K. Ferry, *Phys. Rev. Lett.* 56, 1295 (1986).
263. A. M. Kriman, M. J. Kann, D. K. Ferry, and R. Joshi, *Phys. Rev. Lett.* 65, 1619 (1990).
264. R. P. Joshi and D. K. Ferry, *Phys. Rev. B* 43, 9734 (1991).
265. M. V. Fischetti and S. E. Laux, *J. Appl. Phys.* 78, 1058 (1995).
266. W. J. Gross, D. Vasileska, and D. K. Ferry, *IEEE Electron Device Lett.* 20, 463 (1999).
267. W. J. Gross, D. Vasileska, and D. K. Ferry, *VLSI Design* 10, 437 (2000).
268. D. Vasileska, W. J. Gross, and D. K. Ferry, *Superlattices Microstruct.* 27, 147 (2000).
269. W. J. Gross, D. Vasileska, and D. K. Ferry, *IEEE Trans. Electron Devices* 47, 1831 (2000).
270. K. Tomizawa, "Numerical Simulation of Submicron Semiconductor Devices," Artech House, Norwood, 1993.
271. H. Brooks, *Phys. Rev.* 83, 879 (1951).
272. C. Canali, G. Ottaviani, and A. Alberigi-Quaranta, *J. Phys. Chem. Solids* 32, 1707 (1971).
273. S. Beebe, F. Rotella, Z. Sahul, D. Yergeau, G. McKenna, L. So, Z. Yu, K. Wu, E. Kan, J. McVittie, and R. Dutton, *Int. Electron Devices Meeting* 213 (1994).
274. Medici, "Two-Dimensional Device Simulation Program," Version 1999.2, Avant! Corporation, Fremont, CA, 1999.
275. S. Selberherr, A. Schütz, and H. Pötzl, *IEEE Trans. Electron Devices* 27, 1540 (1980).
276. DESSIS-ISE, ISE TCAD Release 6.0, ISE Integrated Systems Engineering AG, Zürich, 1999.
277. ATLAS User's Manual, 6th Edn., Silvaco International, Santa Clara, CA, 1998.
278. A. Asenov, *IEEE Trans. Electron Devices* 45, 2505 (1998).
279. M. N. O. Sadiku, "Numerical Techniques in Electromagnetics," CRC Press, Boca Raton, FL, 1992.
280. M. Grupen, K. Hess, and L. Rota, "Physics and Simulation of Optoelectronic Devices III" (W. W. Chow and M. A. Osinski, Eds.), Vol. 292, p. 2399. SPIE, 1995.
281. M. Grupen and K. Hess, *Appl. Phys. Lett.* 65, 2454 (1994).
282. K. Hess, J. P. Leburton, and U. Ravaioli, (Eds.), "Hot Carriers in Semiconductors," Plenum, New York, 1996, p. 563.
283. K. E. Meyer, M. Pessot, G. Mourou, R. O. Grondin, and S. N. Chamoun, *Appl. Phys. Lett.* 53, 2254 (1988).
284. W. H. Knox, J. E. Henry, K. W. Goossen, K. D. Li, B. Tell, D. A. B. Miller, D. S. Chemla, A. C. Gossard, J. English, and S. Schmitt-Rink, *IEEE J. Quantum Electron.* 25, 2586 (1989).
285. S. M. El-Ghazaly, R. P. Joshi, and R. O. Grondin, *IEEE Trans. Microwave Theory Tech.* 38, 629 (1990).
286. S. M. Goodnick, S. Pennathur, U. Ranawake, P. Lenders, and V. Tripathi, *Int. J. Num. Model.* 8, 205 (1995).
287. P. R. Smith, D. H. Auston, and M. Nuss, *IEEE J. Quantum Electron.* 24, 255 (1988).
288. J. Son, W. Sha, T. Norris, J. Whitaker, and G. Mourou, *Appl. Phys. Lett.* 63, 923 (1993).
289. K. A. Remley, A. Weisshaar, V. K. Tripathi, and S. M. Goodnick, *VLSI Design* 8, 407 (1998).
290. R. Tsu and L. Esaki, *Appl. Phys. Lett.* 22, 562 (1973).
291. J. Faist, F. Capasso, D. L. Sivco, C. Sirtori, A. L. Hutchinson, and A. Y. Cho, *Science* 264, 553 (1994).
292. C.-J. Sheu and S.-L. Jang, *Solid-State Electron.* 44, 1819 (2000).
293. M. Städele, B. R. Tuttle, and K. Hess, *J. Appl. Phys.* 89, 348 (2001).
294. A. A. Demkov, X. Zhang, and D. A. Brabold, *Phys. Rev. B* 64, 125306/1-4 (2001).
295. J. P. Bird, R. Akis, D. K. Ferry, Y. Aoyagi and T. Sugano, *J. Phys.* 9, 5935 (1997).
296. D. Pfannkuche and R. R. Gerhardts, *Phys. Rev. B* 44, 13132 (1991).
297. C. S. Lent, P. D. Tougaw, W. Porod and G. H. Bernstein, *Nanotechnology* 4, 49 (1993).
298. C. S. Lent, P. D. Tougaw, and W. Porod, *Appl. Phys. Lett.* 62, 714 (1993).
299. A. O. Orlov, I. Amlani, G. H. Bernstein, C. S. Lent, and G. L. Snider, *Science* 277, 928 (1997).
300. C. L. Gardner, *SIAM J. Appl. Math.* 54, 409 (1994).
301. L. de Broglie, *C. R. Acad. Sci. Paris* 183, 447 (1926).
302. L. de Broglie, *C. R. Acad. Sci. Paris* 184, 273 (1927).
303. E. Madelung, *Z. Phys.* 40, 322 (1926).
304. D. Bohm, *Phys. Rev.* 85, 166 (1952).
305. D. Bohm, *Phys. Rev.* 85, 180 (1952).
306. C. Dewdney and B. J. Hiley, *Found. Phys.* 12, 27 (1982).
307. G. J. Iafrate, H. L. Grubin, and D. K. Ferry, *J. Phys.* 42, 307 (1981).
308. E. Wigner, *Phys. Rev.* 40, 749 (1932).
309. D. K. Ferry and J.-R. Zhou, *Phys. Rev. B* 48, 7944 (1993).
310. P. Feynman and H. Kleinert, *Phys. Rev. A* 34, 5080 (1986).
311. C. L. Gardner and C. Ringhofer, *Phys. Rev. E* 53, 157 (1996).
312. C. Ringhofer and C. L. Gardner, *VLSI Design* 8, 143 (1998).
313. D. K. Ferry, *Superlattices Microstruct.* 27, 61 (2000).
314. D. Vasileska and S. S. Ahmed, *IEEE Trans. Electron Devices* 52, 227 (2005).
315. C. Ringhofer, S. Ahmed, and D. Vasileska, *J. Comput. Electron.* 2, 113 (2003).
316. C. Ringhofer, C. Gardner, and D. Vasileska, *Int. J. High Speed Electron. Syst.* 13, 771 (2003).
317. Y. Omura, S. Horiguchi, M. Tabe, and K. Kishi, *IEEE Electron. Device Lett.* 14, 569 (1993).

382. C. Jacoboni, A. Bertoni, P. Bordone, and R. Brunetti, *Math. Comp. Simul.* 55, 67 (2001).
383. L. Shifren, C. Ringhofer, and D. K. Ferry, *IEEE Trans. Electron Devices* 50, 769 (2003).
384. H. Kosina, M. Nedjalkov, and S. Selberherr, *J. Comp. Electron.* 2, 147 (2003).
385. L. Shifren and D. K. Ferry, *Phys. Lett. A* 285, 217 (2001).
386. L. Shifren, private communication.
387. J. R. Barker, *J. Phys. C* 6, 2663 (1973).
388. D. K. Ferry, A. M. Kriman, H. Hida, and S. Yamaguchi, *Phys. Rev. Lett.* 67, 633 (1991).
389. J. R. Barker and D. K. Ferry, *Phys. Rev. Lett.* 42, 1779 (1979).
390. L. Shifren and D. K. Ferry, *Phys. Lett. A* 306, 332 (2002).
391. G. Czycholl and B. Kramer, *Solid State Commun.* 32, 945 (1979).
392. D. J. Thouless and S. Kirkpatrick, *J. Phys. C* 14, 235 (1981).
393. P. A. Lee and D. S. Fisher, *Phys. Rev. Lett.* 47, 882 (1981).
394. D. S. Fisher and P. A. Lee, *Phys. Rev. B* 23, 6851 (1981).
395. B. Kramer and J. Mašek, *J. Phys. C* 21, L1147 (1988).
396. J. Mašek and B. Kramer, *Z. Phys. B* 75, 37 (1989).
397. H. U. Baranger, D. P. DiVincenzo, R. A. Jalabert, and A. D. Stone, *Phys. Rev. B* 44, 1063 (1991).
398. R. Mezzener, Ph.D. thesis, Arizona State University, 1988.
399. Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 44, 8399 (1991).
400. Y. Takagaki and D. K. Ferry, *J. Appl. Phys.* 72, 5001 (1992).
401. Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 45, 12152 (1992).
402. Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 46, 15218 (1992).
403. Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 47, 9913 (1993).
404. Y. Takagaki and D. K. Ferry, *Phys. Rev. B* 48, 8152 (1993).
405. P. C. Martin and J. Schwinger, *Phys. Rev.* 115, 1342 (1959).
406. J. Schwinger, *J. Math. Phys.* 2, 407 (1961).
407. L. V. Keldysh, *Sov. Phys.—JETP* 20, 1018 (1965).
408. J. Rammer and H. Smith, *Rev. Mod. Phys.* 58, 323 (1986).
409. A. P. Jauho, in "Granular Nanoelectronics" (D. K. Ferry, J. R. Barker, and C. Jacoboni, Eds.), Plenum Press, New York, 1991.
410. D. K. Ferry, D. Vasileska, and H. L. Grubin, *Int. J. High Speed Electron.* 11, 363 (2001).
411. I. Prigogine and R. Balescu, *Physica* 25, 281 (1959).
412. P. Danielewicz, *Ann. Phys. (NY)* 152, 239 (1984).
413. D. Semkat, D. Kremp, and M. Bonitz, *J. Math. Phys.* 41, 7458 (200).
414. K. Morawetz, M. Bonitz, V. G. Morozov, G. Röpke, and D. Kremp, *Phys. Rev. E* 63, 020102 (2001).
415. I. Knezevic and D. K. Ferry, *Physica E* 19, 71 (2003).
416. J. R. Barker, *J. Phys. C* 7, 245 (1981).
417. R. Bertoncini, A. M. Kriman, and D. K. Ferry, *Phys. Rev. B* 41, 1390 (1990).
418. R. Bertoncini, A. M. Kriman, and D. K. Ferry, *Phys. Rev. B* 44, 3655 (1991).
419. S. Datta, *J. Phys.* 2, 8023 (1990).
420. R. Lake and S. Datta, *Phys. Rev. B* 45, 6670 (1992).
421. R. Lake, G. Klimeck, R. C. Bowen, and D. Jovanovic, *J. Appl. Phys.* 81, 7842 (1997).
422. R. Venugopal, Z. Ren, S. Datta, M. S. Lundstrom, and D. Jovanovic, *J. Appl. Phys.* 92, 3730 (2002).
423. M. P. Anantram, S. Datta, and V. Xue, *Phys. Rev. B* 61, 14219 (2000).
424. Y. Meir, N. S. Wingreen, and P. A. Lee, *Phys. Rev. Lett.* 66, 3048 (1991).
425. A. P. Jauho, N. S. Wingreen, and Y. Meir, *Phys. Rev. B* 50, 5528 (1994).
426. D. Mamaluy, M. Sabathil, and P. Vogl, *J. Appl. Phys.* 93, 4628 (2003).
427. M. Sabathil, D. Mamaluy, and P. Vogl, *Semicond. Sci. Technol.* 19, S137 (2004).
428. D. Mamaluy, A. Mannargudi, and D. Vasileska, *J. Comput. Electron.* 3, 45 (2004).
429. D. Mamaluy, D. Vasileska, M. Sabathil, and P. Vogl, *Semicond. Sci. Technol.* 19, S118 (2004).
430. D. Mamaluy, M. Sabathil, D. Vasileska, T. Zibold, and P. Vogl, *Phys. Rev. B* 71, 245321 (2005).
431. http://falcon.ecn.purdue.edu:8080/mosfet/10nmstructure.pdf.
432. L. Hedin and S. Lundqvist, *Solid State Phys.* 23, 1 (1969).
433. M. Städele, J. A. Majewski, P. Vogl, and A. Görling, *Phys. Rev. Lett.* 79, 2089 (1997).
434. G. C. Liang, A. W. Ghosh, M. Paulson, and S. Datta, *Phys. Rev. B* 69, 115302 (2004).
435. A. E. Cárdenas, R. D. Coalson, and M. Kurnikova, *Biophys. J.* 79, 80 (2000).

# CHAPTER 2

# Process Simulation for Silicon Nanoelectronic Devices

## Wolfgang Windl

*Department of Materials Science and Engineering, The Ohio State University, Columbus, Ohio, USA*

## CONTENTS

## 1. INTRODUCTION

### 1.1. Rise of Modeling Through Falling Dimensions

The silicon-based metal oxide semiconductor field effect transistor (MOSFET) is at the heart of today's semiconductor industry (Fig. 1) [1]. Gate, insulating oxide (between gate and channel), and channel form a capacitor. A voltage on the gate attracts a conductive layer of charge in the channel region, which connects source and drain—areas, where a high

concentration of "dopant" impurities injects charge carriers like electrons—electrically and switches a source-drain current on. Thus, the MOSFET acts as a switch, which can be turned on and off by the gate voltage. Individual switches can be connected to form the basic building blocks for circuit design, which are in turn used to form microprocessors and memory chips. Since the switching speed of MOSFETs increases with shrinking dimensions, which, in turn, also allows packing more MOSFETs on a chip of a given area, the semiconductor industry has managed to constantly improve the performance of computers by continuously scaling a more or less unchanged device geometry, resulting in performance doubling approximately every 18 months. By now, the industry has entered the nanoelectronics era, where nearly all the critical device dimensions including the feature size, which is the dimension of the smallest feature that is created on the wafer during patterning, are (in part significantly) below 100 nm with revolutionary changes necessary to continue further performance improvement. As an example, Fig. 2 shows the decrease in thickness of the gate oxide (between gate



Figure 2. Range of $SiO_2$ equivalent gate dielectric thickness ($T_{ox}$) as a function of year of first product shipment. Predictions for the future are taken from the 2003 International Technology Roadmap of Semiconductors [2]. The line is an exponential fit to guide the eye. Also shown are typical sizes for protein molecules and single atoms for comparison.

and channel, see Fig. 1) with time [2]. Current devices have gate oxides on the order of 1.2 nm.

The by now nanoscale dimensions of the traditional devices require several significant changes in the device engineering approach, since the increasing importance of single atoms and their random distributions as well as quantum effects visible in the device performance require a departure from the traditional picture, which is based on classical continuum theory. Extensive, time-consuming and costly engineering work is currently going on at semiconductor companies and in academia to further push the miniaturization process of semiconductor devices toward its physical limits and to ultimately define revolutionary new paradigms. In this situation, modeling and simulation has become an increasingly critical cost and especially time saving component of integrated-circuit technology development, provided it is accurate enough.

The physical limits will be reached for a number of device dimensions and materials properties within a few years. State-of-the-art transistors in industry are currently at the 90 nm node [3, 4], while transistors with gate lengths of 6 nm [5], comprising just a few dozen atoms, have been demonstrated. Although this represents a technological tour de force, it will be progressively difficult to continue downscaling at this rate, as increasing sheet resistance and solubility limits for dopants (which we will elaborate in Section 1.2) next to other factors such as gate oxide reliability, quantum tunneling, excessive power dissipation, and interconnect delays start hampering the performance of such devices [6, 7].

As the field of electronic devices ventures further into the nanometer regime, the current strategy consists of two parts. For the near future, the evolutionary "top-down" approach is continued, where traditional device structures are mostly kept. The MOSFET scaling problems are addressed by introducing new materials in places where the traditional ones fail [8] and by improving device design, packaging, and processing. In parallel, "bottom-up" synthetic chemical approaches of assembling nanodevices and circuits directly from their molecular constituents [9] are explored.

Independent of the chosen architecture or materials, one of the grandest challenges of electronic devices at the nanoscale is their fabrication or processing. In conventional semiconductor technology, processing is performed by only six different process steps, diffusion, photolithography, etch, ion implantation, thin film deposition and polishing (Fig. 3) [10]. While this technology has successfully been applied through the history of traditional semiconductor device fabrication and is transferable in principle to many nanoelectronic devices, especially the rapidly increasing cost of fabrication motivates exploration of entirely new paradigms.

In the following, we will discuss as our focus example ultrashallow junction formation, a prominent end-of-the-roadmap problem from the front-end (before the first metallization step) of the fabrication in more detail, which has pioneered to a large fraction the development of new and improved modeling techniques that have taken process simulation from a purely phenomenological continuum approach to more and more physics-based atomic-scale methodologies suitable for nanoelectronics devices. This will be complemented by a discussion of molecular electronics devices in the next section (including especially the use of carbon nanotubes as channel materials in field-effect transistors), which have largely driven a similar development in the field of electron transport (or "device") simulation by establishing an atom-based framework as addition or replacement for continuum-based approaches.

## 1.2. The Ultrashallow-Junction Problem

Device miniaturization does not simply mean that the channel length is shortened and the electrons, traveling with the same drift velocity as before, can make their way from source to drain faster through a shorter channel for a MOSFET device such as the one depicted in Fig. 1. Instead, device scaling requires that all vertical and lateral dimensions of the transistor be scaled [6]. As an example, Fig. 4 shows that the electrical insulation between source and drain in the "off-state" (i.e., with a voltage applied between source and drain, but none at the gate) is much better with a shallow source and drain. This can be seen from the simulated plotted electrical potential values for two MOSFET structures with the only difference being the depth of source and drain (usually referred to as "junction depth," as described next). In

Oxidation        Photoresist      Mask-Wafer         Exposed        Photoresist
(Field oxide)    Coating     Alignment and Exposure  Photoresist    Develop

Oxide            Photoresist      Oxidation          Polysilicon    Polysilicon
Etch             Strip            (Gate oxide)       Deposition     Mask and Etch

Ion              Active           Nitride            Contact        Metal
Implantation     Regions          Deposition         Etch           Deposition and
                                                                    Etch

**Figure 3.** Fabrication steps for top-down lithography-based processing from [10]. Reprinted with permission from [10], G. Bourianoff, *Computer* 36, 44 (2003). © 2001, Prentice Hall.

Fig. 4(a), the potential goes nearly to the background value in the middle of the channel, which means a charge carrier gets little Coulomb acceleration and would not easily travel from source to drain. In Fig. 4(b) with a deep source and drain, the electrical field will be strong, causing a considerable leakage current (i.e., amount of charge to flow between source and drain), which makes it hard to switch the transistor off at all in the worst case.

Since we still need a large enough drive current to flow through the device in the on-state of the MOSFET, we need to increase the dopant (and thus the charge carrier) concentration in the source and drain areas when making them shallower, since the current is proportional to the carrier concentration. The common measure for the amount of doping for a given



**Figure 4.** Color contour plot in arbitrary units of the electrical potential for (a) shallow and (b) deep source/drain profiles. The transistors are in the "off" state with no bias on the gate, source, and channel and a high bias on the drain.

junction depth is the "sheet resistance," which is inversely proportional to the product of carrier concentration and junction depth.

Common dopant atoms are elements from group III and V of the periodic system like boron, arsenic, or phosphorous. Substituting a silicon atom (group IV) with an arsenic atom, for example, leads to a situation where four of the arsenic atom's five valence electrons participate in the strong covalent bonds with the silicon atom, whereas the fifth one is weakly bound with a binding energy of only about 32 meV. Since this energy is very close to the thermal energy per atom at room temperature, the arsenic atom is easily ionized with the fifth electron becoming a free charge carrier. Similarly, a boron atom (group III) is missing an electron to form four covalent bonds with silicon and can easily accept an electron from a neighboring bond, thus creating a hole. Group V doped silicon is usually referred to as $n$-type, whereas group-III doping leads to $p$-type silicon. Usually, both $n$-type and $p$-type channel devices are used in a complimentary fashion (resulting in Complimentary Metal Oxide Semiconductor (CMOS) devices) to achieve small geometries and low power consumption. In $n$-type ($p$-type) material, electrons (holes) are the majority charge carriers.

The source and drain regions are doped to be of the same type, $p$ or $n$, while the channel region is doped to be the opposite type of the source and drain. This results in the source and drain regions being electrically isolated from one another when no voltage is applied to the transistor, provided source and drain are shallow enough. However, when a voltage is applied to the gate in a way that it attracts the source/drain charge (positive for $n$-type, negative for $p$-type) and drives the channel majority carriers away from the gate, a conducting layer of charge can form, electrically connecting the source and drain regions. A voltage between source and drain will then lead to current flow through the device.

Although the common dopant atoms, which are the ones that can be easily ionized at room temperature, have different solubilities (Fig. 5), the contributed maximum carrier concentration (or the "active fraction" of the dopant atoms) is fairly constant for all of them at approximately $2 \times 10^{20}$ cm$^{-3}$ [11]. The inactive dopant atoms do not occupy substitutional sites but form different structures with changed bonding situations, where all electrons are participating in the bonding and do not contribute to the conduction. An example for this is a "$B_2I$" complex (Fig. 6; $I$ stands for Si self-interstitial, meaning there is one more atom than lattice sites in the cluster), where a dumbbell with two boron atoms replaces a silicon lattice atom, leaving all silicon atoms in fourfold coordination with three bonds for each boron atom [12]. Thus, in the best case, the doping distribution consists of a box profile with a concentration of $2 \times 10^{20}$ cm$^{-3}$, which immediately goes to zero at the end (or "junction depth") of the doped source and drain regions. This is of course an idealization, since real dopant profiles are created by ion implantation, where the leading edge is not perfectly sharp (Fig. 7). The projections of the International Technology Roadmap for Semiconductors (ITRS, [2]) for



Figure 5. Solid solubility limits of dopants in Si as a function of temperature (values taken from Ref. [11]).

**Figure 6.** "B,I" complex, where a dumbbell consisting of two boron atoms has replaced a silicon lattice atom. Reprinted with permission from [12], X. Y. Liu et al., *Appl. Phys. Lett.* 77, 2018 (2000); 77, 4064 (2000). © 2000, American Institute of Physics.

future performance requirements ask for junction depth versus sheet resistance values which require to exceed the practical solubility limit of $2 \times 10^{20}$ cm$^{-3}$ and to make the profile as steep as possible (Fig. 8).

Combining the continuing miniaturization while increasing the dopant concentration leads to the "ultrashallow junction problem" of approaching the solubility limit for the dopant atoms in silicon and even having to exceed it by processing techniques that stabilize the metastable supersaturated solid solution. This field has seen intense research during the past decade and has resulted in many interesting findings, like the identification of sub-microscopic nano precipitates of dopant phases, which deactivate the dopants electrically, or nanoscale processing techniques that stabilize metastable solid solutions. Many of the techniques currently used for the modeling and simulation of nanoelectronics devices have their roots in applications for end-of the-roadmap microelectronic devices. We thus have chosen key applications and methods from such work as application examples in this article.



**Figure 7.** Real dopant distribution of a $2 \times 10^{15}$ cm$^{-2}$ boron dose ion-implanted with an energy of 10 keV, in comparison to the ideal case of a box profile with the maximum active concentration of $2 \times 10^{20}$ cm$^{-3}$.

Figure 8. ITRS requirements [2] for junction depth versus sheet resistance as a function of year. The shaded area requires exceeding the dopant solubility limit. The separation line between shaded and white area is calculated from the box profile in Fig. 7.

## 1.3. Molecular Electronics

From its literal meaning, molecular electronics is defined as electronics whose properties are determined by the chemical, physical, and electronic structures of molecules. This definition is very broad; it includes conductive polymers and even the insulating polymer layers on metal wiring. In common use, molecular electronics refers to molecular structures whose characteristic features are on the nanoscale and that contain between one and a few thousand molecules.

Molecules are naturally small, and their abilities of selective recognition and binding can lead to cheap fabrication using self-assembly. In addition, they offer tunability through synthetic chemistry and control of their transport properties due to their conformational flexibility. Remarkable progress in this field has been made in the last few years, as researchers have developed ways of growing, addressing, imaging, manipulating, and measuring small groups of molecules connecting metal leads. Several prototype devices such as conducting wires, insulating linkages, rectifiers, switches, and transistors have been demonstrated [13].

In parallel, there has been significant theoretical activity toward developing the description of nonequilibrium transport through molecules [14]. It is hard to say whether in time such devices could conceivably complement the current silicon-based integrated circuit (IC) industry or generate entirely new areas of applicability. It is clear, however, that in any case we will need to develop models to describe large-bias transport through ultrasmall devices, whether based on silicon or molecules. The fundamental challenge for modeling of molecular electronic involves determining the structure/property relationships for electronic transport (intra- or intermolecular) through a junction containing one or a few molecules as the transport medium and with either two (source/drain) or three (source/gate/drain) electrodes. Such transport structures are a dominant theme of contemporary molecular electronics. A very good summary of the current developments in this field can be found in Ref. [15].

Theoretical ideas for interpreting the current/voltage characteristics in molecular junctions began to appear in the 1990s. Some early approaches used the Landauer formulation [16], originally developed for semiconductor devices. In this simple picture, transport through molecular junctions is interpreted in terms of elastic scattering, and the conductance is given as the product of the quantized unit of conductance $(12.9 \text{ k}\Omega)^{-1}$ and a transmission coefficient describing how effective a molecule is in scattering the incoming electron from the upstream lead into the downstream lead. More powerful formulations, including the nonequilibrium Green's function (NEGF) approach, were also introduced in the mid-1990s [17, 18].

Because the conductance through a molecule depends very sensitively on the correct description of its contact structure with its leads, an atom-by-atom description of the system and fully atomistic process modeling becomes necessary. To cover the necessary time scales, methods beyond straightforward molecular dynamics (MD) simulations such as accelerated molecular dynamics (which will be described in Section 2.2.3) are necessary. This is still a very young field with few existing publications (except for *ad hoc* relaxed contact structures) and thus is not the focus of the current article, although the discussed methods can be applied to molecular systems in a straightforward manner. An example of a carbon nanotube device with titanium leads with contacts optimized with *ab initio* accelerated dynamics is shown in Fig. 9(d) (side view, Fig. 9[b]), which is very different from the metastable structure that a plain relaxation would predict (Fig. 9[c]). The change in the structure affects the contact resistance and thus the conductivity considerably (Fig. 9[e]), with the optimum conduction for an intact nanotube with just rim contact to the titanium leads (Fig. 9[a]) [19]. Because of the large numerical requirements of the calculation, the studied system was rather small, and the periodic boundary conditions and the short overlap between nanotube and leads prevent the nanotube in Fig. 9(c) from completely opening. A more realistic optimized structure is shown in Fig. 9(f).



Figure 9. Carbon nanotube molecular electronic devices with titanium leads in different contact configurations. (a) Carbon nanotube in end contact with titanium leads; (b) lateral view of carbon nanotube in side contact with titanium leads; front view of contact region of (c) relaxed and (d) temperature-accelerated dynamics annealed side-contact structure; (e) current versus voltage characteristics of structures (a), (c), and (d). (f) large-scale contact structure between carbon nanotube and Ti leads. Reprinted with permission from [19], K. Ravichandran et al., unpublished.

# 2. PROCESS SIMULATION THEORY FROM MICRO- TO NANOELECTRONIC DEVICES

During the fabrication of nanoelectronic devices, most mechanisms that lead from the raw materials to the finished product can be described by only two processes, diffusion and reaction (in addition to the radiation effects due to ion implantation, which is not a central focus of the current chapter). In the following, we will discuss their traditional treatment in modeling of microelectronic systems, which are mostly concerned with solid-state diffusion, solid-solid and solid-liquid reactions. We will describe the methodology for dopant and defect diffusion and reaction in a host semiconductor material to introduce the important concepts. As a footnote, the diffusion-reaction picture seems to be more general than previously thought, since a recent application of reaction-diffusion kinetics to the modeling of oxide growth rates [20] was found to be more successful for the modeling of growth of ultrathin gate oxides from dry oxidation than the more *ad hoc* classical Deal-Grove model [21] and its extensions [22, 23].

## 2.1. Micro- and Nanoscale Process Simulation within the Continuum Approximation

### 2.1.1. Diffusion-Reaction Equations

A totally rigorous physical model for the redistribution of impurities and defects in semiconductors does not exist and would probably be too complicated to be implemented into continuum-based process simulation programs. Among the various simplifications, the so-called methodology of diffusion-reaction equations was shown to be very efficient for the numerical simulation of diffusion phenomena. It was developed by Yoshida et al. [24] to describe the diffusion of phosphorus and can be generalized to design a set of coupled partial differential equations from some assumptions about possible reactions between intrinsic point defects and impurities. The following discussion follows closely the description of Pichler [25].

A basic assumption of all diffusion-reaction schemes is that interactions between impurities and point defects lead to distinguishable defect configurations. Each of them is characterized by its diffusion coefficient and charge state. Substitutional impurities are generally assumed to be immobile. Interactions are approximated by binary reaction schemes as discussed below. All possible charge states of the point defects and their possible reactions have to be considered individually at first. Charge states of intrinsic as well as of extrinsic point defects are assumed to result from binary reactions with charge carriers which are assumed to act as quasi-particles.

Development of a model starts with defining the point defects and complexes considered within the model and the possible reactions between them. As an example, the important case of the diffusion of substitutional impurities (dopants) will be considered in the following. Basically all interactions between such substitutional impurity atoms $X_s$, interstitial impurity atoms $X_i$, vacancies $V$, and self-interstitials $I$ fall into the three categories of reactions:

1. An interstitial impurity reacts with a vacancy and becomes substitutional according to the reaction

$$X_i + V \rightleftharpoons X_s \tag{1}$$

This reaction is called "Frank-Turnbull reaction." The reverse reaction is also known as "dissociative mechanism."

2. A self-interstitial displaces a substitutional impurity atom to an interstitial site according to the reaction

$$I + X_s \rightleftharpoons X_i \tag{2}$$

which is known as "kick-out reaction." The reverse reaction is known as "Watkins replacement mechanism."

3. A vacancy or self-interstitial forms a mobile pair with a substitutional impurity according to the reactions

$$V + X_s \rightleftharpoons XV \tag{3}$$

$$I + X_s = XI \tag{4}$$

A comparison of the pair reaction with self-interstitials and Eq. (2) shows that both reactions have the same structure. Because the atomistic form of the mobile complex is only of secondary importance in diffusion-reaction equations, only pair reactions with vacancies present a complementary third reaction to the Frank-Turnbull reaction and the kick-out reaction.

Since charge has to be conserved in the reactions, electrons or holes may also appear at the right-hand sides of the reaction equations. In addition to reactions between extrinsic and intrinsic point defects, reactions between intrinsic point defects may also be considered. The number of considered species (defects and complexes) can be freely chosen to include processes deemed important for the simulation problem at hand. In conjunction with the somewhat hand-waving definition of the reaction constants, discussed in the following section, one can truly speak of a fundamentally phenomenological theory.

As an example to demonstrate the methodology and to introduce the concepts, a fictitious system, consisting of negatively charged and neutral impurities on substitutional sites $X_s^-$ and $X_s^{||}$, negatively charged interstitial impurities $X_i$, and positively charged self-interstitials $I^+$ will be considered in the following. Their concentrations will be denoted by $C_{X_s^-}$, $C_{X_s^{||}}$, $C_{X_i}$, and $C_I$. The holes involved in the reactions will be symbolized by $h^+$, their concentration by $p$. Between these species, the reactions

$$X_s^- + h^- \underset{k_1^r}{\overset{k_1^f}{\rightleftharpoons}} X_s^{||}$$

$$X_s + I^i \underset{k_2^r}{\overset{k_2^f}{\rightleftharpoons}} X_i + h^+ \tag{5}$$

are assumed. The symbols $k_i^f$ and $k_i^r$ stand for the rate constants in forward and reverse direction, respectively. For each of the point defects considered, a continuity equation

$$\frac{\partial C}{\partial t} = -\nabla \cdot \mathbf{J} + G - R \tag{6}$$

needs to be provided, where $\mathbf{J}$ is the flux or diffusion current of the species under consideration, defined as the number of respective atoms passing a unit area within a unit time interval, and $G$ and $R$ stand for the generation and loss rates (often called generation and recombination terms) [25]. In the absence of sinks and sources ($G = R = 0$) or when $G - R = 0$, Eq. (6) is known as Fick's Second Law. In the absence of driving forces, which are all influences which make the jump frequency to depend on the direction of the jump [26], one can use Fick's first law,

$$\mathbf{J} = -D\nabla C \tag{7}$$

to calculate the flux from diffusivity $D$ and the gradient of the concentration. The presence of a driving force $\mathbf{F}$ adds a second term to Eq. (7).

$$\mathbf{J} = -D\nabla C + DC\frac{\mathbf{F}}{k_B T} \tag{8}$$

where $k_B$ is Boltzmann's constant and $T$ is the temperature of the system. For the example of a charged defect/atom with charge $q$ in the presence of an electrostatic field $\mathbf{E}$ (which can be due to external bias or to a charge distribution in the system), the driving force would be given by $\mathbf{F} = q\mathbf{E}$.

---

This limits the useful range of the model to moderate doping temperature regimes, since otherwise $p$ needs to be multiplied by the activity coefficient $\gamma$. [25].

As an example for generation and loss, the first reaction in Eq. (5) describes the generation of a $X_i^0$ while consuming a $X_i$ and a $h^-$. Thus, the continuity equation for species $X_i^0$ needs to contain a generation term, which depends on the availability of $X_i$ and $h^-$, and the equation for $X_i$ needs to contain a corresponding loss term. How do these terms look?

In general terms, the $i = 1, \ldots, N$ allowed reactions between the reactants $A, B, C, D, \ldots$ can be written as

$$a_i A + b_i B + \cdots \underset{k_i^r}{\overset{k_i^f}{\rightleftharpoons}} c_i C + d_i D + \cdots \tag{9}$$

where the stoichiometric number of species $A$ in reaction $i$ is denoted by $a_i$. By definition, stoichiometric numbers appearing on the left-hand side are negative. The extent of a reaction can be described by a reaction constant $\zeta$ [27]. For our purposes, we define $\zeta$ as the number of completed reaction events as described in Eq. (9) per volume $V$. Thus, when only one reaction has to be considered, the change in the concentration of species $A$ in a closed, constant volume is simply given by

$$dC_A = a\, d\zeta \tag{10}$$

When several parallel reactions have to be considered, the change in the concentration of any of the species $A$

$$dC_A^{\text{reactions}} = \sum_i a_i\, d\zeta_i \tag{11}$$

is the sum of the changes associated with the particular reactions $i$. Taking the derivative with respect to time of Eq. (11) relates the generation and loss terms of Eq. (6), $G_i - R_i$, to the reaction rates $d\zeta_i/dt$, since the concentration change due to reactions in Eq. (11) excludes the diffusion term $-\nabla \cdot \mathbf{J}$,

$$\frac{dC_A^{\text{reactions}}}{dt} = G_i - R_i = \sum_i a_i\, \frac{d\zeta_i}{dt} \tag{12}$$

How can we determine the reaction rates $d\zeta_i/dt$? Besides the well-known mass action law for equilibrium concentrations of products and reactants [25], there is also a kinetic law of mass action [28–31], which relates reaction rates to the concentrations of the reactants and thus allows to formulate the generation and loss terms $G - R$ from Eq. (12) explicitly. In our notation, the kinetic law of mass action reads

$$\frac{d\zeta_i}{dt} = k_i^f C_A^{-a} C_B^{-b} \cdots - k_i^r C_C^c C_D^d \cdots \tag{13}$$

We have now all necessary expressions in place and can return to setting up the continuity equations for our example system. Assuming substitutional defects to be immobile as postulated above, we use Eqs. (8), (12), and (13) to formulate

$$\frac{\partial C_{X_i^0}}{\partial t} = k_1^f C_{X_i}\, p - k_1^r C_{X_i^0} \tag{14}$$

$$\frac{\partial C_{X_i}}{\partial t} = -k_1^f C_{X_i}\, p + k_1^r C_{X_i^0} - k_2^f C_{X_i} C_I + k_2^r C_{X_i}\, p \tag{15}$$

$$\frac{\partial C_{X_i}}{\partial t} = \nabla \left[ D_{X_i} \left( \nabla C_{X_i} - C_{X_i} \frac{e\mathbf{E}}{k_B T} \right) \right] + k_2^f C_{X_i} C_I - k_2^r C_{X_i}\, p \tag{16}$$

$$\frac{\partial C_I}{\partial t} = \nabla \left[ D_I \left( \nabla C_I - C_{X_i} \frac{e\mathbf{E}}{k_B T} \right) \right] - k_2^f C_{X_i} C_I + k_2^r C_{X_i}\, p \tag{17}$$

---

The form presented in Eq. (13) is strictly only valid for ideally dilute solutions. For higher concentrations, the concentrations should be multiplied by the activity coefficient $\gamma_i$ ($= 1$ for dilute solutions) which is used to introduce non-ideal effects and will, in general, depend on local parameters like the concentrations of all species in the system as well as on macroscopic parameters like temperature and pressure.

It should be kept in mind that point-defect-impurity pairs and interstitial impurity atoms considered here as individual point defects are just abstractions for the complicated interactions of intrinsic point defects and impurities. Also, the inclusion of the hole concentration in the product $k_2^s C_{X_s} p$ is not mandatory. It serves just to keep $k_2^s$ Fermi-level-independent.

The resulting system of coupled partial differential equations can be solved directly with one of the general-purpose process simulation tools which became customary or with any general partial differential equation solver. Examples for such simulation codes incude PROPHET [32], ALAMODE [33], TAURUS-PMEI [34], and FLOOPS [35]. In general, however, it is often desirable to reduce the number of equations, in order to understand the macroscopic behavior of the system or to speed up solving the equations. The latter is due to the fact that the fastest reaction determines the maximum allowable time step for solving the partial differential equation system. Eliminating the fast reactions thus by physical argu-ments or equilibrium relationships allows to use a longer time step. Because of the much higher mobility of charge carriers than of atoms, equilibrium between the charge states of a point defect will be established on a time scale on which concentration changes due to diffusion or reactions are negligible. On such a time scale, the generation of $X_s^0$ from $X_s^-$ (first term in Eq. [14]) will be equal to the loss of $X_s^0$ to form $X_s^-$ (second term in Eq. [14]), meaning that Eq. (14) is equal to zero. Equation (14) then gives for $C_{X_s}$

$$C_{X_s} = \frac{k_1^i}{pk_1^f + k_1^r} C_{X_s}$$

(18)

using the total concentration $C_{X_s} = C_{X_s^0} + C_{X_s}$ of substitutional dopant atoms. The assump-tion of local equilibrium between the charge states also means that the first two terms in Eq. (15) are zero. With that and Eq. (18), Eqs. (14) and (15) can be combined into one continuity equation,

$$\frac{\partial C_{X_s}}{\partial t} = \frac{\partial \left(C_{X_s^0} + C_{X_s}\right)}{\partial t} = \frac{k_1^r k_2^i}{pk_1^f + k_1^r} C_{X_s} C_I + k_2^s C_{X_s} p$$

(19)

Another customary assumption valid at least for long diffusion times (but inadmissible for short diffusion times [25]) is that the mobile complexes and the substitutional atoms are also in local equilibrium. This means the sum of generation and loss term (terms 3 and 4 in Eq. [15]) is also zero, which leaves us with just two continuity equations from Eqs. (16) and (17) and allows us to express the concentration of interstitial impurities $C_{X_I}$ in the form

$$C_{X_I} = \frac{k_1^r k_2^i C_I}{k_2^s p(pk_1^f + k_1^r) + k_1^r k_2^i C_I} C_X$$

(20)

using the total concentrations of impurities $C_X = C_{X_s} + C_{X_I}$ and self-interstitials $C_I = C_{I^+}$, and to combine Eqs. (16) and (19) into a continuity equation of the form

$$\frac{\partial C_X}{\partial t} = \nabla \left\{ D_X \left[ \nabla \left( \frac{k_1^r k_2^i C_X C_I}{k_2^s p(pk_1^f + k_1^r) + k_1^r k_2^i C_I} \right) - \frac{k_1^r k_2^i C_X C_I}{k_2^s p(pk_1^f + k_1^r) + k_1^r k_2^i C_I} \frac{eE}{k_B T} \right] \right\}$$

(21)

With that, all continuity equations for extrinsic point defects were finally combined into one nonlinear diffusion equation for the total concentration. In fact, the physical assumptions contained implicitly in such macroscopic diffusion models and their range of validity can be made clear by comparison with diffusion-reaction models.

Typical assumptions for the intrinsic point defects are equilibrium, a time-dependent but depth-independent oversaturation as for oxidation-enhanced diffusion, or even a function of space and time, as for implantation-enhanced diffusion. The last simplification often found n the literature for dopant diffusion is that the concentrations of mobile extrinsic complexes are much smaller than the concentration of substitutional atoms. Using the equilibrium condition for interstitial with substitutional impurities, this means $C_{X_I}/C_{X_s} = k_2^i C_I /(k_2^s p) \ll 1$ and that the second term in the denominator of the right-hand side (r.h.s.) of Eq. (20) can be neglected, which leads to $C_X \propto C_X C_I$. It should be noted that this last simplification

leads to a break-down of the resulting equation in all situations where point defects are present in large oversaturation (where the neglected term would guarantee limiting $C_V$ to no more than $C_V$).

In addition to the simple reactions between intrinsic and extrinsic point defects, reactions between extrinsic point defects have to be considered at high concentrations. Typical examples are the formation of ion or defect/ion pairs, that is complexes of a donor and an acceptor impurity or defect at adjacent lattice sites (see, e.g., Section 3.1.2), or of clusters of intrinsic point defects or impurities of the same kind (see, e.g., Section 3.2.2). Ion pairing and the formation of a cluster of two impurities can be incorporated straightforwardly into the diffusion-reaction approach via binary reactions of, for example, a mobile donor complex and a substitutional acceptor ion in case of the ion pair. Clusters with a higher number of atoms can be included, in principle, by considering the complete chain of cluster sizes up to the largest one. Their evolution with time is then determined by a system of binary reactions. However, the expenses in terms of computer resources are considerable and they are usually inadmissible for precipitates which may consist of thousands of intrinsic point defects or impurity atoms. Therefore, the *ad hoc* formation of extended defects or clusters of a certain size is often assumed and formulated as a reaction. While this is certainly a possible simplification, it has to be noted that the absolute values of the forward and backward rate constants can no longer be derived from kinetic considerations. Their ratio, on the other hand, is still determined by the law of mass action.

### 2.1.2. Reaction Rate Constants

As described in Section. 2.1.1, the reaction-diffusion scheme needs as input parameters diffusion constants $D$ and reaction rate constants $k$. The latter again are usually approximated as functions of the diffusion constants, as we will see. We want to discuss the theory of the reaction rate constants in the present section, while leaving the calculation of the diffusivities from atomistic hopping mechanisms for Section 2.1.3. An excellent review of the reaction rate theory can be found in [25].

Von Smoluchowski [36] presented in 1917 the first theory of rate constants within a diffusion-reaction scheme, applied to the coagulation in colloidal solutions. According to his theory, the rate of coagulation is determined mainly by the diffusion of the reactants toward each other. As soon as the distance between the reactants is as low as a capture radius $a_c$, the reaction is assumed to take place immediately.

It is evident that this mechanism will lead predominantly to binary reactions in the form

$$A + B \underset{k^r}{\overset{k^f}{\rightleftharpoons}} C \tag{22}$$

Assuming the reaction to be diffusion-limited, the von-Smoluchowski approach leads to a rate constant given by

$$k^f = 4\pi a_c (D_A + D_B) \tag{23}$$

The symbols $D_A$ and $D_B$ stand for the diffusion coefficients of the reacting species. The theory presented by von Smoluchowski contained implicitly limiting statistical assumptions. However, based on a more solid statistical basis, Waite [37] came to virtually the same conclusion. The capture or reaction radius $a_c$ for reaction between point defects has been usually assumed to be on the order of the distance between two substitutional silicon atoms $(a_c \sim 2.5$ Å$)$, but recent atomistic simulations have shown that it might be about two to three times larger, as discussed in Section 3.1.2. For reactions between point defects and extended defects, an increased capture radius is sometimes introduced which reflects the geometry of the extended defect (see Section 2.2.2) . The treatment of von Smoluchowski and Waite is valid only if, at least, one of the reacting species is electrically neutral. When both reacting species are electrically charged with charges $q_1$ and $q_2$, Coulombic attraction and repulsion has to be taken into consideration. As demonstrated by Debye [38], this leads to a modification of the rate constant in the form

$$k^f = 4\pi a_c (D_A + D_B) \frac{E_{Coul}(a_c)/(k_B T)}{\exp[E_{Coul}(a_c)/(k_B T)] - 1} \tag{24}$$

with the Coulomb energy

$$E_{Coul}(a_c) = \frac{q_1 q_2}{4\pi\varepsilon_0\varepsilon_r a_c} \tag{25}$$

The physically correct value for the relative dielectric constant $\varepsilon_r$ that might be different from the macroscopic value on this extreme nanoscale ($a_c$ is on the order of a few Å) has been a long-standing question [25]. However, it can be determined from atomistic first-principles calculations to be more or less unmodified from the macroscopic value and will be discussed in Section 3.1.2.

In case there is a energy barrier higher than the diffusion barriers to overcome before the reaction can happen, it makes sense to not longer assume that the reaction will happen instantaneously once the reactants are closer than the capture radius. Instead, one assumes a thermally activated process [25, 39],

$$k^f = 4\pi a_c(D_1 + D_B)\exp\left(\frac{-\Delta G}{k_B T}\right) = 4\pi a_c(D_A + D_B)\exp\left(\frac{-\Delta S}{k_B T}\right)\exp\left(\frac{-\Delta H}{k_B T}\right) \tag{26}$$

where $\Delta G$ stands for the barrier in excess of the Gibbs free energy of diffusion of the fastest reactant. The energy barrier might be an enthalpy barrier or an entropy barrier. The meaning of an entropy barrier is that the entropy of the system is lowered during the reaction. Similar modifications can be made for the case of two charged reactants, Eq. (24).

The situation becomes slightly more complicated when charge states and charge carrier concentrations are involved. Let us consider the ionization reaction of an acceptor,

$$A^0 + e \underset{k^r}{\overset{k^f}{\rightleftharpoons}} A^- \tag{27}$$

For an order-of-magnitude estimate, the electron will be assumed to be a quasi-defect to which Waite's theory can be applied. Then, the change in the concentration of neutral acceptors is given by

$$\frac{\partial C_{A^0}}{\partial t} = -k^f C_{A^0} n - k^r C_A = -\frac{C_{A^0}}{\tau} + \cdots \tag{28}$$

Therein, $\tau$ is the characteristic time constant for the charging process. Using Waite's theory, an estimate for $\tau$ can be given in the form

$$\tau = \frac{1}{4\pi D_e a_c n} \tag{29}$$

The diffusion coefficient of the electrons can be estimated via the Einstein relation $\mu = Dq/(k_B T)$ [40, 41] from the mobility given, for example, by Sah et al. [42]. From ab initio calculations of As-$V$ interactions, the electronic reaction radius is at least 1 nm as will be shown in Section 3.1.2. Assuming a reaction radius of 2 nm and the electron concentration to correspond to the intrinsic concentration leads to time constants of less than one picosecond at temperatures of 500°C and above. Thus, it is in general assumed that steady state between the charge states of a defect is established quasi-instantaneously in comparison to process times.

When charge states are involved in binary reactions, as in the reaction equation

$$A^k + B^m \underset{k^r}{\overset{k^f}{\rightleftharpoons}} C^n + (z_k + z_m - z_n)e \tag{30}$$

additional care has to be exercised in estimating the reaction constants. According to the kinetic reaction theories discussed above, the forward reaction rate depends just on the reactants and corresponds to the total concentration of defects $C$ forming per unit time by this reaction. How many of the defects $C$ form in a particular charge states is not predicted. The most straightforward assumption seems to be that charge is conserved so that the pairs forming have charge state $z_k + z_m$. But this is just as artificial as any other assumption since it is not even guaranteed that such a charge state exists for defect $C$. However, at least at

elevated temperatures, it can be assumed that steady state between the charge states $n$ of the products $C$ is established more or less instantaneously. Then, the concentration of pairs forming in a specific charge state per unit time is given by the product of the total rate times the Fermi-level dependent probability that the pair exists in this charge state. It is convenient to use only the total concentrations of reactants and products then, but when individual forward reaction constants have to be given explicitly, they can be written in the form

$$k^l_{kmn} = k^j_{km} \frac{C_{l^n}}{\sum_{\bar{n}} C_{l^n}}$$                                (31)

The rate constants $k^l_{km}$ therein can be estimated as we have discussed.

Waite's theory in the form of Eq. (23) and its extensions are valid only for a uniform distribution of the reacting species. Such conditions are not expected under conditions of, e.g., electron irradiation or ion implantation where vacancies and self-interstitials may occupy closely correlated sites. Extensions to include spatial correlation effects were discussed, e.g., by Peak and Corbett [43].

Having quantified the forward reaction constant $k^f$, we need to quantify also the backward reaction constant $k^r$. However, it is often not possible to derive the reverse reaction constant from kinetic theories. Instead, one uses the fact that, in steady state, the forward and backward rates have to be equal so that $d\zeta_i/dt$ in Eq. (13) vanishes. This allows writing the backward reaction constant in terms of the forward reaction constant and the concentrations in steady state,

$$k^r = k^f \left( \frac{C_A C_B}{C_C} \right)_{\text{steady state}}$$                       (32)

where the concentrations again need to be multiplied by the corresponding activity coefficients in the case of non-dilute concentrations. For the temperature dependence of the backward reaction constant, $k^f$ will contribute with the enthalpy of diffusion of the faster-migrating reactant plus eventual reaction barriers, and $C_A C_B/C_C$ with the binding energy. The sum of these contributions, that is the barrier, which has to be overcome thermally to dissociate the pair, is usually called "dissociation energy."

## 2.1.3. Macroscopic Diffusion Constants from Atomic Hopping

Besides the formation energies, which control the equilibrium concentrations and which are reasonably easy to calculate (at least for elemental systems), the diffusion constants are the central parameters for diffusion reaction systems. For simple systems, only a few species (concentrations) need to be considered in the equation system, and a fit of the kinetic parameters to experimental diffusion data can help to identify them. For nanoscale device systems, on the other hand, where the concentration and nature of many different nanoclusters and molecules become important, this becomes a daunting, if not impossible task. In these cases, atomistic calculations can help to determine the kinetic constants as a function of temperature.

On the atomic scale, diffusion happens when defects, impurities or their complexes perform hops between the different equivalent positions of the crystal structure. Equivalently, reactions, where atoms perform hops, change positions, and create new complexes, are also dominated by atomic hopping. Whereas it is in some cases trivial to extract the macroscopic diffusion constants (such as in the case of a migrating vacancy in an otherwise perfect crystal), it is in general a difficult task to statistically combine the jump rates into the correct macroscopic diffusion constant, and a complete and general recipe how to do it was not available in the literature until 2001 [44]. The complicated nature of this process is not always obvious. We will show in the following that in case there is one step with a much higher energy barrier than all the other hopping events of a diffusion path, the simplified assumption that the highest energy barrier controls the complete diffusion mechanism is reasonable. However, when the barriers start to be comparable, this is of course not the case anymore [45].

A very simple example illustrates the need for a more general treatment. Consider diffusion along an alternating chain, illustrated in Fig. 10, where long-range diffusion along
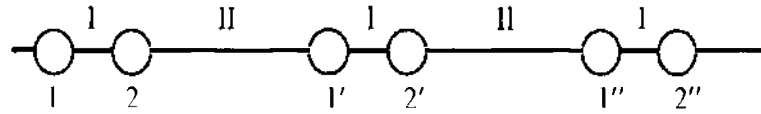
**Figure 10.** Diffusion along an ideal alternating chain. The sites are successively labeled 1, 2, 1', 2', and so on. The hops between sites are of two types, with rates $p_I$ and $p_{II}$. Long-range diffusion occurs in series, so that the resulting diffusivity is limited by the slower of rates $p_I$ and $p_{II}$. Reprinted with permission from [44], M. S. Daw et al., *Phys. Rev. B* 64, 045205 (2001). © 2001, APS.

the chain must combine two hopping rates ($p_I$ and $p_{II}$) in series. A steady current $j$ flowing along the chain requires that the differences in concentrations on the sites ($c_i$) obey $j = (c_2 - c_1)/p_I = (c_{1'} - c_2)/p_{II} = (c_{1'} - c_1)/p_{eff}$ which gives $p_{eff} = p_I p_{II}/(p_I + p_{II})$, thus combining the rates in series. As expected, for strongly different rates like $p_I \gg p_{II}$, our result finds that $p_{eff} \sim p_{II}$, i.e., the slower rate determines the overall diffusion, which is of course not the case anymore when the two rates start to be comparable. Nevertheless, the "non-Arrhenius" behavior of the diffusivity in such a case can still be surprising [45].

In the following, we will outline the general theory of [44] that describes how to combine the atomic hopping events into macroscopic diffusion constants. We assume that there are well-defined states (i.e., energy valleys) which the system can occupy. The concept of state here is quite general. A state can be simply a defect, such as a vacancy, located on a particular lattice site, so that a state is distinguished entirely by specifying its location: this is clearly the case for vacancies in elemental metals.

More generally, however, the defect can have other degrees of freedom. For example, the primitive cell in Si has a two-atom basis. In that case, one must specify on which basis site the defect is located, in addition to which primitive cell. Or, for another example, a vacancy in Si might undergo a Jahn-Teller distortion (where electrons in defect states in the band gap can lower their energy by a symmetry breaking, which adds a little (tetragonal) strain energy to the system, but lowers the position of the occupied defect state), which has three independent orientations. In that case, one must specify orientation in addition to location. An even more complex example would be a cluster of impurity atoms which migrate together—the internal degrees of freedom in this case can be numerous, as we will show.

We assume furthermore that the basic assumptions of transition state theory (Section 2.2.1) hold and the rate of jumping from one state to another is controlled by the difference in energy from the initial valley to the saddlepoint. In the presence of stress, the energy difference, of course, must account for the additional work done against an external stress field. It is convenient to introduce a *local* reference, which is the perfect system subjected to the local value of the external stress. A defect is introduced into the local reference, and changes are measured in the system relative to that reference. We follow the convention of Ref. [44], where the word "creation" has been used to label values calculated with respect to this local reference. To keep the following description simple, we do not consider any trapping, reaction, or dissociation of the defects; we take only hops which maintain the unity of the defect.

An example for defects with a simple basis is an unreconstructed vacancy in a metal having a primitive cell with one atom (e.g., simple cubic, body-centered cubic, and face-centered cubic). In this case, the state is specified completely by designating the spatial location of its primitive cell, $A$. We can describe the system by hopping rates between the different sites $A$ and $B$.

$$p[A \rightarrow B] = p_0 \exp\left[-\beta\left(E_{c(s)}^{|AB|} - E_{c(n)}^{|}\right)\right]$$
$$\equiv \frac{M_{|AB|}}{S_I} \tag{33}$$

with

$$M_{|A|} = M_{|B \cdot i|} \equiv p_0 \exp\left(-\beta E_{c(s)}^{|AB|}\right) \tag{34}$$

$$S_I \equiv \exp\left(-\beta E_{c(n)}^{|}\right) \tag{35}$$

where $\beta = 1/(k_B T)$, $p_0$ is a basic hop rate, and $E_{c(n)}^{|}$ and $E_{c(s)}^{|AB|}$ are the creation energies of the valley at position $A$ and the saddlepoint energy between $A$ and $B$, respectively. In these

definitions, we separate the "solubility factor" $S$ of the defect in the valley, $S_A$, from the "mobility factor" $M_{[AB]}$.

We will begin the treatment by assuming a uniform host, and then extend the results at the end to a nonuniform host. In a uniform host, all valleys have the same energy, so that $E_{(r)}^A = E_{(r)}$ and $S_A = S$. The saddlepoint energies $E_{(s)}^{[AB]}$ do not depend on the absolute positions $A$ and $B$, but only on the relative position. (The saddlepoint energies are not all the same because the host is not assumed to be isotropic. In the presence of stress, even a cubic crystal is not isotropic.)

The concentration $c_A$ on site $A$ develops in time by

$$\dot{c}_A = -c_A \sum_B p[A \to B] + \sum_B c_B p[B \to A] \tag{36}$$

We can make use now of the translation symmetry of the host by expanding the solutions into a Bloch form,

$$c_B = \int_{BZ} u(\mathbf{k}, t) \exp(i\mathbf{k} \cdot \mathbf{R}_B) d^3k \tag{37}$$

Matching Fourier components of the rate equation (Eq. [36]) then gives

$$\dot{u}(\mathbf{k}, t) = \gamma(\mathbf{k}) u(\mathbf{k}, t)/S \tag{38}$$

**with**

$$\gamma(\mathbf{k}) = \sum_B M_{[0B]}[\exp(i\mathbf{k} \cdot \mathbf{R}_{B0}) - 1] \tag{39}$$

where the site 0 has been chosen arbitrarily and $\mathbf{R}_{JI} \equiv \mathbf{R}_J - \mathbf{R}_I$. Because we are looking for the longtime and macroscopic, that is, long-wavelength evolution of the system, we expand $\gamma(\mathbf{k})$ in powers of $\mathbf{k}$,

$$\gamma(\mathbf{k}) = -\frac{1}{2} \text{Tr}\left[ \left( \sum_B M_{[0B]} \mathbf{R}_{B0} \otimes \mathbf{R}_{B0} \right) \cdot (\mathbf{k} \otimes \mathbf{k}) \right] + O(k^4) \tag{40}$$

(The dyad $\mathbf{a} \otimes \mathbf{b}$ defines a matrix with components $a_i b_j$.) There are no terms of order unity because, by number conservation, $\gamma(\mathbf{k} = 0) = 0$. There are no terms linear in $\mathbf{k}$ because $\sum_B M_{[0B]} \mathbf{R}_{B0} = 0$ (in a solid with a simple basis, for every neighbor there is an opposite neighbor which cancels).

When we substitute the series expansion of $\gamma(\mathbf{k})$ (Eq. [40]) into the equations of motion (Eq. [38]) and transform back to real space, each power of $\mathbf{k}$ in $\gamma$ will be associated with a spatial derivative. Thus the second derivative of $\gamma(\mathbf{k})$ at $\mathbf{k} = 0$ will have physical significance in the resulting diffusion equation. We therefore define a "solid permeability tensor"

$$\mathbf{P} = -\frac{1}{2} \left. \frac{\partial^2 \gamma(\mathbf{k})}{\partial \mathbf{k} \partial \mathbf{k}} \right|_{\mathbf{k}=0} \tag{41}$$

and the resulting diffusion equation in the general case is

$$\dot{c}(\mathbf{x}, t) = \nabla \cdot \left[ \mathbf{P}(\mathbf{x}, t) \cdot \nabla \left( \frac{c(\mathbf{x}, t)}{S(\mathbf{x}, t)} \right) \right] \tag{42}$$

where the diffusion tensor $\mathbf{D}$ is obtained from the microscopic hop parameters by

$$\mathbf{D} = \frac{1}{S} \sum_B M_{[0B]} \mathbf{R}_{B0} \otimes \mathbf{R}_{B0} \tag{43}$$

In the case of a uniform host, the hop rates and the solubility are independent on position and are not affected by the gradient operator.

The explicit introduction of the solid solubility factor $S$ is very helpful because it illustrates clearly that the equilibrium condition is $c(\mathbf{x}, t) \propto S(\mathbf{x}, t)$ and also that a gradient in the

solubility factor acts as a driving force for diffusion [46]. The solubility is related to the local chemical potential by $\mu = k_B T \ln(c/S)$. Also, we have introduced the tensor quantity P, which is the product of the diffusivity and solubility factor. In analogy to gaseous and liquid systems, we have chosen to call this the "solid permeability factor"; in an anisotropic medium, the permeability factor in general is a tensor quantity. The solid solubility factor depends only on the valley energy (Eq. [35]), the solid permeability factor depends only on the saddlepoint energy (Eqs. [34], [40], and [41]), and the solid diffusivity depends on the migration energy (difference between saddlepoint and valley; Eq. [43]).

Clearly the choice of a local reference for the energies discussed previously does not affect the diffusivity (which depends on differences in energy), but does affect the solubility factor. In comparing the relative solubility factors of a defect at two different (stressed) locations, one must calculate the energy required to insert the defect into each (stressed) location. Thus, using the stressed but otherwise perfect Si as a local reference is natural for the diffusion problem.

When dealing with systems with a degenerate basis, some care is required beyond the previously discussed case of the simple basis. In this category are crystals with primitive cells having more than one atom in the basis, or defects with internal (e.g., orientational) degrees of freedom. For example, the diamond and hexagonal close-packed (hcp) structures have a two-atom basis. Also, complex defects, such as a dumbbell self-interstitial or vacancy with Jahn-Teller distortion, have an orientation which must be specified in addition to the site.

We identify now a lattice site by $\{Aa\}$, with its cell index $A$ and the index $a$ which denotes the state within the cell. The number of states in the basis is $N_{\text{states}}$. For a symmetric vacancy in diamond, as an example, $N_{\text{states}} = 2$, because there are two sites in the primitive cell. However, if we admit Jahn-Teller distortions, then on each of two atomic sites a vacancy can have one of three orientations, so that in this case $N_{\text{states}} = 6$. It is easy to imagine more complex defects (clusters, for example) where the number of states could be quite large.

The rate of jumping from state $\{Aa\}$ to state $\{Bb\}$ through the saddlepoint $[AaBb]$ is now given by

$$p[Aa \rightarrow Bb] = p_0 \exp\left(-\beta\left(E_{c(s)}^{[AaBb]} - E_{c(v)}^{Aa}\right)\right)$$
$$= M_{[AaBb]}/S_{Aa} \tag{44}$$

with

$$M_{[AaBb]} = M_{[Bb,Aa]} \equiv p_0 \exp\left(-\beta E_{c(s)}^{[AaBb]}\right) \tag{45}$$

$$S_{Aa} \equiv \exp\left(-\beta E_{c(v)}^{Aa}\right) \tag{46}$$

The concentration of state $a$ in cell $A$ develops in time according to

$$\dot{c}_{Aa} = -c_{Aa} \sum_{Bb} p[Aa \rightarrow Bb] + \sum_{Bb} c_{Bb} p[Bb \rightarrow Aa] \tag{47}$$

Making use of the translation symmetry of the host yields

$$c_{Bb} = \sum_{\mathbf{k}} u_b(\mathbf{k}, t) \exp(i\mathbf{k} \cdot \mathbf{R}_{Bb}) \tag{48}$$

Matching Fourier components of the rate equation then gives

$$\dot{u}_a(\mathbf{k}, t) = \sum_b \Gamma_{ab}(\mathbf{k}) u_b(\mathbf{k}, t)/S_b \tag{49}$$

with

$$\Gamma_{ab}(\mathbf{k}) = \sum_B M_{[0aBb]} \exp(i\mathbf{k} \cdot \mathbf{R}_{B0a}) - \delta_{ab} \sum_{Bc} M_{[0aBc]} \tag{50}$$

where the cell $A = 0$ has been chosen arbitrarily. Note that number conservation is expressed as

$$\sum_a \Gamma_{ab}(\mathbf{k} = 0) = 0 \tag{51}$$

The complete dynamics of the system are contained in $\Gamma_{ab}(\mathbf{k})$, which is a symmetric rate matrix of size $N_{states} \times N_{states}$. The eigenvalues are the rate constants of the relaxation process. At least one eigenvalue vanishes at $\mathbf{k} = 0$ because of number conservation (Eq. [51]). We expect generally that only one eigenvalue vanishes for $\mathbf{k} = 0$. This is because a conserved quantity is associated with each vanishing eigenvalue at $\mathbf{k} = 0$. We anticipate that the only conserved quantity associated with diffusion will be the total defect number (we have assumed that the defects diffuse intact). For small $\mathbf{k}$, we have then only one relevant mode (that is, the slowest). All the other modes are fast and correspond to short-range relaxation among the members of the primitive cell.

Consider, for example, the alternating chain discussed in the introduction. In this system, the primitive cell consists of two sites, so that $\Gamma$ is a $2 \times 2$ matrix:

$$\Gamma(\mathbf{k}) = \begin{pmatrix} -p_I - p_{II} & p_I e^{ik(x_2 - x_1)} + p_{II} e^{ik(x_2 - x_1 - a)} \\ p_I e^{ik(x_1 - x_2)} + p_{II} e^{ik(x_1 - x_2 + a)} & -p_I - p_{II} \end{pmatrix} \tag{52}$$

where $x_1$ and $x_2$ are the site positions within the unit cell, $a$ is the periodicity of the chain, and $p_I$ and $p_{II}$ are the rates associated with transitions over the two different saddlepoints. The two eigenvalues are $\gamma_\pm = -p_I - p_{II} \pm \sqrt{p_I^2 + p_{II}^2 + 2p_I p_{II} \cos(ka)}$. The relaxation mode associated with $\gamma_-$ (the "optical mode") has finite lifetime even when $k = 0$. This corresponds to a relaxation within the primitive cell which occurs very quickly to bring the sub-lattice into equilibrium. The rate $\gamma_+$ of the other ("acoustic") mode vanishes at $k = 0$ (that is, small, long-wavelength deviations from uniformity take a very long time to relax) and so this is the only relevant mode. The leading order term for small $k$ is $\gamma_+ \sim -\frac{1}{2}a^2 k^2 p_I p_{II}/(p_I + p_{II})$, which, when we transform to real-space, will become the operator $P \frac{\partial^2}{\partial x^2}$. Thus, $P = (a^2/2)p_I p_{II}/(p_I + p_{II})$, as we expect for processes in series.

In general, $N_{states}$ is potentially large, and it is very difficult to obtain an analytical form for the relevant eigenvalue. However, because we need only the behavior near $\mathbf{k} = 0$ perturbation theory can be used to obtain the permeability, as discussed in depth in [44].

Finally, we note that for a degenerate basis, the appropriate solubility factor is the average solubility over all states; that is, the proper solubility to use in the diffusion equation (Eq. [42]) is

$$\overline{S}_d \equiv \frac{1}{N_{states}} \sum_a S_{da} = \frac{1}{N_{states}} \sum_a \exp(-\beta E_{c(a)}^{(a)}) \tag{53}$$

Since modern nanoelectronic CMOS device structures use very often strained channel materials to enhance the charge carrier mobilities and make the device faster, it is often necessary to specify how stress affects the valley and saddlepoint energies in the hop rates (Eq. [44]). For this, we will assume that the host, in the absence of stress, is uniform.

When a defect is created, the solid changes shape from its original condition. In linear elasticity, the change in the shape of a volume can be expressed as a real, symmetric tensor. For example, a sphere is distorted into an ellipsoid, and the difference can be described in complete generality (within linear elasticity) by three principal values $\Omega_{ci}$ and axes $\hat{t}_i$.

$$\Omega_c = \Omega_{c1}\hat{t}_1 \otimes \hat{t}_1 + \Omega_{c2}\hat{t}_2 \otimes \hat{t}_2 + \Omega_{c3}\hat{t}_3 \otimes \hat{t}_3 \tag{54}$$

where "$c$" denotes again the term "creation" as defined in [44]. The symmetry of the shape change is determined by the symmetry of the defect, i.e., the principal axes will be symmetry axes of the defect. If the defect has an orientation (for example, a dumbbell self-interstitial or a Jahn-Teller distorted vacancy), the principal axes will be directed accordingly, but the set of principal values will be the same for all orientations. A defect with cubic symmetry will have three degenerate principal values, of course. A defect with an orientation may have one eigenvalue unequal to the other two.

We will call the creation energy in the absence of external stress $E_{c[r]}^{[u]}(0)$. In the presence of an external stress, the creation energy must include the work required to distort the solid in opposition to that stress, so that

$$E_{c[r]}^{[u]}(\sigma) = E_{c(r)}^{[u]}(0) + \text{Tr}(\Omega_{c(r)}^{[u]} \cdot \sigma) \tag{55}$$

where the stress is evaluated locally. We have assumed here that all of the defect internal states are energetically degenerate in the absence of stress. The stress can break the degeneracy, depending on the orientation of the principal axes of $\Omega^a_{c(v)}$ relative to the stress tensor. These values of $E_{c(v)}(\sigma)$ are used in determining the solubility factor $S$ (Eqs. [46] and [53]).

Similarly, when the system is at the saddlepoint $[AaBb]$, the shape is different from the reference condition (the perfect, stressed lattice), which is represented by the saddlepoint volume tensor $\Omega^{|AaBb|}_{c(s)}$ and the energy at the saddlepoint becomes, under stress,

$$E^{|AaBb|}_{c(s)}(\sigma) = E^{|AaBb|}_{c(s)}(0) + \text{Tr}(\Omega^{|AaBb|}_{c(s)} \cdot \sigma) \tag{56}$$

where the reference state for the saddlepoint volume is the perfect, stressed crystal, just as it was for the valley volume. The values of $E_{c(s)}(\sigma)$ are used in determining the permeability tensor $\mathbf{P}$ (Eqs. [41], [45], and [50]).

Often the effects of hydrostatic stress (pressure) are separated from those due to deviatoric stress. The deviatoric stress is defined as the traceless part of the stress,

$$\sigma = p\mathbf{1} + \sigma_{dev}$$
$$p \equiv \frac{Tr[\sigma]}{3} \tag{57}$$

Similarly, for the volume tensor,

$$\Omega = (\Omega_h/3)\mathbf{1} + \Omega_a$$
$$\Omega_h \equiv Tr[\Omega] \tag{58}$$

where $\Omega_h$ would be identified as the total scalar volume change and $\Omega_a$ is the traceless part of $\Omega$. We can see that the work against the external stress has two terms: one couples the pressure to the total volume change, and the other couples the deviatoric stress to the anisotropic part of the saddlepoint volume,

$$\text{Tr}(\sigma \cdot \Omega) = p\Omega_h + \text{Tr}(\sigma_{dev} \cdot \Omega_a) \tag{59}$$

From this it is clear that the permeability factor in general will have an overall scalar factor which depends on the pressure and the isotropic saddlepoint volume ($\exp(-\beta p\Omega_{c(s)h})$). The anisotropic part of the saddlepoint volume, along with the deviatoric part of the stress tensor, will determine the anisotropic part of the permeability tensor (sum over terms involving $\exp(-\beta \; Tr[\Omega_{c(s)a} \cdot \sigma_{dev}])$).

Specific examples of the permeability tensor, including the effects of stress, are worked out and displayed in Section 3.3.

## 2.1.4. Limitations of the Continuum Approach

Diffusion-reaction schemes fail when the impurity concentration exceeds about $2 \times 10^{21}$ cm$^{-3}$. For such concentrations, sharp increases of the diffusivities of germanium, tin, arsenic, and antimony were reported (see Chapter 4 in Ref. [25]). Following the suggestion of Mathiot and Pfister [47], they are usually interpreted within the percolation theory of Stauffer [48] to arise from the proximity of the dopant atoms which reduces the formation and migration enthalpies of vacancies in their vicinity and which leads to an enhanced diffusion of dopants within the percolation cluster.

Most commercial process simulation tools until very recently had only implemented the most simple diffusion models using all or most of the simplifications discussed in Section 2.1.1, which break down under many conditions typical of nanofabrication. There, metastable dopant distributions are created by tight thermal budgets with steep temperature ramps and short soak times to beat the solubility limits of the dopants, thus increasing the importance of small intermediate clusters of dopants and defects, which, under more traditional conditions, would not be seen once the thermal equilibrium had been reached and the expected equilibrium precipitates had formed. The density of these "nonequilibrium" point defect and dopant-defect clusters is often quite low and their gradients are very steep, such that their representation by a continuous average concentration becomes questionable. As it

is currently understood the number of different configurations of point defects and dopants that need to be accounted for is rather large. Their binding and activation energies vary significantly, in particular for the small clusters with less than a dozen atoms. If the physically essential set of many equations considering all these defect-dopant clusters is treated, which is usually necessary to capture the nanometer-scale effects essential to nanoelectronic devices, besides the increased numerical effort, a very large number of parameters needs to be determined, which only in very few cases and the availability of many experimental results can be determined by a traditional straightforward fit [49].

Even if a simplification of the equation system would be desired for such annealing processes far from equilibrium, the questions about the "physically essential" set of partial differential equations (for *ad-hoc* forming larger clusters) cannot be answered easily. In the past few years, vast progresses in the field of atomistic and especially first-principles calculations have been found to help with both the equation selection and the parameter determination (see Sections 2.3 and 3.1), which is one of the seminal areas that has opened the door for successful process simulations of nanoscale electronic devices. Nevertheless, such an approach is of course not very flexible and requires a long start-up time for new materials. Both the fitting or calculation of the large number of parameters required by such systems and the solution of these large, highly nonlinear and tightly coupled systems of PDEs, as well as ordinary differential equations (ODEs) and algebraic equations require a significant amount of computational resources and to a large degree cut into the major advantages of the continuum approach, which are its computational efficiency and numerical stability.

Under these circumstances, atomistic process modeling starts to become much more attractive and has in fact recently been incorporated even into commercial process-simulation platforms [50]. While we expect the well-established PDE models to be much more efficient than the Monte Carlo models for one-dimensional (1D) and two-dimensional (2D) simulations, the sub-100 nm devices have a number of inherently three-dimensional effects, and therefore, increasingly require three-dimensional (3D) simulation. The need for 3D simulation and shrinking device sizes drastically reduces the CPU time gap between the two approaches. The recently reported CPU times for the kinetic Monte Carlo diffusion simulations [50, 51] look very promising compared to 3D PDE simulations for realistically complex sets of equations on sufficiently fine meshes.

## 2.2. Atomistic Process Simulation

Besides the limitations to the continuum approach outlined in Section 2.1.4, there is one more reason that makes atomistic modeling increasingly attractive and considerably drove its development. This reason consists in the considerable statistical variations due to the small number of dopant atoms in deca-nanometer device structures [52]. Figure 11 shows a 50-nm
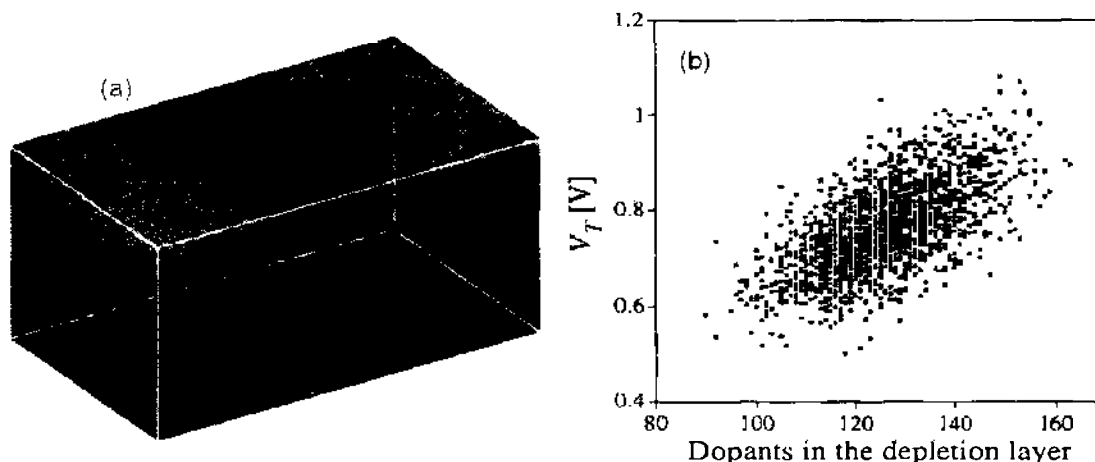


Figure 11. (a) 50-nm MOSFET with randomly distributed dopant atoms in the substrate and in the source and drain regions. (b) Simulated threshold voltages $V_T$ for statistical distributions of the dopant atoms versus the number of dopant atoms in the depletion region. Reprinted with permission from [52]. A. Asenov. "SISPAD 2001 Proceedings." Springer, Berlin. 2001. p. 162. © 2001, Springer.

MOSFET, which, at the current rate of miniaturization, could go into production within the next decade. For such a device, the depletion region of the channel—which is the area where the majority of the current flows when the device is on—contains approximately between 90 and 160 dopant atoms, where the variation is of statistical nature due to the implantation process. As Fig. 11 shows, the number of dopant atoms in the depletion region of course has a significant influence on the threshold voltage, which is the minimum gate voltage required to switch the device (i.e., the current between source and drain) on. However, even when the number of dopant atoms were constant, their statistical distribution still causes a large variation in the threshold voltage of easily 30% for a given number of dopant atoms. Thus, overall the threshold voltage for nominally identical devices on a chip can vary from 0.5 to ~1.1 V. Under these conditions, it is hard to reliably turn all transistors of a given structure on and off, which, for example, could require architectures that can function despite a certain number of statistically distributed malfunctioning transistors. To quantify such effects, an atomistic description of the doping process becomes indispensable.

Different techniques are available on the atomic length scale, ranging from the fastest, but least reliable kinetic Monte Carlo methods to the slowest, but most accurate molecular dynamics (MD) methods. The latter one might be the most tempting to try, since it is a well established and (when using reliable models for the interatomic forces) accurate methodology to determine the time evolution of an atomic (or nearly atomic) scale system, including the chemical reactions and/or diffusion hops in that system. As described in Section 2.1.1, this is the goal of process simulation for nanoscale devices.

Thus, why does not everybody just use straightforward MD for nanoscale process simulation? The answer is that due to the fact that MD simulates all the lattice atoms and, more importantly, that it has to use an almost constant time step on the order of femtoseconds $(10^{-15}$ s) to be able to follow the trajectories of the atoms (which complete a Brownian-motion "orbit" around their equilibrium position in $10^{-13} - 10^{-12}$ s), it cannot simulate the length and more importantly time scales involved in typical technological processing steps (tens to hundreds of nanometers for seconds to hours). The problem is that the transitions of interest are typically many orders of magnitude slower than vibrations of the atoms, so a direct simulation of the classical dynamics ends up being of little use.

This "rare event" problem is best illustrated by an example. A typical, low activation energy for a chemical reaction or diffusion event is 0.5 eV. Such an event can occur thousands of times per second at room temperature and would typically be important in the time evolution of the system. But, the atoms vibrate on the order of $10^{10}$ times before a sufficiently large fluctuation of thermal energy occurs in the right direction for a transition to take place. A direct classical dynamics simulation which necessarily has to faithfully track all this vibrational motion would take thousands of years of computer calculations on the fastest present day computer before a single transition can be expected to occur. Thus, meaningful studies of chemical reactions and/or diffusion cannot be carried out by simply simulating the classical dynamics of the atoms. It is essential to simulate the system on a much longer time scale. This time scale problem is one of the important challenges in atomic-scale process simulation for nanoscale systems.

## 2.2.1. Transition State Theory

The dynamical evolution of an infrequent (or rare) event system consists of vibrational excursions within a potential basin, punctuated by occasional transitions between basins; these transition events are infrequent in the sense that the average time between events is many vibrational periods. If a bottleneck region through which the system must pass in order to make the transition can be identified, the so-called transition state, then transition state theory (TST) [53–59] can be used to separate the time scales between vibrations and transitions and statistically calculate the average amount of time the system would spend in a given state.

In the transition state theory approximation, the classical rate constant for escape from state $A$ (where the system has stayed for many periods of thermal motion) to some adjacent state $B$ is taken to be the equilibrium flux through the dividing surface between $A$ and $B$ (Fig. 12). Although on average such crossings are infrequent, successive crossings can sometimes occur within just a vibrational period or two; these are termed correlated dynamical events.
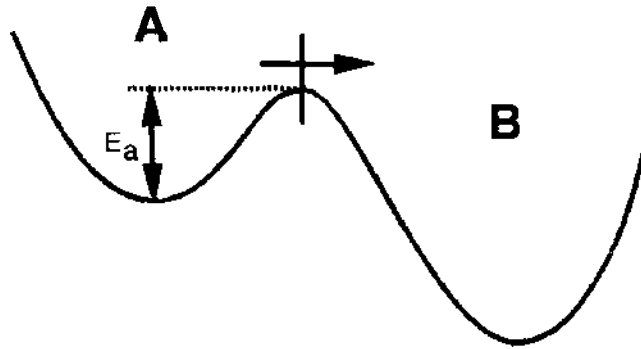
**Figure 12.** A two-state system that illustrates the definition of the transition state theory rate constant as the outgoing flux through the dividing surface bounding state *A*. Reprinted with permission from [60], A. F. Voter et al., *Annu. Rev. Mater. Res.* 32, 321 (2002). © 2002, Annual Reviews.

An example would be a double jump of a vacancy in silicon. For the following, it is sufficient to know that such correlated events exist, although in more or less all of the methods presented below, we assume that they do not occur (this is the primary assumption of transition state theory), which is actually a very good approximation for many solid-state diffusive processes [60]. Apart from assuming no recrossing events and possibly the Born-Oppenheimer approximation for electronic-structure methods, transition-state theory requires only one additional basic assumption, which is that the rate is slow enough so that a Boltzmann distribution is established and maintained in the reactant state. If there are no correlated dynamical events, the transition state theory rate is the exact rate constant. In this description, state *A* and *B* can be the initial and final atomic configuration of a diffusion or reaction event, respectively, and would be characterized by all coordinates of all atoms in the system.

To explain the equilibrium flux concept of transition state theory in more detail, imagine that for a two-state system we run a long classical trajectory, weakly coupled to a heat bath to guarantee on average constant system temperature. We run the trajectory long enough to establish equilibrium, visiting both states an extremely large number of times. By examining this trajectory, we could accurately determine the fraction, $\chi_A$, of the time it spends in state *A* and the number of crossings, per unit time, of the dividing surface. The transition state theory rate constant for escape from *A*, $k_{A \to}^{TST}$, would then be the number of crossings that are exiting state *A* divided by $\chi_A$.

The real beauty of transition state theory, though, is that because this flux is an equilibrium property of the system, we do not need to propagate a trajectory, but can simply compute the transition state theory rate constant from

$$k_{A \to}^{TST} = \left\langle \left| \frac{dx_1}{dt} \right| \delta(x_1 - q) \right\rangle_A \tag{60}$$

Here the angular brackets indicate the grand canonical ensemble average, i.e., for some property $P(\mathbf{r}, \mathbf{p})$.

$$\langle P \rangle = \frac{\int\int P(\mathbf{r}, \mathbf{p}) \exp[-\beta H(\mathbf{r}, \mathbf{p})] d^3r d^3p}{\int\int \exp[-\beta H(\mathbf{r}, \mathbf{p})] d^3r d^3p} \tag{61}$$

where $\mathbf{r}$ and $\mathbf{p}$ include all atoms and have $3N$ components each and $\beta = 1/(k_B T)$. The subscript *A* in Eq. (60) indicates that the configuration space integrals are restricted to the space belonging to state *A* (eliminating the need to divide by $\chi_A$), and the dividing surface, for simplicity here, is at $x_1 = q$, involving only the reaction coordinate $x_1$ ($x_1 \in \mathbf{r}$). If the effective mass *m* of the reaction coordinate is constant over the dividing surface, Eq. (60) reduces to a simpler ensemble average over configuration space only [58].

$$k_{A \to}^{TST} = \sqrt{\frac{2k_B T}{\pi m}} \langle \delta(x_1 - q) \rangle_A \tag{62}$$

The essence of this expression, and of transition state theory, is that the Dirac delta function picks out the probability of the system being at the dividing surface, relative to everywhere

it can be in state $A$. Note that there is no dependence on the nature of the final state $B$. Evaluating Eq. (62) to find the transition state theory rate for a given temperature is relatively straightforward [60], but for the purpose of nanoscale process simulation. is hardly used without the addition of the harmonic approximation.

**Harmonic transition state theory.** In general, the free energy of the transition state needs to be evaluated to estimate the rate. However, for solid systems, where the atoms are relatively tightly held in place by their neighbors, it is often possible to assume that the transition state can be characterized by a few saddle points on the potential energy rim surrounding the initial state basin and that the partition function of the system (denominator of r.h.s. of Eq. (61)) near each saddle point and near the energy minimum can be approximated by a product of harmonic partition functions. In this case, transition state theory simplifies to the harmonic transition state theory (hTST) and the rate of escape, $k$, through each of the saddle point regions can be related to properties of the initial state energy minimum and the saddle point [61, 62] as

$$k^{\mathrm{hTST}} = \nu_0 \exp\left(-\frac{E_a}{k_B T}\right) \tag{63}$$

with

$$\nu_0 = \frac{\prod_i^{3N} \nu_i^{(m)}}{\prod_i^{3N-1} \nu_i^{(s)}} \tag{64}$$

and

$$E_a = E_{(s)} - E_{(m)} \tag{65}$$

As in Section 2.1.3. $E_{(s)}$ is the energy of the saddle point, $E_{(m)}$ is the energy of the local minimum corresponding to the initial state, and the $\nu_i$ are the corresponding normal mode frequencies. All the quantities can be evaluated directly from the potential energy surface without dynamical calculations as described in Section 2.3. Entropic and thermal effects are included through the harmonic partition functions.

With the use of transition state theory or harmonic transition state theory, a long time scale simulation consists of identifying states of the system and finding the mechanism and rate of transitions from a current state to a new state. The key thing is to find the relevant mechanism and not just assume a mechanism. Often, preconceived notions of the transition mechanism have turned out to be incorrect. One example is the mechanism of adatom diffusion on Al(100) where a two atom concerted displacement process turned out to have a lower barrier than the simple hop mechanism [63]. When harmonic transition state theory is used. the most challenging part of the calculation is the search for the low lying saddle points without knowledge of the possible final states. A typical simulation system includes a hundred atoms or many more, which means that saddle points need to be found in a space with at least several hundred degrees of freedom. The large number of degrees of freedom makes this a challenging problem. In Section 2.3. we will review various approaches that have been taken to address this issue.

### 2.2.2. The Kinetic Monte Carlo Approach

Kinetic Monte Carlo (KMC) [64–68] is the most economical atomistic method to extend the time scale of simulations far beyond the vibrational time scale. In combination with the widely used Monte Carlo implantation models, it has been shown to offer the possibility to directly investigate statistical variations for devices with a small number of dopant atoms in the active regions such as the one from Fig. 11 and has already been implemented into commercially available TCAD platforms [50].

In a kinetic Monte Carlo simulation. it is assumed that a list of "all" possible transitions for "each possible" initial state is available. Then, a random number can be used to choose one of the transition processes and evolve the system to a new state. The probability of choosing a certain transition is proportional to its rate, $r_i$. On average. the amount of time that would

have elapsed in order for this process to occur is usually (although not exclusively) assumed to be

$$\Delta t = \frac{1}{\sum r_i} \tag{66}$$

which is independent of the chosen transition. It may also be important to include the appropriate distribution of escape times. For random uncorrelated processes, this is a Poisson distribution. If $\mu$ is a random number from 0 to 1, the elapsed time for a particular transition is given by

$$\Delta t = \frac{-\ln \mu}{\sum r_i} \tag{67}$$

The system is then advanced to the final state of the chosen transition and the process is repeated.

The kinetic Monte Carlo method is thus a purely event-driven technique, that is, it simulates only infrequent events (e.g., diffusion hops) at random, with probabilities according to their respective event rates. In this way it self-adjusts the time step as the simulation proceeds, just to be able to account for the fastest event present at that time. The advantage is that no thermal motion of the atoms is simulated at all, which does not contribute to the structural evolution of the system.

The drawback is that in a traditional kinetic Monte Carlo simulation, all transitions that can ever occur in the system, along with their rates, must be known before the simulation starts. Ideally, the rates are estimated from some description of the atomic interactions [69] such as an empirical interaction potential or *ab initio* calculations, but the problem is to know in advance the mechanism of the relevant transitions for each possible configuration of the atoms. The requirement of knowing and tabulating the relevant transitions ahead of time limits the method to simple systems, or to more complicated systems without a complete event catalog (and thus a limited reliability of the simulation results). Systems that can undergo complicated transitions involving several atoms, such as the examples in Section 3.2.3, or where atoms do not sit at lattice sites are extremely difficult to model with traditional KMC. The selection of the event catalog and the corresponding parameters is similar to the task of defining a sufficient set of diffusion-reaction equations and determining their parameters (as discussed in Section 2.1.1). Like for the continuum technique, this is the most severe bottleneck of the kinetic Monte Carlo technique, since an incomplete or wrong event catalog is a severe hazard for the validity of the simulation results. However, assisted by plentifully available experimental data and the recent increase in *ab initio* data that helped to better understand the importance of the different reactions and hopping mechanisms as well as to define the corresponding kinetic parameters, kinetic Monte Carlo descriptions of the thermal evolution of a system upon annealing [70] have been found to be quite successful for process simulations in nanoscale CMOS structures.

Figure 13 illustrates the concept of the kinetic-Monte Carlo approach for applications in the realm of nanoscale CMOS devices [70]. The figure shows a high resolution TEM view [71] of a silicon sample with a {311} extended defect embedded in the silicon atomic rows. A {311} defect is an extended cluster of self-interstitials added into the lattice in a way that all atoms have fourfold coordination, which gives these defects a low-formation energy [72]. In this situation, the kinetic Monte Carlo technique simulates only the atoms belonging to point or extended defects (represented as circles on the TEM view) and ignores all "perfect" lattice atoms and their migration and reactions with each other. In this example, the infrequent events (on a time scale of maybe nanoseconds) consist of isolated point defects jumping to a neighboring position where then a reaction might happen. At even longer time intervals (e.g., every millisecond) a point defect would be emitted from the extended defect. To simulate this, the kinetic-Monte Carlo simulation starts with timesteps of $10^{-9}$ s to follow the isolated defect atoms. The fast-moving point defects disappear very quickly (being caught by the extended defects), leaving only the extended defects, which would allow to raise the timestep automatically to $10^{-3}$ s, thus allowing for a very fast simulation of the system evolution.
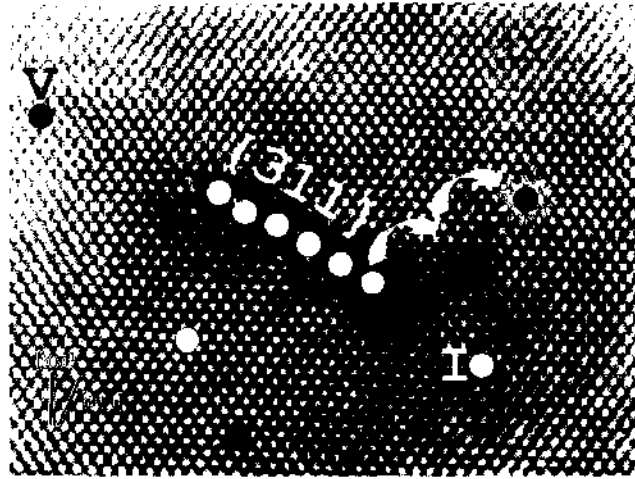
Figure 13. TEM view of an extended {311} defect in silicon [71]. Circles represent the (defect) atoms simulated by the KMC method. As an example, the figure illustrates the emission of an interstitial ($I$) from the {311} defect and its eventual recombination with one of the vacancies ($V$) present in the sample. After Fig. 3. Reprinted with permission from [73]. M. E. Law et al., *Mater. Res. Soc. Bull.* 25, 45 (2000). © 2000. Materials Research Society.

A key component in a KMC simulator is the event manager or "scheduler," which is the procedure responsible for selecting the random events according to their event rates. Figure 13 illustrates the selection procedure for a configuration consisting of three vacancies ($V$), 2 interstitials ($I$) and one {311} defect. If, for example, we assume that the vacancy and interstitial jump rates are 1000 s$^{-1}$ and 100 s$^{-1}$, respectively, and the {311} emission rate is 10 s$^{-1}$, we would expect on average during each second $3 \times 1000 = 3000$ vacancy hops, $2 \times 100$ interstitial hopes, and 10 {311} emission events, or a total of 3210 events. In addition, we have to pick a $V$ with a probability proportional to 3000/3210, an $I$ with a probability proportional to 200/3210, and the {311} with a probability proportional to 10/3210.

A KMC simulator consists of a 3D simulation box, of dimensions ranging from tens of nanometers to a few microns, where point defects can be created and allowed to jump and interact. Each defect or dopant atom is represented by a point, with coordinates $(x, y, z)$ and an interaction radius $a_c$ (equivalent to the capture radius in continuum simulations, see Eq. (23) in Section 2.1.2) in a 3D simulation box. Besides the straightforward simulation of diffusion hops, reactions are simulated by either the transformation of existing point defects to a different species or by agglomerating jumping point defects to form a larger reaction product like the {311} shown in Fig. 13 [73]. Extended defects can be modeled replicating the actual geometry of the defect when known: vacancies can form spherical clusters (voids), interstitials can grow elongated stripes ({311}'s), or planar dislocation loops and stacking faults. Modeling the actual geometry improves the accuracy of the emission and capture rates, and can become essential in the proximity of the surface. A free surface is treated as a particular type of extended defect, and is included in the scheduler with $I$ and $V$ emission rates derived from the surface area and the formation energy of the corresponding point defect. For extended defects, the interaction or capture region can be assumed to be the superposition of the capture spheres of the constituent particles [70]. Finally, each type of capture process can have an associated capture energy barrier. The calculation of realistic capture radii, which have been found to be considerably longer than traditionally assumed, is discussed in Section 3.1.2.

For the event rates, one usually uses the expression of harmonic transition state theory from Eq. (63). As an example, the interstitial jump rate within harmonic transition state theory would be given by

$$ J_{\text{HTST}}^I = 6 \frac{D_0}{\lambda^2} \exp\left( -\frac{E_{\text{mig}}^I}{k_B T} \right) \tag{68} $$

where $D_0$ is the diffusivity prefactor, $\lambda$ is the hopping length, and $E_{\text{mig}}^I$ is the interstitial migration energy.

In addition to point defects, a full KMC simulator needs to implement models for a variety of extended defect types like the above mentioned surfaces, clusters and complexes (e.g., for the formation of voids or interstitial clusters, see Section 3.4.1) and thus runs into similar problems as we had discussed for continuum models, which are to identify the important complexes, their possible and probable reactions with other defects, and the necessary parameters to describe them. As an example for the difficulty to reliably determine such complex parameters, Fig. 14(a) is a plot of the activation energy for the emission of an interstitial from interstitial clusters of different sizes, as fitted to experimental measurements [74]. However, atomic-level *ab-initio* calculations (Fig. 14(b)) did not find the sharp minima for four and eight-atom clusters, but more a smooth dependence of the formation energy on the cluster size [75, 76]. This disagreement leaves a large uncertainty concerning the values of the clustering energies, since both approaches extract the parameters in phase spaces with large numbers of degrees of freedom. On the contrast, fitting a large set of experimental annealing data [77] to determine the formation energies of boron-interstitial complexes (BIC's) $B_m I_n$ led to results that agreed very well with *ab initio* calculations ([12], see Section 3.2.2.2).

## 2.2.3. Accelerated Dynamics Methods

The just discussed uncertainties in event tables and kinetic rates and thus the often decreased credibility of kinetic Monte Carlo simulations are in many cases very unsatisfactory and make the use of methods desirable that can follow the long-time trajectory of a system more reliably. The motivation that led to the development of accelerated molecular dynamics methods becomes particularly clear when we try to understand the dynamical evolution of what we term complex infrequent event systems, where we simply cannot guess where the state-to-state evolution might lead. The underlying mechanisms may be too numerous, too complicated, and/or have an interplay whose consequences are unpredictable. While in very simple systems we can raise the temperature to make diffusive transitions occur on an MD-accessible time scale, in more complex systems this strategy will cause the system to travel down a different path in state space. Ultimately, this will lead to a completely different kind of system, making it impossible to address the questions that the simulation was attempting to answer.

Often, even systems that seem very simple can turn out to be in this complex class. For example, until 1990, we did not know, nor did we expect, that an adatom on the fcc(100) surface diffused in any way other than by a simple hop mechanism. Thus, in a Monte-Carlo simulator to model physical vapor deposition of Al interconnects, the hop mechanism would be the only diffusion mechanism in the event table of an isolated surface adatom. The exchange mechanism [63, 78, 79], involving the adatom and a substrate atom (see Fig. 15) is now known to be the preferred mechanism for diffusion on fcc(100) surfaces for many metals (e.g., Al, Pd, Pt, and Au), and the surface science community has since discovered a
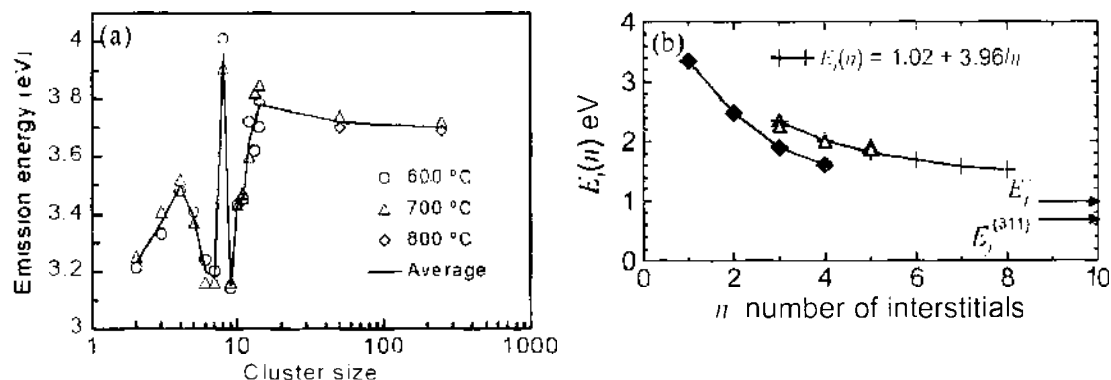


Figure 14. (a) Emission energies of interstitial clusters as a function of the size, as derived from experimental measurements at three different temperatures. Reprinted with permission from [74]. N. E. B. Cowern et al., *Phys. Rev. Lett.* 82, 4460 (1999). © 1999, APS. (b) Dependence of the interstitial-cluster formation energy (per interstitial) on the number of interstitials, $n$. The formation energy of an isolated interstitial, $n = 1$, is 3.35 eV. The compact cluster (diamonds) is only initially more stable than the elongated (triangles) cluster. The trend equation displayed in the figure shows the growing stability of the elongated cluster for increasing number of interstitials. In particular, for "infinite" interstitial defects, the formation energies for the chain defects is only 1.02 eV. Reprinted with permission from [75]. J. Kim et al., *Phys. Rev. Lett.* 84, 503 (2000). © 2000, American Physical Society.
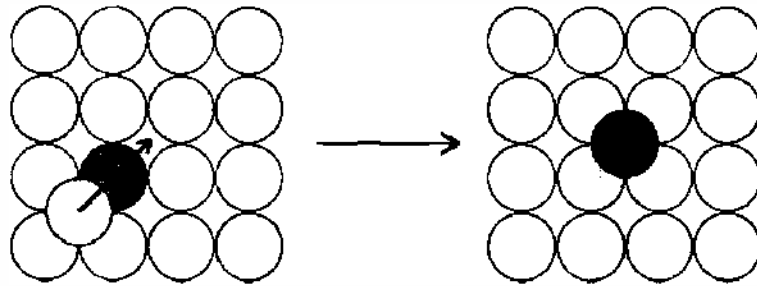
**Figure 15.** Adatom exchange mechanism on an fcc(100) surface. This, as opposed to a jump, is the preferred adatom diffusion mechanism on a number of fcc metals, demonstrating the complexity of even simple atomistic systems. Reprinted with permission from [60], A. F. Voter et al., *Annu. Rev. Mater. Res.* 32, 321 (2002). © 2002, Annual Reviews.

large variety of multiple-atom concerted diffusion mechanisms [80–87]. Many, if not most, materials problems fall into this complex infrequent-event system category.

In the past few years, several so-called accelerated-dynamics methods have been developed that have helped significantly to decrease the time scale problem. For systems in which the long-time dynamical evolution is characterized by a sequence of activated events like diffusive events or reactions), these methods can extend the accessible time scale by orders of magnitude relative to direct MD, while retaining full atomistic detail.

The class of these "accelerated dynamics methods" includes a variety of approaches, such as parallel replica dynamics, hyperdynamics, and temperature accelerated dynamics [60], which we will discuss in the following. The conceptual link between these methods is that the system trajectory, caught in its current state, is stimulated to find an appropriate path for escape more quickly than it would with direct MD. As in MD, no *a priori* information about what this escape path might look like is imposed on the procedure; the trajectory simply finds its own way out of the state. As our focus is on methods that can extend the MD simulation time in an accurate way for nanoelectronics processing problems, we make no attempt to discuss the large body of related work involving enhanced sampling methods and approximate dynamical approaches (e.g., for macromolecule systems). Excellent reviews on accelerated-dynamics methods can be found in [88] and [60].

***Parallel replica dynamics.*** The parallel replica method [89] is the simplest and most accurate of the accelerated dynamics techniques, with the only assumption being that of infrequent events obeying first-order kinetics (exponential decay); that is, for any time greater than $\tau_{corr}$ after entering a state, the probability distribution function for the time of the next escape is given by

$$p(t) = k \exp(-kt) \tag{69}$$

where $k$ is the rate constant for escape. Starting with an $N$-atom system in a particular state (basin), the entire system is replicated on each of $M$ available parallel or distributed processors. After a short dephasing stage ($\Delta t_{deph}$), during which the velocities of the atoms are periodically randomized to eliminate correlations between replicas, each processor carries out an independent constant-temperature MD trajectory for the entire $N$-atom system, thus exploring the phase space within the particular basin $M$ times faster than a single trajectory would. Whenever a transition is detected on any processor, all processors are alerted to stop. The simulation clock is advanced by the accumulated trajectory time summed over all replicas, that is, the total time spent exploring phase space within the basin until an escape pathway is found. It is readily shown [89] that this procedure gives an escape time that is correctly drawn from the distribution in Eq. (69), even if the processor speeds are inequivalent, allowing the use of heterogeneous clusters or widely distributed machines running parallel replica dynamics as low-level background processes [90].

The parallel replica method also correctly accounts for correlated dynamical events (that is, there is no requirement that the system obeys transition state theory), unlike the other two methods presented here. This is accomplished by allowing the trajectory that made the transition to continue on its processor for a further amount of time $\Delta t_{corr}$, during which recrossings or follow-on events may occur. The simulation clock is then advanced by $\Delta t_{corr}$, the final state

is replicated on all processors, and the whole process is restarted. This overall procedure then gives exact state-to-state dynamical evolution because the escape times obey the correct probability distribution: nothing about the procedure corrupts the relative probabilities of the possible escape paths, and the correlated dynamical events are properly accounted for.

The efficiency of the method is limited by both the dephasing stage, which does not advance the system clock, and the correlated event stage, during which only one processor accumulates time. Thus, the overall efficiency will be high when

$$\frac{t_{\text{escape}}}{M} \gg \Delta t_{\text{deph}} + \Delta t_{\text{corr}} \tag{70}$$

An example for the use of the parallel replica method for the study of self-interstitial clustering in silicon is discussed in Section 3.4.2.

**Hyper-MD.** Historically, the first approach to accelerate MD simulations was to decrease the probability that the system is found in an initial state by adding a repulsive potential energy, a bias potential, to the actual interaction potential in such a way as to increase the probability of finding the system at a transition state [91]. The derivation of the method requires that the system obeys transition state theory; that is, it assumes there are no correlated events. The bias potential must be zero at all the dividing surfaces, and the system must still obey transition state theory for dynamics on the biased potential. A trajectory on this modified potential, while relatively meaningless on vibrational time scales, evolves correctly from state to state at an accelerated pace. Moreover, the accelerated time is easily estimated as the simulation proceeds. For a regular MD trajectory, the time advances at each integration step by $\Delta t_{\text{MD}}$, the MD time step (e.g., $\sim 1$ fs). In hyperdynamics, the time advance at each step is $\Delta t_{\text{MD}}$ multiplied by an instantaneous boost factor, the inverse Boltzmann factor for the bias potential at that point, so that the total time after $n$ integration steps is

$$t_{\text{hyper}} = \sum_{j=1}^{n} \Delta t_{\text{MD}} \exp\left(\frac{\Delta V[\mathbf{r}(t_j)]}{k_B T}\right) \tag{71}$$

The ideal bias potential should give a large boost factor, should have low computational overhead (although more overhead is acceptable if the boost factor is very high), and should to a good approximation meet the requirements given above. This is challenging because we want, as much as possible, to avoid utilizing any prior knowledge of the dividing surfaces or the available escape paths. The bias potentials in the first hyperdynamics paper [92] were based on the lowest eigenvalue ($\epsilon_1$) of the Hessian matrix, $\{\partial^2 V/\partial x_i \partial x_j\}_{ij}$. The bias potential was made positive for regions where $\epsilon_1 > 0$, and zero elsewhere, exploiting the fact that $\epsilon_1$ is positive near the bottom of a basin and negative at saddle points. For a periodic two-dimensional example system, this gave substantial boosts (in the thousands when $k_B T$ was $\sim 1/20$ of the barrier height) and excellent accuracy, even when some recrossings were present. Since this bias potential required a diagonalization of the full $3N$-dimensional Hessian at every time step, it became prohibitively expensive when the number of atoms was increased beyond a few tens of atoms [60]. An improved bias potential with a better saddle-point detection algorithm applicable to lower-temperature simulations was later developed [93], which produced a boost factor of 8310 (with a computational overhead of $\sim 30$) in a 221.2 $\mu$s simulation of room-temperature diffusion of a ten-atom cluster of silver on Ag(111), a system with 70 moving atoms.

The bias-potential idea can be easiest demonstrated with the very intuitive simple flat boost potentials (Fig. 16), as suggested in Refs. [94, 95]. They require much less effort for the potential construction but run quickly out of steam for systems with more than a few degrees of freedom (= 3 × number of moving atoms). As an example, the boost factor is already below two for a system with only 20 degrees of freedom [88]. This is because for systems with a larger number of degrees of freedom, the probability of finding the potential energy of the system below a saddle point energy becomes vanishingly small, since each additional degree of freedom brings in $\frac{1}{2}k_B T$ of kinetic energy distorting the potential energy. Despite the large total kinetic energy of larger systems, it is still of course a rare event to find enough energy focused in the right degree of freedom to bring the system through the saddle point region.
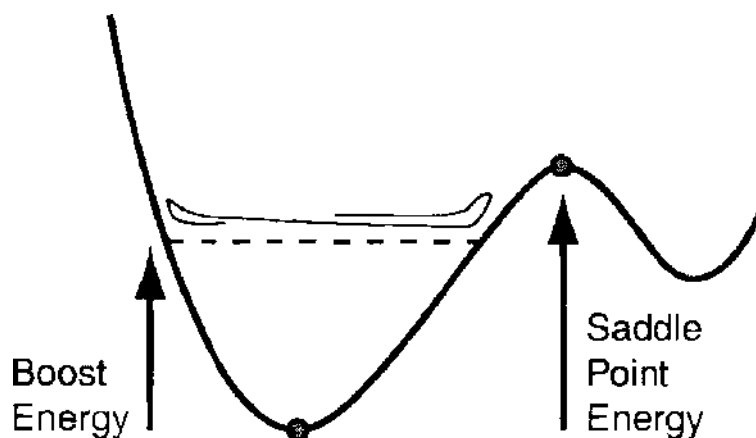
***Temperature accelerated dynamics (TAD).*** Another, perhaps simpler approach for identifying transitions is to increase the rate of rare events by simply heating the system. If the atoms have more energy, they will more likely undergo transitions. One should, however, not expect the favored transition mechanism to be the same at the higher temperature. This situation is illustrated in Fig. 17. The temperature dependence of the rate of two possible processes the system can undergo from a given initial state is shown. One of the processes has a low activation energy and small prefactor while the other has a high activation energy and large prefactor. At low temperature, the low barrier process will have a higher rate and dominate the dynamics. At high temperature, entropy becomes more important and the process with higher prefactor dominates even though the energy barrier is higher.

An even more serious problem arises when the thermodynamic state of the system changes as it is heated up high enough to make transitions observable in a short, classical dynamics simulation. An example of this is diffusion of admolecules on an ice surface. If an ice slab is heated up to a temperature at which the diffusion events occur frequently enough, the surface melts and the diffusion mechanism becomes quite different from the low temperature, long time scale diffusion mechanism.

High temperature dynamics can, however, in many non-pathological cases be used to search for the relevant mechanism if many searches are carried out. In the method proposed in Ref. [96], the system is evolved at a high temperature $T_{high}$ in each basin (while the temperature of interest is some lower temperature $T_{low}$). Whenever a transition out of the basin is detected, the saddle point for the transition is found, e.g., using the nudged elastic band method (see Section 2.3.2). The trajectory is then reflected back into the basin and continued. This basin-constrained molecular dynamics procedure generates a list of escape paths and attempted escape times for the high-temperature system. Assuming that transition state theory holds and that the system obeys first-order kinetics, the probability distribution for the first-escape time for each mechanism is again an exponential, Eq. (69). Because harmonic transition state theory gives an Arrhenius dependence of the rate on temperature (Eq. [63]), depending only on the static barrier height, we can then extrapolate each escape time observed at $T_{high}$ to obtain a corresponding escape time at $T_{low}$ that is drawn correctly from the exponential distribution at $T_{low}$. This extrapolation, which requires knowledge of the saddle point energy, but not the pre-exponential factor, can be illustrated graphically in an Arrhenius-style plot $(\ln(1/t)$ versus $1/T)$, as shown in Fig. 17. The event with the shortest time at low temperature is the correct transition for escape from this basin. Because the extrapolation can in general cause a reordering of the escape times, a new shorter-time event may be discovered as the basin-constrained MD is continued at $T_{high}$. If we make the additional approximation that there is a minimum preexponential factor, $\nu_{min}$, which bounds from below all the pre-exponential factors in the system, we can define a time at which
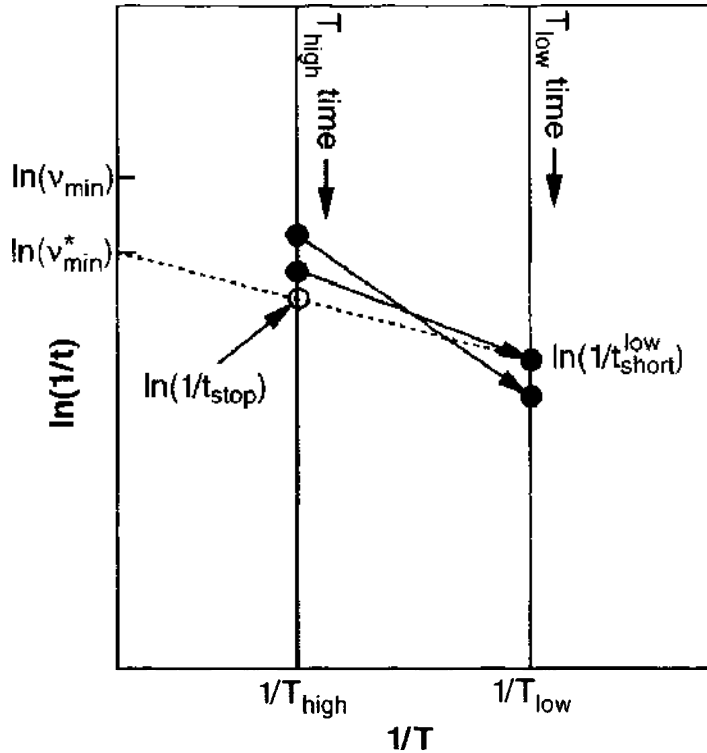
**Figure 17.** Schematic illustration of the temperature-accelerated dynamics method. Progress of the high-temperature trajectory can be thought of as moving down the vertical time line at $1/T_{high}$. For each transition detected during the run, the trajectory is reflected back into the basin, the saddle point is found, and the time of the transition (solid dot on left time line) is transformed (arrow) into a time on the low-temperature time line. Plotted in this Arrhenius-like form, this transformation is a simple extrapolation along a line whose slope is the negative of the barrier height for the event. The dashed termination line connects the shortest-time transition recorded so far on the low temperature time line (solid dot) with the confidence-modified minimum pre-exponential ($\nu_{min}^* = \nu_{min}/\ln(1\delta)$) on the $y$-axis. The intersection of this line with the high-$T$ time line gives the time ($t_{stop}$, open circle) at which the trajectory can be terminated. With confidence $1 - \delta$, we can say that any transition observed after $t_{stop}$ could only extrapolate to a shorter time on the low-$T$ time line if it had a pre-exponential lower than $\nu_{min}$. Reprinted with permission from [60]. A. F. Voter et al., *Annu. Rev. Mater. Res.* 32, 321 (2002). © 2002, Annual Reviews.

the basin-constrained MD trajectory can be stopped, knowing that the probability that any transition observed after that time would replace the first transition at $T_{low}$ is less than $\delta$. This stop time is given by

$$t_{high}^{stop} = \frac{\ln(1/\delta)}{\nu_{min}} \left( \frac{\nu_{min} \, t_{low}^{short}}{\ln(1/\delta)} \right)^{t_{low}/t_{high}}$$ (72)

where $t_{low}^{short}$ is the shortest transition time at $T_{low}$. Once this stop time is reached, the system clock is advanced by $t_{low}^{short}$, the transition corresponding to $t_{low}^{short}$ is accepted, and the TAD procedure is started again in the new basin. The average boost in TAD can be dramatic when $T_{high}/T_{low}$ is large. Any anharmonicity error at $T_{high}$ transfers to $T_{low}$; a rate that is twice the Vineyard harmonic rate owing to anharmonicity at $T_{high}$ will cause the transition times at $T_{high}$ for that pathway to be 50% shorter, which in turn extrapolate to transition times that are 50% shorter at $T_{low}$. If the Vineyard approximation is perfect at $T_{low}$, these events will occur at twice the rate they should. This anharmonicity error can be controlled by choosing a $T_{high}$ that is not too high. An excellent in-depth discussion of TAD can be found in [60].

### 2.2.4. Kinetic Monte Carlo Simulations without Lattice Approximation and Predefined Event Table

Although both the kinetic-Monte Carlo approach and accelerated-dynamics methods show great promise for the area of nanoscale process simulation, both have—mostly complimentary—advantages and disadvantages. As we said, the necessity of a predefined event table for the computationally very efficient Monte Carlo approach creates the danger

of an unphysical system evolution. The accelerated-dynamics methods, on the other hand, guarantee a physically correct trajectory and thus system evolution, but on the other hand are computationally less efficient and still have a hard time to reach the process relevant times $\gg \mu s$ (except for special systems with extremely high barriers). A very interesting approach to combine the two approaches to benefit from their advantages while minimizing their disadvantages has been proposed by Henkelman and Jónsson [97].

In their method, classical dynamics are not used in any form but instead the system is pushed up the potential energy surface using the so-called dimer method (see Section 2.3.3) to find saddle points. The rate of transitions through the vicinity of each saddle point is then estimated within harmonic transition state theory, creating an on-the-fly (more or less complete) event catalog, and a kinetic Monte Carlo method is used to simulate the evolution of the system over long time scales [98]. This method is easy to implement and, compared to existing methods, may require less computational time for small systems. While the use of harmonic transition state theory means that it can only be applied to solids, this method is still applicable to glasses and other amorphous solids, where traditional (lattice) Monte Carlo methods have difficulties.

Henkelman's and Jónsson's method is illustrated for a two-dimensional model potential in Fig. 18. The system is started at a potential minimum, $A$. When a new state is visited, a swarm of dimer searches (around 10–50 or so) is sent out from the vicinity of the potential energy minimum. In this example, ten random displacements from the position of the minimum were chosen as starting points of dimer searches. Figure 18(a) shows the path of the ten dimer searches. In this calculation, four distinct saddle points (*) were found. The system is then quenched on either side of each saddle point in order to verify that it lies on a
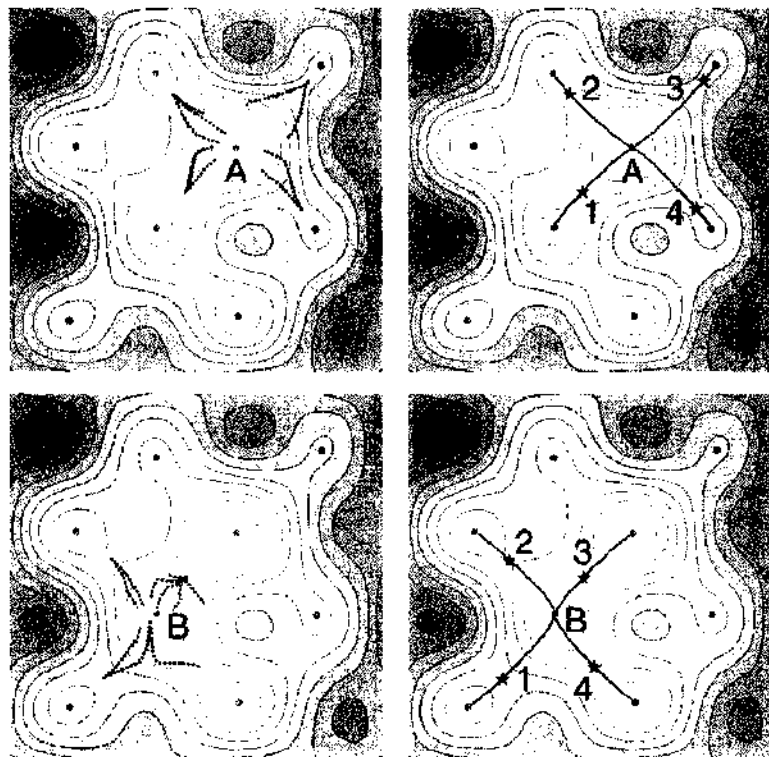


Figure 18. Application of the long time scale simulation method to a model two-dimensional potential surface. The system is initially in state $A$. (a) Ten dimer searches are started from random positions around the minimum. They converge on four distinct saddle points (two of the searches practically overlap). (b) The system is then made to slide down the minimum energy path on either side of the saddle points which are indicated with *. Here all four saddle points have a minimum energy path starting at the initial state minimum $A$, but this does not have to be the case. The rate of each process is then calculated using harmonic transition state theory. A process is chosen at random using the kinetic Monte Carlo algorithm. In this case, process 1 gets chosen. The system is moved to the final state of this process, to minimum $B$. (c) Dimer searches are run from the new minimum, again four distinct saddle points are found. (d) Minimum energy paths are traced out, and the process repeated. Reprinted with permission from [97], G. Henkelman, Ph.D. thesis, University of Washington, 2001. © 2001.

minimum energy path (shown in gray) from the given initial state minimum. All the saddle points found in this case did connect back to the initial minimum. If not, the saddle point is discarded from the list of possible transitions. In the same way as described in Section 2.2.2, a transition is chosen from the list, the system is advanced to the final state of that transition, and the time interval associated with the transition is added to the accumulated time. In this example, transition 1, which corresponds to the lowest barrier was chosen, and the system was advanced to state *B*. From the new minimum the process is repeated. New dimer searches are sent out (Fig. 18[c]), the saddle points verified (Fig. 18[d]) and then one is chosen for the next transition.

Because the novelty of this method, no applications within the field of nanoelectronic process modeling exist yet, with the exception maybe of the simulation of the evolution of radiation damage in silicon carbide [99].

## 2.3. Kinetic Parameters from Atomistic Simulations

After having discussed the different techniques for the evolution of a nanoelectronic system with time, the most important thing missing is the calculation of the transition rates for chemical reactions and diffusion. These calculations are usually performed within the limits of transition-state theory, as described in Section 2.2.1. Since atoms in crystals are usually tightly packed and the relevant temperatures are low compared with the melting temperature, the harmonic approximation to transition state theory can typically be used in studies of diffusion and reactions in crystals [59], and the problem reduces to determining the activation energy (Eq. [65]) and prefactor (Eq. [64]) of an Arrhenius-type expression (Eq. [63]). An excellent review of the techniques to determine these quantities can be found in [97], whose approach we are following in this section.

The most challenging part in this calculation is the search for the relevant saddle points. After a saddle point has been found, one can follow the gradient of the energy downhill, both forward and backward, and map out the Minimum Energy Path (MEP), thereby establishing what initial and final state the saddle point corresponds to. The identification of saddle points ends up being one of the most challenging tasks in theoretical studies of transitions in process modeling.

The minimum energy path is frequently used to define a "reaction coordinate" [100] for transitions. The minimum energy path may have one or more minima in between the end-points corresponding to stable intermediate configurations. The minimum energy path will then have two or more maxima, each one corresponding to a saddle point. Assuming a Boltzmann population is reached for the intermediate (meta)stable configurations, the overall rate is determined by the highest energy saddle point. It is, therefore, not sufficient to find a saddle point, but rather one needs to find the highest saddle point along the minimum energy path, in order to get an accurate estimate of the rate from harmonic transition state theory.

Many different methods have been presented for finding minimum energy paths and saddle points [101, 102]. Since a first order saddle point is a maximum in one direction and a minimum in all other directions, methods for finding saddle points invariably involve some kind of maximization of one degree of freedom and minimization in other degrees of freedom. The critical issue is to find a good and inexpensive estimate of which degree of freedom should be maximized. Below, we give an overview of several commonly used methods in studies of transitions in solids.

### 2.3.1. The Drag Method

The simplest and perhaps the most intuitive method of all is what we will refer to as the drag method. It actually has many names because it keeps being reinvented. One degree of freedom, the drag coordinate, is chosen and is held fixed while all other $(D-1)$ degrees of freedom are relaxed; that is, the energy of the system is minimized in a $(D-1)$ dimensional hyperplane. In small, stepwise increments, the drag coordinate is increased and the system is dragged from reactants to products. The maximum energy obtained on the way is taken to be the saddle point energy. Sometimes, a guess for a good reaction coordinate is used as the choice for the drag coordinate. This could be, for example, the distance between two atoms that initially form a bond which ends up being broken. In the absence of such an

intuitive choice, the drag coordinate can be simply chosen to be the straight line interpolation between the initial and final state. This is a less biased way and all coordinates of the system then contribute in principle to the drag coordinate.

This second approach is illustrated in Fig. 19. For its implementation, Henkelman et al. [97] have suggested to invert the force acting on the system along the drag coordinate and to use the velocity Verlet algorithm [103] with a projected velocity to simulate the dynamics of the system. The velocity projection is carried out at each time step and ensures that only the component of the velocity parallel to the force is included in the dynamics. When the force and projected velocity point in the opposite direction (indicating that the system has gone over the energy ridge), the velocity is zeroed. Such a projected velocity Verlet algorithm is also an efficient and simple minimization algorithm for many of the methods discussed here.

The problem with the drag method is that both the intuitive, assumed reaction coordinate and the unbiased straight line interpolation can turn out to be bad reaction coordinates. They may be effective in distinguishing between reactants and products, but a reaction coordinate must do more than that. A good reaction coordinate should give the direction of the unstable normal mode at the saddle point. Only then does a minimization in all other degrees of freedom bring the system to the saddle point.

Figure 19 shows a simple case where the drag method fails. As the drag coordinate is incremented, starting from the initial state, $R$, the system climbs up close to the slowest
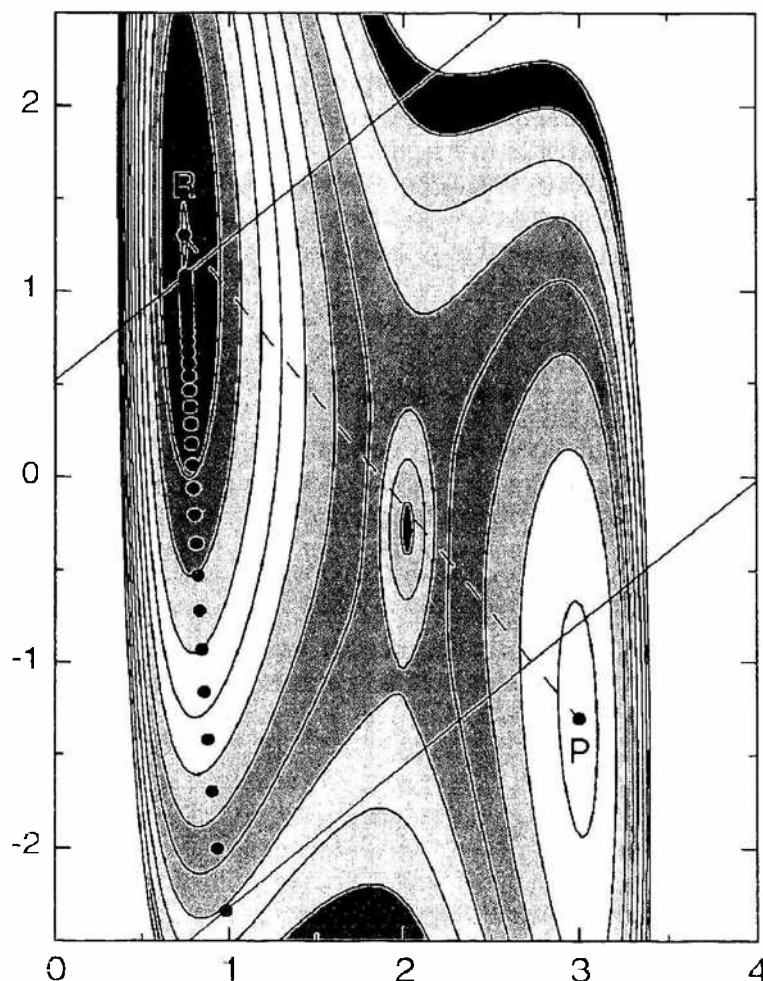


Figure 19. The drag method. A drag coordinate is defined by interpolating from $R$ to $P$ with a straight line (dashed line). Starting from $R$, the drag coordinate is increased stepwise and held fixed while relaxing all other degrees of freedom in the system. In a two-dimensional system, that means relaxation along a line perpendicular to the $P$-$R$ vector. The solid lines show the first and last relaxation line in the drag calculation. The final location of the system after relaxation is shown with filled circles. As the drag coordinate is increased, the system climbs up the potential surface close to the slowest ascent path, reaching a potential larger than the saddle point, and then, eventually, slipping over to the product well. In this simple test case, the drag method cannot locate the saddle point. Reprinted with permission from [97], G. Henkelman, Ph.D. thesis, University of Washington, 2001. © 2001.

ascent path. After climbing high above the saddle point energy, the energy contours eventually stop confining the system in this energy valley and the system abruptly snaps into an adjacent valley (the product valley in the case of Fig. 19). The system is never confined to the vicinity of the saddle point because the direction of the drag coordinate is at a large angle to the direction of the unstable normal mode at the saddle point. While there certainly are cases where the drag method works, there are also many examples where it does not work.

As an example, the originally suggested interstitial-assisted diffusion mechanism for boron in silicon had been identified by the drag method with an intuitively chosen reaction coordinate [104, 105]. This lead to the prediction of an initial kick-out event with subsequent long-range diffusion of the boron atom as an interstitial. Using the nudged-elastic band method (Section 2.3.2) instead, Windl et al. [106] could show that the long-range interstitial diffusion is highly improbable and that instead an immediate kick-in event should be the most probable event after the boron has been kicked out (see Section 3.2.1).

## 2.3.2. The Nudged Elastic Band Method

In the Nudged Elastic Band (NEB) method [102, 107, 108] a string of replicas (or "images") of the system are created and connected with springs in such a way as to form a discrete representation of a path from the reactant configuration, $R$, to the product configuration, $P$. Initially, the images may be generated along the straight line interpolation between $R$ and $P$. An optimization algorithm is then applied to relax the images down towards the minimum energy path. The NEB method is unique among the methods discussed here in the sense that it does not just give an estimate of the saddle point, but also a more global view of the energy landscape, for example, by showing whether more than one saddle point is found along the minimum energy path.

The string of images can be denoted by $[\mathbf{R}_0, \mathbf{R}_1, \mathbf{R}_2, \ldots, \mathbf{R}_N]$, where the endpoints are fixed and given by the initial and final states, $\mathbf{R}_0 = \mathbf{R}$ and $\mathbf{R}_N = \mathbf{P}$. Again, these are vectors of length $3n$ for a system with $n$ atoms, containing the three spatial coordinates of all atoms. The $(N - 1)$ intermediate images are adjusted by the optimization algorithm. The most straightforward approach would be to construct an object function

$$s(\mathbf{R}_1, \ldots, \mathbf{R}_N) = \sum_{i=1}^{N-1} E(\mathbf{R}_i) + \sum_{i=1}^{N} \frac{k}{2}(\mathbf{R}_i - \mathbf{R}_{i-1})^2 \tag{73}$$

and minimize with respect to the intermediate images, $\mathbf{R}_1, \ldots, \mathbf{R}_N$. This mimics an elastic band made up of $(N - 1)$ beads and $N$ springs with spring constant $k$. The band is strung between the two fixed endpoints. The problem with this formulation is that the elastic band tends to cut corners and gets pulled off the minimum energy path by the spring forces in regions where the minimum energy path is curved. Also, the images tend to slide down towards the endpoints, giving lowest resolution in the region of the saddle point, where it is most needed [102].

Both the corner-cutting and the sliding-down problems can be solved easily with a force projection. This is what is referred to as 'nudging'. The reason for corner-cutting is the component of the spring force perpendicular to the path, while the reason for the down-sliding is the parallel component of the true force coming from the interaction between atoms in the system. Given an estimate of the unit tangent to the path at each image (which will be discussed later), $\hat{\tau}_i$, the force on each image should only contain the parallel component of the spring force, and perpendicular component of the true force,

$$\mathbf{F}_i = -\nabla E(\mathbf{R}_i)|_{\perp} + \mathbf{F}_i^s \cdot \hat{\tau}_i \hat{\tau}_i \tag{74}$$

where $\nabla E(\mathbf{R}_i)$ is the gradient of the energy with respect to the atomic coordinates in the system at image $i$, and $\mathbf{F}_i^s$ is the spring force acting on image $i$. The perpendicular component of the gradient is obtained by subtracting out the parallel component

$$\nabla E(\mathbf{R}_i)|_{\perp} = \nabla E(\mathbf{R}_i) - \nabla E(\mathbf{R}_i) \cdot \hat{\tau}_i \hat{\tau}_i \tag{75}$$

In order to ensure equal spacing of the images (when the same spring constant, $k$, is used for all the springs), even in regions of high curvature where the angle between $\mathbf{R}_i - \mathbf{R}_{i-1}$ and

$R_{i+1} - R_i$ deviates significantly from 180°, the spring force should be evaluated as

$$F_i^* = k(|R_{i+1} - R_i| - |R_i - R_{i-1}|)\hat{\tau}_i \tag{76}$$

We now discuss the estimate of the tangent to the path. In the original formulation of the NEB method, the tangent at an image $i$ was estimated from the two adjacent images along the path, $R_{i+1}$ and $R_{i-1}$. The simplest estimate is to use the normalized line segment between the two,

$$\hat{\tau}_i = \frac{R_{i+1} - R_{i-1}}{|R_{i+1} - R_{i-1}|} \tag{77}$$

but a slightly better way is to bisect the two unit vectors,

$$\hat{\tau}_i = \frac{R_i - R_{i-1}}{|R_i - R_{i-1}|} + \frac{R_{i+1} - R_i}{|R_{i+1} - R_i|} \tag{78}$$

and then normalize $\hat{\tau} = \hat{\tau}/|\hat{\tau}|$. This latter way of defining the tangent ensures the images are equispaced even in regions of large curvature. A possible kinkiness in the path can be eliminated by defining the tangent of the path at an image $i$ by the vector between the image and the neighboring image with higher energy [109].

To start the NEB calculation, an initial guess is required. Usually a simple linear interpolation between the initial and final point is adequate. When multiple minimum energy paths are present, the optimization leads to convergence to the minimum energy path closest to the initial guess. In order to find the optimal minimum energy path in such a situation, some sampling of the various minimum energy paths needs to be carried out. for example a simulated annealing procedure, or an algorithm which drives the system from one minimum energy path to another, analogous to the search for a global minimum on a potential energy surface with many local minima [110].

In order to obtain an estimate of the saddle point and to sketch the minimum energy path, it is important to interpolate between the images of the converged elastic band. In addition to the energy of the images. the force along the band provides important information and should be incorporated into the interpolation. By including the force, the presence of intermediate local minima can often be extracted from bands with as few as three images. The interpolation can be done with a cubic polynomial fit to each segment $[R_i, R_{i+1}]$ in which the four parameters of the cubic function can be chosen to enforce continuity in energy and force at both ends [88].

The NEB method has been applied successfully to a wide range of atomic-scale process simulation problems, usually implemented in a plane-wave density functional theory code such as VASP [111], for example bulk diffusion of impurity atoms in silicon such as boron [106] (see Section 3.2.1), nitrogen [112] (see Section 3.1.1), carbon [113], or phosphorous [114]. and diffusion processes at and near semiconductor surfaces [115].

### 2.3.3. Henkelman Dimer Method

When the final state of a transition is not known, the search for the saddle point is more challenging. A climb up from the initial state to the saddle point is more difficult than might at first appear. It is not sufficient to just follow the direction of slowest ascent. The two-dimensional test problem illustrated in Fig. 19 is an example of that. However, adding information from second derivatives (also known as force constants. which determine the vibrational normal modes of the system) can significantly help to guide the climb toward saddle points. These so-called mode-following methods have become widely used in studies of small molecules and clusters. Their disadvantage is that they require the second derivatives of the energy with respect to all the atomic coordinates, i.e., the full Hessian matrix, and that the matrix needs to be diagonalized to find the normal modes, an operation that scales as $N^3$. The evaluation of second derivatives is often very costly, for example in plane wave based density-functional theory calculations. Even when simple empirical potentials are used, a practical limit can be reached for as little as a couple of hundred atoms [116].

The recently introduced "Dimer Method" by Henkelman and Jónsson [117] has the essential qualities of the mode following methods, but only requires first derivatives of the energy and no diagonalization. It can therefore be applied to plane wave DFT calculations and to large systems with several hundred atoms, as shown in Ref. [118]. The method involves two replicas of the system, a 'dimer', as illustrated in Fig. 20. The dimer is used to transform the force in such a way that optimization leads to convergence to a saddle point rather than a minimum. The dimer consists of two images (replicas) separated from their common midpoint $\mathbf{R}$, which denotes initially a (local) energy minimum of the system, by a distance $\Delta R$. The vector $\hat{\mathbf{N}}$ which defines the dimer orientation is a unit vector pointing from one image at $\mathbf{R}_2$ to the other image at $\mathbf{R}_1$. When a transition state search is launched from an initial configuration, with no prior knowledge of what $\hat{\mathbf{N}}$ might be, a random unit vector is assigned to $\hat{\mathbf{N}}$ and the corresponding dimer images are formed,

$$\mathbf{R}_1 = \mathbf{R} + \Delta R \hat{\mathbf{N}} \quad \text{and} \quad \mathbf{R}_2 = \mathbf{R} - \Delta R \hat{\mathbf{N}} \qquad (79)$$

Initially, and whenever the dimer is moved to a new location, the forces acting on the dimer and the energy of the dimer are evaluated. These quantities are calculated from the energy and the force ($E_1$, $\mathbf{F}_1$, $E_2$, and $\mathbf{F}_2$) acting on the two images. The energy of the dimer $E = E_1 + E_2$ is the sum of the energy of the images. The energy and the force acting on the midpoint of the dimer are labeled as $E_0$ and $\mathbf{F}_R$ and are calculated by interpolating between the images. The force $\mathbf{F}_R$ is simply the average force $(\mathbf{F}_1 + \mathbf{F}_2)/2$. The energy of the midpoint is estimated by using both the force and the energy of the two images. $E_0$ can be related to the forces and energies of the two replicas using the finite difference formula for the curvature $C$ of the potential along the dimer,

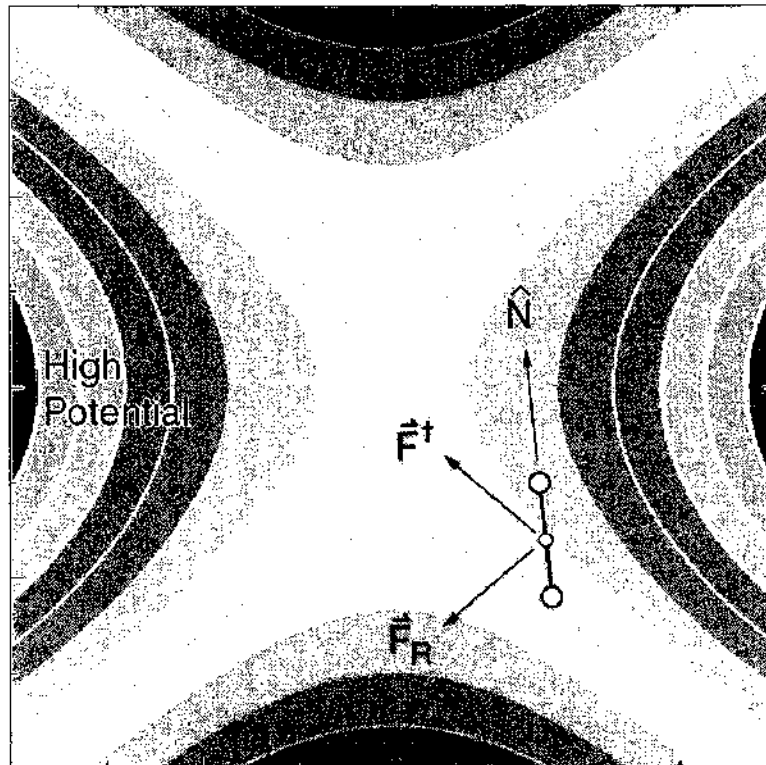$$C = \frac{(\mathbf{F}_2 - \mathbf{F}_1) \cdot \hat{\mathbf{N}}}{2\Delta R} \qquad (80)$$



Figure 20. The calculation of the effective force in the Dimer method. A pair of images, spaced apart by a small distance, on the order of 0.1 Å, is rotated to minimize the energy. This gives the direction of the lowest frequency normal mode. The component of the force in the direction of the dimer is then inverted and the minimization of this effective force leads to convergence to a saddle point. No reference is made to the final state. Reprinted with permission from [88], G. Henkelman and H. Jónsson, J. Chem. Phys. 115, 9657 (2001). © 2001, American Institute of Physics.

and the Taylor expansion of $E$.

$$E = 2E_0 + C(\Delta R)^2 \tag{81}$$

All the properties of the dimer are derived from the forces and energy of the two images. There is no need to evaluate energy and force at the midpoint between the two images. This is important for minimizing the total number of force evaluations required to find saddle points.

There are two parts to each dimer move. The first part is dimer rotation to minimize $E$. Since for that $E_0$ and $\Delta R$ are constant, Eq. (81) shows that the dimer energy, $E$, is minimal along the minimum of the curvature $C$. If the dimer is free to rotate, the forces acting on the two images will pull the dimer to the lowest curvature mode. This is done by defining a rotational force which is the difference in the force on the two images. Minimizing the energy of the dimer with respect to this rotational force aligns the dimer with the lowest curvature mode (this feature was used by Voter in his construction of bias potentials in hyperdynamics [119]). A modified Newton's method can be used to make this rotation efficient [88]. An
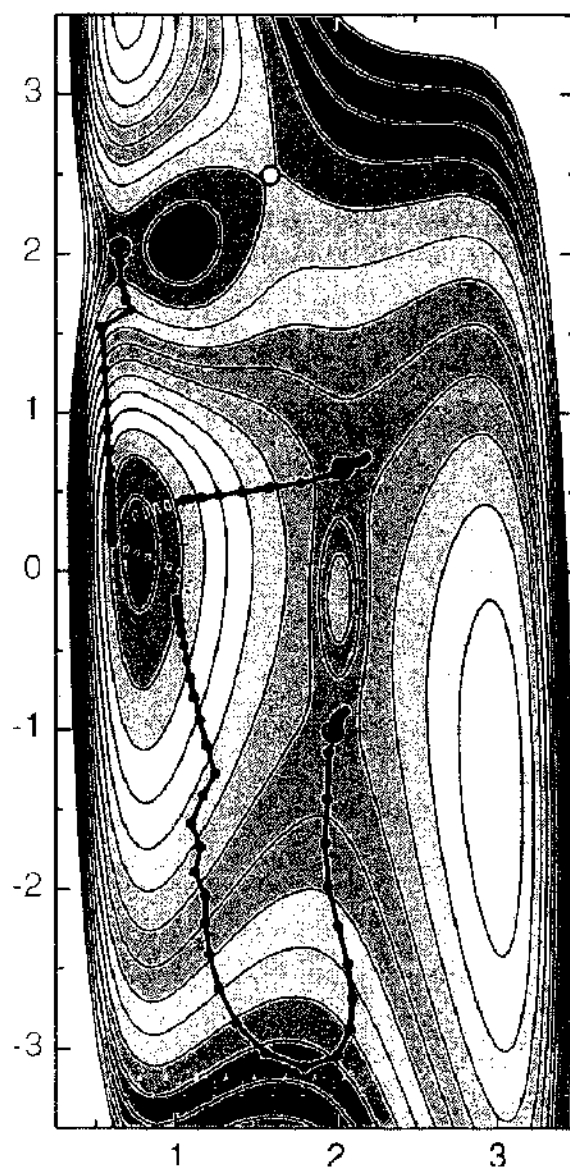


Figure 21. Application of the dimer method to a two-dimensional test problem. Three different starting points are generated in the reactant region by taking extrema along a high temperature dynamical trajectory. From each one of these, the dimer is first translated only in the direction of the lowest mode, but once the dimer is out of the convex region a full optimization of the effective force is carried out at each step (thus the kink in two of the paths). Each one of the three starting points leads to a different saddle point in this case. Reprinted with permission from [88], G. Henkelman and H. Jonsson, J. Chem. Phys. 115, 9657 (2001). © 2001. American Institute of Physics.

important aspect of the dimer method is that it only requires the first derivative of the energy, not the second derivatives.

The second part of the algorithm is translation of the dimer. A first order saddle point on a potential surface is at a maximum along the lowest curvature direction and a minimum in all other directions. In order to converge to a saddle point, the dimer is moved up the potential along the lowest curvature mode, and down the potential in all other directions. This is done by defining an effective force on the dimer, in which the true force due to the potential acting at the center of the dimer has the component along the dimer inverted. Minimizing with respect to this effective force moves the dimer to a saddle point. Usually, minimization using the conjugate gradient method works well.

Figure 21 shows a dimer calculation for the two-dimensional test problem. The initial configurations for the dimer searches were taken from the extrema of a short high temperature molecular dynamics trajectory (shown as a dashed line). The three initial points are different enough that the dimer searches converge to separate saddle points. In general the strategy for the dimer method is to try many different initial configurations around a minimum, in order to find the saddle points that lead out of that minimum basin.

## 3. APPLICATION EXAMPLES

In the following, we want to illustrate the different methods for nanoscale process modeling that were discussed in Section 2 with applications to front-end processing problems within the field of silicon nanoelectronic devices. The discussed examples include

- the calculation of kinetic parameters from atomistic simulations (Section 3.1), especially the
  - diffusion of nitrogen (Section 3.1.1), which illustrates the use of the nudged-elastic band method (Section 2.3.2), the linking of atomistic hopping to macroscopic diffusion constants (Section 2.1.3) and especially the calculation of diffusion prefactors (Eq. [64]);
  - calculation of capture radii (Section 3.1.2), where we describe in detail the basics of a kinetic Monte Carlo program, which is used (Section 2.2.2) to extract the capture radii needed for reaction rate constants (Section 2.1.2) from *ab initio* energetics and hopping parameters determined using the nudged-elastic band method (Section 2.3.2);
- the multiscale modeling of diffusion and deactivation of boron in silicon (Section 3.2), linking atomistic *ab initio* calculations for diffusion and reaction coefficients to continuum modeling; in detail,
  - in Section 3.2.1, the nudged-elastic band constant (Section 2.3.2) is used to determine the diffusivity of boron;
  - in Section 3.2.2, the binding energies for the reaction rate constants (Section 2.1.2) are calculated from first principles;
  - in Section 3.2.3, the resulting reaction kinetics are examined with the dimer method (Section 2.3.3), and
  - in Section 3.2.4, all calculated parameters are implemented into a diffusion reaction equation system (Section 2.1.1).
- the multiscale modeling of stress-mediated diffusion (Section 3.3), especially important for strained-channel devices;
- atomistic modeling of the formation of extended defects in ion-implanted silicon (Section 3.4); specifically, a
  - kinetic Monte Carlo study of the formation of vacancy clusters (or voids) (Section 3.4.1) and
  - accelerated dynamics simulations of interstitial-cluster growth (Section 3.4.2).

### 3.1. Kinetic Parameters from Atomistic Simulations

The advent of efficient and reliable *ab initio* codes in conjunction with the transition-theory based methods described in Section 2.3 have made it possible to include the nanoscale physics into the previously existing continuum and kinetic-Monte Carlo methods and to

produce predictive, physical nanoscale simulation tools. In the following, we want to demonstrate this with the help of a few examples from recent work.

### 3.1.1. Ab Initio Identification of the Nitrogen Diffusion Mechanism in Silicon

Nitrogen doping of silicon has become a process of increasing importance because of its various effects on the formation of extended defects in silicon like the complete suppression of void formation in float-zone processed crystals [120]. Additionally, nitrogen is known to increase the mechanical strength of silicon by locking dislocations [121] and, when implanted with a sufficient dose, to reduce the oxidation rate [122]. In partial explanation of the above effects, it has been shown that the $N_2$ pair readily forms complexes with vacancies, both in a metastable and immobile $N_2V$ configuration, and in the very stable $N_2V_2$ configuration which can either form from the reaction $N_2V + V \rightarrow N_2V_2$ or $N_2 + V_2 \rightarrow N_2V_2$ [123, 124]. The $N_2V_2$ complex has further been shown to attract oxygen into a stable configuration, indicating the possibility of further oxygen precipitate nucleation based on the $N_2$ pair.

As early as the investigations of Stein [125] in 1985 it was concluded from isotope shifts that the nitrogen dimer configuration is prevalent at room temperature. Jones et al. [126] confirmed this conclusion using a combination of spectroscopic and ab initio investigations. Except for the work of Gali et al. [127], who calculated a binding energy of 1.73 eV, all theoretical investigations agree on a rather high binding energy between 3.67 and 4.3 eV [124, 128, 129]. In contrast, the primary mechanism of diffusion of nitrogen in general and of the nitrogen pair in particular, have been controversially discussed in the literature.

In a variety of experimental investigations based on the out-diffusion of nitrogen from doped substrates [130] or on the in-diffusion of nitrogen from the ambient [131-134] the profiles obtained were interpreted in terms of diffusion of $N_2$ as an entity. The diffusion data obtained are shown in Fig. 22 and can be described best by a diffusion constant of

$$D_{N_2} = 35 \exp\left(-\frac{2.34 \text{ eV}}{k_B T}\right) \frac{\text{cm}^2}{\text{s}} \tag{82}$$

with the 90% confidence interval for the activation energy ranging from 2.01 to 2.77 eV. Profiles after ion implantation [135, 136] are considerably more complex and the possibility of a catalytic effect of oxygen on nitrogen diffusion has been reported recently [137]. In the analysis of Adam et al. [138], nitrogen dimers were not taken into consideration nor apparently needed to obtain an excellent description of the experimental profiles. A later analysis of Voronkov and Falster [139] was again based on nitrogen dimers as the prevalent defect, but it was concluded that they would diffuse via dissociation and diffusion of the monomers rather than as an entity. Uncertainties remained about this mechanism because the binding energy of $N_2$ obtained from their analysis, estimated to be between 2.24 eV and 2.9 eV, is considerably smaller than the bulk of the estimates from theoretical work. Concerning nitrogen interstitial diffusion, Schultz and Nelson calculated the migration barrier for nitrogen interstitials to be 0.4 eV for a split-interstitial configuration, resulting in very fast diffusion [140]. Furthermore, a DFT calculation of the nitrogen pair diffusion using the nudged-elastic band method (Section 2.3.2) had resulted in a predicted barrier between 2.9 and 3.3 eV, outside of the confidence range given above [141].

A recent theoretical re-examination of the $N_2$ diffusion problem finally seems to be able to allow a reconciliation of the different experimental findings for $N_2$ diffusion. Using again DFT calculations and the nudged-elastic band method, Stoddard et al. [112] found a new diffusion path by examining more minimum energy paths, resulting in a barrier of $\sim 2.36$ eV, well within the defined confidence range (Fig. 22). The different stages along this new diffusion path are illustrated in Fig. 23. The two nitrogen atoms move disjointedly, with the upper atom moving through Fig. 23[a]–[c] before the lower atom follows (Fig. 23[d]–[f]), while the highest energy configurations are those where the nitrogen atoms are the farthest separated. Due to the use of the climbing-image NEB [142], which pushes the highest energy point to the saddle point, Fig. 23(c) should represent the saddle point of the migration event with an activation energy of 2.36 eV (see Fig. 24), in excellent agreement with the experimental fit.
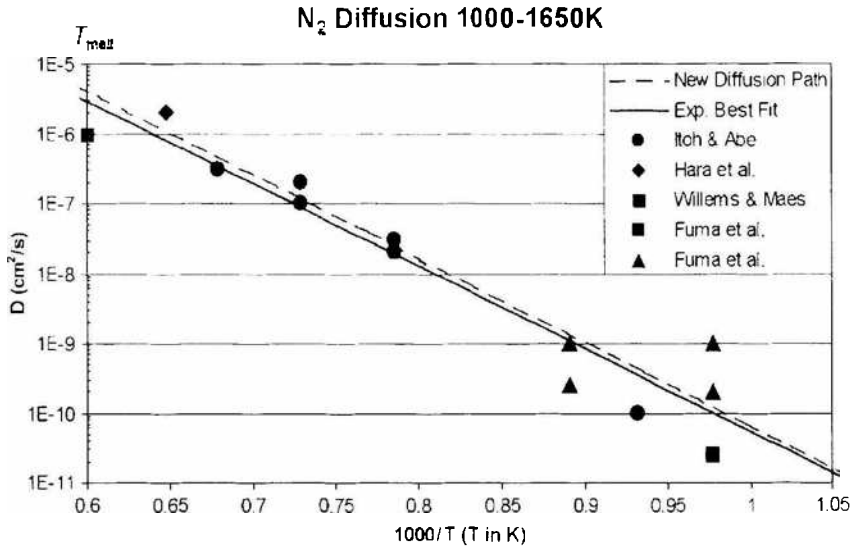
**Figure 22.** Experimental data for the $N_2$ diffusion constant in silicon with best-fit line and theoretical diffusivity. Reprinted with permission from [112], N. G. Stoddard et al., *Phys. Rev. Lett.* 95, 025901 (2005). © 2005, American Physical Society.

The diffusion coefficient prefactor, $D_0$, is determined by calculating vibrational frequencies for the saddle point and the ground state, as well as the entropies of formation and configuration,

$$D_0 = \frac{1}{12} \nu_0 d^2 p \, \exp\left(-\frac{S}{k_B}\right) \tag{83}$$

where $d$ is the jump distance, $p$ is the multiplicity of jump paths, $S$ is the combined entropy and $\nu_0$ is the jump frequency. Within harmonic transition-state theory, the jump frequency is calculated as the ratio of the product of the $\Gamma$-point frequencies of the ground state supercell over the product of all real $\Gamma$-point frequencies of the saddle point cell as defined in Eq. (64). The entropy of configuration is just the natural log of the number of possible defect configurations. Within the quasiharmonic approximation, the vibrational entropy $S$ of



**Figure 23.** Low activation barrier diffusion series for a nitrogen pair (in purple) moving through the silicon lattice (gray) in (100) projection. The first and last configurations are equivalent, and the corresponding energetics are provided in Fig. 24. The atoms each move, but make their jumps at different times. The highest energy configuration (c) is also the point of maximum separation for the nitrogen pair. Reprinted with permission from [112], N. G. Stoddard et al., *Phys. Rev. Lett.* 95, 025901 (2005). © 2005, American Physical Society.
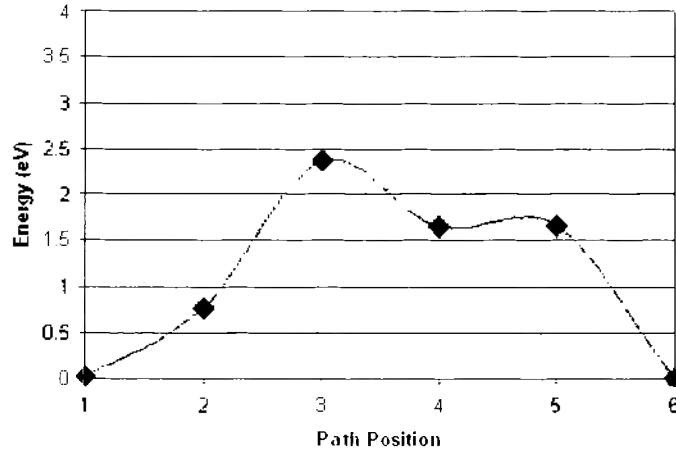
Figure 24. Energy profile for the $N_2$ diffusion path with a migration barrier of 2.36 eV. The points correspond to the atomic configurations shown in Fig. 23. The atoms each move, but make their jumps at different times. The highest-energy configuration (c) is also the point of maximum separation for the nitrogen pair. Reprinted with permission from [112]. N. G. Stoddard et al., *Phys. Rev. Lett.* 95, 025901 (2005). © 2005, American Physical Society.

an atomic configuration can be calculated by [143, 144]

$$S = k_B \sum_i \sum_f \left\{ \frac{\hbar\omega_{if}}{2k_BT} \coth\left( \frac{\hbar\omega_{if}}{2k_BT} \right) - \ln\left[ 2\sinh\left( \frac{\hbar\omega_{if}}{2k_BT} \right) \right] \right\}$$  (84)

where the summation goes over the degrees of freedom, $f$, for each atom, $i$, and $\omega_{if}$ is the characteristic frequency. The entropy of formation can be calculated from the entropy of the $N_2$ defect, N interstitial defect and perfect silicon by

$$S_f = 2S_N - S_{N_2} - S_{perfect}$$  (85)

The temperature dependence of the entropy (Eq. [84]) does not completely cancel out and has a weak effect on the effective migration barrier. For the diffusion path of Fig. 23, the full diffusion constant was calculated, including the entropy temperature dependence, in the temperature range of 800–1400°C. Using an Arrhenius fit where the prefactor is temperature independent, the data is best described in this temperature range by

$$D_{N_2} = 67 \exp\left( -\frac{2.38 \text{ eV}}{k_BT} \right) \frac{cm^2}{s}$$  (86)

see Fig. 22. For low temperatures (300°–700°C), values of $D_0 = 117 \text{ cm}^2/\text{s}$ and $E_a = 2.42$ eV better fit the temperature dependence of the theoretical diffusion constant. This Arrhenius fitting approach will be used throughout. Since the corrections from [145] that were applied to correct for the errors caused by the DFT band gap problem only amounted to ~0.02 eV difference in the migration barrier, the values in the calculated diffusion constant might be a very good estimate for the physical reality.

Stoddard et al. [112] also found a 0.44 eV migration barrier for single-interstitial diffusion of N in Si in good agreement with previous work [140], but a calculation of the temperature-dependent prefactor determined that the theoretical data is best fit by an Arrhenius expression of

$$D_N = 1.7 \exp\left( -\frac{0.56 \text{ eV}}{k_BT} \right) \frac{cm^2}{s}$$  (87)

in the temperature range from 800 to 1600 K. At low temperatures, the migration barrier is better fitted with a value of 0.50 eV. With these values, the interstitial nitrogen will diffuse at least five orders of magnitude faster than nitrogen pairs. Given the strong binding energy of nitrogen pairs and the fast diffusion of nitrogen interstitials, we expect that all of the nitrogen concentration will be paired or complexed within a very short time. Subsequent diffusion will be limited by nitrogen pair diffusion.

### 3.1.2. Combined Kinetic Monte Carlo/Ab Initio Determination of Reaction Capture Radii

We have discussed in Sections 2.1.2 and 2.2.2 that continuum models as well as Monte Carlo simulations of diffusion both employ a so-called capture radius, which in principle defines the distance where two species of atoms or reactant complexes start to feel their mutual presence and begin to react which each other. While used for a long time, this interaction or capture radius was not well characterized and usually was assumed to be on the order of the nearest-neighbor distance (for Si, 2.35 Å) or average spacing in the lattice $((5 \times 10^{22} \, \text{cm}^{-3})^{-1/3} = 2.71$ Å) [146]. However, the simulations of [147], which we want to summarize in the following, showed that defect interactions are more long ranged, extending up to at least eighth-nearest neighbor distance.

For a quantitative evaluation of the capture radii, Beardmore et al. calculated within density-functional theory the interaction potentials for several defect pairs and determined the most reasonable approximation for diffusion barriers between valleys of different energy, using the Vienna *Ab initio* Simulation Package (VASP) [111] and large simulation cells. Those results were used as input for KLMC simulations to determine the capture radii for the defect pairs.

Before the work of Beardmore et al. [147], all *ab-initio* work for vacancy-assisted diffusion in Si had been performed within 64-atom supercells, although those have been shown several years ago to be of insufficient size for defects involving vacancies [148]. Therefore, Beardmore et al. calculated the interaction potentials between defects by structural relaxation in a 216 atom ($3 \times 3 \times 3$ unit cells) volume. For each set of calculations, the two defects were placed in the cell, sampling all possible separation distances and orientations. Each initial configuration was relaxed and the energy and final configuration saved. For saddle point calculations, the nudged-elastic band method (Section 2.3.2) implemented in VASP was used.

The results of the *ab initio* calculations are plotted, relative to the energy at infinite separation, in Fig. 25 and show the existence of long range interactions between defects in silicon, which extend up to at least eighth-nearest neighbors. The interactions are dependent on separation distance, but also on the direction of displacement with respect to crystal orientation and the orientation of the split interstitial defects. To obtain a pair potential, the lowest energy at each distance was used. In the case of $I$-B and $I$-$V$ interactions, the tetrahedral instead of the $\langle 110 \rangle$ split interstitial gave the lowest energy at a given separation, which results from charge transfer to the B atom or $V$, respectively, who both prefer the negative over the neutral charge state (see, e.g., Ref. [106]).

It is interesting to note that for distances larger than $\sim$10 Å, the electronic minimization often, depending on the initial guess of the charge density, resulted in two different energy states for the same structure. For the lower one, charge depletion around the As atom and accumulation around the $V$ was found, which was not true for the higher-energy result. This result can be interpreted in the sense that for distances up to $\sim$10 Å, the excess electron from As is always within the electron capture radius of the $V$, whereas for larger distances,
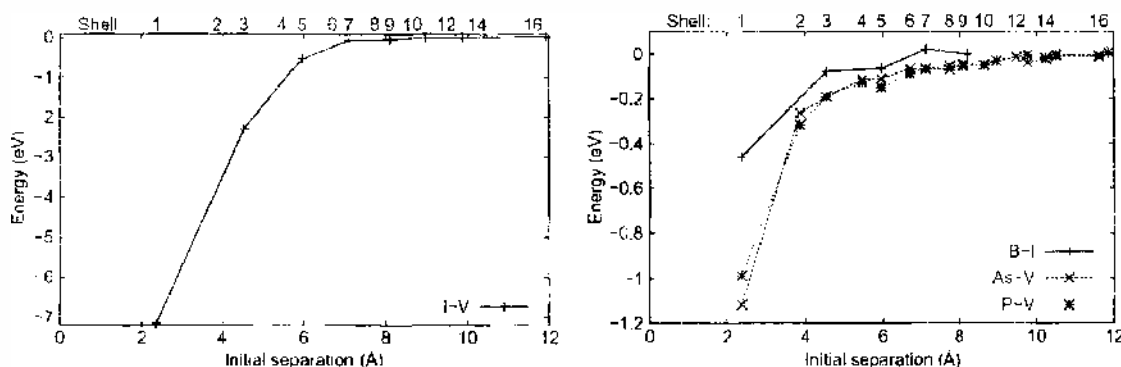


**Figure 25.** Calculated interaction energy as a function of separation for $I$-$I$, B-$I$, As-$I$, and P-$I$. The corresponding interaction shells are also shown. Reprinted with permission from [147], K. M. Beardmore et al., in "Computational Nanoscience and Nanotechnology." Applied Computational Research Society. Cambridge, 2002, p. 251. © 2002, Nanoscience and Technology Institute.

this is not necessarily true. Thus, these calculations give a minimum estimate of the electron capture radius for charge transfer from As to $V$ of $\sim 10$ Å. Knowledge of the electronic capture radius is necessary, for example, to estimate the characteristic time constant for charge transfer processes, see Eq. (29) in Section 2.1.2.

As described in Section 2.2.2, KMC parameters such as interaction energies are obtained from experiment and increasingly from *ab initio* calculations. Beardmore et al. [147] used KMC simulations with *ab initio* interaction and kinetic parameters for the determination of the capture radii [149]. The initial conditions for such capture radii calculations involve two defects, each either a vacancy, interstitial, or dopant randomly located on a 3D lattice in a given volume.

Depending on the interaction potential between a defect and another defect or dopant atom, the initial and final state of a mobile defect can have different energy levels, with no easy estimate for the saddle point energy in-between. Two approximations for such saddle points have been suggested in the past, which we will discuss for the case of an isolated vacancy in otherwise perfect silicon.

The first approximates the saddle-point energy by the sum of the free-vacancy migration barrier, $E_1^0$, and the energy of the higher-laying valley energy, $E_m = \max(E_i, E_f) + E_V^0$ [150]. The other suggests to add the free-vacancy migration energy to the average of the two valley energies, $E_m = (E_i + E_f)/2 + E_1^0$ [151]. While the second approximation might sound more reasonable at first glance, it has the inherent danger to end up with a negative migration barrier in cases where $E_1^0 < (E_i - E_f)/2$.

Beardmore et al. found with a nudged-elastic band ab-initio calculation for vacancy diffusion from the second- to the third-neighbor site that the saddle point is close to the energetically higher third neighbor site, and that its energy can be approximate well by the first approximation from above [150]. For vacancy diffusion from nearest to second-neighbor distance, however, they found a nearly vanishing diffusion barrier only slightly higher than the second-neighbor energy, which is considerably smaller than second-neighbor energy plus free-vacancy migration energy. Thus, the second approximation, which predicts no barrier in this case, seems to be the better approximation. Thus, neither approximation seems to be really universally valid, and a detailed understanding of the barriers seems to be necessary if one wants to use the most realistic kinetic parameters in a KMC simulation.

Since the diffusion step from nearest to second neighbor is the most crucial one for vacancy-assisted dopant diffusion, Beardmore et al. chose the second approximation in their work. Lumping the migration energy into the lattice migration rate $\nu_m$, the hopping frequency is given by $\nu = \nu_m \exp[(E_i - E_f)/(2k_B T)]$, where $E_i - E_f$ is the change in the system energy due to a transition.

All interactions used in determining system energies are assumed to be pairwise additive. The total binding energies are described by $E_X(i) = \sum_{Y,j} E_{XY}(j) \times N_Y(j)$, where $E_X(i)$ is the sum of all defect-defect binding energies for a point defect $X$ at site $i$ with other point defects $Y$ within the range of interaction. $N_Y(j)$ is the number of $Y$ type defects at $j$-th nearest neighbor distance from site $i$ and $E_{XY}(j)$ is the binding energy between $X$ and $Y$ type defects at separation $j$ from Fig. 25.

Within the von Smoluchowski approach outlined in Section 2.1.2 (Eq. [23]), the interaction between defects $A$ and $B$ is described by the reaction rate constant

$$k' = 4\pi a_c (D_A + D_B) \tag{88}$$

where $a_c$ is the capture radius. Since KLMC simulations explicitly include the silicon lattice structure, the capture rate $k'$ can be calculated directly, based on the interaction potential obtained from *ab initio* calculations. Using initial random defect displacements within a simulation box of volume $\Omega$, the capture rate is simply given by

$$k' = \frac{\Omega}{\langle \tau \rangle} \tag{89}$$

where $\langle \tau \rangle$ is the average capture time. Equation (88) is then inverted to give the capture radius.

Using the *ab initio* energies from Fig. 25, the capture radius calculated in Ref. [147] was approximately 7 Å for $I$-$V$ recombination when interactions up to sixth nearest neighbors (as had been done in previous work [149]) are taken into account. This was the same as previously calculated using interaction potentials obtained by empirical MD simulation [149], although the interaction potentials themselves differ significantly. Hence, it appears that the interaction range, not the shape of the potential is what determines the capture radius. However, the *ab initio* results indicate an attractive potential up to the 11th neighbor shell. Using the eleventh neighbor shell as a cutoff, Beardmore et al. found a capture radius of 8.3 Å for $I$-$V$ and 4.6 Å for B-$I$. All these simulations were performed at 900°C.

For both As-$V$ and P-$V$, Beardmore et al. examined the temperature dependence of capture radii up to 1700 K.

On the basis of the 18 shell range potentials, the capture radii are fit very well in the temperature range of 300 K $< T <$ 1700 K by $a_c = a_0 + a_1 \exp(-\sqrt{T/T_1})$, where $a_0 = 2.369$ Å is the silicon bond length. The fit resulted in $a_1 = 28.93$ Å and $T_1 = 339$ K for As-$V$ and $a_1 = 23.15$ Å and $T_1 = 499$ K for P-$V$ with 900°C values of 6.9 Å and 7.4 Å, respectively.

Thus, it was found that the capture radius increases slightly as the temperature is reduced, as might be expected since at low $T$ the capture probability increases for the weaker binding energies seen at larger distances. Overall, the important result of Ref. [147] was that the calculated capture radii seem far larger than those currently used in most continuum diffusion models due to the relatively long range of the interaction potential.

An interesting question in this context concerns the type of interaction between dopants and defects, such as between an As atom and a vacancy. Assuming ionization of As to As$^+$ while the vacancy accepts the extra electron, the interaction should be—apart from some close-range deviations—purely Coulombic. If this were the case, the effective value of the dielectric constant at the capture radius can be determined, which is important to determine reaction rate constants for two charged reactants as stated in Eq. (25). For this purpose, Windl [152] fitted the Madelung-corrected Coulomb energy to the interaction energies calculated in Ref. [147] within the local density approximation (LDA),

$$Z_M \cdot E_{Coul} = -\frac{e^2}{\varepsilon r} + \Delta E \tag{90}$$

where $Z_M$ denotes the Madelung constant for an As$^+$-$V$ pair at distance $r$ [153], $e$ is the electron charge, and $\varepsilon$ (dielectric constant) and $\Delta E$ (possible energy offset to determine the absolute value of the relative interaction energies) are fitting parameters. The Madelung constant corrects for an artifact of the supercell approximation, which determines in an electronic-structure calculation the interaction energy of a periodic array of As$^+$ atoms and $V$, not the interaction energy for an isolated pair [153].

The resulting fit is shown in Fig. 26. The black dots are the original calculation results, arbitrarily shifted to result in an 18th neighbor interaction energy of zero. The black dashed line is a straightforward fit with fixed $\Delta E = 0$. The red dashed line is the result of a fit with variable $\Delta E$, resulting in a downward shift of 0.09 eV and a dielectric constant of 12.1, which is in excellent agreement with experimental (12.1) and theoretical (12.8) values of bulk Si [154]. The fitted curve (red dashed line in Fig. 26) starts to match the calculated interaction energies (red dots in Fig. 26, shifted to their correct absolute values, and red solid line as a guideline to the dots) at ~5 Å, which is smaller than the calculated capture radius at typical annealing temperatures (for example, ~7 Å at 900°C). Thus, this calculations find that the appropriate dielectric constant for Eq. (25) is the bulk Si value of 12.1.

## 3.2. *Ab Initio* to Continuum Modeling of Diffusion and Deactivation of Boron

In this section, we want to demonstrate how a physical multiscale process model can be constructed for the evolution of nanoscale clusters of—in our example—boron in silicon under heat treatment. The modeling examines the atomic-scale detail of the process from first principles calculations and eventually predicts the evolution of macroscopic amounts of material using a reaction-diffusion type continuum model, thus covering the two different
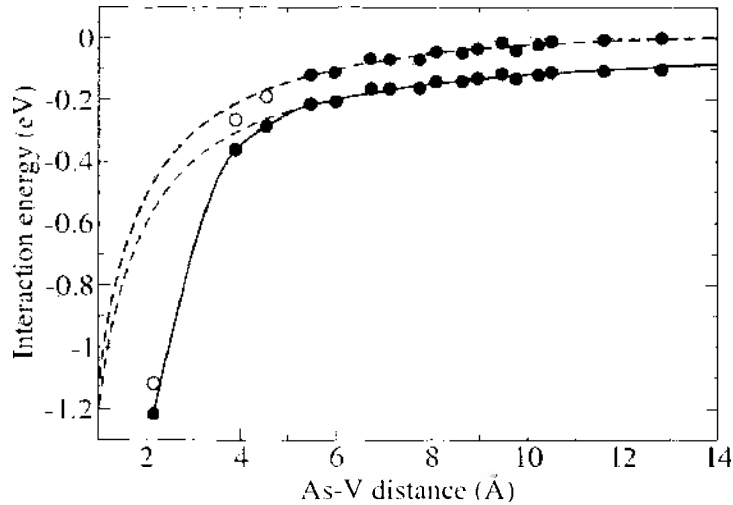
**Figure 26.** Calculated As⁻-I⁻ interaction energies from Ref. [147] on an arbitrary absolute scale (black dots) and shifted (red dots) to the correct absolute values for pure Coulomb interaction between As⁻ and V⁻ following a Coulomb fit. The black dashed line is a Coulomb fit of the black dots, the red dashed line one to the shifted red dots and thus the "true" Coulombic interaction. The red solid line is a spline fit to the red dots to determine where red dashed and red solid lines start to match. Reprinted with permission from [152]. W. Windl (unpublished).

length scales most important for nanoscale CMOS devices: the arrangement of the atoms, which controls the electrical properties of the device, and the influence of processing conditions on the final material. The following summarizes the results of [155].

We chose boron for this example since it is currently the most widely used acceptor dopant in silicon devices and displays very interesting nanoscale structure formation during fabrication. The energetic ions of the ion implantation process damage the host material and build a supersaturation of defects in Si, which impair the device performance. Post-implant annealing is used to heal the implant damage, while activating the dopant atoms electrically. The supersaturation of defects after the implant leads to excessive transient enhanced diffusion (TED) of the implanted B, assisted by mobile interstitials, during the annealing cycle. Excessive diffusion shifts the junction between n and p-type material to undesired locations. On the other hand, the interactions between mobile dopant-defect pairs, mobile defects, and dopants cause the formation of boron clusters, which are immobile, and deactivate the B atoms well below the solid solubility limit [156]. Here, deactivation means that boron does not contribute to the electrical conductivity because, for example, it forms a number of bonds commensurate with the number of its valence electrons.

## 3.2.1. Boron Diffusion in Silicon

Boron diffuses nearly exclusively with the help of Si self-interstitials [157]; that is, the mobile entity is thought to be a B atom paired with an I. As concerning the diffusion mechanism, ab initio modeling had suggested a few years ago that a kick-out mechanism with long-range low-barrier interstitial migration would be the dominant mechanism [105], in contrast to previous perception. In that work, diffusion saddle point configurations had been guessed or estimated by the drag method without systematic advancement of the atoms (see Section 2.3.1) [105]. As we have seen in Section 2.3.1, such methods are often not reliable, especially in cases where the diffusion involves the concerted motion of more than one atom [102]. Therefore, Windl et al. [106] re-examined the minimum-energy barrier diffusion path for I-assisted, charge-state dependent B diffusion using the nudged elastic band method (Section 2.3.2) implemented into VASP [111].

Following Windl et al. [106], the pair with the lowest formation energy in the neutral case, BT⁰ (Fig. 27[c]; T denotes a tetrahedral self-interstitial), has a formation energy of 2.8 (2.5) eV + E_f with respect to the lowest-energy B charge state. B⁻, with a binding energy of 0.9 (0.6) eV relative to the most stable dissociation products I⁺ and B⁻ (here and in the following, numbers without brackets denote generalized-gradient approximation
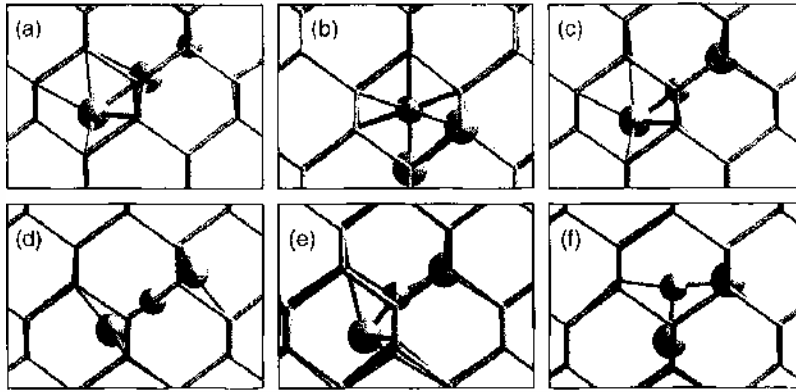
Figure 27. B (black ball)-I pairs in Si (gray balls and sticks). Reprinted with permission from [155]. W. Windl. IEICE Trans. Electron. E86C, 269 (2003). © 2003, Oxford University Press.

[GGA] results, whereas brackets denote LDA values). This binding energy is very similar to the sheer Coulomb attraction of a positive and a negative point charge (0.6 eV). For the +1 charged system, the same configuration, $BT^+$, has the lowest formation energy with a binding energy of 1.0 (0.8) eV with respect to the dissociation products $B^-$ and $T^{++}$. For the −1 charged system, two dumbbell-like interstitials with [110] and [100] orientation have the lowest formation energies, $BX^-$ (Fig. 27[e]) and $BS^-$ (Fig. 27[f]), respectively. $BX^-$ has the lowest total energy with a binding energy of 0.5 (0.3) eV with respect to B− and $X^0$.

In the neutral case, Windl et al. found the $BT^0$ pair to migrate via the $BS^0$ (Fig. 27[f]) to an $BH^0$ (Fig. 27[h]) interstitial by a buckling of the Si-B-$T$ triple dumbbell with a migration barrier of 0.2 (0.4) eV, which is a kick-out event. A cartoon of this diffusion mechanism is shown in Fig. 28[a]. The diffusion path between two neighboring $H$ sites also contains the $S$ interstitial, from where another $BT^0$ configuration can be accessed without barrier, or another $H$ site can be reached over a barrier of 0.1 (0.1) eV. This suggests an immediately following $BH^0 \rightarrow BT^0$ step to be the most probable event after the $BT^0 \rightarrow BH^0$ portion of the diffusion step, which predicts an immediate kick-in event without long-range interstitial diffusion.

For systems with positive charge, a one-step process $BT^+ \rightarrow BT^+$ with no intermediate metastable interstitial position was identified, a bond-centered interstitial $BB^+$ (Fig. 27[e]) as saddle point, and a migration barrier of 0.8 (1.2) eV (see Fig. 28[b]). However, there



(a)

(b)

(c)

Figure 28. Cartoon of diffusion paths for (a) neutral and positive, (b) alternative positive, and (c) negative charge states of BI pairs in (110) projection. The small balls are boron atoms, the larger balls and bond vertices are silicon atoms. Reprinted with permission from [155]. W. Windl. IEICE Trans. Electron. E86C, 269 (2003). © 2003, Oxford University Press.

is a second, competing process, especially for the LDA calculations, which is similar to the neutral diffusion mechanism (Fig. 28[a]), but has $BH^+$ as the saddle point, with a migration barrier of 1.0 (1.3) eV, slightly higher in energy, but with a larger average hopping length and more possible paths [106]. For negatively charged systems, there is a $BX^- \to BS^- \to BX^{--}$ path with an intermediate metastable $BS^-$ configuration and a migration barrier of 0.6 (0.5) eV (shown in Fig. 28[c]). Overall, diffusion and deactivation models derived from these calculations have given excellent results in comparison to experiment [12, 106].

## 3.2.2. Boron Deactivation in Silicon

The implant-anneal cycle can cause the formation of boron precipitates which immobilize and deactivate the boron atoms well below the solid solubility limit. From the observation of the trapping of interstitials by these precipitates, it has been concluded that they consist of boron-interstitial clusters (BICs). In the following, we discuss a systematic study of BIC energetics including the influence of charges and a careful structure minimization within the BIC phase space [12].

### 3.2.2.1. Substitutional Boron Clusters.   Substitutional boron clusters consist of B atoms substituting Si lattice atoms and has been studied in detail by Windl [155]. Although calculations for substitutional clusters might seem to be a straightforward task, care needs to be taken to recover the correct configuration. Traditionally, substitutional clusters were defined to consist of a nearest-neighbor assembly of B atoms, varying the possible configurations under this constraint. In the case of substitutional diatomic carbon clusters, however, the lowest-energy structure has not been found for a nearest-neighbor arrangement, but for the third-neighbor distance between the carbon atoms, where they occupy opposing corners of an hexagonal ring in the [110] plane, which allows for maximum stress relief [158]. Therefore, Windl generalized the concept of substitutional clusters in the sense to search for the minimum-energy configuration without constraining the distance between the B atoms in the Si cell in any way except for the finite size of the supercell (64 atoms).

In the case of the $B_2$ cluster, Fig. 29 shows relaxed energies of two substitutional B atoms at different distances. The previously exclusively studied nearest-neighbor $B_2$ cluster is clearly the energetically most unfavorable cluster, 0.8 eV higher than two B atoms far away from each other. A shallow minimum is found for the third-neighbor configuration, which seems to be analogous to the above described case of C in Si [158]. However, subsequent calculations
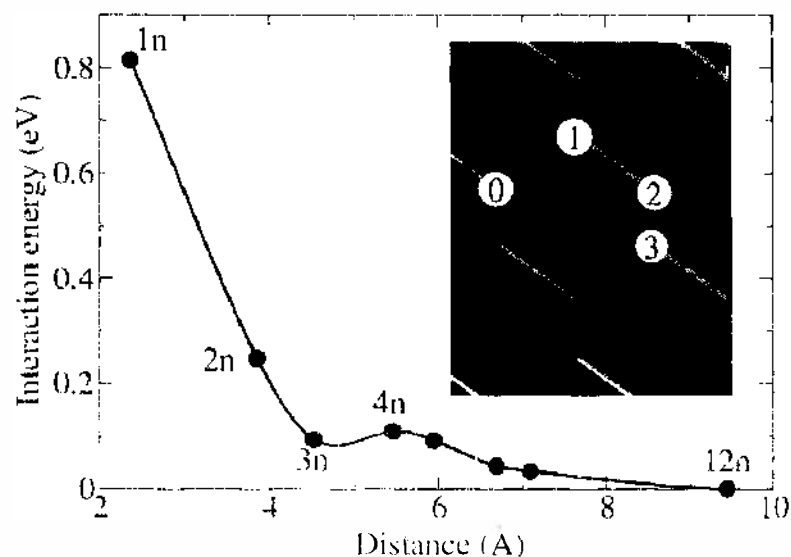


Figure 29. Energy of two substitutional B atoms in a 64-atom Si supercell for relaxed geometries as a function of the distance between the B atoms, relative to the lowest calculated energy at the 12th-neighbor distance. The line is a spline interpolation to guide the eye. The inset shows the position of first, second, and third neighbors relative to atom 0. Reprinted with permission from [155], W. Windl, *H. J. C. E. Trans. Electron.* E86C, 269 (2003). © 2003, Oxford University Press.

with larger cells (216 atoms) indicate that the third-neighbor minimum might be an artifact of the 64-atom supercell. Therefore, substitutional diboron clusters seem to be unstable with respect to the energy of dilute boron atoms.

In order to study if substitutional clusters with more than two boron atoms are possible, Windl [155] looked into the situation of four boron atoms in more detail. Since here, no single parameter such as the distance in the $B_2$ case can describe the arrangement easily, a first-principles Monte Carlo search was performed.

The initial configuration was a 64-atom $2 \times 2 \times 2$ supercell of the conventional unit cell which contains 60 Si atoms and 4 B atoms in random locations. After relaxing the structure and recording the energy, $E_1$, a Si and a B atom are randomly selected and swapped. If the new relaxed energy, $E_2$, is less than $E_1$, the move is accepted, otherwise, the Metropolis criterion is employed to accept or reject the move by comparing a random number between 0 and 1 with $\exp[(E_1 - E_2)/k_B T]$. Five hundred Monte Carlo steps have been run each for the neutral and double negative charge state, which had been identified previously to have the lowest energy in a nearest-neighbor configuration [12].

The resulting energy distribution, which is very similar for both charge states except for a constant energy shift, is shown in Fig. 30. All $B_4$ clusters with nearest-neighbor B arrangements have an energy of 1.3 eV or more higher than the lowest configuration. Again, they form the upper end of the energy distribution, as was the case for the nearest-neighbor $B_2$ pair. The configuration with maximum distance between all B atoms has an intermediate energy, 0.7 eV higher than the lowest-energy configurations. The latter consist of atomic arrangements where two pairs of B atoms sit in third-neighbor positions each. Thus, we conclude that in the case of extremely high B concentration without precipitation of new phases, substitutional nearest-neighbor B clusters would be highly unlikely, but an ordering-effect in third-neighbor positions, similar to the case of C in Si [158, 159], might be possible.

Overall, substitutional B clusters should not play a significant role in the deactivation of B, which suggests that the substitutional side of a reaction-diffusion model can be well covered by single substitutional B atoms.

**3.2.2.2. Boron-Interstitial Clusters.** Liu et al. [12] examined a number of BICs $B_n I_m$ with $n, m < 4$, as well as $B_{12} I_7$ (which had been studied theoretically before without presenting formation energy values [160]) and single B atoms in {311} defects. For their calculations, they used the DFT code VASP [111] within both LDA and GGA and 64-atom supercells.
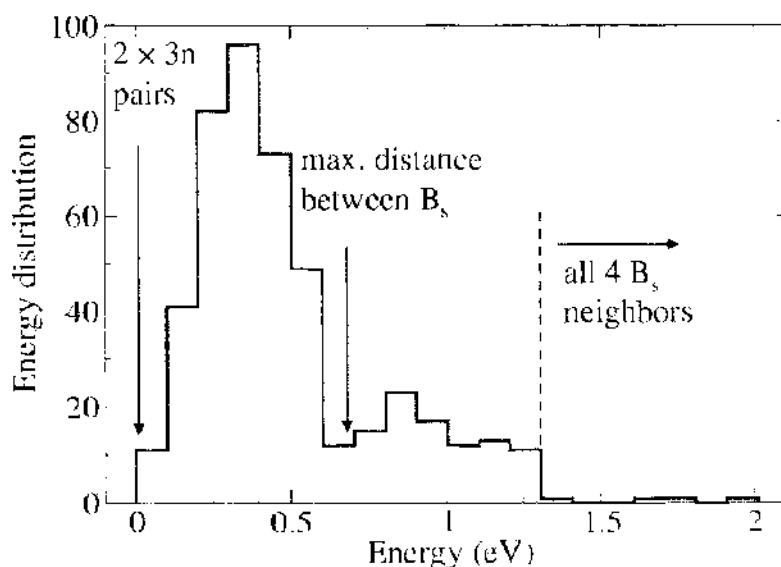


Figure 30. Energy distribution of all 500 steps during Monte Carlo simulation at 1200 K for a $Si_{60}B_4$ cell, where Si and B positions are swapped. The energies of the different nearest-neighbor $B_4$ clusters have been added by hand, since they did not appear during the simulation and their nucleation due to the strong repulsion between two B atoms is highly unlikely. Reprinted with permission from [155], W. Windl, IEICE Trans. Electron. E86C, 269 (2003). © 2003, Oxford University Press.

To decrease the threat of finding a high-energy local instead of the global minimum, they started for each cluster from many different initial configurations that were structurally relaxed.

Figure 31 summarizes the results from Liu et al. [12]. The figure shows the structure, LDA and GGA energies as well as fitted energies from Ref. [49]. The reference states for the cluster energy $E_{cl}$ are $B^-$ and $I^0$, that is, for $B_nI_m$, $E_{cl} = E_{B_nI_m} - n(E_B - E_{Si_{bulk}}) - m(E_{I^0} - E_{Si_{bulk}}(N+1)/N)$, where $N+m$ is the total number of atoms in the supercell. Displayed are the most stable structures and energies found for the lowest-energy charge states at midgap. Except for the substitutional clusters, all clusters have negative formation energy.

### 3.2.3. Dimer Study of Boron Clustering in Silicon

The vastly used von Smoluchowski model for reaction constants (Section 2.1.2) assumes that reactions are diffusion limited. For the case of cluster formation, this means that the barrier for formation of a cluster from smaller components is just the migration barrier of the mobile species. Equivalently, this assumption says that the barrier for cluster dissociation is the sum of the cluster binding energy with respect to the smaller components and the migration barrier of the mobile fragment that diffuses away. Uberuaga et al. [161] have tested this assumption by looking at the atomic mechanisms responsible for cluster formation using the dimer method (described in Section 2.3.3). Their work focused on two clusters, $B_2I$ and $B_3I_2$.

On the basis of the results of continuum modeling discussed in Section 3.2.4, one of the most dominant B clusters seems to be $B_3I$. Therefore, Uberuaga et al. ran 10 dimers from this cluster to study its formation and decay mechanism. Of the interesting results, two led to simple exchange processes of one B atom with another. One of the dimer runs led to the beginning of the dissociation of the cluster. When a second dimer is started from the basin found from the first run, a full dissociation of $B_3I^-$ to $B_2^-+ BI^+$ is found (the charge states are assumed due to the results of Section 3.2.2).

Since formation of $B_2^-$ is unlikely due to the strong repulsive interaction between neighboring boron atoms (Fig. 29), this is an improbable route for the formation of $B_3I^-$ clusters. Therefore, Uberuaga et al. looked at a larger cluster, $B_3I_2$, as an intermediate structure that potentially could decay into $B_3I$.

Of the dimers run from $B_3I_2$, four locked on to zero-curvature modes, two were reorientations of the original cluster (identical in structure, but rotated with respect to the original

geometry), one resulted in a distorted structure, and one lead to the beginning of the dissociation of the cluster.

The original cluster, which involves three B atoms in a trigonal structure bonded to a Si atom, can rotate so that the boron atoms are bonded to a different Si atom in the same tetrahedron. This does not lead to net diffusion of the cluster, but rather just a reorientation. This process has a barrier of 0.8 eV.

The other interesting process found led to the dissociation of the cluster. This reaction involved an intermediate state in which the three B atoms formed a linear chain in the Si crystal which is degenerate in energy with the original cluster. The minimum energy path for this reaction is shown in Fig. 32. The barrier to go from one structure to the other is 1.48 eV. Overcoming a second barrier of 1.58 eV leads to the products $B_2I$ and $BI$ with an energy of 1.42 eV above the original cluster. This already agrees well with the infinite separation energy of these individual products of 1.5 eV (Section 3.2.2). The reverse barrier, of joining $B_2I$ and $BI$ to form $B_3I_2$, is only 0.16 eV, which is smaller than the diffusion barrier of $BI$ in bulk Si [106]. Thus, the formation of the $B_3I_2$ cluster from $B_2I$ and $BI$ is shown to be diffusion limited, in excellent agreement with the postulate of the von Smoluchowski approximation for reaction rate constants (Section 2.1.2).

## 3.2.4. Link to Continuum Modeling

As shown in Section 2.1.2, within the von Smoluchowski approximation, reaction barriers can be calculated from the difference of the total formation energy of the reactants and the total formation energy of the products, which is the binding energy, plus the migration energy of the mobile reactant. In the previous section, we have shown that this model seems to be an excellent approximation at least for the case of boron-interstitial clusters (see Section 3.2.3). In the following, we will describe the formulation of a diffusion-reaction model for the case of BIC formation and decay.

As an example, let us consider the reaction

$$B_2I^0 + I^0 \rightarrow B_2I_2^0 \tag{91}$$

Following Section 2.1.1, the two partial differential equations (PDEs) describing this reaction can be formulated as

$$\frac{dC_{B_2I}}{dt} = -K^1_{B_2I_2} C_{B_2I} C_I + K^r_{B_2I_2} C_{B_2I_2} \tag{92}$$



Figure 32. Minimum energy path for the breakup of $B_3I_2$ into $B_2I$ and $BI$. The reverse barrier for the formation of $B_3I_2$ is only 0.2 eV. Reprinted with permission from [161]. B. P. Uberuaga et al., *Phys. Status Solidi* 233, 244 (2002). © 2002, Wiley-VCH.

and

$$\frac{dC_{\mathrm{B},I_2}}{dt} = K^{\mathrm{f}}_{\mathrm{B}_2 I_2} C_{\mathrm{B}_2 I} C_I - K^{\mathrm{r}}_{\mathrm{B}_2 I_2} C_{\mathrm{B}_2 I_2} \tag{93}$$

Interaction terms for reactions with other clusters can be added to these equations in an analogous way, and the whole PDE system can finally be solved with a PDE sover like ALAMODE [162]. For the equations of the mobile species, a Fick term is added to describe their diffusion to end up with a continuity equation as shown in Eq. (6). The forward and reverse reaction coefficients are described in the standard way (Section 2.1.2) by

$$K^{\mathrm{f}}_{\mathrm{B}_2 I_2} = 4\pi a_c D_{I^0} \tag{94}$$

$$K^{\mathrm{r}}_{\mathrm{B}_2 I_2} = a_c / C_{\mathrm{Si}} \; D_{I^0} \exp\!\left(-\frac{E_b}{k_B T}\right) \tag{95}$$

where $a_c$ is the capture radius of the reaction which can be calculated as described in Section 3.1.2, $D_{I^0}$ is the diffusivity of the mobile species $I^0$ in the reaction (calculation of diffusivities has been demonstrated for the case of nitrogen in Section 3.1.1 and for boron in Section 3.2.1), $C_{\mathrm{Si}}$ is the atom density in Si, and $E_b$ is the binding energy of the reactants in the cluster from Section 3.2.2.

In the current formulation, the reaction barrier is split into the binding energy part and the migration energy of the mobile species, which is hidden in the Arrhenius expression of the interstitial diffusivity. Despite of many theoretical investigations, there is considerable uncertainty in the migration barrier of the Si self-interstitial. The suggested values vary between 0.1 eV [163] and 1.4 eV [164] or even bigger and corresponding prefactors vary over many orders of magnitude.

To quantify the reaction coefficients, we start from the energy states of the species involved in the reaction of Eq. (91), which are according to Tables 1 and 2 [165] within GGA $E(\mathrm{B}_2 I^0) = -2.02$, $E(I^0) = 0$, and $E(\mathrm{B}_2 I^0_2) = -3.30$ eV, respectively. Considering the forward reaction in Eq. (91), we see that the sum of the formation energies on the l.h.s. is considerably higher than that on the r.h.s. and thus would take place whenever a self-interstitial becomes available. This explains *a posteriori* why the forward reaction (Eq. [95]) has no reaction barrier beyond the migration energy of the self-interstitial. The reverse reaction, in contrast, has to overcome a higher barrier with a binding energy of $E_b = -2.02$ eV $+3.30$ eV $= 1.28$ eV and a reaction barrier of $E_r = E_b + E_m(I^0)$. In an analogous way, PDEs and coefficients can be calculated for the whole system.

Liu et al. formulated in this way a continuum model system from the results in Fig 31 and combined it with a well-tested four-stream model for intrinsic B diffusion. For $I$ clustering, they mainly used the model from Ref. [166].

Table 1. Energy state values (eV) of boron clusters with interstitials. The midgap values are taken.

| Lowest energy state | GGA | LDA | The fitted values from Ref. [49] |
|---|---|---|---|
| $\mathrm{B}_2 I^0$ | −2.02 | −1.25 | −0.6 |
| $\mathrm{B}_2 I$ | −3.00 | −2.27 | −3.6 |
| $\mathrm{B}_4 I^-$ | −2.10 | −1.33 | −4.6 |
| $\mathrm{B} I^0_2$ | −2.49 | −2.28 | −2.0 |
| $\mathrm{B}_2 I^0_2$ | −3.30 | −2.56 | −3.0 |
| $\mathrm{B}_3 I^0_2$ | −4.18 | −3.21 | −5.8 |
| $\mathrm{B}_4 I^-_2$ | −5.52 | −4.31 | −6.8 |
| $\mathrm{B} I_3$ | −4.83 | −4.10 | |
| $\mathrm{B}_2 I^0_3$ | −6.05 | −5.20 | −3.0 |
| $\mathrm{B}_3 I_3$ | 6.60 | 5.67 | −7.8 |
| $\mathrm{B}_4 I_3$ | −7.06 | −5.90 | −9.0 |
| $\mathrm{B}_2 I_4$ | 9.25 | 7.89 | −10.9 |
| $\mathrm{B}_3 I_4$ | −25.08 | | |

*Source*: Reprinted with permission from [165]. X. Y. Liu and W. Windl, J. Comp. Elec. (in press), 2005, Springer.

Table 2. Energy state values (eV) of Si interstitials and boron-interstitial pairs. The midgap values are taken.

| Energy state | GGA | LDA | The fitted values from Ref. [49] |
|---|---|---|---|
| $X^0(I)$ | 0 | 0 | 0 |
| $T^{-}(I)$ | 0.25 | 0.38 | |
| $BI^{-}$ | −0.57 | −0.38 | −0.6 |
| $BI$ | −0.44 | −0.30 | |

*Source:* Reprinted with permission from [165]. X. Y. Liu and W. Windl. *J. Comp. Elec.* (in print).

The results for B activation are shown in Fig. 33 [167]. Both LDA and GGA models correctly predict inverse annealing at low temperatures due to the beginning formation of $B_3I_2$ and $B_4I_2$. Because of the low energy for $B_2I_3$, however, the model predicts in contrast to previous work that also the decay of this high-$I$ content cluster contributes significantly to the activation process at higher temperatures. Thus, it might be desirable to include clusters with higher $I$ content into future work.

LDA predicts too much activation too soon, whereas GGA results in too strong clustering as compared to experiment. In order to improve agreement with experiment, Liu et al. used a genetic algorithm to refit the parameters to a large number of SIMS data for different annealing conditions. The fit result was mostly independent of the starting parameters and resulted in parameter values which lay with a few exceptions between the LDA and GGA values, which therefore might be considered as upper and lower boundaries for the clustering energies. No activation data were used for the fit. Nevertheless, the experimental data in Fig. 33 are well predicted, suggesting that the key physics might be described reasonably well within this model.

## 3.3. Multiscale Modeling of Stress-Mediated Diffusion

Since further miniaturization of current MOSFET nanoscale devices becomes increasingly difficult due to the physical limitations and problems outlined in Section 1.1, other solutions to enhance the switching speed have been developed. One very successful approach in this area is the use of a strained channel material. The strain is usually created by growing a thin layer of Si on a relaxed SiGe film (Fig. 34[a]) [168]. Ge has a lattice constant that



Figure 33. Simulated and experimental [167] B activation after a 40 keV, $2 \times 10^{14}$ cm$^{-2}$ B implant for 30 min anneals at varying temperatures. Dashed line: LDA clustering energies; dot-dashed line: GGA energies; solid line: re-fitted to SIMS. Reprinted with permission from [155]. W. Windl. *IEICE Trans. Electron.* E86C, 269 (2003). © 2003. Oxford University Press.

Figure 34. (a) MOSFET structure with a strained Si channel caused by an underlying layer of SiGe. (b) Theoretical prediction of the mobility enhancement of electrons and holes from strained Si channels. (c) Strained buried SiGe quantum well PMOS structure. All figures. Reprinted with permission from [168], N. Collaert and P. Verheyen, *Strained Si and SiGe devices*. IMEC-ASD. http: www.imec.be/wwwinter/processing/asd/activities/strained.shtml. © IMEC.

is 4% larger than that of Si and is completely soluble in Si. Thus, the strained Si layer is under tensile and usually biaxial strain. Strained Si channels have been predicted to allow for enhanced mobilities for both electrons and holes (Fig. 34[b]). Strained Si devices work well for NMOS devices, where source and drain are $n$-type and mostly electrons carry the charge through the channel.

For PMOS devices ($p$-type source and drain with mostly hole transport through the channel), buried Si/SiGe devices have been found to perform better [169]. By introducing a SiGe quantum well beneath the Si/SiO$_2$ interface where the charge carriers can assume lower energy and thus would preferably travel. they can flow from source to drain with minimal interference of surface roughness and interface scattering. The intrinsic higher mobility of strained SiGe for holes introduces an extra advantage. The structure as such consists of a strained SiGe layer epitaxially grown on a Si substrate, followed by a Si cap layer which is needed to allow for the buried channel operation. but also for the growth of the gate oxide since the direct oxidation of SiGe is problematic (Fig. 34[c]) [168].

From the processing point of view, it has been assumed traditionally that the major effect of substrate stresses were dislocation formation and response [170, 171], whereas stress effects on diffusion were thought to be negligible [172]. With the reduction of gate lengths and the use of more exotic gate materials for modern nanoelectronic MOSFET structures. stress-mediated diffusion became a more prevalent component in determining the final dopant profile and subsequent device performance. On the experimental side. contradictory results for the qualitative influence of stress on diffusion especially in the case of boron further motivated a fundamental investigation of stress-mediated diffusion. While the measurements of Aziz et al. suggested for example enhanced boron diffusivity under compressive pressure [173–175], other work found retarded diffusion in that case [172, 176–178].

Most older theoretical work on stress-mediated diffusion assumed hydrostatic stress in the substrate. However, stresses caused by dislocations. thermal processes and geometric effects

all add to a complex stress state under a multi-layered gate stack—where the gate itself acts as a stress concentrator—with magnitudes approaching the material strength even at low temperatures [1, 172, 179, 180].

In this section, we apply the theory described in Section 2.1.3 to two examples of dopant diffusion in Si which are important for microelectronics. The first one, the neutral vacancy, is complicated by the existence of the Jahn-Teller distortion and the dependence of that distortion on stress and charge state [181]. The second example is the diffusion of a B-self-interstitial (B-$I$) pair, and the third one is the "ring-mechanism" for vacancy-assisted diffusion, which is believed to exist for Sb and (at least partially) As.

In all of the examples treated here, symmetry will dictate that two of the eigenvalues of all of the volume tensors will be degenerate. In the case of such degeneracy, we find it convenient to represent the volume tensor of a defect with orientation along direction $\hat{\mathbf{d}}$ by

$$\Omega_c = \Omega_{cl}\hat{\mathbf{d}} \otimes \hat{\mathbf{d}} + \Omega_{ct}(\mathbf{I} - \hat{\mathbf{d}} \otimes \hat{\mathbf{d}})$$

with "longitudinal" ($\Omega_{cl}$) and "transverse" (the doubly degenerate $\Omega_{ct}$) values. We can then describe the volume tensor by two parameters, either the combination $\Omega_{cl}$ and $\Omega_{ct}$, or

$$\Omega_{ch} \equiv \Omega_{cl} + 2\Omega_{ct}$$

$$\Omega_{ra} \equiv \Omega_{cl} - \Omega_{ct}$$

The latter pair measure the overall (scalar) volume and the anisotropic part.

We also note the jump direction is not generally the same as the symmetry axis of the saddlepoint, although for simple defects, such as the vacancy treated here, the two *are* the same (the nearest-neighbor hop also defines the symmetry axis of the saddlepoint). However, this is not true for the second example, the B$I$ pair.

### 3.3.1. Stress Effects on Vacancy Diffusion in Silicon

As noted above, the vacancy in Si may or may not undergo a Jahn-Teller distortion, depending on the charge state *and* the stress [44, 181, 182]. In the presence of a Jahn-Teller distortion, the vacancy in the valley has three possible orientations. Each orientation is symmetric around a $\langle 100 \rangle$ axis, so that the two transverse volumes are equal. The three orientations combined with two lattice sites in the primitive cell makes $N_{states} = 6$.

The solubility factor for the case with Jahn-Teller distortion is

$$S = e^{-\beta p \Omega_{ceh}}$$

$$\times \left[ e^{-\beta\Omega_{ra}(-2\sigma_{xx}+\sigma_{yy}+\sigma_{zz})} + e^{-\beta\Omega_{ra}(\sigma_{xx}-2\sigma_{yy}+\sigma_{zz})} \right.$$

$$\left. + e^{-\beta\Omega_{ra}(\sigma_{xx}+\sigma_{yy}-2\sigma_{zz})} \right]/3$$

where $p = Tr[\sigma]/3$. The case of no Jahn-Teller distortion is obtained by taking $\Omega_{cra} \to 0$, in which limit the solubility factor becomes $S = e^{-\beta p \Omega_{ceh}}$.

The vacancy hops by a jump to a nearest-neighbor site; the saddle point configuration has a symmetry axis along a {111} direction [44], and once again the two transverse volumes are equal (see Fig. 35). Defining

$$\alpha \equiv 2\beta\sigma\Omega_{c(s)a}/3$$

$$\delta \equiv e^{-(\alpha_{xy}+\alpha_{yz}+\alpha_{xz})} + e^{-(+\alpha_{xy}-\alpha_{yz}-\alpha_{xz})}$$

$$+ e^{-(-\alpha_{xy}+\alpha_{yz}-\alpha_{xz})} + e^{-(-\alpha_{xy}-\alpha_{yz}+\alpha_{xz})}$$

we find the components of the permeability tensor to be

$$P_{xx} = 2e^{-\beta p \Omega_{ceh}}[\cosh(2\alpha_{xy}) + \cosh(2\alpha_{zx})]/\delta$$

$$P_{xy} = -2e^{-\beta p \Omega_{ceh}}\sinh(2\alpha_{xy})/\delta \qquad (96)$$

The other components can be obtained by cyclic permutation of the Cartesian components $x, y, z$. The Jahn-Teller distortion of the vacancy in its equilibrium position has no effect on the permeability tensor.
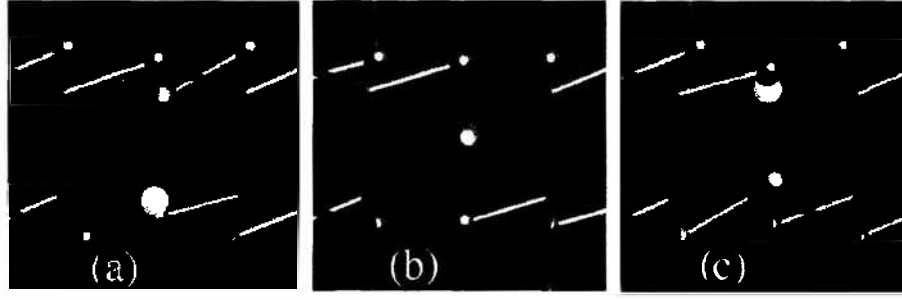
Figure 35. Projection of a Si crystal containing a migrating Jahn-Teller distorted vacancy on the (121) plane. The migrating Si atom (which moves in the opposite direction as the vacancy) is shown as a large dark ball, the vacancy in (a) and (c) by a small white ball. The small dark balls show the other Si atoms surrounding the two involved vacancy sites. (a) is the initial configuration. (b) the saddle point with a threefold symmetry around the migration direction (111) which lies in the paper plane, and (c) the final configuration. Reprinted with permission from [44]. M. S. Daw et al., *Phys. Rev. B* 64, 045205 (2001). © 2001, American Physical Society.

## 3.3.2. Effect of Stress on Boron Diffusion in Silicon

We have discussed boron diffusion, which is assisted by Si self-interstitials, in Section 3.2.1. Fig. 28[a] shows the diffusion mechanism for the B$I$ pair in $p$-type silicon (groundstate in Figs. 27[c] and initial and final structures in 28[a]), where the intermediate hexagonal interstitial (Fig. 27[b] and center of Fig. 28[a]) is a saddle point for the positive, and a local minimum close to the saddle point for the neutral charge state, and is in the following assumed to be the dominant saddle point for the purpose of examining stress effects in $p$-type Si.

### 3.3.2.1. Stress-Dependent Diffusion Coefficient for Interstitial-Assisted Boron Diffusion. The B$I$ pair results in more complex forms for the solubility and permeability factors than a simple point defect such as the vacancy discussed in the previous section. The defect in the valley has a $\langle 111 \rangle$ symmetry axis, with the B$I$ bond aligned so that the (substitutional) B lies along the line between the $I$ and a Si lattice atom (Figs. 27[c]; 28[a], initial and final structures) [183]. The three-fold symmetry around the $\langle 111 \rangle$ axis fixes the two transverse volumes as equal.

There are four orientations for the B$I$ pair on each site, and two lattice sites per primitive cell, so that $N_{states} = 8$. The solubility factor for the boron-self-interstitial pair is:

$$S = e^{-\beta p\Omega_{c(nh)}} \left( e^{\omega_{yz} - \omega_{zx} - \omega_{xy}} + e^{\omega_{xy} + \omega_{zx} + \omega_{yz}} \right.$$

$$\left. + e^{+\omega_{yz} - \omega_{zx} + \omega_{xy}} + e^{+\omega_{yz} + \omega_{zx} - \omega_{xy}} \right)/4$$

where $\omega \equiv 2\beta\Omega_{c(n)a}\sigma/3$.

The migration in $p$-type Si occurs when the $I$ pushes into the lattice site occupied by the B, which is displaced to a nearby hexagonal or quasi-hexagonal site (Fig. 27[b]); because the hexagonal interstitial is a saddle point for the positive, and a local minimum close to the saddle point for the neutral charge state [183], we assumed it to be the dominant saddle point in $p$-type Si. The quasi-hexagonal site has a $\langle 111 \rangle$ symmetry axis. The B in the quasi-hexagonal site is surrounded by 6 Si atoms, anyone of which may now be displaced by the B, leading again to a B-$I$ pair. The resulting hops can be 1NN, 2NN, or even 3NN. There are a total of 768 paths that contribute to the reduced rate matrix (Eq. [50]). The resulting permeability tensor is

$$\mathbf{P} = c_1 \rho + c_2 \chi \tag{97}$$

with

$$c_1 = \zeta/(4\lambda_x\lambda_y\lambda_z)$$

$$\rho_{xx} = (1 + \lambda_y^2\lambda_z^2 + \lambda_z^2\lambda_x^2 + \lambda_x^2\lambda_y^2)$$

$$\rho_{xy} = -5(1 - \lambda_y^2\lambda_z^2 - \lambda_z^2\lambda_x^2 + \lambda_x^2\lambda_y^2)/11$$

$$c_2 = \zeta/(132\lambda_x\lambda_y\lambda_z(1 + \lambda_y^2\lambda_z^2 + \lambda_z^2\lambda_x^2 + \lambda_x^2\lambda_y^2))$$

$$\chi_{\iota\iota} = -(1 + \lambda_1^2\lambda_2^2 - \lambda_2^2\lambda_1^2 - \lambda_1^2\lambda_1^2)^2$$

$$\chi_{\iota\nu} = -1 + 2\lambda_2^2\lambda_\nu^2 + \lambda_1^4\lambda_2^4 + \lambda_2^4\lambda_1^4 - \lambda_1^4\lambda_1^4 - 2\lambda_1^2\lambda_1^2\lambda_2^4 \tag{98}$$

$$\zeta \equiv \exp\left(-\beta p\Omega_{(s)h}\right)$$

$$\lambda_z \equiv \exp\left(2\beta\Omega_{(s)a}\sigma_{\nu\nu}/3\right)$$

$$\Omega_{\iota\iota\nu h} \equiv \Omega_{c(s)l} + 2\Omega_{r(s)l}$$

$$\Omega_{\iota(s)h} \equiv \Omega_{c(s)l} - \Omega_{\iota(s)l}$$

### 3.3.2.2. Diffusion Model and Atomistic Calculation of the Parameters.

In the seminal work by Laudon et al. [184], the first multiscale modeling approach was presented to model effects of anisotropic stress on diffusion in nanoscale MOSFET structures. A standard four-stream diffusion-reaction model for B was used as a starting point, consisting of the concentration fields for $I$ (self-interstitials, mobile), $V$ (vacancies, mobile), B (substitutional B, immobile), and B$I$ (B- interstitial pair, mobile). Conjecturing that the point-defects ($I$ and $V$) are in equilibrium with a free surface and that the $I$ concentration is independent of the B concentration, the four-stream model was simplified to an effective one-stream model,

$$\frac{\partial C_B}{\partial t} = \nabla \cdot [\mathbf{P}_B^{\text{eff}}\nabla(C_B/S_B^{\text{eff}})] \tag{99}$$

where $C_B$ is the B concentration, $\mathbf{P}_B^{\text{eff}}$ is the effective solid permeability tensor, and $S_B^{\text{eff}}$ is the solid solubility factor [44]. We have demonstrated such a simplification process and the arguments that are used for its justification for the derivation of Eq. (21) in Sec. 2.1.1. In the hydrostatic case where the stress tensor $\sigma$ is given by $\sigma = p\,\mathbf{Id}$, the permeability (which is a scalar now) as a function of the pressure $p$ is given by

$$P_B^{\text{eff}} = P_{BI} = D_B^0\exp\left(-\frac{\varepsilon_s^{BI} + p\Omega_s^{BI}}{k_B T}\right) \tag{100}$$

where $D_B^0$ is the diffusivity prefactor for intrinsic B diffusion, which was calculated from first principles within the harmonic Vineyard method [185], $\varepsilon_s^{BI}$ is the creation energy for the B$I$ pair at the saddle point as introduced in Eq. (35), and $\Omega_s^{BI}$ is the corresponding creation volume from Eq. (54). The corresponding solubility is

$$S_B^{\text{eff}} = \exp\left[-\frac{\varepsilon_v^B + \varepsilon_{at} + p(\Omega_v^B + \Omega_{at})}{k_B T}\right] \tag{101}$$

where $\varepsilon_v^B$ is the creation volume for substitutional B in its ground state or "valley," $\varepsilon_{at}$ is the total energy per atom of the perfect Si cell. $\Omega_v^B$ is the corresponding creation volume for substitutional B, and $\Omega_{at}$ the volume per atom of perfect Si. For the anisotropic case with general stress tensor, the expressions are more complicated.

The general stress dependence is given by the respective creation volume tensors which are calculated by the length changes $\Delta L_\alpha$ between the defective cell and the perfect Si cell with lattice parameters $L_\alpha$, $\Omega = (\Omega)_{\alpha\beta} = \delta_{\alpha\beta}\varepsilon_{\alpha\gamma_\lambda}\Delta L_\alpha L_\gamma L_\delta$ in a principal-axes system. Because the single elements of the volume tensors are hard to separate experimentally and had not been known previously, they were calculated from first principles in Ref. [184], although several ab initio investigations had examined the (scalar) hydrostatic pressure dependence of diffusion before [182, 186, 187].

The values proposed in Ref. [184] were $\varepsilon_r^B + \varepsilon_{at} = -7.14$ eV and $\varepsilon_s^{BI} = -3.39$ eV, which resulted in a net activation energy of $E_a = -(\varepsilon_v^B + \varepsilon_{at}) + \varepsilon_s^{BI} = 3.75$ eV in good agreement with experiment [188]. Furthermore, the creation volumes were found to be $(\Omega_s^{BI})_{\alpha\beta} = \delta_{\alpha\beta}[9.5\delta_{\alpha x} - 3.8(\delta_{\alpha y} + \delta_{\alpha z})]$ Å$^3$ in a principal-axis system with the $x$-axis parallel to the (111) direction and $\Omega_v^B + \Omega_{at} = 2.4\,\mathbf{Id}$ Å$^3$, respectively, with scalar values for the hydrostatic case of $\text{Tr}(\Omega_s^{BI}) = 1.9$ Å$^3$ and $\text{Tr}(\Omega_v^B + \Omega_{at}) = 7.2$ Å$^3$. The corresponding net scalar activation volume was $V_a = -\text{Tr}(\Omega_v^B + \Omega_{at}) + \text{Tr}(\Omega_s^{BI}) = -5.3$ Å$^3$, in good agreement with previous experiments [174] of $-3.4$ Å$^3$ and isotropic ab initio calculations [187] of $-3.1$ Å$^3$. Although

the hydrostatic value is small, there is considerable anisotropy in the permeability volume tensor which can have a significant effect on diffusion under anisotropic strain. These results suggest that in the considered equilibrium case B diffusion is *enhanced* by compressive pressure. This is mostly a solubility effect and, in consequence, due to the fact that the point defects were assumed to be in equilibrium with a free surface.

### 3.3.2.3. Results for metal gate MOSFET structure.

In Ref. [184], a titanium nitride (TiN) metal gate integration on a sub-quarter micron *p*-MOSFET was used as a demonstration example of the stress-diffusion phenomena where anomalous stress-dependent boron diffusion had been recently discovered at Motorola [1]. Electrically measured lateral diffusion results indicated an enhancement in boron diffusion with increasing gate stress. Figure 36 shows an example scanning electron microscope image of the gate along with a feature scale stress model.

A finite-element method was used to predict high-temperature feature scale stresses in the Si substrate under the TiN metal gate. Since the temperature dependence of elastic properties for most materials in a gate stack, except for Si, is unknown, empirical data are used to calibrate the high-temperature stress simulations. A full 3D finite-element model had to be used, since plane-strain and plane-stress 2D reductions are insufficient in the area of interest near the Si surface. A 100 Å SiO$_2$ film was modeled over the gate in order to reduced the potential singularity found in finite-element peeling stress results near free traction boundaries [180]. The TiN gate has a high tensile stress at anneal temperatures (1025°C), resulting in compressive horizontal stresses directly under the gate and large compressive and tensile stress concentrations just under and outside the gate edge, respectively.

Resulting 3D stress tensors from the finite-element model were passed through nodal data to a stress diffusion solver based on the partial differential equation solver described in [189]. This solver employed the gradient-weighted moving finite element method which uses a continuously moving mesh that adapts to the evolving solution. The diffusion equation implemented was

$$S_B^{cH} \frac{\partial w}{\partial t} = \nabla \cdot P_B^{cH} \nabla w + \nabla w \cdot P_B^{cH} \nabla w \qquad (102)$$

where the transformed variable $w = \log(C_B/S_B^{cH})$ was introduced to achieve a relative accuracy in the concentration tail comparable to that in the high-concentration regions. Galerkin equations were obtained by minimizing the residual of this equation with respect to a $1/S_B^{cH}$ weighted $L^2$ norm. The stress tensor, upon which $S_B^{cH}$ and $P_B^{cH}$ depend, was simply obtained by interpolating the finite-element computed stress field onto the moving mesh. Using the described procedure, post anneal diffusion profiles for the TiN metal gate as well as for a reference stress-free gate were calculated.

The resulting profiles for a 5 keV implanted B profile diffused at 1025°C for 10 s can be seen in Fig. 37. Because of stress effects, an 8% change in $L_{eff}$ for a 250-nm $L_{drawn}$ device was predicted. Equivalently, a 30% change in $L_{eff}$ was predicted for a 65-nm $L_{drawn}$ device. These numbers were quantitatively in good agreement with the experimental findings.



Figure 36. Scanning electron microscope image [1] and finite-element model of the TiN gate stack of [1]. Reprinted with permission from [184], M. Laudon et al., *Appl. Phys. Lett.* 78, 201 (2001). © 2001, American Institute of Physics.

**Figure 37.** Boron concentration contours for a 10 s, 1025°C. boron diffusion (source/drain + extension). Comparison of the TiN metal gate stressed case (heavy line) to a unstressed (light line) solution show a 8% difference in lateral diffusion for a 250-nm $L_{gate}$ gate. Reprinted with permission from [184]. M. Laudon et al., *Appl. Phys. Lett.* 78, 201 (2001). © 2001, American Institute of Physics.

### 3.3.3. Ring Mechanism for Vacancy-Assisted Diffusion in Silicon

As the most complicated case discussed to date, Daw and Windl have shown how the macroscopic diffusivity depends on the microscopic hopping for dopants migrating via the commonly assumed ring mechanism for vacancy diffusion [190]. The diffusion of Sb and (at least partially) As in silicon has been suggested to occur via the ring mechanism [153]. In the following, we assume As as the dopant, but of course any other vacancy-assisted diffuser can be substituted for it. The vacancy mechanism involves a number of vacancy hops beyond the simple exchange between impurity and vacancy to result in net mass transport. In particular, it is clear that the $V$ must—for the minimum-distance path—move around one of the (110) quasi-hexagonal rings in order to allow the As to accomplish a macroscopic hop. Indeed, from the energetics in [153], we see that the $V$ hopping around the ring limits the rate of macroscopic As diffusion. We can use the theoretical procedure of Daw et al. [44] to calculate the solubility and permeability factors (and hence the diffusivity) for this complex diffusion mechanism.

First, we enumerate the states of the As$V$ (recall that the theoretical treatment from Section 2.1.3 requires that the As$V$ moves as a unit, if even a loosely bound one). Each state must be identified as belonging to a particular unit cell; then we must distinguish between various states within each unit cell. In the present case we find it convenient to describe the general state of the As$V$ defect by a pair of positions. (1) the *absolute* position of the As atom, and (2) the position of the $V$ *relative* to the As. For each As position, there are four positions for $V$ as 1NN, 12 as 2NN, and 12 as 3NN. There are two As positions within the diamond unit cell. The total number of states of the As$V$ defect per unit cell of Si is then 56.

The solubility factor for the ring mechanism is easy to anticipate. The As$V$ pair at 1NN has energy $\varepsilon_1$, at 2NN it has $\varepsilon_2$, and at 3NN it has $\varepsilon_3$. There are three times as many 2NN sites as 1NN, and the same for 3NN to 1NN. Thus,

$$S = (s_1 + 3s_2 + 3s_3)/7 \tag{103}$$

with $s_i = \exp(-\beta\varepsilon_i)$. According to Ref. [153], the 1NN position is easily the most stable in the case of As, so that at temperatures below $(\varepsilon_2 - \varepsilon_1)/k_B \sim 8000$ K, the solubility factor is dominated by the As$V$ energy at 1NN.

The permeability is constructed from the hops as described in Section 2.1.3. The rate matrix is built from the terms describing the hops which occur among the various As$V$ states. Using translation symmetry, we can block-diagonalize the rate matrix to pieces on the order of the size of the number of states per unit cell. Thus, the reduced rate matrix for this problem will be of size 56 × 56. There are two types of hops to consider. (1) $V$ hops

around the ring, leaving As fixed, and (2) $V$ interchanges with As (1NN relationship before and after).

The As site and its 1NN, 2NN, and 3NN sites are arranged in six-fold quasihexagonal rings. For the $V$ to traverse a ring (while the As remains fixed) requires it to move from a 1NN site to 2NN, 2NN to 3NN, 3NN to a different 2NN and on to a new 1NN site. From each 1NN site, the $V$ has several possible rings that it could use to return to another 1NN site. The $V$ hops are then of two basic types, from 1NN to 2NN (or vice versa), which we designate "H12" and from 2NN to 3NN (or vice versa), called "H23."

In addition to the hops where the As remains fixed, we need to include the jump where the As interchanges with a 1NN $V$, which we designate as "H01." In all, there are 152 hops to be included in the rate matrix: 8 are H01, 48 are H12, and 96 are H23. The resulting permeability tensor is isotropic, with the (scalar) permeability constant

$$P = \frac{a^2}{56} \cdot \frac{p_{01}\, p_{12}\, p_{23}}{2\, p_{01}\, p_{12} + 4\, p_{12}\, p_{23} + 3\, p_{01}\, p_{23}} \tag{104}$$

with $p_{ij} = \exp(-\beta \varepsilon_{ij})$. $\varepsilon_{ij}$ is the energy of the saddlepoint between valley $i$ and valley $j$, referenced to the same energy as the valley states.

The structure of $P$ is not surprising. It reveals that the As$V$ diffusion is governed by all three hop types (H01, H12, and H23) in series. The effective inverse rate (inverse of $P$) is proportional to the sum of the inverse rates; that is

$$\frac{1}{P} \sim \frac{4}{p_{01}} + \frac{3}{p_{12}} + \frac{2}{p_{23}}$$

In the case that one of the three saddle points is much higher than the others, it will dominate. According to Ref. [153], H23 is higher in energy than the other two, so that at temperatures below $(\varepsilon_{23} - \varepsilon_{12})/k_B \sim 12,000$ K the permeability is dominated by the hops between 2NN and 3NN sites.

The diffusion constant is $P/S$. Using the numbers from Ref. [153] for all the equilibrium states and saddle points, the permeability is dominated in the relevant temperature range by $p_{23}$ and the solubility factor by $s_1$, so the diffusivity has a simple Arrhenius form with the energy $E_d \approx \varepsilon_{23} - \varepsilon_1$, which is about 0.9 eV, consistent with the results noted in Ref. [153].

**Effect of stress on As V diffusion by the ring mechanism.** Stress effects are accounted for by knowing the volume tensors at each of the valleys and saddlepoints. There are three different types of valleys, and fundamentally we would need three volume tensors to describe the effects of stress on the solubility. Likewise, there are three different types of saddlepoints, and three volume tensors to describe the effects of stress on the permeability. However, this leads to such complicated intermediate expressions that simplifying assumptions become necessary to make the computation tractable. Fortunately, for the temperature range of interest, we are able to make such simplifying assumptions:

1. The solubility factor is dominated by the As$V$ energy at 1NN, so we will include only the volume tensor at 1NN. This defect has a $\langle 111 \rangle$ symmetry axis, which means only two parameters are required for its description (e.g., one longitudinal and one transverse).
2. The permeability factor is dominated by the H23 saddlepoint, so we will include only the volume tensor at this position. The H23 saddlepoint is formed by placing the $V$ between 2NN and 3NN from the As atom. The separation between the As and $V$ is sufficiently large at this point that we will approximate the volume tensor by the sum of the two volume tensors for infinite separation (that is, the As valley volume tensor and the volume tensor for $V$ at its saddlepoint). The saddlepoint for $V$ migration has a $\langle 111 \rangle$ symmetry axis, and the As valley volume is isotropic (no Jahn-Teller distortion), so the resulting volume tensor has two parameters required for its description.

Because of these simplifying assumptions, the permeability and solubility factors are both determined by volumes with $\langle 111 \rangle$ symmetry.

This implies that uniaxial stress along a $\langle 100 \rangle$ axis does not break the isotropic symmetry of the permeability tensor, whereas for uniaxial stress along $\langle 110 \rangle$ and $\langle 111 \rangle$, the permeability

tensor is anisotropic. Furthermore, this means that we need to consider only deviatoric (shear) stresses of the general type (using the diamond cube axes as reference)

$$
\begin{pmatrix}
0 & \sigma_{xy} & \sigma_{zx} \\
\sigma_{xy} & 0 & \sigma_{yz} \\
\sigma_{zx} & \sigma_{yz} & 0
\end{pmatrix}
\tag{105}
$$

The solubility factor is then

$$
S = (\bar{s}_1 + 3s_2 + 3s_3)/7
\tag{106}
$$

with $s_2$ and $s_3$ the same as without stress (see Eq. [103]) and

$$
\begin{aligned}
\bar{s}_1 = s_1 \exp[-\beta p\Omega_h^r]\frac{1}{4}\big( & \exp[2\beta(-\sigma_{xy} + \sigma_{yz} + \sigma_{zx})\Omega_a^r/3] \\
& + \exp[2\beta(\sigma_{xy} + \sigma_{yz} - \sigma_{zx})\Omega_a^r/3] \\
& + \exp[2\beta(\sigma_{xy} - \sigma_{yz} + \sigma_{zx})\Omega_a^r/3] \\
& + \exp[2\beta(-\sigma_{xy} - \sigma_{yz} - \sigma_{zx})\Omega_a^r/3]\big)
\end{aligned}
\tag{107}
$$

The permeability factor depends only on the off-diagonal components of the stress tensor (in the coordinate frame where the cube axes are the coordinate axes). The components of the permeability factor are

$$
\begin{aligned}
P_{xx} = P_{yy} &= P_{zz} = (\xi_x^2 + \xi_y^2 + \xi_z^2 + \xi_x^2\xi_y^2\xi_z^2)/(4\xi_x\xi_y\xi_z) \\
P_{xy} &= (\xi_x^2 + \xi_y^2 - \xi_z^2 - \xi_x^2\xi_y^2\xi_z^2)/(12\xi_x\xi_y\xi_z) \\
P_{zx} &= (\xi_x^2 - \xi_y^2 + \xi_z^2 - \xi_x^2\xi_y^2\xi_z^2)/(12\xi_x\xi_y\xi_z) \\
P_{yz} &= (-\xi_x^2 + \xi_y^2 + \xi_z^2 - \xi_x^2\xi_y^2\xi_z^2)/(12\xi_x\xi_y\xi_z)
\end{aligned}
\tag{108}
$$

where

$$
\begin{aligned}
\xi_x &\equiv \exp(-2\beta\Omega_x^a\sigma_{yz}/3) \\
\xi_y &\equiv \exp(-2\beta\Omega_y^a\sigma_{zx}/3) \\
\xi_z &\equiv \exp(-2\beta\Omega_z^a\sigma_{xy}/3)
\end{aligned}
\tag{109}
$$

This expression for **P** has factored out of it the value at zero stress (Eq. [104]) and also the hydrostatic component (which is $\exp(-\beta\Omega_x^h p)$, where $p$ is the pressure).

## 3.4. Atomistic Modeling of Extended Defect Formation in Ion-Implanted Silicon

The complexity of the atomic level interaction mechanisms involved in today's semiconductor device processing is a subject of concern for the microelectronics industry, as had been already pointed out in the 1997 SIA Roadmap [191]:

> Continuum physics models are no longer sufficient below 100 nm. Tools are needed for the physical and chemical processes at an atomic level.

Furthermore, the discreteness of the channel dopants has been shown to give rise to both a shift of the threshold voltage (as compared with the prediction from the continuum approach) and to an asymmetry in drain current upon interchanging the source and drain [192] in addition to the statistical effects [52] discussed in Section 2.1.4.

When studying discreteness and statistical effects for small nanodevices, only an atomistic description of the system can provide the necessary results. Jaraiz et al. have shown that such a task is well within the reach of present-day atomistic simulation tools, especially within the kinetic Monte Carlo approach [70]. They demonstrated the atomic-level simulation of a 50 nm $n$-MOSFET structure, which included the 40-nm S/D extension and half the channel region (see Fig. 38). The simulated process flow consisted of a 5 keV, $10^{14}$ As cm$^{-2}$ S/D extension implant, a subsequent 70 keV, $10^{13}$ BF$_2$ cm$^{-2}$ SPI implant and a 10s, 950°C anneal.

**Figure 38.** Fifty-nanometer N-MOSFET fabricated as described in the text. The simulated region is indicated in the figure. Reprinted with permission from [70]. M. Jaraiz, in "The Encyclopedia on Materials Science and Technology" (K. H. J. Buschow, R. W. Cahn, M. C. Flemings, B. Ilschner, E. J. Kramer, and S. Mahajan, Eds.). p. 7861. Pergamon, Boston, 2001. © 2001, Elsevier.

Figure 39 shows a snapshot of the simulation after 25 ms annealing time at which the average size of {311} interstitial clusters reaches its maximum. Figure 39(a) and 39(b) are plan views of the {311} defects and As implanted atoms, respectively. Figure 39(c) and 39(d) are cross-sectional views. Assuming non-amorphizing conditions, microvoids form near the surface, whereas interstitials are pushed beyond the As implant range due to the heavy As ion mass. The Si interstitials generated by the arsenic implant induce TED for the boron atoms already present in the channel and deplete the channel to the S/D transition region. This effect is known as the reverse short channel effect with detrimental consequences in deep submicron devices.

Since we already have discussed atomic-level simulation work on defect-assisted diffusion, including transient enhanced diffusion, for boron (Section 3.2) and arsenic (Section 3.1.2), and have also discussed in the latter section the basic technicalities of the kinetic Monte Carlo approach, we want to concentrate in this section on the evolution of extended defect clusters during annealing, which consist of vacancy clusters or voids and interstitial clusters, which, before forming dislocation loops, assume the shape of rodlike {311} defect clusters. We discuss both the use of the kinetic Monte Carlo technique (for void formation, see Section 3.4.1) and accelerated dynamics (for {311} formation, see Section 3.4.2). Atomistic methods are preferable in this case, since especially for interstitial clusters, not only their concentration (which we could model as well within the continuum approach), but also their shape is important because of the strongly anisotropic growth mechanism, as discussed in Section 3.4.2.

### 3.4.1. Void Formation in Ion-Implanted Silicon

The increasingly stringent requirements on the quality of crystalline semiconductor substrates has continued to fuel fundamental research aimed at optimal control of defect diffusion, nucleation, and growth. Commercial silicon wafers currently are grown under vacancy-rich conditions to avoid the aggregation of self-interstitials into dislocation loops and other microdefects such as the oxidation-induced stacking-fault ring [193], all of which are devastating to microelectronic device yield and performance. However, excess vacancies also can aggregate during crystal growth to form (111)-oriented octahedral vacancy clusters (voids) [194] of up to several hundred nanometers in diameter [195].

Equivalently, high-dose implantation of heavy ions at elevated temperatures (to prevent amorphization) can lead to the agglomeration of vacancies in the form of voids, since the large momentum of these heavy ions drives the silicon atoms deeper into the substrate, leaving behind a vacancy rich region near the surface, where the voids are formed.

Voids can act as gettering centers because their internal surface behaves like a clean silicon surface [11]. Voids have also been linked to the degradation of microelectronic device performance, particularly gate oxide integrity. Thus, understanding the vacancy clustering process is important to improve the control of void formations, especially in view of the potentially increasing use of heavier ions for the improved formation of ultrashallow junctions in present-day semiconductor devices [2].

Figure 39. Simulated atomistic configurations of the 50-nm N-MOSFET device of Fig. 38, after a short annealing of 25 ms. (a) and (b) are plan views of the {311} defects and As implanted atoms, respectively, while (c) and (d) are their respective cross sectional views. Reprinted with permission from [70], M. Jaraiz, in "The Encyclopedia on Materials Science and Technology" (K. H. J. Buschow, R. W. Cahn, M. C. Flemings, B. Ilschner, E. J. Kramer, and S. Mahajan, Eds.), p. 7861. Pergamon, Boston, 2001. © 2001, Elsevier.

For a successful modeling of void growth, an energy model is needed to identify the binding energy of vacancies in voids of varying sizes $N$. Since voids can grow to diameters of many nanometers [195], an analytical model for the binding energy as a function of void size is desirable. Such models have been proposed, for example, by Prasad and Sinno [196], on the basis of molecular dynamics simulations, and by Jaraiz et al. [197], based on scaling arguments. The binding energy for a vacancy cluster (void) of size $N$ is usually defined as

the energy difference between two configurations for $N$ vacancies, a single vacancy plus a size-$(N - 1)$ void, and a size-$N$ void.

In Prasad and Sinno's work, a powerlaw fit demonstrates $N^{2/3}$ (exponent $= 0.64$) scaling behavior across the entire size range considered. These results indicate that a phenomenological representation of cluster energies; that is, $E_N^f = \sigma(T)N^{2/3}$ is appropriate even for very small clusters, where $\sigma(T)$ is an effective surface energy. For the effective surface energy of larger clusters, they found $\sigma(T) = (1.05 - 1.69 \times 10^{-4}T)$ J/m$^2$.

Jaraiz et al. [197] have assumed the same $N^{2/3}$ scaling for the void energy with the argument that the energy should be proportional to the number of vacancies at its surface. Taking 3.65 eV as the vacancy formation energy at the free surface ($N \to \infty$), and a prefactor to fit the experimental activation energy for the divacancy of 1.2 eV, they approximated the vacancy binding energy by [197]

$$E_{bind}^v = 3.65 + 4.9[(N - 1)^{2/3} - N^{2/3}]$$ (110)

However, this expression carries no temperature dependence, which Prasad and Sinno have shown to be crucial for a good description of, for example, the experimental void aggregation onset temperature [196], where the reduction of the effective surface free energy at high temperatures due to the significant entropic contribution was an important ingredient for the correct simulation of void growth.

The assumption of a temperature-independent cluster binding energy might by the reason why in Ref. [197], the simulation of post-implant void growth after a 100 keV, $10^{16}$ cm$^{-2}$ As$^+$ implant [198] predicts a depth distribution of voids in good agreement with experiment, but an average void size smaller than in the experiment. However, this is difficult to say since there is considerable uncertainty about the temperature in the experimental work. Another possible explanation suggested in Ref. [199] is that smaller clusters should have a larger surface energy due to their curvature and, therefore, are less stable, which is consistent with the findings of Prasad and Sinno [196], where the surface energy $\sigma(T)$ is found to be constant for cluster sizes $N > 100$, but increases (smoothly) from the $N = 100$ value of 1.1 J/m$^2$ to around 1.3 J/m$^2$ for very small clusters.

Figure 40(a) shows a cross-sectional view of the 100 keV, $10^{16}$ cm$^{-2}$ As$^+$ implant [198]. Beam heating with a high flux ion beam was used to prevent amorphization. A 45 nm region can be identified, extending from the surface, which contains a high density of voids followed by a band of dislocations extending to about 200 nm. The KMC simulation, shown at the same scale in the bottom figure, predicts the formation of voids and interstitial clusters within the same depth ranges under those implant conditions, but with a different size distribution as discussed above, using the expression for vacancy clustering (or void formation) binding energy from Eq. (110).

## 3.4.2. Accelerated Dynamics Simulations of Interstitial-Cluster Growth

Besides void formation through vacancy clustering, self-interstitials can in a similar fashion agglomerate subsequent to ion implantation and form {311} rod-like interstitial clusters and dislocation loops. Once they have formed, they control the dopant diffusion in ion-implanted silicon by providing traps for and sources of mobile interstitials [71, 74, 156]. Understanding the energetic and dynamical properties of interstitial defects is thus an essential step in accurately modeling the time evolution of dopant profiles [74, 200]. Recent theoretical studies of interstitial defects [75, 201–204] elucidate their growth path: interstitial clusters become more stable by adding more interstitials. This is consistent with the thermal behavior of boron transient enhanced diffusion with or without large-scale interstitial defects [71, 156]. However, the microscopic understanding of the critical steps leading to the cluster growth is limited. Cost for atomistic simulations grows prohibitively for the increasing complexity of the system. Furthermore, the long time scale involved with the diffusion and nucleation is not easily achieved by conventional molecular dynamics simulations involving many degrees of freedom.

Although the use of kinetic Monte Carlo simulations for the study of {311} formation is well capable of demonstrating the discrete nature, geometry and orientational disorder of

Figure 40. Cross-sectional TEM of a beam-heated, 100 keV, $10^{16}$ cm$^{-2}$ As$^+$ implant [198], compared with the corresponding KMC simulation, drawn to the same scale (units in nm). Reprinted with permission from [199]. M. Jaraiz et al., Mater. Sci. Semicond. Proc. 3, 59 (2000). © 2000, Elsevier.

the defects in space [70] (see Fig. 41), a more detailed description is necessary to understand the growth mechanism of the defects, which is responsible for their shape and size, and thus produce a sensible event catalog for the KMC simulations.

In order to model the necessary longer time scale with atomic resolution, Birner et al. [205] have applied the parallel-replica method (Section 2.2.3) to study the growth of silicon interstitial clusters. The silicon-silicon interaction is described using a classical potential [206]



Figure 41. (a) Plan view TEM micrograph [156] and (b) KMC simulation of a 40 keV 5 × $10^{13}$ cm$^{-2}$ silicon implant after a 30 s annealing at 800 °C, exhibiting {311} defects. Reprinted with permission from [70]. M. Jaraiz, in "The Encyclopedia on Materials Science and Technology" (K. H. J. Buschow, R. W. Cahn, M. C. Flemings, B. Ilschner, E. J. Kramer, and S. Mahajan, Eds.), p. 7861. Pergamon, Boston. 2001. © 2001, Elsevier.

of the modified embedded atom method (MEAM) form. The cluster growth was modeled by (1) coalescence of randomly distributed interstitials for small clusters and (2) mobile-interstitial capture by an $n$-interstitial cluster. Local minimum configurations and transition states extracted from collected trajectories of many parallel replica runs provide microscopic models for the interstitial-diffusion and interstitial-trapping processes. Using between 4 and 32 processors (dephasing time 1–2 ps), simulations times as long as 0.1 $\mu$s were achieved, allowing accurate estimates of dynamical quantities such as diffusion constants of mobile interstitial defects.

**Small interstitial clusters: $n = 2 - 4$.** Randomly distributed isolated interstitials coalesce into small clusters. For example, di-interstitial $I_2$ clusters are invariably formed when two interstitials are located within the third-nearest neighbor sites of a silicon lattice. At 500 and 800 K, four interstitials in a 512-atom unit cell form compact clusters, after making several interstitial jumps. With the used potential, the simulations predict a single-interstitial diffusion constant of $D_I = 4.5 \times 10^{-6}$ cm$^2$s$^{-1}$ exp($0.2$ eV$/(k_B T)$).

$I_2$ clusters are found to be as mobile as $I_1$'s with a migration energy of 0.1 eV. This is consistent with the *ab initio* and tight-binding predictions that $I_2$'s are important mobile components in interstitial-supersaturated silicon [202]. The formation of a tri-interstitial cluster $I_3$ occurs in two steps: (1) formation of $I_2$ and (2) its subsequent capture of an interstitial. Among the local minimum structures identified is the compact $I_3$ with $T_d$ symmetry [75, 204].

Metastable precursors dominate the formation of an $I_4$ cluster. Figure 42 shows the local minima found during a parallel replica run at 800 K of the $I_4$ formation. Within 0.2 ns, a metastable $I_4$ cluster is formed, releasing 5.6 eV. Once a metastable $I_4$ cluster is formed, transitions between related structures and diffusion of the cluster occur for 2 ns. After about 300 transitions, the cluster falls into the groundstate configuration ($D_{2d}$ symmetry) whose core structure is shown in Fig. 42.

About 1 eV is released when the ground-state $I_4$ is formed by bonding rearrangements from metastable precursors. Both $I_2$ and $I_3$ clusters are precursors of an $I_4$ cluster, consistent with interstitial cluster growth models of small clusters [74, 200].

The classically simulated $I_4$ ground state is in good agreement structurally and energetically with *ab initio* calculations [75, 203]: a low 1.4 eV formation energy versus the 1.5 eV *ab initio* result. Both experiment [74] and calculations [75, 201, 203] predicted the $I_4$-$D_{2d}$ to be extremely stable. Indeed tens of nanoseconds additional simulation at both 500 and 800 K find no lower-energy state.



Figure 42. Formation energy of a four-interstitial defect during a parallel-replic run at 800 K and the core structure of the ground $I_4$-$D_{2d}$ cluster. The starting configuration consists of four randomly distributed $I$-interstitials in a 512-atom cell. After 20 interstitial jumps, a metastable $I_4$ precursor is formed, releasing 5.4 eV (off the scale). Once the transition to $I_4$-$D_{2d}$ occurs, no more transitions are observed within 10 ns. Reprinted with permission from [205] S. Birner et al., Solid State Commun. 120, 279 (2001). © 2001, Elsevier

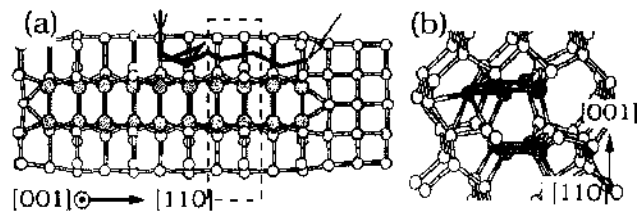Figure 43. Core structures of $I_5$ local minima identified by parallel-replica runs. The captured interstitial is shown as solid circles. Note that (a)-(c) contain an interstitialcy or split interstitial. Using a MEAM potential, (a) is the ground structure. Metastable structures (b)-(d) are typically found at the moment of interstitial capture. Gray atoms denote the $I_4-D_{2d}$ structure that remains intact after formation of $I_5$ by an interstitial capture. Reprinted with permission from [205]. S. Birner et al., *Solid State Commun.* 120, 279 (2001). © 2001, Elsevier.

### Intermediate interstitial clusters: n = 5 – 20.

For the simulation of interstitial-cluster growth larger than $n = 4$, a different approach was used. Starting configurations combined (1) an interstitial cluster $I_n$ with (2) an $I_1$ or $I_2$ placed at a distance far from the existing cluster. The final cluster $I_{n+1}$ or $I_{n+2}$ is then used as the initial cluster for subsequent simulations. These initial configurations are chosen to simulate the assumption that every reaction can be separated into subsequent binary reactions as outlined in Section 2.1.1, an approximation also used in kinetic Monte Carlo methods (Section 2.2.2).

The simulations of [205] find in general that the cluster growth occurs by (1) interstitial capture at the cluster boundaries and (2) subsequent local rearrangements of the core atoms to lower-energy structures.

As an example, the starting configuration of a compact $I_4$ cluster captures after several interstitial hops an additional interstitial present in the initial structure. Figure 43 shows four low-lying core structure of $I_5$ identified during the parallel-replica simulations. With the classical MEAM potential, Fig. 43(a) is the ground state which is found most frequently independent of the initial configurations. The metastable structures Fig. 43(b)-43(d) are characteristic transient structures found at the early stage of the interstitial trapping process. Interestingly, the $I_4$ core structure remains intact in $I_5$. Further interstitial injection into the cell containing a compact $I_5$ results in the formation of compact clusters as large as $n = 20$.



Figure 44. Schematics for interstitial trapping by an elongated cluster: (a) an interstitial (solid circle) is placed at the nearest hexagonal site; and (b) an extra chain is formed by concerted motions by an interstitial and the atoms at the end. The gray atoms denote 〈110〉 interstitialcies constituting a chain. The transition from (a) to (b) releases 2 eV. Reprinted with permission from [205]. S. Birner et al., *Solid State Commun.* 120, 279 (2001). © 2001, Elsevier.

Figure 45. (a) Schematics for interstitial trapping by an elongated cluster. An interstitial is first trapped in the middle of the chain. The solid line is a trajectory of the interstitial making jumps in the [110] chain directon. Eventually, the interstitial is incorporated into the chain end indicated by the arrow on the right. (b) Local mininum configuration, when the interstitial is located in the dotted region of (a). Reprinted with permission from [215], S. Birner et al., *Solid State Commun.* 120, 279 (2001). © 2001, Elsevier.

Using *ab initio* calculations, the found state configuration changes with size: compact for small cluster and, for larger clusters, interstitial chains forming extended {311} defects [5]. Chainlike, elongated clusters are local minimum structures of any interstitial cluster larger than $n = 3$. The stability of interstitial chains is achieved by maximizing the ratio of fourfold coordinated atoms to the number of interstitials incorporated into a cluster. The elongated clusters are basic building blocks of extended {311} defects [75, 207, 208]. Therefore, the growth mechanism of elongated clusters is an important step in understanding the formaton of extended {311} defects.

The role of the shape of interstitial traps is investigated by performing parallel replica simulations using chain-like clusters as nuclei for cluster growth. Figure 44 shows the chain structure before and after interstitial capture. By concerted motions, an interstitial can be easily trapped at the chain end. In Figure 44, the interstitial chain extends from $n = 5$ to 6. This extension releases about 2 eV. Many trapping paths for various trajectories of an incoming interstitial were found. Once a metastable elongated cluster is formed, the transition from an elongated shape to a compact one would require bond rearrangements involving many atoms, thus becoming kinetically inaccessible within typical simulation times of 1–10 ns.

Although the most energetically favorable trapping occurs at the chain ends, the active capture sites extend along the entire interstitial chain. Figure 45 illustrates a growth mech-anism for an elongated cluster by an interstitial capture in the middle of a chain, folloved by interstitial diffusion steps to the chain end. Seven-member rings surrounding a chain are efficient interstitial trap sites with an interstitial-binding energy of 1 eV. The interstitial cap-tured in the seven-member rings diffuses by making random jumps in the chain direction. Eventually, the interstitial settles at the chain end, releasing an extra 1 eV.

The interstitial capture mechanisms of Figs. 44 and 45 lead to the elongation of interstitial chains. Furthermore, they present possible paths for the growth of planar {311} defects. When two interstitials captured at seven-member ring sites interact during the random junps in the chain direction, they form a stable, immobile structure attached to an existing chain. This provides new trap nuclei for an additional interstitial chain parallel to the existng interstitial chain.

## 4. SUMMARY

In this article, we have discussed how the established process simulation toolset of "Technology Computer Aided Design" (TCAD), which was developed for the simulatior of microelectronic systems, needs to evolve in order to be capable of simulating nanoelectronic devices. Due to the nanometer dimensions of current and especially future devices, rew effects based on the atomic-level structure need to be included into the modeling, whch until recently consisted nearly exclusively of classical continuum models.

We have started by summarizing the traditional continuum diffusion reaction mocels (Sec. 2.1) and have described novel work relating the atomic hopping mechanisms to he macroscopic kinetic parameters (with and without applied stress field, Section 2.1.3) and calculation of those parameters with atomic simulations, especially *ab initio* techniqies (Section 2.3). We have argued that this atomic-level detail, when included into the continuum modeling, can lead to efficient models that are valid well into the nanoregime. Where thse

models finally fail, atomic-level process simulation becomes necessary (Section 2.2). As an introduction into this field, we have first reviewed transition state theory, which is the basis of most approaches to describe the kinetics of a system and becomes essential on the atomic level (Section 2.2.1). We then have discussed the currently predominant atomistic process simulation methods, which are the kinetic Monte Carlo (Section 2.2.2) and accelerated dynamics methods (Section 2.2.3) and have suggested the use of a novel technique by Henkelman and Jónsson [97], which is an event-catalog-free kinetic Monte Carlo approach (Section 2.2.4), as a possible technique to combine the best of both worlds for nanoscale process simulation.

For the calculation of kinetic parameters necessary for continuum and KMC methods, we have described the evolution from the (by now mostly outdated) drag method (Section 2.3.1) to the novel, much more accurate transition-state theory based methods, especially the nudged-elastic band method (Section 2.3.2) and the dimer method (Section 2.3.3).

Finally, we have demonstrated the application of these methods to the field of nanoscale process simulation and the impact that they can have. This started from calculation of kinetic parameters for continuum modeling (Section 3.1)—we discussed nitrogen (Section 3.1.1) and boron diffusion (Section 3.2.1) and capture radii in silicon (Section 3.1.2) as examples—and continued with the definition of a continuum model for boron nanoclustering (Section 3.2) and examples for the effect of stress on different diffusion mechanisms (Section 3.3). Finally, atomic-level process simulation has been demonstrated for the growth of nanoscale defect clusters after ion implantation, using the kinetic Monte Carlo approach to study nanovoid formation (Section 3.4.1) and accelerated dynamics for self-interstitial cluster growth (Section 3.4.2).

Although these examples are front-end processes from the fabrication of silicon nanoelectronic devices, the discussed concepts are quite general and are applicable to a wide variety of processes with focus on the ones involving temperature-activated diffusion and reaction events. In summary, it looks as if the path of process simulation towards a fully atomistic description thus has been paved, although the definition of a trustworthy standard methodology seems still to be a task of the future.

## ACKNOWLEDGMENTS

## REFERENCES

1. B. Maiti, P. J. Tobin, C. Hobbs, R. I. Hegde, F. Huang, D. L. O'Meara, D. Jovanovic, M. Mendicino, J. Chen, D. Connelly, O. Adetutu, J. Mogab, J. Candelaria, and L. B. La, *Proc. IEDM* 1998, 781 (1998).
2. The International Technology Roadmap for Semiconductors, 2003 ed. (Semiconductor Industry Association, San Jose, CA, 2003); http://public.itrs.net/.
3. G. Bourianoff, *Computer* 36, 44 (2003).
4. S. Thompson, N. Anand, M. Armstrong, C. Auth, B. Arcot, M. Alavi, P. Bai, J. Bielefeld, R. Bigwood, J. Brandenburg, M. Buehler, S. Cea, V. Chikarmane, C. Choi, R. Frankovic, T. Ghani, G. Glass, W. Han, T. Hoffmann, M. Hussein, P. Jacob, A. Jain, C. Jan, S. Joshi, C. Kenyon, J. Klaus, S. Klopcic, J. Luce, Z. Ma, B. McIntyre, K. Mistry, A. Murthy, P. Nguyen, H. Pearson, T. Sandford, R. Schweinfurth, R. Shaheed, S. Sivakumar, M. Taylor, B. Tufts, C. Wallace, P. Wang, C. Weber, and M. Bohr, in "IEDM Technical Digest" Institute of Electrical and Electronics Engineers, Piscataway, NJ, 2002, p. 61.
5. B. Doris, M. Ieong, T. Kanarsky, Y. Zhang, R. A. Roy, O. Dokumaci, Z. Ren, F.-F. Jamin, L. Shi, W. Natzle, H.-J. Huang, J. Mezzapelle, A. Mocuta, S. Womack, M. Gribelyuk, F. C. Jones, R. J. Miller, H.-S. P. Wong, and W. Haensch, in "IEDM Technical Digest," Institute of Electrical and Electronics Engineers, Piscataway, NJ, 2002, p. 267.
6. P. A. Packan, *MRS Bull.* 25, 18 (2000).
7. D. J. Frank, R. H. Dennard, E. Nowak, P. M. Solomon, Y. Taur, and H.-S. P. Wong, in "Proceedings of the IEEE 89," Institute of Electrical and Electronics Engineers, Piscataway, NJ, 2001, p. 259.
8. M. Lundstrom, *Science* 299, 210 (2003).
9. A. W. Ghosh, P. S. Damle, S. Datta, and A. Nitzan, *MRS Bull.* 29, 391 (2004).
10. M. Quirk and J. Serda, "Semiconductor Manufacturing Technology." Prentice Hall, New York, 2001, Fig. 9.1.
11. J. D. Plummer, M. D. Deal, and P. B. Griffin, "Silicon VLSI Technology." Prentice Hall, Upper Saddle River, NJ, 2000.

12. X. Y. Liu, W. Windl, and M. P. Masquelier, *Appl. Phys. Lett.* 77, 2018 (2000); 77, 4064 (2000).
13. For a recent review, see A. Nitzan and M. A. Ratner, *Science* 300, 1384 (2003).
14. A. W. Ghosh and S. Datta, *J. Comput. Electron.* 1, 515 (2002), and references therein.
15. C. R. Kagan, M. A. Ratner, Eds., *MRS Bull.* 29, 376 (2004).
16. V. Mujica, M. Kemp, and M. A. Ratner, *J. Chem. Phys.* 101, 6849 (1994).
17. S. Datta, "Electronic Transport in Mesoscopic Systems," Cambridge University Press, Cambridge, UK, 1995.
18. W. Tian, S. Datta, S. H. Hong, R. Reifenberger, J. I. Henderson, and C. P. Kubiak, *J. Chem. Phys.* 109, 2874 (1998).
19. K. Ravichandran, L. Fonseca, and W. Windl (to be published).
20. C. Krzeminski et al., in "Front-End Models for Silicon Future Technology, Public Final Report," FRENDTECH Consortium, IST Project 2000-30129 (2004), p. 47.
21. B. E. Deal and A. S. Grove, *J. Appl. Phys.* 36, 3770 (1965).
22. H. Z. Massoud, J. D. Plummer, and E. A. Irene, *J. Electrochem. Soc.* 132, 2685 (1985); 132, 2693 (1985).
23. C. J. Han and C. R. Helms, *J. Electrochem. Soc.* 134, 1297 (1987).
24. M. Yoshida, E. Arai, H. Nakamura, and Y. Terunuma, *J. Appl. Phys.* 45, 1498 (1974).
25. P. Pichler, "Intrinsic Point Defects, Impurities, and Their Diffusion in Silicon," Springer, Berlin, 2004, Chap. 1.5.
26. J. R. Manning, "Diffusion Kinetics for Atoms in Crystals," Princeton, Van Nostrand, 1968.
27. T. De Donder, *Bull. Classe Sci. Acad. R. Belg.* 24, 15 (1938).
28. J. N. Bronsted, *Z. Phys. Chem.* 102, 169 (1922).
29. J. N. Brönsted, *Z. Phys. Chem.* 115, 337 (1925).
30. R. Haase, *Z. Phys. Chem. Neue Folge* 128, 225 (1981).
31. R. Haase, *Z. Phys. Chem. Neue Folge* 153, 217 (1987).
32. M. Pinto, D. M. Boulin, C. S. Rafferty, R. K. Smith, W. M. Coughran, I. C. Kizilyalli, and M. J. Thoma, in "IEDM Technical Digest," IEEE, Piscataway 1992, p. 923.
33. D. W. Yergeau, E. C. Kan, M. J. Gander, and R. W. Dutton, in "6th International Conference on Simulation of Semiconductor Devices and Processes (H. Ryssel and P. Pichler, Eds.), SISDEP '95" p. 66, Springer, Wien, 1995.
34. TAURUS-MEDICI data sheet, Synopsis, Inc., 2003; http://www.synopsys.com/products/mixedsignal/taurus/device_sim_ds.html.
35. M. E. Law, Florida Object Oriented Process Simulator, University of Florida, 1991–2004; http://www.tec.ufl.edu/~flooxs/.
36. M. von Smoluchowski, *Z. Phys. Chem.* 92, 129 (1917).
37. T. R. Waite, *Phys. Rev.* 107, 463 (1957).
38. P. Debye, *Trans. Electrochem. Soc.* 82, 265 (1942).
39. T. R. Waite, *J. Chem. Phys.* 28, 103 (1958).
40. W. Nernst, *Z. Phys. Chem.* 2, 613 (1888).
41. A. Einstein, *Ann. Phys.* 17, 549 (1905).
42. C. T. Sah, P. C. H. Chan, C.-K.Wang, R. L.-Y. Sah, K. A. Yamakawa, and R. Lutwack, *IEEE Trans. Electron Devices* 28, 304 (1981).
43. D. Peak and J. W. Corbett, *Phys. Rev. B* 5, 1226 (1972).
44. M. S. Daw, W. Windl, N. N. Carlson, M. Laudon, and M. P. Masquelier, *Phys. Rev. B* 64, 045205 (2001).
45. P. Ramanarayanan, B. Srinivasan, K. J. Cho, and B. M. Clemens, *J. Appl. Phys.* 96, 7095 (2004).
46. F. C. Larche and J. W. Cahn, *Acta Metall.* 30, 1835 (1982).
47. D. Mathiot and J. C. Pfister, *J. Appl. Phys.* 66, 970 (1989).
48. D. Stauffer, *Phys. Rep.* 54, 1 (1979).
49. L. Pelaz, G. H. Gilmer, H.-J. Gossmann, C. S. Rafferty, M. Jaraiz, and J. Barbolla, *Appl. Phys. Lett.* 74, 3657 (1999).
50. N. Strecker, V. Moroz, and M. Jaraiz, "Proceedings of the 2002 International Conference on Computational Nanoscience" (2002), p. 247.
51. M. Jaraiz, P. Castrillo, R. Pinacho, I. Martin-Bragado and J. Barbolla, "SISPAD 2001 Proceedings," Springer, Berlin, 2001, p. 10.
52. A. Asenov, "SISPAD 2001 Proceedings," Springer, Berlin, 2001, p. 162.
53. H. Eyring, *J. Chem. Phys.* 3, 107 (1935).
54. E. Wigner, *Trans. Faraday Soc.* 34, 29 (1938).
55. J. C. Keck, *Adv. Chem. Phys.* 13, 85 (1967).
56. P. Pechukas, in "Dynamics of Molecular Collisions" (W. Miller Ed.), Plenum, New York, 1976, Part B.
57. D. Chandler, *J. Chem. Phys.* 68, 2959 (1978).
58. A. F. Voter and D. Doll, *J. Chem. Phys.* 80, 5832 (1984).
59. A. F. Voter and D. Doll, *J. Chem. Phys.* 82, 80 (1985).
60. A. F. Voter, F. Montalenti, and T. C. Germann, *Annu. Rev. Mater. Res.* 32, 321 (2002).
61. C. Wert and C. Zener, *Phys. Rev.* 76, 1169 (1949).
62. G. H. Vineyard, *J. Phys. Chem. Solids* 3, 145 (1957).
63. P. J. Feibelman, *Phys. Rev. Lett.* 65, 729 (1990).
64. A. B. Bortz, M. H. Kalos, and J. L. Lebowitz, *J. Comput. Phys.* 17, 10 (1975).
65. D. T. Gillespie, *J. Comput. Phys.* 22, 403 (1976).
66. D. T. Gillespie, *J. Phys.* 81, 2340 (1977).

67. D. T. Gillespie, *J. Comput. Phys.* 28, 395 (1978).

68. G. H. Gilmer, *Science* 208, 335 (1980).

69. A. F. Voter, *Phys. Rev. B* 34, 6819 (1986).

70. M. Jaraiz, in "The Encyclopedia on Materials Science and Technology" (K. H. J. Buschow, R. W. Cahn, M. C. Flemings, B. Ilschner, E. J. Kramer, and S. Mahajan, Eds.), p. 7861, Pergamon, Boston, 2001.

71. D. J. Eaglesham, P. A. Stolk, H. J. Gossmann, and J. M. Poate, *Appl. Phys. Lett.* 65, 2305 (1994).

72. J. N. Kim, J. W. Wilkins, F. S. Khan, and A. Canning, *Phys. Rev. B* 55, 16186 (1997).

73. M. E. Law, G. H. Gilmer, and M. Jaraiz, *Mater. Res. Soc. Bull.* 25, 45 (2000).

74. N. E. B. Cowern, G. Mannino, P. A. Stolk, F. Roozemboom, H. G. A. Huizing, J. G. M. van Berkum, F. Cristiano, A. Claverie, and M. Jaraiz, *Phys. Rev. Lett.* 82, 4460 (1999).

75. J. Kim, F. Kirchhoff, J. W. Wilkins, and F. S. Khan, *Phys. Rev. Lett.* 84, 503 (2000).

76. D. A. Richie, J. Kim, S. A. Barr, K. R. A. Hazzard, R. Hennig, and J. W. Wilkins, *Phys. Rev. Lett.* 92, 045501 (2004).

77. L. Pelaz, M. Jaraiz, G. H. Gilmer, H. J. Gossmann, C. S. Rafferty, D. J. Eaglesham, and J. M. Poate, *Appl. Phys. Lett.* 70, 2285 (1997).

78. G. L. Kellogg and P. J. Feibelman, *Phys. Rev. Lett.* 64, 3143 (1990).

79. C. Chen and T. T. Tsong, *Phys. Rev. Lett.* 64, 3147 (1990).

80. P. Stoltz and J. K. Norskov, *Phys. Rev. B* 48, 5607 (1993).

81. M. Villarba and H. Jónsson, *Phys. Rev. B* 49, 2208 (1994).

82. R. Wang and K. A. Fichthorn, *Surf. Sci.* 301, 53 (1994).

83. J. C. Hamilton, M. S. Daw, and S. M. Foiles, *Phys. Rev. Lett.* 74, 2760 (1995).

84. Z.-P. Shi, Z. Zhang, A. K. Swan, and J. F. Wendelken, *Phys. Rev. Lett.* 76, 4927 (1996).

85. F. Montalenti and F. Ferrando, *Phys. Rev. Lett.* 82, 1498 (1999).

86. T. R. Linderoth, S. Horch, L. Petersen, S. Helveg, E. Laegsgaard, I. Stengaard, and F. Besenbacher, *Phys. Rev. Lett.* 82, 1494 (1999).

87. O. S. Trushin and A. Ala Nissila, *Phys. Rev. B* 62, 1611 (2000).

88. G. Henkelman and H. Jónsson, *J. Chem. Phys.* 115, 9657 (2001).

89. A. F. Voter, *Phys. Rev. B* 57, 985 (1998).

90. M. Shirts and V. S. Pande, *Science* 290, 1903 (2000).

91. E. K. Grimmelman, J. C. Tully, and E. Helfand, *J. Chem. Phys.* 74, 5300 (1981).

92. A. F. Voter, *J. Chem. Phys.* 106, 4665 (1997).

93. A. F. Voter, *Phys. Rev. Lett.* 78, 3908 (1997).

94. M. M. Steiner, P. A. Genilloud, and J. W. Wilkins, *Phys. Rev. B* 57, 10236 (1998).

95. S. Pal and K. A. Fichthorn, *Chem. Eng. J.* 74, 77 (1999).

96. M. R. Sorensen and A. F. Voter, *J. Chem. Phys.* 112, 9599 (2000).

97. G. Henkelman, Ph.D. thesis, University of Washington, 2001.

98. G. Henkelman and H. Jonsson, *Mater. Res. Soc. Symp. Proc.* 677, AA8.1 (2001).

99. F. Gao, G. Henkelman, W. J. Weber, L. R. Corrales, and H. Jónsson, *Nucl. Instrum. Meth. B* 202, 1 (2003).

100. R. Marcus, *J. Chem. Phys.* 45, 4493 (1966).

101. M. McKee and M. Page, in "Reviews in Computational Chemistry" (K. B. Lipkowitz and D. B. Boyd, Eds.), Vol. 4, VCH, New York, 1993.

102. H. Jónsson, G. Mills, and K. W. Jacobsen, in "Classical and Quantum Dynamics in Condensed Phase Simulations" (B. J. Berne, G. Ciccotti, and D. F. Coker, Eds.), p. 385, World Scientific, Singapore, 1998.

103. H. C. Andersen, *J. Chem. Phys.* 72, 2384 (1980).

104. C. S. Nichols, C. G. Van de Walle, and S. T. Pantelides, *Phys. Rev. B* 40, 5484 (1989).

105. J. Zhu, T. D. de la Rubia, L. H. Yang, C. Mailhiot, and G. H. Gilmer, *Phys. Rev. B* 54, 4741 (1996); J. Zhu, *Comput. Mater. Sci.* 12, 309 (1996).

106. W. Windl, M. M. Bunea, R. Stumpf, S. T. Dunham, and M. P. Masquelier, *Phys. Rev. Lett.* 83, 4345 (1999).

107. G. Mills and H. Jónsson, *Phys. Rev. Lett.* 72, 1124 (1994).

108. G. Mills, H. Jónsson, and G. K. Schenter, *Surf. Sci.* 324, 305 (1995).

109. G. Henkelman and H. Jónsson, *J. Chem. Phys.* 113, 9978 (2000).

110. N. Mousseau and G. T. Barkema, *Phys. Rev. E* 57, 2419 (1998).

111. G. Kresse and J. Hafner, *Phys. Rev. B* 47, 558 (1993); 49, 14251 (1994); G. Kresse and J. Furthmüller, *Comput. Mater. Sci.* 6, 15 (1996); *Phys. Rev. B* 54, 11169 (1996). G. Kresse and J. Hafner, *J. Phys.* 6, 8245 (1994).

112. N. G. Stoddard, P. Pichler, G. Duscher, and W. Windl, *Phys. Rev. Lett.* 95, 025901 (2005).

113. C. L. Liu, W. Windl, L. Borucki, S. F. Lu, and X. Y. Liu, *Appl. Phys. Lett.* 80, 52 (2002).

114. X. Y. Liu, W. Windl, K. M. Beardmore, and M. P. Masquelier, *Appl. Phys. Lett.* 82, 1839 (2003).

115. B. P. Uberuaga, M. Levskovar, A. P. Smith, H. Jónsson, and M. Olmstead, *Phys. Rev. Lett.* 84, 2441 (2000).

116. N. P. Kopsias and D. N. Theodorou, *J. Chem. Phys.* 109, 8573 (1998).

117. G. Henkelman and H. Jonsson, *J. Chem. Phys.* 111, 7010 (1999).

118. G. Henkelman, G. Jóhannesson, and H. Jónsson, in "Progress on Theoretical Chemistry and Physics" (S. Schwartz, Ed.), Kluwer Academic, New York, 2000.

119. G. T. Barkema and N. Mousseau, *Phys. Rev. Lett.* 77, 4358 (1996).

120. T. Abe, H. Harada, and J.-I. Chikawa, in "Defects in Semiconductors II" (S. Mahajan and J. W. Corbett, Eds.), p. 1, MRS Symposia Proceedings 14, Materials Research Society, Pittsburgh, 1983.

121. K. Sumino, I. Yonenaga, M. Imai, and T. Abe, *J. Appl. Phys.* 54, 5016 (1983).

122. W. J. M. J. Josquin and Y. Tamminga, *J. Electrochem. Soc.* 129, 1803 (1982).

123. A. Karoui, F. S. Karoui, G. Rozgonyi, M. Hourai, and K. Sueoka, J. Electrochem. Soc. 150, G771 (2003).
124. H. Kageshima, A. Taguchi, and K. Wada, Appl. Phys. Lett. 76, 3718 (2000).
125. H. J. Stein, in "Thirteenth International Conference on Defects in Semiconductors" (L. C. Kimerlirg and J. M. Parsey, Jr., Eds.), p. 839. Metallurgical Society of AIME, 1985.
126. R. Jones, S. Öberg, F. Berg Rasmussen, and B. Bech Nielsen, Phys. Rev. Lett. 72, 1882 (1994).
127. A. Gali, J. Miro, P. Deák, C. P. Ewels, and R. Jones, J. Phys. 8, 7711 (1996).
128. H. Sawada and K. Kawakami, Phys. Rev. B 62, 1851 (2000).
129. A. Karoui, F. Sahtout Karoui, G. A. Rozgonyi, M. Hourai, and K. Sueoka, in "Semiconductor Silicon 2002" (H. R. Huff, L. Fabry, and S. Kishino, Eds.), p. 670. Electrochemical Society Proceedings 2002-2 Electrochemical Society, 2002.
130. T. Itoh and T. Abe, Appl. Phys. Lett. 53, 39 (1988).
131. A. Hara, T. Fukuda, T. Miyabo, and I. Hirai, Appl. Phys. Lett. 54, 626 (1989).
132. G. J. Willems and H. E. Maes, J. Appl. Phys. 73, 3256 (1993).
133. N. Fuma, K. Tashiro, K. Kakumoto, and Y. Takano, Mater. Sci. Forum 196–201, 797 (1995).
134. N. Fuma, K. Tashiro, K. Kakumoto, and Y. Takano, Jpn. J. Appl. Phys. Part 1 35, 1993 (1996).
135. L. S. Adam, M. E. Law, K. S. Jones, O. Dokumaci, C. S. Murthy, and S. Hegde, J. Appl. Phys. 87, 2282 (2000).
136. R. S. Hockett, Appl. Phys. Lett. 54, 1793 (1989).
137. G. Mannino, V. Privitera, S. Scalese, S. Libertino, E. Napolitani, P. Pichler, and N. E. B. Cowern, Eletrrchem. and Solid State Lett. 7, G161 (2004).
138. L. S. Adam, M. E. Law, O. Dokumaci, and S. Hegde, J. Appl. Phys. 91, 1894 (2002).
139. V. Voronkov and R. Falster, Solid State Phenom. 95–96, 83 (2004).
140. P. A. Schultz and J. S. Nelson, Appl. Phys. Lett. 78, 736 (2001).
141. H. Sawada, K. Kawakami, A. Ikari, and W. Ohashi, Phys. Rev. B 65, 075201 (2002).
142. G. Henkelman and H. Jónsson, J. Chem. Phys. 113, 9901 (2000).
143. D. Chandler, in "Introduction to Modern Statistical Mechanics." Oxford University Press, New York, 1987, p. 70.
144. R. LeSar, R. Najafabadi, and D. J. Srolovitz, Phys. Rev. Lett. 63, 624 (1989).
145. W. Windl, Phys. Status Solid (B) 241, 2313 (2004).
146. M. M. Bunea, Ph.D. thesis, Boston University, 2000.
147. K. M. Beardmore, W. Windl, B. P. Haley, and N. Gronbech-Jensen, in "Computational Nanoscence and Nanotechnology," Applied Computational Research Society, Cambridge, 2002, p. 251.
148. M. J. Puska, S. Pöykkö, M. Pesola, and R. M. Nieminen, Phys. Rev. B 58, 1318 (1998).
149. M. M. Bunea, P. Fastenko, and S. T. Dunham, Mater. Res. Soc. Proc. 568, 135 (1999).
150. S. M. Hu, Phys. Status Solidi B 60, 595 (1973).
151. S. T. Dunham and C. D. Wu, J. Appl. Phys. 78, 2362 (1995).
152. W. Windl (unpublished).
153. O. Pankratov, H. Huang, T. D. de la Rubia, and C. Mailhiot, Phys. Rev. B 56, 13172 (1997).
154. W. Windl, Ph.D. thesis, University of Regensburg, CH Publishing, Regensburg, 1995, p. 89.
155. W. Windl, IEICE Trans. Electron. E86C, 269 (2003).
156. P. A. Stolk, H.-J. Gossmann, D. J. Eaglesham, D. C. Jacobson, and J. M. Poate, Appl. Phys. Ler. 66, 568 (1995); L. H. Zhang, K. S. Jones, P. H. Chi, and D. S. Simons, Appl. Phys. Lett. 67, 2025 (1995).
157. A. Ural, P. B. Griffin, and J. D. Plummer, J. Appl. Phys. 85, 6440 (1999).
158. M. A. Meléndez-Lira, J. D. Lorentzen, J. Menéndez, W. Windl, N. G. Cave, R. Liu, J. W. Christiansen, N. D. Theodore, and J. J. Candelaria, Phys. Rev. B 56, 3648 (1997).
159. W. Windl, A. A. Demkov, and O. F. Sankey, "Theory of Strain and Electronic Structure of Si₁C₁ and Si₁₋ₓ₋ᵧGeₓC₁ Alloys, in Silicon-Germanium Carbon Alloys: Growth, Properties and Applications (Optoelectronic Properties of Semiconductors and Superlattices)" (S. T. Pantelides and S. Zollner, Eds.), Chap 8. Taylor and Francis, New York, 2002.
160. J. Yamauchi, N. Aoki, and I. Mizushima, Phys. Rev. B 55, 10245 (1997).
161. B. P. Uberuaga, G. Henkelman, H. Jónsson, S. T. Dunham, W. Windl, and R. Stumpf, Phys. Status Solidi 233, 244 (2002).
162. D. Yergeau, E. C. Kan, M. J. Gander, and R. W. Dutton, in "Proceedings of the 6th International Conference on Simulation of Semiconductor Devices and Processes (SISDEP'95), Erlangen" (H. Ryssel and P. Pichler, Eds.), p. 151. Springer, Wien, 1995.
163. W. C. Lee, S. G. Lee, and K. J. Chang, J. Phys. 10, 995 (1998).
164. M. Tang, L. Colombo, J. Zhu, and T. Diaz de la Rubia, Phys. Rev. B 55, 14279 (1997).
165. X. Y. Liu and W. Windl, J. Comp. Elec. (in print).
166. C. S. Rafferty, G. H. Gilmer, M. Jaraiz, D. Eaglesham, and H.-J. Gossmann, Appl. Phys. Lett. 68, 2395 (1996).
167. A. Mokhberi, P. B. Griffin, and J. D. Plummer, E. Paton, S. McCoy, and K. Elliott, IEEE T. Electron. Dev. 49, 1183 (2002).
168. N. Collaert and P. Verheyen, Strained Si and SiGe devices, IMEC-ASD, http://www.imec.be/wwwinter/processing asd activities/strained.shtml.
169. N. Collaert, P. Verheyen, K. De Meyer, R. Loo, and M. Caymax, Solid-State Electron. 47, 1173 (2003).
170. P. Kuo, Appl. Phys. Lett. 66, 580 (1995).
171. Y. Zaitsu, K. Osada, T. Shimizu, S. Matsumoto, M. Yoshida, F. Arai, and T. Abe, Mater. Sci. Forum 196–201, 1891 (1995).
172. H. Park, K. S. Jones, J. A. Slinkman, and M. E. Law, J. Appl. Phys. 78, 3664 (1995).

*173.* M. Aziz. *Defect Diffusion Forum* 153–155, 1 (1998).

*174.* Y. Zhao, M. Aziz, H.-J. Gossmann, S. Mitha, and D. Schiferl. *Appl. Phys. Lett.* 75, 941 (1999).

*175.* Y. Zhao, M. Aziz, H.-J. Gossmann, S. Mitha, and D. Schiferl. *Appl. Phys. Lett.* 74, 31 (1999).

*176.* S. Chaudhry and M. E. Law, *J. Appl. Phys.* 82, 1138 (1997).

*177.* K. Osada Y. Zaitsu, S. Matsumoto, M. Yoshida, E. Arai, and T. Abe, *J. Electrochem. Soc.* 142, 202 (1995).

*178.* Y. Todokoro and I. Teramoto, *J. Appl. Phys.* 49, 3527 (1978).

*179.* E. Suhir, *J. Appl. Mechanics* 55, 143 (1988); D. B. Bogy, *J. Appl. Mechanics* 35, 460 (1968).

*180.* B. Gu and P. E. Phelan, *Appl. Supercond.* 6, 19 (1998).

*181.* E. Tarnow, *J. Phys.* 5, 1863 (1993).

*182.* A. Antonelli, E. Kaxiras, and D. J. Chadi, *Phys. Rev. Lett.* 81, 2088 (1998).

*183.* W. Windl, M. M. Bunea, R Stumpf, S. T. Dunham, and M. P. Masquelier, *Phys. Rev. Lett.* 83, 4345 (1999).

*184.* M. Laudon, N. N. Carlson, M. P. Masquelier, M. S. Daw, and W. Windl, *Appl. Phys. Lett.* 78, 201 (2001).

*185.* B. P. Uberuaga, R. Stumpf, W. Windl, H. Jonsson, and M. P. Masquelier (unpublished); this work also finds that stress effects on prefactors are small compared to exponential effects, whence we neglect the former.

*186.* A. Antonelli and J. Bernholc, *Phys. Rev. B* 40, 10643 (1989).

*187.* B. Sadigh, T. J. Lenosky, S. K. Theiss, M.-J. Caturla, T. Diaz de la Rubia, and M. A. Foad, *Phys. Rev. Lett.* 83, 4341 (1999).

*188.* Y. M. Haddara, B. T. Folmer, M. E. Law, and T. Buyuklimanli, *Appl. Phys. Lett.* 77, 1976 (2000).

*189.* N. N. Carlson and K. Miller, *SIAM J. Sci. Computing* 19, 728 (1998); 19, 766 (1998).

*190.* M. S. Daw and W. Windl, *Stress dependence of multi-stream dopant diffusion* (unpublished).

*191.* National Technology Roadmap for Semiconductors (NTRS) 1997. Semiconductor Industry Association (1997), p. 188.

*192.* H. S. Wong and Y. Taur, *Proc. Int. Electron Dev. Meeting Proc.* 1993, 705 (1993).

*193.* E. Dornberger, J. Virbulis, B. Hanna, R. Hoelzl, E. Daub, and W. Von Ammon, *J. Cryst. Growth* 229, 11 (2001).

*194.* M. Itsumi, H. Akiya, T. Ueki, M. Tomita, and M. Yamawaki, *Jpn. J. Appl. Phys. Part 1* 35, 812 (1996).

*195.* T. Sinno and R. A. Brown, *J. Electrochem. Soc.* 146, 2300 (1999).

*196.* M. Prasad and T. Sinno, *Appl. Phys. Lett.* 80, 1951 (2002).

*197.* M. Jaraiz, L. Pelaz, E. Rubio, J. Barbolla, G. H. Gilmer, D. J. Eaglesham, H. J. Gossmann, and J. M. Poate, *Mater. Res. Soc. Symp. Proc.* (N. E. B. Cowern, D. C. Jacobson, P. B. Griffin, P. A. Packan, and R. P. Webb, Eds.), p. 43. Materials Research Society, Warrendale, Pennsylvania, 1998.

*198.* O. W. Holland and C. W. White, *Nucl. Instrum. Meth. B* 59/60, 353 (1991).

*199.* M. Jaraiz, E. Rubio, P. Castrillo, L. Pelaz, L. Bailon, J. Barbolla, G. H. Gilmer, and C. S. Rafferty, *Mater. Sci. Semicond. Proc.* 3, 59 (2000).

*200.* M. Jaraiz, G. H. Gilmer, J. M. Poate, and T. D. de la Rubia, *Appl. Phys. Lett.* 68, 409 (1996).

*201.* N. Arai, S. Takeda, and M. Kohyama, *Phys. Rev. Lett.* 78, 4265 (1997).

*202.* J. Kim, F. Kirchhoff, W. G. Aulbur, J. W. Wilkins, and F. S. Khan, *Phys. Rev. Lett.* 83, 1990 (1999).

*203.* B. J. Coomer, J. P. Gross, R. Jones, S. Oberg, and P. R. Briddon, *J. Phys.* 13, L1 (2001).

*204.* A. Bongiorno, L. Colombo, F. Cargnoni, C. Gatti, and M. Rosati, *Euro. Phys. Lett.* 50, 608 (2000).

*205.* S. Birner, J. Kim, D. A. Richie, J. W. Wilkins, A. F. Voter, and T. J. Lenosky, *Solid State Commun.* 120, 279 (2001).

*206.* T. J. Lenosky, B. Sadigh, E. Alonso, V. V. Bulatov, T. Diaz de ka Rubia, J. Kim, A. F. Voter, and J. D. Kress, *Model. Simul. Mater. Sci. Eng.* 8, 825 (2000).

*207.* S. Takeda, *Jpn. J. Appl. Phys.* 30, L639 (1991).

*208.* M. Kohyama and S. Takeda, *Phys. Rev. B* 46, 12305 (1992).

# CHAPTER 3

# Electron Transport in Nanostructured Systems— *Ab Initio* Study

**Yoshiyuki Kawazoe, Hiroshi Mizuseki,**
**Rodion Belosludov, Amir Farajian**

*Institute for Materials Research, Tohoku University, Sendai, Japan*

## CONTENTS

# 1. INTRODUCTION

Recently, nanotechnology (clusters, molecules, and supramolecules on a scale of 1–100 nm) has attracted considerable attention from scientists, as it certainly will play an important role in a broad range of technological fields such as semiconductor manufacturing, materials science, chemistry, biology, genetics, and medicine [1]. Two basic approaches for creating nanostructures are the top-down approach and the bottom-up approach. The top-down approach involves molding or etching materials into smaller components. This approach has traditionally been used in making parts for computers and electronics. The bottom-up approach involves assembling structures atom-by-atom or molecule-by-molecule.

In parallel with the progress of more effective fabrication technologies, theoretical study of promising molecular and cluster structures based on quantum mechanical calculations is also a key factor in the design of new nanomaterials with desirable physical and chemical characteristics. When scientists perform research, they should carefully consider what benefits can be achieved from their investigations. To apply the nano to different technologies, experimentalists should make much effort to bring new ideas and produce nanocomposites with desired properties. Theoreticians can aid in this search by narrowing the field through accurate prediction of the chemical and physical properties of different nanomaterials. Thus, the nano is a field where theory can play a significant role hand-in-hand with experiment. Computational materials science is rapidly becoming an essential tool for investigation of a variety of physical and chemical properties of nanomaterials. The value of simulation can be evaluated based on its ability to rapidly and accurately predict the properties of novel functional nanomaterials in a more cost-effective way than is experimentally possible. Instead of synthesizing and testing a large number of potential candidates, it is possible to use the combinatorial computational chemistry approach to screen a large number of candidates, including very expensive elements. Experiments only need to be performed on a small number of the most promising candidates [2]. Using powerful computers and highly accurate methods, we can accelerate the realization of novel nanomaterials and propose these materials for different applications.

The basic principle of first-principles simulations is accurately determining the total energy of an investigated system. There is a cost scale to computational materials science because so many physical properties are related with the total energies. While just one piece of the theoretical tool is necessary to calculate all the physical properties that are related to the total energies, completely different pieces of experimental tools are required to measure each class of physical properties of a material [3]. This represents an enormous advantage of computational methods over experimental measurements. Simulations are easy to perform, even for very complex systems; often their complexity is no worse than the complexity of the physical description. As the capabilities of computers expand, simulations of many-body systems will be able to treat more complex physical systems at higher levels of accuracy. Thus, the ultimate result of an extremely wide range of scientific and engineering applications will undoubtedly be profound.

Despite a remarkable miniaturization trend in the semiconductor industry, in the next 10–15 years, conventional Si-based microelectronics likely faces fundamental limitations when feature lengths shrink below 10 nm. Thus, molecular electronics has attracted considerable attention from scientists and the semiconductor industry as a "postsilicon technology" for future applications and trends in advanced computer electronics [4]. The main challenge in molecular electronics is to establish that single molecules or a finite number of self-assembled molecules can perform all the basic functions of conventional electronics components such as wires, diodes, and transistors [5]. There have been significant advances in the fabrication and demonstration of molecular electronic wires [6, 7], molecular diodes, and two-terminal electrical switches made from single molecules [8, 9]. In parallel with the progress of more effective fabrication technologies, theoretical studies are also an important part of the molecular electronics because they can directly propose and design novel nanodevices as well as enrich the experimental intuition. For example, the molecular electronics dates to the original theoretical work in which Aviram and Ratner have revealed the possibility of an organic molecule functioning as a molecular rectifier [10].

Our group developed the TARABORD code enables to simulate the transport properties based on first-principles calculations [11]. Despite the fact that the first-principles methods are either computationally expensive, they are important because they can be used to propose novel nanodevices as well as enrich experimental intuition. In this program, the nonequilibrium Green's function (NEGF) approach for first-principles modeling of current-voltage characteristics of molecular electronics devices has been used. The molecular device is modeled on atomic level and the electronic structures of the studied systems will be described using the density functional theory (DFT). This approach includes the following steps:

1. Dividing the systems into electrode and scattering region.
2. Determining the one electron DFT Hamiltonian and overlap matrices.
3. Setting up the NEGF, determing the charge density.
4. Calculating the effective potential.

Hamiltonian and the overlap matrices corresponding to the gold contacts, the surface Green's functions describing the semi-infinite electrodes attached to the molecules from the left and right sides are derived. These surface Green's functions together with the Hamiltonians and overlap matrices of the molecule, and the molecule-electrode part, are then used to determine the conductance of the system using the Landauer approach [12]. The location of the Fermi level related to the molecular level is directly calculated, taking into consideration the charge transport between molecule and metal contacts, to obtain an accurate value of the conductance.

# 2. THEORETICAL INVESTIGATION OF TRANSPORT IN NANOMETER-SCALE SYSTEMS

In this section, we discuss a theoretical model based on the nonequilibrium Green's function approach in order to calculate transport properties of nanometer-scale systems. The method is presented for the particular case of contact-molecule-contact (i.e., double-terminal systems). However, our approach is general and can be applied equally well to three- and multiterminal cases. Using this approach, one is able to characterize and predict the transport properties of basic nano- and molecular electronic components.

The main functional parts in nano- and molecular electronics applications are a few nanometers across. Two examples are depicted in Fig. 1. Figure 1(a) shows a polythiophene molecule and parts of the gold nanocontacts to which the molecule is attached. The nanocontacts are semi-infinite in extent, that is, are bounded only by the molecule sandwiched between them. In Fig. 1(b), part of an infinitely long carbon nanotube is shown. The nanotube is deposited on a double-crystal substrate, hence its "left" and "right" parts are subject to different doping effects. Applying bias will cause current to flow through these systems.

Making use of Landauer's formalism [13], it is rather straightforward to calculate the quantum conductance of an open contact-molecule-contact system, in which the contacts are semi-infinite. If we consider semi-infinite quasi–one-dimensional (Q1D) nanocontacts, then

**(a)**



**(b)**



Figure 1. Two typical examples of nano- and molecular electronic components. (a) A polythiophene molecule attached to two gold nanocontacts via sulfur atoms. (b) A carbon nanotube deposited on a double-crystal substrate.

the band structures of the nanocontacts will also play an important role in calculating conduction. This issue has practical importance as with the advancement of nanolithographic techniques, contacts whose cross sections lie in the nanometer-scale are experimentally realized nowadays. When two such nanocontacts sandwich a molecule as in Fig. 1(a), the conduction of the whole system can be calculated by deriving the probability for an excitation well inside the "left" nanocontact to travel first to the left contact's "surface," then across the middle junction onto the "surface" of the "right" contact, and finally well within the right contact. The problem of calculating the transition probability is conveniently solved by calculating the "surface Green's functions" of the left and right contacts, and by "attaching" them to the Green's function of the middle junction [14].

For the system of Fig. 1(b), for example, the left and right contacts are the left and right parts of the nanotube located "far enough" from the interface of the substrate crystals, and the middle junction part is the middle part of the nanotube that joins the left side to the right side. By mentioning "far enough" we mean that the effects of the disturbance caused by the interface of the substrate crystals are negligible for the left and right parts of the nanotube. These effects, however, are not negligible for the middle junction part. In other words, the middle part is long enough to effectively screen the disturbance [15].

In the sections to follow, we first briefly explain our general approach as implemented in the computer program TARABORD [11] and then present typical applications of the method. The method presented here is independent of the particular approach used for calculating the electronic structure of the system. Although the electronic structure data are the necessary prerequisite input of our method, the method is capable of using any level of accuracy. That is, for example, *ab initio* and tight-binding descriptions can both be used. This will be clarified in the following sections.

## 2.1. Calculating Transport Using Nonequilibrium Green's Function Approach

The nonequilibrium surface Green's function matching (NSGFM) employed here for transport calculations is a generalization of the surface Green's function matching approach [14, 16] that was previously applied to some nanotube junctions under equilibrium conditions, that is, without external bias [17]. Both the equilibrium and nonequilibrium surface Green's function matching formalisms are independent of the particular description of the electronic structure of the system. In particular, they can be applied to *ab initio* or semiempirical model Hamiltonians. Using tight-binding description, we have previously applied this method to various carbon nanotube systems [15, 18–21]. Generalization to *ab initio* description with nonorthogonal basis makes it possible to calculate transport characteristics of any desired contact-molecule-contact system with high accuracy [11]. For example, we use *ab initio* and tight-binding modelings to calculate transport characteristics of a polythiophene-based molecular device, as well as some nanotube-based devices.

The NSGFM method is schematically shown in Fig. 2. It is seen that an excitation far inside the left contact travels to its surface after passing through successive layers. After tunneling across the middle junction, the excitation travels from the surface of the right contact to layers far inside. We make use of the NSGFM method in order to calculate the conductance and I-V characteristics of an open system that consists of a general finite system (e.g., the functional molecule) attached on its left and right hand to two semi-infinite contacts or electrodes. The method applies to the general case of nonorthogonal basis. In Fig. 2, the left and right contacts are divided into successive layers. Each layer interacts only with its nearest-neighboring layers and includes the minimum possible number of unit cells of the electrode. Moreover, the middle molecular junction is assumed to have direct interactions only with the surface layers of the left and right contacts.

As is clear from Fig. 2, the general system that we would like to consider for transport calculations is open. This means that the contacts, that are responsible for applying bias to the functional middle part and to let the current pass through the system, are semi-infinite. The mathematically open system corresponds to the physical case of macroscopic contacts attached to a nanometer-scale functional part.

Figure 2. Schematic presentation of the nonequilibrium surface Green's function matching approach.

The transport calculation starts with obtaining the Hamiltonian and overlap matrices of the system in any localized basis set. As mentioned before, the exact method and level of accuracy of electronic structure calculation used to provide the Hamiltonian and overlap matrices is independent of the transport calculation. Using these matrices, we first calculate the transfer matrices corresponding to the left and right electrodes [22, 23]. Next, the surface Green's functions of the left and right electrodes are calculated [11, 14, 16] that effectively close the open contact-molecule-contact system, as depicted in Fig. 3.

The surface Green's functions, together with the Hamiltonian and overlap matrices coupling the middle molecular junction to the left and right electrodes, are then used to calculate the self-energies $\Sigma_l$ and $\Sigma_R$ of the left and right contacts. The conductance of the system, $\Gamma$, is then obtained through [11, 24]

$$\Gamma(E, V) = Tr[\Gamma_R G \Gamma_L G^\dagger],$$

where $\Gamma_{L,R} = \Sigma_{L,R} - \Sigma_{L,R}^\dagger$ and $G$ is the total Green's function of the system projected to the molecular junction space [11, 14, 16]. The conductance $\Gamma(E, V)$ indeed depends on both the (carrier) energy $E$ and the bias voltage $V$, as all the matrices used in the definition of conductance depend on $E$ and $V$.

The current $I(V)$ is then obtained by integrating $\Gamma(E, V)$ taking into account the Fermi-Dirac distributions of the left and right contacts [25].

In the following sections, we will apply the implementation of the method in the computer program TARABORD [11] to some particular nanosystems, namely, doped nanotube junctions as well as nanoelectromechanical switches and sensors based on bent nanotubes. These applications are based on semiempirical tight-binding modeling of the systems. We will also show typical applications using *ab initio* descriptions of the system.



Figure 3. The closed system corresponding to the open system of Fig. 2, obtained via calculating the surface Green's functions of the electrodes.

## 2.2. Negative Differential Resistance and Rectifying Behavior of Doped Nanotube Junctions

Carbon nanotubes are probably the most-studied materials for various nanotechnolog' applications, including nanoelectronics. The possibility of doping carbon nanotubes with akali or halogen atoms has been the subject of experiments [26, 27] and theoretical works [28, 29]. Doping carbon nanotubes with donor and acceptor atoms, as depicted in Fig. 4, is one way to produce a "nanotube junction," that is, a nanotube in which the valence bands of he left and right sides are shifted with respect to each other [18, 30].

In Fig. 4, charge transfer to/from the nanotube is responsible for this shift. Another possibility is the one shown in Fig. 1(b), where the two substrates on which the nanotube is deposited are assumed to have different work functions. This will cause the charge transfer to/from the left and right sides of the nanotube to be different and a shift of the valence bands to be present.

We use a simple one orbital per atom tight-binding Hamiltonian to model the nanotube [18]. The on-site energies are determined by the external bias and the initial shifts of the valence bands, and the hopping terms are a nonzero constant (taken as unit of energy) only for nearest-neighboring sites. As an example to illustrate the effect of doping, in the present work, initial shifts 0.2 and 0.1 and 0.3 and 0.0 are assumed in the calculations for the left and right parts of the metallic and semiconducting tubes, respectively. These values determine the initial position of chemical potentials with respect to the density of states when here is no external bias. The external bias will then determine the relative shifts of the new chemical potentials determined by the dopings [18].

The $I(V)$ characteristics at temperature $T = 0$ of (4,4) and (3,3) armchair tubes, as well as those of (7,0), (5,0), and (3,0) zigzag tubes are depicted in Figs. 5(a) and 5(b) These curves are obtained by assuming an external step potential at the junction, without any self-consistent calculation of the potential drop. For the (3,0) tube, however, the curve obtained from a self-consistent treatment is also given for comparison. One can notice that the current from the self-consistent calculation is larger than the current of the step potential calculation. The smoothening effect of the self-consistent treatment of the (3,0) tube is of the same order as that of a larger tube (7,0). Therefore, we expect that the self-consistent calculations of the current for larger tubes also differ slightly from the step potential ones.

It is seen from Fig. 5(a) that for the armchair case, there are regions of negative differential resistance (NDR). They arise as a result of the reduction of conduction under certain bias voltages that cause the channels with different rotational symmetries to coincide [18]. Tunneling between such conduction channels are prohibited. After the bias voltage reaches the width of the (pseudo) gap of these metallic tubes, other channels start to conduct, and hence the current increases. We can notice an enhancement of the NDR as the tube radius gets smaller. The NDR feature has a wide range of applications, including amplification, logic, and memory, as well as fast switching [31]. The unique character of the NDR of metallic nanotubes is that its mechanism, that is, the selection rule involved, is a direct consequence of the rotational symmetry of carbon nanotubes and is different from the mechanism responsible for NDR in Esaki diodes and resonant tunneling structures. As the main cause of NDR is the rotational symmetry selection rule between the eigenstates of bulk systems on left and right, we do not expect the assumption of sharp potential drop to modify qualitatively the results of our calculations.

Next, we consider the case of zigzag tubes in Fig. 5(b). Although there exist regions of NDR for the (3,0) metallic zigzag tube, the $I(V)$ characteristics of (7,0) and (5,0)



Figure 4. A doped nanotube junction generated by inserting donor and acceptor atoms, for example, $X$ and $I$, inside a nanotube.

**Figure 5.** $I(V)$ characteristics of some doped nanotube junctions of metallic (a) and semiconducting (b) nanotubes. $I$ is in units of $[(2e/h)$ hopping] and $V$ is in units of [hopping] [18].

semiconducting zigzag tubes do not contain regions of NDR. However, there exists a region of zero current for nonzero potential differences in the $I(V)$ characteristic of these semiconducting tubes, which is asymmetric with respect to the sign of the bias. This arises from the asymmetric modification of the gap of these tubes due to the initial dopings. Because of the asymmetry of its $I(V)$ characteristic, the junction of doped semiconducting tubes can function as a rectifier in much the same way that an ordinary p-n junction is used for rectification, provided that the difference between the left and right initial dopings is large enough.

## 2.3. Nanoelectromechanical Sensors and Switches Based on Bent Nanotubes

The interrelationship of the electronic and mechanical properties of carbon nanotubes [32–38] gives rise to natural speculations for possible applications. It has been shown both experimentally [39] and theoretically [40, 41] that the reversible bending of nanotubes can be used to alter their conduction, which, in turn, may be used in nanoelectromechanical switch/sensor applications. For calculating $I(V)$, we again use a tight-binding formulation but employ of a four orbital per atom parameterization introduced by Xu et al. for carbon [42]. Our aim is to investigate the interrelationship of mechanical response and transport properties of typical metallic and semiconducting carbon nanotubes [19]. We use the general formalism introduced earlier, in order to derive the $I(V)$ characteristics of the bent nanotubes.

To obtain the relaxed structures under bending and to calculate transport properties, we consider parts of a (6,6) armchair and a (10,0) zigzag nanotubes that contain 972 and 940 carbon atoms, respectively. The diameters of the (6,6) armchair and (10,0) zigzag nanotubes are 8.1 and 7.8 Å, respectively. The lengths of these nanotube portions are 98 Å. Considering the large number of atoms in the systems that makes ab initio geometry optimization formidable, and taking into account the disadvantage of using classical potentials due to ignoring hybridizations, we choose the four orbital per atom tight-binding approach of Xu et al. for carbon [42], both to obtain the optimized geometries and to calculate the electronic/transport properties. Geometry optimizations are performed via the O(N) density-matrix electronic structure calculation method of Li et al. [43], combined with the Broyden



**Figure 6.** Optimized geometries of (6,6) (a) and (10,0) (b) nanotubes under different bending angles [19].

Figure 7. $I(V)$ characteristics of the bent nanotube structures of Fig. 6 [19].

minimization scheme [44], within the previously mentioned tight-binding approach. The optimization of the bent structures proceeds as follows. For successive bending angles, while fixing eight carbon rings (96 atoms for the armchair case, and 80 atoms for the zigzag case) at each end of the nanotube, the structure is optimized such that the maximum force acting on the unconstrained atoms becomes less than 0.05 eV/Å. The results of geometry optimizations are depicted in Fig. 6.

For calculating conductance, two semi-infinite perfect nanotubes (i.e., "leads"), are assumed to be attached to the two ends of the bent region. When a bias voltage $V$ is applied to the bent region, the chemical potentials of the two leads are assumed to be pinned at $V/2$ and $-V/2$. This determines the shifts of the band structures and density of states of the leads. This is achieved by shifting the on-site elements of the tight-binding Hamiltonians of the two leads by $V/2$ and $-V/2$. Within the bent region, however, the on-site shifts are determined by the potential drop pattern. As for the functional form of the potential drop, which is necessary in calculating $I(V)$ characteristics, we assume a linear drop across the bent region. Although self-consistent calculation of the accurate pattern of the potential drop is possible in principle [15, 18], applying this approach to the bent nanotubes of this study is currently formidable because of the large number of carbon atoms involved. Our assumption of linear drop is justified by the observation that, as we shall see shortly, the deformations within the bent nanotubes are distributed more or less uniformly.

The $I(V)$ curves are depicted in Fig. 7. From this figure, it is evidently seen that as the bending angle increases, the current passing through the bent armchair (6,6) tube decreases, while that of the zigzag (10,0) tube increases. This difference is caused by the fact that the localized states produced as a result of bending can increase the tunneling probability for the zigzag tube but not for the armchair tube, whose transmission is restricted by the limited number of channels of the leads within their pseudogap [19].

The results of Fig. 7 can be used for designing nanoelectromechanical sensors and switches, as there is evidently a correspondence between the bending angle and the current that passes through the nanotubes at constant bias. In particular, the current through the semiconducting nanotube can be switched on or off depending on the bending angle.

# 3. STRUCTURAL, ELECTRONIC, AND TRANSPORT PROPERTIES OF "ENAMEL" MOLECULAR WIRE

## 3.1. Basic Concepts of "Enamel" Molecular Wire

The wire is a very important component in molecular electronics because it can be used as a connection between a metal electrode and other functional molecules, such as a molecular diode or transistor, to create complex molecular circuits. It is also important that the molecular wire should have metallic characteristics. Among the different candidates for molecular wire, the conducting polymers are very attractive materials for several reasons: they can be

synthesized with highly controlled length; their electrical conductivity can be control over the full range from insulator to metal by chemical or electrochemical doping [45]; and they can be chemically bonded, without changing of their electronic properties, with other functional molecules to create circuits with more complex functionality. To prevent the possible interaction between different molecular wires it would be better if a single polymer chain were to be encapsulated into a bulky insulated structure and hence forming a molecular enamel wire. The "molecular enamel wire" concept, in which insulators are placed around a conducting center, was first proposed by Wada et al. [46]. The authors also suggested that the "molecular enamel wire" would be one of the key concepts for realizing a high-performance molecular supercomputer [46].

One of the possible approaches for using isolated molecular wires for the realization of this concept is the formation of a self-assembled supramolecular complex between the conducting polymer and cyclic cyclodextrin (CD) molecules as shown in Fig. 1. CDs offer not only the roles of filters and sieves but also offer the subnanoscale of special chemical reaction fields to make a supramolecular complex not formed by natural interactions, such as Van der Waals attractive force and ionic force. The cavity size of CD can be regulated by the number of D-glucose units in each CD molecule (6, 7, and 8 for $\alpha$-, $\beta$-, and $\gamma$-CD, respectively) and a molecular tube can be created by cross-linking adjacent $\alpha$-CD units using a hydroxypropylene bridge [47]. The inner diameters ($D$, see Fig. 8[b]) are 4.5, 7.0, and 8.5 Å for $\alpha$-, $\beta$-, and $\gamma$-CDs, respectively. The depths ($d$) are nearly the same length of 7.0 Å [48]. Recently, the formation of such inclusion complexes between a conducting polymer, polyaniline with emeraldine base, and CDs was realized [49, 50]. Atomic force microscopy (AFM) and scanning tunneling microscopy (STM) observations indicated the formation of an inclusion complex in which the polymer is fully covered by $\beta$-CD molecules [49] and also a molecular nanotube of cross-linking $\alpha$-CD molecules [50]. The experimental results indicated that, in such molecular nanotube, the conformation of the polyaniline chain remains rodlike (all trans conformation).

In this section, we have discussed the structural, electronic, and transport properties of various polymer-CDs inclusion complexes. Moreover, the effect of metal contacts on the structural and electronic properties of polymer chains in CD complexes has been analyzed. The structure of the doped polymer wires in CD complexes as well as the conductance of the doped conducting polymer has been also investigated in order to understand the possibility



**Figure 8.** (a) Chemical form of D-glucopyranose. (b) Shape of a cyclodextrin which is synthesized from D-glucopyranoses. Diameter $D$ is the inside diameter at the middle of cyclodextrin and $d$ is the depth of it. (c) Schematic diagram of inclusion complex formation of CDs and conducting polymer.

of using this inclusion complex for active molecular wire interconnections. The aim of this study is to demonstrate the possibility of the realization of isolated "enamel" molecular wires using quantum mechanical simulations.

## 3.2. Configuration of Polymers into Cyclodextrin Molecular Nanotubes

The control of the structural order of conjugated polymers is important for realization of molecular wires with good electrical conductivity because the carrier mobility and hence the electrical conductivity is limited by their structural disorder. It is important to know what configuration of polyaniline (PANI) fragment is formed in CDs because the source of conductivity for a conjugated conducting polymer is a set of $\pi$-type molecular orbitals that lie above and below the plane of the molecule when it is in a planar or near-planar conformation. To obtain the optimized polymer structure inside CD molecules, the two-layered "own N-layered integrated molecular orbital and molecular mechanics" (ONIOM) method [51] has been applied. In this hybrid method, the structure of seven thiophene monomers is treated quantum mechanically (Hartree-Fock [HF] and various density functional methods), while the remainder of the system (two $\beta$-CDs) is treated by a semiempirical method.

The polymer fragment has been optimized in free space using both full optimization to find the lowest energy structure and partial optimization while maintaining the planar configuration of the PANI fragment as shown in Fig. 9(a) and 9(b), respectively [52]. In the case of full optimization, the imaginary frequencies are absent, which means that it is a local minimum. In this configuration, the adjacent benzene rings have a dihedral angle of 90°. This would reduce the extent of $\pi$-orbital overlap between adjacent rings break up the electron channels and decrease the conductivity of the molecular wire. The planar structure is higher in energy by 38.65 kcal/mol at HF/6-31G⁺ level compared with the most stable structure. Figure 9(c) shows the transmission spectrum of PANI in both planar and most stable configurations under zero bias. The transport calculations show that in the case of bent geometry the conductivity decrease due to a reduction the extent of $\pi$-orbital overlap between adjacent benzene rings. Thus, in the case of most stable PANI configuration, the $G(E_F)$ is found to be $0.0009G_0$ (where $G_0 = 2e^2/h$ is a unit of quantum conductance). The conductivity increases for planar conformation and PANI has a $G(E_F)$ of approximately $0.006G_0$, which is six times larger than that of nonplanar one. The small conductance values of both PANI fragments

(a) $L = 39.2$ Å                                    (b) $L = 46.9$ Å



(c)        Conductance $(2e^2/h)$



E (eV)

**Figure 9.** Structural analysis of PANI fragments: (a) the most stable configuration and (b) the planar configuration in free space. (c) Conductance of the bend (dashed line) and planar (solid line) PANI structures. The Fermi level (vertical dotted line) has been chosen as zero energy.

near Fermi level indicate that they are in semiconductor state in agreement with the large value of calculated energy gap between the highest-occupied (HOMO) and the lowest-unoccupied (LUMO) orbitals of these oligomers. The imaginary frequencies are found for the planar configuration and correspond to the combination of out-of-plane (bent) vibrations of benzene rings. Therefore, it is necessary to apply external forces for stabilization of the planar configuration.

The geometry of the PANI has been also optimized in $\beta$-CDs and the cross-linking $\alpha$-CDs. Figure 10(a) and 10(b) show the respective optimized structures [52]. The structure of the PANI in $\beta$-CDs lies higher in energy by 7.10 kcal/mol as compared with the most stable configuration of the same PANI fragment in free space. In the case of the cross-linking $\alpha$-CD molecular nanotube, this energy difference is found to be a 20.91 kcal/mol. Moreover, the length of polymer chain in this case (45.5 Å) is very close to the length of planar conformation of polymer chain (46.9 Å). These results indicate that the configurations of PANI in the cross-linking $\alpha$-CD molecular nanotube are closer to the planar structure of PANI in free space than the configuration of PANI in $\beta$-CDs. It has also been found in the both cases that there is no charge transfer between polymer fragment and CDs and hence the interaction between these molecules has a noncovalent character.

Because of the steric hindrance, the backbone of PANI is not planar and the lack of strict planarity has led to local distortions, the so-called ring-twist distortions. However, in the case of a heterocyclic polymer, such polythiophene (PT), these distortions are much smaller, and hence a near-planar structure with high degree of polymerization that is well-ordered and stable in the air at room temperature can be synthesized [53]. Therefore, the structural properties of polythiophene inside CD tubes have been also investigated [54]. The optimized structures of PT fragments in molecular nanotubes of cross-linking $\alpha$-CDs and in $\beta$-CDs are shown in Figs. 11(a) and 11(b), respectively. The results of calculations indicate that the configurations of trapped PT are close to the planar structure in the cases of $\beta$-CDs and molecular nanotubes of cross-linking $\alpha$-CDs. The structures of PT fragment in $\beta$-CDs and molecular nanotubes of cross-linking $\alpha$-CDs lie at energies higher by 0.296 and 0.723 kcal/mol, respectively, than the configuration of five monomer units in free space. It is interesting to note that the different results have been obtained in the case of the $\alpha$-CD-polythiophene inclusion complex. The configuration of the PT fragment lies at energy higher by 2.525 kcal/mol than the most stable configuration of five-monomer units. Moreover, the $\alpha$-CD molecules are distorted and lose structural order along the polymer chain.

The results of the structural optimizations show that the near-planar configurations of the various polymer chains are formed in a molecular nanotube of cross-linking $\alpha$-CDs due to weak interactions such as Coulomb and van der Waals interactions between the host framework of the CDs and polymer chains.

## 3.3. Electronic Properties of "Enamel" Molecular Wires

To understand the electron transport through polymer chain in CDs, we have analyzed the spatial extent of the frontier orbital, which provides a strategy by which the transport properties of these systems can be understood [52, 54, 55]. Analysis of the molecular orbital energy diagrams for the configurations of PANI in the cross-linking $\alpha$-CD nanotube and $\beta$-CDs host frameworks (Fig. 12[a] and 12[b], respectively) shows that the lowest unoccupied (LUMO and LUMO + 1) orbitals (as well as LUMO + 2) are located on the polymer fragment and

(a) $L = 43.9$ Å          (b) $L = 45.5$ Å



**Figure 10.** Structural analysis of PANI fragments: (a) in $\beta$-CDs and (b) in molecular nanotube of cross-linking $\alpha$-CDs.

**Figure 11.** Structural analysis of PANI fragments: (a) in molecular nanotube of cross-linking α-CDs and (b) in β-CDs.

their contours are similar to those in the case of the planar configuration of PANI in free space. The same results have been observed in the case PT fragment [54]. Figures 12(c) and 12(d) indicate that the electronic structure of PT is almost the same as that of the most stable conformation of the PT fragment in free space. Moreover, in all cases, there is no overlap of electron density between CDs and the polymer fragments and the charge transfer between the polymer fragment and host framework of CDs has not been found. This reveals that the inclusion complex based on the cross-linking α-CD nanotube and β-CD host frameworks



**Figure 12.** Contour of the lowest-unoccupied orbitals LUMO and LUMO + 1 of (a) PANI-cross-linking α-CDs (b) PANI-β-CDs, (c) PT-cross-linking α-CDs, and (e) PT-β-CDs.

can be used as molecular enamel wire. These results are important from a practical point of view because the electron will transport only through the polymer chain and there is no current leakage across the CD molecules.

To estimate the conductance through molecules attached to the gold electrode, it is important to understand the molecular-metal contact influence on the electronic structure of the molecule because the wires connect to the metal contact through the alligator clips in the experiment. It is known for S-terminated organic molecules that the most probable adsorption side on the gold surface is an equilateral triangle (threefold adsorption side) in which a chemically bonded sulfur atom overlaps equally with three gold atoms [56]. Therefore, we optimized the PT fragment connected to Au contacts, which was approximated by three and six gold atoms in order to estimate the structural properties of PT-gold contacts [57, 58]. This metal cluster configurations are the minimal models for a one- and two-layer configuration of the threefold side of the Au(111) surface, respectively. For a conjugated molecular wire, it can be assumed that the lowest-unoccupied orbitals that span the length of the molecule are responsible for the electron transport and these orbitals can be termed as the conduction MOs. Figure 13 shows the HOMO and LUMO + K (K values are taken up to the first conduction MO) energy values with and without metal atoms attached. In the case of the free molecular wire, the HOMO and LUMO orbitals are delocalized over the polythiophene fragment. In the case of attachment of six gold atoms [57], the LUMO orbital of the free molecule resembles the LUMO + 5 of the metal-polythiophene complex due to the localization of LUMO, LUMO + 1, LUMO + 2, LUMO + 3 and LUMO + 4 on the metal clusters. The energy difference between the HOMO and conduction MO in the case of the metal-wire complex (2.86 eV) is slightly larger than that in the case of the free molecular wire (2.38 eV), as shown in Fig. 13. It was also found that the structure of the PT fragment is not changed significantly even for the thiophene ring nearest to the surface. These results indicate that the metal-molecule interaction is localized to the interfacial region and hence the metal contacts do not interfere with the molecular orbitals of the long molecular wires.

## 3.4. Transport Properties of "Enamel" Molecular Wires

To accurately estimate the conductance through various molecular devices, it is desirable to include metal atoms as part of a device because the surface effect will be automatically accounted for self-consistent calculations of conductance. Our previous results also showed that the conductance characteristics depend on the model used for the electrodes [59]. Moreover, the size of these clusters should be sufficient to eliminate the interaction between the end of the metal cluster and thiophene oligomer. Therefore, the PT fragments with and without $\beta$-CDs connected to two Au$_{22}$ gold clusters have been considered as shown in Fig. 14. The Au$_{22}$ clusters arranged in the FCC(111) geometry for right and left contacts have been selected. These clusters have 6 layers (with 3, 3, 5, 3, 3 and 5 atoms in each layer,



Figure 13. Schematic MO diagrams of the most stable configurations of the PT fragments: (a) without gold and (b) connected with two gold (Au) clusters.

**Figure 14.** Models used for transport calculations: (a) PT fragment encapsulated in $\beta$-CDs with two gold contacts and (b) PT fragment connected to gold contacts.

respectively), which are just extended configuration of $Au_n$ contacts used in the previous calculations [58]. Figure 15(a) shows the electron density of states (DOS) of these devices at zero bias voltage. The first large peak was observed in the energy region between −12 and −6 eV. The molecular orbital analysis in this energy region for both devices shows that the orbitals are located mainly at the gold clusters (see Fig. 15[c]) and hence are apparently irrelevant for electron transport. This is in agreement with results of the conductance calculations as shown in Fig. 15(b). There are no conductance peaks in this energy region for both PT and PT-CD complexes. It has also been found that the conducting channels near the Fermi energy in the case of the polymer chain trapped in CDs are practically the same as those in the case of the single polymer fragment. Analysis of the molecular orbital for the PT-CD complex connected to the gold contacts in that energy region (between −5 and 5 eV) shows that these orbitals span the length of the whole polymer chain is the same as the orbitals of the isolated PT fragment (see Fig. 15[d]). Therefore, these conductance channels correspond to electron transport through the polymer chain even in the case of CDs-polymer inclusion complex because there are not localized on CD molecules. The small differences can be explained by weak noncovalent interactions between polymer and CD molecules which slightly affected on geometry of PT fragment inside CDs molecular nanotube [58]. The conductance channel in the high energy region (20–33 eV) is related to the conductance of CDs and it can be activated only under very high applied bias. The small



**Figure 15.** (a) Density of states (DOS) and (b) conductance at zero bias voltage for the isolated PT fragment (solid line) and encapsulated in $\beta$-CDs (dashed line). Contour of selected MOs (c and d) for PT fragment and for PT fragment in $\beta$-CDs.

value of conductance near Fermi level indicates that the polythiophene fragment is in semi-conductor state in agreement with electronic calculations of that oligomer which shows a large HOMO-LUMO energy gap (2.38 eV). Therefore, to use these inclusion complexes as molecular wires, the metallic state of polymer inside CDs should be realized.

## 3.5. Doping of Polymer inside Molecular Nanotube

To realize the concept of molecular enamel wire, it is also necessary to understand the stability and the electronic properties of the conducting polymers in the metallic state when they are encapsulated within molecular nanotubes. Charge injection onto conjugated polymers can be achieved by charge-transfer redox chemistry, such as oxidation ($p$-type doping) or reduction ($n$-type doping) [45]. Schematically it can be described by the following equations:

1. $p$-type doping

$$(\pi - \text{polymer})_n + 3/2 n y(I_2) \rightarrow [(\pi - \text{polymer})^{+y}(I_3^-)_y]_n \qquad (1)$$

2. $n$-type doping

$$(\pi - \text{polymer})_n + n y[\text{Na}^+(C_{10}H_8)^-] \rightarrow [(\text{Na}^+)_y(\pi - \text{polymer})^{-y}]_n + n y(C_{10}H_8)^0 \qquad (2)$$

It is well known through experiment that protonation by the acid-base chemistry leads to an internal redox reaction and the conversion from semiconductor (the emeraldine base; EB) to metal (the emeraldine salt; ES) [45]. The five monomers in free space have been optimized using both full optimization to find the lowest energy structure, and partial optimization while maintaining the planar configuration of PANI with ES [54]. It is found that the lowest-energy structure of ES with five monomer units has a total spin of S = 1, which indicates the existence of two unpaired spins. The HOMO-LUMO energy difference is significantly reduced as compared with the same energy difference for EB which has the same number of benzene rings. This indicates the transition of PANI from semiconducting to metallic state. It is also found that by using the cross-linking CD molecular nanotubes one can stabilize the near-planar configuration of the metallic form of PANI. Analysis of molecular orbital energy diagrams for the configuration of EB in molecular nanotube shows that the single-occupied-molecular orbitals (SOMO, SOMO + 1) as well the lowest-unoccupied (LUMO + 2 and LUMO + 3) orbitals (Fig. 16) are located on the polymer fragment and their contours are similar to those in the case of the planar configuration of ES. These orbitals are located on polymer chain and hence the CDs can be used as insulator between different single molecular wires [54]. Therefore, the present theoretical results provide support for the "molecular



Figure 16. Contour of the selected molecular orbitals of ES of PANI fragment in a molecular nanotube of cross-linking α-CDs: (a) SOMO; (b) SOMO + 1; (c) LUMO + 2; and (d) LUMO + 3.

(a)



(b)  Conductance ($2e^2/h$)



E(eV)

Figure 17. (a) Model used for transport calculations for Na doped PT fragment and (b) the conductance for the Na doped (solid line) and undoped PT fragment (dashed line). The Fermi level (vertical dotted line has been chosen as zero energy.

enamel wire" concept [46] and hence indicate a high possibility of realizing this concept in molecular electronics.

It is difficult to realize the $p$-type doping of iodine ($I_3$) molecules into PT inside the CDs molecular nanotube because of their large size. Therefore, we consider the $n$-type doping of PT using Na atoms. First, the geometry optimization of a Na-doped PT fragment was performed [60]. The charge transfer is observed from Na to the polymer chain (0.81$e$ per Na atom) where the charge goes to the inner thiophene rings. This results in a significant geometry modification of the inner rings that are located closer to the Na atoms. Thus, there is an interchange of the single C—C and double C=C bonds as compared with the undoped case. This leads to a local transformation to a quinoidlike structure of PT which has metallic characteristics. The value of the HOMO-LUMO energy difference is found to be 1.35 eV, which is significantly reduced as compared with that of the undoped case (2.38 eV). The changes in the geometry and electronic structure affect the conductance properties of the PT fragment [60]. Figure 17 shows the transmission spectrum before and after Na doping under zero bias. It has been found that the conducting channel near the Fermi energy in the case of the Na-doped polymer chain is nearly twice as large ($G(E_F) = 0.225G_0$) as that in the case of the undoped polymer fragment ($G(E_F) = 0.116G_0$). Moreover, in the energy

(a)                                    (b)



(c)



Figure 18. Optimized geometries of Na-doped PT fragment: (a, b) in $\beta$-CDs and (c) in molecular nanotube of cross-linking $\alpha$-CDs

**Figure 19.** (a) Structural analysis, (b) contour of the HOMO, and (c) LUMO + 2 of Na-doped PT fragment in molecular nanotube of cross-linking α-CDs.

region near the Fermi level, there are more peaks in the case of the Na-doped PT fragment, and these conductance channels correspond to the electron transport through the polymer chain. The results indicate the increase in the conductivity of the doped polymer chains, which is in agreement with many experimental results related to metallic behavior of doped conjugated polymers [45].

The possibility of doping a polymer chain inside the CD molecular nanotube has been also investigated. Therefore, the structures of *n*-type doped PT fragments in various inclusion complexes based on CD molecules were also studied [61]. For the β-CDs, two initial structures have been selected in which the Na atoms are located inside and outside the CD molecules (Fig. 18a and 18b, respectively). In the first case, after optimization strong deformation of the PT fragment is observed. Moreover, one of the Na atoms is moved outside the CD molecules. In the second case, the doping atoms remain closer to their initial positions but the distance between the CD molecules increases. This indicates that doping of PT in the case of the β-CDs is difficult because one must control the separation distance between the β-CD molecules to prevent a large deformation of the PT chain. However, in the case of a molecular nanotube of cross-linking α-CDs (Fig. 18[c]), such control can be easily realized because the CD molecules are connected by chemical bonds. As in the case of doped PT in free space, the geometry modifications in the inner rings of PT provoke, along the carbon path, the interchange of the single- and double-like bonds. The bonds between rings decrease from 1.436 to 1.380 Å, as shown in Fig. 19(a), which means a local transformation of PT from aromatic to quinoid form as in the case of doping PT without covering by CD molecules. Analysis of the molecular orbital energy diagrams (Fig. 19[b] and 19[c]) for the configurations of doped PT in the CDs host framework shows that the HOMO and LUMO+2 orbitals (as well as LUMO and LUMO + 1) are located on the polymer fragment and their contours are similar to those in the case of the planar configuration of doped PT in free space.

Thus, for the formation of an isolated metallic single-polymer chain, it is necessary to control the separation distance between CD molecules. This can be easily achieved in the case of cross-linking α-CDs because the α-CD molecules are connected by a hydroxypropylene bridge with highly controlled length. The present theoretical results in combination with the experimental data provide support for the "molecular enamel wire" concept and hence can suggest this supramolecular system as a good candidate for realizing this concept in molecular electronics.

## 4. PORPHYRIN- AND PHTHALOCYANINE-BASED DEVICES

The porphyrin molecule is also a promising material for future nanoelectronics applications since it can be used as a building block in various molecular devices. Moreover, the porphyrin-based complexes would be potential candidates for other applications, such as

materials for novel nonlinear optic devices, infrared detectors, spintronic devices, and photochemical energy conversion [62]. Thus, it is very important to investigate the structures and electronic properties of different metal-containing porphyrin molecules by using first-principles calculations.

In this section the structural, electronic, and transport properties of various porphyrins-based complexes have been presented. The effect of metal doping on the structural and electronic properties of porphyrin molecule has been analyzed. The structures of the porphyrin wires and their electronic properties have been also investigated in order to understand the possibility of using this complex as molecular rectifier. Moreover, the stability of phthalocyanine-fullerene supramolecular structure has been also studied.

## 4.1. Electronic and Transport Properties of Metal Porphyrins

The idea to use porphyrin as a building block in molecular wires has been supported by the recent discovery of different porphyrin arrays having rigid geometric structures and stability in air [63]. Using such arrays, it is possible to control the $\pi$ orbital delocalization, which is desirable for molecular wire applications. Moreover, such porphyrin polymers can be doped by different metals that will also affect on electron transport through the porphyrin chain. Several theoretical studies related to the transport properties of porphyrin molecular wires have been performed [64–66]. In those studies, the different types of conjugated porphyrin molecular wires both physically deposited on Al contacts [64, 65] and chemically bonded to Au contacts [66] were investigated. It was shown that the conductivity depends on the type of conjugated connection between the porphyrin monomers as well as on the type of metal-molecule contact. Recently, we have studied the effect of metal doping on the electronic and conductance properties of porphyrin using different types of metal atoms [60, 67].

We consider a porphyrin molecule with two thiophenyl substituents at its right and left side, respectively, because thiol-type sulfur atom has widely been used as the alligator clip for the connection to Au electrodes. The structure and electronic properties both undoped, so-called free-based (FB) and metal-doped porphyrins have been examined (see Fig. 20). Seven first-row transition-metal complexes (from Cr to Zn) of 5,15-di-(4-thiophenyl)-porphyrin (abbreviated as MDTP), have been selected for the investigation. The geometry optimization of MDTPs has been performed by using of density functional theory (DFT) couplet with B3PYL exchange-correlation functional and 6-311G basis set. Moreover, the most stable spin configurations of MDTPs (M = Cr, Mn, Fe, Co, and Ni) have been examined and the only most stable spin configurations have been selected for the consideration. It is well-known fact that the density functional approximation (for example, BLYP) cannot reproduce the energy of excited state and always gives the underestimated value as compared with experimental one. The introduction of hybrid functional in the DFT formalism gives a larger value as compared with pure DFT one. Moreover, the B3LYP functional within the time-dependent density functional theory have been shown to produce low-lying excitation energies that are in excellent agreement with experiment for porphyrins [68]. It has been also found that the optimization geometry of zinc tetraphenylporphyrine (ZnTPP) is better reproduce the experimental structure by B3LYP functional than by BLYP functional. The 6-311G and 6-31G(d)



Figure 20. Structure of metal 5,15-di-(4-thiophenyl)-porphyrin (MDTP) (M = Cr, Mn, Fe, Co, Ni, Cu, and Zn).

basis sets have been used. The 6-311G basis set is commonly used for the prediction of the electronic structure of the first transition-row complexes, especially Fe [69], and therefore, it is mostly used in this study. For comparison, the 6-31G(d) basis set has been also applied in order to verify the obtained calculated results. All calculations have been performed by using Gaussian 98 set of programs [70].

It has been found that the dihedral angle between porphyrin ring and thiophenyl group is nearly perpendicular. Although it is well known that the phenyl ring at the mesoposition prefers perpendicular conformation, the results for thiophenyl group may be interesting from the experimental point of view, especially the energy difference between planar and perpendicular conformations. For instance, to accurately measure the transport through these molecular structures, they should be properly adsorbed on metal contacts. Figure 21 shows the energy difference of ZnDTP as related to the most stable structure as a function of the rotation angle of the thiophenyl substituents. As it would be expected, from this figure it is clear that the near perpendicular conformation of mesosubstituents is the most stable, while the planar conformation has the highest energy with the energy difference about 10.3 eV. The energetically favorable perpendicular conformation leads to the destruction of the $\pi-\pi$ interaction between the porphyrin core and mesosubstituents and thus to the reduction of the electron transport properties in this system. It is also interesting to note that there is no significant change in energy, when the torsion angle between mesophenyl group and porphyrin core varied from 70° to 90°. On the other words, the perpendicular conformation is flexible and can be easily changed under external perturbation, such as temperature.

The calculated M-N distances in CoDTP and NiDTP are close to 1.98 Å that shorter than those in ZnDTP, MnDTP, and CrDTP, which are around 2.07 Å. The bonds between carbon and nitrogen atoms in pyrrole rings are practically the same for different MDTPs. The agreement between the calculated and available experimental data is quite good [71–74] with the largest deviation of 0.04Å. Accordingly, it is confirmed that the selected calculation method is accurately predicts the structural properties of studied molecules. The most stable spin configuration for each of MDTPs has been determined by the comparison of the total energy of optimized MDTP structures with the different spin configurations. First, we have calculated the most stable spin configuration of FeTPP using B3LYP/6-311G level and compared with available theoretical and experimental data to make sure that the selected method can accurately predict the magnetic state of metal in metalporphyrin complexes. The most stable spin configuration of Iron tetraphenylporphyrine (FeTPP) has total spin $S = 1$ (Symbol $S$ indicates the total spin of the system, which is the difference between the number of spin-up and spin-down electrons). The state with $S = 2$ lies 0.51 eV above most stable spin configuration and in comparison with a magnetic susceptibility measurement that yielded a value of 0.62 eV [75]. The lowest closed-shell state ($S = 0$) lies 0.98 eV above the ground state. The recently theoretical results are showed the same magnetic configuration of FeTPP. The most stable spin configuration ($S = 1$) and the energy difference between this



**Figure 21.** Relative energy of ZnDPT as a function of dihedral angle between porphyrin plane and thiophenyl groups.

spin state and other configurations ($S = 2$ and $S = 0$) are found to be 0.75 eV and 1.15 eV, respectively [76]. Predicted the most stable spin configurations for Fe, Co, and Ni complexes are in an agreement both with the theoretical [76] and experimental results [71, 77, 78] obtained for the metal-containing mesotetraphenyl porphyrins (MTPP). Table 1 shows the difference in total energy for each of the spin configurations of MDTP as related to the total energy of the most stable configuration defined as 0. CuDTP and ZnDTP have the only one spin configuration, $S = 1/2$ and $S = 0$, respectively. In the cases of Cr, Mn, Fe, Co, and NiDTPs, the most stable spin configurations are $S = 2$, 5/2, 1, 1/2, and 0, respectively. The smaller energy difference between the lowest electronic configurations for the triplet and quintet states in FeDTP (0.21 eV) as compared with that in FeTPP (0.75 eV) [76], can be explained by the reduction of the molecular symmetry from $D_{4h}$ in FeTPP to $C_{2h}$ in FeDTP. The calculated large-energy difference between the different spin configurations in the Mn, Fe, and Cr complexes suggest that these compounds can be used as a building blocks in the spintronic devices.

The conductance properties of metal 5,15-di-(4-thiophenyl)-porphyrin (MDTP, where M = Zn or Ni) have been investigated and the model of calculations has been shown in Fig. 22[a]). The transport calculations show that the conductivities are very low for all cases (Fig. 22[b]). The other interesting observations in this figure are the increasing conductivities of the MDFTs [60]. Thus in the case of DTP, the $G(E_f)$ is found to be $0.0003G_0$, which is an almost negligible value in agreement with previous theoretical results [66]. The conductivity increases for MDTP. ZnDTP has a $G(E_f)$ of approximately $0.0012G_0$, which is only four times larger than that of DTP, yet still very small. A larger value has been calculated for NiDTP ($0.0177G_0$).

To understand the effect of metal-ligand interactions on the conductance properties of MDTP, the $3d$ orbital splitting of different central metal atoms (Zn and Ni) by a porphyrin ligand field (DTP) has been analyzed. There are two possible schemes of $3d$ orbital splitting by the square-planar ligand field. The first is an in-plane splitting scheme, which reflects a dominant in-plane interaction between the metal $3d$ ($d_x 2 \,_y2$ and in part $d_z2$) orbitals and the porphyrin $\sigma$-donor orbitals. The second is an out-of-plane splitting scheme, which reflects a dominant out-of-plane interaction between the metal $3d$ ($d_{xz}$, $d_{yz}$, and $d_{xy}$) orbitals and $\pi$-donor or acceptor orbitals of the porphyrin ligand. It has been found, by analysis of the electronic structures, that the $3d$ orbital splitting has a predominantly in-plane character in the case of ZnDTP (Fig. 23[a]), while it has a predominantly out-of-plane character in the case of NiDTP (Fig. 23[b]). Thus Zn has strong $\sigma$-type interactions with the porphyrin core, while Ni has strong $\pi$-type interactions with the porphyrin core. Moreover, in the case of ZnDTP, the orbitals with the metal contribution lie lower than those in the case of NiDTP and are all occupied. In the case of NiDTP, $d_x 2 \,_y2$ is an unoccupied orbital (LUMO + 2) while $d_\pi$ and $d_z2$ are close to HOMO. The absence of electrons in the $d_x2_-y2$ orbital results in a shorter distance between Ni and N, and therefore, the desirable overlap between the orbitals of porphyrin ring is larger in this case as compared with that in the case of the ZnDTP, where all the $3d$ orbitals of the metal are filled. All of these results account for the higher conductivity of NiDTP as compared with that of ZnDTP.

The results imply that incorporation of the Ni into the porphyrin enhance the transport properties larger than incorporation of the Zn. Moreover, the thiophenyl groups

Table 1. Relative energy (eV) of MDTPs as calculated for the various spin states. Zero value indicates the energetically most stable spin configuration.

| Structure | $S = 0$ | $S = 1$ | $S = 2$ |
|-----------|---------|---------|---------|
| CrDTP | 2.84 | 1.97 | 0 |
| FeDTP | 1.53 | 0 | 0.21 |
| NiDTP | 0 | 0.51 | — |

| | $S = 1/2$ | $S = 3/2$ | $S = 5/2$ |
|-----------|-----------|-----------|-----------|
| MnDTP | 0.97 | 0.16 | 0 |
| CoDTP | 0 | 0.20 | — |

Figure 22. (a) Models used for transport calculations metal-porphyrins and (b) conductance for DTP(dotted line with circles). ZnDTP (dashed line with black squares), and NiDTP (solid line). The Fermi level (vertical dotted line) has been chosen as zero energy.

conformation leads to the elimination of the $\pi$ conjugation in whole MDTP structures and other connections of porphyrin to electrodes should be proposed in order to realize the molecular wire for the interconnections between electronic devices.

## 4.2. Electronic Properties of Porphyrin Wires

The realization of molecular-scale rectifying function is one of the most important and fundamental requirements in molecular electronics [10]. Aromatic molecules have $\pi$-conjugation structure and hence electrons can flow easily through the $\pi$-orbital overlapping between adjacent blocks. However, it is possible to increase or decrease the $\pi$-electron density by substituting for the different functional groups in an aromatic system and thereby creating acceptor ($p$-type) and donor ($n$-type) molecular subunits. Moreover, a rectifier could be designed using the combination of these two molecular subunits between two electrodes, in which electrons can flow from cathode to the acceptor and from donor to the anode [79, 80]. To realize a rectifying function by using this scheme, the HOMO and LUMO orbitals have to localize on the donor and acceptor parts, respectively.

Porphyrin possesses good electron-donating properties because of its large easily ionized $\pi$-electron system and various metal porphyrins are available. Moreover, a long molecular



Figure 23. Schematic diagram of $3d$ metal orbital splitting in the square-planar ligand field: (a) ZnDTP and (b) NiDTP.

wire of fully conjugated porphyrin polymer has been reported by experimental group [63], and the theoretical study on the electronic properties of this polymer has been performed by tight-binding method [81]. According to these reports, HOMO-LUMO gap of fully conjugated porphyrin polymer is much smaller than that in usual conjugate polymer. This is a good feature for molecular electronics applications, especially, for a molecular wire, which need good conductance properties. Moreover, we have proposed that a rectifier diode can be created by combining two metal porphyrin molecules with different transition metal atoms. In order to investigate the electron transport properties through this polymer, we have performed the molecular orbital analysis of various porphyrin oligomers and suggested a strategy by which the rectifying properties of the porphyrin polymer can be understood [82].

Many configurations of porphyrin polymers have been already synthesized [63]. Figure 24(a), 24(b), and 24(c) indicate the different structures investigated in this study, fully, partially, and nonconjugated free-base porphyrin polymers, respectively. To evaluate the effect of the molecular structure on a localization of their frontier orbitals, we have investigated the porphyrin polymer fragment based on four monomers. The metal-metal junction has been formed in this porphyrin fragment by inserting of two transition metals (see Fig. 25).

The results of the orbital spatial distribution in a porphyrin polymer obtained by HF/6-311G are shown in Figs. 26 and 27 and Table 2. In Table 2, "No" means that these molecules did not show a rectifying function. It is assumed that the unoccupied orbitals provide channels for electron conduction through the molecules. The energy difference of the lowest-unoccupied levels between a donor and an acceptor has been used to estimate a criterion (potential drop) of a rectifying function. The potential drop in a vacuum can be explained as the difference in the LUMO energies between the donor and acceptor molecules when they are widely separated ($\Delta E_{LUMO} = E_{LUMO}(\text{donor}) - E_{LUMO}(\text{acceptor})$) [80]. It is clearly seen from Fig. 25 that a full planar structure (fully conjugated) does not exhibit the rectification properties (except for the case of Cr-Cu in which the empty porphyrin responses for the donor function). Moreover, similar results are obtained for the partially conjugated fragments. HOMO and LUMO for partially conjugated fragments are delocalized on the entire whole system. A rectifier is reported in the case of a D(donor)-$\pi$-A(acceptor) structure [83], even though, this porphyrin polymer does not exhibit the localized frontier orbital. However, a nonconjugated chain displays rectifying features (Cr-Cu and Zn-Fe). HOMO and LUMO + 5 for Cr-Cu in the nonconjugated polymer form the localized donor side (Cu porphyrin) and LUMO forms the localized acceptor side (Cr porphyrin). Consequently, these results together with previous report [80] indicate that the geometry of spacer plays an important role in localizing the frontier orbitals. The localization of frontier orbitals can predict a rectifying function for a device, but the value of rectification should be estimated by a combination the Green's function approach and first-principles calculations.



Figure 24. Structures of (a) fully conjugated, (b) partially conjugated, and (c) nonconjugated free base porphyrin chains.

**Figure 25.** Four porphyrin monomers, arranged as a metal-metal junction in porphyrin chain.



**Figure 26.** Contour of the selected molecular orbitals for the Zn-Fe pair in fully conjugated porphyrin fragment. (a) HOMO, (b) LUMO, and (c) LUMO + 1.



**Figure 27.** Contour of the selected molecular orbitals for the Cr-Cu pair in nonconjugated porphyrin fragment. (a) SOMO, (b) LUMO, and (c) LUMO + 5.

Table 2. The energy difference $\Delta E_{H-L\text{ MO}}$ and $\Delta E_{HOMO-LUMO}$ of various porphyrin chains.

| Structure | $\Delta E_{H-L\text{ MO}}$ (eV) | $\Delta E_{HOMO-LUMO}$ (eV) |
|---|---|---|
| Fully conjugated(Cr-Cu) | 1.59 | 5.37 |
| Fully conjugated(Zn-Zn) | No | 4.30 |
| Fully conjugated(Zn-Fe) | No | 4.32 |
| Fully conjugated(Zn-Ni) | No | 4.35 |
| Partially conjugated(Zn-Zn) | No | 5.13 |
| Partially conjugated(Zn-Ni) | No | 5.13 |
| Nonconjugated(Cr-Cu) | 0.54 | 7.36 |
| Nonconjugated(Zn-Zn) | No | 6.18 |
| Nonconjugated(Zn-Fe) | 0.11 | 6.16 |
| Nonconjugated(Zn-Ni) | 0.02 | 6.18 |

## 4.3. Photovoltaic Materials Based on Organic Molecule—Fullerene Mixture

Conjugated polymers emerged in the mid-eighties to early nineties, and were developed for a wide range of optoelectronic applications, such as organic transistors, light-emitting diodes, and solar cells. The current general trend in research and development of photovoltaic elements is aimed at producing lower-cost devices. Solar cells based on conjugated polymers alone have been disappointing because of their low-quantum efficiencies. However, an encouraging breakthrough in the development of highly efficient materials has been achieved by mixing electron-donor–type polymers with suitable electron acceptors [84, 85]. This dual molecule approach, for example, using a conjugated polymer/fullerene mixture. has been successful and is well documented [86–89]. Many fullerene-based organic mixture have been proposed as potential materials for organic photovoltaic devices, with their electrochemical and photoelectrochemical properties measured under light illumination. Organic molecule based photovoltaic elements have attracted much attention as a replacement for "inorganic semiconductor" and offer the possibility of cheap, easy-to-produce photovoltaic energy from light. Consequently, an organic molecule/fullerene mixture is therefore a potential material candidate for a photovoltaic cell due to its large and flexible absorption combined with electrical properties similar to the early photo effects of natural photosynthesis [86]. especially solar cells harvesting the infrared part of the solar spectrum.



Figure 28. Chemical structure of phthalocyanine ($M = 2H$, free base phthalocyanine, H2-Pc: $M = $ Zn, zinc phthalocyanine. Zn  Pc).

**Figure 29.** Optimized structure of zinc phthalocyanine-fullerene supramolecule (Zn – Pc + C$_{60}$).

The chemical structure of phthalocyanine is shown in Fig. 28. Phthalocyanine possesses good electron-donating properties due to their large easily ionized $\pi$-electron system, whereas fullerene is good $\pi$-electron acceptor that can be connected with other organic molecules. In other word, phthalocyanine has good electron-donating properties and fullerene is a good $\pi$-electron acceptor. Recently, the synthesis of these supramolecules has been reported [89]. Contrary to other organic molecule-fullerene based supramolecules that were synthesized for photovoltaic applications and have a sigma bond between the polymer and the fullerene [86], the above supramolecule have van der Waals bond instead of a sigma bond. Here, we discuss the electronic structure of the phthalocyanine-fullerene supramolecule after geometry optimization.

The optimized geometries and energetics of all the structural variables have been obtained using first-principles calculations. These calculations have been performed using the Gaussian98 program [70] with Hartree-Fock (HF) theory and suitable basis set. Since the treated fullerene/phthalocyanine supramolecular complex consist of a large number of atoms, the small basis set (3–21 G) has been selected to save computation time. The molecular structure has been energetically optimized to reach the stable structure. After optimization, the electron spatial distributions of different molecular orbitals have been analyzed. The analysis of



**Figure 30.** Comparison of the five highest-occupied and five lowest-unoccupied orbital levels. Pc. ZnPc-C$_{60}$ supramolecule. C$_{60}$.

the molecular orbital provided a strategy for understanding the photovoltaic and electrochemical properties (such as charge separation) of various molecular compositions.

Figure 29 shows the optimized structures of phthalocyanine-fullerene supramolecules using the HF/3–21G level. The planar structure of the free base phthalocyanine is retained in this complex (see Fig. 29[a]), whereas a slight bending is observed in the case of zinc phthalocyanine. Zinc atom is protruding from plane of the phthalocyanine molecule in a fullerene direction. The distance between the zinc atom and the nearest carbon atom in fullerene is 2.54 Å. Zinc atom is the bridge site between the six-member ring of the absorbed fullerene. The charge transfer (0.20$e$) from zinc phthalocyanine to the fullerene while there is no charge transfer to free base phthalocyanine. Figure 30(a) and 30(b) illustrate a comparison of the five highest occupied and five lowest unoccupied orbital levels for the optimized structure of the zinc phthalocyanine, fullerene, and the zinc-phthalocyanine-fullerene supramolecule. It is interesting to note that while the LUMO energy levels of the supramolecule compare well with the LUMO energy level of fullerene, the HOMO energy levels of the supramolecule are close to the HOMO energy levels of the phthalocyanine. The same tendency of molecular orbital localizations for other organic molecules-fullerene supramolecular system has been also observed [90].

## 5. CONCLUSIONS

In this chapter, we introduce recent development in theoretical development on the first principles treatment of transport properties in nanostructured materials. The subject attracts much attention not only theoretically but also experimentally and in industries, because of the expected ending of the present day silicon technology and emerging new era of nanotechnology based on these new quantum mechanical world of transport of electron as the minimum unit of information carrier.

## REFERENCES

1. C. P. Poole, Jr. and F. J. Owens, "Introduction to Nanotechnology." Wiley, New York, 2003.
2. R. V. Belosludov, S. Takami, M. Kubo, and A. Miyamoto, in "Combinatorial Materials Synthesis" (X.-D. Xiang and I. Takeuchi, Eds.), p. 363. Marcel Dekker, New York, 2003.
3. M. C. Payne, M. P. Teter, D. C. Allan, T. A. Arias, and J. D. Joannopoulos. *Rev. Mod. Phys.* 64, 1045 (1992).
4. C. Joachim, J. K. Gimzewski, and A. Aviram, *Nature (London)* 408, 541 (2000).
5. M. A. Reed and J. M. Tour, *Sci. Am.* 282, 86 (2000).
6. M. A. Reed, C. Zhou, C. J. Muller, T. P. Burgin, and J. M. Tour, *Science* 278, 252 (1997).
7. L. A. Bumm, J. J. Arnold, M. T. Cygan, T. D. Dunbar, T. P. Burgin, L. Jones II, D. L. Allara, J. M. Tour, and P. S. Weiss, *Science* 217, 1705 (1996).
8. J. Chen, M. A. Reed, A. M. Rawlett, and J. M. Tour, *Science* 286, 1550 (1999).
9. R. M. Mertzger, *Acc. Chem. Res.* 32, 950 (1999).
10. A. Aviram and M. A. Ratner, *Chem. Phys. Lett.* 29, 277 (1974).
11. A. A. Farajian, R. V. Belosludov, H. Mizuseki, and Y. Kawazoe, *Thin Solid Films* (2005), in press.
12. R. Landauer, *Phys. Lett. A* 85, 91 (1981).
13. R. Landauer, *IBM J. Res. Dev.* 32, 306 (1998).
14. M. C. Munoz, V. R. Velasco, and F. Garcia-Moliner, *Prog. Surf. Sci.* 26, 117 (1987).
15. A. A. Farajian, K. Esfarjani, and M. Mikami, *Phys. Rev. B* 65, 165415 (2002).
16. F. Garcia-Moliner and V. R. Velasco, "Theory of Single and Multiple Interfaces." World Scientific, Singapore, 1992.
17. L. Chico, L. X. Benedict, S. G. Louie, and M. L. Cohen, *Phys. Rev. B* 54, 2600 (1996).
18. A. A. Farajian, K. Esfarjani, and Y. Kawazoe, *Phys. Rev. Lett.* 82, 5084 (1999).
19. A. A. Farajian, B. I. Yakobson, H. Mizuseki, and Y. Kawazoe, *Phys. Rev. B* 67, 205423 (2003).
20. A. A. Farajian, H. Mizuseki, and Y. Kawazoe, *Physica E* 22, 675 (2004).
21. A. A. Farajian, B. I. Yakobson, H. Mizuseki, and Y. Kawazoe, *Int. J. Nanosci.* 3, 131 (2004).
22. M. P. Lopez Sancho, J. M. Lopez Sancho, and J. Rubio, *J. Phys. F* 14, 1205 (1984).
23. M. P. Lopez Sancho, J. M. Lopez Sancho, and J. Rubio, *J. Phys. F* 15, 851 (1985).
24. S. Datta, "Electronic Transport in Mesoscopic Systems." Cambridge University Press, Cambridge, 1995.
25. M. Buttiker, Y. Imry, R. Landauer, and S. Pinhas, *Phys. Rev. B* 31, 6207 (1985).
26. O. Stephan, P. M. Ajayan, C. Colliex, P. Redlich, J. M. Lambert, P. Bernier, and P. Lefin, *Science* 266, 1683 (1994).
27. L. Grigorian, G. U. Sumanasekera, A. L. Loper, S. Fang, J. L. Allen, and P. C. Eklund, *Phys. Rev. B* 58, 4195 (1998).

28. Y. Miyamoto, A. Rubio, X. Blase, M. L. Cohen, and S. G. Louie, *Phys. Rev. Lett.* 74, 2993 (1995).
29. A. A. Farajian, K. Ohno, K. Esfarjani, Y. Maruyama, and Y. Kawazoe, *J. Chem. Phys.* 111, 2164 (1999).
30. K. Esfarjani, A. A. Farajian, Y. Hashi, and Y. Kawazoe, *Appl. Phys. Lett.* 74, 79 (1999).
31. K. Esfarjani, A. A. Farajian, Y. Kawazoe, and S. T. Chui, *J. Phys. Soc. Jpn.* 74, 515 (2005).
32. A. Rochefort, D. R. Salahub, and P. Avouris, *Chem. Phys. Lett.* 297, 45 (1998).
33. A. Rochefort, P. Avouris, F. Lesage, and D. R. Salahub, *Phys. Rev. B* 60, 13824 (1999).
34. M. B. Nardelli and J. Bernholc, *Phys. Rev. B* 60, R16338 (1999).
35. S. Paulson, M. R. Falvo, N. Snider, A. Helser, T. Hudson, A. Seeger, R. M. Taylor, R. Superfine, and S. Washburn, *Appl. Phys. Lett.* 75, 2936 (1999).
36. P. E. Lammert, P. Zhang, and V. H. Crespi, *Phys. Rev. Lett.* 84, 2453 (2000).
37. D. Tekleab, D. L. Carroll, G. G. Samsonidze, and B. I. Yakobson, *Phys. Rev. B* 64, 035419 (2001).
38. D. Bozovic, M. Bockrath, J. H. Hafner, C. M. Lieber, H. Park, and M. Tinkham, *Appl. Phys. Lett.* 78, 3693 (2001).
39. T. W. Tombler, C. Zhou, L. Alexseyev, J. Kong, H. Dai, L. Liu, C. S. Jayanthi, M. Tang, and S. Y. Wu, *Nature (London)* 405, 769 (2000).
40. L. Liu, C. S. Jayanthi, M. Tang, S. Y. Wu, T. W. Tombler, C. Zhou, L. Alexseyev, J. Kong, and H. Dai, *Phys. Rev. Lett.* 84, 4950 (2000).
41. A. Maiti, A. Svizhenko, and M. P. Anantram, *Phys. Rev. Lett.* 88, 126805 (2002).
42. C. H. Xu, C. Z. Wang, C. T. Chan, and K. M. Ho, *J. Phys. Condens. Matter* 4, 6047 (1992).
43. X.-P. Li, R. W. Nunes, and D. Vanderbilt, *Phys. Rev. B* 47, 10891 (1993).
44. K. Ohno, K. Esfarjani, and Y. Kawazoe, "Computational Materials Science from Ab Initio to Monte Carlo Methods." Springer, Berlin, 1999.
45. A. J. Heeger, *J. Phys. Chem. B* 105, 8475 (2001).
46. Y. Wada, M. Tsukada, M. Fujihira, K. Matsushige, T. Ogawa, M. Haga, and S. Tanaka, *Jpn. J. Appl. Phys.* 39, 3835 (2000).
47. A. Harada, J. Li, and M. Kamachi, *Nature* 356, 325 (1994).
48. M. L. Bender and M. Komiyama, "Cyclodextrin Chemistry." Springer-Verlag, Tokyo, 1978.
49. K. Yoshida, T. Shimomura, K. Ito, and R. Hayakawa, *Langmuir* 15, 910 (1999).
50. T. Shimomura, T. Akai, T. Abe, and K. Ito, *J. Chem. Phys.* 116, 1753 (2002).
51. S. Humbel, S. Sieber, and K. Morokuma, *J. Chem. Phys.* 105, 1959 (1996).
52. R. V. Belosludov, A. A. Farajian, Y. Kikuchi, H. Mizuseki, and Y. Kawazoe, *Comp. Mater. Sci.* (in press).
53. M. Kobayashi, J. Chen, T. C. Chung, F. Moraes, A. J. Heeger, and F. Wudl, *Synth. Metals* 9, 77 (1984).
54. R. V. Belosludov, H. Sato, A. A. Farajian, H. Mizuseki, K. Ichinoseki, and Y. Kawazoe, *Jpn. J. Appl. Phys.* 42, 2492 (2003).
55. R. V. Belosludov, H. Mizuseki, K. Ichinoseki, and Y. Kawazoe, *Jpn. J. Appl. Phys.* 41, 2739 (2002).
56. N. B. Larsen, H. Biebuyck, E. Delamarche, and B. Michel, *J. Am. Chem. Soc.* 119, 3107 (1997).
57. R. V. Belosludov, H. Sato, A. A. Farajian, H. Mizuseki, and Y. Kawazoe, *Thin Solid Films* 438–439, 80 (2003).
58. R. V. Belosludov, A. A. Farajian, H. Mizuseki, K. Ichinoseki, and Y. Kawazoe, *Jpn. J. Appl. Phys.* 43, 2061 (2004).
59. A. A. Farajian, R. V. Belosludov, H. Mizuseki, and Y. Kawazoe, *Physica E* 18, 253 (2003).
60. R. V. Belosludov, A. A. Farajian, H. Baba, H. Mizuseki, and Y. Kawazoe, *Jpn. J. Appl. Phys.* 44, 2823 (2005).
61. R. V. Belosludov, H. Sato, A. A. Farajian, H. Mizuseki, and Y. Kawazoe, *Mol. Cryst. Liq. Cryst.* 406, 195 (2003).
62. N. Robertson and C. A. McGowan, *Chem. Soc. Rev.* 32, 96 (2003).
63. A. Tsuda and A. Osuka, *Science* 293, 79 (2001).
64. K. Tagami and M. Tsukada, *Jpn. J. Appl. Phys.* 42, 3606 (2003).
65. K. Tagami, M. Tsukada, T. Matsumoto, and T. Kawai, *Phys. Rev. B* 67, 245324 (2003).
66. K. Tagami and M. Tsukada, *e-J. Surf. Sci. Nanotechnol.* 1, 45 (2003).
67. Y. Kikuchi, R. V. Belosludov, H. Baba, H. Mizuseki, and Y. Kawazoe, *Mol. Simul.* 30, 929 (2004).
68. K. A. Nguyen, P. N. Day, and P. Pachter, *J. Chem. Phys.* 110, 9135 (1999).
69. V. N. Nemykin, N. Kobayashi, V. Y. Chernii, and V. K. Belsky, *Eur. J. Inorg. Chem.* 3, 733 (2001).
70. Gaussian 98, Revision A.11.1 (Gaussian, Pittsburgh, PA, 2001).
71. J. P. Collman, J. L. Hoard, N. Kim, G. Lang, and C. A. Reed, *J. Am. Chem. Soc.* 97, 2676 (1975).
72. P. Madura and W. R. Scheidt, *Inorg. Chem.* 15, 3182 (1976).
73. A. L. Maclean, G. J. Foran, B. J. Kennedy, P. Turner, and T. W. Hambley, *Aust. J. Chem.* 49, 1273 (1996).
74. E. B. Fleischer, C. K. Miller, and L. E. Webb, *J. Am. Chem. Soc.* 86, 2342 (1964).
75. P. D. W. Boyd, A. D. Buckingham, R. F. McMecking, and S. Mitra, *Inorg. Chem.* 18, 3585 (1979).
76. M.-S. Liao and S. Scheiner, *J. Chem. Phys.* 117, 205 (2002).
77. W. C. Lin, *Inorg. Chem.* 15, 1114 (1976).
78. J. Subramanian, in "Porphyrins and Metalloporphyrins" (K. M. Smith, Ed.), p. 555. Elsevier Scientific, Amsterdam, 1975.
79. J. C. Ellenbogen and J. C. Love, *Proc. IEEE* 88, 386 (2000).
80. C. Majumder, H. Mizuseki, and Y. Kawazoe, *J. Phys. Chem. A* 105, 9454 (2001).
81. T. G. Pedersen, T. B. Lynge, P. K. Kristensen, and P. M. Johansen, *Thin Solid Films* 477, 182 (2005).
82. H. Mizuseki, R. V. Belosludov, A. A. Farajian, N. Igarashi, and Y. Kawazoe, *Mater. Sci. Eng. C* 25, 718 (2005).
83. R. M. Metzger, *Chem. Rev.* 103, 3803 (2003).
84. G. Yu and A. J. Heeger, *J. Appl. Phys.* 78, 4510 (1995).

85. J. J. M. Halls, C. A. Walsh, N. C. Greenham, E. A. Marseglia, R. H. Friend, S. C. Moratti, and A. B. Homes, *Nature* 376, 498 (1995).
86. D. M. Guldi, *J. Phys. Chem. B* 109, 114321 (2005).
87. G. Yu, J. Gao, J. C. Hummelen, F. Wudl, and A. J. Heeger, *Science* 270, 1789 (1995).
88. S. E. Shaheen, C. J. Brabec, N. S. Sariciftci, F. Padinger, T. Fromherz, and J. C. Hummelen, *Appl. Phys. Lett.* 78, 841 (2001).
89. T. Nojiri, M. M. Alam, H. Konami, A. Watanabe, and O. Ito, *J. Phys. Chem. A* 101, 7943 (1997).
90. H. Mizuseki, N. Igarashi, R. V. Belosludov, A. A. Farajian, and Y. Kawazoe, *Jpn. J. Appl. Phys.* 42, 2503 (2003).

# CHAPTER 4

# Single-Electron Functional Devices and Circuits

## Takashi Morie,[1] Yoshihito Amemiya[2]

[1]Graduate School of Life Science and Systems Engineering,
Kyushu Institute of Technology, Kitakyushu, Japan
[2]Department of Electrical Engineering, Hokkaido University, Sapporo, Japan

## CONTENTS

# 1. INTRODUCTION

This chapter describes various functional devices and circuits based on single-electron operation, as well as their possible applications. Single-electron circuits composed of nanostructures are promising in the construction of ultimately high-density VLSI (very large scale integration) systems after the era of CMOS (Complementary Metal-Oxide-Semiconductor) technology. The principles described in this chapter are so general that they can be applied in various materials. The nanoelectronic devices can be made of semiconductors, metals, organic molecules, and so on. In particular, silicon devices are most important for the practical point of view.

Regarding the theory of single-electron operation and its application to devices and circuits, excellent reviews and books have been published [1, 2]. The fabrication technologies of single-electron devices will not be described in this chapter. Those using silicon technology are described in some reviews [3].

It is noted that the operation principles described in this chapter are completely different from so-called quantum computing. Quantum computing is a sort of super-parallel computing based on the superposition of quantum states, which is never understood by the concepts of classical mechanics and is never realized in the classical electronic devices such as CMOS (see Section 5.7.1).

In contrast, the principles described here are based on the quantum mechanical theory, but one can have some intuitive classical pictures in which an electron can be treated as a classical particle. As long as the movement of an electron can be considered as that of a pure particle, the single-electron device can be considered classical. In the quantum mechanical system, the condition that an electron existing over a tunnel junction can be considered as a particle is that tunnel resistance $R_T$ is much larger than the quantum resistance $R_Q = h/e^2 = 25.8\text{k}\Omega$, where $h$ is the Plank constant and $e$ $(= 1.6 \times 10^{-19}$ C) is the elementary charge. If $R_T < R_Q$, the existence probability of an electron spreads beyond the tunnel junction, and an electron cannot be considered as a particle anymore.

The single-electron operation can be considered as an extreme case of the operation in the MOS devices. As the device becomes smaller, the number of electrons employed for electronic operation becomes smaller. In a MOS transistor with a design rule of 0.1 $\mu$m, the number of electrons is on the order of $10^4$. In contrast, in nanoelectronic devices, fewer than several tens of electrons are used for operation. The macroscopic operation principles, for example, controlling a current by applying a voltage, cannot be applied to controlling a group of few electrons because of their large fluctuation.

This chapter is organized as follows. In Section 2, principles of single-electron operation are described and the basic single-electron circuits are reviewed. In Section 3, some open issues for practical application of the single-electron devices are discussed. In Section 4, an overview of possible information-processing architectures for single-electron systems is addressed. In Section 5 various functional circuits using single-electron operation are described. Finally, Section 6 concludes this chapter.

# 2. PRINCIPLES FOR SINGLE-ELECTRON OPERATION AND BASIC CIRCUITS

## 2.1. Coulomb Blockade

The Coulomb blockade based on the Coulomb repulsion effect between electrons is the most basic principle that controls a single electron. Let us assume a conduction island that is isolated electrically from other parts, and its capacitance is $C$. When a charge $Q$ is stored in this island, the charging energy $E$ is given by $Q^2/2C$. For the macroscopic capacitance such as 10 fF ($= 10^{-14}$ F), which is a typical value of the gate capacitance of a submicron MOS transistor, $E$ is on the order of $10^{-24}$ J when $Q$ is only one electron charge ($-e$). This energy corresponds to a temperature of 0.1 K because the thermal energy at temperature $T$ K is expressed as $k_B T$, where $k_B$ is the Boltzmann constant. Thus, the effect of charging is negligible compared with thermal energy at room temperature. However, if the island is very small, for example, $C = 10$ aF is assumed, $E$ for one electron can be comparable to the thermal energy at 100 K, and thus the movement of an electron is limited. This is called the Coulomb blockade phenomenon.

The movement of an electron is determined by the total energy of the whole electron system. To use this limited movement of electrons effectively, basic single-electron devices have two electric components: tunnel junctions and normal capacitors. Both components have very small capacitance, less than 10 aF, to cause Coulomb blockade phenomena at room temperature or moderate low temperature. The difference between the two components is that electrons can tunnel the tunnel junction without applying a voltage, but they cannot tunnel the normal capacitor with the normal operation voltage.

To meet the above condition in the tunnel junction, when the junction is made of a dielectric insulator such as silicon oxide, the thickness $d$ of the junction must be less than 2-3 nm. The capacitance $C$ of normal dielectric material is determined by $C = \varepsilon S/d$, where $\varepsilon$ is the dielectric constant of this material and $S$ is the area of the junction, provided the fringe effect is negligible. Even if the insulator is made of vacuum, that is, $\varepsilon = \varepsilon_0 = 8.85$ [pF/m], to obtain $C < 10$ aF, $S$ should be less than 3400 nm$^2$ for $d = 3$ nm. This is the reason why nanostructures are required for single-electron operation.

On the basis of the Coulomb blockade described above and the Coulomb repulsion effect, various single-electron devices and circuits have been proposed. In the following text, some important and useful circuits are reviewed.

## 2.2. Single-Electron Box

The simplest single-electron device is the so-called single-electron box consisting of a tunnel junction and a normal capacitor, as shown in Fig. 1a.

The tunnel junction capacitance $C_j$ and the normal capacitance $C_o$ are assumed to have charges $Q_j$ and $Q_o$, respectively. The number of electrons that tunnel through the junction along the direction depicted in the figure is represented by $n$. The free energy of this system is expressed by

$$F(n) = \frac{Q_j^2}{2C_j} + \frac{Q_o^2}{2C_o} - qV \tag{1}$$

where the first two terms of the right hand represent static electric energy of the junction and the capacitor and the third term represents the work by the voltage supply with a voltage $V$.



Figure 1. Single-electron box (a) and its characteristics (b).

The following conditions must be satisfied:

$$\frac{Q_i}{C_i} + \frac{Q_o}{C_o} = V \tag{2}$$

$$-Q_i + Q_o = en \tag{3}$$

$$q = en + Q_i \tag{4}$$

where the initial charge at the conduction island is assumed to be zero. From Eq. (1) and Eq. (4), one obtains

$$F(n) = \frac{1}{2C_\Sigma}(en)^2 - \frac{C_o}{C_\Sigma} Ven \tag{5}$$

where $C_\Sigma \equiv C_i + C_o$ is the total capacitance at the conduction island, and the term unrelated to $n$ is omitted.

When $V$ increases, if $F(n + 1) < F(n)$, then an electron tunnels the junction. Using Eq. (5), this condition is given by

$$V \sim \frac{1}{C_o}\left(ne - \frac{e}{2}\right) \tag{6}$$

The number of electrons stored in the island as a function of $V$ is shown in Fig. 1b.

See Ref. [4] for a more detailed explanation, including temperature dependence.

## 2.3. Single-Electron Transistor

A well-known and useful single-electron device is a single-electron transistor (SET) [5]. One can construct a conduction island by one normal capacitor and two tunnel junctions, as shown in Fig. 2a and b. The normal capacitance $C_g$, which is called a gate capacitance, and two tunnel junction capacitances $C_1$ and $C_2$ are assumed to have charges $Q_g$, $Q_1$, and $Q_2$, respectively. The number of electrons that tunnel through junctions 1 and 2 along the directions depicted in the figure are represented by $n_1$ and $n_2$. The total free energy of this system is expressed by

$$F(n_1, n_2) = \frac{Q_1^2}{2C_1} + \frac{Q_2^2}{2C_2} + \frac{Q_g^2}{2C_g} - qV_b - Q_gV_g \tag{7}$$

where the first three terms of the right hand represent static electric energy of the junctions and the gate and the fourth and fifth terms represent the work by the bias and gate voltage supplies. The following conditions must be satisfied:

$$\frac{Q_1}{C_1} + \frac{Q_2}{C_2} = V_b \tag{8}$$

$$\frac{Q_g}{C_g} + \frac{Q_2}{C_2} = V_g \tag{9}$$

$$-Q_1 + Q_2 - Q_g = e(n_1 - n_2) \tag{10}$$

$$q = en_1 - Q_1 \tag{11}$$

where the initial charge at the conduction island is assumed to be zero. From Eq. (7) and Eq. (11), one obtains

$$F(n_1, n_2) = \frac{1}{2C_\Sigma}[C_g V_g + e(n_1 - n_2)]^2 - en_1 \frac{C_2 + C_g}{C_\Sigma} V_b - en_2 \frac{C_1}{C_\Sigma} V_b \tag{12}$$

where $C_\Sigma \equiv C_1 + C_2 + C_g$ is the total capacitance at the conduction island and the terms unrelated to $n_1$ and $n_2$ are omitted.

The conditions under which the Coulomb blockade is effective at both junctions are as follows:

$$F(n_1 \pm 1, n_2) > F(n_1, n_2) \tag{13}$$

$$F(n_1, n_2 \pm 1) > F(n_1, n_2) \tag{14}$$

thus,

$$\frac{1}{C_2 + C_g}\left(Q_0 - \frac{e}{2}\right) < V_b < \frac{1}{C_2 + C_g}\left(Q_0 + \frac{e}{2}\right) \tag{15}$$

$$-\frac{1}{C_1}\left(Q_0 + \frac{e}{2}\right) < V_b < -\frac{1}{C_1}\left(Q_0 - \frac{e}{2}\right) \tag{16}$$

where $Q_0 \equiv C_g V_g + e(n_1 - n_2)$. These conditions define diamond regions at the $V_g$-$V_b$ plane, as shown in Fig. 2c. When one changes $V_g$ with a constant $V_b$ along the dashed line shown in Fig. 2c, the operational point enters and exits the Coulomb blockade regions repeatedly, and thus the current from the voltage supply $V_b$ changes periodically. This means that gate voltage $V_g$ can modify the current, and the SET can operate as a field-effect transistor (FET), such as MOS-FETs. The periodicity of current modulation by $V_g$ is a unique property of SETs, unlike MOS-FETs. By changing the bias point, a SET can have a different characteristics (i.e., complementary operation as in CMOS circuits can be achieved by using single structure SETs [5, 6]). Therefore, various CMOS-like logic circuits composed of SETs have been proposed. Such CMOS-like circuits will be described in Section 5.1.

It is noted that switching can achieved by a very small charge $\Delta Q$ (shown in Fig. 2c), which can be much less than the elementary charge $e$. This means that a SET can be used as a charge sensor with a very high sensitivity [2, 7].

## 2.4. Single-Electron Turnstile and Single-Electron Pump

Unlike the SET described above, a single-electron turnstile [8] and single-electron pump [9] are electron-transfer devices that can transfer electrons one by one in a desired direction. The operation of these circuits can also be understood by a similar analysis, such as described above.

The single-electron turnstile consists of four tunnel junctions and a gate terminal, as shown in Fig. 3a. The operation principle is described using a diagram that defines Coulomb blockade regions, as shown in Fig. 3b. By setting voltages $U_g$ and $V$, the operation point of the circuit can set at point $A$ in Fig. 3b. If the operation point is moved to point $B$, the



(a)                                    (b)

Figure 3. Single-electron turnstile (a) and its characteristics (b).

(a)                                                          (b)

Figure 4. Single-electron pump (a) and its characteristics (b).

number of excess electrons at the center node $n$ increases by one. Then, if the operation point is returned to point $A$, $n$ decreases by one. This process means that electrons move from $C_4$ to $C_1$ one by one.

A single-electron pump is another type of electron-transfer device. The equivalent circuit and the diagram that explains the operation are shown in Fig. 4. Using two-phase voltage signals, the phase transition shown in Fig. 4b is realized. The numbers of electrons at nodes $n_1$ and $n_2$ are indicated in Fig. 4b by $(n_1, n_2)$. According to the phase combination of two-phase voltages, the clockwise or counterclockwise transition is achieved: $(0,0) \rightarrow (0,1) \rightarrow (1,0) \rightarrow (0,0)$ or $(0,0) \rightarrow (1,0) \rightarrow (0,1) \rightarrow (0,0)$. As can be seen in the latter transition, the unique feature of this circuit is that it can pump up electrons from the node with lower potential to that with higher potential. The name of "pump" is derived from this feature.

## 2.5. Coulomb Repulsion Effect in Quantum Dot Circuits

Coulomb repulsion effects between charges confined in coupled quantum dots can be used for information processing. Using the quantum-dot circuit, a cellular-automaton model can be implemented, called a quantum-dot cellular automaton (QCA).

A cellular automaton (CA) is an information-processing model known as a universal computer. In the usual CA model, the processing units (PUs) having internal binary states (0/1) are arranged in a two-dimensional array, and only neighboring PUs interact. The dynamics are defined in discrete time, and the internal state of each PU at the next time step is determined only by the internal states of the neighboring PUs.

The detail of the QCA model and logic circuits using it will be described in Section 5.2. Other circuits using the Coulomb repulsion effects will be described in Sections 5.8.5, 5.8.6, and 5.9.

# 3. OPEN ISSUES FOR PRACTICAL APPLICATION OF SINGLE-ELECTRON DEVICES

## 3.1. Operation Temperature

One of the drawbacks of single-electron devices and circuits is that they can operate only at very low temperatures (e.g., 30 K for a total capacitance of 0.1 aF). This is because the charging energy in these circuits is directly related to the tunnel junction capacitance.

As described in Section 2.1, to observe the Coulomb blockade phenomenon at room temperature, the capacitance $C$ related to the conduction island must be on the order of 1 aF. However, to guarantee a high-enough reliability in large-scale integrated systems consisting of large number of SETs operated at room temperature, $C$ should be much smaller than this value. The reliability of SET circuits was analyzed on the basis of error rate characteristics [10]. For room-temperature operation, $C$ must be less than 0.01 aF, which corresponds to a capacitance between dots with a diameter less than the atomic scale.

To overcome this problem, another approach to operating a single-electron device at room temperature is described in Section 5.8.6, which actively uses thermal noise for proper operation.

## 3.2. Random Background Charges

Effects of charges randomly induced on isolated islands surrounded by tunnel junctions, as shown in Fig. 5, which are referred to as random background-charge effects, are considered to be a serious problem for circuit design [1, 2]. Random background charges are mainly induced by defects or impurities located within the oxide barriers [6], and these defects cannot be entirely removed.

In single-electron transistors, even if the amounts of charges induced on an island are large, an excess over $\pm e/2$ is compensated for by electrons tunneling to or from other islands or electrodes. Therefore, background charges distribute over a range of $[-e/2, e/2]$ with a uniform probability. However, it has been estimated that the background charge margins in CMOS-like logic circuits are much smaller than $|e/2|$, typically $0.03e$ [11, 12].

Even if almost no islands have any background charge, even a few background charges may cause fatal errors in the operation of a whole circuit in multistage logic circuits. In small-scale circuits, this problem can be solved by individual biasing for each island, but this solution cannot be used for large-scale integrated circuits.

Fortunately, it has been reported that background charges in silicon single-electron devices seem to be smaller than those in metallic single-electron devices [13, 14]. Furthermore, architecture-level solutions to reduce the background-charge problem will be proposed in the following sections.

## 3.3. Higher-Order Tunneling (Cotunneling)

Higher-order tunneling, which is often called cotunneling [15] or macroscopic quantum tunneling [16–18], is a phenomenon in which plural electrons are involved, as shown in Fig. 6. The probability that the $n$th order tunneling occurs is proportional to $(R_Q/R_T)^n$. As described in Section 1, single-electron devices and circuits should be designed to guarantee that $R_T \gg R_Q$. Therefore, higher-order tunneling seldom occurs, but second-order tunneling becomes the major leakage factor in Coulomb blockade regions where ordinary first-order tunneling is prohibited. Thus, higher-order tunneling arises as a problem in devices that use charge confinement by Coulomb blockade phenomena such as single-electron memories. An effective method for reducing the effect of higher-order tunneling is using multitunnel junctions.

## 3.4. Injection of a Single Electron to Nanodot Arrays

In quantum-dot (or nanodot) circuits described in Section 2.5 and some circuits and applications that will be described in Section 5, introduction of a single electron to a nanodot array is an essential assumption.

One might consider using Fowler-Nordheim (F-N) tunneling when a high-voltage pulse is applied to the injection node. Because of statistical tunneling processes, a short high-voltage pulse will not ensure a single electron injection. When a short high-voltage pulse is applied to many nanodot arrays, a single electron can be injected into some arrays but not into others. Only the arrays with an electron will operate properly. The success rate of single-electron injection may be a crucial issue.

To ensure single-electron injection, the Coulomb blockade effect could be used. However, the blockade would hold only while the electron stayed at the injected dot. When it moves to one of the adjacent dots, another electron can tunnel in it. This tunneling will depend on the tunnel resistance to the adjacent dots. This is another time constant issue that must be

isolated              background
island                charge



Figure 5. Random background charge issue.

(a) 1st-order tunneling (Coulomb blockade)    (b) 2nd-order tunneling (cotunneling)

**Figure 6.** Expamle of cotunneling process.

taken into account in the circuit design. To prevent the injected electron from moving to the other dot immediately. an additional electrode to hold the potential may be necessary near the injection dot. Some experimental results have been reported about restricted electron injection in nanocrystalline floating-dot MOSFET devices [19].

# 4. OVERVIEW—INFORMATION PROCESSING ARCHITECTURES FOR SINGLE-ELECTRON SYSTEMS

## 4.1. Information-Processing Models for Single-Electron Systems

One of the goals in single-electron circuit technology is the development of computing systems that can perform information processing by using single-electron phenomena. There are two basic ways that we can use single-electron phenomena for information processing. One is to construct transistor-like devices based on single-electron phenomena and imitate existing silicon LSI systems, using the devices as analogs of MOS-FETs. The problem with this approach is that integrated circuits composed of new transistor-like devices would have a hard time competing with the original silicon LSI, which is a product of well-established, mature technologies. The more promising approach is to reconsider the procedure for implementing information processing and take up a way different from that of existing LSIs. From this approach, one will be able to construct novel computing devices that make good use of the properties of single-electron phenomena.

Various processing methods for information processing (or procedures for solving a given problem) known for now are shown in Fig. 7. The procedure used in existing LSI systems is as follows (shown by boldface in the figure).

1. Devise an algorithm for processing given information
2. Execute each step of the algorithm under von Neumann–type computer architecture
3. In the execution, express the algorithm in the form of a sequence of Boolean operations on digital functions
4. Describe each operation in terms of Boolean expression (a combination of logic operators AND, OR, and NOT)
5. Implement the Boolean expression with binary logic gates made from MOS FETs.

This conventional procedure, Neumann-Boolean computing, has become the mainstream in information processing. However, there are many potential ways of processing that are different from Neumann-Boolean computing, as shown in Fig. 7. Let us call these processing ways *unconventional computing*. The process of unconventional computing is much more sophisticated than that of Neumann-Boolean computing; consequently, can provide the possibility of solving problems that are intractable for Neumann-Boolean computing. Unfortunately, unconventional computing is not easy to implement on LSIs because the CMOS



**Figure 7.** Various methods of information processing.

transistor gate—a Boolean logic device by nature—is the only device we can use at this time for constructing practical LSIs. However, we will be able to create novel, functional information-processing LSIs based on unconventional computing if we can develop single-electron devices suitable for implementing non-Neumann, non-Boolean operations.

The theme of this chapter is to consider developing such single-electron processing systems based on unconventional computing. Challenges to this interesting subject have already been reported with leading examples. Section 5 will review these vanguard attempts at developing new-paradigm systems for information processing.

## 4.2. Design Strategy for Single-Electron Functional Circuits

A strategy for overcoming the difficulties described in Section 3 and for achieving large-scale integration of single-electron devices is summarized in Fig. 8 [20].

Single-electron devices should be used for massively parallel processing with a huge number of devices because of their large packing density and ultralow power dissipation. The operation speed is improved by parallel processing. Adopting a few logic stages, a small fanout, and regularity or repeatability in the circuit architecture overcomes the interconnection difficulty and lowers the background-charge effects.

Parallel operation and such circuit architecture may make deep logic processing more difficult. Therefore, the use of ultrasmall CMOS devices is essential. Single-electron devices should be used for simple functional circuits with few logic stages, and ultrasmall CMOS devices are used for multistage logic circuits.

To overcome the design difficulty, it is a good idea to use large output capacitance as a buffer for each circuit component. Here, "large" means a capacitance value that can store a few tens of electrons. Because of this output buffer, the operation of the circuit component is not affected by the following stage circuit, and thus modularized circuit design is applicable. The output capacitance is also used for an interface between SET and CMOS circuits. The information generated by massively parallel processing in SET circuits is collected and integrated into the output buffer and is transferred to the CMOS circuits. The gate capacitance of an ultrasmall MOS device can be used as the buffer capacitance.

To reduce the sensitivity to capacitance and background charges, new information-processing concepts and models are required in addition to solutions regarding the averaging of information or redundancy configurations.

The functionalities created by the single-electron operations are as follows:

- Multiinputs circuits using multigate configuration in single-electron devices
- Nonmonotone input-output functions resulting from multipeak (oscillatory) characteristics of single-electron operation
- Stochastic operation resulting from stochastic tunneling phenomena
- Energy minimization.

Various circuits and architectures using these functionalities are reviewed in the next section.



Figure 8. Single-electron circuit system design strategy. Reprinted with permission from [20], T. Morie et al., J. Nanosci. Nanotechnol. 2, 343 (2002). © 2002, American Scientific Publishers.

## 5. SINGLE-ELECTRON FUNCTIONAL DEVICES AND CIRCUITS

### 5.1. CMOS-Like SET Logic Circuits

#### 5.1.1. CMOS-Like SET Logic

Using SETs described in Section 2.3, CMOS (complementary metal-oxide-semiconductor)-like circuits can be constructed [6]. Because of periodical characteristics between the gate voltage and the tunneling current in a SET, p-type and n-type MOS transistor operations can be realized by setting different bias voltages at the additional gate (back-gate) electrode, as shown in Fig. 9. Such a circuit can use the same circuit configuration as the present CMOS logic circuits, and thus various studies about this approach have been reported [6, 11, 12, 21]. For example, a CMOS-like SET inverter and an exclusive-OR (XOR) logic gate are shown in Figs. 10 and 11, respectively.

In this case, the value of load capacitance is important. To construct multistage logic circuits, a large load capacitance is required. If small capacitance is used, the whole multistage circuits must be considered as one quantum system, and therefore the circuit design will be much different from the conventional CMOS circuit design. In the following text, how to design such circuits, including plural SETs, is described.

#### 5.1.2. CMOS-Like SET Circuit Design

A SET circuit has a number of island nodes that are interconnected by means of tunnel junctions. Its internal state is determined by the configuration of electrons (i.e., the pattern in which the excess electrons are distributed among the nodes) and is expressed by a set of the numbers of excess electrons on the nodes. The circuit changes its state through tunneling in response to the input and thereby changes its output voltage as a function of the input.

A SET circuit changes its state to decrease its free energy; hence, the circuit operates as an organic whole. Therefore, any SET circuit has to be designed taking into consideration the global stability of the whole circuit. Because a SET circuit has complex internal states, a "guide map" is needed to grasp the overall situation of the circuit. The guide map for this purpose is known as the stability diagram.

The stability diagram is the diagram that depicts the internal states of a SET circuit in a multidimensional space of circuit variables. Its concept is as follows. The bias voltage setting criteria necessary to maintain a circuit in a stable state are given as a combination of inequalities [e.g., Eqs. (15) and (16)], each of which indicates a condition for maintaining the circuit energy minimum. The inequalities involve all the circuit parameters as variables, so the stable bias region for a given state forms a hypersolid surrounded by a number of hypersurfaces on a variable space.

For instance, a Tucker's inverter in Fig. 10 has 11 variables—two voltage variables $V_{in}$ and $V_{dd}$ and nine capacitance variables. Accordingly, the stable bias region is drawn in an 11-dimensional space. Each hypersurface corresponds to a threshold condition for tunneling through a junction in one direction; therefore, if the number of tunnel junctions is $N$, then $2N$ hypersurfaces exist. For a different state of the circuit, there is a different set of hypersurfaces that determines a hypersolid of a stable region.

The stability diagram is a map that illustrates all the hypersolids that represent all possible states of the circuit. Looking at a stability diagram, we can see the changes of the internal



Figure 9. Complementary operation of SET: (a) device structure, (b) p-type MOS-like operation, (c) n-type MOS-like operation.

Figure 10. CMOS-like inverter using SETs, which is often called Tucker's inverter.

states, the stability, and the output values, as functions of the circuit variables. To illustrate a stability diagram on a sheet of paper, we have to reduce the diagram to a two-dimensional representation. For this purpose, we select two of the variables and assume the others to be constant. In general, it is convenient in designing a circuit to choose the input voltages for the circuit as the variables. If the circuit has one and only one input, it is advisable to use as the other variable the voltage of another voltage source in the circuit. In such a two-dimensional stability diagram with two voltage variables, the hypersurfaces and the hypersolids are reduced to lines and polygons. For examples of a stability diagram with two voltage variables, see Figs. 2c, 3b, 4b, 9b, and 9c.

**5.1.2.1. Drawing Stability Diagrams**  A stability diagram can be calculated analytically for a simple SET circuit composed of a few junctions. However, a circuit of greater complexity is difficult to calculate on paper; therefore, simulation by computer is needed. The way of drawing the stability diagram for a given circuit is as follows. First, a current state of



Figure 11. CMOS-like XOR gate using SETs.

the circuit (i.e., a current set of the numbers of excess electrons on nodes) is assumed, and then a set of values of the circuit variables is also assumed. After that, the energy change of the circuit is calculated for each possible tunneling. If all the energy change is incremental, no tunneling occurs, and it can be considered that the current state is stable under the set of circuit variables. If the energy change is reduced for one or more tunnelings, then it can be judged that the current state is unstable under the set of circuit variables. The above procedure is named a trial. After the judgment for a trial, each variable is changed slightly, and then another trial calculation is repeated for the new set of circuit variables. By scanning the whole variable space, the stability diagram can be drawn for the assumed state of the circuit. The same sequence is repeated for the other states. As the result of these iterations, the whole stability diagram can be obtained.

As a sample circuit, let us take up the Tucker's inverter in Fig. 10 and assume for simplicity that $C_{j1} = 1$ aF, $C_{j2} = 2$ aF, $C_1 = 8$ aF, $C_2 = 7$ aF, and $C_{out} = 24$ aF. Then the stability diagram can be drawn on a two-dimensional $V_{in}$-$V_{dd}$ plane ($V_{in}$: input voltage, $V_{dd}$: power voltage), as shown in Fig. 12. The circuit has three island nodes (L, M, and N), and its internal state is expressed by a set of the numbers ($l, m, n$) of excess electrons stored on the three nodes. The stable regions take on various configurations. Most regions overlap with one another. In a set of electron numbers ($l, m, n$), $n$ is a main factor of determining the output voltage, and $l$ and $m$ change the output voltage slightly.

In the range of the voltage variables of Fig. 12, the numbers of $m = -1$, 0, and produce output voltages of about 6 V, 0 V, and −6 mV, whereas $l$ and $n$ modify the output voltage by 0.7 mV or less in the same way. In Fig. 12, the approximate output voltage for each state is shown by putting the letters H, L, or LL before the electron-number set; for example, H(0, −1, 0) indicates a state of high-output voltage (about 6 mV), L(0, 0, 1) of low-output-voltage (about 0 mV), and LL(0, 1, 0) of much lower output voltage (about −6 mV).

The stability diagram offers an insight into the functions that can be obtained from the circuit. If we set $V_{dd}$ to 7 mV and operate the circuit on segment A-B, we can use the circuit as an inverter with the transfer curve shown in Fig. 13a. The inverter has an unstable region on its transfer curve for intermediate values of the input, but this is not a problem for use of binary-logic applications. In contrast, if we set $V_{dd}$ to 5.8 mV to operate the circuit on segment C-D, we can use the same circuit as a Schmidt trigger. In the bistable region denoted by H(0, −1, 0) and L(0, 0, 0), the circuit maintains the same state as just before entering this region. Therefore, the threshold voltage of the circuit is higher for an increasing input and



Figure 12. Stability diagram for the circuit shown in Fig. 10. For the capacitance parameters, see the text. The shaded region is an unstable region where tunneling occurs repeatedly and the circuit state alternates between two or more different states.

**Figure 13.** Transfer curves of the circuit shown in Fig. 10. The output voltage is plotted as a function of the input voltage for two different values of $V_{dd}$: (a) $V_{dd} = 7$ mV (the circuit operates as an inverter) and (b) $V'_{dd} = 5.9$ mV (the circuit operates as an inverting Schmidt trigger circuit).

lower for a decreasing input; this condition results in a hysteresis characteristic, as shown in Fig. 13b. The width of the hysteresis region can be set up as desired by adjusting $V_{dd}$.

**5.1.2.2. Designing NAND Logic Gates**  A NAND logic gate can be constructed as follows: Let us assume a circuit configuration for a two-input NAND gate by analogy with a CMOS gate, based on the model of the SET inverter circuit proposed by Tucker [6]. The configuration, illustrated in Fig. 14, consists of a pull-up tree (a parallel connection pair of SET transistors) and a pull-down tree (a merged series pair of SET transistors), with an output capacitance $C_0$. The circuit has seven tunnel junctions (the junction capacitances are $C_{j1}$ through $C_{j7}$), four input capacitors ($C_1$ through $C_4$), and four bias capacitors ($C_5$ through $C_8$), with a power voltage $V_{dd}$. The circuit accepts two input voltages ($V_1, V_2$) and produces the corresponding output voltage $V_{out}$.

The next work is to determine an optimum set of capacitance parameters so that the circuit will produce a NAND logic operation. The optimum parameter set can be determined in theory by calculating and scrutinizing the stability diagram of the circuit, but this is impossible in practice because this circuit has too many circuit variables. The circuit has 19 circuit variables (16 capacitance parameters plus 2 input voltages and 1 power voltage), so the stability diagram will be drawn in a 19-dimensional space of circuit variables. Empirically,



**Figure 14.** CMOS-like SET-NAND gate.

if a diagram has eight or more dimensions, it takes an enormous amount of time to search the whole diagram (even by computer), and it is virtually impossible to detect an optimum point in the diagram. Therefore, the number of circuit variables has to be reduced.

One way to reduce the number of circuit variables is to divide a circuit into smaller sub-circuit, as Tucker did in designing a SET inverter [6]. Here, let us divide the circuit into two subcircuits, the pull-up tree and the pull-down one, and analyze each tree circuit separately. For a NAND logic operation. the tree circuits are required to have the following character-istics: First, when both input voltages are $V_{dd}$ (i.e., $V_1 = V_2 = V_{dd}$), the pull-up tree is always stable (switch "off"), and the pull-down tree is unstable (switch "on") if $V_{out} > 0$ and stable if $V_{out} = 0$. (If this specification is satisfied, the output node will be discharged through the pull-down tree to reach a zero voltage, so the output will settle down at $V_{out} = 0$.) Second, when either or both inputs are 0, the pull-down circuit is always stable, and the pull-up circuit is unstable for $V_{out} < V_{dd}$ and stable for $V_{out} = V_{dd}$. (If this is satisfied, the output node will be charged up to reach a voltage of $V_{dd}$, so the output will become stable at $V_{out} = V_{dd}$.)

For further reduction of the number of variables, let us make three assumptions: the pull-up tree consists of two identical SET transistors (i.e., $C_{j1} = C_{j3}$, $C_{j2} = C_{j4}$, $C_1 = C_2$, and $C_5 = C_6$); in the pull-down tree, two input capacitances are equal to each other ($C_3 = C_4$) and two bias capacitances are also equal ($C_7 = C_8$); and two junction capacitances $C_{j5}$ and $C_{j6}$ are equal to each other ($C_{j5} = C_{j6}$). The last assumption for $C_{j5}$ and $C_{j6}$ is groundess, but we still assume it for reducing the variables. Through these assumptions, the number of capacitance parameters has been reduced to four for each tree circuit.

Our work is to determine an optimum set of capacitance parameters for each tree circuit that produces the required pull-up and pull-down characteristics. This can be performed by scrutinizing the stability diagram of each tree circuit. To simplify the task, we here adopt an approximation: In calculating the stability diagrams, we consider the output voltage $V_{out}$ as a variable. Strictly speaking, this is inaccurate, because $V_{out}$ is dependent on other variables. However, if the output capacitance is large, $V_{out}$ can be considered a variable in calculating the stability diagrams. This is because, for large-output capacitance, a value of $V_{out}$ is kept almost constant during the change of the circuit state induced by one-electron tunneling. For simplicity, we set the power supply voltage $V_{dd}$ at 6.67 mV (no special reason for this value; any other positive voltage can be assumed). We draw the stability diagram of each tree circuit in the seven-dimensional space of circuit variables (four capacitance parameters plus three voltages, $V_1$, $V_2$, and $V_{out}$) and then search the diagrams for an optimum point that produces the required characteristics of the tree circuits. For this power supply voltage, several optimum sets of capacitance parameters can be found. A sample set is $C_{j1} = C_{j3} = 1$ aF, $C_{j2} = C_{j4} = 2$ aF, $C_{j5} = C_{j6} = 2$ aF, $C_{j7} = 1$ aF, $C_1 = C_2 = 5$ aF, $C_3 = C_4 = 3$ aF, $C_5 = C_6 = 7.4$ aF, and $C_7 = C_8 = 7.7$ aF, ($V_{out} = 6.67$ mV, as mentioned above).

For reference, a part of the calculated stability diagram is illustrated in Fig. 15a–f on a plane of the two voltage variables ($V_{out}$ and $V_1$ or $V_2$). The output capacitance $C_0$ can be set at any value on the condition that it is sufficiently larger than the other capacitance parameters and that the value of $C_0 V_{dd}$ is a multiple of the elementary charge. An example value of $C_0$ is 240 aF.

**5.1.2.3. Constructing Adders**    A full adder can be constructed with the NAND gates, as shown in Fig. 16. The output capacitance is set to 240 aF. In each NAND gate, the input capacitance is sufficiently smaller than the output capacitance, so several gates can be connected to a preceding gate without affecting the operation of the preceding gate. The add operation of the full adder, for a tunnel resistance of 200 kΩ, is shown in Fig. 17. The figure depicts the waveforms for adder inputs $A$ and $B$, a carry input $C_{in}$, an adder output $Sum$, and a carry output $C_{out}$. The average add time is 10 ns.

## 5.2. QCA-Based Logic Circuits

As described in Sec. 2.5, using the Coulomb repulsion effect in coupled quantum dots, quantum-dot cellular automata (QCA) can be constructed [22].

Let us consider quantum-dot arrays as shown in Fig. 18 [23]. Here, four or more quantum dots are coupled and form a unit cell, in which two electrons are confined. The two electrons

**Figure 15.** Operation diagram of CMOS-like SET-NAND gate.

can exist at the four vertexes (and the center position in cells A and B) of the square shape cell, but at the steady state only the two diagonal arrangements of the electron positions (polarization) are possible because of the Coulomb repulsion effect. Therefore, they can be assigned logical "1" and "0", respectively. There are various quantum-dot configurations for constructing a cell, as shown in Fig. 18. The polarization response by cell–cell interaction depends on the number of tunneling paths and the distance between dots. By the quantum-mechanical analysis, cells A and B exhibit better polarization response than cells C and D [23].

As the first step for constructing logic circuits, binary-information transmission lines can be realized by capacitively-coupled one-dimension cells. QCA lines (or wires) with various configurations are shown in Fig. 19 [24–26].

Using QCA, basic logic gates can be constructed [25]. A digital inverter using QCA is shown in Fig. 20. A QCA majority logic gate with three inputs is shown in Fig. 21a. Using this gate, AND and OR gates can be constructed, as shown in Fig. 21c. There are many reports about circuit architecture, design, and fabrication technology for QCA circuits [22, 27–31].

## 5.3. Single-Electron Logic Circuits Based on the Binary Decision Diagram

### 5.3.1. Representing Digital Functions by Binary Decision Diagrams

The binary decision diagram (BDD) is a way of representing digital functions by using a directed graph instead of a Boolean expression [32, 33]. It can represent any digital function



**Figure 16.** Full-adder circuit.

**Figure 17.** Timing diagram of the full-adder circuit.



**Figure 18.** Quantum-dot arrays for quantum cellular automata (QCA). Circles indicate quantum dots, and the gray dots indicate those occupied by an electron. The solid line connecting two dots indicates a tunneling path of the electrons. All cells A to D can express logical 0 and 1 by different electron position states. Adapted with permssion from [23]. P. D. Tougaw et al., *J. Appl. Phys.* 74, 3558 (1993). © 1993, American Institute of Physics.



**Figure 19.** Interconnection lines using QCA. (a) bent line with a right angle. (b) branch lines. (c) crossing lines. Adapted with permission from [25]. P. D. Tougaw and C. S. Lent, *J. Appl. Phys.* 75, 1818 (1994). © 1994, American Institute of Physics.

**Figure 20.** QCA inverter. Adapted with permission from [25]. P. D. Tougaw and C. S. Lent. *J. Appl. Phys.* 75, 1818 (1994). © 1994, American Institute of Physics.

and provides a concise representation for most digital functions encountered in logic-design applications.

As an example, consider the four-variable digital function represented by the Boolean equation given in Fig. 22a. This function can also be represented by the BDD shown in Fig. 22b. A BDD is a graph consisting of many nodes and two terminals, with each node labeled by a variable; in this example, each node is represented by a circle labeled by a variable $X_i (i = 1, 2, 3, 4)$, and a terminal is represented by a square labeled 0 or 1. Each node in a BDD has two branches labeled 1 and 0 and is connected to its adjacent nodes with the branches.

In determining the value of the function for a given set of variable values, we enter at the root and proceed down to a terminal. At each node, we follow the branch corresponding to the value of the variable; that is, we follow the 1 branch if $X_i = 1$ and the 0 branch if $X_i = 0$. For a given set of variable values, there is one and only one path from the root to either terminal. The value of the function is equal to the value of the terminal we reach; the function is 1 if we reach the 1-terminal and 0 for the 0-terminal. Some other examples of the BDD are shown in Fig. 23, together with the corresponding Boolean equations.

### 5.3.2. Implementing a BDD with Single-Electron Devices

As described in the previous section, when an input (a set of variable values) is presented to a BDD, a path through the BDD can be traced from the root to either terminal. In this condition, a signal that is injected in the root can travel along the path to reach a 1- or 0-terminal. Therefore, we can determine the value of the function by observing which terminal the signal reaches. This signal is called a messenger.

A BDD is composed of many identical interconnected nodes, so the node is the unit element of a BDD. The function of the element is simple two-way switching controlled by an input variable. To implement this function, we can use many physical effects that change the course of a traveling messenger (electron, photon, single flux quantum, etc.) in response to an input. The advantage of BDD logic systems is that they can make use of even a physical effect that is useless for constructing transistor-like devices.

Asahi and others proposed implementing the BDD node function by means of single-electron circuits [34–36]. Their unit element, a BDD device, illustrated in Fig. 24, consists

**Figure 21.** QCA majority logic gate. (a) Configuration. (b) logic gate notation. (c) AND and OR logic realization by the QCA majority logic gate. Adapted with permission from [25]. P. D. Tougaw and C. S. Lent. *J. Appl. Phys.* 75, 1818 (1994). © 1994, American Institute of Physics.

(a) Boolean equation

(b) Binary decision diagram

$$f = (\overline{X1}\,X2 + X.\ \overline{X2}\,)\,\overline{X3}\,\overline{X4}$$
$$+ (\overline{X1}\,\overline{X2} + X1.\,X2\, )\,X3.X4$$
$$+ X3.\overline{X4}$$



**Figure 22.** Representation of digital functions.

of four tunnel junctions $(J_1, \ldots, J_4)$ and three capacitors $(C_1, \ldots, C_3)$ and is driven by a voltage clock $(\phi)$. It has entry branch $(A)$ and two exit branches $(D, E)$. Voltage input $X$ (and its complement $\overline{X}$), specifying the value of a variable, is applied to island $B$ (and $C$) through capacitor $C_2$ (and $C_3$); $X$ is an appropriate positive voltage (and $\overline{X}$ is a negative one) if the variable value is 1, and $X$ is negative ($\overline{X}$ is positive) if the variable value is 0.

The node device receives a messenger electron from a preceding device through the entry branch and sends the electron to a following device through either exit branch that corresponds to the binary value of the input. The path of the electron transport is $A \to B \to D$ (the 1 branch) if input $X$ is positive and $A \to C \to E$ (the 0 branch) if input $X$ is negative.

## 5.3.3. Constructing BDD Logic Circuits

A logic circuit can be constructed by connecting many BDD devices in a cascade manner to build the tree of a BDD graph, as illustrated in Fig. 25. Each BDD device corresponds to a node of the graph and operates as a two-way switch for the transport of a messenger electron. The entire system is composed of the BDD graph circuit, an electron injector, and an output circuit, as illustrated in Fig. 26. At the start of operation, a messenger electron is injected from the electron injector into the root node and then transferred by four clocks $(\phi_0, \ldots, \phi_3)$ through the BDD graph circuit to a terminal. The messenger electron travels

(a) $f = X1+X2+X3+X4$     (b) $f = (X1+X2)(X3+X4)$     (c) odd parity



**Figure 23.** Examples of BDD representation.

**Figure 24.** Unit element for single-electron BDD logic circuits. Its function is to provide two-way switching for electron transport. Reprinted with permission from [35], N. Asahi et al., *IEEE Trans. Electron Devices* 44, 1109 (1997). © 1997. IEEE.

along the BDD path specified by a given set of voltage inputs ($X_1, X_2, \ldots$). The value of the logic is determined by observing which terminal the messenger electron reaches. The output circuit detects the arrival of a messenger electron at the 1-0 terminals and produces the corresponding binary voltage output. The messenger electron ejected from the terminal is returned to the root node through the feedback loop to be used for successive logic operations. Two or more messenger electrons can be used in a BDD circuit; we can put a messenger electron on every transfer stage (a subcircuit unit that is driven by a set of four clock pulses) to produce pipelined operation for increasing throughput of processing.

## 5.3.4. Logic Operation of BDD Circuits

This section illustrates the logic operation of BDD circuits, with an example of computer-simulated results for a 4-bit adder reported in Ref. [36]. The adder accepts two 4-bit adder inputs (an addend and an augend) and produces the corresponding 4-bit sum output and



**Figure 25.** Unit devices cascaded to build the tree of a BDD graph. Dashed lines represent a path of a messenger-electron transfer. Reprinted with permission from [35], N. Asahi et al., *IEEE Trans. Electron Devices* 44, 1109 (1997). © 1997. IEEE.

**Figure 26.** Schematic view of a BDD-circuit block, consisting of a BDD circuit (a tree circuit and terminals), an output circuit, and an electron injector. Reprinted with permission from [35], N. Asahi et al., *IEEE Trans. Electron Devices* 44, 1109 (1997). © 1997, IEEE.

1-bit carry output. The two adder inputs are hereafter represented by binary numbers $a_3a_2a_1a_0$ and $b_3b_2b_1b_0$, the adder output by $s_3s_2s_1s_0$, and the carry output by $c_3$, where the value of each element $a_3$ through $c_3$ is either 1 or 0.

The operation of the adder can be represented by using the set of BDD graphs illustrated in Fig. 27. Each bit of the outputs ($s_0$ through $s_3$ and $c_3$) is produced by the corresponding BDD graph, which contains the value of the input bits ($a_0$ through $a_3$ and $b_0$ through $b_3$) as node variables. In each node, the exit branch on the right side shows a 1 branch and the exit branch on the left shows a 0 branch.

To illustrate the single-electron circuits that implement the BDD graphs, Fig. 28 defines symbols for two devices—the BDD device and a buffer. The buffer is a subcircuit consisting of a tunnel junction and a capacitor. Its function is to set up a dummy node and holds a messenger electron for one clock period—an indispensable timing element for constructing the single-electron BDD circuits.



**Figure 27.** A set of BDD graphs representing 4-bit addition.

(a) BDD device  (b) Buffer

**Figure 28.** Symbols of the BDD device and a buffer.

The circuits designed for the adder are illustrated in Fig. 29. In each circuit, a root node is indicated by the boxed word "Root," and terminal nodes are indicated by the boxed numerals 1 and 0. The configuration of each circuit was obtained by simply replacing the nodes in the BDD graphs of Fig. 27 with BDD devices. Factors to keep in mind in designing the circuit include the following:

1. To transfer and circulate the messenger electron in the circuit, the four-phase clock ($\phi_0, \ldots, \phi_3$) is applied to node devices and buffers. The phase shift of the clock is $\phi_0 = 0, \phi_1 = -\pi/2, \phi_2 = -\pi, \phi_3 = -3\pi/2$. The bit signals of the adder inputs are applied in sequence to the BDD devices such that the bit signal for a BDD device is applied synchronously with the clock pulse for the consecutive BDD device (or the consecutive buffer).

2. In actual systems, an input-bit signal will be applied to all the BDD circuits simultaneously with the clock. For successful operation in such a situation, the buffers (dummy nodes) have to be set up on appropriate points on the path to ensure that a messenger electron will arrive at each BDD device at the correct time.

3. Two or more messenger electrons can be used in a BDD circuit. Putting a messenger electron on every transfer stage (a subcircuit that is driven by a set of four clock pulses) produces a pipelined operation that improves the processing throughput. In this instance, three messenger electrons can be put in each BDD circuit.

Figure 30 shows the simulated operation of the designed BDD adder for several sets of input data. The parameters used in this example are junction capacitance = 10 aF and tunnel resistance = 100 k$\Omega$ for each tunnel junction, capacitance = 10 aF for each capacitor, and temperature = 0 K. Three messenger electrons were put in each BDD circuit. The figure plots the waveforms for clock pulse $\phi_0$ (the other clocks are omitted), input-bit pulses



**Figure 29.** Circuit configuration of the 4-bit adder. Reprinted with permission from [36], N. Asahi et al., *IEICE Trans. Electron.* E81-C, 49 (1998). © 1998, IEICE.

**Figure 30.** Operation of the 4-bit BDD adder. A set of input and output data is marked by a gray line. Reprinted with permission from [36], N. Asahi et al., *IEICE Trans. Electron.* E81-C, 49 (1998). © 1998, IEICE.

($a_3$ through $b_0$), and output charges ($s_0$ through $c_3$, the charges on the 1-terminal nodes of the four BDD circuits, normalized to the electron charge). The adder produces an output data flow in response to the input data flow, as a linear systolic array; in the figure, a set of input and output data is marked by a dashed line. One logic operation (from accepting $a_3$ to producing the corresponding output) requires 8 ns, but the pipelined processing can produce an output data every period of 4 ns.

## 5.3.5. Logic Circuits Based on the Shared BDD

As described in Fig. 22, a BDD has a path from the root to either terminal for a given set of variable values. Therefore, we can also determine the value of the digital function by tracing the path upward from the 1-terminal to the root to check whether the 1-terminal is

connected to the root. If the connection is established, the function is logical 1, and if not, the function is logical 0. This upward tracing leads us to another form of single-electron BDD circuits, as shown in Sections. 5.3.6–5.3.8.

In the upward tracing, each node in BDD acts as a switching element with two entry branches (1-branch and 0-branch) and an exit, as shown in Fig. 31a. In each node in BDD graphs (e.g., in Fig. 31b and c), one entry branch corresponding to the value of the variable ($x_i$) is connected to the exit, and the other branch is disconnected (e.g., if variable $x_i$ is logical 1, the 1-branch is connected to the exit and the 0-branch is not connected). To determine the value of the logic function, we check whether the 1-terminal is connected to the root; if the connection is established, the function is logical 1, and if not, the function is logical 0. Most digital systems contain multiple output functions that are closely related to one another, and these functions can be represented by a single graph with multiple roots (one root for each function), as shown in Fig. 31d (a combination of two BDD graphs in Fig. 31b and c). This kind of BDD is called a shared BDD. Figure 32 shows another example of a shared BDD, a representation of 2-bit addition, which has three roots for three output bits (a 2-bit sum output and a 1-bit carry output).

## 5.3.6. Constructing Shared-BDD Circuits Combined with the Upward Tracing

Yamada and others proposed implementing the node function for the upward tracing by means of tunneling gates [37]. The tunneling gate is a tunnel junction with a gate electrode that controls electron transport through the tunnel junction (Fig. 33a). It accepts a binary gate voltage as input and transmits electrons through the junction if the gate voltage is logical 1 (the junction is on) or transmits no electrons if the gate voltage is logical 0 (the junction is off). The tunneling gates can be made with quantum-dot devices, as shown later in Section 5.3.8.

Figure 33b shows the BDD device consisting of two tunneling gates and a ground capacitor joined together at the exit node. A binary gate voltage (and its complement), specifying the value of an input variable $x_i$, is applied to tunneling gate labeled $x_i$ (and gate labeled $\overline{x_i}$) to control the tunneling transport of the junction. If variable $x_i = 1(\overline{x_i} = 0)$, the BDD device transports electrons from the 1-branch to the exit, and if $x_i = 0(\overline{x_i} = 1)$, it transports electrons from the 0-branch to the exit.

Any combinational logic can be implemented by combining the BDD devices to build a BDD graph circuit. In operating the circuit, we apply input gate voltages to the tunneling gates and inject electrons into the circuit from the 1-terminal and then observe at each root whether the electrons flow out or not. As an example, Fig. 34 shows a sample configuration that implements the shared BDD given in Fig. 31c. The circuit is designed on the basis of the following principles:

1. In composing the circuit, the 0-terminal and related branch connections in BDDs are unnecessary and removed.
2. A gradient potential from the 1-terminal to the roots must be appropriately established so that, in every BDD device, electrons can flow from the entry branches to the exit.



Figure 31. Examples of the binary decision diagram (BDD).

augend $a_1$ $a_0$
addend $a_1$ $b_0$

sum: $s_1$ $s_0$
carry: $c_1$



**Figure 32.** Shared-BDD representation for 2-bit addition.



**Figure 33.** BDD device consisting of tunneling gates. Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.
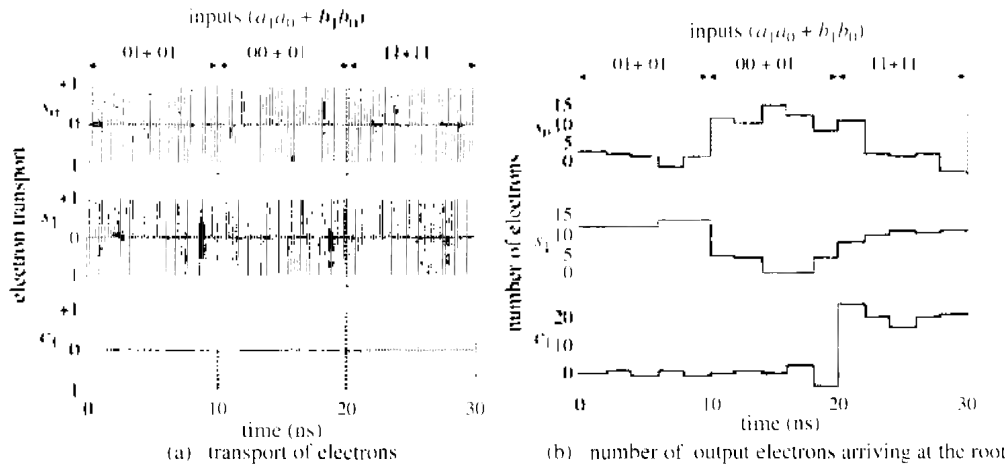


**Figure 34.** BDD circuit implementing the shared BDD graph in Fig. 31(c). The ground capacitance for each node is omitted. Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.

This can be achieved by inserting dummy tunnel junctions in electron paths so that every path from the 1-terminal to a root will have the same number of tunnel junctions.

The circuit accepts two voltage inputs (specifying variables $x_1$ and $x_2$) and produces two corresponding outputs ($f_1$ and $f_2$). The input voltages are applied to the tunneling gates labeled with variables $x_1$, $x_1$, $x_2$, and $x_2$ to control (turn on and off) electron transport through the tunneling gates. The nonlabeled junctions are dummy tunnel junctions that establish a potential gradient through electron paths. Each node has a ground capacitance, but for simplicity the capacitance is omitted in the figure.

To operate the circuit, all the roots are grounded and a negative power voltage is applied to the 1-terminal. Electrons are injected from the power voltage into the circuit and transported toward each root along the paths specified by the variables. The logical value of an output is 1 if electrons can reach the corresponding root, and 0 if they cannot. The circuit is analogous in operation to pass-transistor circuits composed of MOSFETs, and it can be considered to be an ultra-low-power version of a pass-transistor circuit. In the circuit, the electron flow (and therefore power consumption) is regulated by the Coulomb blockade and can be set to a far smaller quantity than can be in ordinary pass-transistor circuits. The number of electrons that flow through the circuit during one logic operation can be reduced to tens or so by adjusting the tunnel junction and the ground capacitances (see Figs. 35 and 36).

The logic operation of the sample circuit is shown in Fig. 35. This is a result simulated with a set of device parameters: tunneling gate junction capacitance = 10 aF, dummy tunnel junction capacitance = 10 aF, ground capacitance of a node = 20 aF, and tunneling gate resistance = 1 M$\Omega$ in the "on" state and 10 G$\Omega$ in the "off" state. The voltage of the 1-terminal was set to 3.5 mV, and the temperature was assumed to be absolute zero. In the figure, four input combinations ($x_1x_2$ = 00, 10, 01, 11) are applied in sequence to the gates of the circuit at 10-ns intervals. Each impulse in the figure signifies the arrival of an electron at a root. Thus, a total of 15 electrons flowed out of root $f_1$ during the period from 10 to 20 ns. The circuit produces the expected correct outputs. The number of output electrons, and therefore the power consumption of the circuit, can be controlled by adjusting the 1-terminal voltage.

## 5.3.7. Logic Operation of a 2-Bit Adder Consisting of a Shared-BDD Circuit

As an example of a larger system, this section shows a circuit that implements the shared BDD for 2-bit addition given in Fig. 32. The designed circuit is shown in Fig. 37. The circuit accepts two 2-bit binary inputs [augend ($a_1a_0$) and addend ($b_1b_0$)] and produces the corresponding 2-bit sum output ($s_1s_0$) and a 1-bit carry output ($c_1$).

Computer simulation shows that the designed circuit perform adds correctly for all possible input combinations. Some of the simulation results are shown in Figs. 36 and 38. (The



Figure 35. Operation of the BDD circuit in Fig. 34 (simulation). Each impulse signifies the arrival of an electron at a root. The expected logical values for each output bit are written in the figure. Temperatures is assumed to be 0 K. Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.

**Figure 36.** Operation of the adder circuit at 0 K (simulation). Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.

same device parameters as for the preceding circuit were used. The power voltage of the 1-terminal was set to 10 mV for the results shown in Fig. 36 and 20 mV for Fig. 36). Figure 36 shows the results for zero temperature (0 K). Each impulse in the figure signifies the arrival of an electron at a root. The circuit produces the expected correct outputs.

Figure 38a shows the results at a temperature of 20 K. The circuit produces noisy outputs as a result of thermal agitation (i.e., a countercurrent or countertransport of electrons from a root into the circuit is frequently observed, shown by negative impulses in the figure). Nevertheless, the correct outputs can be retrieved by counting the net number of output electrons. The results are shown in Fig. 38b. In this example, the counting was repeated



**Figure 37.** Shared BDD circuit for a 2-bit adder. The ground capacitance for each node is omitted. Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.

Figure 38. Operation of the adder circuit at 20 K (simulation). Reprinted with permission from [37], T. Yamada et al., *Jpn. J. Appl. Phys.* 40, 4485 (2001). © 2001, Institute of Pure and Applied Physics.

every 2 ns (e.g., a total of nine electrons flowed out of root $s_0$ during the period from 12 to 14 ns). Using a longer period for the counting enables the operation of the circuit at higher temperatures.

### 5.3.8. Shared-BDD Integrated Circuits Fabricated by a Semiconductor Process Technology

Kasai and others developed single-electron subsystems based on the shared BDD, with the aim of composing a microprocessor [38, 39]. These subsystems are fabricated with a sophisticated process technology, Hexagon NanoProcess. In the first step of this technology, a hexagonal network of GaAlAs-GaAs heterostructure nanowires is formed on an insulating GaAs layer (the left figure in Fig. 39a). A quantum wire that consists of two-dimensional electron gas (2DEG) goes through the nanowire (the left figure in Fig. 39a). In the second step, a Schottky gate consisting of metal is formed on the wire to construct a tunneling gate (the right figure in Fig. 39a). The Schottky gate is wrapped around the wire, and a depletion layer or a potential barrier extends from the gate into the quantum wire. This potential barrier divides the quantum wire into two separate 2DEG regions, and a tunnel junction is formed between the two 2DEG regions.

A notable characteristic of this tunnel junction is that electron transport between the two 2DEG regions can be modulated by controlling the voltage of the gate to regulate the width of the potential barrier between the 2DEG regions. Electrons will be transferred by tunneling from one 2DEG region to the other if the gate is set to a low negative voltage, whereas no electrons will be transferred if the gate is set to a highly negative voltage. In this way, the tunnel junction can operate as the tunneling gate that turns electron transport on and off. With these tunneling gates, BDD devices can be constructed as shown in Fig. 39b. The device in the left figure consists of two tunneling gates, and the device in the left figure uses two single-electron transistor switches consisting of four tunneling gates. The on–off conductance and the two-way switching of the BDD devices are shown in Fig. 40.



Figure 39. (a) Hexagonal network of quantum wires and a tunneling gate with a Schottky gate (WPG), and (b) two structures of BDD devices, i.e., a quantum-wire (QWR) device consisting of two tunneling gates and a single-electron-transistor (SE) device consisting of four tunneling gates. Reprinted with permission from [39], S. Kasai et al., in "Proceedings of the 2003 Asia-Pacific Workshop on Fundamentals and Application of Advanced Semiconductor Devices," 2003, p. 177. © 2003, IEICE.

Figure 40. On-off conductance in (a) tunneling gate and (b) single-electron-transistor switch measured at various temperatures, and (c) two-way switching of the BDD device measured at 1.7 K and 297 K. Reprinted with permission from [39], S. Kasai et al., in "Proceedings of the 2003 Asia-Pacific Workshop on Fundamentals and Application of Advanced Semiconductor Devices," 2003, p. 177. © 2003, IEICE.

With the tunneling-gate BDD devices, several subsystems were constructed based on the shared BDD graphs shown in Fig. 41. To implement with Hexagon NanoProcess, these BDD graphs are drawn on a hexagonal-network chart and designed so that they have no intersection of branches. In Hexagon NanoProcess, the subsystem circuits were fabricated on a GaAlAs-GaAs hexagonal nanowire network that was formed by chemically etching a AlGaAs/GaAs heterostructure layer on an insulating GaAs substrate. To make QWR BDD devices, Schottky gate electrodes were attached on the nanowires by using EB lithography, metal deposition, and lift-off process. The width of the nanowire was 300–500 nm, and the gate length was 300 nm. As an example of fabricated subsystems, a 2-bit adder is shown in Fig. 42; the panel a shows a SEM image, and panel b shows input–output waveforms of 2-bit add operation for all possible input combinations. Using Hexagon NanoProcess will enable us to construct single-electron LSIs with large-scale integration of 25–45 million BDD devices/cm$^2$.

## 5.4. Single-Electron Majority Logic Circuits

### 5.4.1. A Short Sketch of Majority Logic

The majority logic is a way of implementing digital operations in a manner different from that of Boolean logic [40, 41]. It represents and manipulates digital functions on the basis of the principle of majority decision instead of using Boolean logic operators AND, OR, and their complements. The logic process of majority logic is much more sophisticated than that of Boolean logic, so the majority logic is more powerful for implementing a given digital function with a smaller number of logic gates.

The prospects for the practical applications of majority logic wholly depend on the feasibility of a logic device suitable for majority logic. In the late 1950s, several computer systems based on the majority logic architecture were developed and constructed for practical use by using a nonlinear-reactance device called the parametron, a majority logic device that uses the phenomenon of parametric phase-locked oscillation. After these developments, however,

**Figure 41.** Shared BDD graphs designed to construct several subsystems for a single-electron 2-bit processor. Reprinted with permission from [39], S. Kasai et al., in "Proceedings of the 2003 Asia-Pacific Workshop on Fundamentals and Application of Advanced Semiconductor Devices," 2003. p. 177. © 2003, IEICE.

majority logic had to leave the stage because the transistor gate circuit—a Boolean logic device by nature—came to be the dominant device in digital electronics. However, majority logic can be expected to make a comeback with the development of single-electron technology. This is so because, as shown in the following sections, the single-electron circuit will provide functional properties that can be well used for implementing majority logic operations.

### 5.4.2. Unit Function of Majority Logic

The unit function of majority logic is to determine the output state by means of the majority vote of input states. The logic element. a majority logic gate, or simply a majority gate has an odd number of binary inputs and a binary output. It produces an output of 1 if the majority of the inputs is 1, and produces an output of 0 if the majority is 0. The function of a three-input gate is shown in Fig. 43a, together with the logic symbol of the gate. When, for instance, the three inputs are 0, 1, and 1. the output is 1 (fourth row in the table); when the inputs are 1, 0, and 0, the output is 0 (fifth row in the table). (For further details on

**Figure 42.** Single-electron 2-bit adder fabricated with Hexagon NanoProcess technology. (a) An SEM and (b) input-output waveforms compared with theoretical ones, measured at room temperature. Reprinted with permission from [39], S. Kasai et al., in "Proceedings of the 2003 Asia-Pacific Workshop on Fundamentals and Application of Advanced Semiconductor Devices," 2003, p. 177. © 2003, IEICE.

majority logic, see Refs. [40] and [41].) Any digital function can be implemented using a combination of majority gates and inverters. A majority gate can have five or more inputs, but three-input gates suffice for the construction of any logic system.

Figure 43b gives the complement of the three-input majority function. The gate devices described in the following sections produce this complementary function. Any digital function can be implemented using only the gates of the complementary majority function; using inverters jointly makes the logic system design more concise. The logic symbol given in Fig. 43b is used for a complementary majority gate.

Majority logic provides a concise implementation for most functions encountered in logic design applications. As an example, Fig. 44a shows the implementation of a full adder, and Fig. 44b shows a 4-bit adder consisting of four full adders connected into a cascade manner. In the figures, an inverter is represented by a segment on a connection branch, according to the conventional flow-diagram description of majority logic. A full adder is composed of only three gates with two inverters. In contrast, a Boolean-based implementation requires a larger circuit with seven or eight gate elements (about 25–30 MOSFETs).

The essential function for the majority gate operation is to ascertain whether the majority of inputs is 1 or 0. This can be performed by calculating the mean of the inputs (each input is 1 or 0) and comparing the mean value with 0.5; if the mean value is larger than 0.5,



**Figure 43.** Three-input majority gates: (a) majority-logic gate and (b) complementary majority-logic gate.

Figure 44. Full adder and a 4-bit adder consisting of majority gates and inverters. Each symbol denotes a complementary function gate. An inverter is represented by a segment on a connection branch.

the majority can be considered to be 1, and if it is smaller than 0.5, then the majority can be considered to be 0. A promising way of implementing this function with single-electron circuits is to use the technique of charge-balancing capacitor summation, with a threshold device for producing the corresponding 1/0 output. The following sections give an outline of two majority gates: One uses a Tucker's inverter, and the other uses a single-electron box as a threshold device.

## 5.4.3. Constructing a Majority Gate with a Tucker's Inverter

**5.4.3.1. Gate Circuit** Iwamura and others proposed a majority gate that used a Tucker's inverter as a threshold device [42]. Figure 45(a) illustrates the gate with an example of a three-input configuration. The gate consists of an input capacitor array (six capacitors $C$) for input summation and an inverter subcircuit (four tunnel junctions $C_{j1}-C_{j4}$ and three capacitors $C_1-C_3$) for threshold operation. The gate accepts three input voltages $V_1$, $V_2$, and $V_3$ and produces the corresponding output voltage $V_{out}$. The configuration of this circuit is modeled on the single-electron inverter proposed by Tucker [6], illustrated in Fig. 45b; that is, starting with the inverter configuration of Fig. 45b, each of the two input capacitances ($C_0$ in Fig. 45b) are divided into three equal capacitances, then the divided capacitances ($C$ in Fig. 45a) are connected to three input terminals to create the majority gate. The gate input capacitance $C$ is set at one-third of the inverter input capacitance $C_0$; that is, $C = C_0/3$ (the other parameters are set to be the same as those of the original inverter).

This gate produces the complementary majority function as follows: The input nodes $P$ and $Q$ of the inverter subcircuit are coupled to each input $V_1$, $V_2$, or $V_3$ through each input capacitance $C$, so the potential of each input node is changed in proportion to the mean value of the inputs. Therefore, if two or three inputs are 1, the gate circuit is in the same state as that of an inverter that is applied an input of larger than 0.5, and so the corresponding output will be 0. If two or three inputs are 0, the gate state is the same as the inverter state for an input of smaller than 0.5, therefore the output will be 1.



Figure 45. Single-electron majority logic gate circuit: (a) circuit configuration and (b) Tucker's inverter circuit. Reprinted with permission from [42]. H. Iwamura et al., IEICE Trans. Electron. E81-C, 42 (1998). © 1998. IEICE.

To achieve the correct gate operation, the inverter subcircuit must be designed so that it will be stable throughout its transfer curve except a high–low transition point to exhibit the characteristic of a step inverter (Fig. 46). This is important for majority gate application because, in the majority gate, the inverter subcircuit will frequently receive an intermediate input value between 1 and 0. Selecting appropriate parameter values can give an inverter the step characteristic. A sample set of parameters for a three-input gate is $C = 1$ aF, $C_{j1} = C_{j4} = 1$ aF, $C_{j2} = C_{j3} = 2$ aF, $C_1 = 9$ aF, $C_2 = 9$ aF, $C_3 = 24$ aF, and $V_{dd} = 65$ mV.

This parameter set gives the transfer curve shown in Fig. 47, simulated under the condition of all the three input terminals being applied the same input voltage. A slight hysteresis is observed at the high–low transition but is not a problem in the gate circuit. It s worthy of note that, unlike CMOS inverters, single-electron step inverters do not produce short-circuit current even for intermediate values of input. This is quite convenient for low-power operation.

**5.4.3.2. Subsystem Design: Adders**   One of the simple subsystems is a full adder (Fig. 44a) consisting of three majority gates and two inverters. This circuit operates as summarized in Fig. 48, simulated with parameters same as those given above and a tunnel junction resistance of 100 kΩ. Temperature was assumed to be 0 K. In the figure, the three inputs (adder input $A$, adder input $B$, carry input $C_{in}$) are applied in sequence as 000, 001, 010, 011, 100, 101, 110, 111), and the corresponding two outputs (sum output $Sum$, carry output $C_{out}$) is observed as (00, 10, 10, 01, 10, 01, 01, 11).

Some explication is needed for the operation speed of the circuit. Electron tunneling is in general a probabilistic phenomenon, so the operation delay time in a single-electron circuit is not a fixed value but differs in each operation event. (In simulation, the probabilistic characteristic is taken into account by use of random numbers, and Fig. 48 shows the result for a given set of random numbers.) It is therefore the usual practice for single-electron circuits to represent the operation speed by the mean of operation delay time. Following this practice shows that the mean add time of the full adder is 0.3 ns.

Figure 49 shows the operation of the circuit for a 4-bit adder (Fig. 44b). The output waveforms for the add operation of 1111 + 0001 are plotted. In this figure, both of the 4-bit inputs (an addend and an augend) are initially set at 0000; then, at time = 0, one input is turned to 1111 and the other to 0001. The operation speed is limited by the carry delay, and the mean add time for the operation of 1111 + 0001 is 1.2 ns.

## 5.4.4. Constructing a Majority Gate with a Single-Electron Box

**5.4.4.1. Single-Electron Box as a Threshold Device**   Oya and others proposed a majority gate that used a single-electron box instead of a Tucker's inverter [43, 44]. Their gate uses a balanced pair of single-electron boxes or an irreversible single-electron box as a threshold device that compares the output of a capacitor summation array with 0.5 to produce the 1/0 output.

Figure 50a shows the majority gate with an irreversible single-electron box [44]. It consists of two identical tunnel junctions $C_j$ connected in series, a bias capacitor $C_1$, and a bias



Figure 46. Transfer curve of a step inverter required for the inverter subcircuit.

**Figure 47.** Transfer curve of the majority gate circuit. Simulated under the condition of all the three input terminals being applied the same input voltage. For the circuit parameters, see the text. Reprinted with permission from [42]. H. Iwamura et al., *IEICE Trans. Electron.* E81-C, 42 (1998). © 1998, IEICE.



**Figure 48.** Add operation of the single-bit full adder (simulation): the waveforms of a sum output *Sum* and a carry output *Cout* for the eight input combinations. Reprinted with permission from [42]. H. Iwamura et al., *IEICE Trans. Electron.* E81-C, 42 (1998). © 1998, IEICE.



**Figure 49.** Operation of the 4-bit adder (simulation): the output waveforms (a sum output *Sum* and a carry output *Cout*) for the add operation of 0001 + 1111. Reprinted with permission from [42]. H. Iwamura et al., *IEICE Trans. Electron.* E81-C, 42 (1998). © 1998, IEICE.

Figure 50. Irreversible single-electron box: (a) circuit configuration with a double tunnel junctions, (b) static number $n$ of excess electrons on node 1 as a function of bias voltage $V_d$, (c) voltage at node 1 as a function of $V_d$.

voltage $V_d$. It has an island node 1, at which excess electrons are stored. At the ow temperatures at which the Coulomb-blockade effect occurs, the number $n$ of excess electrons takes a value such that the electrostatic energy in the circuit (including the bia; voltage source) is locally minimized. The value of $n$ is 0 at $V_d = 0$, and it changes with $V_d$ because of electron tunneling between node 1 and the ground through junctions $C_j$ via intermediate node 2. In this circuit, $n$ is a hysteretic staircase function of $V_d$, as shown in Fig. 50b; $n$ changes from 0 to 1 when $V_d$ is increased above threshold $V_2$ [$= e(1 + C_L/C_j)/(2C_L)$; $e$ is the elementary charge] and returns from 1 to 0 when $V_d$ is decreased below threshold $V_1$ [$= e(1 - C_L/C_j)/(2C_L)$]. Because of this discrete changes in $n$, the voltage at node 1 is a hysteretic sawtooth function of $V_d$, as shown in Fig. 50c. For this operation, the irreversible single-electron box is called a single-electron trap.

A similar hysteretic function can be obtained using a multijunction box—a single-electron box that has three or more tunnel junctions connected in series. A multijunction box can therefore also be used to construct the majority gate. In an $N$ junction trap, the thresholds are given by $V_2 = e[1 + (N - 1)C_L/C_j]/(2C_L)$ and $V_1 = e[1 - (N - 1)C_L/C_j]/(2C_L)$.

### 5.4.4.2. Gate Circuit Configuration

A majority gate can be constructed with an irreversible single-electron box, as illustrated in Fig. 51a with a three-input configuration. The majority gate consists of a double-junction box ($C_L$ and two junctions $C_j$), three input capacitors $C$, and an output capacitor $C$ (the input and output capacitors are set at equal capacitance). Three input voltages, $V_1$, $V_2$, and $V_3$, are applied to node 1 through the input capacitors. The input capacitors form a voltage-summing network and produce the mean of their inputs on node 1. The double-junction box then produces the complementary majority-logic output on the same node, as illustrated later, and the output is sent to a succeeding gate through the output terminal. Binary logic values 1 and 0 are represented by a positive voltage and a negative voltage of equal amplitude.

The operation of the majority gate is as follows. We first ground the output terminal and apply the input voltages and then increase bias voltage $V_d$ to an appropriate excitation value, $V_{ex}$. The voltage at node 1 reaches a positive value determined by $V_{ex}$ and input voltages $V_1$, $V_2$, and $V_3$. If the node voltage exceeds a threshold, an electron will tunnel from the ground to node 1 via intermediate node 2: consequently, the node voltage will turn negative.



Figure 51. Majority-gate device: (a) simulated circuit configuration and (b) voltage at node 1 as a hysteretic function of $V_d$ with four sets of inputs as parameters. The dashed lines show the node voltage at inputs $V_1 = V_2 = V_3 = 0$ V. Reprinted with permission from [44]. E. Oya et al., IEEE Trans. Nanotechnology 2, 15 (2003). © 2003, IEEE.

In contrast, if the node voltage does not reach the threshold, it will remain positive. We then retrieve the node voltage as an output. For successful majority-logic operation, we set excitation voltage $V_{ex}$ to $e[1 + (C_L + 4C)/C_j]/(2C_j)$ so that electron tunneling will take place only if two or three inputs are a logical 1 (or only if the mean of three input voltages is positive). After exciting the gate, we decrease bias $V_d$ to a holding voltage $e/(2C_L)$, ground the input terminals (or set $V_1$, $V_2$, and $V_3$ to 0 V), and then observe the voltage at node 1, the output voltage. The output voltage (therefore input voltage for a succeeding gate) is $-e/(2C_L + 8C + C_j)$, or logical 0, if two or three inputs are 1 (electron tunneling takes place), and it is $e/(2C_L + 8C + C_j)$, or logical 1, if two or three inputs are 0 (no electron tunneling takes place).

### 5.4.4.3. Majority Logic Operation

Figure 51b shows the gate operation, with simulated results for a sample set of parameters, $C_L = 2$ aF, $C_j = 20$ aF, $C = 2$ aF, and zero temperature. In this example, a logical 1 is represented by a voltage of 4 mV and a logical 0 by $-4$ mV. The figure shows the voltage at node 1 as a function of $V_d$ with a set of three inputs as a parameter set (the output terminal is grounded). In the figure, for example, "inputs (1, 1, 0)" means that two inputs are set to 4 mV (logical 1) and that one input is set to $-4$ mV (logical 0). With increasing $V_d$, the node potential increases to a maximum, then drops to a negative because of electron tunneling. The threshold value of $V_d$, at which the electron tunneling takes place, depends on the sum of the inputs; in this example, the threshold is 56 mV for inputs (1, 1, 0) and 64 mV for inputs (1, 0, 0). To operate the device as a majority gate, $V_d$ is increased to an excitation value of 60 mV (indicated by $V_{ex}$ in the figure). The logic output, or the voltage at node 1, was retrieved through the output terminal. To do this retrieving, we decreased bias $V_d$ to a holding value of 40 mV (indicated by $V_{ss}$), grounded the three input terminals, and checked the voltage of node 1. The node voltage was 4 mV (quiescent point $A$ on the dashed line) when the output was 1 and $-4$ mV (quiescent point $B$) when the output was 0.

Figure 52 illustrates the logic operation for all input combinations, with simulated results for $C_L = 2$ aF, $C_j = 20$ aF, and $C = 2$ aF; tunnel junction conductance $= 5$ $\mu$S; and zero temperature. The bias voltage $V_d$ is the two-step clock pulse shown in the upper plot in the figure; first, $V_d$ is set to an excitation voltage $V_{ex}$ of 60 mV, and then it is set to a holding voltage $V_{ss}$ of 40 mV. Three inputs ($V_1$, $V_2$, and $V_3$) are applied synchronously with the bias clock. They are the rectangular pulses (4 mV for logical 1 and $-4$ mV for logical 0) in the middle plot in Fig. 52. In the figure, the four sets of inputs (1, 1, 1), (1, 1, 0), (1, 0, 0), and



Figure 52. Simulated majority-logic operation of the gate device. Reprinted with permission from [44], T. Oya et al., *IEEE Trans. Nanotechnology* 2, 15 (2003). © 2003, IEEE.

(0. 0. 0) were applied in sequence. Depending on the majority of the inputs, the voltage at node 1 changes from 0 to positive (1-valued) or negative (0-valued) (bottom plot in Fig. 52). For a 0 output, the output voltage initially goes high for an instant with the rise in $V_d$ and then turns negative as electron tunneling takes place. The output established in each clock cycle is maintained after the input pulses are turned off until $V_d$ returns to zero (duration $T$ in the bottom plot).

### 5.4.4.4. Subsystem Design: Full Adder

Any logic function can be implemented by combining identical gates into a cascade configuration, with the output capacitor of one gate acting as the input capacitor of the following gate. An example is illustrated in Fig. 53(a). This majority gate is bilateral, so a three-phase clock is used to control the signal-flow direction. To do this, the gate circuits are divided into three groups, and each group is excited in turn by one phase of the three clock signals. $\phi_1$ to $\phi_3$, as shown in Fig. 53b. For instance, in Fig. 53a. the leftmost gate (and every fourth gate thereafter) belongs to the first group and is excited by the $\phi_1$-phase clock; the middle gate (and every fourth gate thereafter) belongs to the second group and is exited by the $\phi_2$-phase clock; the rightmost gate (and every fourth gate thereafter) belongs to the third group and is excited by the $\phi_3$-phase clock. The phases of the three clock signals overlap so that the output of a stage will be established while the preceding stage is maintaining its output during its holding period; signals are thus transmitted from one gate to the next. For successful interstage coupling, the duration of the overlap has to be longer than the excitation period, as shown in Fig. 53b. The signal-flow direction is determined by the relative timing of the three phases; in Fig. 53a, it is rightward.

Figure 54a illustrates an implementation of the full adder shown in Fig. 44. The adder accepts three inputs, augend $A$, addend $B$, and carry input $C_{in}$ (and their complements $\overline{A}$, $\overline{B}$. and $\overline{C}_{in}$); it then produces the corresponding carry output $C_o$ and sum output $S_o$. The core of the adder consists of majority gates 2, 3, and 7. The other gates (1, 4, and 5) act as a delay buffer to transfer the signal from the preceding stage to the following stage with the correct clock timing. Gate 6 is a fan-out buffer that transmits signals from the following stage to the succeeding two stages. The inputs are taken in while clock $\phi_1$ is in the excitation period. Carry output $C_o$ is produced when $\phi_2$ goes high and is retrieved by gate 7 while $\phi_2$ is in the holding period. Sum output $S_o$ is produced when $\phi_3$ goes high again. The delays between the inputs and the carry output are two-thirds of a clock period, and that between the inputs and the sum output is one clock period.

A simulated result is shown in Fig. 54b. Three clock cycles of the output waveforms of carry $C_o$ and sum $S_o$ are shown: the device parameters are as given in the previous section. Two sets of inputs $(A, B, C_{in}) = (1, 1, 0)$ and $(0, 0, 1)$ were sequentially entered, and the correct outputs $(C_o, S_o) = (1, 0)$ and $(0, 1)$ were produced in response. Multibit adders can be constructed by combining the full adders into a cascade configuration.

Figure 55a shows the configuration of a 4-bit adder consisting of four full adders connected into a cascade configuration. In this adder, signals flow leftward with clocks $\phi_1$ through $\phi_3$. The circuit accepts two 4-bit inputs, augend $X_3X_2X_1X_0$ and addend $Y_3Y_2Y_1Y_0$ (and the complement of each input bit); it then produces the corresponding carry output $C_3$ and 4-bit sum output $S_3S_2S_1S_0$. Each bit signal of the inputs is applied synchronously with the corresponding clock timing. as indicated by codes $\phi_1$ to $\phi_3$ in the figure; for example. input

(a)                                                    (b)

**Figure 54.** Full adder: (a) circuit configuration and (b) simulated output waveforms of carry $C_o$ and sum So. waveforms of $\phi_3$-clock and the inputs. Clocks $\phi_1$ and $\phi_2$ are not shown. Reprinted with permission from [44]. T. Oya et al., *IEEE Trans. Nanotechnology* 2, 15 (2003). © 2003. IEEE.

bits $X_0$ and $Y_0$ (and bias inputs 1 and 0) are applied when $\phi_3$ goes high, and $X_1$ and $Y_1$ are applied when $\phi_2$ goes high. Output bits $S_0$, $S_1$, $S_2$, $C_3$, and $S_3$ are produced in this order with clocks $\phi_3$, $\phi_2$, $\phi_1$, $\phi_3$, and $\phi_3$. The delay between the first inputs ($X_0$ and $Y_0$) and the last output $S_3$ is a three-clock period. (Using shift resistors as well will provide modified adders that can accept all input bits simultaneously and produce all output bits simultaneously.)

Figure 55b illustrates the operation of the adder, with four clock cycles of output waveforms for $C_3$ and $S_3$ (the waveforms of $S_0$ through $S_2$ are omitted because of limited space). The device parameters were as given in Section 5.4.2. In the figure, for instance, the output in the period from 155 to 195 ns ($C_3 = 1$ and $S_3 = 0$) is the response to inputs augend 1111 and addend 0001, which were applied sequentially from 35 to 150 ns.

## 5.5. Boltzmann-Machine Neural Networks Using the Stochastic Nature of Single-Electron Tunneling

### 5.5.1. Outline of Boltzmann Machines

The Boltzmann machine is a kind of recurrent neural network that can solve various problems in areas such as combinatorial optimization, classification, and learning. It consists of a large network of processing units called neurons that are interconnected bidirectionally with



(a)                                                    (b)

**Figure 55.** 4-bit ripple carry adder: (a) circuit configuration and (b) simulated output waveforms of carry bit $C_o$ and most significant sum bit $S_3$. waveforms of $\phi_3$-clock and the inputs. Clocks $\phi_1$ and $\phi_2$ are not shown.

signal connections having various connection weights. Each neuron receives input signals from every other neuron and sends output signals to every other neuron. The neuron has a binary output state and changes its state in response to the inputs, according to a stochastic transition rule. All neurons operate in parallel, and each one adjusts its own state to the states of all the others; consequently, the whole network converges into an optimal configuration. The structure of mathematical problems such as combinatorial optimization can be mapped onto the structure of a Boltzmann machine by determining the connection weights between the neurons. In this way, finding the solution to a problem can be reduced to finding the optimal configuration of the Boltzmann machine. The unique and important feature of Boltzmann machines is their stochastic neuron operation combined with simulated annealing algorithms. This allows Boltzmann machines to reach a globally optimal configuration (and thereby an optimal solution) without falling into local minimum configurations. For a detailed explanation, see Refs. [45] and [46].

A Boltzmann machine LSI circuit for practical use must integrate thousands of neurons on a chip. The crucial problem in developing such LSIs is how to implement the generation of randomness for the stochastic neuron operation. Every neuron has to have its own randomness because stochastic independence between the neurons is required. Electronic circuits that are currently available for generating randomness—such as the thermal noise amplifier and the random bit generator [47]—consist of many device elements and, consequently, require a large area. They, therefore, cannot be used for large-scale Boltzmann machine LSIs, and so there is a need for a novel device for constructing Boltzmann machine LSIs.

The following sections describe an attempt to implement Boltzmann machines on electronic circuits, using the single-electron circuit technology. The single-electron circuit shows stochastic behavior in its operation because of the probabilistic nature of electron tunneling, so it can therefore be successfully used for implementing Boltzmann machines with simple construction. The following sections explain the neuron function required for the Boltzmann machine, implementation of Boltzmann machines with single-electron circuits, and the problem-solving operation of a single-electron Boltzmann machine network.

### 5.5.2. Function of Neurons Required for Boltzmann Machine Operation

A Boltzmann machine consists of a network of many identical neurons interconnected with each other. The configuration of the network is illustrated in Fig. 56. The output of each neuron is fed back into the inputs of other neurons, and each neuron exchanges signals with others to update its own output. Notation $W_{ij}$ is the connection weight to neuron $i$ from neuron $j$, $T_i$ is the threshold connection weight to neuron $i$ from a bias that is fixed at the value of 1, and $x_i$ is the output of neuron $i$. Connection weights $W_{ij}$ and $T_i$ can be given at any desired value under the restrictions that $W_{ij} = W_{ji}$ and $W_{ii} = 0$. The output $x_i$ of each



Figure 56. Concept of the Boltzmann machine. It is a recurrent network consisting of many identical neurons with signal connections. Each neuron produces a binary stochastic output.

neuron is binary (i.e., $x_i = 1$ or $-1$), and a set of neuron states ($x_1$, $x_2$, $x_3$, etc.) is called the state of the network.

In this network, each neuron $i$ takes a weighted sum of inputs according to the following equation:

$$S_i = \sum_i W_{ij} x_i + T_i \tag{17}$$

The neuron produces an output of a $1/-1$ bit stream in response to the weighted sum $S_i$, following the logistic-sigmoid probability function given by

$$f(S_i) = 1/[1 + \exp(S_i/C)] \tag{18}$$

where $f(S_i)$ is the probability for generation of an output 1 ($x_i = 1$) in the bit stream. A control parameter $C$ regulates the dependence of the probability on $S_i$. It is decreased slowly from a large value to a very small value during the simulated-annealing process, as shown later. Through this process, a Boltzmann machine network changes its state to minimize the "energy" function $E$ defined by

$$E = -\frac{1}{2} \sum_i \sum_j W_{ij} x_i x_j - \sum_i T_i x_i \tag{19}$$

By our setting appropriate connection weights $W_{ij}$ and $T_i$, the energy function of the network can be related to the objective function (cost function) of a given optimization problem. This way, the solution to the problem can be found simply by observing the final state that the network reaches.

### 5.5.3. Creating Neuron Device with a Single-electron Circuit

The function of the Boltzmann machine neuron can be implemented with a single-electron digital oscillator that generates an output of $1/-1$ bit stream. In a single-electron digital oscillator, the duration of an output 1 (or an output $-1$) will fluctuate randomly because of the probabilistic nature of electron tunneling. This oscillator will enable us to produce an output of random $1/-1$ bit stream required for the Boltzmann-machine neuron operation. For the complete function of the neuron, the digital oscillator must be designed so that the probability of an output 1 can be modulated by a signal input [$S_i$ in Eq. (17)], according to the logistic-sigmoid probability function regulated by a control input [$C$ in Eq. (18)].

Amemiya and others proposed constructing such oscillators by using a Tucker's inverter and modifying its circuit configuration [48, 49]. Figure 57 shows their oscillator or a neuron device [49]. The oscillator consists of four tunnel junctions ($C_{j1}$ through $C_{j4}$) and seven



Figure 57. Configuration of the unit-neuron circuit for single-electron Boltzmann machines. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

capacitors ($C_1$ through $C_7$), with two bias voltages (i.e., a positive voltage $V_{dd}$ and a negative voltage $-V_{ss}$). Offset voltage $V_h$ is for adjusting the operating point of the circuit. The circuit, with the appropriate set of parameters given later, operates as an astable multivibrator or a square-wave oscillator. That is, the circu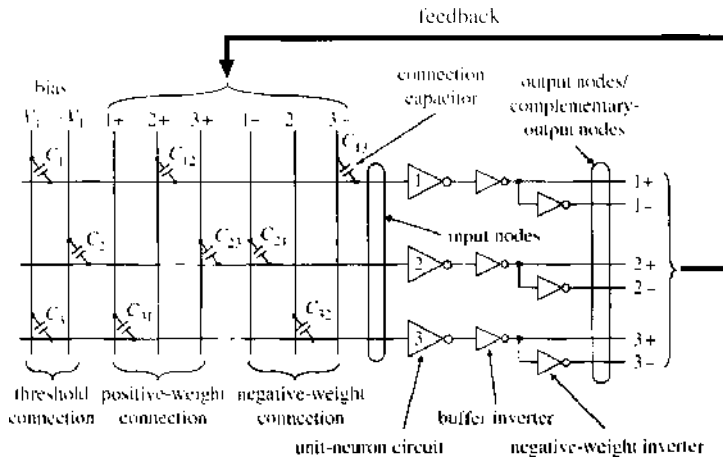it repeats a cycle of transferring one electron from $-V_{ss}$ to $V_{dd}$ by making an electron tunneling from one node to another in the following sequence: node $L \rightarrow V_{dd}$ (tunneling 1), node $M \rightarrow$ node $L$ (tunneling 2), $-V_{ss} \rightarrow$ node $N$ (tunneling 3), and node $N \rightarrow$ node $M$ (tunneling 4). Output voltage $x_i$ of the circuit is nearly equal to $V_{dd}$ (an output 1) during a period between tunneling 2 and tunneling 4 because, in this period, an electron is extracted from node $M$ (output node) and, consequently, node $M$ is charged positive. In the remaining period, the output is nearly equal to $-V_{ss}$ (an output $-1$). The time that the node $M$ is charged positive depends on the waiting time for tunnelings, so the output randomly alternates between 1 and $-1$ according to the probabilistic fluctuation in the waiting time for tunneling. The probability for an output 1 can be controlled by external voltage $S_i$. Thus, the circuit accepts voltage input $S_i$ (the weighted sum of inputs) and produces the corresponding voltage output $x_i$ in a form of random $1/-1$ bit stream; this circuit is hereafter called an unit neuron circuit. [This neuron circuit does not include a subcircuit for calculating the weighted sum of inputs $S_i$. For calculating $S_i$ according to Eq. (17), an additional capacitor subcircuit will be used, as described in Section 5.5.4.]

To create the stochastic neuron operation, the circuit parameters have to be set so that the circuit can operate under oscillating conditions. To determine the optimum parameters, the stability diagram of the circuit is used as a guide map. (A stability diagram illustrates the internal states of a single-electron circuit in a multidimensional space of circuit variables—namely, the voltages of powers and inputs, and the capacitances of tunnel junctions and capacitors.) An example set of the capacitance parameters for the unit-neuron circuit is

$$C_{i1} = C_{i4} = 1 \text{ aF}$$

$$C_{i2} = C_{i3} = 2 \text{ aF}$$

$$C_1 = C_2 = 12 \text{ aF}$$

$$C_3 = C_4 = 4 \text{ aF}$$

$$C_5 = C_6 = 10 \text{ aF}$$

$$C_7 = 24 \text{ aF}$$

(20)

## 5.5.4. Operation of Unit-Neuron Circuit

The internal state of the unit-neuron circuit is expressed by the numbers of excess electrons ($l$, $m$, $n$) on the three nodes ($L$, $M$, and $N$) in the circuit. The stability diagram of the unit-neuron circuit is drawn in a four-dimensional space of four voltage variables ($S_i$, $V_h$, $V_{dd}$, and $-V_{ss}$). In Fig. 58a-d, a part of the diagram is illustrated on a plane of two voltage variables, input $S_i$ and offset $V_h$, assuming the capacitance parameters given by Eq. (20). The two white regions are stable regions, in which the circuit stabilizes at internal states (0, $-1$, 0) and (0, 0, 0); state (0, $-1$, 0) produces a positive output voltage (an output 1), whereas state (0, 0, 0) produces negative output voltage (an output $-1$). The output state for each internal state is illustrated by putting a letter ($H$ for an output 1, or $L$ for an output $-1$) before the electron number set. The dark regions are unstable regions in which electron tunneling frequently occurs and, consequently, the circuit alternates between two or more internal states to output a random $1/-1$ bit stream. The width of the unstable region can be controlled by regulating bias voltages $V_{dd}$ and $-V_{ss}$ as shown in the Figs. 58a-d.

The line $PQ$ in each figure shows the operating line on which the unit-neuron circuit has to be operated. Increasing input $S_i$ moves the operating point from $H(0, -1, 0)$ region to $L(0, 0, 0)$ region on line $PQ$, and this changes the probability for generation of an output 1. The control parameter for the probability function can be changed by regulating bias voltages $V_{dd}$ and $-V_{ss}$ to change the width of the unstable region. In regulating $V_{dd}$ and $-V_{ss}$, the value of offset $V_h$ has to be adjusted simultaneously so that the operating point for

**Figure 58.** Stability diagrams of the unit-neuron circuit given in Fig. 57, plotted on a plane of input voltage $S_i$ and offset voltage $V_b$. For capacitance parameters, see the text. The shaded regions are unstable regions. Figure 58 (a) through (d) correspond to a gradual increase in $V_{dd}$ and $V_{ss}$. The value of ($V_{dd}$, $-V_{ss}$) is: (a) (2.72 mV, $-3.29$ mV), (b) (2.74 mV, $-3.30$ mV), (c) (2.80 mV, $-3.36$ mV), (d) (2.87 mV, $-3.43$ mV). Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

zero inputs ($S_i = 0$) will be situated on the center line of the unstable region and, thereby, the probability for an output 1 will exactly be 0.5 at zero inputs. A set of ($V_{dd}$, $-V_{ss}$, and $V_b$) is hereafter called the control-parameter set.

Figure 59a–b show the operation of the unit-neuron circuit, simulated with the capacitance parameters given by Eq. (20), tunnel resistances of 100 kΩ for junctions $C_{j1}$ and $C_{j4}$ and 5 MΩ for $C_{j2}$ and $C_{j3}$, and zero temperature. The results are for the control-parameter set of (2.80 mV, $-3.36$ mV, $-0.87$ mV) that corresponds to the operating line $PQ$ in Fig. 58c. The output voltage waveform (a random $1/-1$ bit stream) for two instance values of the input voltage is plotted: (a) input $S_i = 1$ mV (point $Y$ in Fig. 58c) and (b) $S_i = -1$ mV (point $X$ in Fig. 58c). The probability of an output 1 can be changed by input $S_i$. The state of a negative



**Figure 59.** Output voltage waveforms for the unit-neuron circuit with the control-parameter set of (2.80 mV, $-3.36$ mV, $-0.87$ mV). Simulated for two input voltages: (a) $S_i = 1$ mV and (b) $S_i = -1$ mV. For circuit parameters, see the text. Temperature is 0 K. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

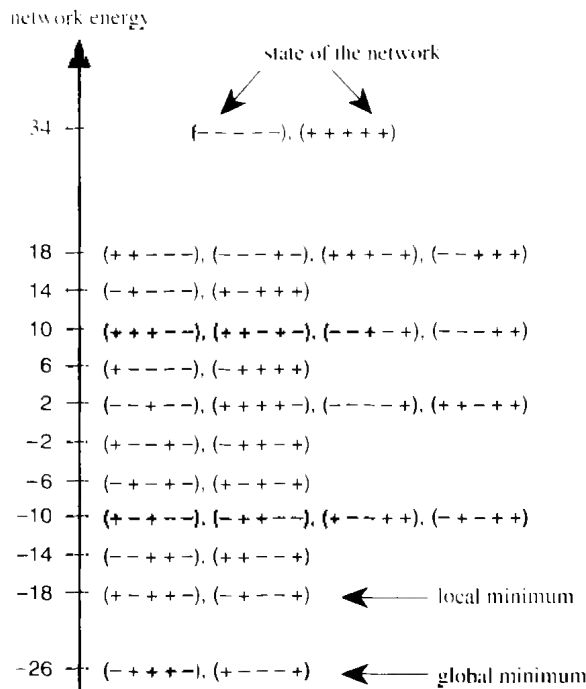output (output $-1$) is dominant for a negative value of $S_i$ [Fig. 59a], whereas the state of a positive output (output $1$) is dominant for a positive value of $S_i$ (Fig. 59b). Intermediate states are also generated (i.e., states $(0, -1, 1)$ and $(-1, 0, 0)$ in Fig. 59a), but this is not a problem because their duration is always short regardless of the input voltage value. In this example, the circuit changes its internal state in a cycle of: $L(0, 0, 0) \rightarrow L(-1, 0, 0) \rightarrow H(0, -1, 0) \rightarrow H(0, -1, 1) \rightarrow L(0, 0, 0)$.

The probability of an output $1$ is illustrated in Fig. 60a as a function of input $S_i$, for various control-parameter sets ($V_{dd}$, $-V_{ss}$, $V_h$). The probability function required for the stochastic neuron can be obtained. The probability can be controlled by regulating the control-parameter set; the number ($1$ through $4$) for each curve indicates a specific control-parameter set that is required for producing the characteristic of the curve. Figure 60b illustrates a diagram for setting the control-parameter set ($V_{dd}$, $-V_{ss}$, $V_h$), where the numbers $1$ through $4$ indicate the sets for producing curve $1$ through $4$ in Fig. 60a (e.g., the set for curve 3 can be obtained for $V_{dd} = 2.80$ mV, $-V_{ss} = -3.36$ mV, and $V_h = -0.87$ mV). In curve 1, the circuit acts as a simple threshold element without stochastic operation, which corresponds to the condition of $C = 0$ in Eq. (18). [Strictly speaking, the obtained characteristic is somewhat different from that of Eq. (18)—the characteristic is a line sigmoid function rather than a logistic sigmoid.]

### 5.5.5. Designing Boltzmann Machine Networks

The Boltzmann machine can be constructed by combining the unit-neuron circuits into a network. The overall configuration of the network circuit is illustrated in Fig. 61. The network consists of a number of the unit-neuron circuits with buffer inverters, negative-weight inverters, and connection capacitors. The buffer inverter is added to each unit-neuron circuit for intensifying the power of load drivability. Hereafter, the voltage of buffer-inverter outputs ($1+$, $2+$, $3+$, etc.) is called the output of neurons. The negative-weight inverters produce voltage signals ($1-$, $2-$, $3-$, etc.) that are complementary to the neuron outputs ($1+$, $2+$, $3+$, etc.); the complementary signals are used for obtaining negative weight connections. The buffer inverters and the negative-weight inverters are shown in Fig. 62. Both the output and its complement of each neuron circuit are fed back into inputs for other unit-neuron circuits. The connection between two neurons is established by a coupling capacitor $C_{ij}$. The threshold for each neuron is set up by positive bias voltage $V_1$ (or by negative voltage $-V_2$) with a coupling capacitor $C_i$.

Connection weights between neurons can be set at any desired values by choosing the capacitances of the coupling capacitors. Each weight [$W_{ij}$ and $T_i$ in Eq. (17)] is given by

$$|W_{ij}| = C_{ij}/(C_{ij} + C_i)$$
$$|T_i| = C_i/(C_{ij} + C_i)$$

(21)

For a positive weight ($W_{ij} > 0$), the coupling capacitor is connected with the input node of neuron $i$ and the output node ($1+$, $2+$, $3+$, etc.) of neuron $j$, and for a negative weight ($W_{ij} < 0$), with the input node and the complementary-output node ($1-$, $2-$, $3-$, etc.). The threshold coupling is made between the input node and positive-bias node $V_1$ if $T_i > 0$ and



Figure 60. Probability function of the unit-neuron circuit. For the device parameters, see the text: (a) the probability of generating an output 1 for various control-parameter sets, and (b) a diagram for setting the control-parameter set. The curves #1 through #4 in (a) are obtained for the control-parameter sets #1 through #4 in (b). Reprinted with permission from [49]. T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001. IOP Publishing.

feedback



**Figure 61.** Overall configuration of the single-electron Boltzmann machine network. The neuron outputs (1+, 2+, 3+, etc.) and the complementary outputs (1−, 2−, 3−, etc.) are fed back to become the inputs for the unit-neuron circuits. The connection between two neurons is established by coupling capacitor $C_{ij}$, and the threshold input for each neuron is set by coupling capacitor $C_i$. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

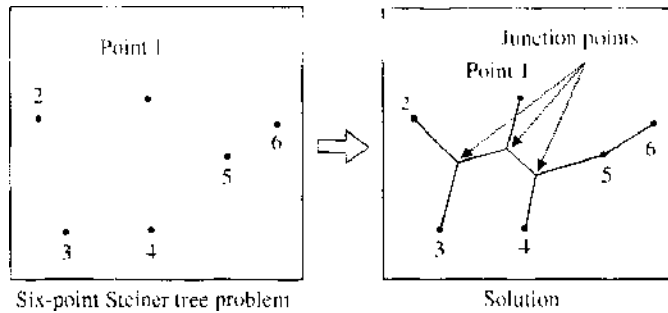between the input node and the negative-bias node $-V_2$ if $T_i < 0$. The capacitances $C_{ij}$ and $C_i$ have to be set at such values that the symmetry in connection (i.e., $W_{ij} = W_{ji}$ for all $i$, $j$) can be established. Under these conditions, the network circuit will operate as a complete Boltzmann machine.

## 5.5.6. Problem Solving Operation of the Network Circuit

Setting appropriate values for coupling capacitances can implement various optimization problems on the network circuit. As an example, this section shows a sample network circuit that solves an instance of the max-cut problem.

**5.5.6.1. Implementing the Max-Cut Problem on the Network Circuit** The max-cut problem is stated as follows: given a graph $G = (V, E)$ with positive weights on the edges, find a partition of the vertices $V = \{1, 2, \ldots, n\}$ into two disjoint sets $V_0$ and $V_1$ such that the sum of the weights of the edges that have one endpoint in $V_0$ and one endpoint in $V_1$ is maximal. To formulate the objective function for this problem, we here define a number of variables. Let $d_{ij}$ be the weight associated with the edge $i$, $j$ (by definition, $d_{ij} = d_{ji}$) and let



**Figure 62.** Tucker's inverter used for the buffer inverters and the negative-weight inverters, together with its circuit parameters (tunnel resistance = 100 kΩ for all of the four junctions). Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

$x_i$ be a $1/-1$ variable defined as

$$x_i = 1(\text{if } i < V_1)$$
$$= -1(\text{if } i < V_0) \tag{22}$$

then the max-cut problem can be formulated as

*maximize*

$$\sum_i \sum_j \frac{d_{ij}}{4}(x_i - x_j)^2 \tag{23}$$

which can be rewritten, by using $x_i^2 = x_j^2 = 1$, as

*minimize*

$$-\frac{1}{2}\sum_i \sum_j d_{ij} x_i x_j \tag{24}$$

The max-cut problem can then be implemented by designing a network circuit such that the output of each neuron represents each variable $x_i$. As an instance, this section takes up a weighted graph given in Fig. 63a and design a network circuit whose structure is isomorphic to the graph. To implement this problem instance, the network circuit with five neurons is prepared, and with this circuit, vertex $i$ of the problem graph is represented by $i$th neuron ($i = 1$ through 5). The required connection weights $W_{ij}$ between the neurons can be determined as in Fig. 63b by comparing the objective function given by Eq. (24) with the energy function given by Eq. (19). From weight values $W_{ij}$, a set of the coupling capacitances for the circuit construction can be determines by using Eq. (21). The result is given in Fig. 63c. The network circuit with this coupling capacitance set will hereafter be called the sample network.

### 5.5.6.2. Energy Function and Local Minima in the Sample Network
The internal state of the sample network is expressed by a set of five neuron outputs $(x_1, x_2, x_3, x_4, x_5)$, where $x_i$ is 1 or $-1$. For simplicity, represent the set by a code of signs such as $(+ + - - -)$, where $+$ denotes $x_i = 1$ and $-$ denotes $x_i = -1$. The sample network has 32 possible internal states. The value of the energy function calculated from Eq. (19) is plotted in Fig. 64 for all the states. States $(- + + + -)$ and $(+ - - - +)$ are the global minimum and represent the correct solution to the problem (i.e., the maximal cut for the problem graph of Fig. 63a is given by two disjointed sets of vertices 1, 5 and 2, 3, 4). The network can change its internal state through the transition of a Hamming distance of 1. From the second-lowest states $(- + + - +)$ and $(+ - - + -)$ to the global minimum states, there is no transition path of a Hamming distance of 1: therefore, these two states act as a local minimum.

## 5.5.7. Problem-Solving Operation of the Sample Network
For problem solving, it is essential that, starting with a given initial state, the network circuit should converge to its global minimum energy states. To observe the behavior of the sample network, the state transition of the network was simulated by computer.



Figure 63. Instance of the max cut problem and the corresponding connection weights and coupling capacitances: (a) a weighted graph for the problem instance, (b) connection weights $W$ between neurons, and (c) coupling capacitances for the network circuit. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

**Figure 64.** Energy diagram for the sample network circuit corresponding to the problem instance of Fig. 63(a). The notation, such as $(+ + - + -)$, denotes the set of five neuron outputs. States $(- + + + -)$ and $(+ - - - +)$ correspond to the global minimum that represents the correct solution to the problem. The second lowest states $(+ - + + -)$ and $(- + - - +)$ correspond to local minima. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

The behavior of the sample network (simulated results) is as follows. The first is an operation without an simulated annealing process; that is, all the unit-neuron circuits were set at the condition of simple threshold (curve 5 in Fig. 60a). The result is shown in Fig. 65. The network was initially set at state $(+ + + + +)$, then it was allowed to change its state without restraint. After some transition time, the circuit stabilized into a final state. This procedure, a trial, was repeated many times using a different series of random numbers; the results of three trials are illustrated in the figure. The network was sometimes able to converge into the global minimum state $(- + + + -)$ or $(+ - - - +)$ (as shown by number 1), but it frequently became stuck in the local minimum state $(+ - + + -)$ or $(- + - - +)$ and could not reach the global minimum (as shown by numbers 2 and 3).



**Figure 65.** State transition in the sample network without the annealing (computer simulation). The results of three trials are plotted. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

The second result is an operation with simulated annealing. In the annealing, the control-parameter set for the unit-neuron circuits was gradually changed with the advance in time, according to an appropriate schedule. In this experimentation, we changed the control-parameter set gradually, following a cooling schedule given by $V_{dd} = 2.72 + 0.88$ $\exp(-t/50)$ [mV] ($-V_{ss}$ and $V_b$ were also changed in accordance with the curves in Fig. 60b), where $t$ is time in a unit of nanoseconds. The simulation result of the circuit operation is illustrated in Fig. 66. The circuit was initially set at state (+ + + + +), then was allowed to change its state under the annealing operation. The result for a trial is plotted in the figure. The circuit successfully reached the global minimum state. In this way, the correct solution to the max-cut problem can be found.

## 5.6. Analog Computation Based on the Energy-Minimizing Behavior of Single-Electron Circuits

### 5.6.1. What is Analog Computation?

Analog computation is a way of processing that solves a mathematical problem by applying an analogy of a physical system to the problem. To solve the problem in this way, you prepare an appropriate physical system and represent each problem variable by a physical quantity in the system. If the mathematical relations between the physical quantities are analogous to those of the problem, then you can find the solution to the problem by observing the behavior of the system and measuring the corresponding physical quantities. A way of processing based on this principle is called analog computation.

The analog computation is quite different from the commonly used Neumann-Boolean computation. In Neumann-Boolean computation, we first devise an algorithm (a set of instructions for finding the solution to a problem), then execute each step of the algorithm in the manner of Boolean operation, under Neumann computing architecture. In contrast, analog computation is concerned with no symbolic Boolean operation; instead it uses the properties of a physical system to perform the mathematical operations required for the solution. An important feature of analog computation is concurrency or parallelism in computing, and thereby it can provide the possibility of solving complex problems in a short time.

### 5.6.2. An Example of Analog Computation: Solving the Steiner Tree Problem by Using a Soap-film System

Consider the following problem (Fig. 67). Connect $n$ points on a plane with a graph of minimum overall length, using additional junction points if necessary. This is a combinatorial problem called the Steiner tree problem. Plainly expressed, the problem is "to connect $n$ cities by a road network of minimum total length."



Figure 66. State transition in the sample network under the annealing (computer simulation). The result of a trial is plotted. The network can successfully reach the global minimum state. Reprinted with permission from [49], T. Yamada et al., *Nanotechnology* 12, 60 (2001). © 2001, IOP Publishing.

**Figure 67.** Steiner tree problem. Connect given points on a plane with a graph of minimum overall length. This is difficult to solve using existing computers because it requires enormous computing time.

This problem is intractable for digital computation. There are many possible graphs with junction points, and we must examine all the possible ones to find the minimum solution. The number of computational steps required increases exponentially with the number $n$ of original points. Indeed, the Steiner tree problem belongs to the class of NP-hard problems (nondeterministic polynomial–time hard problems). Except for inefficient exponential-time procedures, no algorithm is known for the solution. This problem therefore requires enormous computing time to solve and is virtually unsolvable for large values of $n$.

Nevertheless, there is an ingenious analog-computation method that can quickly solve the problem. This method uses soap films to make a physical system analogous to the problem (Fig. 68) [50]. Prepare two parallel glass plates and insert $n$ pins between the plates to represent the points; then dip the structure into a soap solution and withdraw it. The soap film will connect the $n$ pins in the minimum Steiner-tree graph. The computing process is parallel and instantaneous, so the solution can be obtained in very short time regardless of the number $n$ of the pins.

In this analog computation, the energy-minimizing principle is well utilized for problem solving. Any physical system changes its configuration to decrease its total energy. In liquids at rest, the relevant energy components are the gravitational potential energy and the surface energy. The latter is dominant in a thin soap film, and so a soap-film system changes its configuration to minimize its total area, and therefore its length, and thereby its surface energy.

Strictly speaking, it is not possible to be certain, in this system, that the absolute minimum solution can always be obtained. Depending on the angle at which the system is withdrawn from the soap solution, the soap-film network sometimes assumes topologies different from the optimum one that gives the minimum network length. This arises from the property that many local minima exist in the energy-topology relation of a soap film. Even in such cases, however, the networks obtained are always nearly equal to the minimum one. Hence it can be said that the system works well in general. For other examples of analog computation, see Refs. [51] and [52].

### 5.6.3. Single-Electron Circuit for Solving the Colorability Problem

It is interesting to speculate what analog computations are possible using the properties of single-electron devices. A way is to make use of the property that the single-electron circuit



**Figure 68.** Soap film solution to the Steiner tree problem. The problem can be quickly solved by utilizing the equilibrium of a soap-film system.

changes its state to decrease its free energy. By constructing a single-electron circuit such that its free energy function is related to the objective function of a given combinatorial problem, we will be able to solve the problem simply by observing what state the circuit will settle down. As an example, this section introduces an analog computation device proposed by Tokuda and others [53]. The device is a single-electron circuit that solves a combinatorial problem, the three-colorability problem, in an analog computation manner.

### 5.6.3.1. Three-Colorability Problem

Consider the following problem: Can the countries on a given map be colored with three colors such that no two countries that share a border have the same color (Fig. 69)? This is called the three-colorability problem and is difficult to solve for a map with many countries. There are colorable maps and uncolorable ones, but we cannot tell the colorability of a given map before examining all the possible colorings. (The problem is quite easy if we can use four colors, because it has been proven that four colors suffice for any map.) The three-colorability problem belongs to the class of NP-complete problems (nondeterministic polynomial–time complete) and is intractable for digital computation because only exponential-time algorithms are known for the solution.

This problem is reduced to graph coloring, as shown in Fig. 70. Any map can be converted into a corresponding dual graph by reducing each country to a vertex and drawing an edge between two vertices if the corresponding two countries share a border. Coloring the map is then equivalent to coloring the graph, under the rule that two vertices connected by an edge cannot have the same color. The following text describes a way of solving the three-colorability problem by using single-electron circuits. The task for solution is first to construct a single-electron circuit analogous to a given map for the problem and then to solve the problem by using the circuit.

### 5.6.3.2. Implementing a Dual Graph By Using a Single-Electron Circuit

Let us consider the map of Fig. 69a as an example and construct the analogous single-electron circuit for problem solving. The map can be converted into the dual graph of Fig. 70, so our task is to construct a single-electron circuit that is analogous to the graph.

To represent a vertex of the dual graph, a triangular subcircuit illustrated in Fig. 71a is used. It consists of three identical tunnel junctions $(C_j)$ connected in series to form a ring with three nodes $(A_1, A_2, A_3)$. One excess electron $(e^-)$ is put in the subcircuit, and it occupies one of the three nodes. A ground capacitance $C_0$ exists between each node and ground. Hereafter, this subcircuit with the excess electron is called a triangle, and the excess electron is simply called an electron. We define that the three nodes of the triangle represent three differing colors (e.g., $A_1$ represents red, $A_2$ blue, and $A_3$ green), and that the vertex is colored in the color of the node occupied by the excess electron (e.g., the vertex is colored green if the electron is on node $A_3$).

The first task for solution is to implement two vertices $A$ and $B$ connected by an edge. This is done by coupling two triangles in the manner illustrated in Fig. 71b, using a coupling capacitor $C$ to connect each two nodes that represent the same color. (A ground capacitance exists for each node, but it is not illustrated for simplicity.) The free energy in this coupling

Figure 69. Coloring of a given map with three colors. (a), (b) Colorable maps. The numbers 0, 1, and 2 represent three colors. (c) An uncolorable map. The trial solution shown fails on reaching region F.

**Figure 70.** Dual graph for the map in Fig. 69(a). Each vertex is colored in one of three colors. Two vertices connected by an edge cannot have the same color.

circuit is equal to the electrostatic energy and takes a large value for a state in which two electrons occupy same-color nodes to face each other (e.g., occupying nodes $A_1$ and $B_1$); therefore, the circuit will tend to avoid such single-color states to decrease its energy. In consequence, two electrons in the coupled two triangles will occupy two nodes that represent differing colors (e.g., if the electron in triangle $A$ is on node $A_1$, then the electron in triangle $B$ will be on a node of differing color, either $B_2$ or $B_3$). Connecting six triangles in series produces a circuit analogous to a subgraph consisting of six points with five lines, as shown in Fig. 71c.

A complete circuit analogous to the graph of Fig. 70 can be obtained by connecting six triangles, using 24 coupling capacitors $C$, as illustrated in Fig. 72. (A ground capacitance for each node is omitted in the illustration.) Electrons in neighboring two triangles occupy differing-color nodes, if possible, to reduce the total energy of the circuit; this satisfies the requirement of the three-colorability problem.

It should be stressed that this procedure of constructing analogous circuits can be applied to every other map. For any problem map given, the corresponding analogous circuit can be constructed by combining identical triangles and coupling capacitors.

### 5.6.3.3. Solving the Problem by Using the Constructed Circuit

The three-colorability problem asks whether a given map is colorable, and the answer is either yes or no. This problem is solved with the analogous circuit in the following procedure. Put the circuit in an initial state (any state will do), let the circuit settle down to its equilibrium state with the minimum electrostatic energy, and then check to see whether two electrons in any coupled triangle subcircuits are on nodes of differing colors. If they are, the answer is yes, and colors of the occupied nodes indicate the coloring in which the map can be colored. If they are



**Figure 71.** Construction of a single-electron circuit analogous to the three-colorability problem. (a) A triangular subcircuit representing a vertex of the graph. (b) A circuit analogous to the two connected vertices A and B (for simplicity, the ground capacitances are omitted in illustration). E. Tokuda et al., *Analog Integrated Circuit and Signal Processing* 24, 41 (2000). © 2000, Springer Science and Business Media.

**Figure 72.** Analogous circuit for solving the three-colorability problem for the graph in Fig. 70 (or for the map in Fig. 69(a)).

not, the answer is no. (As for the circuit of Fig. 72, the answer will be yes because the circuit is for the colorable map of Fig. 69a.) This solution is based on the following two principles. First, electrostatic energy in the analogous circuit has a large value when two electrons face each other at neighboring nodes of the same color. Therefore, the circuit will change its state to minimize the number of such electron pairs and, if possible, to reduce such pairs to zero. Second, "a map is colorable" is equivalent to "in the analogous circuit, at least one arrangement of electrons exists such that no two electrons face each other at same-color nodes." (Let us call a state of such electron arrangement the satisfaction state.) In contrast, a circuit for an uncolorable map has no such satisfaction state.

In the minimum-energy state, the circuit for a colorable map is in the satisfaction state, and it will be found that electron pairs in any coupled triangles are on dots of differing colors. In a circuit for an uncolorable map, no satisfaction state can be attained, so it will be found that one or more electron pairs occupy the dots of the same color.

A similar solution using single-electron circuits should exist for other NP-complete problems because every NP-complete problem belongs to the same class, and one can be converted into another.

### 5.6.4. Simulating Circuit Operation of Problem Solving

For the problem solving, it is essential that, starting with a given initial state, analogous circuits should settle down to their minimum-energy state. Unfortunately, analogous circuits in general have many states of locally minimum energy, as shown later, and therefore it cannot be certain, as things stand, that the circuit can achieve the state of globally minimum energy without getting stuck in the local minima. To make the circuit converge exactly to the minimum energy state, the annealing method must be used to operate the analogous circuits successfully. With this method, the global-minimum state will be obtained in most cases, thereby offering the correct solution to the problem. The details will be described in the following sections, with computer simulation.

**5.6.4.1. Energy Function and Local Minima in Analogous Circuits** Electrostatic energy of single-electron circuits is a function of the electron arrangement in the circuit. The analogous circuit given in Fig. 72 has an energy function shown in Fig. 73. The horizontal axis in the figure indicates the number of electron arrangement—one number corresponds to

**Figure 73.** Electrostatic energy versus electron arrangement for the circuit of Fig. 72.

one arrangement of electrons; 729 arrangements are possible because three possible arrangements exist for an electron in each of the six triangles. (In calculation, the circuit parameters were assumed as the tunnel junction capacitance $C_j = 100$ aF, the coupling capacitance $C = 100$aF, and the ground capacitance $C_0$ of each node $= 1$ aF.)

The energy of this circuit becomes the minimum for several specific electron arrangements (the arrangements of numbers 147, 209, 303, 427, 521, and 583; indicated by solid arrows in the figure), which correspond to the satisfaction states representing the correct solution to the graph (or map) coloring. However, the energy function also has many local minima with energy values close to that of the minimum-energy states. It is therefore not possible to be certain that the circuit can always achieve the correct solution without getting stuck in the local minima. (The electron arrangements of numbers 1, 365, and 729, indicated by dashed arrows, correspond to states of monochromatic coloring; i.e., coloring of the graph [or map] with a single color. These states have the maximum energy value.)

The local minima prevent us from finding the solution to the problem. This is shown in Fig. 74, a simulated operation of the analogous circuit given in Fig. 72. The circuit was initially set at the state of monochromatic coloring (any state will do) and then was left changing its state without restraint. After some transition time, the circuit stabilized in a final state; the results of two trials are illustrated in the figure. The circuit may reach the global-minimum energy state by chance, but in most cases, it gets stuck in a local minimum and cannot reach the global-minimum state.

**5.6.4.2. Annealing Operation Method** A way of overcoming the local-minimum difficulty is to operate the circuit with annealing. The annealing process consists of the following four steps:

1. Put the analogous circuit into a heat bath and set the circuit at an initial state (any state will do)
2. Increase initially the temperature of the heat bath to a maximum value at which the circuit changes its state or electron arrangement randomly
3. Carefully decrease the temperature of the heat bath until the circuit arranges its electrons in an equilibrium state (or until the circuit reaches convergence)
4. Check the final arrangement of electrons in the circuit to see whether or not the circuit is in the satisfaction state.

Figure 74. State transition in the analogous circuit of Fig. 72 (computer-simulated). The results of two trials are plotted; in most cases, the circuit gets stuck in a local minimum. The transition of electron arrangement it both trials is illustrated with a set of six triangles that represents the analogous circuit. A dot on each triangle represents an electron on each triangle subcircuit.

If the lowering of the temperature is done slowly enough, the analogous circuit can reach thermal equilibrium at each temperature, and it therefore can approach the global-minimum state with a decrease in temperature. Therefore, the solution to the problem can be obtained by observing the final state of the circuit.

### 5.6.4.3. Convergence to the Minimum Energy State through the Annealing    Figure 75

shows the circuit converging to the global-minimum state with annealing. The cooling schedule (a decrement function for lowering the temperature in annealing) used is a natural cooling given by $T = T_0 \exp(-\rho t)$, where $T$ is the temperature, $T_0$ is an initial value of the temperature, $\rho$ is a cooling-speed coefficient, and $t$ is time. (Parameters $T_0$ and $\rho$ govern the convergence of a circuit during annealing. The values for successful convergence can be inferred by experience from the size of a given analogous circuit.) The circuit was initially set at an monochromatic-coloring state and then was left changing its state under the natural cooling given by $T_0 = 0.5$ mK and $\rho = 0.03$ $\mu s^{-1}$). The result for a trial is plotted in the figure, and the global-minimum state is obtained successfully. State 3 in the figure is an entrance to a transition path that leads to local minima (see state 3a in Fig. 74), but the circuit was able to escape this state because of thermal excitation resulting from annealing. With annealing,



Figure 75. State transition in the analogous circuit with annealing (computer simulation). The transition of electron arrangement is also illustrated. The circuit can successfully reach the global-minimum state.

almost every trial results in the global minimum successfully. In this way, we can find the global minimum state of analogous circuits, and hence the solution to given problems.

## 5.7. Quantum Hopfield Networks Using the Cotunneling Effect

### 5.7.1. Parallel Computation Using Quantum Mechanics

Introducing quantum mechanics into computation may produce the capability for massive parallel processing. The quantum generalization of the Turing machine (or Neumann-Boolean computing), known as the quantum Turing machine [54], or so-called quantum computing, is an example. The quantum Turing machine can perform ultra-high-speed computation because it can accept a coherent superposition of many input data and perform a computation on every input datum simultaneously. This concurrency or parallelism can be used to quickly solve several important problems, such as factoring and discrete logarithms, that are intractable for the classical Turing machine, and therefore for existing computers. Several approaches have been proposed to implement the quantum Turing machine.

Is this type of quantum effect exclusive to the Turing machine? Probably not. There are various other computation models besides the Turing machine (Neumann-Boolean computing), as shown in Fig. 7, and it is likely that the parallelism of each model can be enhanced with the application of quantum mechanics. This section takes up the Hopfield network as an example and shows that quantum parallelism can be obtained in this computation models as well. A quantum version of Hopfield networks, the quantum Hopfield network, may provide novel computation devices that solve combinatorial problems without being troubled by the local-minimum difficulty.

### 5.7.2. The Hopfield Network—A Computation Model for Solving Combinatorial Problems

The Hopfield network is a computation model for solving combinatorial optimization problems. It makes use of the operation of a specific recurrent network (hereafter, we call the recurrent network itself a Hopfield network). The concept of a Hopfield network is shown in Fig. 76. The network consists of threshold elements and connections. The connection weights $W_{ij}$ and $\theta_i$ can be given any desired value, with the restrictions that $W_{ij} = W_{ji}$ and $W_{ii} = 0$. The outputs $V_i$ of the threshold elements $i$ wrap around to become the inputs to the network. Each threshold elements $i$ produces an output 1 if the weighted sum of inputs $(W_{ij}X_j + \theta_i)$ is positive, and an output 0 if the weighted sum of inputs is negative. The state of the network is defined as a set of the outputs $V_i$ of the threshold elements. The point of this network is that, starting at a given initial state, it changes its state to minimize the value of the energy function defined by

$$E = -\frac{1}{2}\sum_{ij} W_{ij}V_iV_j - \sum_i \theta_i V_i \tag{25}$$



Figure 76. Configuration of the Hopfield network.

By adjusting the connection weights, we can relate the energy function of the network to the objective function of a given optimization problem. In this way, we can find the solution to the problem simply by observing the final state that the network reaches. For details, see Refs. [55] and [56]. (The Hopfield network uses the same configuration of network as that of Boltzmann machine in Section 5.5. The difference is that the Hopfield network uses neither stochastic neuron nor simulated annealing algorithm.)

The computation in the Hopfield network is quite different from the commonly-used computation methods. In the common solutions to combinatorial problems, we cannot obtain the solution to a problem until examining all the possible combinations of problem variables; therefore, computing time that is required increases exponentially with the size of the problem. In contrast, in the Hopfield network, a given problem is mapped onto the network itself and is solved quickly through concurrent or parallel operation of all the elements in the network. The Hopfield network, therefore, has the possibility of solving combinatorial problems in a short time, regardless of the size of the problem. This parallelism may provide an efficient way of solving difficult combinatorial problems such as NP-hard problems, which are often encountered in engineering fields but take enormous computing time to solve with digital computers.

Unfortunately, it is not possible to be certain that the correct solution can always be obtained. This is because the Hopfield network in general has many states of locally minimum energy in addition to the globally minimum state. In most cases the network will get stuck in a local minimum, and a solution will not be reached. The computation model of Hopfield networks is based on the premise that the final state of the network can be considered as globally minimum in energy, and without this premise we cannot be convinced that an obtained result is the correct solution. This local-minimum difficulty is an inevitable drawback in the Hopfield network and has limited the application field of the Hopfield network.

The local-minimum difficulty is a natural result of the fact that each event of state transition in the threshold elements is independent of others. The threshold elements update their output states irrelevantly, with no mutual correlation, and consequently the network can make at a time only a limited state transition with a Hamming distance of 1 (i.e., a transition corresponding to the output change of one threshold element). Under these conditions the network cannot escape from a local minimum, even if there are other possible states with lower energy, because an output change of any one threshold element will increase the network energy. This is inevitable as far as we are tied to the classical concept of the Hopfield network.

The way of overcoming this difficulty is to create a special network in which two or more threshold elements can change their outputs simultaneously in a form of coherent combination. In such a network, a transition with a larger Hamming distance (2 or more) can occur, and consequently the global minimum state can be achieved without being troubled by local minima. Such a network can be constructed with single-electron circuits.

### 5.7.3. Constructing a Hopfield Network with Single-Electron Circuits

The single-electron circuit changes its state to decrease its free energy. Akazawa and others proposed making use of this behavior to construct a quantum Hopfield network [57]. Components of the network are shown in Fig. 77. A tunnel junction with an excess electron is used as a threshold element (Fig. 77a); we define the state of the tunnel junction as 1 if the electron is on the right of the junction, and as '0' if it is on the left. The connection between two tunnel junctions can be established by a pair of coupling capacitors; the connection weight can be set positive or negative by choosing the layout of the capacitor coupling (Fig. 77b). The overall configuration of the network is illustrated in Fig. 78. (An excess electron is also set on each bias node.) A ground capacitance exists between each node and a ground (not illustrated for simplicity). A sample set of capacitance parameters is given in the figure. Starting from a given initial position, the circuit changes its internal state (the arrangement of electrons) to minimize its free energy. In this circuit, the free energy is equal to the electrostatic energy and is given in the form of

$$E = A - \frac{1}{2} \sum_{ij} B_{ij} N_i N_j - \sum_i C_i N_i \qquad (26)$$

**Figure 77.** Single-electron Hopfield network: (a) tunnel junction as a threshold element and (b) positive connections and negative ones. Reprinted with permission from [57], M. Akazawa et al., *Analog Integrated Circuit and Signal Processing* 24, 51 (2000). © 2000, Springer Science and Business Media.

where $N_i$ (either 1 or 0) is the state of each tunnel junction, and $A$, $B_{ij}$, and $C_i$ are coefficients. The value of each coefficient can be set at any desired value by selecting the connection pattern and the capacitance of the tunnel junctions, connection capacitors, and ground capacitors. This equation is in essence the same as the energy function [Eq. (25)] of the Hopfield network. In this way we can be certain that the circuit will operate as a Hopfield network.

The internal state of the circuit is expressed by a set of the states of the tunnel junctions. For the sample network circuit shown in Fig. 78, the internal state is expressed as $(N_1, N_2, N_3)$. Figure 79 illustrates the relative energy values of possible internal states (see Fig. 80 for the exact values of each state). The global minimum is state (1, 1, 1). The solid arrows in the figure indicate the possible occurrence of a state transition resulting from one tunneling event, corresponding to the transition with a Hamming distance of 1, which also occurs in classical Hopfield networks (assumed zero temperatures and no energy excitation). For such transitions, states (1, 0, 0) and (0, 1, 0) have several incoming paths but no outgoing path; therefore, these two states seem to be local minima.

However, as described in the following section, the single-electron Hopfield network has quantum properties unlike its classical relatives. Therefore, the said states (1, 0, 0) and (0, 1, 0) do not act as a local minimum state, and the network never gets stuck in these states.

## 5.7.4. Quantum Operation in Single-Electron Hopfield Networks

The transition from state (1, 0, 0) or (0, 1, 0) to state (1, 1, 1) is a transition with a Hamming distance of 2, which can occur only when two threshold elements (tunnel junctions) change their states simultaneously, with mutual correlation. Such a transition is nonexistent in a classical sense, but in the single-electron Hopfield network, it can actually occur through a quantum effect known as the cotunneling phenomenon (or higher-order tunneling).

Cotunneling is a phenomenon in which two or more tunneling events occur simultaneously in the form of quantum coherent combination. In single-electron circuits, two or more



**Figure 78.** Sample configuration of the single-electron Hopfield network. The state of the network is represented by $(N_1, N_2, N_3)$. Reprinted with permission from [57], M. Akazawa et al., *Analog Integrated Circuit and Signal Processing* 24, 51 (2000). © 2000, Springer Science and Business Media.

**Figure 79.** Electrostatic-energy diagram showing possible transitions for the sample network of Fig. 78. The state of the network is expressed as $(N_1, N_2, N_3)$. A solid arrow shows a transition by one tunneling event (Hamming distance is 1). A dashed arrow shows a higher-order transition (Hamming distance is 2). which cannot occur in classical Hopfield networks. Reprinted with permission from [57]. M. Akazawa et al., *Analog Integrated Circuit and Signal Processing* 24. 51 (2000). © 2000, Springer Science and Business Media.

tunnelings can occur simultaneously through the cotunneling if such an event decreases the energy of the circuit. This enables transitions with a larger Hamming distance (2 or more) in single-electron Hopfield networks. Through this phenomenon, the sample circuit in Fig. 78, for example, can change its state from (1, 0, 0) and (0, 1, 0) to the global minimum (1, 1, 1), as shown by the dashed arrows in Fig. 79; thus, the local-minimum difficulty disappears. This kind of Hopfield network is called the quantum Hopfield network. In the quantum Hopfield network, it is certain that, starting at a given initial state, the global minimum state can always be established.

Figure 80 illustrates the behavior of the single-electron quantum Hopfield network, with the results simulated for the sample circuit given in Fig. 78. The cotunneling phenomenon was taken into account by the method of Refs. [18, 58]; the tunnel resistance of the junctions is set at 200 kΩ, and temperature is assumed to be 0 K. The circuit is initially set at maximum energy state (0, 0, 1) and then is left changing its state without restraint. After some transition time, the circuit stabilizes in a final state. This procedure is called a trial, and the results of three trials are shown in the figure.

In every trial, transitions with a Hamming distance of 1 are observed at first; they are denoted by numbers 1 through 6 in the figure. These transitions correspond to those



**Figure 80.** State transition in the single-electron Hopfield network of Fig. 78 (computer-simulated). The results of three trials are plotted. The circuit finally achieves the global-minimum state. Reprinted with permission from [57]. M. Akazawa et al., *Analog Integrated Circuit and Signal Processing* 24, 51 (2000). © 2000, Springer Science and Business Media.

observed in classical Hopfield networks. The network can sometimes converge to minimum energy state (1, 1, 1) through transitions with a Hamming distance of 1 (as shown by numbers 2 and 5 in the figure). but it usually got stuck in the intermediate states (1, 0, 0) or (0, 1, 0) (as represented by numbers 1 and 4 and numbers 3 and 6). The two states would act as local minima if this network were a classical Hopfield network. The situation is, however, quite different in the single-electron Hopfield network. After some waiting time, we can observe the transition from the states (1, 0, 0) or (0, 1, 0) to the global-minimum state (1, 1, 1). This transition is a transition with a Hamming distance of 2 that is induced by the cotunneling phenomenon.

In more complex networks, state transitions with a larger Hamming distance (3 or more) will be required for convergence, but it is certain that these transitions surely occur after some waiting time. In general, single-electron Hopfield networks can always reach the global-minimum energy state, starting at a given initial state. Making use of this property will enable us to develop novel computation devices that solve combinatorial problems without being troubled by the local-minimum difficulty. (Strictly speaking, there is an open question of whether it would be practical to built single-electron network circuits that can generate the cotunneling event frequently enough to deal with any given complicated problems.)

## 5.8. Single-Electron Circuits for Stochastic Associative Processing

### 5.8.1. Concept of Stochastic Associative Processing

Associative processing or associative memory is a function that extracts a pattern similar to the input key pattern from the stored patterns. All the conventional associative memories achieve deterministic association; the same input pattern leads to the same association result. In associations performed in the human brain, however, different outputs are often obtained from the same input. This may be described as chaotic behavior in highly nonlinear systems. As another model of such associative processing, this section describes an unconventional associative processing scheme, that utilizes the stochastic property of single-electron devices effectively.

There are some associative processing models also in the artificial neural network field. One is the associatron, a historical neural network model of associative processor [59], and another famous one is related to Hopfield networks [55, 60]. However, these associative processing models are not considered in this section, but the conventional digital associative processing architecture is used.

In the conventional associative processing, the input pattern is compared with all of the stored patterns, and the stored pattern most similar to the input is deterministically extracted. In contrast, a stochastic associative processor does not always extract the most similar pattern. Instead, the more similar pattern is extracted with a larger association probability that depends on the similarity of the pattern to the input. When this concept is applied to digital data processing, its mathematical formalization is as follows: Let an associative processing device be defined as one that extracts similar patterns to the key pattern $\Psi$ from stored patterns $\Phi_k$ ($k = 1, \ldots, M$), where $\Psi$ is given by the external system and $\Phi_k$ are stored in the device. All patterns consist of $N$ bit binary data:

$$\Psi = \{\xi_j: \ j = 1, \ldots, N\} \tag{27}$$

$$\Phi_k = \{\zeta_{jk}; \ j = 1, \ldots, N, \ k = 1, \ldots, M\} \tag{28}$$

$$\xi_j \in \{0, 1\}, \ \zeta_{jk} \in \{0, 1\} \tag{29}$$

Let $\nu_k$ be defined as the number of unmatched bits between $\Psi$ and $\Phi_k$, which is referred to as a Hamming distance. In addition, let $k_1$ be defined as the suffix of the most similar $\Phi_k$ to $\Psi$, which means that $\nu_{k_1} \leq \nu_k$, $\forall k \neq k_1$, and $k_2$ as that of the second, and so on: that is,

$$\nu_{k_1} \leq \nu_{k_2} \leq \nu_{k_3} \leq \cdots \tag{30}$$

The associative processor described in this section stochastically extracts $\Phi_{k_1}$. It sometimes extracts $\Phi_{k_2}$ or $\Phi_{k_3}$ and so on. If we define the probabilities of extracting $\Phi_k$ as $P_k$, we can

expect

$$P_{k_1} \geq P_{k_2} \geq P_{k_3} \geq \cdots \qquad (31)$$

By repeating numerous extraction trials. we can obtain the order of $k$ in similarity (i.e., $k_1, k_2, k_3, \ldots$).

This concept offers an approach to intelligent information processing that differs from the conventional deterministic approach. Useful and unique examples of this type of processing are sequential stochastic association [61] and clustering for vector quantization [62].

## 5.8.2. Stochastic Associative Processor Architecture Using Nanostructures

On the basis of the processing model described above, a processor architecture is described. A pattern datum is often referred to as a word in digital processing. Therefore, the associative processor has $M$ word-comparators (WCs) and an extractor. as shown in Fig. 81. Each WC has $N$ bit-comparators (BCs). Bit-level comparisons between $\Psi$ and $\Phi_k$ are performed at all BCs in each WC in parallel. Each WC unifies the results from the $N$ BCs and the extractor compares all the outputs of WCs and extracts the most similar word. The function that extracts the most similar data is achieved by a winner-take-all (WTA) circuit.

The bit-comparator can be realized by using an XOR or exclusive-NOR (XNOR) logic operation, where the former outputs logical 1 only if the both inputs have the different logical bits, and the latter outputs 1 only if the both inputs have the same logical bit In these two cases, the output of the WC represents the Hamming distance $v_k$ or $N - v_i$, respectively. The associative processor extracts the stored pattern having the shortest Hamming distance.

To achieve stochastic associative processing, random fluctuation is added somewhere in the circuit. If stochastic behavior of the single-electron devices is used for such random



Figure 81. Architecture of associative memory. Reprinted with permission from [64]. M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998. IEICE.

fluctuation, a BC is a suitable part, because BCs are regularly arranged in the circuit, as described in Section 4.2.

Figure 82 shows an architecture of a stochastic associative processor that uses nano-structures [63]. The BCs consist of nanostructures and perform bit-comparison with random fluctuation. Thus, the input pattern is stochastically compared with the stored patterns by WCs. The comparison result by each WC is expressed as the total number of electrons released from the BCs that belong to the WC. The electrons are collected at each capacitor $(C_{o1}, C_{o2}, \ldots, C_{oM})$. The results of all of the WCs are fed into the WTA circuit, which deterministically extracts the stored pattern evaluated as that most similar to the input. The WTA circuit can be constructed by CMOS devices.

In the following sections, some single-electron circuits and nanostructures for BCs and WCs with random fluctuation are described.

### 5.8.3. Word-Comparator Using SET Logic Circuits

**5.8.3.1. Circuit Configuration** In the WC circuit described here [64], if two inputs are equal: $\xi_j = \zeta_{jk}$, then the binary output of the $j$th BC is 0. Otherwise it stochastically oscillates between {0, 1}. Figure 83 shows an example of the WC including BC circuits composed of SET circuits. The BC circuit consists of two-input CMOS-like SET inverter *INV* and two-input SET switch *SW*. Here, let $I_{jk}$ be defined as the output current of the BC associated with electrons passing through tunnel junction *TJ*. By adjusting the device parameters, this circuit operates as follows.

Input voltages $V_a$ and $V_b$ correspond to data bits $\xi_j$ and $\zeta_{jk}$, respectively. If $V_a = V_b$, then the voltage of node $P$, $V_P$, is always $\overline{V_a}$. Output current $I_{jk}$ is zero in this case. If $V_a \neq V_b$, then an electron goes in and out of node $P$ at random; that is, $V_P$ oscillates. Current $I_{jk}$ oscillates accompanying electrons passing through *TJ* when $V_P = V_a$. Figure 84 schematically shows this situation. The frequency and duty-ratio of the oscillations are related to the probability of the electron transition and depend on the device parameters.

Then, all current outputs $I_{jk}$ of BCs are summed up at capacitor $C_o$, and each WC outputs the result as a voltage. Because capacitor $C_o$ is shorted at intervals of sampling time $t_s$ by switch $SW_1$, output voltage $V_k$ at the end of each interval is

$$V_k = \frac{1}{C_o} \sum_{j=1}^{N} \int_0^{t_s} I_{jk}(t)dt \qquad (32)$$



Figure 82. Architecture of the stochastic associative processor.

Figure 83. Word comparator using single-electron circuit. Reprinted with permission from [64]. M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.

where $C_o$ is set proportionally to $I_s$ to obtain an appropriate voltage of $V_k$. If $\Psi = \Phi_k$; that is, $\xi_j = \zeta_{jk}$, $\forall j$, then obviously $V_k = 0$. When $\Psi \neq \Phi_k$, $\forall k$, to extract the most similar pattern, there are the following two matching methods.

**5.8.3.2. Voltage-Domain Matching** If we set $I_s$ longer than the average period of the oscillation in $V_p$, $V_k$ is statistically proportional to the number of unmatched bits, $\nu_k$, in the voltage domain. Thus, the order of similarity in $\Phi_k$, $k_j$, expressed in Eq. (30), can statistically be obtained by comparing $V_k$ with the ramping reference voltage, as in conventional analog sorting circuits.

The fluctuation of $V_k$ decreases when making $I_s$ long. If $I_s$ is long enough, the order of similarity can be obtained almost deterministically, which is the same as with ordinary associative processor. In contrast, if $I_s$ is set in the order of the average period of the oscillation, the fluctuation of $V_k$ is very large, and association becomes stochastic.



Figure 84. Schematic figure for explaining the relation between $V_p$ and $I_k$, where $V_a = 1$ is assumed. Reprinted with permission from [64]. M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.

**5.8.3.3. Time-Domain Matching** If $t_s$ is set shorter than the average period of the oscillation in $V_p$, $V_k$ also oscillates, and there exist sampling intervals during which $V_k = 0$. The average span between the sampling events at which $V_k = 0$ statistically increases with increases in $\nu_k$ because $V_k = 0$ only if $I_{jk} = 0$ for all $j$. Thus, if the output of the $k_t$th WC, $V_{k_t}$, becomes zero earliest in the consecutive sampling intervals, our associative processor extracts $\Phi_{k_t}$. The probability that $\Phi_{k_t}$ is most similar to $\Psi$ is a maximum in this case.

Let us estimate the time dependence of the probability of detecting $V_k = 0$ in the sampling events with a parameter of $\nu_k$. Consider a BC with different inputs and a typical $V_p$ change as shown in Fig. 85, where $V_a = 1$ is assumed; therefore, $I_{jk} \neq 0$ when $V_p = 1$, and vice versa. Here, $T$ and $a$ are defined as the average period and the ratio of the interval when $V_p = 0$ in oscillation of $V_p$, respectively.

When a sampling event starts within the time span of $aT - t_s$, $I_{jk}$ is always zero in this sampling interval. Because the relation between $V_p$ oscillation and sampling timing is random, the probability that $I_{jk}$ is always zero in a sampling interval is statistically estimated at

$$prob_{jk} = \frac{aT - t_s}{T} = a - \frac{t_s}{T}$$ (33)

In a WC, the probability of detecting $V_k = 0$ is statistically estimated at

$$prob_k = \left(a - \frac{t_s}{T}\right)^{r_k}$$ (34)

The probability that $V_k$ becomes zero at time $t$ for the first time is

$$prob_k(t) = \left[1 - \left(a - \frac{t_s}{T}\right)^{r_k}\right]^{t/t_s}\left(a - \frac{t_s}{T}\right)^{r_k}$$ (35)

Thus, the probability that $V_k$ becomes zero at least once by time $t$ is

$$Prob_k(t) = \sum_{j=0}^{n-1}\left[1 - \left(a - \frac{t_s}{T}\right)^{r_k}\right]^{j}\left(a - \frac{t_s}{T}\right)^{r_k}$$ (36)

$$n = t/t_s \geq 1$$ (37)

Figure 86 shows the time dependence of $Prob_k(t)$ with a parameter of $\nu_k$, where $T = 36$ ns, $a = 0.5$, $t_s = 8$ ns. Thus, it is confirmed from Fig. 86 that the order of similarity in $\Phi_k$, $k_t$ expressed in Eq. (30), is stochastically obtained in the time domain.

**5.8.3.4. Simulation Results** Figure 87 shows the SET circuit simulation results about the dependence of $I_{jk}$ on $V_{a,b}$ in the BC. It can be seen that $I_{jk}$ oscillates randomly when $V_a \neq V_b$. The black areas indicate that $I_{jk}$ oscillates at a high frequency, which is the effect of electrons passing through the switch $SW$. In this simulation, parameters were $V_{dd} = 6.5$ mV, $V_{bias} = 2.3$ mV, $C_a = 55$ aF, $C_b = 55$ aF, $C_{bias} = 10$ aF, $C_1 = 6$ aF, $C_2 = 9$ aF, $C_{s1} = 10$ aF, $C_{s2} = 9$ aF, $C_{t1} = 1$ aF, $C_{t2} = 2$ aF, $C_L = 22$ aF, $R_{t1} = 45$ M$\Omega$, and $R_{t2} = 1$ M$\Omega$, and the ambient temperature was 1 mK.

Figure 88 shows the waveforms of $I_{jk}$ when $\xi_j \neq \zeta_{jk}$, $j = 1, 2, 3$. and detection timing depending on $\nu_k$. Here, $t_s = 10$ ns is assumed. If only the first bit ($j = 1$) is unmatched, the



**Figure 85.** Schematic figure introducing the probability of detecting $I_{jk} = 0$, where $V_a = 1$ is assumed. Reprinted with permission from [64]. M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.

**Figure 86.** Time dependence of $Prob_k(t)$ with a parameter of $v_k$. Reprinted with permission from [64], M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.



**Figure 87.** Simulation results about the dependence of $I_{jk}$ on $V_{a,b}$ in the bit-comparator. Reprinted with permission from [64], M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.

**Figure 88.** Waveforms of $I_{jk}$ and detection timing depending on $\nu_k$. Reprinted with permission from [64], M. Saen et al., *IEICE Trans. Electron.* E81-C, 30 (1998). © 1998, IEICE.

first time when $V_k = 0$ is $T_1$. If the first and second bits ($j = 1, 2$) are unmatched, the first time when $V_k = 0$ is $T_2$. If the first to third bits ($j = 1, 2, 3$) are unmatched, the first time when $V_k = 0$ is $T_3$. We can apparently see that $T_1 < T_2 < T_3$.

As described above, a WC including BCs with random fluctuation can be constructed by using SET circuits.

## 5.8.4. Word-Comparator Circuit Using Nanodots on a MOSFET Gate

### 5.8.4.1. Simple Single-Electron XNOR Gate for a Bit-Comparator
Another type of WC that has a simpler structure than that in the previous section is described here. A nano-structure image is also described [63].

As shown in Fig. 11 (Section 5.1), an XOR gate can easily be constructed by using the SET inverter. Furthermore, if dynamic operation is assumed, an XOR gate can be constructed by using a SET and a capacitor, as shown in Fig. 89a. Here, $V_{Co}$ is assumed to be reset at a certain voltage level before operation. Electrons pass through the SET to the capacitor $C_o$ only when $V_a = V_b = L$ or $H$, although the output voltages in both cases are not exactly equal, as shown in Fig. 89b. To equalize both output voltages, a complementary configuration of SETs (CS) is introduced, as shown in Fig. 89c. Figures 89d and 89e show the waveforms of $V_{Co}$, where $C_o$ is 500 aF and 5000 aF, respectively. In both simulations, the operations are repeated 50 times when the same inputs are applied, and all waveforms are represented by the gray lines in the figures.

Because of the stochastic event in each electron tunneling, the rise time of voltage $V_{Co}$ fluctuates in every operation, even if the same inputs are applied. The fluctuation shown in Fig. 89e is smaller than that shown in Fig. 89d because the stochastic characteristic is eventually averaged as the number of electrons accumulating in the capacitor increases. Thus, the XNOR gates with larger $C_o$ operate more deterministically, and those with smaller $C_o$ operate more stochastically. Larger $C_o$ obviously leads to a longer circuit-delay time, and a $C_o$ that is too small makes integrated circuit design difficult. The appropriate capacitance value depends on the application.

Although operation temperature of 1 mK was assumed in the simulations, it was verified that the circuit with the same parameters as shown in Fig. 89a can operate properly up to around 3 K. Moreover, it was also verified that the circuit can operate up to around 100 K if the capacitance values except $C_o$ is reduced to one-tenth, where thermal noise is used for stochastic operation.

**Figure 89.** Dynamic exclusive-NOR gate using a SET and a capacitor: (a) basic circuit configuration and (b) transient behavior of $V_{co}$; (c) the gate using complementary configuration and (d)(e) its stochastic behavior (simulation results of 50 trials).

To avoid misunderstanding, it should be noted that the circuits shown in Fig. 89 are not logic gates in the usual sense of this term. It is difficult to use them in multistage logic circuits because they should be externally reset after each clock cycle; they cannot sustain arbitrary logic levels for long time because of the cotunneling effect; and their output voltage is lower than the input voltage.

**5.8.4.2. Stochastic Associative Processing Circuit Using SETs**  A BC with random fluctuation can be constructed by using the CS shown in Fig. 89c. The WC is constructed by connecting nodes $N_c$ of the plural CS' with the common output capacitor $C_o$. The number of CS' that pass electrons to the capacitor $C_o$ increases as the Hamming distance decreases, resulting in a shorter rise time for the voltage $V_{Co}$. However, even when the Hamming distance is at its shortest, the output voltage does not always increase rapidly because of the stochastic characteristics in CS' as shown in Fig. 89d or e.

In this architecture, the extractor (WTA) circuit selects the word comparator output that reaches the threshold voltage most rapidly, as shown in Fig. 90, which means that it selects the stored pattern having the shortest Hamming distance. Because the output capacitance



**Figure 90.** Dynamic winner-take-all operation.

$C_o$ is large enough (e.g., 500 aF), the input stage of the WTA circuit can be constructed by sub-100-nm MOS devices, whose gate capacitance is less than 1 fF.

A device structure image of the WC is shown in Fig. 91. Isolated islands of SETs are regularly arranged on a capacitor plate, which corresponds to an electrode of $C_o$, and also a gate electrode of an ultra-small MOS device, where the gate capacitance acts as $C_o$.

### 5.8.4.3. Simulation Results

The basic association operation was confirmed by simulation of a digit pattern association. Each stored pattern consisted of seven segments and represented a digit number of $\{0, 1, \ldots, 9\}$. Because the ON/OFF state of each segment corresponded to bit data, the stored data consisted of 10 vectors, each of which had seven binary elements. The simulation of the system was performed as follows: input bit data $V_a$ and stored data $V_i$ ($i = 1, 2, \ldots, 9$) were supplied to BCs as voltage signals; operation of the WC that consisted of single-electron devices was simulated by the Monte Carlo simulator; and operation of the WTA circuit was simulated as an ideal black box that simply selected a winner from outputs of the word comparators.

Two examples of output-voltage changes in the word comparators when the input pattern was 5 are shown in Fig. 92a. Because of the stochastic property, different patterns become winners; one is the pattern most similar to the input pattern 5, and the other is a second similar pattern 6. Figure 92b shows the association probabilities of all stored patterns where the input pattern is 5, and Fig. 92c is another simulation result, where the input pattern is not included in the stored patterns. Both simulation results show that the association probability of the stored pattern increases as the Hamming distance from the input pattern decreases.

## 5.8.5. Bit-comparator Using Coulomb Repulsion in Nanodots

Two types of nanodot circuits measuring a Hamming distance using the Coulomb repulsion effect are described. They have structures where a nanodot array is arranged on a gate electrode of an ultrasmall MOSFET. The device parameters for successful operation are clarified from Monte-Carlo simulation.

### 5.8.5.1. Bit-Comparison Principle Using Coulomb Repulsion in Nanodots

Let us assume a string of nanodots, as shown in Fig. 93a, put an electron $e_M$ at one of the three dots $D_1$, and represent a bit (0 or 1) of the input and stored data by whether an electron is put at each end dot $D_2$ or not. When the corresponding bits of both data are matched, because Coulomb repulsion is symmetric, electron $e_M$ is stabilized at the center; otherwise, it is off-center. To detect the position of electron $e_M$, there are two possible detection circuits: circuit A that detects the center position, and circuit B that detects the off-center position, as shown in Fig. 93b and 93c, respectively. By the Coulomb repulsion effect, the bit-matching result reflects whether electron $e_R$ tunnels to node $N_o$. Consequently, electrons whose number is equal to that of the matched bits (circuit A) or that of the unmatched bits (circuit B) are accumulated in $C_o$. To ensure the stabilizing processes, control voltage $V_g$ is used.

A three-dimensional arrangement of nanodots realizing these circuits is shown in Fig. 94, where common terminals to the ground and control voltages are omitted. The capacitance $C_o$ corresponds to a gate capacitance of an ultrasmall CMOS transistor.

### 5.8.5.2. Simulation Results

Single-electron circuit simulation was performed, where parasitic capacitance between ground and nanodots $C_g$, and that between the second-neighbor



Figure 91. Device structure image of the word comparator.

**Figure 92.** Simulation results: (a) output voltage changes in word comparators, (b),(c) association probablities for stored patterns, where the association is repeated 100 times.



**Figure 93.** Principle of bit-comparison using Coulomb repulsion (a), and two word-comparator circuits (b)(c).

**Figure 94.** 3-D structure images of a word-comparator.

nanodots $C_d$ are considered. Because it was found that circuits A and B have almost the same characteristics, only the simulation results for circuit A are described here.

The operation temperature ranges for feasible capacitance values are shown in Fig. 95a. The upper limit of operation temperature gradually lowers with increasing $C_g$ and $C_d$. To operate the circuit at higher temperature, one has to scale down all the values of capacitance and, at the same time, scale up the applied voltages. For room-temperature operation, tunnel junction capacitance of 0.01 aF is required, as shown in Fig. 95b, although this value is very difficult to realize.

Setting the margin of $C_o$ is shown in Fig. 95c as a function of the word length (the number of the connected BCs). There exist minimum values of $C_o$ for the correct operation, and the values increase as the word length increases. This is because some electrons $e_R$ cannot tunnel to node $N_o$ because of the Coulomb blockade effect by other $e_R$s. If the parasitic capacitance is negligible, there do not exist the upper limits of $C_o$. However, $C_o$ should be as small as possible because the sensitivity to one electron in the output voltage $e/C_o$ decreases with increasing $C_o$.



**Figure 95.** Simulation results for circuit-A: (a) operation temperature range with one BC, (b) output voltage changes at room temperature, and (c) setting margin for $C_o$.

## 5.8.6. Multinanodot WC Circuit and Structure Using Thermal-Noise Assisted Tunneling

The nanodot WC circuits described in the previous section operate only at very low temperature for practical junction capacitance. The multinanodot WC circuit and structure described in this section can operate even at room temperature with a junction capacitance around 0.1 aF by using tunneling processes assisted by thermal noise [20]. In the stochastic associative processing operation, the association probability distribution can be controlled by changing the detection timing of the electron position.

### 5.8.6.1. Circuit and Structure

Let us assume nanodot structures constructed on a MOS transistor gate electrode, as shown in Fig. 96. Each nanodot structure consists of a pair of one-dimensional dot arrays: $A_r$ ($D_{r1}$, $D_{r2}$, $D_{r3}$) and $A_h$ ($D_1, \ldots, D_n, D_c, D_n, \ldots, D_1$), where $n$ is the number of dots at a side of $A_h$, and it should be more than four for the proper operation described below. The array $A_h$ has dot $D_c$ outside of each end. The capacitance $C_0$ corresponds to the gate capacitance of an ultrasmall MOS transistor. A bit (1 or 0) of the input and stored data is represented by whether or not an electron is placed at each end dot $D_c$. (Alternatively, an appropriate voltage corresponding to a bit may be applied directly at $D_c$). Bias voltages are applied to the plate $P_c$ over $D_c$ ($V_{pc}$), to the nodes outside of $D_c$ ($V_c$), and to the back gate of the MOS transistor ($V_{bg}$).

An electron $e_M$ is introduced at the center dot $D_c$ of the one-dimensional array $A_h$. Electron $e_M$ can move along array $A_h$ through tunnel junctions $C_t$, but it cannot move to either of $D_c$s or to $D_{r3}$ through normal capacitors $C_1$ or $C_2$. Each nanodot structure works as an XNOR logic gate (BC) with random fluctuation, as explained below.

### 5.8.6.2. Stabilization Process

By applying appropriate bias voltages $V_{pc}$, $V_c$, and $V_{bg}$, the profile of the total energy as a function of the position of $e_M$ along the one-dimensional array $A_h$ has a minimal value at $D_c$, as shown in Fig. 97. For 1-1 state, where electrons are placed at



Figure 96. Multinanodot circuit and a structure image. Reprinted with permission from [20]. T. Morie et al., J. Nanosci. Nanotech. 2, 343 (2002) © 2002, American Scientific Publishers.

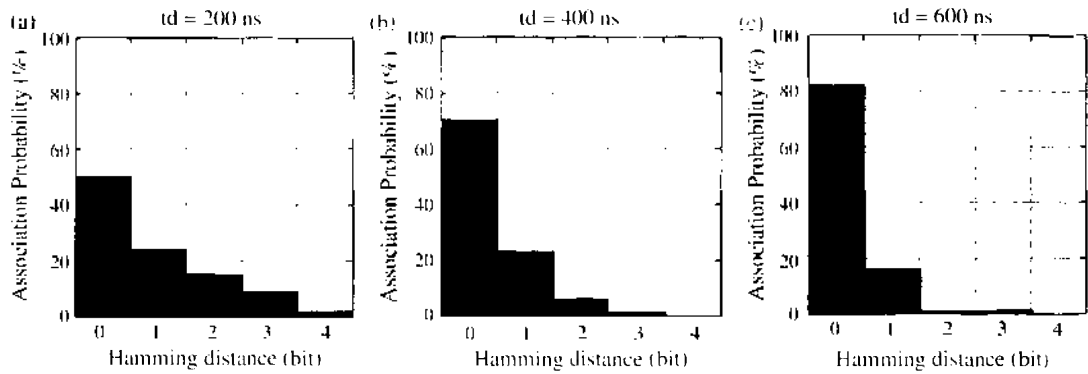Figure 97. Schematics of total energy profile of 1-D dot-array structure. Reprinted with permission from [20], T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002). © 2002. American Scientific Publishers.

both $D_c$s. the energy at $D_1$ rises, and thus $e_M$ is most strongly stabilized at the center position. Therefore, the difference between 0-0 state and 1-0 (or 0-1) state is important for correct BC operation. In the two states, the energy profile has another minimal value at $D_1$. The energy barrier height for $e_M$ located at $D_c$ is approximately determined by the total capacitance for $e_M$ and bias voltages. The greater number of serial capacitance connections causes higher energy barriers, and the energy differences can be much larger than the thermal energy at room temperature, even if the tunnel junction capacitance $C_j$ is around 0.1 aF.

The energy barrier at the 0 side in the 1-0 (0-1) state becomes lower than that in 0-0 state because of the Coulomb repulsion force of the electron placed at the opposite $D_c$, as shown in Fig. 97. Thus, $e_M$ in the 1-0 (0-1) state can more easily overcome the barrier when assisted by thermal noise at nonzero temperature. and it then moves to $D_1$ at the 0 side. As a result. there exists a certain time span $t_0$ within which $e_M$ in 1-0 (0-1) state moves to $D_1$, whereas $e_M$ in 0-0 state stays at $D_c$.

**5.8.6.3. Detection Process**   After spending $t_0$, the vertical dot array $A_v$ detects whether or not $e_M$ stays at $D_c$. by changing the bias voltages if necessary. Only if $e_M$ stays at $D_c$ is array $A_v$ polarized and an electron induced at the gate electrode of $C_g$. (To achieve stable polarization, at least three dots are required in $A_v$). The total number of electrons induced at the gate electrode is proportional to the number of matched bits; this reflects the gate voltage $V_g$, and it can be measured by the source-drain current of the MOS transistor. Thus, the Hamming distance can be measured by this MOS transistor with nanostructure arrays.

If this detection process starts just after $t_0$, the most accurate bit comparison operation is achieved, although some statistical fluctuation remains. However, if the detection timing ($t_d$) is shifted from $t_0$, an arbitrary amount of fluctuation can be introduced in the bit comparison result. Thus, controlled stochastic association can be achieved, which is necessary to apply the stochastic association model to various types of intelligent information processing effectively.

The bit-comparator circuit composed of nanodot arrays works only at nonzero temperatures because at 0 K, $e_M$ can never escape from $D_c$, the valley of the energy profile. Furthermore, the time span $t_0$ depends on the operating temperature. Conversely, for a given $t_0$, an appropriate amount of thermal noise is required; if the thermal noise is too small, electron $e_M$ in 1-0 (0-1) state cannot escape from $D_c$, and thus no bit comparison is achieved. In contrast, if thermal noise is too large, $e_M$ in both the 0-0 state and the 1-0 (0-1) state escapes from $D_c$, and thus the two states cannot be distinguished. In this sense, it can be considered that this circuit utilizes a stochastic resonance effect [65] by thermal noise.

### 5.8.6.4. Simulation Results

We analyzed the proposed circuit shown in Fig. 96 by using a Monte Carlo single-electron simulator, where the tunnel junction capacitance $C_j$ is 0.1 aF and tunnel resistance $R_t$ is 5 M$\Omega$, and other parameters are shown in Fig. 96. In this case, the dot diameter is assumed to be around 1 nm. The bias voltages applied were $V_{pc} = 0$ V, $V_t = 1.15$ V, $V_{bg} = 0$ V for the $e_M$ stabilization process, and $V_{pc} = 0$ V, $V_r = 1.8$ V, $V_{bg} = 3$ V for the $e_M$ position detection process.

Figure 98 shows the total energy profiles at the 0 state side of $A_h$ as a function of the position of $e_M$, where the energy when $e_M$ is located at $D_c$ is defined as zero. It is confirmed that the barrier height for $e_M$ at $D_c$ is larger than the thermal energy at room temperature (i.e., 26 meV), and the barrier height in 1-0 (0-1) state is lower than that in 0-0 state.

Figure 99 shows the relationship between operating temperature and time ($t_M$) required until $e_M$ moves to $D_1$. The closed and open circles indicate $t_M$ in many trials at the 1-0 state and 0-0 state, respectively. Because the moving process assisted by thermal noise is purely stochastic, $t_M$ scatters over a wide range. However, time span $t_0$, defined in the previous section, can be determined from these simulation results.

To determine $t_0$ precisely, we measured $t_M$ for 100 simulations for 1-0 and 0-0 states with different seeds for random number generation. The 100 data obtained for $t_M$ were sorted in increasing order and numbered from 1 to 100. The assigned number indicates the number of electrons that move to $D_1$ within the corresponding $t_M$. Therefore, the relationship between $t_M$ and the assigned number can approximately be considered as the probability that $e_M$ moves to $D_1$ as a function of time. The results obtained at room temperature (300 K) are shown in Fig. 100. The optimum $t_0$ is obtained as the time having the smallest overlap between the two states; it is about 1 $\mu$s. It should be noted here that $t_0$ depends on tunnel resistance $R_t$. If lower tunnel resistance is available, $t_0$ becomes shorter.

From Fig. 100, we can obtain the probability of wrong detection; that is, the conditional probability that a given 1-0 state is detected as a 0-0 state or vice versa. For example, when the detection timing is $10^{-7}$ s, the probability that $e_M$ moves to $D_1$ at 1-0 state is only 20%. This means that wrong detection occurs with a probability of 80%.

By using this effect, we can add fluctuation to the bit comparison operation. However, in the above case, it must be noted that fluctuation can be added only in the 1-0 state. Bit-matched (1-1 and 0-0) states always answer correctly. Therefore, a comparison between patterns with a shorter Hamming distance is performed more deterministically. This means that a stochastic association operation cannot be achieved. An easy way to overcome this difficulty is to reverse the input bit pattern. Although this leads to a deterministic comparison between patterns with a longer Hamming distance, such patterns are seldom associated, and thus it hardly affects the stochastic association operation.



Figure 98. Energy profiles for electron $e_M$ at the 0 state side in 1-D dot array structures. Reprinted with permission from [20], T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002), © 2002, American Scientific Publishers.

**Figure 99.** Relation between operating temperature and time when $c_M$ moves to $D_1$. Reprinted with permission from [20]. T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002). © 2002. American Scientific Publishers.

Figure 101a–c show association probability distributions as a function of the Hamming distance for some $t_d$. In these simulations, the input pattern was (1,1,1,1), and the stored reference patterns were (1,1,1,1), (1,1,1,0), (1,1,0,0), (1,0,0,0), and (0,0,0,0). These reference patterns were reversed when they were applied to the multinanodot circuit, for the reason described above. With the number of electrons indicating the results of Hamming distance evaluation with fluctuation, the reference pattern having the smallest evaluation result became the winner. If two or more reference patterns had the same number of electrons, we determined the winner stochastically. The simulations were repeated 100 times with different seeds for random number generation. The number of trials when a given reference pattern becomes a winner is approximately proportional to the probability that it is associated. The simulation results shown in Fig. 101 confirm that as $t_d$ becomes further apart from $t_0$ (= 1 $\mu$s in this case), the association probability distribution becomes flatter. Thus, the association probability distribution is controlled by changing $t_d$.



**Figure 100.** Probability that $c_M$ moves to $D_1$ as a function of detection timing. Reprinted with permission from [20]. T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002). © 2002. American Scientific Publishers.

**Figure 101.** Association probability distribution as a function of Hamming distance for various detection timing $t_d$. Reprinted with permission from [20], T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002). © 2002, American Scientific Publishers.

Figure 102 shows the time dependence of voltage $V_o$ as a parameter of the Hamming distance for a 4-bit word comparator at room temperature, where the voltage for a distance of 0 bits is defined as 0 V. The voltage changes are proportional to the Hamming distance, and the voltage difference per bit is larger than 1 mV, which is large enough to detect with a CMOS circuit.

## 5.8.7. Some Remarks about Nanostructures for Stochastic Associative Processors

In these architectures described above, if it is difficult to represent 1 bit by 1 nanodot, a redundant architecture can be applied easily, in which plural BCs represent 1 data bit. Such an architecture based on a majority decision principle has advantages of robustness against effects of random background charge.

For realizing these nanostructures described above, the basic technology of nanocrystalline floating-dot MOSFET devices, which are closely related to these nanostructures, has been reported [19, 66, 67]. Fabrication technology using self-organization may also be applied [68]. Furthermore, well-controlled self-assembly processes using molecular manipulation technology, especially that using DNA [69], would be utilized to fabricate the nanostructures.

Although the use of digital data is assumed in the above sections, analog data can be treated in the same circuit by using pulse-width modulation (PWM) signals, which have a digital amplitude and an analog pulse width [70]. Instead of a Hamming distance, a Manhattan distance, the summation of the absolute value of difference, is evaluated by using these



**Figure 102.** Time dependence of voltage $V_o$ as a parameter of the Hamming distance. Reprinted with permission from [20], T. Morie et al., *J. Nanosci. Nanotech.* 2, 343 (2002). © 2002, American Scientific Publishers.

nanostructures. The clustering algorithm using stochastic association for vector quantization [62] can use this distance evaluation approach.

## 5.9. A Multinanodot Floating-Gate MOSFET Circuit for Spiking Neuron Models

### 5.9.1. Neural Network Models

To realize brain-like information processing functions, such as association, perception, and recognition, one very important and challenging approach is to mimic brain functions and structures. In the brain, a neuron receives many electric impulses via a few thousand synapses, and it outputs spike pulses. A typical neuron has three parts: dendrites, a soma, and an axon. Pulse signals called spikes are fed into the dendrites via synapses, the effects of the inputs are gathered up at the soma, and a spike pulse is output from the axon [71]. To analyze and apply neuronal information processing, it is essential to model real neurons. In creating neuron models, however, there are different tradeoffs between faithful simulation of biological reality and minimization of computational expense.

Until the mid-1990s, pulse-rate coding models, which use analog values as averages of pulse events, were studied intensively in both artificial neural networks [72] and their VLSI implementation. The approaches to analog VLSI implementation treat such analog values directly in the voltage or current domain [73–75], whereas the approaches to digital VLSI implementation represent such analog values only by a set of digital bits [76–78]. Another approach, called pulse-stream neural networks or pulse-density modulation (PDM), also focuses on the rate of pulse events [79, 80].

In contrast, since the mid-1990s, computational neuroscientists have focused on more-realistic spiking neuron models, which treat spike pulses directly [71, 81]. In principle, the computational power of spiking neuron models is superior to that of the conventional rate coding models [71, 82].

In this section, a single-electron circuit based on a multinanodot floating-gate MOS device for spiking neurons is described [83].

### 5.9.2. Spiking Neuron Models

In spiking neuron models, information is represented by spatiotemporal patterns in spike pulse trains. A simple spiking neuron model, called the Spike Response Model (SRM) [71], is shown in Fig. 103. A spike pulse inputted to a neuron via a synapse generates a postsynaptic potential (PSP). There are two types of synapses: excitatory and inhibitory. Respectively,



**Figure 103.** Schematic of a simple spiking neuron model, the Spike Response Model (SRM). Reprinted with permission from [83], T. Morie et al., *IEEE Trans. Nanotechnology* 2, 158 (2003). © 2003, IEEE.

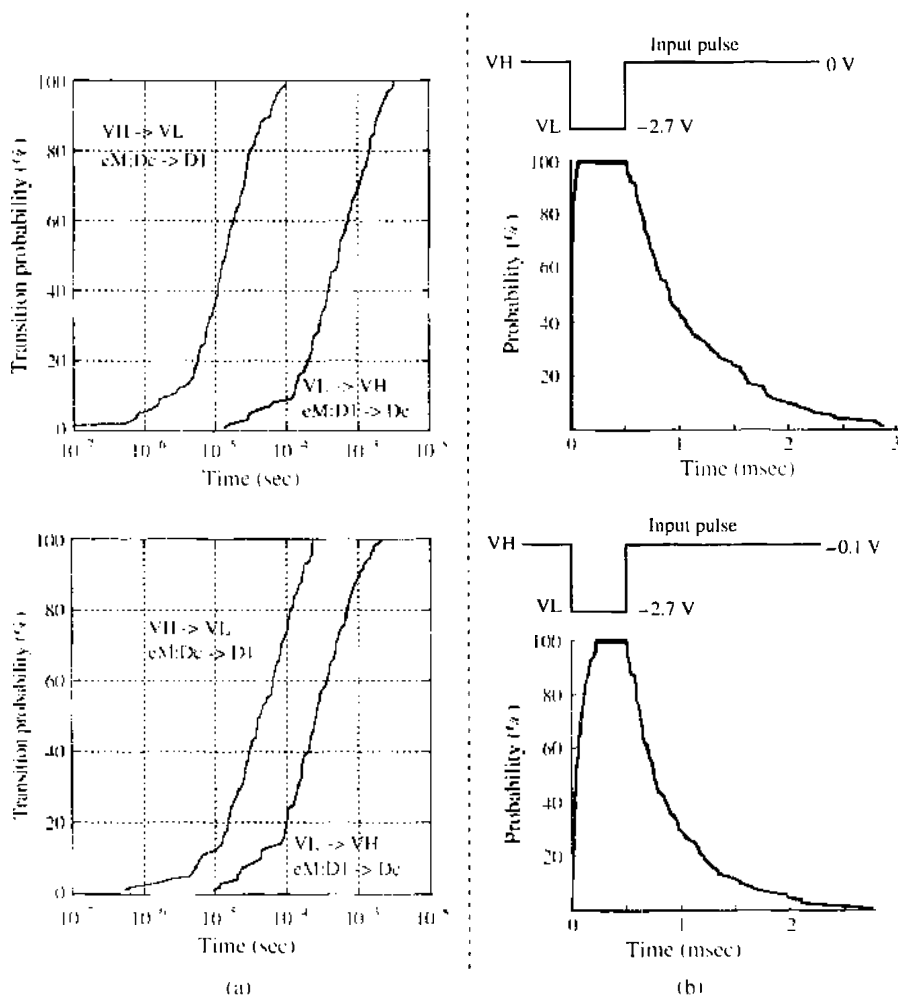these have positive and negative synaptic weights and generate an excitatory PSP (EPSP) and an inhibitory PSP (IPSP). The PSP temporarily increases or decreases according to whether the synaptic connection is excitatory or inhibitory, respectively. The typical time course of a PSP is approximated by a so-called $\alpha$-function $\propto x \exp(-x)$, where $x = (t - t_0 - \Delta^{ax})/\tau_s$, $t - t_0 > \Delta^{ax}$, $t_0$ is a firing time, $\Delta^{ax}$ is the axonal transmission delay, and $\tau_s$ is a time constant [71]. In this function, the rise time and the fall time are related to each other; however, for various applications it is desirable for them to be independently changeable.

The neuron's internal potential $I(t)$ is equal to the spatiotemporal summation of all PSPs generated by the input spike pulses. If $I(t)$ exceeds a certain threshold ($th$), the neuron emits a spike pulse and $I(t)$ is reset to the resting level. There exists a refractory period after firing, in which the neuron cannot fire, even if many spikes are inputted.

The property of the spiking neural network consisting of SRM neurons depends on the relation between the average of interspike intervals $t_{ISI}$ and the PSP decay time constant $\tau_{PSP}$. As shown in Fig. 104, if $\tau_{PSP} \gg t_{ISI}$, then PSPs are integrated during $\tau_{PSI}$, in which case the neuron acts as an *integrator*. In contrast, if $\tau_{PSP} \ll t_{ISI}$, only those PSPs generated by spike pulses inputted within the time interval $\tau_{PSP}$ are integrated. In that case, the neuron acts as a coincidence detector of spike timing. That is, the neuron detects whether or not plural input spikes arise within a certain time span comparable to $\tau_{PSP}$.

The latter type of neuron enables higher-order and faster intelligent information processing functions such as the brain might perform. It is indicated that the computational power of networks composed of spiking neurons is superior to that of conventional analog neural networks based on the rate coding models [71]. For example, spiking neurons can detect temporal patterns irrespective of a common additive constant, and they can compute weighted summations. Furthermore, they can approximate any continuous functions more efficiently.

The spiking neural networks can process various types of information by using various PSP decay constants. This means that rate-coding-type integrators (conventional analog neural networks) and coincidence detectors can coexist in the same network. The two features can be combined for various applications. Therefore, $\tau_{PSP}$ must be controlled in the VLSI implementation of SRM neurons.

### 5.9.3. Multinanodot Floating-Gate MOS Device for Spiking Neurons

Because one real biological neuron has a few thousand synapses, the key for very large scale integration of neural networks is the design of a small synapse circuit. Thus, the realization



Figure 104. Schematic explaining that the spiking neural network property depends on the relation between the average of interspike intervals$t_{ISI}$ and the PSP decay time constant $\tau_{PSP}$. (a) $\tau_{PSP} \gg t_{ISI}$. (b) $\tau_{PSP} \ll t_{ISI}$. Reprinted with permission from [83], T. Morie et al., *H.El. Trans. Nanotechnology* 2, 158 (2003). © 2003, IEEE.

of the synapse functions should be more focused. The synapse circuit of SRM neurons has to implement a generation of controllable PSPs, as well as a synaptic weighting.

A nanodot circuit and structure for realizing the synapse of the SRM model is shown in Fig. 105, which is very similar to the multinanodot structure shown in Fig. 96 in Section 5.8.6. The circuit parameters used in the following simulations are indicated in this figure. These parameters are assumed as an example in the case of the minimum size for nanometer-scale structures. Multinanodot structures are constructed on a MOSFET gate electrode. Each nanodot structure consists of a one-dimensional dot array: $A_h$ $(D_1, \ldots, D_n, D_c, D_{n'}, \ldots, D_1)$, where $n(= n')$ is the number of dots at a side of $A_h$. In the following simulations, $n = 5$ is assumed. The center dot $D_c$ is capacitively coupled with the MOS gate via one nanodot $D_v$. The capacitance $C_a$ corresponds to the gate capacitance of the MOSFET. It is assumed that only one electron $e_M$ exists in the array $A_h$. Electron $e_M$ can move along array $A_h$ only through tunnel junctions, and it cannot move outside of $A_h$ through normal capacitors $C_1$ or $C_2$.

Spike pulses are inputted at nodes $IN$ that are capacitively coupled to one end of $A_h$. To achieve a high probability that $e_M$ exists at center dot $D_c$ when no spike pulse is inputted, appropriate bias voltages are applied to the top plate $P_c$ $(V_p)$, to the end electrode $P_c$ $(V_c)$, and to the back gate of the MOS transistor $(V_{bg})$. For example, assume that $V_p = V_{bg} = 0$ V and $V_c \leq 0$ V. The baseline voltage of the input signal $V_{ll}$ is also set at 0 V or slightly less than 0 V. Thus, the profile of the total energy as a function of the position of $e_M$ along one-dimensional array $A_h$ has two peaks, as shown in Fig. 106, because of the charging energy of $e_M$ itself. The minimal values of the energy are located at $D_c$, $D_1$, and $D_1$.



Figure 105. Single-electron circuit and nanostructure for the SRM neuron. The circuit parameters used in the simulations are also indicated. Reprinted with permission from [83], T. Morie et al., *IEEE Trans. Nanotechnology* 2, 158 (2003). © 2003, IEEE.

Figure 106. Schematics of the total energy profile in 1-D nanodot array and simulation results, where $I_p = V_c =$ $V_{hc} = V_{tt} = 0$ V, and $V_t = -2.7$ V. Reprinted with permission from [83], T. Morie et al., *IEEE Trans. Nanotechnology* 2. 158 (2003). © 2003, IEEE.

Figure 106 also shows simulation results. The energy barrier height for $e_M$ staying at $D_c$ is approximately determined by the total capacitance for $e_M$ and the bias voltages. The energy differences can be larger than the thermal energy at room temperature if the tunnel junction capacitance $C_j$ is around 0.1 aF.

Thus, in the stationary state without input pulses, electron $e_M$ almost always stays at the center dot $D_c$, although thermal noise sometimes causes $e_M$ to move to edge position $D_1$ and back to $D_c$. When a spike pulse is inputted, if it has appropriate pulse width and amplitude voltage, $e_M$ moves quickly to $D_1$. Then, after the pulse signal ceases, $e_M$ moves slowly back to the center dot $D_c$ because of thermal-noise-assisted tunneling through the energy barrier between nanodots. This behavior of $e_M$ creates a PSP. The energy barrier can be larger than the thermal energy at room temperature, and thus the one-dimensional nanodot array can operate at room temperature.

The position of $e_M$ affects the gate voltage $V_o$ through dot $D_c$ (i.e., when $e_M$ stays at center dot $D_c$, $V_o$ is reduced because of capacitive coupling between $D_c$ and the gate electrode). Because the above process occurs stochastically, $V_o$ fluctuates because of the position of $e_M$. However, if one adopts a configuration in which the same input pulse is applied to a number of arrays, or if one uses a low-pass filter in the detection part, an analog PSP is obtained

**Figure 107.** Neuron circuit consisting of a differential-pair including n-MOSFETs with nanodot synapses. Inputs for positive synaptic weights are fed into $IN^+$, and those for negative weights are fed into $IN^-$. The typical number of these inputs is on the order of 100, and the size of each n-MOSFET's gate electrode is around 20 nm by 5000 nm. Reprinted with permission from [83], T. Morie et al., *IEEE Trans. Nanotechnology* 2, 158 (2003). © 2003, IEEE.



**Figure 108.** Electron transition probability as a function of time with baseline voltages $V_{H} = 0$ and $-0.1$ V: (a) transition probability of $e_{M}$ between $D$ and $D_1$ as a function of time, and (b) time dependence of the probability that $e_{M}$ stays at $D_1$ when a pulse is applied. Reprinted with permission from [83], T. Morie et al., *IEEE Trans. Nanotechnology* 2, 158 (2003). © 2003, IEEE.

from the MOSFET. In the former case, this means that the interconnection wires for nput can be much wider than the interval between nanodot array sets. This condition is preferable for VLSI fabrication technology. Furthermore, by changing the circuit parameters such as $V_c$ and $V_p$ locally, the profile of PSP generated by each input pulse can be controlled.

PSPs generated in plural nanodot arrays are summed on the gate capacitance of the MOSFET. If one adopts models that ignore the difference in the transmission delay caused by the spatial distribution of synapses on dendritic trees, excitatory and inhibitory PSPs are separately integrated, each by a different MOSFET having a nanostructure on the gate and then the inhibitory contribution is subtracted from the excitatory one. A differentia-pair circuit can be used effectively for this subtraction and for the thresholding operation, as shown in Fig. 107.

The averaging process can also reduce various nonidealities in nanostructures, such as additional charge effects resulting from background offset, parasitic or surplus charges and the effect of device parameter fluctuation. In single-PSP operation, the addition of chrges comparable to the elementary charge causes a fatal error. However, even if some dot arrays do not work correctly by additional charges, the averaging process by many dot array: can decrease the effect. Thus, the proposed device and circuit will operate successfully even if the fabrication technology is not yet fully mature.

### 5.9.4. Simulation Results

The proposed circuit shown in Fig. 105 was analyzed by using a Monte Carlo single-electron simulator. In the simulations, the tunnel junction capacitance $C_j$ was assumed to be 01 aF, which means that the dot diameter is around 1 nm. The tunnel resistance $R_j$ was assumed to be 5 M$\Omega$. The operation temperature was set at 300 K.

Figure 108a shows the transition probability of $e_M$ between $D_c$ and $D_1$ as a function of time according to an input voltage change. Those data were obtained in the same wiy as for obtaining Fig. 100 in Section 5.8.6. By replotting Fig. 108a, one can obtain the time



Figure 109. Time dependence of MOS gate voltage (V): (a) one-input case. (b) two-input case. (A): electron $e_M$ comes back to center-dot (D). (X): first electron comes back to $D_c$. (Y): second electron comes back o $D_c$; (c) four-input case. a fitted α-function curve is indicated by the solid line. Reprinted with permission froa [83], T. Morie et al., IEICE Trans. Nanotechnology 2, 158 (2003). © 2003. IEEE.

dependence of the probability that $e_M$ stays at $D_1$ when a pulse is applied, as shown in Fig. 108b. The time course of PSP $(V_o)$ is approximately proportional to this probability. The time constant of PSP can be controlled by the baseline voltage of the input $(V_H)$, as shown in Fig. 108.

Figure 109a–c shows the time dependence of $V_o$ when a pulse is applied in the case of one, two and four nanodot arrays, respectively. Voltage $V_o$ fluctuates because of the position of the electron, but the average voltage roughly represents a PSP (i.e., the exponential decay characteristic of a PSP is realized in $V_o$). In Fig. 109c, a fitted $\alpha$-function curve is indicated by a solid line. Although an $\alpha$-function expresses only a typical PSP time course, the fit is reasonably good. It is verified from Fig. 109b and 109c that the summation of plural inputs can also be realized.

Thus, the function of the synapse part, which is the generation of arbitrary PSPs, can be realized by using a MOSFET with nanostructures. The decay constant of PSPs is controlled by the input baseline voltage $V_H$ and the bias voltage $V_c$, which affect the potential barrier height. The functions of the soma part in the neuron, such as thresholding and refractoriness, are realized by MOS circuits, including the base MOSFET of nanostructures.

# 6. CONCLUSION

This chapter reviewed single-electron devices and circuits, described various new operation principles and algorithms utilizing single-electron operation, and illustrated various functional circuits implementing such algorithms. Most of these circuits are not based on conventional CMOS Boolean computing, but they can coexist with conventional digital systems. Only sophisticated nanotechnology may construct these functional circuits, and at the same time, they should be promising applications of post-CMOS nanoelectronics.

# REFERENCES

1. H. Grabert and M. H. Devoret, (Eds.), "Single Charge Tunneling," Plenum Press, New York, 1992.
2. K. K. Likharev, Proc. IEEE 87, 606 (1999).
3. Y. Takahashi, Y. Ono, A. Fujiwara, and H. Inokawa, J. Phys.: Condens. Matter 14, R995 (2002).
4. D. Esteve, in "Single Charge Tunneling" (H. Grabert and M. H. Devoret, Eds.), p. 109, Plenum Press, New York, 1992.
5. K. K. Likharev, IEEE Trans. Magn. 23, 1142 (1987).
6. J. R. Tucker, J. Appl. Phys. 72, 4399 (1992).
7. L. S. Kuzmin, P. Delsing, T. Claeson, and K. K. Likharev, Phys. Rev. Lett. 62, 2539 (1989).
8. L. J. Geerligs, V. F. Anderegg, P. A. M. Holweg, J. E. Mooij, H. Pothier, D. Esteve, C. Urbina, and M. H. Devoret, Phys. Rev. Lett. 64, 2691 (1990).
9. H. Pothier, P. Lafarge, R. F. Orfila, C. Urbina, D. Esteve, and M. H. Devoret, Physica B 169, 573 (1991).
10. S. Shimano, K. Masu, and K. Tsubouchi, Jpn. J. Appl. Phys. 38, 403 (1999).
11. R. H. Chen, A. N. Korotkov, and K. K. Likharev, Appl. Phys. Lett. 68, 1954 (1996).
12. A. N. Korotkov, R. H. Chen, and K. K. Likharev, J. Appl. Phys. 78, 2520 (1995).
13. N. Zimmerman, W. H. Huber, A. Fujiwara, and Y. Takahashi, Appl. Phys. Lett. 79, 3188 (2001).
14. A. Fujiwara, S. Horiguchi, M. Nagase, and Y. Takahashi, Jpn. J. Appl. Phys. 42, 2429 (2003).
15. T. M. Eiles, G. Zimmerli, H. D. Jensen, and J. M. Martinis, Phys. Rev. Lett. 69, 148 (1992).
16. L. J. Geerligs, D. V. Averin, and J. E. Mooij, Phys. Rev. Lett. 65, 3037 (1990).
17. D. V. Averin and A. A. Odintsov, Phys. Lett. A 140, 251 (1989).
18. D. V. Averin and Y. V. Nazarov, in "Single Charge Tunneling" (H. Grabert and M. H. Devoret, Eds.), p. 217, Plenum Press, New York, 1992.
19. A. Kohno, H. Murakami, M. Ikeda, S. Miyazaki, and M. Hirose, Jpn. J. Appl. Phys. 40, L721 (2001).
20. T. Morie, T. Matsuura, M. Nagata, and A. Iwata, J. Nanosci. Nanotech. 2, 343 (2002).
21. N. Kuwamura, K. Taniguchi, and C. Hamaguchi, IEICE Trans. Electron. J77-C-II, 221 (1994).
22. "QCA home page, http://www.nd.edu, qcahome," (2004).
23. P. D. Tougaw, C. S. Lent, and W. Porod, J. Appl. Phys. 74, 3558 (1993).
24. C. S. Lent and P. D. Tougaw, J. Appl. Phys. 74, 6227 (1993).
25. P. D. Tougaw and C. S. Lent, J. Appl. Phys. 75, 1818 (1994).
26. C. S. Lent and P. D. Tougaw, J. Appl. Phys. 75, 4077 (1994).
27. P. D. Tougaw and C. S. Lent, Jpn. J. Appl. Phys. 34, 4373 (1995).
28. P. D. Tougaw and C. S. Lent, J. Appl. Phys. 80, 4722 (1996).
29. W. Porod, J. Franklin Institute-Eng. Appl. Math. 334B, 1147 (1997).
30. A. O. Orlov, I. Amlani, G. H. Bernstein, C. S. Lent, and G. L. Snider, Science 277, 928 (1997).

31. C. S. Lent and B. Isaksen, *IEEE Trans. Electron Devices* 50, 1890 (2003).
32. S. B. Akers, *IEEE Trans. Comput.* C-27, 509 (1978).
33. R. E. Bryant, *IEEE Trans. Comput.* C-35, 677 (1986).
34. N. Asahi, M. Akazawa, and Y. Amemiya, *IEEE Trans. Electron Devices* 42, 1999 (1995).
35. N. Asahi, M. Akazawa, and Y. Amemiya, *IEEE Trans. Electron Devices* 44, 1109 (1997).
36. N. Asahi, M. Akazawa, and Y. Amemiya, *IEICE Trans. Electron.* E81-C, 49 (1998).
37. T. Yamada, Y. Kinoshita, S. Kasai, H. Hasegawa, and Y. Amemiya, *Jpn. J. Appl. Phys.* 40, 4485 (2001).
38. S. Kasai, M. Yumoto, T. Tamura, and H. Hasegawa, in "Conference Digest of 61st Device Research Conference (DRC)," 2003, p. 97.
39. S. Kasai, M. Yumoto, T. Tamura, and H. Hasegawa, in "Proceedings of the 2003 Asia-Pacific Workshop on Fundamentals and Application of Advanced Semiconductor Devices," 2003, p. 177.
40. S. Amarel, G. Cooke, and R. O. Winder, *IEEE Transactions Electronic Computers* 13, 4 (1964).
41. A. R. Meo, *IEEE Transactions Electronic Computers* 15, 606 (1966).
42. H. Iwamura, M. Akazawa, and Y. Amemiya, *IEICE Trans. Electron.* E81-C, 42 (1998).
43. T. Oya, T. Asai, T. Fukui, and Y. Amemiya, *J. Nanosci. Nanotech.* 2, 333 (2002).
44. T. Oya, T. Asai, T. Fukui, and Y. Amemiya, *IEEE Trans. Nanotechnology* 2, 15 (2003).
45. G. E. Hinton and T. J. Sejnowski, in "Parallel Distributed Processing," MIT Press, Cambridge MA, 1986, Chap. 7.
46. E. Aarts and J. Korst, "Simulated Annealing and Boltzmann Machines: A Stochastic Approach to Combinatorial Optimization and Neural Computing," John Wiley and Sons, 1989.
47. J. Alspector, J. W. Gannet, S. Harber, M. B. Parker, and R. Chu, in "IEEE Proc. of Int. Symp. Circuits and Systems," 1990, p. 1058.
48. M. Akazawa and Y. Amemiya, *Appl. Phys. Lett.* 70, 670 (1997).
49. T. Yamada, M. Akazawa, T. Asai, and Y. Amemiya, *Nanotechnology* 12, 60 (2001).
50. C. Y. Isenberg, "The Science of Soap Films and Soap Bubbles," Tieto Ltd., 1978.
51. A. K. Dewdny, *Scientific American* 250, 15 (1984).
52. A. K. Dewdny, *Scientific American* 252, 12 (1985).
53. E. Tokuda, N. Asahi, T. Yamada, and Y. Amemiya, *Analog Integrated Circuit and Signal Processing* 24, 41 (2000).
54. D. Deutsch, *Proc. Royal. Soc. Lond. A* 400, 97 (1985).
55. J. J. Hopfield, *Proc. Natl. Acad. Sci. USA* 79, 2554 (1982).
56. J. J. Hopfield and D. W. Tank, *Biol. Cybern.* 52, 141 (1985).
57. M. Akazawa, E. Tokuda, N. Asahi, and Y. Amemiya, *Analog Integrated Circuit and Signal Processing* 24, 51 (2000).
58. H. D. Jensen and J. M. Martinis, *Phys. Rev. B* 46, 13407 (1992).
59. K. Nakano, *IEEE Trans. Syst. Man. Cybern.* SMC-2, 380 (1972).
60. J. J. Hopfield, *Proc. Natl. Acad. Sci. USA* 81, 3088 (1984).
61. T. Yamanaka, T. Morie, M. Nagata, and A. Iwata, *IEICE Trans. Electron.* E84-C, 1723 (2001).
62. T. Morie, T. Matsuura, M. Nagata, and A. Iwata, in "Advances in Neural Information Processing Systems" (T. G. Dietterich, S. Becker, and Z. Ghahramani, Eds.), Vol. 14, p. 1115, MIT Press, Cambridge, MA, 2002.
63. T. Yamanaka, T. Morie, M. Nagata, and A. Iwata, *Nanotechnology* 11, 154 (2000).
64. M. Saen, T. Morie, M. Nagata, and A. Iwata, *IEICE Trans. Electron.* E81-C, 30 (1998).
65. A. R. Bulsura and L. Gammaitoni, *Physics Today* 49, 39 (1996).
66. S. Tiwari, F. Rana, H. Hanafi, A. Hartstein, E. F. Crabbé, and K. Chan, *Appl. Phys. Lett.* 68, 1377 (1996).
67. R. Ohba, N. Sugiyama, J. Koga, K. Uchida, and A. Toriumi, in "Ext. Abs. of Int. Conf. on Solid State Devices and Materials (SSDM)," 2000, p. 122.
68. S. Huang, G. Tsutsui, H. Sakaue, S. Shingubara, and T. Takahagi, *J. Vac. Sci. Technol. B* 18, 2653 (2000).
69. R. A. Kiehl, in "Extended Abstracts, 4th International Workshop on Quantum Functional Devices (QFD)," 2000, p. 49.
70. A. Iwata and M. Nagata, *IEICE Trans. Fundamentals.* E79-A, 145 (1996).
71. W. Maass and C. M. Bishop, (Eds.) "Pulsed Neural Networks," MIT Press, Cambridge, MA, 1999.
72. D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, "Parallel Distributed Processing," MIT Press, Cambridge, MA, 1986.
73. R. E. Howard, D. B. Schwartz, J. S. Denker, R. W. Epworth, H. P. Graf, W. E. Hubbard, L. D. Jackel, B. L. Straughn, and D. M. Tennant, *IEEE Trans. Electron Devices* ED-34, 1553 (1987).
74. T. Shima, T. Kimura, Y. Kamatani, T. Itakura, Y. Fujita, and T. Iida, *IEEE J. Solid-State Circuits* 27, 868 (1992).
75. T. Morie and Y. Amemiya, *IEEE J. Solid-State Circuits* 29, 1086 (1994).
76. M. Yasunaga, N. Masuda, M. Yagyu, M. Asai, K. Shibata, M. Ooyama, M. Yamada, T. Sakiguchi, and M. Hashimoto, *IEEE J. Solid-State Circuits* 28, 106 (1993).
77. Y. Kondo, Y. Koshiba, Y. Arima, M. Murasaki, T. Yamada, H. Amishiro, H. Shinohara, and H. Mori, in "IEEE Int. Solid-State Circuits Conf. Dig." 1994, p. 218.
78. O. Saito, K. Aihara, O. Fujita, and K. Uchimura, in "IEEE Int. Solid-State Circuits Conf. Dig." 1998, p. 94.
79. A. F. Murray and L. Tarassenko, "Analogue Neural VLSI—A Pulse Stream Approach," Chapman & Hall, London, UK, 1994.
80. Y. Hirai and M. Yasunaga, in "Proc. Int. Conf. on Neural Information Processing (ICONIP)," 1996, p. 1251.
81. W. Maass, *Neural Networks* 10, 1659 (1997).
82. W. Maass, in "Advances in Neural Information Processing Systems" (M. C. Mozer, M. I. Jordan, and T. Petsche, Eds.), Vol. 9, p. 211, The MIT Press, 1997.
83. T. Morie, T. Matsuura, M. Nagata, and A. Iwata, *IEEE Trans. Nanotechnology* 2, 158 (2003).

# CHAPTER 5

# Modeling of Single-Electron Transistors for Efficient Circuit Simulation and Design

## YunSeop Yu

*Department of Information and Control Engineering,*
*Hankyong National University, Anseong, Kyeonggi-do, Republic of Korea*

## SungWoo Hwang

*Department of Computer and Electronics Engineering, Korea University, Sungbuk,*
*Seoul, Republic of Korea, and Institute of Quantum Information*
*Processing and Systems, University of Seoul, Jeonnong, Dongdaemun,*
*Seoul, Republic of Korea*

## Doyeol Ahn

*Institute of Quantum Information Processing and Systems, University of Seoul,*
*Jeonnong, Dongdaemun, Seoul, Republic of Korea*

## CONTENTS

# 1. INTRODUCTION

Modern microelectronics strongly tends toward scaling down of electronic device size for the development of ultra-large-scale integrated circuits (ULSIs) [1–4]. Meal-oxide-semiconductor field-effect transistors (MOSFETs) have been the most prevalent electron devices for ULSI applications. Because of higher integration of MOSFETs, the/ lead to an increase in performance speed due to smaller delay times in information exciange. In the early years of the 21st century, the scaling of complementary-MOSFETs (CMOSFETs) entered the deep sub-50-nm regime [5]. in which fundamental limits of CMOSFETs and technological challenges with regard to the scaling of CMOSFETs are encountered 4, 6]. On the other hand, it has opened up new possibilities based on quantum-mechanical efects [7]. Therefore, it is essential to introduce a new device having an operation principle that is effective in small dimensions, such as down-scaling existing devices, molecular devices. mesoscopic



**Figure 1.** Flow chart for the SPICE.

Figure 2. Single-electron circuits consisting of tunnel junctions, capacitors, and voltage sources.

devices with quantum effects, superconducting devices with tunneling effects, single-electron devices (SEDs) with single-electron effects, and so on, and thus provide a new functionality beyond that attainable with CMOSFETs [8–14]. Among the new devices, single-electron devices [8, 15–17] are an attractive candidate for their potential room temperature application to very high density memory and logic circuits with conventional silicon very-large-scale integrated circuit (VLSI) processing [18–28] techniques because they can retain their scalability even on an atom scale and, moreover, they can control the motion of even a single electron. Therfore, the ULSI consisting of SEDs will have the attributes of extremely high integration and extremely low power consumption.

However, a number of issues have to be overcome that at present pose crucial obstacles for implementing single-electron transistor (SET) logic gates. The maximum voltage gain of an SET, which is defined as the ratio of the gate to the drain capacitance, is very small, usually less than one or slightly more than one [29, 30]. The SETs have low current driving capability, which degrades device performance, as it takes a long time to charge up the large interconnection capacitance connected to the output node of a device. One of the approaches to overcome these inherent disadvantages of SETs is to construct a hybrid circuit consisting of MOSFETs, which have a high gain, high output resistance, and high applicable voltages and can thus supplement the SET. Such a hybrid circuit has already been demonstrated experimentally [31–35], and in order to evaluate the merits of this approach thoroughly, powerful simulation tools are necessary.

Most of the single-electron circuit simulators [36–39] have procedures that calculate the charge states of all the Coulomb islands altogether to take into account the interaction between neighboring Coulomb islands. The pioneering SET circuit simulator, the MOSES [36], the SIMON [37], and the KOSEC [38], use the Monte Carlo method to obtain the average number of electrons in Coulomb islands. Other important simulators, the SENECA [39], directly solves the master equation for the population probability of Coulomb islands. Although both methods provide accurate simulation results, they have to deal with a huge amount of matrix calculation and require a huge amount of computation time because the evolution of all the Coulomb islands in the entire circuit has to be monitored [40, 41].

To the best of our knowledge, it is very difficult to include the simulation of complementary metal-oxide-semiconductor (CMOS) devices and to predict the performance of SET-CMOS hybrid circuits with the above-mentioned pioneering simulators. Therefore, various types of SET-CMOS hybrid circuit simulators for efficient circuit simulation have been developed [41–51].

In this chapter, the modeling of single-electron circuits and SETs for efficient circuit simulation and proposed circuits for logic applications are explained. In Section 2, the simulation method of the conventional simulation program such as SPICE and the basic assumptions used with SPICE are illustrated. In Section 3, we explain the simulation method of single-electron circuits consisting of SEDs with the procedure that calculates the charge states of all the Coulomb islands together to take into account the interaction between neighboring

**Figure 3.** Simplified flow chart of the Monte Carlo method for the simulation of single-electron circuits. $t$ is simulation time. $t_{sim}$ is the simulation end time, and $t_{step}$ is the simulation time step.

Coulomb islands. In Section 4, when an SET circuit is applied to SPICE, the features of SET circuits are illustrated. In Section 5, the minimum size (capacitance) condition of the interconnection, where each SET can be handled independently in the circuit simulation at the DC and the transient, is investigated using the single-electron inverter (SEI) circuit as an example. For those regimes where each SET can be handled independently, the various types of compact-modelings of the SET to efficiently simulate SET circuits are introduced in Section 6. In Section 7, comparisons between each SET modeling are performed. In Section 8, various types of SET circuits and SET-MOSFET hybrid circuits are introduced.

## 2. SIMULATION METHOD OF CONVENTIONAL CIRCUITS

In this section, simulation methods of the conventional circuit simulation program such as SPICE are illustrated, and the basic assumptions used with SPICE are introduced.

**Figure 4.** Partition of state space into a frequent state domain and a rare state domain, and contribution of rare states and events by stepping through the event tree. Black solid arrows mark transitions due to frequent events, black dashed arrows mark transitions due to rare events, and blue dashed arrows mark transitions due to rare or frequent events. Reprinted with permission from [37], C. Wasshuber et al., *IEEE Trans. Comp. Aided Design* 16, 937 (1997). © 1997, IEEE.

## 2.1. Basic Equations for Conventional Circuit Simulation

A conventional circuit simulator such as SPICE is based on the Kirchhoff's current law (KCL) equation at each node. Figure 1 summarizes the flow chart for the SPICE simulation. When there are $N$ nodes in the circuit, the $N \times 1$ nodal voltage vector $\mathbf{V}$ (collection of the voltages of all the nodes in the circuit) can be obtained by solving the following matrix equation.

$$\mathbf{YV} = \mathbf{J} \tag{1}$$



**Figure 5.** Circuit diagram for a single-electron inverter consisting of two SET's in series.

**Figure 6.** The voltage transfer characteristics and the current of the lower SET of the inverter when $I_g = 0.03$ V, $C_d = C_s = 1.6$ aF, $C_g = 3.2$ aF, $C_l = 1.6$ aF, $R_d = R_s = 100$ MΩ, and $T = 30$ K. Reprinted with permission from [42], Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

where $\mathbf{Y}$ is the $N \times N$ nodal admittance matrix and $\mathbf{J}$ is the $N \times 1$ nodal current vector. As shown in the flow chart, $\mathbf{V}$ is calculated by the initial guess and successive iteration and Newton-Raphson convergence [52, 53]. During the successive iteration, the elements of $\mathbf{Y}_n$, $\mathbf{V}_n$, and $\mathbf{J}_n$ (the subscript $n$ means $n$th iteration) are related by the following formula:

$$[\mathbf{Y}_n]_{ij} = \frac{\partial I_i^n}{\partial [\mathbf{V}_n]_j} \tag{2}$$

$$[\mathbf{J}_n]_i = \sum_k^{all\ nodes} \left\{ \frac{\partial I_i^n}{\partial [\mathbf{V}_n]_k} [\mathbf{V}_n]_k \right\} - I_i^n \tag{3}$$

where $I_i^n$ is the sum of all the currents flowing out of the $i$th node (the superscript $n$ means $n$th iteration). The SPICE itself automatically provides equivalent circuits so that $I_i^n$ can be calculated at each node.

## 2.2. Basic Assumptions

In the case of the conventional circuit, the compact simulators such as SPICE [54] are used to simulate the characteristics of the given circuit topology. In these simulators, two assumptions are implicitly used to build the model. The first assumption is that once the parameter



**Figure 7.** The voltage transfer characteristics and the current of the lower SET of the inverter when $I_g = 0.03$ V, $C_d = C_s = 1.6$ aF, $C_g = 3.2$ aF, $C_l = 500$ aF, $R_d = R_s = 100$ MΩ, and $T = 30$ K. Reprinted with permission from [42], Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

**Figure 8.** (a) The output voltage swing, $V_{out}^{swing}$ and (b) the average output level, $V_{out}^{avg}$ as a function of $C_l$ ($V_B =$ 0.03 V, $C_d = C_s = 1.6$ aF, $C_g = 3.2$ aF, $R_d = R_s = 100$ MΩ, and $T = 30$ K). Reprinted with permission from [42], Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

of the isolated transistor is determined from the device simulator or other modeling tools, they can be used in the whole circuit. The current–voltage characteristic of the device is approximated by a linear function of the device scale. The second assumption is that the $I-V$ characteristics of the device are affected by neighboring transistors only through the charges of the terminal voltages of those transistors. The interaction between each adjacent device is usually overlooked.

## 3. SIMULATION METHODS OF SINGLE-ELECTRON CIRCUITS

To simulate a single-electron circuit (with $n_a$ nodes) consisting of voltage sources and capacitors including tunnel junctions as shown in Fig. 2, the node voltages and the charges stored in the capacitors are calculated and are related by

$$CV = Q \tag{4}$$

where $C$ is the $n_a \times n_a$ capacitance matrix of the circuit network, $V$ is the $n_a \times 1$ vector of node voltages, and $Q$ is the $n_a \times 1$ vector of node charges. Here, $n_a = n_i + n_b$, where $n_i$ and $n_b$ are the number of islands in the circuit and nodes connected to the voltage sources or the



**Figure 9.** Average difference between the output voltages ($V_{out}$) obtained from two types of calculations when the input voltage ($V_{in}$) of the inverter is swept from 0 to 0.03 V in three different $t_r$s ($C_d = C_s = 1.6$ aF, $C_g = 3.2$ aF, $R_d = R_s = 100$ MΩ, $V_B = 0.03$ V, and $T = 30$ K). The result $V_{ME}$ is the output voltage calculated by solving the time-dependent master equation considering the overall probability distribution of three Coulomb islands. The result $V_{K+I}$ is the output voltages calculated by solving two time-dependent master equations of the lower and the upper SETs independently and applying the Kirchhoff's law at the interconnection.

(a)

(b)

```
.macro SET 1 2 3 Cg=1.0e-18 Cd=1.0e-18 Cs=1.0e-18
.param
+pi=3.1415926535897932846
+CF1=40          ; 2C_G/e
+CI2=0.2e-9
+CR1=300e+6
+CR2=100e+6
+CVp=0.02
+CP1=0
+Csum='Cd+Cs+Cg'

V2 5 3 DC CVp
V3 7 3 DC '-CVp'
RG 2 3 100G
RR1 1 3 R= 'CR1+CR2*cos(CF1*pi*V(2,3)+CP1)'
RR2 1 4 R= 'CVp/(CI2-2*CVp/(CR1+CR2*cos(CF1*pi*V(2,3)+CP1)))'
RR3 1 6 R= 'CVp/(CI2-2*CVp/(CR1+CR2*cos(CF1*pi*V(2,3)+CP1)))'
CGS 2 3 'Cg*Cs/Csum'
CGD 2 1 'Cg*Cd/Csum'
CDS 1 3 'Cd*Cs/Csum'
D1 4 5 DIODE
D2 7 6 DIODE
.MODEL DIODE D(N=0.01)
.eom
```

**Figure 10.** (a) The equivalent circuit and (b) the SPICE macro model code of the SET for the macromodeling. Reprinted with permission from [42]. Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.



**Figure 11.** The current-voltage characteristics of an SET in Fig. 7 at various gate biases when $C_g = C = 1.6$ aF, $C = 3.2$ aF, $R_1 = R = 100$ MΩ, and $T = 30$ K. Reprinted with permission from [42]. Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

**Figure 12.** The current–voltage characteristics of the SET in Fig. 7 at various temperatures when $C_g = C_s = 1.6$ aF. $C_j = 3.2$ aF. $R_g = R_s = 100$ MΩ, and $T = 30$ K. The open symbols are Monte Carlo results, and the filled symbols are obtained from the proposed macromodel. Reprinted with permission from [42], Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

ground, respectively. Ordering all the biased nodes before the island nodes, the following equation in terms of submatrices and vectors is rewritten [43].

$$\begin{pmatrix} \mathbf{Q_I} \\ \mathbf{Q_B} \end{pmatrix} = \begin{pmatrix} \mathbf{C_{II}} & \mathbf{C_{IB}} \\ \mathbf{C_{BI}} & \mathbf{C_{BB}} \end{pmatrix} \begin{pmatrix} \mathbf{V_I} \\ \mathbf{V_B} \end{pmatrix} \tag{5}$$

where $\mathbf{Q_I}$ is the $n_i \times 1$ vector of island charges, $\mathbf{Q_B}$ is the $n_b \times 1$ vector of biased node charges. $\mathbf{V_I}$ is the $n_i \times 1$ vector of island voltages, $\mathbf{V_B}$ is the $n_b \times 1$ vector of biased node voltages. $\mathbf{C_{II}}$ is the $n_i \times n_i$ matrix of island–island capacitances, $\mathbf{C_{IB}}$ is the $n_i \times n_b$ matrix of island-biased capacitances, and so on (note that $\mathbf{C_{BI}} = \mathbf{C_{IB}^t}$). Then, because $\mathbf{V_B}$ and $\mathbf{Q_I}$ are known, $\mathbf{Q_B}$ and $\mathbf{V_I}$ have to be calculated by [43].

$$\begin{pmatrix} \mathbf{Q_B} \\ \mathbf{V_I} \end{pmatrix} = \begin{pmatrix} \mathbf{C_{BB}} - \mathbf{C_{BI}} \cdot \mathbf{C_{II}^{-1}} \cdot \mathbf{C_{IB}} & \mathbf{C_{BI}} \cdot \mathbf{C_{II}^{-1}} \\ -\mathbf{C_{II}^{-1}} \cdot \mathbf{C_{IB}} & \mathbf{C_{II}^{-1}} \end{pmatrix} \begin{pmatrix} \mathbf{V_B} \\ \mathbf{Q_I} \end{pmatrix} \tag{6}$$

To simulate the tunneling of electrons from island to island, the rates of all possible tunnel events have to be determined. The normal tunnel rate of a tunnel junction is given by [15–17]

$$\Gamma = \frac{\Delta F}{e^2 R_T \left[ 1 - \exp(\frac{\Delta F}{k_B T}) \right]} \tag{7}$$

where $\Delta F$ is the change in free energy after and before a tunneling event, $e$ is the electron charge, $R_T$ is the tunnel resistance, and $k_B T$ is the thermal energy. The free energy is defined

**Table 1.** The macromodel parameters at various temperatures. The parameters CFI (= 40) and CPI (= 0) are temperature independent.

| Temp. Para. | 10 K | 30 K | 77 K | 100 K | 300 K |
|---|---|---|---|---|---|
| CI2 | 0.2 n | 0.2 n | 0.25 n | 0.27 n | 0.35 n |
| CR1 | 1.35 G | 300 M | 168 M | 147.5 M | 118 M |
| CR2 | 1.15 G | 100 M | 14 M | 4.5 M | 0 |

Reprinted with permission from [42], Y. S. Yu et al., *IEEE Trans. Electron. Devices* 46, 1667 (1999). © 1999, IEEE.

**Figure 13.** The circuit diagram and the linearized equivalent circuit of a four-terminal SET for the seninumerical modeling. The symbols, $C_d$, $C_g$, $C_b$, $C_{gd}$, $C_{gb}$, and $C_{bd}$ are the drain-source, the gate-source, the backgate-source, the gate-drain, the gate-backgate, and the backgate-drain capacitance, respectively. Reprinted with permission from [49], Y. S. Yu et al., *Electron. Lett.* 38, 850 (2002). © 2002, IEE.

by the difference in electrostatic energy stored in the circuit ($U$) and the work done by the voltage sources ($W$) and is given by [37]

$$F = U - W = \frac{1}{2} (Q_B^T, \; V_I^T) \begin{pmatrix} V_B \\ Q_I \end{pmatrix} - W \tag{8}$$

with

$$W = \sum_m \int V_m(t) i_m(t) dt \tag{9}$$

where $V_m(t)$ and $i_m(t)$ are the voltage of the $m$th voltage sources and the current through $m$th voltage source, respectively.

Single-electron tunneling is a stochastic phenomenon, and its theory can only predict probability rates of the possible tunneling events [15–17]. Therefore, two simulation approaches are used for single-electron devices and circuits. One is based on a Monte Carlo method [55], and the other is based on a master equation [15]. Monte Carlo method simulates possible scenarios of electron tunneling between the islands. This is a very convenient method for studying the typical behavior of the electrons in a device. However, this method is done many times to simulate the transport of electrons through the network and especially, in case when very rare tunneling events such as co-tunneling take place against the background of much more frequent events (single-electron tunneling), it is difficult to resolve the rare tunneling events by the Monte Carlo method. Master equation is a description for underlying Markov process [56] of electrons tunneling between the islands, and thus the circuit occupies different states. The set of all possible states of the circuit is needed. A state is defined by the set of voltages of voltage sources and charge distribution in the circuit. This method can include the tunneling phenomena through a single tunnel junction as well as the tunneling phenomena through multiple tunnel junctions (co-tunneling). Therefore, the probabilities of all possible states of the circuit (each characterized by a particular charge configuration) are calculated by the master equation method. This method has the necessity to store information about all states of the circuit. To rectify the requirement to a large extent by taking into account the hierarchy of the tunneling rates, only rates higher than

**Figure 14.** (a) The drain ($I_d$) and the source current ($I_s$), and (b) the gate ($I_g$) and the backgate current ($I_b$) of the SET in Fig. 10 when $V_{gs}$ varies linearly from 0 to 0.1 V in $t_r = 50$ ps, and (c) the drain current when $V_{gs}$ varies as shown in the inset ($V_{ds} = 0.1$ V, $V_{bs} = 0.075$ V, $C_d = 0.2$ aF, $C_s = 0.1$ aF, $C_g = 0.8$ aF, $C_b = 0.7$ aF, $R_d = R_s = 1$ MΩ, and $T = 30$ K). The results from our DC model and the Monte Carlo results are also shown. The symbols denote the results of our SPICE model and the lines denote the results of the master equation or the Monte Carlo.

a certain threshold value are considered at a time. Two simulation approaches will be the subject of detailed discussion in the next sections.

## 3.1. Monte Carlo Method

To use the Monte Carlo method for the simulation of single-electron circuits was first proposed and implemented by N. Bakhvalov et al. [57]. Other groups adopted this method later [58, 59]. Figure 3 shows the simplified flow chart of the Monte Carlo method for the simulation of single-electron circuits. Summarizing the Monte Carlo method, it starts with all possible tunnel events, calculates their probabilities, and chooses one of the possible events randomly, weighted according to their probabilities.

Simulation procedure is as follows. Given a tunneling rate $\Gamma$ for a tunneling event, the probability that a tunnel event out of state 0 happens at $t$ and not earlier is given by

$$P_0(t) = e^{-\Gamma t} \qquad (10)$$

Here, the tunneling events are considered as independent and exponentially distributed processes. Therefore, time duration to the tunneling events time is determined as [60, 61]

$$\Delta t = -\frac{1}{\Gamma} \ln(r) \qquad (11)$$

where $r$ is an evenly distributed random number from the interval $0 \le r < 1$. $\Delta t$ is calculated for each tunneling event starting from the present state. The tunneling event with the shortest $\Delta t$ of all (shorter than the simulation time step $t_{step}$) is chosen as the one that actually occurs.

**Figure 15.** Equivalent circuit of the SET in Fig. 13 for two-state approximation of the full-analytical modeling [50]. Reprinted with permission from [50], H. Inokawa et al., *IEEE Trans. Electron Devices* 50, 455 (2003). © 2003, IEEE.

Then, the state of the circuit (node voltages and node charges) is updated. Consequently, the free energy changes too, and one has to calculate all possible tunneling rates again. This procedure is followed repeatedly, which is performed many times to simulate the transport of electrons through the network.

```
SUBCKT ASET  1 2 3 4
  +Cs=2E-18            : Capacitance of junction 1 (source)
  +Cd=2E-18            : Capacitance of junction 2 (drain)
  +Rs=2E6              : Resistance of junction 1 (source)
  +Rd=2E6              : Resistance of junction 2 (drain)
  +Cg=8E-18            : Capacitance of gate
  +Cb=0                : Capacitance of back-gate
  +C0=0                : Self Capacitance of the island
  +Q0=0                : Offset charge in units of e
  +TSET=4.2            : Temperature

.PARAM E=1.60217733E-19      : Electronic charge
.PARAM KB=1.380658E-23       : Boltzmann's constant
.PARAM Csum='Cs+Cd+Cg+Cb+C0'   : The total capacitance of the SET
.PARAM NT='2*KB*TSET*CSUM/(E^2)'    : Normalized temperature 2T*kB*Csum/e^2
.PARAM Rsum='2*Rs*Rd/(Rs+Rd)'
.PARAM NR='(Rd-Rs)/(Rs+Rd)'

A1 5 1 V='2*Cg*V(3,1)/E + 2*Cb*V(4,1)/E - (Cg+Cb+Cs-Cd)*V(2,1)/E - 1'
A2 8 1 V= IF COS(PI*(V(5,1)+1)/2) > 0
  +      THEN ASIN(SIN(PI*(V(5,1)+1)/2))/PI
  +      ELSE -ASIN(SIN(PI*(V(5,1)+1)/2))/PI
A21 6 1 V='V(5,1) - 2*(V(5,1)+1)/2 - V(8,1)'
A3 7 1 V='Csum*V(2,1)/E'
A4 2 1 I='E*(1-NR^2)*(V(6,1)^2-V(7,1)^2)*SINH(V(7,1)/NT)/(4*Rsum*Csum/
  +      ((V(6,1)*SINH(V(6,1)/NT)-V(7,1)*SINH(V(7,1)/NT)
  +      + NR*(V(7,1)*SINH(V(6,1)/NT)-V(6,1)*SINH(V(7,1)/NT)))'
A5 2 1 I=E*(1-NR^2)*(V(6,1)+2)^2-V(7,1)^2)*SINH(V(7,1)/NT)/(4*Rsum*Csum/
  +      ((V(6,1)+2)*SINH((V(6,1)+2)/NT)-V(7,1)*SINH(V(7,1)/NT)
  +      + NR*(V(7,1)*SINH((V(6,1)+2)/NT)-(V(6,1)+2)*SINH(V(7,1)/NT)))'
A6 2 1 I='E*(1-NR^2)*(V(6,1)-2)^2-V(7,1)^2)*SINH(V(7,1)/NT)/(4*Rsum*Csum/
  +      ((V(6,1)-2)*SINH((V(6,1)-2)/NT)-V(7,1)*SINH(V(7,1)/NT))
  +      + NR*(V(7,1)*SINH((V(6,1)-2)/NT)-(V(6,1)-2)*SINH(V(7,1)/NT)))'

.ENDS ASET
```

**Figure 16.** Source code of SPICE subcircuits implemented to SmartSpice for two-state approximation of the full-analytical modeling [50, 73]. Reprinted with permission from [50], H. Inokawa et al., *IEEE Trans. Electron Devices* 50, 455 (2003). © 2003, IEEE.

**Figure 17.** $I_d$-$I_g$ characteristics of asymmetric SETs calculated by two-state approximation of the full-analytical modeling (line) and the reference simulator (symbols) for $R_d/R_s$ of 1 MΩ/19 MΩ (open diamonds) and 19 MΩ/1 MΩ (open circles). Other parameters are $V_{gs} = 26.7$ mV, $C_d = C_s = C_g = 1$ aF, $C_e = 0$. and $T = 18.6$ K. Reprinted with permission from [50], H. Inokawa et al., *IEEE Trans. Electron Devices* 50, 455 (2003). © 2003, IEEE.

## 3.2. Master Equation Method

Monte Carlo method becomes impractical for the study of rare tunnel events (co-tunneling) taking place against the background of much more frequent tunnel events operation (single-electron tunneling) of the system, because their events typically differ by several orders of magnitude, making the sampling time extremely long in the Monte Carlo method. In order to analyze rare events, the master equation method for the simulation of single-electron circuits was proposed and implemented by Pothier et al. [62], Jensen and Martinis [63], Fonseca [39, 64], and others [40, 43].

The single-electron tunneling based on the "orthodox" theory was derived from first-order perturbation theory [15–17]. However, in the Coulomb blockade regime, where the first-order tunnel rate is very low, or at zero temperature even zero, higher order processes may become important. The co-tunneling is the quantum tunneling phenomena through multiple tunnel junctions at a time. The co-tunneling effect is a quantum mechanical effect that allows electrons to tunnel via an intermediate virtual state, where normal tunneling would be impossible or very unlikely due to missing thermal energy. For finite temperature, the



**Figure 18.** Coulomb staircase ($I_d$-$I_{ds}$ characteristics) of an asymmetric SET calculated by two-state approximation of the full-analytical modeling (lines) and the reference simulator (symbols) for $V_{gs} = 0$ (open diamonds) and $V_{gs} = e/2C_g$ (= 80.1 mV) (open circles). Other parameters are $R_d = 1$ MΩ, $R_s = 19$ MΩ, $C_d = 0.1$ aF, $C_s = 1.9$ aF, $C_g = 1$ aF, $C_e = 0$. and $T = 18.6$ K. Reprinted with permission from [50]. H. Inokawa et al., *IEEE Trans. Electron Devices* 50, 455 (2003). © 2003, IEEE.

**Figure 19.** Equivalent circuit of the SET to the left of Fig. 13, based on the multistate (11 states) approximation of the full-analytical SET model [51]. Reprinted with permission from [51]. G. Lientschnig et al., *Jpn. J Appl. Phys.* 42, 6467 (2003). © 2003, Institute of Pure and Applied Physics.

rate of an $N$th-order "inelastic" co-tunneling process is given by the following expression [65, 66]

$$\Gamma^{(N)} = \frac{2\pi}{\hbar}\left(\prod_{i=1}^{N}\frac{\hbar}{2\pi e^2 R_{Ii}}\right)\int_{0}^{\infty} S^2(\omega_1,\ldots,\omega_{2N})\delta\left(\Delta E_N + \sum_{i=1}^{2N}\omega_i\right)\prod_{i=1}^{2N}[1 - f(\omega_i)]d\upsilon_i \quad (12)$$

where $\Delta E_N = E_N - E_0$ is the change in the electrostatic energy during co-tunnding from the initial state 0 to the final state $N$, $f$ is Fermi function, $\delta$ is the Dirac delta function, and the factor of tunneling matrix element $S$ is given by

$$S(\omega_1,\ldots,\omega_{2N}) = \sum_{perm\{k_1,\ldots,k_N\}}\prod_{k=1}^{N-1}\frac{1}{\varepsilon_k} \quad (13)$$

where each permutation of the tunneling configuration perm $\{k_1,\ldots,k_N\}$ gives rise to a tunneling sequence with the intermediate energies $\varepsilon_k$ given by [65]

$$\varepsilon_k = \Delta E_k + \sum_{i=1}^{2k}\omega_i \quad (14)$$

where $\Delta E_k = E_k - E_0$ is the change in the electrostatic energy during co-tunnding from initial state 0 to intermediate state $k$, and $\omega_{2k-1} + \omega_{2k}$ is the energy of one of the $k$ electron–hole excitations created between states 0 and $k$. A different mechanism, so called "elastic" co-tunneling [67], is also possible, but its rate is negligibly low in metallic structur [65].

Denoting a state with a particular charge configuration as **P**, the corresponding master equation is given by the following either scalar or matrix form [15].

$$\frac{dP_i}{dt} = \sum_i \Gamma_{ji}P_j - \sum_i \Gamma_{ij}P_i \quad \text{or} \quad \frac{d\mathbf{P}}{dt} = \Gamma\mathbf{P} \quad (15)$$

where $P_i$ is the occupation probability of state $i$. $(\Gamma)_{ij} = \Gamma_{ji}$ is the sum of transition rate from $j$ state to $i$ state, and $(\Gamma)_{ii} = -\sum_{j\neq i}\Gamma_{ij}$. The stationary case of Eq. (15). $\Gamma\mathbf{P} = 0$, is a system of linear equations that may be solved by a multitude of numerical algorithns [68]. In the transient case, which is a system of ordinary linear first-order homogeneous differential equations, if the simulation time step $\Delta t$ is small enough on the scale on which the external

voltages vary, the rate matrix to be constant in each of these intervals may be assumed so that the solution of the master equation can be immediately given by [65]

$$\mathbf{P}(t + \Delta t) = \exp(\Gamma \Delta t)\mathbf{P}(t) \tag{16}$$

where the exponential of the rate matrix is defined by the corresponding series expansion for a scalar argument. The Padé approximation [68] is used to calculate the matrix $\exp(\Gamma \Delta t)$ [37, 39, 63, 64]. This approximation gives an absolute error for $\exp(A)$ less than $10^{-10}$ for $\|A\| < 1$. For $\Delta t$, a value less than the inverse of the maximum rate occurring in the rate matrix is used.

In principle, the number of charge states in a single-electron circuit is infinite because it takes into account all possible charge states resulting from both classical and co-tunneling transitions in some range of the involved parameters. Handling of all possible states and transitions between them is impossible, and a large number of charge states have very low possibilities $P$ (below a certain threshold $P_{th}$). To solve this problem, charge states having probabilities below $P_{th}$ are ignored [39].

The current through an external electrode $n$ consists of the sum of two parts where the first term is the current due to tunneling through the junctions connected directly to the

```
.SUBCKT SET  1 2 3 4  PARAMS:  :source drain gate back-gate
  +Cs=1E-18              : Capacitance of junction 1
  +Cd=1E-18              : Capacitance of junction 2
  +Rs=1E5               : Resistance of junction 1
  +Rd=1E5               . Resistance of junction 2
  +Cg=1E-18              : Capacitance of gate 1
  +Cb=0                 : Capacitance of gate 2
  +C0=0                 : Self Capacitance of the island
  +Q0=0                 : Offset charge in units of e
  +TEMP=4.2             : Temperature

.PARAM PI=3.1415926535897932        : Pi constant
.PARAM E=1.60217733E-19             : Electronic charge
.PARAM KB=1.380658E-23              : Boltzmann's constant
.PARAM CSUM={Cs+Cd+Cg+Cb+C0}        : The total capacitance of the SET
.PARAM T={TEMP*CSUM*5.3785467E14}   : Normalized temperature 5.3785467E14 = kB/e^2
.PARAM RN1={Rs/(Rs+Rd)}             : Normalized resistance of junction 1
.PARAM RN2={Rd/(Rs+Rd)}             : Normalized resistance of junction 2

.FUNC Q(v1,v2,v3,v4) { (Cg*v3+Cb*v4+Cs*v1+Cd*v2)/E+Q0} : Definition of a charge term in units of e
.FUNC VN(v) { CSUM*v/E }                          : The normalized voltage
.FUNC GAMMA(u) { IF(T==0,IF(u<0,-u,0),IF(u==0,T,u/(EXP(u/T)-1))) } : The rate function
.FUNC ROUND(x) { x-IF(cos(PI*x)>0,arcsin(sin(PI*x))/PI,-arcsin(sin(PI*x))/PI) } : The round() function
.FUNC N_OPT(v1,v2,v3,v4) { ROUND(-Q(v1,v2,v3,v4)+(CSUM/E)*(v1*RN2+v2*RN1)) } : The most probable charge on the island in
units of e
*********** the rates for the four tunnel events **********
.FUNC R1L(n,v1,v2,v3,v4) {GAMMA(0.5 - n - Q(v1,v2,v3,v4) + VN(v1))/RN1}
.FUNC R1R(n,v1,v2,v3,v4) {GAMMA(0.5 + n + Q(v1,v2,v3,v4) - VN(v1))/RN1}
.FUNC R2L(n,v1,v2,v3,v4) {GAMMA(0.5 + n + Q(v1,v2,v3,v4) - VN(v2))/RN2}
.FUNC R2R(n,v1,v2,v3,v4) {GAMMA(0.5 - n - Q(v1,v2,v3,v4) + VN(v2))/RN2}
* determine the relative probabilities, charge state N_OPT is initially assumed to have a relative probability equal to one
.FUNC PN_1(n,v1,v2,v3,v4) {(R1L(n,v1,v2,v3,v4)+R2R(n,v1,v2,v3,v4))/(R1R(n-1,v1,v2,v3,v4)+R2L(n-1,v1,v2,v3,v4))}
.FUNC PN_2(n,v1,v2,v3,v4) { PN_1(n,v1,v2,v3,v4)*
         +R1L(n-1,v1,v2,v3,v4)+R2R(n-1,v1,v2,v3,v4))/(R1R(n-2,v1,v2,v3,v4)+R2L(n-2,v1,v2,v3,v4))}
.FUNC PN_3(n,v1,v2,v3,v4) { PN_2(n,v1,v2,v3,v4)*
         +R1L(n-2,v1,v2,v3,v4)+R2R(n-2,v1,v2,v3,v4))/(R1R(n-3,v1,v2,v3,v4)+R2L(n-3,v1,v2,v3,v4))}
.FUNC PN_4(n,v1,v2,v3,v4) { PN_3(n,v1,v2,v3,v4)*
         +R1L(n-3,v1,v2,v3,v4)+R2R(n-3,v1,v2,v3,v4))/(R1R(n-4,v1,v2,v3,v4)+R2L(n-4,v1,v2,v3,v4))}
.FUNC PN_5(n,v1,v2,v3,v4) { PN_4(n,v1,v2,v3,v4)*
         +R1L(n-4,v1,v2,v3,v4)+R2R(n-4,v1,v2,v3,v4))/(R1R(n-5,v1,v2,v3,v4)+R2L(n-5,v1,v2,v3,v4))}
.FUNC PN1(n,v1,v2,v3,v4) {(R2L(n,v1,v2,v3,v4)+R1R(n,v1,v2,v3,v4))/(R2R(n+1,v1,v2,v3,v4)+R1L(n+1,v1,v2,v3,v4))}
.FUNC PN2(n,v1,v2,v3,v4) { PN1(n,v1,v2,v3,v4)*
         +R2L(n+1,v1,v2,v3,v4)+R1R(n+1,v1,v2,v3,v4))/(R2R(n+2,v1,v2,v3,v4)+R1L(n+2,v1,v2,v3,v4))}
.FUNC PN3(n,v1,v2,v3,v4) { PN2(n,v1,v2,v3,v4)*
         +R2L(n+2,v1,v2,v3,v4)+R1R(n+2,v1,v2,v3,v4))/(R2R(n+3,v1,v2,v3,v4)+R1L(n+3,v1,v2,v3,v4))}
.FUNC PN4(n,v1,v2,v3,v4) { PN3(n,v1,v2,v3,v4)*
         +R2L(n+3,v1,v2,v3,v4)+R1R(n+3,v1,v2,v3,v4))/(R2R(n+4,v1,v2,v3,v4)+R1L(n+4,v1,v2,v3,v4))}
.FUNC PN5(n,v1,v2,v3,v4) { PN4(n,v1,v2,v3,v4)*
         +R2L(n+4,v1,v2,v3,v4)+R1R(n+4,v1,v2,v3,v4))/(R2R(n+5,v1,v2,v3,v4)+R1L(n+5,v1,v2,v3,v4))}
```

**Figure 20.** Source code of SPICE subcircuits implemented to SmartSpice for multistate (11 states) approximation of the full-analytical modeling [51]. Reprinted with permission from [51]. G. Lientschnig et al., *Jpn. J. Appl. Phys.* 42, 6467 (2003). © 2003, Institute of Pure and Applied Physics.

```
.FUNC PSUM(n,v1,v2,v3,v4) { PN_5(n,v1,v2,v3,v4)+PN_4(n,v1,v2,v3,v4)+PN_3(n,v1,v2,v3,v4)+PN_2(n,v1,v2,v3,v4)
                            +PN_1(n,v1,v2,.3,v4)+1+PN1(n,v1,v2,v3,v4)+PN2(n,v1,v2,v3,v4)+PN3(n,v1,v2,v3,v4)
                            +PN4(n,v1,v2,v3,v4)+PN5(n,v1,v2,v3,v4) }
*********  calculate the current from source to drain *********
.FUNC CUR(n,v1,v2,v3,v4) { PN_5(n,v1,v2,v3,v4)*(R1R(n-5,v1,v2,v3,v4)-R1L(n-5,v1,v2,v3,v4))
                          +PN_4(n,v1,v2,v3,v4)*(R1R(n-4,v1,v2,v3,v4)-R1L(n-4,v1,v2,v3,v4))
                          +PN_3(n,v1,v2,v3,v4)*(R1R(n-3,v1,v2,v3,v4)-R1L(n-3,v1,v2,v3,v4))
                          +PN_2(n,v1,v2,v3,v4)*(R1R(n-2,v1,v2,v3,v4)-R1L(n-2,v1,v2,v3,v4))
                          +PN_1(n,v1,v2,v3,v4)*(R1R(n-1,v1,v2,v3,v4)-R1L(n-1,v1,v2,v3,v4))
                          +(R1R(n,v1,v2,v3,v4)-R1L(n,v1,v2,v3,v4))
                          +PN1(n,v1,v2,v3,v4)*(R1R(n+1,v1,v2,v3,v4)-R1L(n+1,v1,v2,v3,v4))
                          +PN2(n,v1,v2,v3,v4)*(R1R(n+2,v1,v2,v3,v4)-R1L(n+2,v1,v2,v3,v4))
                          +PN3(n,v1,v2,v3,v4)*(R1R(n+3,v1,v2,v3,v4)-R1L(n+3,v1,v2,v3,v4))
                          +PN4(n,v1,v2,v3,v4)*(R1R(n+4,v1,v2,v3,v4)-R1L(n+4,v1,v2,v3,v4))
                          +PN5(n,v1,v2,v3,v4)*(R1R(n+5,v1,v2,v3,v4)-R1L(n+5,v1,v2,v3,v4)) }

.FUNC CURRENT(n,v1,v2,v3,v4 { I*CUR(n,v1,v2,v3,v4)/(CSUM*PSUM(n,v1,v2,v3,v4)*(R1+R2)) }
***********  ******  ****  calculate the island voltage *****************************
.FUNC VOLT(n,v1,v2,v3,v4) { PN_5(n,v1,v2,v3,v4)*(n-5+Q(v1,v2,v3,v4))
                           +PN_4(n,v1,v2,.3,v4)*(n-4+Q(v1,v2,v3,v4))
                           +PN_3(n,v1,v2,.3,v4) (n-3+Q(v1,v2,v3,v4))
                           +PN_2(n,v1,v2,.3,v4)*(n-2+Q(v1,v2,v3,v4))
                           +PN_1(n,v1,v2,.3,v4)*(n-1+Q(v1,v2,v3,v4))
                           +n+Q(v1,v2,v3,v4)
                           +PN1(n,v1,v2,v3,v4)*(n+1+Q(v1,v2,v3,v4))
                           +PN2(n,v1,v2,v3,v4)*(n+2+Q(v1,v2,v3,v4))
                           +PN3(n,v1,v2,v3,v4)*(n+3+Q(v1,v2,v3,v4))
                           +PN4(n,v1,v2,v3,v4)*(n+4+Q(v1,v2,v3,v4))
                           +PN5(n,v1,v2,v3,v4)*(n+5+Q(v1,v2,v3,v4)) }
.FUNC VOLTAGE(n,v1,v2,v3,v4 { (E/CSUM)*VOLT(n,v1,v2,v3,v4)/PSUM(n,v1,v2,v3,v4) }
E1 5 0 VALUE={VOLTAGE(N_OPT(V(1),V(2),V(3),V(4)),V(1),V(2),V(3),V(4))}   ; Voltage of the island
G1 1 2 VALUE={CURRENT(N_OPT(V(1),V(2),V(3),V(4)),V(1),V(2),V(3),V(4))}   ; Current from source to dran
CT1 1 5 {Cs}
CT2 2 5 {Cd}
CGATE1 3 5 {Cg}
CGATE2 4 5 {Cb}
.ENDS SET
```

Figure 20. Continued.

electrode, while the second term is the polarization-induced current (displacement current) due to tunneling through other junctions and the direct influence of voltage changes. The first term is defined simply by [39, 69]

$$I_n^{tun}(t) = e \sum_{i,j,\Delta k} P_i(t) \Gamma_{ij}^{\Delta k}(t) \Delta k_n \tag{17}$$

where $\Gamma_{ij}^{\Delta k}$ is the rate of transition rate from $i$ state to $j$, which leads to a transfer of charge $e\Delta k_n$ into external electrode $n$.

The induced part of the current (displacement current) through the external electrode $n$ can be calculated as [39, 43, 70]

$$I_n^{ind}(t) = \frac{dQ_n^{ind}(t)}{dt} \tag{18}$$

where the induced charge $Q_n^{ind}$ is given by [39]

$$Q_n^{ind}(t) = \sum_{i,m} P_i(t) C_{mn} V_{mn}(t) \tag{19}$$

Here the summation is over all states, all capacitances $C_{mn}$ connected to external electrode $n$, and $V_{mn}$ are the voltages across the capacitances $C_{mn}$.

## 3.3. Comparison Between Monte Carlo Method and Master Equation Method

Excellent dynamic and transient characteristics of single-electron circuits are obtained by using the Monte Carlo method and the master equation method because in the Monte Carlo method, electron tunneling between islands is simulated in a very direct manner, and in the

Figure 21. Comparison of the single-electron transistor simulations performed with multistate (11 states) approximation of the full-analytical modeling (black lines) with the ones performed with a conventional Monte Carlo simulation program (gray lines). In all three simulations, the parameters used were: $V_{b} = 0$ V. $C_{d} = C_{s} = C_{g} = 1$ aF. $C_{o} = 0$, $R_{d} = R_{s} = 100$ kΩ. and $T = 4.2$ K. (a) The current through the SET transistor is plotted as a function of the bias voltage for six gate voltages, $R_{d} = R_{s} = 100$ kΩ and $T = 4.2$ K. (b) The current through the SET transistor is plotted against the gate voltage for bias voltages $V_{b} = 5$ mV to 150 mV in steps of 5 mV. $R_{d} = R_{s} = 100$ kΩ. and $T = 4.2$ K. (c) Current–voltage characteristics are plotted for three different temperatures. $V_{g} = 0$ V. $R_{s} = 100$ kΩ. and $R_{d} = 1$ MΩ. Reprinted with permission from [51]. G. Lientschnig et al., *Jpn. J. Appl. Phys.* 42, 6467 (2003). © 2003, Institute of Pure and Applied Physics.

master equation method, the time-dependent master equation is directly solved. There is one major disadvantage of the Monte Carlo method. In cases when very rare tunneling events of some type take place against the background of much more frequent events of another kind, the Monte Carlo method becomes impractical because of the demand for very long simulation time, but the master equation method can include the code to calculate rare errors due to co-tunneling. There is one major disadvantage of the master equation method. In the master equation method, the number of states that have to be considered becomes often very large. Consequently, the matrix operations involved in calculating the exponential operator are time consuming and the approximations do not converge quickly. Numerical instabilities can easily appear. The master equation method needs to include more than relevant states to correctly simulate the circuits, but the Monte Carlo method is not required to find the relevant states before starting the simulation. To overcome those problems of the Monte Carlo method and the master equation, a new algorithm that combines the Monte Carlo method and the master equation is developed, which is explained in the next subsection.

## 3.4. Method Combined with Monte Carlo and Master Equation

Although the master equation method gives theoretically accurate results, it has many other impracticalities that limit its accuracy and usability. The starting point of the master equation is the set of all relevant states that a circuit will occupy during operation. In order to complete correctly the simulation, many states that would be relevant have to be included, which results in extremely long simulation times and sometimes bad numerical stability.

**Figure 22.** Simplified flow chart to calculate the SET drain currents for quasi-analytical modeling. Reprinted with permission from [47], S. Mahapatra et al., *IEEE Electron. Device Lett.* 23, 366 (2002). © 2002, IEEE.

To overcome this problem, SIMON [37] combines the advantages of the Monte Carlo method and the master equation. All possible states are divided into two subspaces, the frequent state space and the rare state space as shown in Fig. 4 [37]. The Monte Carlo method is used to simulate only the frequent state space, which gives the occupation probabilities $P_i$ of frequent states. The occupation probability $P_i$ is calculated as the ratio of time $T_i$ spent in state $i$ to the total simulation time $T_\Sigma$ and is given by [37]

$$P_i = \frac{T_i}{T_\Sigma} \tag{20}$$



**Figure 23.** Drain current–gate voltage $I_D$ – characteristics and the transconductance $dI_{D}/dV$ – characteristics of the single and double-gate SET calculated by quasi-analytical modeling and the reference simulator at various $V$ . Reprinted with permission from [47], S. Mahapatra et al., *IEEE Electron. Device Lett.* 23, 366 (2002). © 2002, IEEE.

**Figure 24.** $I_{ds}$-$V_{ds}$ characteristics of SETs calculated by the Monte Carlo simulator SIMON (solid lines), the macro-modeling (open squares), the seminumerical modeling (filled circles), and the full-analytical modeling (open triangles) at the DC analysis ($V_{gs} = 0$ V, $C_d = C_s = C_g = 1$ aF, $C_b = 0$, $R_d = R_s = 10$ M$\Omega$, and $T = 30$ K).

with

$$T_\Sigma = \sum_i T_i \tag{21}$$

Instead of waiting for the Monte Carlo simulator to step into the rare state space, which would result in impractically long simulation times, we directly calculate the contribution of events leading to rare states by stepping through the event tree starting at frequent states as shown in Fig. 4. The essential assumption is that the time spent in the rare state $j$, $T_{j,rare} \ll T_\Sigma$, because the rare states cause only a small perturbation to the frequent state probabilities. Here, $T_{j,rare}$ is given by by [37]

$$T_{j,rare} = \frac{\sum_i (T_i \Gamma_{ij})}{\sum_k \Gamma_{jk}} \tag{22}$$

where $T_i \Gamma_{ij}$ is the number of times that the rare state $j$ would be on average visited from state $i$. $\sum_{k \neq j} \Gamma_{kj}$ is the exit rate of state $j$, and thus $1/\sum_{k \neq j} \Gamma_{kj}$ is the average time spent in state $j$ for one visit. Time averages are used for the direct calculation. Actually, the durations



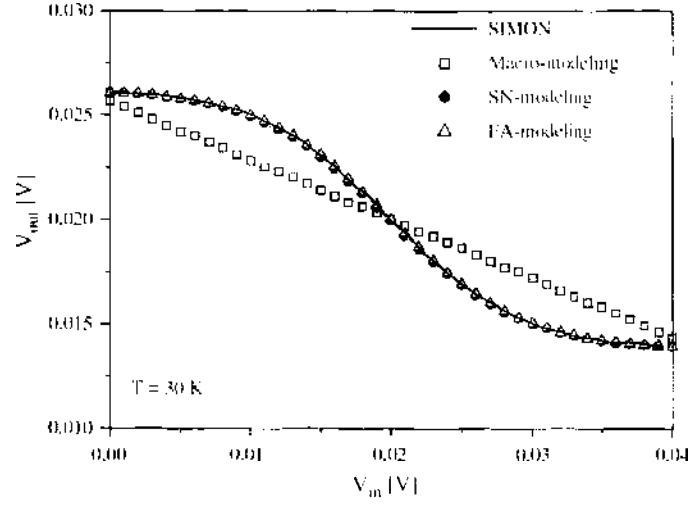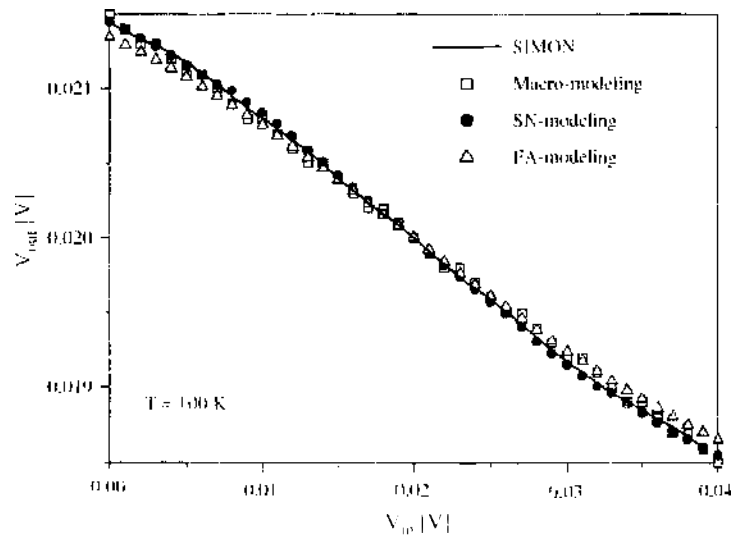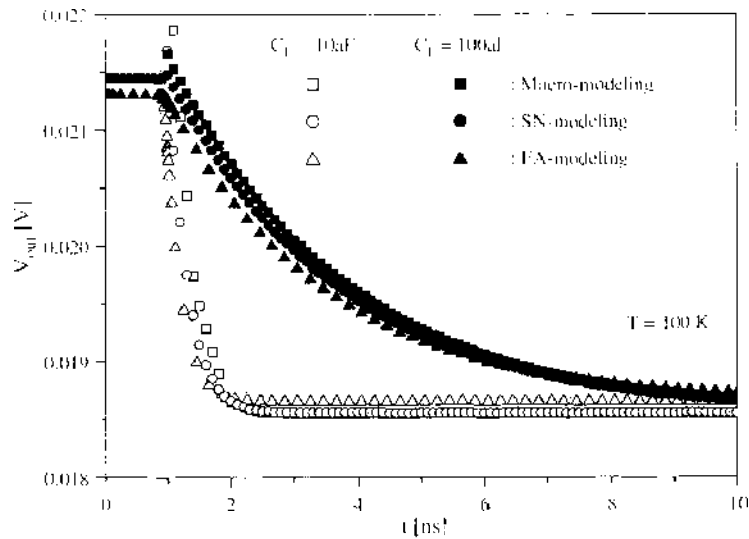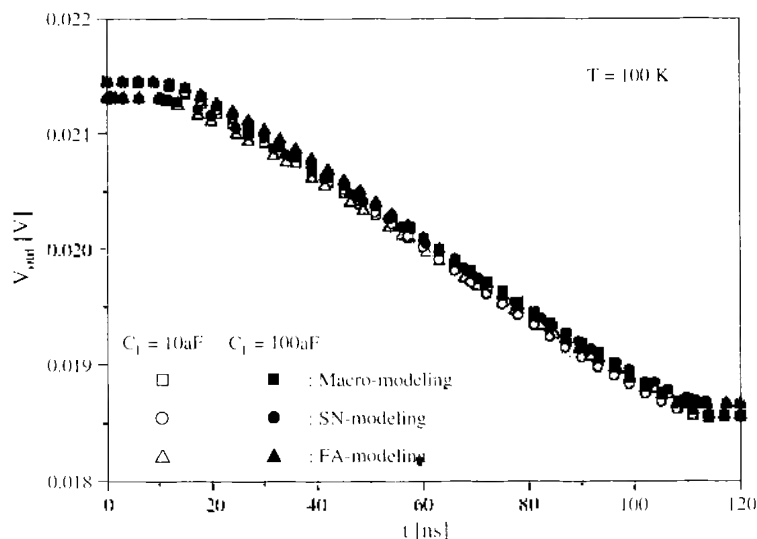**Figure 25.** $I_{ds}$-$V_{gs}$ characteristics of SETs calculated by the Monte Carlo simulator SIMON (solid lines), the macro-modeling (open squares), the seminumerical modeling (filled circles), and the full-analytical modeling (open triangles) at the DC analysis when $T = 30$ K and $T = 100$ K. ($V_{ds} = 26.7$ mV, $C_d = C_s = C_g = 1$ aF, $R_d = R_s = 10$ M$\Omega$, and $C_b = 0$).

are distributed with a Poisson distribution. However, the directly calculated rare state space is only a small perturbation to the frequent state space, which is calculated with a Monte Carlo method where the Poisson distribution of the tunnel durations is fully incorporated. Therfore, the occupation probability for a rare state is given by [37].

$$P_{j,rare} = \frac{T_{j,rare}}{T_\Sigma} = \frac{\sum_i (T_i \Gamma_{ij})}{T_\Sigma \sum_k \Gamma_{jk}} \tag{23}$$

The algorithm follows all possible events starting at frequent states. If a frequent state is encountered, the algorithm terminates that branch, because the state probability is already known. Once the probability of a state is lower than a predefined limit no further descent into following branches from this state is made as shown in Fig. 4.

## 4. FEATURES OF SET CIRCUITS FOR THE SPICE APPLICATION

In the case of the circuits with SETs, it is known that the second assumption in SPICE may not be valid. The terminal currents of the SET are determined from the average charge state of the Coulomb island of the SET. When several SETs are connected, the charge state of the Coulomb island of one SET is strongly affected by the charge states of neighboring islands of other SETs. Therefore, the terminal currents of the SET in the circuit may be different from those of the isolated SET even at the same bias condition. The following numerical example shows such interaction among neighboring Coulomb islands more intuitively. Figure 5 shows a single-electron inverter (SEI), consisting of a series combination of two SETs. There are three Coulomb islands (two from the SETs and one from the interconnection) in the circuit, and the charge states of these three Coulomb islands are correlated. The filled squares and open squares in Fig. 6 show the calculated voltage transfer characteristics and current of the lower SET of the inverter, respectively. The open circles are the currents of an isolated SET when the terminals are biased at the same voltages as the lower SET in the inverter. Even at the same bias condition, the terminal current of the isolated SET is expected to be totally different from that of the SET in the circuit.

On the other hand, when the size of the interconnection is large enough, the interconnection would serve as a reservoir for neighboring SETs rather than a Coulomb island. In that case, the Coulomb islands of SETs became isolated by the interconnection, and the interaction among neighboring SETs and interconnections may not be significant such that applying the conventional circuit simulation techniques to the simulation of SET circuits is possible. Accordingly, simulation of hybrid circuits consisting of SET circuits and conventional circuits is possible.

## 5. REGIME OF THE COMPACT MODELING

In this section, it is illustrated that the $I$-$V$ characteristics of SETs can be calculated independently with one another when the sizes of the interconnections are large enough. Furthermore, using the case of the single electron inverter, the range of this conventional circuit simulation regime is systematically identified.

### 5.1. DC Regime

First of all, the same calculation as that of Fig. 6 has been repeated with a large interconnect capacitance, and the results are shown in Fig. 7. It can be clearly seen that the terminal current of the lower transistor (open symbols) does not change whether the SET is in the circuit or is isolated, only if the terminal bias conditions are the same. In this case, the transfer characteristic of the inverter can be calculated by the $I$-$V$ curves obtained from the isolated SETs and by applying the usual Kirchhoff circuit laws. The filled circles in Fig. 4 are the transfer characteristic obtained from such methods. They are in good agreement with the results resorting to the orthodox theory (filled squares) [36–38], considering the interconnection as a Coulomb island.

The regime of the compact modeling can be identified by studying the discrepancies between the results of the above two methods as a function of the interconnect capacitance,

$C_j$. The simulation parameters are the bias voltage $V_B = 0.03$ V, the tunneling resistance $R_t = 100M\Omega$, the junction tunnel capacitor $C_1 = C_2 = C_3 = C_4 = 1.6$ aF, and the gate normal capacitor $C_g = 3.2$ aF. Figures 8a and 8b summarize the difference in the output voltage swing, $V_{out}^{swing}$, and the average output level, $V_{out}^{avg}$, of the transfer characteristics as a function of $C_j$, respectively. It is clearly seen that the conventional circuit law type method with the $I-V$s of isolated SETs works well when $C_I > 6.25C_j$, where $C_j$ is the value of a typical junction capacitance.

The above criterion will be satisfied in most cases mainly because the interconnect is connected to several gate capacitances of the next stage. Furthermore, the charge states of Coulomb islands are interacting mainly with the nearest neighbors, and the regime of $C_j$ obtained from the inverter can be generally used in other types of circuits.

## 5.2. Transient Regime

Two types of transient simulations are performed for the SEI. The first type of calculation is performed by solving the time-dependent master equation (TME) considering the overall probability distribution of three Coulomb islands (two islands of the lower and the upper SET and the interconnection island) in the circuit, which is used to SENECA [39]. Solving two time-dependent master equations of the lower and the upper SETs independently and applying Kirchhoffs law at the interconnection performs the second type of calculation, which is used to SPICE. Figure 9 shows the maximum difference between the output voltages $V_{out}$ obtained from two types of calculations when the input voltage $V_{in}$ of the inverter is swept from 0 to 0.03 V in three different $t_r$s. The difference between the two results is negligible when the interconnection capacitance $C_I$ is larger than 64 aF and the deviation increases rapidly with the decrease of $C_I$. This result suggests that, as in the steady-state case, each SET can be treated independently even in the transient case when the interconnection capacitance is large enough. However, the value of $C_I$ for the independent treatment of SETs is approximately 10 times larger than that of the steady-state case [42]. This will be more clearly seen in the following sections where the compact models are applied for more complicated circuits.

# 6. COMPACT MODELING OF THE SINGLE-ELECTRON TRANSISTORS

For the application of full conventional simulation techniques to single-electron circuits, a compact modeling is essential for the $I-V$ characteristics of isolated SETs rather than the $I-V$ characteristics obtained from the Monte Carlo method. In this section, the compact modeling for SETs, which is fully compatible to SPICE, is introduced.

## 6.1. Macromodeling

Figure 10 shows the SET equivalent circuit and its macromodel code of an SET for the HSPICE [71] simulation. Symmetric features of the drain-source current-voltage $(I_{ds}-V_{ds})$ characteristics are incorporated with two branches consisting of combinations of resistors, diodes, and voltage sources. They are denoted by $R2/D2/V2$ and $R3/D3/V3$, respectively. The directions of $D2$ and $V2$ are opposite to those of $D3$ and $V3$ to have adequate current flow in both positive and negative drain-source bias. The charging energy, periodically changing as a function of the gate bias, is included in $R1$, $R2$, and $R3$ where the cosine of the gate bias is used. They are expressed as follows [42]:

$$R_1(V_G) = CR1 + CR2\cos(CF1 \cdot \pi \cdot V_G + CP1) \tag{24-1}$$

$$R_2(V_G) = R_3(V_G) = \frac{CVp}{CI2 - 2CVp/R_1(V_G)} \tag{24-2}$$

The parameters, $CF1$, $CVp$, $CI2$, $CP1$, $CR1$, and $CR2$ are used to fit the current-voltage characteristics at various gate biases. The gate-source capacitance $C_{gs}$, the gate-drain capacitance $C_{gd}$, and the drain-source capacitance $C_{ds}$ in the proposed SET equivalent circuit

Table 2. Summary of the differences between simulation results of four types of the compact SET modeling and the Monte Carlo simulation results in the DC analysis [50].

| Type | Low $T$ ($< 0.1\ e^2/k_B C_\Sigma$) | | High $T$ ($\geq 0.1\ e^2/k_B C_\Sigma$) | |
| | Low $V_{ds}$ ($< e/C_\Sigma$) | High $V_{ds}$ ($\geq e/C_\Sigma$) | Low $V_{ds}$ ($< e/C_\Sigma$) | High $V_{ds}$ ($\geq e/C_\Sigma$) |
|---|---|---|---|---|
| Macromodeling | Large | Large | Small | Small |
| Seminumerical modeling | No | No | No | No |
| Full-analytical modeling | No | Large | Small | Large |

Large, small, and no in Table 2 more than 10% deviation, from 1% to 10% deviation, and below 1% deviation, respectively.

are expressed as $C_g C_s/C_\Sigma$, $C_g C_d/C_\Sigma$, and $C_d C_s/C_\Sigma$, respectively, based on the principle of charge conservation. Figure 11 shows the simulated $I-V$ characteristics of an SET at various gate biases. The solid lines are the Monte Carlo calculation results and the open symbols are the macromodel calculation results. With a proper choice of the parameters, $CF1 = 40$, $CVp = 0.02$, $CI2 = 0.2 \times 10^{-9}$, $CP1 = 0$, $CR1 = 300 \times 10^6$, and $CR2 = 100 \times 10^6$, the macromodel calculations reproduce the Monte Carlo calculations reasonably well. Because the $I_{ds}-V_{ds}$ characteristics of SETs strongly depend on $T$, the macromodel parameters are the functions of $T$. Figure 12 shows the $I_{ds}-V_{ds}$ characteristics of the SET in Fig. 11 at various $T$s. Again, with the proper choice of parameters, the macromodel calculations can reproduce the Monte Carlo results reasonably well. The parameter values for the results of Fig. 12 are summarized in Table 1.

## 6.2. Seminumerical Modeling

The equivalent circuit of the SET under the isolated SET approximation is constructed using the same idea as the CMOS equivalent circuit. Figure 13 shows the schematic diagram (left figure) and the linearized equivalent circuit model (right figure) of the SET. The elements in the equivalent circuit are given by [49].

$$
\begin{bmatrix} J_{ds} \\ J_{gs} \\ J_{bs} \end{bmatrix} = \begin{bmatrix} g_{ds} & g_{dmg} & g_{dmb} \\ g_{gmd} & g_{gs} & g_{gmb} \\ g_{bmd} & g_{bmg} & g_{bs} \end{bmatrix} \begin{bmatrix} V_{ds} \\ V_{gs} \\ V_{bs} \end{bmatrix} - \begin{bmatrix} I_d \\ I_g \\ I_b \end{bmatrix} \quad \begin{bmatrix} g_{ds} & g_{dmg} & g_{dmb} \\ g_{gmd} & g_{gs} & g_{gmb} \\ g_{bmd} & g_{bmg} & g_{bs} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial I_d}{\partial V_{ds}} & \dfrac{\partial I_d}{\partial V_{gs}} & \dfrac{\partial I_e}{\partial V_{bs}} \\ \dfrac{\partial I_g}{\partial V_{ds}} & \dfrac{\partial I_g}{\partial V_{gs}} & \dfrac{\partial I_g}{\partial V_{bs}} \\ \dfrac{\partial I_b}{\partial V_{ds}} & \dfrac{\partial I_b}{\partial V_{gs}} & \dfrac{\partial I_b}{\partial V_{bs}} \end{bmatrix} \quad (25)
$$

Figure 26. The drain ($I$) and the source current ($I$) of an SET calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (squares), the seminumerical modeling (circles), and the full-analytical modeling (triangles) when $T$ varies linearly from 0 to 0.1 V in $t = 1\ \mu s$ ($I_0 = 26.7$ mV, $C_g = C_s = C_d = 1$ aF, $C_b = 0$, $R_s = R_d = 10$ M$\Omega$ and $T = 30$ K).

**Figure 27.** The drain ($I_d$) and the source current ($-I_s$) of an SET calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (squares), the seminumerical modeling (circles), and the full-analytical modeling (triangles) when $V_g$ varies linearly from 0 to 0.1 V in $t_r = 1$ ns ($V_{ds} = 26.7$ mV, $C_d = C_s = C_g = 1$ aF, $C_b = 0$, $R_d = R_s = 10$ MΩ, and $T = 30$ K).

The terminal current of the SET is consists of the tunneling current across the tunnel junction and the displacement current. For example, the drain current $I_d$ is expressed as

$$I_d(t_q) = I_d^{tun}(t_q) + I_d^{id}(t_q) + I_d^{ed}(t_q) \tag{26}$$

where $I_d^{tun}(t_q)$ is the tunneling current across the drain tunnel junction, $I_d^{id}(t_q)$ is the internal displacement current, and $I_d^{ed}(t_q)$ is the external displacement current. The first term is given by

$$I_d^{tun}(t_q) = e \sum_{n=-\infty}^{\infty} [\Gamma_d^{an}(n, t_q) - \Gamma_d^{dn}(n, t_q)] P_n(t_q) \tag{27}$$

where

$$\Gamma_d^a(n, t_q) = \frac{1}{e^2 R_d} \frac{\Delta E_d^a(n, t_q)}{1 - \exp\left(-\frac{\Delta E_d^a(n, t_q)}{k_B T}\right)} \tag{28}$$

$$\Delta E_d^a(n, t_q) = \begin{cases} e V_d'(t_q) + e^2 \left(\dfrac{-2n-1}{2C_\Sigma}\right), & a = \text{inc} \\[2ex] -e V_d'(t_q) + e^2 \left(\dfrac{2n-1}{2C_\Sigma}\right), & a = \text{dec} \end{cases} \tag{29}$$

$$V_d'(t_q) = \frac{1}{C_\Sigma}\left[ \sum_{k=s,g,b} (-C_k V_k(t_q)) + (C_\Sigma - C_d) V_d(t_q) \right] \tag{30}$$

Here, $e$ is the positive elementary charge, $n$ is an integer that specifies the number of elementary charges added in the Coulomb island, $R_d$ is the tunneling resistance of the drain tunnel junction, $C_\Sigma = C_d + C_g + C_s + C_b$ is the total capacitance of the Coulomb island of the transistor, $V_k(t_q)$ is the $k$-terminal voltage, $V_d'(t_q)$ is the voltage drop between the drain and the Coulomb island, and $\Delta E_d^a(n, t_q)$ is the electrostatic energy difference when the charge in the Coulomb island is changed from $ne$ to $(n+1)e$ ($a = \text{inc}$) or to $(n-1)e$ ($a = \text{dec}$) due to the tunneling between the drain and the Coulomb island. Finally, $\Gamma_d^a(n, t_q)$ is the corresponding tunneling rate.

The term $P_n(t_q)$ is the probability that the charge in the Coulomb island in an SET equals $ne$ at time $t_q$. Every successive iteration step requires a new $P_n(t_q)$, as $P_n(t_q)$ is a function of the terminal voltages. The sampling time interval $t_s (= t_{q+1} - t_q)$ is discretized

**Figure 28.** The voltage transfer characteristics of an SET calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (open squares), the seminumerical modeling (filled circles), and the full-analytical modeling (open triangles) at the DC analysis ($V_{cc} = 26.7$ mV, $C_j = C_s = C_g = 1$ aF, $C_k = 0$, $R_{tj} = R = 10$ MΩ, and $T = 30$ K).

with the interval of $t$, and the evolution of $P_n(t_q)$ is obtained from the time-dependent master equation [15, 39, 43].

$$P_n(t_q + (v+1)\Delta t) = \Gamma_{n,n-1}(t_q + v\Delta t)\Delta t P_{n+1}(t_q + v\Delta t) + \Gamma_{n,n+1}(t_q + v\Delta t)\Delta t P_{n-1}(t_q + v\Delta t)$$
$$+ [1 - \Gamma_{n+1,n}(t_q + v\Delta t)\Delta t - \Gamma_{n-1,n}(t_q + v\Delta t)\Delta t]P_n(t_q + v\Delta t), \quad (31)$$

$$\sum_n P_n(t_q + v\Delta t) = 1$$

where

$$\Gamma_{n+1,n}(t_q + v\Delta t) = \Gamma_d^{un}(n, t_q + v\Delta t) + \Gamma_s^{un}(n, t_q + v\Delta t),$$

$$\Gamma_{n-1,n}(t_q + v\Delta t) = \Gamma_d^{dec}(n, t_q + v\Delta t) + \Gamma_s^{dec}(n, t_q + v\Delta t), \quad (32)$$

$$v = 0, 1, 2, \ldots\ldots, M, \quad \text{and} \quad M = \text{int}\left(\frac{t_{q+1} - t_q}{\Delta t}\right)$$



**Figure 29.** The voltage transfer characteristics of an SET calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (open squares), the seminumerical modeling (filled circles), and the full-analytical modeling (open triangles) at the DC analysis ($V_{cc} = 26.7$ mV, $C_j = C_s = C_g = 1$ aF, $C_k = 0$, $R_{tj} = R = 10$ MΩ, and $T = 100$ K).

**Figure 30.** The transient response of an SEI calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (open squares), the seminumerical modeling (filled circles), and the full-analytical modeling (open triangles) when $I^*_{ex}$ varies linearly from 0 to 0.04 V in $t_r = 0.1$ ns ($V_{ds} = 26.7$ mV, $C_d = C_s = C_g = 1$ aF, $C_b = 0$, $R_d = R_s = 10$ MΩ, and $T = 100$ K).

In the simulation, a threshold probability $P_{th}$ is set so that only a finite number of probabilities are included [39, 43]. The second term of Eq. (26) is the internal displacement current induced from the potential change of the Coulomb island by the tunneling across the drain and the source tunnel junction [70], and it is given by

$$I_d^{id}(t_q) = -\frac{C_d}{C_\Sigma}[I_s^{tun}(t_q) + I_d^{tun}(t_q)] \tag{33}$$

The third term is the external displacement current induced from the change of the terminal voltages. It is given by

$$I_d^{ed}(t_q) = C_d \frac{dV_d'(t_q)}{dt} \tag{34}$$

Other terminal currents are obtained similarly.



**Figure 31.** The transient response of an SEI calculated by the Monte Carlo simulator SIMON (solid lines), the macromodeling (open squares), the semi-numerical modeling (filled circles), and the full-analytical modeling (open triangles) when $V_{ex}$ varies linearly from 0 to 0.04 V in $t_r = 100$ ns ($V_{ds} = 26.7$ mV, $C_d = C_s = C_g = 1$ aF, $C_b = 0$, $R_d = R_s = 10$ MΩ, and $T = 100$ K).

**Figure 32.** The simulation times of several types of SET circuits in the transient simulation. The x-axis denotes the complexity of the circuit: 1 (single electron inverter). 2 (an SE-NOR gate). 3 (two SE-NOR gates). 4 (an SE-OR gate; three SE-NOR gates).

Finally, the conductance terms are evaluated from the derivatives of the terminal current. For example,

$$g_{ds} = \frac{\partial I_d(t_q)}{\partial V_{ds}(t_q)} = \frac{\partial I_d^{tun}(t_q)}{\partial V_{ds}(t_q)} - \frac{C_d}{C_\Sigma}\left(\frac{\partial I_s^{tun}(t_q)}{\partial V_{ds}(t_q)} + \frac{\partial I_d^{tun}(t_q)}{\partial V_{ds}(t_q)}\right) + \frac{\partial I_d^{cd}(t_q)}{\partial V_{ds}(t_q)} \qquad (35)$$



**Figure 33.** Examples of two-input SET logic gates with three terminal SETs. Schematics of (a) two-input NOR [58]. (b) two-input XOR [58]. and (c) two-input OR gates [74]. Reprinted with permission from [58]. R. Chen et al.. *Appl. Phys. Lett.* 68. 1954 (1996). © 1996. American Institute of Physics. Reprinted with permission from [74]. I. Karafyllidis et al.. *Electron. Lett.* 36. 407 (2000). © 2000. IEE.

The first term $\partial I_d^{tun}(t_q)/\partial V_{ds}(t_q)$ is given by

$$\frac{\partial I_d^{tun}(t_q)}{\partial V_{ds}(t_q)} = e \sum_{n=-\infty}^{\infty} \left\{ \left[ \frac{\partial \Gamma_d^{inc}(n,t_q)}{\partial V_{ds}(t_q)} - \frac{\partial \Gamma_d^{dec}(n,t_q)}{\partial V_{ds}(t_q)} \right] P_n(t_q) + \left[ \Gamma_d^{inc}(n,t_q) - \Gamma_d^{dec}(n,t_q) \right] \frac{\partial P_n(t_q)}{\partial V_{ds}(t_q)} \right\}$$

(36)

where

$$\frac{\partial \Gamma_d^a(n,t_q)}{\partial V_{ds}(t_q)} = \frac{1}{e^2 R_d} \frac{\dfrac{\partial \Delta E_d^a(n,t_q)}{\partial V_{ds}(t_q)} \left[ 1 - \left( 1 + \dfrac{\Delta E_d^a(n,t_q)}{k_B T} \right) \exp\left( -\dfrac{\Delta E_d^a(n,t_q)}{k_B T} \right) \right]}{\left[ 1 - \exp\left( -\dfrac{\Delta E_d^a(n,t_q)}{k_B T} \right) \right]^2}$$

(37)

$$\frac{\partial \Delta E_d^a(n,t_q)}{\partial V_{ds}(t_q)} = \begin{cases} e \dfrac{\partial V_d'(t_q)}{\partial V_{ds}(t_q)}, & a = inc, \\[12pt] -e \dfrac{\partial V_d'(t_q)}{\partial V_{ds}(t_q)}, & a = dec \end{cases}$$

(38)

$$\frac{\partial V_d'(t_q)}{\partial V_{ds}(t_q)} = \frac{(C_\Sigma - C_d)}{C_\Sigma}$$

(39)

$$\frac{\partial P_0(t_q)}{\partial V_{ds}(t_q)} = -P_0^2(t_q) \sum_{n=1}^{\infty} \left\{ \prod_{s=-n+1}^{n} \frac{\Gamma_{s-1,s}(t_q)}{\Gamma_{s,s-1}(t_q)} \sum_{s=-n+1}^{n} \left[ \frac{1}{\Gamma_{s-1,s}(t_q)} \frac{\partial \Gamma_{s-1,s}(t_q)}{\partial V_{ds}(t_q)} - \frac{1}{\Gamma_{s,s-1}(t_q)} \frac{\partial \Gamma_{s,s-1}(t_q)}{\partial V_{ds}(t_q)} \right] \right.$$

$$\left. + \prod_{s=0}^{n-1} \frac{\Gamma_{s+1,s}(t_q)}{\Gamma_{s,s+1}(t_q)} \sum_{s=0}^{n-1} \left[ \frac{1}{\Gamma_{s+1,s}(t_q)} \frac{\partial \Gamma_{s+1,s}(t_q)}{\partial V_{ds}(t_q)} - \frac{1}{\Gamma_{s,s+1}(t_q)} \frac{\partial \Gamma_{s,s+1}(t_q)}{\partial V_{ds}(t_q)} \right] \right\}$$

$$\frac{\partial P_n(t_q)}{\partial V_{ds}(t_q)} = P_n(t_q) \left[ \sum_{s=-n+1}^{0} \left( \frac{1}{\Gamma_{s-1,s}(t_q)} \frac{\partial \Gamma_{s-1,s}(t_q)}{\partial V_{ds}(t_q)} - \frac{1}{\Gamma_{s,s-1}(t_q)} \frac{\partial \Gamma_{s,s-1}(t_q)}{\partial V_{ds}(t_q)} \right) \right.$$

$$\left. + \frac{1}{P_0(t_q)} \frac{\partial P_0(t_q)}{\partial V_{ds}(t_q)} \right], \quad n < 0$$

(40)

$$\frac{\partial P_n(t_q)}{\partial V_{ds}(t_q)} = P_n(t_q) \left[ \sum_{s=0}^{n-1} \left( \frac{1}{\Gamma_{s+1,s}(t_q)} \frac{\partial \Gamma_{s+1,s}(t_q)}{\partial V_{ds}(t_q)} - \frac{1}{\Gamma_{s,s+1}(t_q)} \frac{\partial \Gamma_{s,s+1}(t_q)}{\partial V_{ds}(t_q)} \right) \right.$$

$$\left. + \frac{1}{P_0(t_q)} \frac{\partial P_0(t_q)}{\partial V_{ds}(t_q)} \right], \quad n > 0$$

$$\frac{\partial \Gamma_{n+1,n}(t_q)}{\partial V_{ds}(t_q)} = \frac{\partial \Gamma_d^{inc}(n,t_q)}{\partial V_{ds}(t_q)} + \frac{\partial \Gamma_s^{inc}(n,t_q)}{\partial V_{ds}(t_q)}, \quad \frac{\partial \Gamma_{n-1,n}(t_q)}{\partial V_{ds}(t_q)} = \frac{\partial \Gamma_d^{dec}(n,t_q)}{\partial V_{ds}(t_q)} + \frac{\partial \Gamma_s^{dec}(n,t_q)}{\partial V_{ds}(t_q)}$$

(41)

The term $\partial I_s^{tun}(t_q)/\partial V_{ds}(t_q)$ of Eq. (35) is obtained similarly. Also, the third term $I_d^{cd}(t_q)$ and $\partial I_d^{cd}(t_q)/\partial V_{ds}(t_q)$ of the capacitor network of the linearized equivalent circuit in Fig. 13 are calculated by the iterative companion model of a capacitor with the backward Euler approximation [52]. The capacitor network of the SET shown at the schematic diagram in Fig. 13 has changed from that of the linearized equivalent circuit in Fig. 13, based on the principle of conservation of charge. The capacitance $C_{ab}$ of the linearized equivalent circuit in Fig. 13 is equal to $C_a C_b / C_\Sigma$, where $C_a$ is the capacitance between the Coulomb island and the terminal $a$, and $C_b$ is the capacitance between the Coulomb island and the terminal $b$ of the schematic diagram in Fig. 13. The DC equivalent circuit of an SET can automatically be obtained by erasing all the elements except for $J_{ds}$, $g_{ds}$, $g_{dmg}$, $g_{dmb}$ of the linearized equivalent

circuit in Fig. 13. In the calculation of the terminal current, the internal and the external displacement current do not exist. Therefore, $I_g = I_b = 0$ and $I_d = I_d^{tun} = -I_s^{tun} = -I_s$. The probability $P_n$ is obtained from the steady-state master equation (SME) [16, 40]. The seminumerical SET model is implemented to SmartSpice [72] as a user-defined element.

Figure 14a and 14b show the terminal currents $I_d(t)$, $-I_s(t)$, $I_g(t)$, and $I_b(t)$ of the SET in Fig. 10 when gate-source voltage $V_{gs}$ varies linearly from 0 to 0.1 V in 50 ps, drain-source voltage $V_{ds} = 0.1$ V, and backgate-source voltage $V_{bs} = 0.075$ V. The parameter $t = 0.1$ ps is used. All terminal currents calculated from our transient model (solid symbols) almost exactly match with the results obtained from the direct numerical solution of the master equation (solid line). It suggests that the integration of the transient model into the Smart-Spice is properly achieved. The transient currents also show appreciable difference from DC current (open symbols), which suggests that the time-dependent charging of the Coulomb island and the contribution of the displacement current are correctly included in the transient model. Such current components cannot be accounted for in the DC model. Finally, the DC currents calculated from the DC model match with the DC currents calculated from the Monte Carlo method. It is concurrent with our model, as the DC model is the limited case of the transient model. Figure 14c shows $I_d(t)$ when $V_{gs}$ varies as shown in the inset. The transient model predicts that the hysteretic behavior is originated from the sign change of the displacement current according to the sweep direction of the bias.



Figure 34. Examples of two-input SET logic gates with four terminal SETs. Schematics of (a) an SET inverter. (b) two-input XOR. (c) two-input NAND. and (d) two input NOR gates. Reprinted with permission from [75], M. Y. Jeng et al., *Jpn. J. Appl. Phys.* 36, 6706 (1997). © 1997. Institute of Pure and Applied Physics.

## 6.3. Full-Analytical Modeling

### 6.3.1. Two-State Approximation

Figure 15 shows the equivalent circuit of the SET to the left of Fig. 13, based on the two-state approximation of the full-analytical SET model. Figure 16 shows the source code of SPICE subcircuits implemented to SmartSpice. The assumptions for deriving the full-analytical drain current $I_{d_1}$ in the SET equivalent circuit are that at each given gate voltage, only the two most-probable charging states of the Coulomb island are taken into account, at which the model is sufficiently accurate over a wide drain-source voltage range of $|V_{ds}| \leq e/C_\Sigma$ and even at a relatively high temperature of $0.1e^2/2C_\Sigma k_B$. The model is based on the "orthodox" theory and the steady-state master equation [16, 40]. Under such an assumption, the steady-state master equations having two states, $n$ and $n + 1$, becomes

$$\left[ \Gamma_d^{dec}(n + 1) + \Gamma_s^{dec}(n + 1) \right] P_{n+1} - \left[ \Gamma_d^{inc}(n) + \Gamma_s^{inc}(n) \right] P_n = 0 \qquad (42)$$

where $P_n$ and $P_{n+1}$ are the probabilities that the charge in the Coulomb island in an SET equals $ne$ and $(n + 1)e$, respectively, $\Gamma_d^{dec}(n + 1)$ and $\Gamma_s^{dec}(n + 1)$ are the tunneling rates at the drain junction and the source junction when the charge in the Coulomb island is changed from $(n + 1)e$ to $ne$, respectively, and $\Gamma_d^{inc}(n)$ and $\Gamma_s^{inc}(n)$ are the tunneling rates at the drain junction and the source junction when the charge in the Coulomb island is changed from $ne$ to $(n + 1)e$, respectively. By considering that $P_n + P_{n+1} = 1$, and by using the asymmetry factor $r = (R_d - R_s)/(R_d + R_s)$ and the hyperbolic sine function, drain current $I_n$ can be neatly rearranged and then can be expressed as [45, 48, 50].

$$I_n = \frac{e}{4C_\Sigma R_T} \frac{(1 - r^2)\left( \tilde{V}_{gs}^2 - \tilde{V}_{ds}^2 \right) \sinh\left( \frac{\tilde{V}_{ds}}{\tilde{T}} \right)}{A} \qquad (43)$$

$$A = \left\{ \tilde{V}_{gs} \sinh\left( \frac{\tilde{V}_{gs}}{\tilde{T}} \right) - \tilde{V}_{ds} \sinh\left( \frac{\tilde{V}_{ds}}{\tilde{T}} \right) \right\} + r \left\{ \tilde{V}_{ds} \sinh\left( \frac{\tilde{V}_{gs}}{\tilde{T}} \right) - \tilde{V}_{gs} \sinh\left( \frac{\tilde{V}_{ds}}{\tilde{T}} \right) \right\} \qquad (44)$$

$$\tilde{V}_{gs} = \frac{2C_g V_{gs}}{e} + \frac{2C_b V_{bs}}{e} - \frac{(C_g + C_b + C_s - C_d) V_{ds}}{e} + 2n - 1 \qquad (45)$$

$$\tilde{V}_{ds} = \frac{C_\Sigma V_{ds}}{e} \qquad (46)$$

$$\tilde{T} = \frac{k_B T}{e^2/2C_\Sigma} \qquad (47)$$

Here, $R_T$ is the harmonic mean of the tunneling resistances, $2R_d R_s/(R_d + R_s)$. It should be noted that the dependence of $I_n$ of $\tilde{V}_{gs}$ takes bell-shape characteristics with quasi-exponential decays on both sides and can be represented by only one peak. Because there is one dominant charging state and two almost-equally minor states that determine the leakage current, it is also inaccurate in the Coulomb blockade (CB) regions on both sides. The summation of $I_n$s for different $n$s in the relevant gate-voltage range gives the Coulomb oscillations of the drain current [45] and simultaneously compensates the inaccuracy caused by the insufficient number of charging states considered in the model. The SET model is implemented to SmartSpice as a subcircuit comprising analogue behavior devices [73].

Figure 17 shows the $I_d$-$V_{gs}$ characteristics of asymmetry SETs calculated by the model and the reference simulator. The lines are calculated according to the model and the symbols are calculated by the Monte Carlo simulator SIMON [37]. There is virtually no difference between the results of the model and the reference simulator. Figure 18 reproduces the Coulomb staircase, which is peculiar to asymmetric SETs, for different gate biases. The results calculated according to the model coincide well with the simulated ones, at least in the CB and SET regions, as can be expected under the two-charging-state assumption.

## 6.3.2. Multistate Approximation

Figure 19 shows the equivalent circuit of the SET to the left of Fig. 13, based on the multistate (11 states) approximation of the full-analytical SET model. Here, a stray capacitance $C_0$ to ground and a background charge $Q_0$ are included in the model. The white voltage sources representing the bias (the source and the drain) and gate voltages (the gate and the back gate) are external to the SET SPICE model, and the gray sources are internal to the model. E1 is a voltage source that defines the island voltage in the SET (node 5), and G1 is a current source that specifies the source-drain current from node 1 to node 2. This model is almost similar to the DC model of the seminumerical modeling of a SET if the number of states is fixed to a constant value such as 11 states in this model. Figure 20 shows the source code of SPICE subcircuits implemented to SmartSpice.

In this model, first the most probable charge state, which has the highest probability of the charge state $n$ that the charge in the Coulomb island of an SET equals $ne$, should be determined. The most probable charge state is expressed as [51]

$$n_{opt} = \frac{-(Q_0 + C_s V_s + C_d V_d + C_g V_g + C_b V_b)}{e} + \frac{C_\Sigma (V_s R_d + V_d R_s)}{e(R_s + R_d)} \tag{48}$$

where the total capacitance of the Coulomb island in the SET $C_\Sigma = C_d + C_g + C_s + C_b + C_0$, including $C_0$ in this model. And then the steady-state master equation for the 11 charge states around this most probable charge state is solved to calculate the values for E1 and G1. To determine all probabilities $P_n$s that the charge in the Coulomb island in an SET equals $ne$, the following recursion relation [16, 51] is used.

$$P_{n+1} = P_n \left( \frac{\Gamma_d^{inc}(n) + \Gamma_s^{inc}(n)}{\Gamma_d^{dec}(n+1) + \Gamma_s^{dec}(n+1)} \right) \tag{49}$$

To obtain the current source G1, the average current flowing through the SET in the direction from the source tunnel junction to the drain tunnel junction should be calculated, which is expressed as

$$I_{avg} = e \sum_n P_n (\Gamma_s^{inc} - \Gamma_s^{dec}) == e \sum_n P_n (\Gamma_d^{dec} - \Gamma_d^{inc}) \tag{50}$$

To obtain the voltage source E1, the average voltage in the Coulomb island should be calculated, which is expressed as

$$V_{avg} = \sum_n P_n V_n \tag{51}$$

where $V_n$ is the island voltage in the case when the charge in the Coulomb island of an SET equals $ne$.

Figure 21 compares the simulations of a single-electron transistor using this model (black lines) with the ones using the Monte Carlo simulator SIMON (gray lines) to demonstrate the accuracy of multistate (11 states) approximation of the full-analytical SET model. As it clearly can be seen, the simulations give identical results except for some wiggles in the Monte Carlo simulation lines, which are due to the stochastic character of the Monte Carlo algorithm. In Fig. 21a, the current–voltage $I_{ds}-V_{ds}$ characteristics at various gate voltages $V_{gs}$ are plotted. Figure 21b shows the current–voltage $I_{ds}-V_{gs}$ characteristics at various drain-source bias voltages $V_{ds}$. If the thermal fluctuations ($k_B T$) are larger than the energy it takes to add an electron to the island, the CB is washed out. This is demonstrated in Fig. 21c where the current–voltage $I_{ds}-V_{ds}$ characteristics at three different temperatures are plotted.

## 6.4. Quasi-Analytical Modeling

Quasi-analytical modeling has the similar assumption to two-state approximation of the full-analytical modeling of SET introduced in the previous section: (i) $|V_{ds}| \leq e/C_\Sigma$, and (ii) the interconnect capacitances associated with gate, source, and drain terminals are much larger than the device capacitance.

Figure 22 shows the simplified flow chart to calculate the SET drain currents. The drain current $I_{ds}$ in the SET is modeled via two components: (i) the tunneling current $I_{DST}$, which is independent of temperature, $T$, and (ii) the thermal current, $I_{DSTH}$, which fully includes the effect of temperature as

$$I_{ds} = I_{DST} + I_{DSTH} = \frac{I_D I_s}{I_D + I_s} + I_{DSTH} \tag{52}$$

where $I_{DST}$ is subsequently modeled based on the individual calculations of drain $I_D$ and source $I_s$ tunneling components, which is considered to be proportional to the tunneling rates [46]. As $I_{DST}$ is independent of temperature, tunneling rates are calculated at zero temperature. The analytical linear expressions of $I_D$ and $I_s$ is obtained from the island voltage $V_{island}$ as shown in Fig. 22. Before any electron tunneling has occurred, $V_{island}$ can be expressed as

$$V_{island} = \frac{C_s V_{ds} + C_g V_{gs} + C_b V_{bs} - ne}{C_\Sigma} \tag{53}$$

According to the "orthodox theory," when the voltage drop $V_{di}$ between the drain and the Coulomb island or the voltage drop $V_{is}$ between the Coulomb island and the source becomes larger than $V_\Sigma (= e/2C_\Sigma)$, one electron tunnels-in or tunnels-out from the source or the drain to the island and as a result $V_{island}$ decreases (for tunnel-in) or increases (for tunnel-out) by an amount of $2V_\Sigma$. However, if $V_{di}$ or $V_{is}$ becomes less than $V_s$, no electron tunneling happens and the device enters into the CB region. The two "while" statements in the quasi-analytical algorithm as shown in Fig. 22 are used to modify $V_{island}$ in order to capture the periodic Coulomb oscillation characteristics of SET. Based on this modified value of $V_{island}$, the tunneling current $(I_{DST})$ is modeled as shown in Fig. 22. The overall calculation of the drain current is performed by using the following equations, $I_{ds} = f(V_{gs}, V_{bs}, V_{ds})$, where $f$ is a computing subroutine.

As explained in Fig. 22, the current due to thermal effects is considered as a temperature dependent leakage current $I_{DSTH}$, which can be calculated using the following quasi-empirical equation:

$$I_{DSTH} = \begin{cases} (I_{peak} - I_{DST})\frac{2V_t}{T_\Sigma}, & \text{when } I_{DST} \neq 0 \\ I_{OFF} = \frac{V_t}{2(R_d + R_s)} \ln\left[1 + \exp(1 - \frac{V_s}{mV_t})\right], & \text{when } I_{DST} = 0 \end{cases} \tag{54}$$

where $V_t$ is thermal voltage ($= k_B T/e$), $m$ is fitting coefficient, and $I_{peak} = V_{ds}/2(R_d + R_s)$ [37]. The off-state current $I_{OFF}$, which is a function of temperature and the device-size, has been checked to agree with the results calculated by the Monte Carlo simulator SIMON over two decades of temperature (0.1 K to 10 K) for fixed value of $m = 10$ [47].

Figure 23 shows the drain current-gate voltage $I_{ds}$-$V_{gs}$ characteristics and the transconductance $\partial I_{ds}/\partial V_{gs}$-$V_{gs}$ characteristics of the single and double-gate SET at various $V_{ds}$ to demonstrate the accuracy of the quasi-analytical modeling. The filled circles and the open circles represent the results of the single and double-gate SET simulated by the quasi-analytical modeling, respectively. The solid lines and the dotted lines represent the results of the single and double-gate SET simulated by the SIMON, respectively. The quasi-analytical modeling is further validated in terms of SET transconductance $g_m$: the quasi-linear plot in Fig. 23b shows the ability of the quasi-analytic modeling to accurately describe the first derivative of the drain current.

# 7. COMPARISON BETWEEN EACH COMPACT-MODELING

Figure 24 shows the $I_{ds}$-$V_{ds}$ characteristics of SETs calculated by the macromodeling, the seminumerical modeling, and the full-analytical modeling in the DC analysis. Because the results simulated with the full analytical modeling are considerably the same to those with the quasi-analytical modeling [47], the quasi-analytical modeling will be not compared

Figure 35. Principle of programmable SET logic. (a) Schematic of the programmable SETs with a nonvolatile memory (NVM) node that is a key element of the programmable SET logic. The SET with NVM consists of a quantum dot (QD), tunnel junctions, and a memory node. Here. $C_{dg}$, $C_{mg}$, $C_{dm}$, $I_d$, and $V_g$ are capacitance between the QD and the gate electrode, the capacitance between the memory node and the gate electrode, the capacitance between the QD and the memory node, the drain current, and the gate voltage, respectively. (b) Characteristics of SET with NVM. Initially, the SET with NVM shows the same $I_d$–$V_g$ characteristics as those of the conventional SET (upper figure). The complementary SET (lower figure) is realized after writing operation generating the half-period ($\pi$) phase shift of Coulomb oscillations. For simplicity, $C_{dm}$ is not shown in the schematics of the SETs with NVM. (c) Logical meaning of complementary SET. The operation of the complementary SET is equivalent to that of conventional SET to which logically inverted signal is fed. Reprinted with permission from [19], K. Uchida et al., *IEEE Trans. Electron. Devices* 50, 1623 (2003). © 2003. IEEE.

with other simulation results. The solid lines, the open squares, the filled circles, and the open triangles are calculated by the Monte Carlo simulator SIMON, the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. The simulation results of the macromodeling and the seminumerical modeling match with the Monte Carlo simulation results. The simulation results of the full-analytical modeling match with the Monte Carlo simulation results at the drain-source voltages range of $|V_{ds}| \leq e/C_\Sigma$, but they have considerable differences with the Monte Carlo simulation results at $|V_{ds}| \geq e/C_\Sigma$.

**Figure 36.** Example of programmable SET logic. Schematics of (a) an SET inverters, (b) and (c) two-input SET NOR, and (d) two-input SET AND gates. Reprinted with permission from [19], K. Uchida et al., *IEEE Trans. Electron. Devices* 50, 1623 (2003). © 2003, IEEE.

Figure 25 shows the $I_{ds}-V_{gs}$ characteristics of an SET calculated by the macromodeling, the seminumerical modeling, and the full-analytical modeling in the DC analysis. The solid lines, the open squares, the filled circles, and the open triangles are calculated by the Monte Carlo simulator SIMON, the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. When the temperature $T = 30$ K, the simulation results of the seminumerical modeling and the full-analytical modeling matches with the Monte Carlo simulation results, but the simulation results of the macromodeling show considerable differences with the Monte Carlo simulation results. When the temperature $T = 100$ K, the simulation results of the macromodeling and the full-analytical modeling matches with the Monte Carlo simulation results, but the simulation results of the full-analytical modeling show small differences with the Monte Carlo simulation results. The values of the parameters in the macromodeling are that $CF1 = 12.5$, $CVp = 0.04$, $CI2 = 1.9 \times 10^{-9}$, $CP1 = 0.51$, $CR1 = 190 \times 10^6$, and $CR2 = 150 \times 10^6$ when $T = 30$ K and $CF1 = 12.5$, $CVp = 0.04$, $CI2 = 1.9 \times 10^{-9}$, $CP1 = 0.51$, $CR1 = 51 \times 10^6$, and $CR2 = 12 \times 10^6$ when $T = 100$ K. In Table 2, the differences between simulation results of four types of the compact SET modeling and the Monte Carlo simulation results are summarized in the dc analysis. The seminumerical modeling is the most excellent above other modeling to match with the Monte Carlo simulator over the entire range of the gate voltages, the drain voltages, and the temperature because the seminumerical modeling solves thoroughly the master equation.

Figure 26 shows the drain currents and source currents of SETs when gate-source voltage $V_{gs}$ varies linearly from 0 to 0.1 V in 1 $\mu$ and drain-source voltage $V_{ds} = 26.7$ mV. The solid lines, the squares, the circles, and the triangles are calculated by the Monte Carlo simulator SIMON, the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. The filled symbols and the open symbols are the drain currents $I_d$ and the source currents $-I_s$, respectively. In this slow transient simulation, the drain currents $I_d$ are equal to the source currents $-I_s$ as the DC simulation results in Fig. 25. The simulation results of the macromodeling, the seminumerical modeling, and the full-analytical modeling are as same as the DC simulation results in Fig. 25. Figure 27 shows the drain currents and source currents of SETs calculated by the macromodeling, the seminumerical modeling, and the full-analytical modeling when gate-source voltage $V_{gs}$ varies linearly from 0 to 0.1 V in 1 ns and drain-source voltage $V_{ds} = 26.7$ mV. In this fast transient simulation, the drain currents $I_d$ are different with the source currents $-I_s$. The simulation results of the macromodeling

Figure 37. Examples of multi-input SET logic gates. Schematics of (a) the multi-input NAND gate |75| and (b) the four-input XOR gate [44]. (a) Reprinted with permission from [75], M. Y. Jeng et al., *Jpn. J. Appl Phys.* 36, 6706 (1997). © 1997. JAACC. (b) Reprinted with permission from |44|, K. Uchida et al., *Jpn. J. Appl. Phys.* 38, 4027 (1999). © 1999, Institute of Pure and Applied Physics.

and the full-analytical modeling are different than those of the seminumerical modeling solving accurately the time-dependent master equation.

Figures 28 and 29 show the voltage transfer characteristics (VTC) of the SEI in Fig. 5 when $T = 30$ K and $T = 100$ K, respectively. The solid lines, the squares, the circles, and the triangles are calculated by the Monte Carlo simulator SIMON, the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. Figures 28 and 29 show the same results as shown in Table 2. Figures 30 and 31 show the transient response of the SEI when the input voltage $V_{in}$ varies linearly from 0 to 0.04 V in 0.1 ns and in the 100 ns, respectively. The filled symbols and open symbols are calculated when $C_t = 10$ aF

**Figure 38.** The simulated timing chart of the circuit in Fig. 37b. $V_{CMOS} = 0.8$ V, $V_{SET} = 50$ mV, $C_d = C_s = 0.06$ aF, $C_g = 0.1$ aF, $C_j = 10$ fF, $C_l = 1$ fF, $C_t = 50$ aF, $R_d = R_s = 500$ kΩ, and $T = 293$ K. Reprinted with permission from [44], K. Uchida et al., *Jpn. J. Appl. Phys.* 38, 4027 (1999). © 1999, Institute of Pure and Applied Physics.

and $C_l = 100$ aF, respectively. The squares, the circles, and the triangles are calculated by the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. In the slow transient simulation in Fig. 31, the transient simulation results of the macromodeling, the seminumerical modeling, and the full-analytical modeling are almost the same, but in the fast transient simulation in Fig. 30, the transient simulation results of the macromodeling, the seminumerical modeling, and the full-analytical modeling show a slight difference. Figure 32 shows the simulation time of several types of SET circuits in the transient simulation. The open squares, the filled circles, and the open triangles are calculated by the macromodeling, the seminumerical modeling, and the full-analytical modeling, respectively. The $x$-axis denotes the complexity of the circuits: 1 (single-electron inverter), 2 (a single-electron NOR gate), 3 (two single-electron NOR gates), 4 (three single-electron NOR gates). The simulation time of the full-analytical modeling is similar to that of the macromodeling, but the simulation time of the seminumerical modeling is longer than that of the macromodeling and the full-analytical modeling.



**Figure 39.** Examples of SE-FET hybrid circuit consisting of a CMOSFET amplifier and an SE-FET converter/inverter consisting of one SET (in Fig. 35) and one FET in series. Reprinted with permission from [19], K. Uchida et al., *IEEE Trans. Electron. Devices*, 50, 1623 (2003). © 2003, IEEE.

## 8. SET CIRCUITS AND SE-FET HYBRID CIRCUITS FOR LOGIC APPLICATIONS

In this section, various types of SET circuits and SE-FET hybrid circuits consisting of SET circuits and field effect transistor (FET) circuits for logic applications are introduced. As the first example of proposed SET circuits, an SET inverter (Fig. 5), two-input NOR, two-input XOR, and two-input OR gates as shown in Fig. 33 are proposed by Chen et al. [58] and Karafyllidis [74], which consist of three terminal SETs as shown in Fig. 10a. However, the two-input OR gate in Fig. 33c cannot be simulated by the compact models of an SET (implemented to SPICE) described in Section 6 because the circuit includes three terminal SETs as well as a tunnel junction. As the second example of proposed SET circuits, an SET inverter, two-input XOR, two-input NAND, and two input NOR gates as shown in Fig. 34 are proposed by Jeng et al. [75], which consist of four terminal SETs proposed by Tucker [76] as shown in Fig. 13. As the third example of proposed SET circuits, an SET inverter, two-input SET NOR, and two-input SET AND gates as shown in Fig. 16 are proposed by Uchida et al. [23], in which programmable SETs with a nonvolatile memory (NVM) node (Fig. 35 and Fig. 36) are used and are similar to the complementary SETs proposed by Tucker [76]. As the fourth example of proposed SET circuits, multi-input gates are proposed by Jeng et al. [75], Takahashi et al. [77], and Uchida et al. [44]. Figures 37a and 37b show the multi-input NAND gate and the four-input XOR gate, respectively. The functionality of the multi-input gate SET enables us to make multibit adders with a small number of transistors without any wire crossing [78, 79]. In Fig. 37b, The four-input XOR gate is similar to dynamic



Figure 40. Experimental room-temperature demonstration of programmable SET logic operation of the SET-pMOSFET circuit. (a) Waveform of the SET gate voltage. (b) Output waveform of the SET-pMOSFET circuit. The initial characteristics are indicated by the solid line. The characteristics after applying 8 V writing pulse are indicated by the dashed line. The initial waveform is synchronized with that of the SET gate voltage. On the other hand, the waveform after writing operation is logically inverted to that of the SET gate voltage, demonstrating that function of the SET-pMOSFET circuit can be programmed from a converter to an inverter by using the NVM function. (c) Output waveform of the CMOS inverter. The waveform is logically inverted to that of the SET-pMOSFET circuit. It should be noted that the small output voltage of the SET-pMOSFET circuit is amplified with the CMOS inverter. Reprinted with permission from [19], K. Uchida et al., IEEE Trans. Electron. Devices, 50, 1623 (2003). © 2003, IEEE.

**Figure 41.** Examples of multivalued logics and memories with SE-FET hybrid circuits consisting of one current source, one MOSFET, and one SET. Schematics of (a) the universal literal gate (ULG), (b) the quantizer, (c) 3-bit analog-digital converter (ADC), and (d) the static random-access memory (SRAM) with the SE-FET hybrid circuits. Reprinted with permission from [80], H. Inokawa et al., *IEEE Trans. Electron. Devices* 50, 462 (2003). © 2003, IEEE.



**Figure 42.** (a) Measured voltage transfer ($V_{in}$-$V_{out}$) characteristics of the ULG in Fig. 41a. The $V_{gg}$ is set at 1.08 V to attain a SET drain voltage of about 10 mV. The CC load is realized by a current-mode output of 4.5 nA from a semiconductor parameter analyzer with compliance (voltage limit) of 5 V. (b) Measured quantizer operation of the quantizer in Fig. 41b with $V_{gg}$ = 1.08 V and $I_0$ = 4.5 nA. The frequency of short pulses of CLK, $f_{CLK}$ = 2.5 Hz. Reprinted with permission from [80], H. Inokawa et al., *IEEE Trans. Electron. Devices* 50, 462 (2003). © 2003, IEEE.

**Figure 43.** Example of a multivalued logic with SE-FET hybrid circuits consisting of an SET-resistor inverter and a CMOS inverter [83].

logic family with two cascaded switches. The four-input XOR gate is amplified to the same voltage as the gate voltage swing of SETs with CMOS inverter in order to drive the next gate [31]. This figure shows also an example of a hybrid circuit consisting of SET circuits and FET circuits. Figure 38 shows the simulated timing chart of the four-input XOR gate.

The key advantages of SET are ultrasmall size, ultralow power consumption. and new functionalities such as the Coulomb oscillation and the Coulomb blockade [14-17]. However, there are particular obstacles for implementing SET logic gates. The maximum voltage gain of SET, which is defined as the ratio of the gate to the drain capacitance, is very small, usually less than one or slightly more than one [29, 30]. The SETs have low current driving capability, which degrades device performance, as it takes a long time to charge up the large interconnection capacitance connected to the output node of a device. One of the approaches to overcome these inherent disadvantages of SETs is to construct a hybrid circuit consisting of metal-oxide-semiconductor field-effect transistors (MOSFETs), which have a high gain. high output resistance, and high applicable voltages, and can thus supplement



**Figure 44.** Example of a multivalued logic with SE-FET hybrid circuits using a complementary self-biasing scheme [84].

**Figure 45.** Example of the quantizer with two SETs and one load capacitance. Reprinted with permission from [85], S. Mahapatra et al., *Electron. Lett.* 38, 443 (2002). © 2002, IEE.

the SET. As the first example of proposed SE-FET hybrid circuits, the hybrid circuit (as shown in Fig. 39) consisting of a CMOSFET amplifier and an SE-FET converter/inverter consisting of one SET (in Fig. 35) and one FET in series is proposed by Uchida et al. [23]. The programmable operation of the SE-FET converter/inverter and its amplication (due to the CMOSFET inverter) are experimentally demonstrated (as shown in Fig. 40). As the second example of proposed SE-FET hybrid circuits, a multivalued logics and memories with SE-FET hybrid circuits consisting of one current source, one MOSFET, and one SET (as shown in Fig. 41) are proposed by Takahashi et al. [80]. SETs are very suitable for multivalued logic because the discreteness of the electron charge in the Coulomb island can be directly related to multivalued operation. Figures 41a to 41d show the schematic of the universal literal gate (ULG), the quantizer, 3-bit analog-digital converter (ADC), and the static random-access memory (SRAM) with SE-FET hybrid circuits consisting of current sources, MOSFETs, and SETs, respectively. The operations of the multivalued logics with the SE-FET hybrid circuits are experimentally demonstrated by using the SET fabricated with the PADOX process [81, 82] (as shown in Fig. 42) [80]. Figure 42a shows the measured input-output characteristics of the ULG with the SE-FET hybrid circuits. Figure 42a shows a periodic binary output, which is a unique characteristics of the ULG. The output voltage increases and decreases periodically with the input voltage, reflecting the Coulomb oscillation characteristics of the SET. Figure 42b shows the measured quantizer operation of the



**Figure 46.** Examples of (a) a charged lock loop with the SE-FET hybrid circuits and (b) an SET amplifier with a current bias by a FET. Reprinted with permission from [51], G. Lientschnig et al., *Jpn. J. Appl. Phys.* 42, 6467 (2003). © 2003, Institute of Pure and Applied Physics.

**Figure 47.** (a) Schematic of a single-electron CCD on a silicon-on-insulator wafer. The Si wire-MOSFETs are defined by the fine poly-Si gates that cover a part of the T-shaped Si wire. In addition to the fine gates, the upper poly-Si gate is formed (after the formation of SiO₂ interlayer) to cover a large area (not shown here). The Si wire is surrounded by gate oxide (not shown). A single hole can be stored and transferred. (b) Demonstration of the manipulation of an elementary charge. Procedure for demonstrating CCD operation. (1) One hole is stored in MOSFET 2 and is sensed. (2) The hole is transferred to MOSFET 1. (3) The hole is sensed again. (4 The hole is transferred backwards. (c) Results for the manipulation of a single hole between the two MOSFETs. The two drain currents are plotted for the repeated procedures of the sensing and the single-hole transfer. Owing to the sensing, the numbers of stored holes $[n_h(1), n_h(2)]$ are alternately detected as (0,1) and (1,0). Reprinted with permission from [88], Y. Ono et al., *nature* 410, 560 (2003). © 2003. Nature Publishing Group.

quantizer in Fig. 42b. A trianglular wave was fed to $V_{in}$ and the gate of the transfer gate MOSFET was driven by short pulses of $CLK$. Differential voltage levels in $V_{in}$ were sampled by the transfer gate MOSFET, transferred to the storage node $V_{out}$, and quantized. $V_{out}$ was quantized to levels **a** ~ **f**. Another example of the ULG as shown in Fig. 43 is proposed by Chun et al. [83], which consists of a SET-resistor inverter and a CMOS inverter. As the third example of proposed SE-FET hybrid circuits, the multivalued logic with SE-FET hybrid circuits using a complementary self-biasing scheme (as shown in Fig. 44) are proposed by Song et al. [84], in which SETs have the four terminals proposed by Tucker [76]. The complementary self-biasing method enables the multivalued logic to operate well at higher temperature



**Figure 48.** Schematic diagram of a future possible SE-FET hybrid ULSI architecture [90].

Table 3. Features of various types of SET-FET hybrid circuit simulators.

| Reporter | Year | Method | Accuracy of DC simulation | Transient simulation | |
|---|---|---|---|---|---|
| | | | | Accuracy | Simulation speed |
| Yu [17] | 1999 | Macromodeling | Good except for low $T$ | Good except for fast transient analysis | Fast |
| Amakawa [22] | 1998 | SME Seminumerical modeling | Excellent | Excellent except for fast transient analysis | Fast |
| Kirihara [20] | 1999 | TME Seminumerical modeling | Excellent | Excellent | Slow |
| Yu [19] | 2002 | | | | |
| Uchida [24] | 2000 | Full-analytical modeling | Excellent except for high $T$ and high $V_{th}$ | Excellent except for fast transient analysis | Fast |
| Wang [27] | 2000 | | | | |
| Mahapatra [28] | 2002 | | | | |
| Inokawa [25] | 2003 | | | | |

and higher $V_{in}$ condition than the ULG in Fig. 41a [84]. As the fourth example of proposed SE-FET hybrid circuits, the quantizer as shown in Fig. 45 is proposed by Mahapatra et al. [85], which consists of two SETs and one load capacitance. The quantizer does not require any external sampling signal for sampling baseband signal, and the sampling rate can be controlled by varying the device capacitor. As the fifth example of proposed SE-FET hybrid circuits, a charged lock loop to solve the background charge problem [16, 86, 87] and an SET amplifier with a current bias by a FET as shown in Fig. 46 are proposed by Lientschnig et al. [51]. The charge-clocked loop in Fig. 46a uses feedback to keep the charge in the Coulomb island of an SET constant. In Fig. 46b, the SETs provide the charge sensitivity while the FETs provide the gain and the low output impedence. Also, an FET of the second stage is used to buffer the output of the SET, which increases the speed of the circuit. As the fifth example of proposed SE-FET hybrid circuits, a single-electron CCD as shown in Fig. 47 is proposed by Fujiwara et al. [88], which consists of two lower gates (which partially cover the T-shaped wire branch) constituting two small MOSFETs connected in series and one upper gate (which covers the entire region of the wire) on the Si nanowire fabricated by the electron-beam lithography. Figure 47b illustrates the procedure for demonstrating the single-hole transfer between the two wire-MOSFETs. In the initial state (1), one hole is stored only in MOSFET 2. To read out the number of holes in both the MOSFETs, we measure two sensing currents by setting the sense-gate voltage as high as 0.88 V. As only MOSFET 2 stores a hole, current 1 should be low and current 2 should be high (as shown in Fig. 47c). After the sensing, the sense-gate voltage is decreased to −1 V. This depletes electrons not only from the channel below the front gate but also from neighboring channels. Next, the front-gate voltages are controlled such that the hole potential is as illustrated in (2). Then the hole can be transferred from MOSFET 2 to MOSFET 1. After that, the transferred hole is sensed again (3) and transferred backwards (4), and so on. The results for the single-hole manipulation are shown in Fig. 47c. An SE-FET hybrid pump with a similar structure to the single-electron CCD in Fig. 47 is also proposed by Ono et al. [89], which consists of two gates on the Si nanowire fabricated by the PADOX method. Figure 48 shows the equivalent circuits of the SE-FET hybrid pump. The pump consists of one SET and two ultra-small MOSFETs. This place tiny MOSFETs close to a SET and opens their channels only when we want electrons to pass to (from) the SET. This makes it possible to constitute a Si-based pump.

Uchida et al. recently reported a future possible SE-FET hybrid ULSI architecture as shown in Fig. 48 [90]. In Fig. 48, SET circuits and SE-FET hybrid circuits are used in programmable and low-power circuit block, and CMOSFET circuits are used in high-speed and I/O circuit blocks.

## 9. CONCLUSIONS AND FUTURE WORKS

The SET with very low power consumption, ultrasmall size, and high functional characteristics is promising for future large-scale integrated circuits (LSIs) as information process increases and the size of chips is reduced. However, it has crucial obstacles for implementing

SET logic gates because of its low current drivability, low voltage gain, and low temperature operation. To overcome these problems, the hybrid circuits consisting of SET and MOSFET have been extensively investigated because MOSFET has a high gain. high output resistance, and high applicable voltages, thus it can supplement the SET. Such a hybrid circuit has already been demonstrated experimentally, and in order to evaluate the merits of this approach thoroughly, powerful simulation tools have been developed. The simulation tools are developed on the conventional circuit simulator SPICE. Therefore, the simulation of the SET circuits must satisfy the basic assumptions of the SPICE; the $I$–$V$ characteristics of the device are affected by neighboring transistors only through the charges of the terminal voltages of those transistors. In this chapter, first, considerably accurate simulation methods (the Monte Carlo and the master equation method) for SET circuits were introduced. The simulation methods give considerably exact results, but they consume much simulation time and cannot simulate SE-FET hybrid circuits. Second, the possibility of compact modeling in the single-electron circuit simulation was introduced to simulate SE-FET hybrid circuits and reduce the simulation time. In the possible conditions, various types of SET-CMOS hybrid circuit simulators for efficient circuit simulation have been developed, and their features are summarized in Table 3. Each simulator is based on four types of simulation methods; the macromodeling, the seminumerical modeling, the full-analytical modeling, and quasi-analytical modeling. In DC analysis, the macromodeling is accurate except for the regime of low temperature $T(< 0.1\ e^2/k_B C_\Sigma)$, and the full-analytica modeling and the quasi-analytical modeling are accurate except for the regimes of high temperature $T(> 0.1 e^2/k_B C_\Sigma)$ and high drain-source voltage $V_{ds}(> e/C_\Sigma)$; the seminumerical modeling is accurate in all regimes. In slow transient analysis $(t_r > 1000 R_t C_t$, where $R_t$ and $C_t$ are the tunneling resistance and the tunneling capacitance), the simulation results of the macromodeling and the full-analytical modeling and the quasi-analytical modeling match with those of the seminumerical modeling, but. in fast transient analysis $(t_r < 1000 R_t C_t)$, the macromodeling, the full-analytical modeling, and the quasi-analytical modeling are considerably different than the seminumerical modeling. The seminumerical modeling, however, consumes more simulation time than the macromodeling, the full-analytical modeling, or the quasi-analytical modeling in fast transient analysis. In the future, the modeling that can have the merits of the seminumerical modeling, accurate in all operation regions in DC analysis, and can be fast enough to simulate by using the macromodeling, the full-analytical modeling, or the quasi-analytical modeling in transient analysis, will be required. Finally, various types of SET circuits and SE-FET hybrid circuits were explained.

The SET is a device capable of measuring charge with a charge sensitivity of $3 \times 10^{-6}$ $e/\sqrt{Hz}$ [91], which may be used in applications from very sensitive charge meters and current standards [16, 92–95]. Especially, the radio-frequency SET (RF-SET) uses a tank circuit for impedance transformation and has been shown to operate at frequencies up to 100 MHz with a sensitivity of the order $10^{-5}$ $e/\sqrt{Hz}$ [96–99]. The RF-SET model in SPICE is required to design and simulate the RF circuits consisting of tack circuits, conventional devices, and an RF-SET [100]. However, an AC SET model, especially the RF-SET model, has not been developed yet. Therefore, in the future. the RF-SET model will be developed on SPICE.

## ACKNOWLEDGMENTS

## REFERENCES

1. International Technology Roadmap for Semiconductor (ITRS) 2003 Edition. Available at http://public.itrs.net/Files/2003ITRS/Home2003.htm.
2. R. H. Dennard, F. H. Gaensslen, II. Yu. V. L. Rideout, E. Bassous, and A. R. Leblanc. *IEEE J. Solid-State Circuits* 9, 256 (1974).
3. G. Baccarani, M. R. Wordeman, and R. H. Dennard. *IEEE Trans. Electron. Devices* 31, 452 (1984).
4. D. J. Frank, R. H. Dennard, E. Nowak, P. M. Solomon, Y. Taur, and H.-S. P. Wong. *Proc. IEEE* 88, 259 (2001).

5.  L. Chang, Y.-K. Choi, D. Ha, P. Ranade, S. Xiong, J. Bokor, C. Hu, and T.-J. King, *Proc. IEEE* 91, (2003).
6.  D. Goldhaver-Gordon, M. S. Montemerlo, J. C. Love, G. J. Opiteck, and J. C. Ellenborgen, *Proc. IEEE* 85, 521 (1997).
7.  H. Majima, H. Ishikuro, and T. Hiramoto, "Technical Digest IEDM," Washington, DC, USA, 2001, p. 379.
8.  M. A. Kastner, *Rev. Modern Phys.* 64, 849 (1992).
9.  T. Van Duzer and C. W. Turner, "Principle of Superconducting Circuits," Elsevier, New York, 1981.
10. R. H. Mathews, J. P. Sage, T. C. L. G. Sollner, S. D. Calawa, C.-L. Chen, L. J. Mahoney, P. A. Maki, and K. M. Molvar, *Proc. IEEE* 87, 596 (1999).
11. M. R. Stan, P. D. Franzon, S. C. Goldstein, J. C. Lach, and M. M. Ziegler, *Proc. IEEE* 91, 1940 (2003).
12. B. S. Doyle, S. Datta, M. Doczy, S. Hareland, B. Jin, J. Kavalieros, T. Linton, A. Murthy, R. Rios, and R. Chau, *IEEE Electron. Device Lett.* 24, 263 (2003).
13. J. Kedzierski, M. Ieong, E. Nowak, T. S. Kanarsky, Y. Zhang, R. Roy, D. Boyd, D. Fried, and H.-S. P. Wong, *IEEE Trans. Electron. Devices* 50, 952 (2003).
14. M. Dorojevets and P. Bunyk, *IEEE Trans. Appl. Supercond.* 13, 446 (2003).
15. D. V. Averin and K. K. Likharev, in "Mesoscopic Phenomena in solids," North-Holland, Elsevier Science B. V., Oxford, New York, Tokyo, 1991.
16. H. Grabert and H. Devoret, "Single Charge Tunneling-Coulomb Blockade Phenomena in Nanostructure," NATO ASI Series B: Physics Plenum, New York and London, 1992.
17. K. K. Likharev, *IBM J. Res. Dev.* 32, 144 (1988).
18. Y. Ono, Y. Takahashi, K. Yamazaki, M. Nagase, H. Namatsu, K. Kurihara, and K. Murase, *Appl. Phys. Lett.* 76, 3121 (2000).
19. K. Uchida, J. Koga, R. Ohba, and A. Toriumi, *IEEE Trans. Electron. Devices* 50, 1623 (2003).
20. D. H. Kim, S.-K. Sung, K. R. Kim, J. D. Lee, B.-G. Park, B. H. Choi, S. W. Hwang, and D. Ahn, *IEEE Trans. Electron. Devices* 49, 627 (2002).
21. L. Guo, E. Leobandung, and S. Chou, *Science* 275, 649 (1997).
22. K. Yano, T. Ishii, T. Sano, T. Mine, F. Murai, T. Hashimoto, T. Kobayashi, T. Kure, K. Seki, *Proc. IEEE* 86, 633 (1999).
23. A. Nakajima, T. Futatsugi, K. Kosemura, T. Fukano, and N. Yokoyama, *Appl. Phys. Lett.* 70, 1742 (1997).
24. A. Tilke, R. H. Blick, H. Lorenz, and K. P. Kotthaus, *J. Appl. Phys.* 89, 8159 (2001).
25. B. H. Choi, S. W. Hwang, I. G. Kim, H. C. Shin, Y. Kim, and E. K. Kim, *Appl. Phys. Lett.* 73, 3129 (1998).
26. S. Tiwari, F. Rana, H. Hanafi, A. Harstein, E. F. Crabbe, and K. Chan, *Appl. Phys. Lett.* 68, 1377 (1996).
27. N. Asahi, M. Akazawa, and Y. Amemiya, *IEEE Trans. Electron. Devices* 44, 1109 (1997).
28. H. Inokawa, A. Fujiwara, and Y. Takahashi, *IEEE Trans. Electron. Devices* 50, 462 (2003).
29. R. A. Smith and H. Ahmed, *Appl. Phy. Lett.* 71, 3838 (1997).
30. C. P. Heij and P. Hadley, *Rev. Sci. Instrum.* 73, 791 (2002).
31. A. Ohatha, A. Toriumi, and K. Uchida, *Jpn. J. Appl. Phys.* 36, 1686 (1997).
32. H. Inokawa, A. Fujiwara, and Y. Takahashi, *Appl. Phys. Lett.* 79, 3620 (2001).
33. Z. A. K. Durrani, A. C. Irvine, and H. Ahmed, *IEEE Trans. Electron. Devices* 47, 2334 (2000).
34. K. Nishiguchi, H. Inokawa, Y. Ono, A. Fujiwara, and Y. Takahashi, *IEEE Electron. Device Lett.* 25, 31 (2004).
35. M. Saitoh and T. Hiramoto, "Technical Digest IEDM," 2003, p. 753.
36. A. N. Korotov, R. H. Chen, and K. K. Likharev, *J. Appl. Phys.* 78, 2520 (1995).
37. C. Wasshuber, H. Kosina, and S. Selberherr, *IEEE Trans. Comp. Aided Design* 16, 937 (1997).
38. KOSEC (Korea Single Electron Circuit simulator) developed at the Nanoelectronics Laboratory in Korea University, Seoul, Korea.
39. L. R. C. Fonseca, A. N. Korotov, K. K. Likharev, and A. A. Odintsov, *J. Appl. Phys.* 78, 3238 (1995).
40. S. Amakawa, H. Majima, H. Fukui, M. Fujishima, and K. Hoh, *IEICE Trans. Electron.* E81-C, 21 (1998).
41. M. Fujishima, S. Amakawa, and K. Hoh, *Jpn. J. Appl. Phys.* 37, 1478 (1998).
42. Y. S. Yu, S. W. Hwang, and D. Ahn, *IEEE Trans. Electron. Devices* 46, 1667 (1999).
43. M. Kirihara, K. Nakazato, and M. Wagner, *Jpn. J. Appl. Phys.* 38, 2028 (1999).
44. K. Uchida, K. Matsuzawa, and A. Toriumi, *Jpn. J. Appl. Phys.* 38, 4027 (1999).
45. K. Uchida, K. Matsuzawa, and J. Koga, R. Ohba, S. Takagi, and A. Toriumi, *Jpn. J. Appl. Phys.* 39, 2321 (2000).
46. X. Wang and W. Porod, *Superlatt. Microstruct.* 28, 345 (2000).
47. S. Mahapatra, A. M. Ionescu, and K. Banerjee, *IEEE Electron. Device Lett.* 23, 366 (2002).
48. S.-H. Lee, D. H. Kim, K. R. Kim, J. D. Lee, B.-G. Park, Y.-J. Gu, G.-Y. Yang, and J.-T. Kong, *IEEE Trans. Nanotechnol.* 1, 226 (2002).
49. Y. S. Yu, S. W. Hwang, and D. Ahn, *Electron. Lett.* 38, 850 (2002).
50. H. Inokawa and Y. Takahashi, *IEEE Trans. Electron. Devices* 50, 455 (2003).
51. G. Lientshnig, I. Weymann, and P. Hadley, *Jpn. J. Appl. Phys.* 42, 6467 (2003); available at http://vortex.tn.tudelft.nl/research/set
52. R. M. Kielkowski, "Inside SPICE," McGraw-Hill, New York, 1994.
53. P. Antognetti and G. Massobrio, "Semiconductor device modeling with SPICE," McGraw-Hill, New York, 1988.
54. Available at http://infopad.eecs.berkeley.edu/~icdesign/SPICE.
55. K. K. Likharev, N. S. Bakhvalov, G. S. Kasacha, and S. I. Seryukova, *IEEE Trans. Magn.* 25, 1436 (1989).
56. W. J. Stewart, Introduction to the Numerical Solution of Markov Chains, Princeton University press, Princeton, NJ, 1994.
57. N. S. Bakhvalov, G. S. Kazacha, K. K. Likharev, and S. I. Serdyukova, *Sov. Phys. JETP* 68, 581 (1989).

58. R. H. Chen, A. N. Korotkov, and K. K. Likharev, *Appl. Phys. Lett.* 68, 1954 (1996).
59. S. A. Roy, Simulation Tools for the Analysis of Single Electronic Systems, PhD thesis, University of Glasgow, 1994.
60. M. Kirihara, N. Kuwamura, K. Taniguchi, and C. Hamaguchi, "Ext. Abst. Int. Conf. Solid State Devices Mater." Yokohama, Japan, 1994, p. 328.
61. G. E. Johnson, *Proc. IEEE* 82, 270 (1994).
62. L. J. Geerligs, V. F. Anderegg, P. A. M. Holweg, J. E. Mooij, H. Potier, D. Esteve, C. Urbina, and M. H. Devoret, *Phys. Rev. Lett.* 64, 2691 (1990).
63. H. D. Jensen and J. M. Martinis, *Phys. Rev. B* 46, 13407 (1992).
64. L. R. C. Fonseca, A. N. Korotov, and K. K. Likharev, *J. Appl. Phys.* 79, 9155 (1996).
65. D. V. Averin and Yu. V. Nazarov, *Phys. Rev. Lett.* 65, 2446 (1990).
66. D. V. Averin and A. A. Odintsov, *Phys. Lett. A* 140, 251 (1989).
67. L. I. Glazman and K. A. Matveev, *Sov. Phys. JETP* 71, 1031 (1990).
68. G. H. Goluv and C. F. Van Loan, "Matrix Computations," Johns Hopkins University Press, Baltimore, MD, and London, 1989, p. 557.
69. I. O. Kulik and R. I. Shekhter, *Sov. Phys. JETP* 41, 308 (1975).
70. H. S. Lee, Y. S. Yu, and S. W. Hwang, *J. Korean Phys. Soc.* 33, s266 (1998).
71. *STAR-HSPICE Manual Release* 2001.2, Avanti (2001).
72. *SmartSpice User Manual Vol 2*, Silvaco International (1997).
73. H. Inokawa and Y. Takahashi, "Proceedings 33rd IEEE International Symposium Multiple-Valued Logic," Tokyo, Japan, 2003, p. 259.
74. I. Karafyllidis, *Electron. Lett.* 36, 407 (2000).
75. M. Y. Jeng, Y. H. Jeng, S. W. Hwang, and D. M. Kim, *Jpn. J. Appl. Phys.* 36, 6706 (1997).
76. J. R. Tucker, *J. Appl. Phys.* 72, 4399 (1992).
77. Y. Takahashi, A. Fujiwara, H. Namatsu, K. Kurihara, *Appl. Phys. Lett.* 76, 637 (2000).
78. Y. Ono, K. Y. Amazaki, M. Nagase, S. Horiguchi, K. Shiraishi, and Y. Takahashi, *Microelectron. Eng.* 59, 435 (2001).
79. Y. Ono, K. H. Inokawa, and Y. Takahashi, *IEEE Trans. Nanotechnol.* 1, 93 (2002).
80. H. Inokawa, A. Fujiwara, and Y. Takahashi, *IEEE Trans. Electron. Devices* 50, 462 (2003).
81. Y. Takahashi, H. Namatsu, K. Kurihara, K. Iwadate, M. Nagase, and K. Murase, *IEEE Trans. Electron. Devices* 43, 1213 (1996).
82. S. Horiguchi, *Jpn. J. Appl. Phys.* 40, L29 (2001).
83. M.-J. Chun and Y.-H. Jung, "Proceedings of 3rd IEEE-NANO," San Francisco, CA, 2003, p. 745.
84. K.-W. Song, S. H. Lee, D. H. Kim, K. R. kim, J. Kyung, G. Baek, C.-A. Lee, J. D. Lee, and B.-G. Park, in "Proceedings 33rd IEEE International Symposium Multiple-Valued Logic," San Francisco, CA, 2003, p. 267.
85. S. Mahapatra, A. M. Inoscu, K. Banerjee, and M. Declercq, *Electron. Lett.* 38, 443 (2002).
86. I. I. Abramov and E. G. Novick, *Semiconductor* 35, 489 (2001).
87. H.-O. Muller, M. Furlan, T. Heinzel, and K. Ensslin, *Europhys. lett.* 55, 253 (2001).
88. A. Fujiwara and Y. Takahashi, *Nature* 410, 560 (2001).
89. Y. Ono and Y. Takahashi, *Appl. Phys. Lett.* 82, 1221 (2003).
90. K. Uchida, J. Koga, and A. Toriumi, "Technical Digest ISSCC," 2002, p. 206.
91. A. Aassime, D. Gunnarsson, K. Bladh, P. Delsing, and R. Schoelkopf, *Appl. Phys. Lett.* 79, 4031 (2001).
92. H. Pothier, P. Lafarge, P. F. Orfila, C. Urbina, D. Esteve, and M. H. Devoret, *Physica B* 169, 57 (1991).
93. M. W. Keller, J. M. Martinis, N. M. Zimmerman, and A. H. Steinbach, *Appl. Phys. Lett.* 69, 1804 (1996).
94. Y. Ono, N. M. Zimmerman, K. Yamazaki and Y. Takahashi, *Jpn. J. Appl. Phys.* 42, L1109 (2003).
95. S. Amakawa, H. Mizuta, and K. Nakazato, *J. Appl. Phys.* 89, 5001 (2001).
96. R. J. Schoelkopf, P. Wahlgren, A. A. Kozhevnikov, P. Delsing, and D. E. Prober, *Science* 280, 1238 (1998).
97. T. M. Buehler, D. J. Reilly, R. P. Starrett, S. Kenyon, A. R. Hamilton, A. S. Dzurak, and R. G. Clark, *Microelectron. Eng.* 67, 775 (2003).
98. T. M. Buehler, D. J. Reilly, R. Brenner, A. R. Hamilton, A. S. Dzurak, and R. G. Clark, *Appl. Phys. Lett.* 82, 557 (2003).
99. T. Fujisawa and Y. Hirayama, *Appl. Phys. Lett.* 77, 543 (2000).
100. V. O. Turin and A. N. Korotkov, *Appl. Phys. Lett.* 83, 2898 (2003).

# CHAPTER 6

# Electric Properties of Nanostructures

## K. Palotás, B. Lazarovits, P. Weinberger

*Center for Computational Materials Science,*
*Technical University of Vienna, Vienna, Austria*

## L. Szunyogh

*Center for Computational Materials Science, Technical University of Vienna,*
*Vienna, Austria; and Department of Theoretical Physics, Center for Applied*
*Mathematics and Computational Physics Budapest University of Technology*
*and Economics, Budapest, Hungary*

## CONTENTS

# 1. INTRODUCTION

This contribution is devoted to the theoretical description of the electric properties of nano-structured matter, in particular to structures nanoscaled in two dimensions, namely supported clusters of atoms such as finite chains of atoms embedded in the surface of a metallic substrate or atomic-sized contacts. Because this description is based on a "real space" representation of the so-called Kubo-Greenwood equation, it was felt necessary to give first a proper account of the theoretical background of linear response theory in terms of electric fields. For this reason Section 2 deals quite generally with currently available transport theories. In putting the Kubo-Greenwood equation into a computationally accessible scheme the use of density functional theory and multiple scattering approaches is required. Therefore only after having summarized very shortly the main quantities in a Korringa-Kohn-Rostoker-type realization of multiple scattering (Section 3), practical expressions for evaluating electric properties of nanostructures are introduced (Section 4). Clearly enough the numerical accuracy of such approaches have to be documented before any kind of application to nanosized matter can be given. Unfortunately this kind of numerical "test" leads back to bulk materials, for which the electric properties are well documented, both experimentally and theoretically. However, only the "tests" discussed in Section 5 provide the necessary confidence for the theoretical results presented in Sections 6 and 7 for finite wires and atomic-sized contacts.

Not dealt with in this contributions are systems nanosized only in one dimension such as spin valves or other heterojunctions, as a review of such systems—also based on a Green's function realization of the Kubo-Greenwood equation—only appeared rather recently [1] that discusses in quite some length, for example, properties of the giant magnetoresistance (GMR).

# 2. TRANSPORT THEORIES

In this section, methods describing electric transport in solid matter are reviewed, with emphasis, however, on the Kubo-Greenwood approach. Consider a system of $N$ interacting electrons moving in the electrostatic potential of the nuclei, the effective one-electron (Kohn-Sham-) Hamilton operator is then given by

$$\widehat{H}_0 = -\frac{\hbar^2}{2m} \sum_{i=1}^{N} \mathbf{\nabla}_i^2 + \sum_{i=1}^{N} u_{\text{eff}}(\mathbf{r}; \sigma, \mathcal{M}) \tag{1}$$

where the first term is the kinetic energy operator and the second the effective one-electron potential that in turn depends on the "spin" of the electrons ($\sigma$) as well as on the magnetic configuration of the system ($\mathcal{M}$). As it is well-known the corresponding one-electron Kohn-Sham equation can be written, e.g., in the case of a three-dimensional periodic system as

$$\widehat{H}_0 \phi_{k,\sigma}(\mathbf{r}) = E_{k,\sigma} \phi_{k,\sigma}(\mathbf{r}), \quad k = (\nu, \mathbf{k}) \tag{2}$$

where $\nu$ refers the so-called band index, $\mathbf{k}$ to the momentum, and $E_{k,\sigma}$ and $\phi_{k,\sigma}(\mathbf{r})$ denote the one-electron energies and states, respectively.

## 2.1. Boltzmann Formalism

This kind of theoretical approach assumes the existence of a distribution function $f_{k,\sigma}(\mathbf{r})$ that measures the probability of charge carriers with spin $\sigma$ in state $k$ in the neighborhood of $\mathbf{r}$. The change of $f_k(\mathbf{r})$.

$$f_k(\mathbf{r}) = \sum_\sigma f_{k,\sigma}(\mathbf{r})$$

in time is then described by the famous *Boltzmann equation*.

$$\left(\frac{\partial f_k(\mathbf{r}, t)}{\partial t}\right) + \left(\frac{\partial f_k(\mathbf{r})}{\partial t}\right)_{\text{diffusion}} + \left(\frac{\partial f_k(\mathbf{r})}{\partial t}\right)_{\text{field}} = -\left(\frac{\partial f_k(\mathbf{r})}{\partial t}\right)_{\text{scattering}} \tag{3}$$

in which the various terms correspond to different effects, namely from the left to right: an explicit time dependence, diffusion, the influence of external fields and scattering. Stationarity implies now that the total time dependence of $f_k(\mathbf{r})$ vanishes, see Eq. (3).

For matters of simplicity in the following the $\mathbf{r}$-dependence of the distribution function will be neglected. The local change of electrons resulting from elastic scattering of independent particles can be correlated to the *microscopic scattering probability*,

$$P_{k\sigma, k'\sigma'} = \begin{pmatrix} P_{k\uparrow, k'\uparrow} & P_{k\uparrow, k'\downarrow} \\ P_{k\downarrow, k'\uparrow} & P_{k\downarrow, k'\downarrow} \end{pmatrix} \tag{4}$$

in the following manner

$$\left(\frac{\partial f_k}{\partial t}\right)_{\text{scattering}} = \sum_\sigma \sum_{k'\sigma'} [f_{k',\sigma'}(1 - f_{k,\sigma}) P_{k'\sigma', k\sigma} - (1 - f_{k',\sigma'}) f_{k,\sigma} P_{k\sigma, k'\sigma'}] \tag{5}$$

The first contribution is usually called *scattering-in term* and describes the scattering of electrons from occupied states $(k', \sigma')$ into an empty state $(k, \sigma)$; the second term refers to the reverse process, namely the scattering of an electron from an occupied state $(k, \sigma)$ into empty states $(k', \sigma')$ and is called *scattering-out term*. It therefore seems reasonable to separate the distribution function into two parts,

$$f_{k,\sigma} = f^0_{k,\sigma} + g_{k,\sigma} \tag{6}$$

where

$$f^0_{k,\sigma} = \frac{1}{e^{\beta(E_{k,\sigma} - E_F)} + 1} \tag{7}$$

is the Fermi-Dirac distribution function with $E_{k,\sigma}$ denoting the one-electron energies, see Eq. (2), $E_F$ the Fermi energy, and $\beta = 1/k_B T$ with $k_B$ being the Boltzmann constant and $T$ the temperature. In Eq. (6) $g_{k,\sigma}$ denotes the deviation from the equilibrium distribution function. Making use of the *principle of microscopic reversibility*,

$$P_{k\sigma, k'\sigma'} = P_{k'\sigma', k\sigma}$$

for the microscopic scattering probabilities and the separation of $f_{k,\sigma}$ in Eq. (6), the scattering term in Eq. (5) can be rewritten as

$$\left(\frac{\partial f_k}{\partial t}\right)_{\text{scattering}} = \sum_\sigma \sum_{k'\sigma'} P_{k\sigma, k'\sigma'}(g_{k',\sigma'} - g_{k,\sigma}) \tag{8}$$

Neglecting now in Eq. (3) terms with an explicit time dependence of the distribution function and changes caused by diffusion, that is, keeping only changes of $f_k$ arising from a homogeneous external electric field $\mathbf{E}$, the following expression is obtained,

$$e \sum_\sigma \left(\frac{\partial f^0_{k,\sigma}}{\partial E_{k,\sigma}}\right) \mathbf{v}_{k,\sigma} \mathbf{E} = \sum_{\sigma'} \sum_{k'\sigma'} P_{k\sigma, k'\sigma'}(g_{k',\sigma'} - g_{k,\sigma}) \tag{9}$$

where $\mathbf{v}_{k,\sigma}$ is the velocity of the electrons with spin $\sigma$, which in turn can be related semi-classically to the one-electron energies as follows

$$\mathbf{v}_{k,\sigma} = \frac{1}{\hbar}\frac{\partial E_{k,\sigma}}{\partial \mathbf{k}} \qquad (10)$$

Assuming that $g_{k,\sigma}$ depends linearly on the external electric field, the following *ansatz* can be made,

$$g_{k,\sigma} = -e\left(\frac{\partial f^0_{k,\sigma}}{\partial E_{k,\sigma}}\right)\Lambda_{k,\sigma}\mathbf{E} \qquad (11)$$

where $\Lambda_{k,\sigma}$ is the so-called *mean free path vector* of electrons of spin $\sigma$. The magnitude of $\Lambda_{k,\sigma}$ measures the path of an electron with spin $\sigma$ between two scattering events.

By introducing a so-called *relaxation time* $\tau_{k,\sigma}$, which specifies the time that an electron stays in state $(k, \sigma)$ until the next scattering event (*scattering life time*) occurs as

$$\tau_{k,\sigma}^{-1} = \sum_{k'\sigma'} P_{k\sigma,k'\sigma'} \qquad (12)$$

Eq. (9) can be solved with the ansatz in Eq. (11) to give

$$\Lambda_{k,\sigma} = \tau_{k,\sigma}\left(\mathbf{v}_{k,\sigma} + \sum_{k'\sigma'} P_{k\sigma,k'\sigma'}\Lambda_{k',\sigma'}\right) \qquad (13)$$

This now is a system of coupled integral equations. The different spin-components can be decoupled by ignoring *spin-flip scattering processes*, namely assuming in Eq. (4) that

$$P_{k\uparrow,k'\downarrow} = 0, \qquad P_{k\downarrow,k'\uparrow} = 0$$

such that a relatively simple integral equation is obtained,

$$\Lambda_{k,\sigma} = \tau_{k,\sigma}\left(\mathbf{v}_{k,\sigma} + \sum_{k'} P_{k\sigma,k'\sigma}\Lambda_{k',\sigma}\right) \qquad (14)$$

from which in principle $\Lambda_{k,\sigma}$ can be evaluated.

Due to the neglect of *spin-flip scattering processes*, the total current density can be written as

$$\mathbf{j} = \sum_{\sigma}\mathbf{j}_\sigma = \frac{e}{V}\sum_{k,\sigma}f_{k,\sigma}\mathbf{v}_{k,\sigma} \qquad (15)$$

where $V$ is the volume of the system. The *conductivity tensor* $\underline{\underline{\sigma}}$ at $T = 0$ is then obtained by using *Ohm's law*, $\mathbf{j}_\sigma = \underline{\underline{\sigma}}_\sigma \mathbf{E}$, and Eqs. (6), (11), (15),

$$\underline{\underline{\sigma}} = \sum_{\sigma}\underline{\underline{\sigma}}_\sigma = \frac{e^2}{V}\sum_{k,\sigma}\delta(E_{k,\sigma} - E_F)\mathbf{v}_{k,\sigma}\circ\Lambda_{k,\sigma} \qquad (16)$$

where $\circ$ denotes a dyadic product (resulting in a $3 \times 3$ tensor). The contributions to the total conductivity refer therefore to independent *majority* ($\uparrow$) and *minority* ($\downarrow$) spin channels (*two current model* [2]).

Neglecting the scattering-in term in Eq. (13), the conductivity is finally given by

$$\underline{\underline{\sigma}} = \sum_{\sigma}\underline{\underline{\sigma}}_\sigma = \frac{e^2}{V}\sum_{k,\sigma}\delta(E_{k,\sigma} - E_F)\tau_{k,\sigma}\mathbf{v}_{k,\sigma}\circ\mathbf{v}_{k,\sigma} \qquad (17)$$

where

$$\sum_{k,\sigma}\delta(E_{k,\sigma} - E_F) = n(E_F) \qquad (18)$$

is nothing but the density of states at the Fermi energy. Obviously the conductivity tensor is determined by three factors: the density of states, the velocities, and the relaxation times of the electrons at the Fermi surface. While the first two factors arise from the electronic structure of the system, the last one refers to defects or impurities present in the solid. Moreover, different approximations can be made for the relaxation time in Eq. (17), for example, an isotropic $\tau$, or spin-dependent $\tau_\sigma$, thus resulting in a simple expression for the conductance.

It has to be mentioned that the Boltzmann equation can easily be implemented within traditional bulk bandstructure methods, since in the semi-classical interpretation the velocity is given by the energy dispersion, see Eq. (10). Apart from being a semi-classical theory, the main disadvantages is that in the form of Eq. (17) only ordered bulk systems (three-dimensional cyclic boundary conditions) can be described, as a welldefined Fermi surface is needed and the relaxation times are system-dependent parameters.

## 2.2. Landauer Formalism

The *Landauer-Büttiker approach* [3, 4] is an effective tool to describe transport in mesoscopic systems. Suppose a multiprobe structure consists of a finite region connected to $N_L$ leads, each lead being attached to an ideal "reservoir." The electrons are then scattered in a finite region (*scattering or interaction region*) caused either by disorder or due to a particular geometry. The transport through the scattering region is thought to be completely coherent, no phase breaking is taken into account, and because of assumed low temperatures inelastic scattering processes are supposed to be negligible. The leads are used to inject and drain current or measure voltage, whereas the reservoirs are assumed to fulfil certain conditions: the reservoir for the $n$th lead has to be in equilibrium at a given chemical potential $\mu_n$,

$$\mu_n = E_F + eV_n \tag{19}$$

where $V_n$ is the applied potential and $E_F$ the Fermi energy. Furthermore, a steady-state current flowing from/into the reservoir is supposed not to change $\mu_n$, implying of course large enough reservoirs. Moreover, it is assumed that no additional resistance is produced by the interface between a reservoir and the scattering region. This in turn implies that an electron that enters a reservoir must be scattered inelastically before returning to the coherent scattering region, providing thus a phase-randomization. The current passsing through the $n$th lead can be written as

$$I_n = \sum_{m \neq n} g_{nm} V_m \tag{20}$$

where the sum extends over all leads except the $n$th one, and the $g_{nm}$ are the so-called conductance coefficients of the system.

Introducing incoming and outgoing scattering channels which play the same role as incoming and outgoing Bloch states in scattering theory, the conductance can be expressed in terms of a transmission probability ($T_{ni,mj}$) or $S$-matrix ($S_{ni,mj}$) as

$$g_{nm} = \frac{2e^2}{h} \sum_{ij} T_{ni,mj} = \frac{2e^2}{h} \sum_{ij} |S_{ni,mj}|^2 \tag{21}$$

In this equation $T_{ni,mj}$ is the transmission probability for an electron from an incoming channel $j$ in lead $m$ to an outgoing channel $i$ in lead $n$, the factor 2 accounts for the two spin directions, and the sum has to be carried out over all incoming and outgoing channels in the corresponding leads.

The advantage of using the Landauer formalism is first and foremost seen for two-probe structures, for which only one conductance coefficient $g$ has to be considered such as, for example, in the case of perpendicular transport (current perpendicular to plane-CPP) in layered structures or for quantum point-contacts. The main parameters of a contact refer to the characteristic lengths of the system, namely the *contact diameter* ($d$) and the *mean free path for elastic* ($\lambda_e$) and *inelastic* ($\lambda_i$) *scattering*, that is, the length of an electron's

path between two elastic (inelastic) scattering events. If $d \ll \Lambda_e$, $\Lambda_i$ one speaks of a *ballistic* point-contact, as an electron travels through the contact without any scattering. If $d \gg \Lambda_e$ a point-contact is said to belong to the *diffusive* regime meaning that an electron experiences a lot of elastic scattering when traveling through the contact. In both cases the contact diameter must be much larger than the electron's wavelength.

In the case of *two-probe structures* the conductance can be written according to Eq. (21) as

$$g = \frac{2e^2}{h} \sum_{ij} T_{ij} \tag{22}$$

where $T_{ij}$ is a hermitian matrix. Diagonalized the conductance can be formulated in the eigenchannel basis as

$$g = \frac{2e^2}{h} \sum_{i=1}^{N} T_i \tag{23}$$

where $N$ is the number of conducting channels and the $T_i$ are the real eigenvalues of $T_{ij}$, $0 < T_i < 1$. For an *ideal ballistic point-contact* as well as for the theoretically interesting case of an infinite periodic wire, $T_{ij} = \delta_{ij}$ which implies that the conductance is quantized in units of the so-called conductance quantum $G_0 = 2e^2/h$,

$$g = N_{ch} G_0 \tag{24}$$

Such quantized conductances have been observed by many experimental groups. Within the Landauer approach, the conductance obviously depends on the *number of open eigenchannels* $N_{ch}$, which in turn depends on the sample geometry. This implies that $N_{ch}$ is determined for the entire system by the *narrowest cross section* of a point-contact or a wire.

## 2.3. Kubo Formalism

In the 1950s, Kubo developed a method of evaluating the response of a quantum mechanical system to an external potential, in particular, the current in response to an electric field [5]. To first order, known as linear response, the two quantities are related by a conductivity (Ohm's law), which is given in terms of the equilibrium properties of the system, that is, in the limit of a vanishing field. Moreover, conductance coefficients can be derived from the conductivity, which describes the total current flowing in and out of the system in response to the voltages applied.

### 2.3.1. Linear Response Theory

**2.3.1.1. Linear Response and the Green Function**  Assuming a time-dependent perturbation $\hat{H}'(t)$, the Hamilton operator of the perturbed system is of the form,

$$\hat{H}(t) = \hat{H}_0 + \hat{H}'(t) \tag{25}$$

For a grand-canonical ensemble the density operator of the unperturbed system can be written as

$$\hat{\varrho}_0 = \frac{1}{Z} e^{-\beta \hat{n}_0} \tag{26}$$

with

$$\hat{n}_0 = \hat{H}_0 - \mu \hat{N} \tag{27}$$

where $\mu$ is the chemical potential, $\hat{N}$ the (particle) number operator, and $Z$ is the grand canonical partition function,

$$Z = Tr(e^{-\beta \hat{n}_0}) \tag{28}$$

Because the expectation value of a physical observable $A$, associated with a hermitian operator $\hat{A}$ corresponding to the unperturbed system is given by

$$A_0 = \langle A \rangle = \frac{1}{Z} Tr(\hat{A} e^{-\beta \hat{n}_0}) = Tr(\hat{\varrho}_0 \hat{A}) \tag{29}$$

within the Schrödinger picture the equation of motion for the density operator can be written as

$$ih\frac{\partial\hat{\varrho}(t)}{\partial t} = [\widehat{\mathscr{H}}(t), \hat{\varrho}(t)] \tag{30}$$

where

$$\widehat{\mathscr{H}}(t) = \hat{H}(t) - \mu\hat{N} = \widehat{\mathscr{H}}_0 + \hat{H}'(t) \tag{31}$$

Clearly enough, in the absence of a perturbation, $\hat{\varrho}(t) = \hat{\varrho}_0$. Therefore, partitioning $\hat{\varrho}(t)$ as

$$\hat{\varrho}(t) = \hat{\varrho}_0 + \hat{\varrho}'(t) \tag{32}$$

and making use of the fact that $[\widehat{\mathscr{H}}_0, \hat{\varrho}_0] = 0$, one gets in first order in $\hat{H}'$,

$$ih\frac{\partial\hat{\varrho}'(t)}{\partial t} = [\widehat{\mathscr{H}}_0, \hat{\varrho}'(t)] + [\hat{H}'(t), \hat{\varrho}_0] \tag{33}$$

or, by switching to the interaction (Dirac) picture,

$$\hat{\varrho}_D(t) = \hat{\varrho}_0 + \hat{\varrho}'_D(t), \quad \hat{\varrho}'_D(t) = e^{(i/h)\widehat{\mathscr{H}}_0 t}\hat{\varrho}'(t)e^{-(i/h)\widehat{\mathscr{H}}_0 t} \tag{34}$$

$$ih\frac{\partial\hat{\varrho}'_D(t)}{\partial t} = [\hat{H}'_D(t), \hat{\varrho}_0] \tag{35}$$

This equation has to be solved now for a given initial condition. Turning on the external field adiabatically at $t = -\infty$, implies that the density operator of the system at $t = -\infty$ represents an ensemble of systems in thermal equilibrium, that is,

$$\lim_{t\to-\infty} \hat{\varrho}(t) = \hat{\varrho}_0, \quad \lim_{t\to-\infty} \hat{\varrho}'_D(t) = 0$$

Using this boundary condition for $\hat{\varrho}'_D(t)$ results into the following integral equation

$$\hat{\varrho}'_D(t) = -\frac{i}{h}\int_{-\infty}^{t} dt'[\hat{H}'_D(t'), \hat{\varrho}_0] \tag{36}$$

such that in the Schrödinger picture the density operator can be approximated in first order as

$$\hat{\varrho}(t) \approx \hat{\varrho}_0 - \frac{i}{h}\int_{-\infty}^{t} dt'\, e^{-(i/h)\widehat{\mathscr{H}}_0 t}[\hat{H}'_D(t'), \hat{\varrho}_0]e^{(i/h)\widehat{\mathscr{H}}_0 t} \tag{37}$$

Considering now the time evolution of the physical observable $A(t)$,

$$A(t) = Tr(\hat{\varrho}(t)\hat{A}) = A_0 - \frac{i}{h}\int_{-\infty}^{t} dt'\, Tr(e^{-(i/h)\widehat{\mathscr{H}}_0 t}[\hat{H}'_D(t'), \hat{\varrho}_0]e^{(i/h)\widehat{\mathscr{H}}_0 t}\hat{A})$$

$$= A_0 - \frac{i}{h}\int_{-\infty}^{t} dt'\, Tr([\hat{H}'_D(t'), \hat{\varrho}_0]\hat{A}_D(t)) \tag{38}$$

where $A_0$ is defined in Eq. (29) and the Dirac representation of operator $\hat{A}$ is given by

$$\hat{A}_D(t) = e^{(i/h)\widehat{\mathscr{H}}_0 t}\hat{A}e^{-(i/h)\widehat{\mathscr{H}}_0 t} \tag{39}$$

then by making use of the identity,

$$Tr([\hat{A}, \hat{B}]\hat{C}) = Tr(\hat{A}\hat{B}\hat{C} - \hat{B}\hat{A}\hat{C}) = Tr(\hat{B}\hat{C}\hat{A} - \hat{B}\hat{A}\hat{C}) = Tr(\hat{B}[\hat{C}, \hat{A}])$$

one arrives at

$$\delta A(t) = A(t) - A_0 = -\frac{i}{h}\int_{-\infty}^{t} dt'\, Tr(\hat{\varrho}_0[\hat{A}_D(t), \hat{H}'_D(t')]) \tag{40}$$

Assuming finally that the perturbation $\hat{H}'(t)$ is of the form,

$$\hat{H}'(t) = \hat{B}F(t) \tag{41}$$

where $\hat{B}$ is a *Hermitian operator* and $F(t)$ is a *complex function* (classical field), Eq. (40) transforms to

$$\delta A(t) = -\frac{i}{\hbar} \int_{-\infty}^{t} dt'\, F(t')Tr(\hat{\varrho}_0[\hat{A}_D(t), \hat{B}_D(t')]) \tag{42}$$

which can be written in terms of a *retarded Green function* as,

$$G_{AB}^{ret}(t, t') = -i\Theta(t - t')Tr(\hat{\varrho}_0[\hat{A}_D(t), \hat{B}_D(t')]) \tag{43}$$

or, by introducing a so-called *generalized susceptibility*,

$$\chi_{AB}(t, t') = \frac{1}{\hbar}G_{AB}^{ret}(t, t') \tag{44}$$

as

$$\delta A(t) = \frac{1}{\hbar} \int_{-\infty}^{\infty} dt'\, F(t')G_{AB}^{ret}(t, t') = \int_{-\infty}^{\infty} dt'\, F(t')\chi_{AB}(t, t') \tag{45}$$

Suppose now that the operators $\hat{A}$ and $\hat{B}$ do not depend explicitly on time, then $G_{AB}^{ret}(t, t')$ and $\chi_{AB}(t, t')$ are only functions of the argument $(t - t')$. Consequently, the Fourier coefficients of $\delta A(t)$ can be written as

$$\delta A(\omega) = \frac{1}{\hbar} F(\omega)G_{AB}^{ret}(\omega) = F(\omega)\chi_{AB}(\omega) \tag{46}$$

where

$$X(\omega) = \int_{-\infty}^{\infty} dt\, X(t)e^{i\omega t}, \quad X(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} d\omega\, X(\omega)e^{-i\omega t} \tag{47}$$

applies for any time-dependent quantity $X(t)$.

Because by definition $G_{AB}^{ret}(\omega)$ is analytical only in the upper complex semi-plane (retarded sheet), for a real argument $\omega$ the limit $\varpi \rightarrow \omega + i0$ has to be considered. The *complex admittance* $\chi_{AB}(\omega)$ can therefore be expressed in terms of the retarded Green function as

$$\chi_{AB}(\omega) = \frac{1}{\hbar}G_{AB}^{ret}(\omega + i0) = -\frac{i}{\hbar} \int_0^\infty dt\, e^{i(\omega + i0)t}Tr(\hat{\varrho}_0[\hat{A}(t), \hat{B}(0)]) \tag{48}$$

The occurance of the side-limit $\omega + i0$ in $\chi_{AB}(\omega)$ is usually termed adiabatic switching on of the perturbation as it corresponds to a time-dependent classical field,

$$F'(t) = \lim_{s \to +0}(F(t)e^{st}) \tag{49}$$

### 2.3.1.2. The Kubo Formula   Returning now to Eq. (38).

$$\delta A(t) = -\frac{i}{\hbar} \int_{-\infty}^{t} dt\, Tr([\hat{H}'_H(t'), \hat{\varrho}_0]\hat{A}_H(t)) \tag{50}$$

where the operators are defined within the Heisenberg picture with respect to the unperturbed system and using *Kubo's identity*,

$$\frac{i}{\hbar}[\hat{X}_H(t), \hat{\varrho}] = \hat{\varrho} \int_0^\beta d\lambda\, \dot{\hat{X}}_H(t - i\lambda\hbar).$$

$$\varrho = \frac{e^{-\beta\hat{H}}}{Tr(e^{-\beta\hat{H}})}, \quad \hat{X}_H(t) = e^{i t \hat{H}/\hbar}\hat{X}(t)e^{-i t \hat{H}/\hbar}, \quad \dot{\hat{X}}_H(t) = -\frac{i}{\hbar}[\hat{X}_H(t), \hat{H}] \tag{51}$$

in Eq. (50) finally yields the famous *Kubo formula*:

$$\delta A(t) = -\int_{-\infty}^{t} dt' \int_{0}^{\beta} d\lambda\, Tr(\hat{\varrho}_0 \dot{\hat{H}}'_{II}(t' - i\lambda\hbar)\hat{A}_{II}(t))$$

$$= -\int_{-\infty}^{t} dt' \int_{0}^{\beta} d\lambda\, Tr(\hat{\varrho}_0 \dot{\hat{H}}'_{II}(t')\hat{A}_{II}(t - t' + i\lambda\hbar)) \tag{52}$$

### 2.3.2. The Current-Current Correlation Function

In the case of electric transport, a time-dependent external electric field is applied. Obviously, this induces currents, which in turn creates internal electric fields. Suppose that the total electric field, $E(r, t)$ is related to the perturbation, $\hat{H}'(t)$ in terms of a scalar potential $\phi(r, t)$ as

$$\hat{H}'(t) = \int d^3r\, \hat{\rho}(r)\phi(r, t), \quad E(r, t) = -\nabla\phi(r, t) \tag{53}$$

where $\hat{\rho}(r) = e\psi(r)^+\psi(r)$ is the charge density operator, $\psi(r)$ a field operator and $e$ the charge of an electron. Then the time-derivative of $\hat{H}'_{II}(t)$ can be calculated as follows,

$$\dot{\hat{H}}'(t) = \int d^3r\, \underbrace{\frac{1}{i\hbar}[\hat{H}_0, \hat{\rho}(r)]}_{\frac{\partial\rho(r,t)}{\partial t}\big|_{t, 0}}\phi(r, t) = -\int d^3r\, \nabla\hat{J}(r)\,\phi(r, t)$$

$$= \int d^3r\, \hat{J}(r)\,\nabla\phi(r, t) = -\int d^3r\, \hat{J}(r)\,E(r, t) \tag{54}$$

where the current-density operator is given by

$$\hat{J}(r) = \begin{cases} \dfrac{e\hbar}{2mi}\psi(r)^+(\overrightarrow{\nabla} - \overleftarrow{\nabla})\psi(r), & \text{in non-relativistic case,} \\[2ex] ec\psi(r)^+\,\hat{\alpha}\,\psi(r), & \text{in relativistic case} \end{cases} \tag{55}$$

and the $\hat{\alpha}$ denote Dirac matrices. Note that in Eq. (54) the continuity equation was used and periodic boundary conditions were assumed such that when using Gauss' integration theorem the corresponding surface term vanishes. Making use of Eqs. (52) and (54), the $\mu$th component of the current density can be written as

$$J_\mu(r, t) = \sum_\nu \int d^3r' \int_{-\infty}^{t} dt'\sigma_{\mu\nu}(r, r'; t, t')\,E_\nu(r', t') \tag{56}$$

where the occurring space-time correlation function is given by

$$\sigma_{\mu\nu}(r, r'; t, t') = \Theta(t - t')\int_{0}^{\beta} d\lambda\, Tr(\hat{\varrho}_0\,\hat{J}_\nu(r, 0)\hat{J}_\mu(r', t - t' + i\lambda\hbar)) \tag{57}$$

by which the linear response of the current density at $(r, t)$ in direction $\mu$ is correlated to the local electric field at $(r', t')$ applied in direction $\nu$. Note that in the above equation the current-density operators are assumed to be Heisenberg operators.

Consider now the Fourier components of the electric field,

$$E(q, \omega) = \int d^3r \int_{-\infty}^{\infty} dt\, E(r, t)e^{-iq\cdot r - i\varpi t} \tag{58}$$

$$E(r, t) = \frac{1}{2\pi V}\int d^3q \int_{-\infty}^{\infty} d\omega\, E(q, \omega)e^{iq\cdot r - i\varpi t} \tag{59}$$

where $\varpi = \omega + i0$ and $V$ is the volume of the system. Although $\sigma_{\mu\nu}(r, r'; t, t')$ trivially depends only on $(t - t')$, in general, it is a function of the independent space variables $r$ and $r'$. In cases, however, the current density is an average of the local current density in Eq. (56) over a large enough region, $\sigma_{\mu\nu}(r, r'; t, t')$ can be assumed to be homogeneous in space, that is, $\sigma_{\mu\nu}(r, r'; t - t') = \sigma_{\mu\nu}(r - r'; t - t')$. This usually is the case if $|q|$ is small,

implying that *long-wavelength excitations* are studied. The $(\mathbf{q}, \omega)$ component of the current density per unit volume,

$$J_\mu(\mathbf{q}, \omega) = \frac{1}{V} \int d^3r \int_{-\infty}^{\infty} dt\, J_\mu(\mathbf{r}, t) e^{-i\mathbf{q}\cdot\mathbf{r}+i\omega t} \tag{60}$$

can then be determined from Eqs. (56) and (57),

$$J_\mu(\mathbf{q}, \omega) = \sum_\nu \sigma_{\mu\nu}(\mathbf{q}, \omega) E_\nu(\mathbf{q}, \omega) \tag{61}$$

where $\sigma_{\mu\nu}(\mathbf{q}, \omega)$ is the *wave-vector* and *frequency-dependent conductivity tensor*,

$$\sigma_{\mu\nu}(\mathbf{q}, \omega) = \frac{1}{V} \int_0^\infty dt\, e^{i\omega t} \int_0^\beta d\lambda\, Tr(\hat{\varrho}_0 \hat{J}_\nu(-\mathbf{q}, 0)\hat{J}_\mu(\mathbf{q}, t + i\lambda\hbar)) \tag{62}$$

and

$$J_\mu(\mathbf{q}, t) = \int d^3r\, J_\mu(\mathbf{r}, t) e^{-i\mathbf{q}\cdot\mathbf{r}} \tag{63}$$

In using contour integration techniques one arrives at

$$\sigma_{\mu\nu}(\mathbf{q}, \omega) = \frac{i}{\hbar V} \int_0^\infty dt\, e^{i\omega t} \int_t^\infty dt'\, Tr(\hat{\varrho}_0 [\hat{J}_\mu(\mathbf{q}, t'), \hat{J}_\nu(-\mathbf{q}, 0)]) \tag{64}$$

such that by introducing the below *current-current correlation function*,

$$\Sigma_{\mu\nu}(\mathbf{q}, \omega) = \frac{1}{\hbar V} \int_0^\infty dt\, e^{i\omega t}\, Tr(\hat{\varrho}_0 [\hat{J}_\mu(\mathbf{q}, t), \hat{J}_\nu(-\mathbf{q}, 0)]) \tag{65}$$

the conductivity tensor can finally be expressed as

$$\sigma_{\mu\nu}(\mathbf{q}, \omega) = \frac{\Sigma_{\mu\nu}(\mathbf{q}, \omega) - \Sigma_{\mu\nu}(\mathbf{q}, 0)}{\omega} \tag{66}$$

For a homogeneous system with carrier density $n$ and mass of carriers $m$,

$$-\frac{\Sigma_{\mu\nu}(\mathbf{q}, 0)}{\omega} = i\frac{ne^2}{m\omega} \delta_{\mu\nu} \tag{67}$$

one obtains the phenomenological Drude term for noninteracting particles. Furthermore, the *static limit*, that is, when $\omega \to 0$ and $|\mathbf{q}| \to 0$, is defined as

$$\sigma_{\mu\nu}(\mathbf{q} = 0, \omega = 0) = \lim_{s \to +0} \frac{\Sigma_{\mu\nu}(\mathbf{q} = 0, is) - \Sigma_{\mu\nu}(\mathbf{q} = 0, 0)}{is}$$

$$= \frac{d\Sigma_{\mu\nu}(\mathbf{q} = 0, \omega)}{d\omega}\bigg|_{\omega=0} \tag{68}$$

### 2.3.3. Kubo Formula for Independent Particles

An important special case arises when considering independent particles. Represented in the basis of the eigenfuctions of $\hat{H}_0$ (spectral representation),

$$\hat{H}_0 |n\rangle = \varepsilon_n |n\rangle, \quad \langle m|n\rangle = \delta_{mn}, \quad \sum_n |n\rangle\langle n| = \hat{I} \tag{69}$$

the equilibrium density operator and its matrix elements are given by

$$\hat{\varrho}_0 = \sum_n f(\varepsilon_n)|n\rangle\langle n|, \quad \langle n|\hat{\varrho}_0|p\rangle = f(\varepsilon_n)\delta_{pn} \tag{70}$$

and the thermal average of the current–current commutator can be written as

$$Tr(\hat{\varrho}_0 [\hat{J}_\mu(\mathbf{q}, t), \hat{J}_\nu(-\mathbf{q}, 0)]) = \sum_{nm} \{f(\varepsilon_n) - f(\varepsilon_m)\} e^{(it/\hbar)(\varepsilon_n - \varepsilon_m)}$$

$$\times J_\mu^{nm}(\mathbf{q}) J_\nu^{mn}(-\mathbf{q}) \tag{71}$$

with

$$J_\mu^{nm}(\mathbf{q}) \equiv \langle n|\widehat{J}_\mu(\mathbf{q})|m\rangle \quad \text{and} \quad J_\nu^{nm}(-\mathbf{q}) \equiv \langle m|\widehat{J}_\nu(-\mathbf{q})|n\rangle \tag{72}$$

Substituting Eq. (71) into Eq. (65) then yields

$$\Sigma_{\mu\nu}(\mathbf{q},\omega) = \frac{1}{\hbar V}\sum_{nm}\{f(\varepsilon_n) - f(\varepsilon_m)\}J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})\int_0^\infty dt\, e^{(i/\hbar)(\varepsilon_n - \varepsilon_m + \hbar\varpi)t} \tag{73}$$

The integral with respect to $t$, however, is just the Laplace transform of the identity,

$$\int_0^\infty dt\, e^{[(i/\hbar)(\varepsilon_n - \varepsilon_m + \hbar\omega) - s]t} \underset{(s\to0)}{=} -\frac{e^{[(i/\hbar)(\varepsilon_n - \varepsilon_m + \hbar\omega) - s]}}{(i/\hbar)(\varepsilon_n - \varepsilon_m + \hbar\omega) - s} \tag{74}$$

therefore, Eq. (73) can be transformed to

$$\Sigma_{\mu\nu}(\mathbf{q},\omega) = \frac{i}{V}\sum_{nm}\frac{f(\varepsilon_n) - f(\varepsilon_m)}{\varepsilon_n - \varepsilon_m + \hbar\varpi}J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q}) \tag{75}$$

with $f(\varepsilon)$ being the Fermi-Dirac function. Together with Eq. (64), this now provides a numerically tractable tool to calculate the conductivity tensor. Because

$$\frac{1}{\varepsilon_n - \varepsilon_m + \hbar\varpi} - \frac{1}{\varepsilon_n - \varepsilon_m} = \frac{-\hbar\varpi}{(\varepsilon_n - \varepsilon_m)(\varepsilon_n - \varepsilon_m + \hbar\varpi)}$$

$\sigma_{\mu\nu}(\mathbf{q},\omega)$ can also be written in the following compact form,

$$\sigma_{\mu\nu}(\mathbf{q},\omega) = \frac{\hbar}{iV}\sum_{nm}\frac{f(\varepsilon_n) - f(\varepsilon_m)}{\varepsilon_n - \varepsilon_m}\frac{J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})}{\varepsilon_n - \varepsilon_m + \hbar\varpi}. \tag{76}$$

$$\varpi = \omega + i\delta$$

## 2.3.4. Contour Integrations

$\Sigma_{\mu\nu}(\mathbf{q},\varpi)$ can be evaluated by using contour integration techniques. Considering a pair of eigenvalues, $\varepsilon_n$ and $\varepsilon_m$, for a suitable contour $C$ in the complex energy plane (see Fig. 1) the residue theorem implies that

$$\oint_C dz\frac{f(z)}{(z - \varepsilon_n)(z - \varepsilon_m + \hbar\omega + i\delta)} = -2\pi i\frac{f(\varepsilon_n)}{\varepsilon_n - \varepsilon_m + \hbar\omega + i\delta}$$

$$+ 2i\delta_r\sum_{k=-N_2+1}^{N_3}\frac{1}{(z_k - \varepsilon_n)(z_k - \varepsilon_m + \hbar\omega + i\delta)} \tag{77}$$



Figure 1. Integrations along the contours $C$ (left entry) and $C'$ (right entry). Reprinted with permission from [10], J. Zabloudil et al., "Electron Scattering in Solid Matter," Springer, New York, 2004. © 2004, Springer.

where the $z_k = E_F + i(2k - 1)\delta_T$ are the (fermionic) Matsubara-poles with $E_F$ being the Fermi energy, $\delta_T \equiv \pi k_B T$ and $T$ the temperature. In Eq. (77) it was supposed that $N_1$ and $N_2$ Matsubara-poles in the upper and lower semi-plane lie within the contour $C$, respectively. Equation (77) can be rearranged as follows

$$i\frac{f(\varepsilon_n)}{\varepsilon_n - \varepsilon_m + \hbar\omega + i\delta} = -\frac{1}{2\pi}\oint_{C} dz \frac{f(z)}{(z - \varepsilon_n)(z - \varepsilon_m + \hbar\omega + i\delta)}$$

$$+ i\frac{\delta_T}{\pi}\sum_{k = N_1+1}^{N_1} \frac{1}{(z_k - \varepsilon_n)(z_k - \varepsilon_m + \hbar\omega + i\delta)} \qquad (78)$$

Similarly, by choosing contour $C'$ the following expression,

$$-i\frac{f(\varepsilon_m)}{\varepsilon_n - \varepsilon_m + \hbar\omega + i\delta} = \frac{1}{2\pi}\oint_{C'} dz \frac{f(z)}{(z - \varepsilon_m)(z - \varepsilon_n - \hbar\omega - i\delta)}$$

$$+ i\frac{\delta_T}{\pi}\sum_{k = -N_1+1}^{N_2} \frac{1}{(z_k - \varepsilon_m)(z_k - \varepsilon_n - \hbar\omega - i\delta)} \qquad (79)$$

can be derived. Deforming the contours such that the real axis is crossed at $\infty$ and $-\infty$, $\Sigma_{\mu\nu}(\mathbf{q}, \varpi)$ can be expressed as

$$\Sigma_{\mu\nu}(\mathbf{q}, \varpi) = -\frac{1}{2\pi V}\left\{\oint_{C} dz\, f(z)\sum_{mn} \frac{J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})}{(z - \varepsilon_n)(z - \varepsilon_m + \hbar\omega + i\delta)}\right.$$

$$\left. -\oint_{C'} dz\, f(z)\sum_{mn} \frac{J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})}{(z - \varepsilon_m)(z - \varepsilon_n - \hbar\omega - i\delta)}\right\}$$

$$+ i\frac{\delta_T}{\pi V}\left\{\sum_{k = -N_1+1}^{N_1}\sum_{mn} \frac{J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})}{(z_k - \varepsilon_n)(z_k - \varepsilon_m + \hbar\omega + i\delta)}\right.$$

$$\left. + \sum_{k = -N_1+1}^{N_2}\sum_{mn} \frac{J_\mu^{nm}(\mathbf{q})J_\nu^{mn}(-\mathbf{q})}{(z_k - \varepsilon_m)(z_k - \varepsilon_n - \hbar\omega - i\delta)}\right\} \qquad (80)$$

Consider now the resolvent of the unperturbed Hamilton operator, i.e., of the Kohn-Sham Hamiltonian,

$$\hat{G}(z) = (z\hat{I} - \hat{H})^{-1} \qquad (81)$$

and its adjoint,

$$\hat{G}(z)^\dagger = (z^*\hat{I} - \hat{H})^{-1} = \hat{G}(z^*) \qquad (82)$$

$$\hat{G}(z) = \sum_n \frac{|n\rangle\langle n|}{z - \varepsilon_n} \qquad (83)$$

it is straightforward to rewrite Eq. (80) as

$$\Sigma_{\mu\nu}(\mathbf{q}, \varpi) = -\frac{1}{2\pi V}\left\{\oint_{C} dz\, f(z)\, Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(z + \hbar\omega + i\delta)\hat{J}_\nu(-\mathbf{q})\hat{G}(z))\right.$$

$$\left. -\oint_{C'} dz\, f(z)\, Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(z)\hat{J}_\nu(-\mathbf{q})\hat{G}(z - \hbar\omega - i\delta))\right\}$$

$$+ i\frac{\delta_T}{\pi V}\left\{\sum_{k = -N_1+1}^{N_1} Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(z_k + \hbar\omega + i\delta)\hat{J}_\nu(-\mathbf{q})\,\hat{G}(z_k))\right.$$

$$\left. + \sum_{k = -N_1+1}^{N_2} Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(z_k)\hat{J}_\nu(-\mathbf{q})\hat{G}(z_k - \hbar\omega - i\delta))\right\} \qquad (84)$$

Introducing for matters of convenience the quantity,

$$\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z_1, z_2) = -\frac{1}{2\pi V} Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(z_1)\hat{J}_\nu(-\mathbf{q})\hat{G}(z_2)) \tag{85}$$

$$\tilde{\Sigma}_{\nu\mu}(-\mathbf{q}; z_2, z_1) = \tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z_1, z_2),$$

$$\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z_1^*, z_2^*) = \tilde{\Sigma}_{\nu\mu}(\mathbf{q}; z_1, z_2)^* = \tilde{\Sigma}_{\mu\nu}(-\mathbf{q}; z_2, z_1)^* \tag{86}$$

because of the reflection symmetry for the contours $C$ and $C'$, Eq. (84) can be written as

$$\Sigma_{\mu\nu}(\mathbf{q}, \varpi) = \oint_C dz\, f(z)\, \tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z + \hbar\omega + i\delta, z) - \left(\oint_C dz\, f(z)\, \tilde{\Sigma}_{\mu\nu}(-\mathbf{q}; z - \hbar\omega + i\delta, z)\right)^*$$

$$- 2i\delta_\Gamma \sum_{k=N_2+1}^{N_1} (\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z_k + \hbar\omega + i\delta, z_k) + \tilde{\Sigma}_{\mu\nu}(-\mathbf{q}; z_k - \hbar\omega + i\delta, z_k)^*) \tag{87}$$

## 2.3.5. Integration Along the Real Axis: The Limit of Zero Lifetime Broadening

Deforming the contour $C$ to the real axis such that the contributions from the Matsubara poles vanish and using the relations in Eq. (86), Eq. (87) trivially reduces to

$$\Sigma_{\mu\nu}(\mathbf{q}, \varpi) = \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\{\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; \varepsilon + \hbar\omega + i\delta, \varepsilon + i0) - \tilde{\Sigma}_{\mu\nu}(-\mathbf{q}; \varepsilon + \hbar\omega + i\delta, \varepsilon - i0)\}$$

$$- \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\{\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; \varepsilon - i0, \varepsilon - \hbar\omega - i\delta) - \tilde{\Sigma}_{\mu\nu}(-\mathbf{q}; \varepsilon + i0, \varepsilon - \hbar\omega - i\delta)\} \tag{88}$$

or, by inserting the definition of $\tilde{\Sigma}_{\mu\nu}(\mathbf{q}; z_1, z_2)$,

$$\Sigma_{\mu\nu}(\mathbf{q}, \varpi) = -\frac{1}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\{Tr(\hat{J}_\mu(\mathbf{q})\hat{G}(\varepsilon + \hbar\omega + i\delta)\hat{J}_\nu(-\mathbf{q})\hat{G}^+(\varepsilon))$$

$$- Tr(\hat{J}_\mu(-\mathbf{q})\hat{G}(\varepsilon + \hbar\omega + i\delta)\hat{J}_\nu(\mathbf{q})\hat{G}^-(\varepsilon))$$

$$- Tr(\hat{J}_\mu(\mathbf{q})\hat{G}^-(\varepsilon)\hat{J}_\nu(-\mathbf{q})\hat{G}(\varepsilon - \hbar\omega - i\delta))$$

$$+ Tr(\hat{J}_\mu(-\mathbf{q})\hat{G}^+(\varepsilon)\hat{J}_\nu(\mathbf{q})\hat{G}(\varepsilon - \hbar\omega - i\delta))\} \tag{89}$$

where $\hat{G}^+(\varepsilon)$ and $\hat{G}^-(\varepsilon)$ are the so-called up- and down-side limits of the resolvent

$$\hat{G}^\pm(\varepsilon) = \lim_{\theta \to +0} \hat{G}(\varepsilon \pm i\theta), \quad \hat{G}^\pm(\varepsilon)^\dagger = \hat{G}^\mp(\varepsilon) \tag{90}$$

By taking the limit $\delta \to 0$, Eq. (89) reduces to

$$\Sigma_{\mu\nu}(\mathbf{q}, \omega) = -\frac{1}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\{Tr(\hat{J}_\mu(\mathbf{q})\hat{G}^+(\varepsilon + \hbar\omega)\hat{J}_\nu(-\mathbf{q})\hat{G}^+(\varepsilon))$$

$$- Tr(\hat{J}_\mu(-\mathbf{q})\hat{G}^+(\varepsilon + \hbar\omega)\hat{J}_\nu(\mathbf{q})\hat{G}^-(\varepsilon))$$

$$- Tr(\hat{J}_\mu(\mathbf{q})\hat{G}^-(\varepsilon)\hat{J}_\nu(-\mathbf{q})\hat{G}^-(\varepsilon - \hbar\omega))$$

$$+ Tr(\hat{J}_\mu(-\mathbf{q})\hat{G}^+(\varepsilon)\hat{J}_\nu(\mathbf{q})\hat{G}^-(\varepsilon - \hbar\omega))\} \tag{91}$$

which for $\mathbf{q} = 0$ yields

$$\Sigma_{\mu\nu}(\omega) = -\frac{1}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\{Tr(\hat{J}_\mu\hat{G}^+(\varepsilon + \hbar\omega)\hat{J}_\nu[\hat{G}^+(\varepsilon) - \hat{G}^-(\varepsilon)]) \tag{92}$$

$$+ Tr(\hat{J}_\mu[\hat{G}^+(\varepsilon) - \hat{G}^-(\varepsilon)]\hat{J}_\nu\hat{G}^-(\varepsilon - \hbar\omega))\} \tag{93}$$

### 2.3.6. The Static Limit

In order to obtain the correct zero-frequency conductivity tensor, Eq. (89) has to be used in Eq. (68). Making use of the analyticity of the Green functions in the upper and lower complex semiplanes this then leads to the famous *Kubo-Luttinger formula* [5, 6],

$$
\sigma_{\mu\nu} = -\frac{\hbar}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)
$$

$$
\times Tr\left( \hat{J}_\mu \frac{\partial \widehat{G}^+(\varepsilon)}{\partial \varepsilon} \hat{J}_\nu [\widehat{G}^+(\varepsilon) - \widehat{G}^-(\varepsilon)] - \hat{J}_\mu [\widehat{G}^+(\varepsilon) - \widehat{G}^-(\varepsilon)] \hat{J}_\nu \frac{\partial \widehat{G}^-(\varepsilon)}{\partial \varepsilon} \right) \quad (94)
$$

Integrating by parts yields

$$
\sigma_{\mu\nu} = -\int_{-\infty}^{\infty} d\varepsilon \frac{df(\varepsilon)}{d\varepsilon} S_{\mu\nu}(\varepsilon) \quad (95)
$$

with

$$
S_{\mu\nu}(\varepsilon) = -\frac{\hbar}{2\pi V} \int_{-\infty}^{\varepsilon} d\varepsilon'
$$

$$
\times Tr\left( \hat{J}_\mu \frac{\partial \widehat{G}^+(\varepsilon')}{\partial \varepsilon'} \hat{J}_\nu [\widehat{G}^+(\varepsilon') - \widehat{G}^-(\varepsilon')] - \hat{J}_\mu [\widehat{G}^+(\varepsilon') - \widehat{G}^-(\varepsilon')] \hat{J}_\nu \frac{\partial \widehat{G}^-(\varepsilon')}{\partial \varepsilon'} \right) \quad (96)
$$

which has the meaning of a zero-temperature, energy dependent conductivity. For $T = 0$, $\sigma_{\mu\nu}$ is obviously given by

$$
\sigma_{\mu\nu} = S_{\mu\nu}(E_F) \quad (97)
$$

A numerically tractable expression can be obtained only for the *diagonal elements of the conductivity tensor*, the so-called *Kubo-Greenwood formula* [7, 8] for the dc-conductivity at finite temperatures,

$$
\sigma_{\mu\mu} = -\frac{\hbar}{4\pi V} \int_{-\infty}^{\infty} d\varepsilon \left( -\frac{df(\varepsilon)}{d\varepsilon} \right) Tr(\hat{J}_\mu [\widehat{G}^+(\varepsilon) - \widehat{G}^-(\varepsilon)] \hat{J}_\mu [\widehat{G}^+(\varepsilon) - \widehat{G}^-(\varepsilon)] \quad (98)
$$

which at $T = 0$ temperature obviously can be written as

$$
\sigma_{\mu\mu} = -\frac{\hbar}{4\pi V} Tr(\hat{J}_\mu [\widehat{G}^-(E_F) - \widehat{G}^-(E_F)] \hat{J}_\mu [\widehat{G}^+(E_F) - \widehat{G}^-(E_F)])
$$

$$
= \frac{\hbar}{\pi V} Tr(\hat{J}_\mu \text{Im}\,\widehat{G}^+(E_F) \hat{J}_\mu \text{Im}\,\widehat{G}^+(E_F)) \quad (99)
$$

Recalling finally the spectral resolution of the resolvent,

$$
\text{Im}\,\widehat{G}^+(\varepsilon) = -\pi \sum_n |n\rangle \langle n| \delta(\varepsilon - \varepsilon_n) \quad (100)
$$

it is easy to see that Eq. (99) is identical with the original *Greenwood equation* [7],

$$
\sigma_{\mu\mu} = \frac{\pi\hbar}{V} \sum_{nm} J_\mu^{nm} J_\mu^{mn} \delta(E_F - \varepsilon_n) \delta(E_F - \varepsilon_m) \quad (101)
$$

Equation (94), however, can also be reformulated as follows.

$$
\sigma_{\mu\nu} = \frac{\hbar}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\, Tr\left( \hat{J}_\mu \frac{d\widehat{G}^+(\varepsilon)}{d\varepsilon} \hat{J}_\nu \widehat{G}^+(\varepsilon) + \hat{J}_\mu \widehat{G}^+(\varepsilon) \hat{J}_\nu \frac{d\widehat{G}^+(\varepsilon)}{d\varepsilon} \right)
$$

$$
- \frac{\hbar}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\, Tr\left( \hat{J}_\mu \frac{d\widehat{G}^-(\varepsilon)}{d\varepsilon} \hat{J}_\nu \widehat{G}^-(\varepsilon) + \hat{J}_\mu \widehat{G}^-(\varepsilon) \hat{J}_\nu \frac{d\widehat{G}^-(\varepsilon)}{d\varepsilon} \right)
$$

$$
= \frac{\hbar}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon \left( -\frac{df(\varepsilon)}{d\varepsilon} \right) Tr(\hat{J}_\mu \widehat{G}^+(\varepsilon) \hat{J}_\nu \widehat{G}^-(\varepsilon))
$$

$$
- \frac{\hbar}{2\pi V} \int_{-\infty}^{\infty} d\varepsilon\, f(\varepsilon)\, Tr\left( \hat{J}_\mu \frac{d\widehat{G}^-(\varepsilon)}{d\varepsilon} \hat{J}_\nu \widehat{G}^+(\varepsilon) + \hat{J}_\mu \widehat{G}^+(\varepsilon) \hat{J}_\nu \frac{d\widehat{G}^-(\varepsilon)}{d\varepsilon} \right) \quad (102)
$$

namely in terms of an expression which is similar to that of Baranger and Stone [9], but clearly can be cast into a relativistic form. This expression is of practical relevance reasonable if conductances have to be calculated.

# 3. GREEN'S FUNCTIONS AND SCATTERING PATH OPERATORS

In the following, only a very brief summary of multiple scattering is given. For a detailed treatise on this topic, the reader is refered to a very recent book [10] by some of the authors of the current article that contains also so-called full-potential approaches not considered here.

Suppose the potential in Eq. (1) can be partitioned into non-overlapping, spherically symmetric potentials $V_i$, centered at lattice positions $R_i$, $i = 1 \ldots N$,

$$V(\mathbf{r}) = \sum_{i=1}^{N} V_i(\mathbf{r}_i)\,(\mathbf{r}_i = \mathbf{r} - \mathbf{R}_i) \tag{103}$$

$$V_i(\mathbf{r}_i) = \begin{cases} V_i(r_i) & \text{if } |\mathbf{r}_i| < S_i \\ \text{constant} & \text{otherwise} \end{cases} \tag{104}$$

where $N$ denotes the number of scatterers in the system. For non-overlapping spheres this refers to the so-called *muffin-tin* approach and $S_i$ is called the *muffin-tin radius* of the $i$th sphere. In the so-called atomic sphere approximation (ASA), the spheres are chosen to have the same volume as the Wigner-Seitz cell, thus they overlap slightly, the effect of overlapping, however, is neglected. In the region between the spheres the potential is a constant, commonly set to zero.

## 3.1. Single-Site Scattering

In the absence of effective fields, the *Kohn-Sham-Dirac equation* is of the form [11, 12],

$$\hat{H}|\psi\rangle = \begin{pmatrix} (V(r) + mc^2)\hat{I}_2 & c\hat{\sigma}_r\left(\dfrac{\partial}{\partial r} + \dfrac{1}{r} - \dfrac{1}{r}\beta\hat{K}\right) \\ c\hat{\sigma}_r\left(\dfrac{\partial}{\partial r} + \dfrac{1}{r} - \dfrac{1}{r}\beta\hat{K}\right) & (V(r) - mc^2)\hat{I}_2 \end{pmatrix} |\psi\rangle = W|\psi\rangle \tag{105}$$

where $c$ is the speed of light, $\hat{\sigma}_r = \hat{\mathbf{r}} \cdot \hat{\sigma}$ with $\hat{\mathbf{r}} = \mathbf{r}/|\mathbf{r}|$, $W$ is the total energy of the particle,

$$W^2 = p^2 c^2 + m^2 c^4$$

with $p$ being its momentum and

$$\hat{K} = \hat{\sigma} \cdot \hat{\mathbf{L}} + \hbar\hat{I}_2, \quad \text{and} \quad \hat{\beta} = \begin{pmatrix} \hat{I}_1 & 0 \\ 0 & -\hat{I}_1 \end{pmatrix} \tag{106}$$

The wavefunction $|\psi\rangle$ can be decoupled into two bi-spinors: $|\psi\rangle = |\phi, \chi\rangle$. The total angular momentum operator is defined as $\hat{\mathbf{J}} = \hat{\mathbf{L}} + \hat{\mathbf{S}}$, where $\hat{\mathbf{L}}$ is the angular momentum operator and $\hat{\mathbf{S}} = \frac{\hbar}{2}\hat{\sigma}$ is the spin momentum operator. The eigenfunctions of $\hat{J}^2$ and $\hat{J}_z$,

$$\hat{J}^2|\phi\rangle = \hbar^2 j(j+1)|\phi\rangle, \quad \hat{J}^2|\chi\rangle = \hbar^2 j(j+1)|\chi\rangle, \quad j = l \pm \frac{1}{2} = \frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \ldots,$$

$$\hat{J}_z|\phi\rangle = \hbar\mu|\phi\rangle, \quad \hat{J}_z|\chi\rangle = \hbar\mu|\chi\rangle, \quad \mu = -j, \ldots, j. \tag{107}$$

$$\hat{K}|\phi\rangle = -\hbar\kappa|\phi\rangle, \quad \hat{K}|\chi\rangle = \hbar\kappa|\chi\rangle, \quad \kappa = \mp\left(j + \frac{1}{2}\right)$$

are the so-called spin spherical harmonics,

$$|\kappa, \mu\rangle = |Q\rangle = \sum_{s=\pm\frac{1}{2}} C\left(l, j, \frac{1}{2}\Big|\mu - s, s\right)|l, \mu - s\rangle\Phi_s, \quad |-\kappa, \mu\rangle = |\bar{Q}\rangle \tag{108}$$

where the $C(l, j, 1/2|\mu - s, s)$ denote *Clebsch-Gordan coefficients* [13], $|l, \mu - s\rangle$ *complex spherical harmonics*,

$$\langle l, \mu - s|\hat{\mathbf{r}}\rangle = Y_l^{\mu - s}(\hat{r}) \quad \text{and} \quad \langle \hat{\mathbf{r}}|l, \mu - s\rangle = Y_l^{\mu - s}(\hat{r})^*$$

and the $\Phi_s$ are the following *spinor basis functions* [11],

$$\Phi_{1,2} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \Phi_{-1,2} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \tag{109}$$

It should be noted that in the so-called weak-relativistic limit, the following approach is used,

$$\varepsilon = W - mc^2 = \sqrt{p^2c^2 + m^2c^4} - mc^2 \approx \frac{p^2}{2m} \ll mc^2 \tag{110}$$

### 3.1.1. Free-Space Green's Functions

The nonrelativistic Green's function of a free electron in angular momentum representation can be written as

$$G_0^{nr, \pm}(\varepsilon, \mathbf{r}, \mathbf{r}') = \mp ip \sum_L j_l\left(\frac{pr_<}{\hbar}\right) h_l^\pm\left(\frac{pr_>}{\hbar}\right) Y_L(\hat{r}) Y_L^*(\hat{r}') \tag{111}$$

where $L = (l, m)$, $r_< = \min(r, r')$, $r_> = \max(r, r')$, and $h_l^\pm = j_l \pm in_l$ is a spherical Hankel function with $j_l$ and $n_l$ being spherical Bessel and Neumann functions, respectively [14]. In the relativistic case the Green's function of a free electron in angular momentum representation is of the form

$$G_0^r(\varepsilon, \mathbf{r}, \mathbf{r}') = -ip\frac{\varepsilon + 2mc^2}{2mc^2} \sum_Q [J_Q(\varepsilon, \mathbf{r}) H_Q^+(\varepsilon, \mathbf{r}')^\dagger \theta(r' - r) + H_Q^+(\varepsilon, \mathbf{r}) J_Q^\dagger(\varepsilon, \mathbf{r}') \theta(r - r')] \tag{112}$$

with

$$F_Q(\varepsilon, \mathbf{r}) = \begin{bmatrix} f_l\left(\dfrac{pr}{\hbar}\right)\langle Q|\hat{\mathbf{r}}\rangle \\[2mm] \dfrac{iS_\kappa pc}{\varepsilon + 2mc^2} f_{\bar{l}}\left(\dfrac{pr}{\hbar}\right)\langle \bar{Q}|\hat{\mathbf{r}}\rangle \end{bmatrix} \tag{113}$$

$$F_Q^\dagger(\varepsilon, \mathbf{r}) = \left[f_l\left(\frac{pr}{\hbar}\right)\langle\hat{\mathbf{r}}|Q\rangle, \quad \frac{-iS_\kappa pc}{\varepsilon + 2mc^2} f_{\bar{l}}\left(\frac{pr}{\hbar}\right)\langle\hat{\mathbf{r}}|\bar{Q}\rangle\right] \tag{114}$$

and

$$S_\kappa = \frac{\kappa}{|\kappa|}, \quad \bar{l} = l - S_\kappa$$

with $f_l$ denoting spherical Bessel-, Neumann-, or Hankel functions. A general solution of Eq. (105) can be written as

$$R(\varepsilon, \mathbf{r}) = \sum_Q R_Q(\varepsilon, \mathbf{r}) = \sum_Q \begin{bmatrix} g_\kappa(\varepsilon, r)\langle Q|\hat{\mathbf{r}}\rangle \\ if_\kappa(\varepsilon, r)\langle \bar{Q}|\hat{\mathbf{r}}\rangle \end{bmatrix} \tag{115}$$

a solution outside the bounding sphere as

$$R_Q(\varepsilon, \mathbf{r}) = J_Q(\varepsilon, \mathbf{r}) - ip\sum_Q H_Q^+(\varepsilon, \mathbf{r}) t_{QQ}(\varepsilon) \tag{116}$$

where $t_{QQ}(\varepsilon)$ is usually called *single-site t-matrix*,

$$t_{QQ}(\varepsilon) = \int_s d^3r \int_s d^3r' J_Q^\dagger(\varepsilon, \mathbf{r}') t(\varepsilon, \mathbf{r}', \mathbf{r}) J_Q(\varepsilon, \mathbf{r}) \tag{117}$$

which, for example, for nonmagnetic systems and $\varepsilon \geq 0$ can easily be obtained from the below matching condition

$$g_\kappa(\varepsilon, r) = \cos \delta_\kappa(\varepsilon) j_l\left(\frac{pr}{\hbar}\right) - \sin \delta_\kappa(\varepsilon) n_l\left(\frac{pr}{\hbar}\right) \tag{118}$$

$$cf_\kappa(\varepsilon, r) = pS_\kappa\left[\cos \delta_\kappa(\varepsilon) j_l\left(\frac{pr}{\hbar}\right) - \sin \delta_\kappa(\varepsilon) n_l\left(\frac{pr}{\hbar}\right)\right] \tag{119}$$

at the muffin-tin radius $S$.

$$\tan \delta_\kappa(\varepsilon) = \frac{L_\kappa(\varepsilon, S) j_l(pS/\hbar) - pS_\kappa j_l(pS/\hbar)}{L_\kappa(\varepsilon, r) n_l(pS/\hbar) - pS_\kappa n_l(pS/\hbar)} \tag{120}$$

$$L_\kappa(\varepsilon, S) = \frac{cf_\kappa(\varepsilon, S)}{g_\kappa(\varepsilon, S)} \tag{121}$$

### 3.1.2. Scattering Solutions

The so-called (regular) scattering solutions are defined as follows

$$Z(\varepsilon, r) = \sum_Q Z_Q(\varepsilon, r) = \sum_{Q'Q} R_{Q'}(\varepsilon, r) t_{Q'Q}^{-1}(\varepsilon) \tag{122}$$

$$Z_Q(\varepsilon, r) = \sum_{Q'} J_{Q'}(\varepsilon, r') K_{Q'Q}^{-1}(\varepsilon) + pN_Q(\varepsilon, r) \tag{123}$$

where $K_{QQ'}(z)$ usually is termed *reactance*

$$K_{QQ'}^{-1}(\varepsilon) = t_{QQ'}^{-1}(\varepsilon) - ip\delta_{QQ'} \tag{124}$$

$$K(\varepsilon) = K^\dagger(\varepsilon) \tag{125}$$

## 3.2. Multiple Scattering

### 3.2.1. Scattering Path Operators

The angular momentum representation of the so-called scattering path operator is given by [15]

$$\tau_{QQ'}^{ii}(\varepsilon) = t_{QQ'}^i(\varepsilon)\delta_{ij} + \sum_{k \neq i} \sum_{Q_1 Q_2} t_{QQ_1}^i(\varepsilon) G_{Q_1 Q_2}^{0, ik}(\varepsilon) \tau_{Q_2 Q'}^{kj}(\varepsilon) \tag{126}$$

where the relativistic structure constants

$$G_{QQ'}^{0, ij}(\varepsilon) = \frac{\kappa + 2mc^2}{2mc^2} \sum_s C\left(l, j, \frac{1}{2}\Big|\mu - s, s\right) G_{LL'}^{0, ij}(\varepsilon) C\left(l', j', \frac{1}{2}\Big|\mu' - s, s\right) \tag{127}$$

are obtained from the nonrelativistic ones

$$G_{0}^{m}(\varepsilon, r_i + R_i, r_j' + R_j) = \sum_{L, L'} j_L(\varepsilon, r_i) G_{LL'}^{0, ij}(\varepsilon) j_{L'}^\times(\varepsilon, r_j') \tag{128}$$

In here, the notation $G_{LL'}^{0, ij}(\varepsilon) = G_{LL'}^{0}(\varepsilon, R_i - R_j)$ is applied, $R_i$ and $R_j$ denoting the position vectors of sites $i$ and $j$, and $t_{QQ'}^i(\varepsilon)$ refers to the $t$-matrix at site $R_i$, and

$$G_{LL'}^{0}(\varepsilon, r) = -4\pi i^{(l-l'+1)} p \sum_{L''} C_{LL'}^{L''} i^{-l''} h_{L''}(\varepsilon, r) \tag{129}$$

$$C_{LL'}^{L''} = \int d\hat{r} Y_l(\hat{r}) Y_{l'}^*(\hat{r}) Y_{L''}(\hat{r}) \tag{130}$$

$$j_L(\varepsilon, r) \equiv j_l\left(\frac{pr}{\hbar}\right) Y_L(\hat{r}), \quad j_L^\times(\varepsilon, r) \equiv j_l\left(\frac{pr}{\hbar}\right) Y_L^*(\hat{r}) \tag{131}$$

where the $C_{LL'}^{L''}$ are the so-called Gaunt coefficients.

### 3.2.2. Green's Functions

The Green's function in the case of an ensemble of scatterers can be defined n terms of the scattering path operator and scattering solutions as

$$G(\varepsilon, \mathbf{r}_i, \mathbf{r}'_j) = \sum_{Q, Q'} [Z^i_Q(\varepsilon, \mathbf{r}_i) \tau^{ij}_{QQ'}(\varepsilon) Z^j_{Q'}(\varepsilon, \mathbf{r}'_j)^+ - \delta_{ij}\delta_{QQ'}(\theta(r_i - r'_i)Z^i_Q(\varepsilon, \mathbf{r}_i)I^i_{Q'}\varepsilon, \mathbf{r}_i)^+$$

$$+ \theta(r'_i - r_i)I^i_Q(\varepsilon, \mathbf{r}_i)Z^i_{Q'}(\varepsilon, \mathbf{r}_i)^+)] \tag{132}$$

where $Z^i_Q(\varepsilon, \mathbf{r})$ and $I^i_Q(\varepsilon, \mathbf{r})$ denote the regular and irregular scattering solutons of the Dirac equation in cell $i$ and—as should be recalled—at the muffin-tin radius of he $i$th cell ($S_i$) the following relations have to be satisfied,

$$Z^i_Q(\varepsilon, S_i) = \sum_{Q'}(t^i_{QQ'})^{-1}J_{Q'}(\varepsilon, S_i) - ipH^+_Q(\varepsilon, S_i) \tag{133}$$

$$I^i_Q(\varepsilon, S_i) = J_Q(\varepsilon, S_i) \tag{134}$$

## 3.3. KKR Method for Layered Systems

Layered systems are systems with (at least) two-dimensional translational symmetry. In the case of a surface or an interface, the translational symmetry is "broken" along the direction "perpendicular" to the plane. Suppose such a layered system corresponds to a parent infinite (three-dimensional periodic) system consisting of a simple lattice with only one atom per unit cell, then any lattice site $\mathbf{R}_{pi}$ can be written as

$$\mathbf{R}_{pi} = \mathbf{C}_p + \mathbf{T}_i; \qquad \mathbf{T}_i \in L_2 \tag{135}$$

where $\mathbf{C}_p$ is the "spanning vector" of a particular layer $p$ and the two-dimensonal (real) lattice is denoted by $L_2 = \{\mathbf{T}_i\}$ with in-plane lattice vectors $\mathbf{T}_i$ and where the corresponding set of indices is $I(L_2)$. It should be noted that $\mathbf{C}_p$ not necessarily has to be perpendicular to the planes of atoms, e.g., in a body centered cubic (BCC) lattice for (001)-panes $\mathbf{C}_p = p \cdot a \cdot (\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$, where $a$ is the three-dimensional lattice constant.

The real-space structure constants are now defined by

$$\underline{G}^{ij}_0(\varepsilon) = \underline{G}_0(\varepsilon, \mathbf{R}_{pi} - \mathbf{R}_{qj}) = \underline{G}_0(\varepsilon, \mathbf{C}_p + \mathbf{T}_i - \mathbf{C}_q - \mathbf{T}_j) = \widehat{\underline{G}}^{pq}_0(\varepsilon, \mathbf{T}_i - \mathbf{T}_j)$$

$$= \frac{1}{\Omega_{BZ}} \int_{BZ} d^2k_{\parallel} \widehat{\underline{G}}^{pq}_0(\varepsilon, \mathbf{k}_{\parallel})e^{-i\mathbf{k}_{\parallel}\cdot(\mathbf{T}_i - \mathbf{T}_j)} \tag{136}$$

where $\Omega_{BZ}$ denotes the volume (area) of the two-dimensional Brillouin zone, and the symbol "hat" denotes a layer-indexed quantity; the two-dimensional lattice Fourier transform of the "reciprocal" structure constants is simply given by

$$\widehat{\underline{G}}^{pq}_0(\varepsilon, \mathbf{k}_{\parallel}) = \sum_{\mathbf{T}_i} \widehat{\underline{G}}^{pq}_0(\varepsilon, \mathbf{T}_i)e^{i\mathbf{k}_{\parallel}\cdot\mathbf{T}_i} \tag{137}$$

Introducing the below matrix notation,

$$\underline{t}(\varepsilon) = \{\underline{t}_p(\varepsilon)\delta_{pq}\}, \qquad \widehat{\underline{G}}_0(\varepsilon, \mathbf{k}_{\parallel}) = \{\widehat{\underline{G}}^{pq}_0(\varepsilon, \mathbf{k}_{\parallel})\}, \qquad \underline{\tau}(\varepsilon, \mathbf{k}_{\parallel}) = \{\underline{\tau}_{pq}(\varepsilon, \mathbf{k}_{\parallel})\} \tag{138}$$

the so-called KKR (Korringa-Kohn-Rostoker) equation can be written as

$$\underline{\tau}(\varepsilon, \mathbf{k}_{\parallel}) = (\underline{t}^{-1}(\varepsilon) - \widehat{\underline{G}}_0(\varepsilon, \mathbf{k}_{\parallel}))^{-1} \tag{139}$$

## 3.4. The Screened KKR Method (SKKR)

For systems containing several atoms per unit cell as well as for layered structures severe computational difficulties arise from the long range behavior of the structure constants. For such systems, tight-binding (TB) methods seem to be better suited. However, it can be shown

that by applying a so-called screening transformation, the KKR-method can be transformed into a TB form [16–18]. In Ref. [17], for example, a reference system is suggested with a constant repulsive potential $V^r$, defined within non-overlapping muffin-tin spheres and zero otherwise. In the following, $r$-indexed quantities refer to the reference system and quantities without such an index to those of the physical system, the corresponding *Green's function matrices* being defined by the below Dyson equations

$$\underline{G}(\varepsilon) = \underline{G}_{0}(\varepsilon)(\underline{I} - \underline{t}(\varepsilon)\underline{G}_{0}(\varepsilon))^{-1}, \quad \underline{G}^{r}(\varepsilon) = \underline{G}_{0}(\varepsilon)(\underline{I} - \underline{t}^{r}(\varepsilon)\underline{G}_{0}(\varepsilon))^{-1} \tag{140}$$

Furthermore, in defining the difference of the inverse of the $t$-matrices as

$$\underline{t}_{\Delta}(\varepsilon) = \underline{t}(\varepsilon) - \underline{t}^{r}(\varepsilon) \tag{141}$$

and the screened scattering path operator as

$$\underline{\tau}_{\Delta}(\varepsilon) = \left(\underline{t}_{\Delta}^{-1}(\varepsilon) - \underline{G}^{r}(\varepsilon)\right)^{-1} \tag{142}$$

the unscreened (physical) scattering path operator can be calculated from the screened one by using the following invariance property [17]

$$\underline{G}(\varepsilon) = \underline{t}^{-1}(\varepsilon)\underline{\tau}(\varepsilon)\underline{t}^{-1}(\varepsilon) - \underline{t}^{-1}(\varepsilon) = \underline{t}_{\Delta}^{-1}(\varepsilon)\underline{\tau}_{\Delta}(\varepsilon)\underline{t}_{\Delta}^{-1}(\varepsilon) - \underline{t}_{\Delta}^{-1}(\varepsilon) \tag{143}$$

$$\underline{\tau}(\varepsilon) = \left[\underline{I} - \underline{t}^{r}(\varepsilon)\underline{t}_{\Delta}^{-1}(\varepsilon)\right]\underline{\tau}_{\Delta}(\varepsilon)\left[\underline{I} - \underline{t}_{\Delta}^{-1}(\varepsilon)\underline{t}^{r}(\varepsilon)\right] + \left[\underline{t}^{r}(\varepsilon) - \underline{t}^{r}(\varepsilon)\underline{t}_{\Delta}^{-1}(\varepsilon)\underline{t}^{r}(\varepsilon)\right] \tag{144}$$

In the two-dimensional lattice Fourier transform of the screened scattering path operator,

$$\hat{\tau}_{\Delta, pq}(\varepsilon, \mathbf{k}_{\parallel}) = \left\{\left((\hat{t}_{\Delta})^{-1}(\varepsilon) - \widetilde{\underline{G}}^{r}(\varepsilon, \mathbf{k}_{\parallel})\right)^{-1}\right\}_{pq} \tag{145}$$

however—because of the screening—$\widehat{G}^{r}$ is of block-tridiagonal form, the blocks being related to so-called principal layers that contain $n$ atomic layers ($n \geq 3$). If these layers form the top of a semi-infinite bulk (substrate) or are situated between two semi-infinite bulk systems, a so-called surface Green function method [16] has to be applied to ensure proper boundary conditions. The real-space physical $\tau$-matrix is then obtained first by performing the following Brillouin zone integral,

$$\underline{\tau}_{\Delta}(\varepsilon, \mathbf{R}_{pi} - \mathbf{R}_{qj}) = \frac{1}{\Omega_{BZ}} \int_{BZ} d^{2}k_{\parallel}\hat{\underline{\tau}}_{\Delta, pq}(\varepsilon, \mathbf{k}_{\parallel})e^{-i\mathbf{k}_{\parallel}\cdot(\mathbf{T}_{i} - \mathbf{T}_{j})} \tag{146}$$

and subsequent use of the transformation defined in Eq. (144). It should be noted that in principle for a two-dimensional translational invariant medium the physical real space $\tau$-matrix is in principle defined by

$$\underline{\tau}^{pi, qj}(\varepsilon) = \frac{1}{\Omega_{BZ}} \int_{BZ} d^{2}k_{\parallel}\hat{\underline{\tau}}^{pq}(\varepsilon, \mathbf{k}_{\parallel})e^{-i\mathbf{k}_{\parallel}\cdot(\mathbf{T}_{i} - \mathbf{T}_{j})} \tag{147}$$

## 3.5. The Embedding Technique

A finite *cluster* is defined as a geometrical arrangement of a set of scatterers. Let $\ell$ denote the set of the position vectors of sites in the cluster,

$$\ell \equiv \{\mathbf{R}_{i}\}, \quad i = 1, \ldots, N \tag{148}$$

where $N$ is the number of atoms in the cluster, and $C_{N}$ a corresponding set of site-indices

$$C_{N} \equiv \{i | \mathbf{R}_{i} \in \ell, \quad i = 1, \ldots, N\} \tag{149}$$

Because for the host system (*substrate*) the potential is given by

$$V^{host}(\mathbf{r}) = \sum_{i} V_{i}^{host}(\mathbf{r}_{i}) \tag{150}$$

and for an *embedded cluster* by

$$V^{\text{clus}}(\mathbf{r}) = \sum_i V_i^{\text{clus}}(\mathbf{r}_i), \qquad V_i^{\text{clus}}(\mathbf{r}_i) = \begin{cases} V_i^{\text{host}}(\mathbf{r}_i) & \text{if } i \notin C_N \\ V_i^{\text{imp}}(\mathbf{r}_i) & \text{if } i \in C_N \end{cases} \tag{151}$$

the single site $t$-matrices of the perturbed system are of the form,

$$\underline{t}^i_{\text{clus}}(\varepsilon) = \begin{cases} \underline{t}^i_{\text{host}}(\varepsilon) & \text{if } i \notin C_N \\ \underline{t}^i_{\text{imp}}(\varepsilon) & \text{if } i \in C_N \end{cases} \tag{152}$$

It is important to emphasize that by performing real space embedding, a cluster usually contains the investigated impurity atoms, some sites from the host material (and empty spheres (vacuum) in the case of a surface or a nanocontact). In principle according to the above classification all $V_i^{\text{imp}}$, $\underline{t}^i_{\text{imp}}$ are different.

The KKR-equation for the unperturbed and perturbed systems are then given by,

$$\underline{\tau}^{-1}_{\text{host}}(\varepsilon) = \underline{t}^{-1}_{\text{host}}(\varepsilon) - \underline{G}_0(\varepsilon), \qquad \underline{\tau}^{-1}_{\text{clus}}(\varepsilon) = \underline{t}^{-1}_{\text{clus}}(\varepsilon) - \underline{G}_0(\varepsilon) \tag{153}$$

respectively, with

$$\underline{\tau}(\varepsilon) = \{\underline{\tau}^{ij}(\varepsilon)\}, \qquad \underline{\tau}^{ij}(\varepsilon) = \{\tau^{ij}_{QQ'}(\varepsilon)\}, \qquad \underline{t}(\varepsilon) = \{\underline{t}^i(\varepsilon)\delta_{ij}\}, \qquad \underline{t}^i(\varepsilon) = \{t^i_{QQ'}(\varepsilon)\} \tag{154}$$

Defining the following quantities

$$\Delta\underline{t}^{-1}(\varepsilon) = \underline{t}^{-1}_{\text{host}}(\varepsilon) - \underline{t}^{-1}_{\text{clus}}(\varepsilon) = \{\Delta\underline{t}^i(\varepsilon)^{-1}\delta_{ij}\}$$

$$\Delta\underline{t}^i(\varepsilon)^{-1} = \begin{cases} 0, & \text{if } i \notin C_N \\ \underline{t}^i_{\text{host}}(\varepsilon)^{-1} - \underline{t}^i_{\text{imp}}(\varepsilon)^{-1}, & \text{if } i \in C_N \end{cases} \tag{155}$$

from Eq. (153), one then obtains

$$\underline{\tau}^{-1}_{\text{clus}}(\varepsilon) = \underline{\tau}^{-1}_{\text{host}}(\varepsilon) - \Delta\underline{t}^{-1}(\varepsilon)$$

$$= [\underline{I} - \Delta\underline{t}^{-1}(\varepsilon)\underline{\tau}_{\text{host}}(\varepsilon)]\underline{\tau}^{-1}_{\text{host}}(\varepsilon)$$

$$= \underline{\tau}^{-1}_{\text{host}}(\varepsilon)[\underline{I} - \underline{\tau}_{\text{host}}(\varepsilon)\Delta\underline{t}^{-1}(\varepsilon)] \tag{156}$$

which inverted finally leads to the below *embedding equation*,

$$\underline{\tau}_{\text{clus}}(\varepsilon) = \underline{\tau}_{\text{host}}(\varepsilon)[\underline{I} - \Delta\underline{t}^{-1}(\varepsilon)\underline{\tau}_{\text{host}}(\varepsilon)]^{-1} = [\underline{I} - \underline{\tau}_{\text{host}}(\varepsilon)\Delta\underline{t}^{-1}(\varepsilon)]^{-1}\underline{\tau}_{\text{host}}(\varepsilon) \tag{157}$$

It should be mentioned that for a *single impurity* at site $i_0$ the above embedding equation reduces to

$$\underline{\tau}^{i_0 i_0}(\varepsilon) = \underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon)[\underline{I} - \Delta\underline{t}^{i_0}(\varepsilon)^{-1}\underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon)]^{-1} = \underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon)\widetilde{\underline{D}}^{i_0}(\varepsilon)$$

$$= [\underline{I} - \underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon)\Delta\underline{t}^{i_0}(\varepsilon)^{-1}]^{-1}\underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon) = \underline{D}^{i_0}(\varepsilon)\underline{\tau}^{i_0 i_0}_{\text{host}}(\varepsilon) \tag{158}$$

## 3.6. The Coherent Potential Approximation

### 3.6.1. Configurational Averages

Suppose a binary bulk alloy is of composition $A_c B_{1-c}$ with $c_A = c$ being the concentration of species A and $c_B = (1 - c)$ the concentration of species B. Furthermore, suppose the total number of atoms is $N$ and the number of A atoms and B atoms $N_A$ and $N_B$, respectively,

$$N = N_A + N_B, \qquad N_A = cN, \qquad N_B = (1 - c)N \tag{159}$$

A substitutional binary alloy refers to a system with no positional disorder, all atoms are placed (for matters of simplicity) in the positions of an underlying ideal simple lattice $J$ which is characterized by the set of indices, $I(J)$. Assuming this kind of disorder, the potential can be written as

$$V(\mathbf{r}) = \sum_{i \in I(J)} V_i(\mathbf{r}_i - \mathbf{R}_i) \tag{160}$$

$$V_i(\mathbf{r}_i - \mathbf{R}_i) = \xi_i V_1(\mathbf{r}_i - \mathbf{R}_i) + (1 - \xi_i)V_B(\mathbf{r}_i - \mathbf{R}_i) \tag{161}$$

where $\xi_i$ is an occupational variable such that $\xi_i = 1$ if site $\mathbf{R}_i$ is occupied by species A and $\xi_i = 0$ if this site is occupied by species B. For a completely random alloy the probability for $\xi_i = 1$ is $c_A$ and correspondingly for $\xi_i = 0$ the probability is $c_B$. In Eq. (161), $V_A(\mathbf{r}_i - \mathbf{R}_i)$ and $V_B(\mathbf{r}_i - \mathbf{R}_i)$ are the individual (effective) potentials of species A and B at site $\mathbf{R}_i$, respectively. Then $\{\xi_i \mid i \in I(J)\}$ is one particular arrangement of atoms A and B on the positions of $J$. Such an arrangement is called a *configuration*. Quite clearly for one particular configuration the Kohn-Sham equation,

$$\hat{H}\{\xi_i\}\psi_n(\mathbf{r}, \{\xi_i\}) = \varepsilon_n\{\xi_i\}\psi_n(\mathbf{r}, \{\xi_i\}) \tag{162}$$

where $\hat{H}$ is the Hamiltonian of the system and $n$ labels the eigenstates, can be solved using standard techniques. Observables, however, in general do not map a particular configuration but an average over all configurations. Let $\langle A_{nn'} \rangle$ be the configurationally averaged matrix element of a Hermitian operator $\hat{A}$. Then

$$\langle A_{nn'} \rangle = \sum_{\{\xi_i\}} P(\{\xi_i\})\langle\psi_n\{\xi_i\}|\hat{A}|\psi_{n'}\{\xi_i\}\rangle \tag{163}$$

where $P(\{\xi_i\})$ is the microcanonical probability for a particular configuration $\{\xi_i\}$. In the above equations it was assumed that the occupational probabilities for different sites are independent from each other, that is, that

$$P(\{\xi_i\}) = \prod_i P_i(\xi_i), \qquad \sum_{\xi_i = 0, 1} P_i(\xi_i) = 1$$
$$\langle \xi_i \rangle \equiv P_i(1) = c, \qquad \langle 1 - \xi_i \rangle \equiv P_i(0) = 1 - c \tag{164}$$

Obviously the calculation of averages such in specified Eq. (163) is greatly simplified by directly calculating the configurationally averaged Green function $\langle G^+(\varepsilon, \mathbf{r}, \mathbf{r}') \rangle$ from which typical one-particle physical properties can immediately be obtained.

*Restricted ensemble averages*, denoted by $\langle \ldots \rangle_{(i=\alpha)}$, have the following meaning: in cell $i$ the occupation is fixed to atom $\alpha(\alpha \in \{A, B\})$ and the averaging is restricted to all configurations for the remaining $N - 1$ sites. By using restricted ensemble averages the configurational average is partitioned into two subsets, for which the following condition has to be satisfied,

$$\langle G^+(\varepsilon, \mathbf{r}_i, \mathbf{r}_i) \rangle = \sum_{\alpha \in \{A, B\}} c_\alpha \langle G^+(\varepsilon, \mathbf{r}_i, \mathbf{r}_i) \rangle_{(i=\alpha)} \tag{165}$$

### 3.6.2. The CPA Single-Site Approximation

In the so-called single-site approximation to the coherent potential approximation (CPA), *short-range-order effects* are explicitly excluded. Multiple scattering effects, however, are implicitly included since the single-site approximation is based on the idea of a single scatterer immersed in an average medium, that is, on the very concept of a "mean field theory." From the definition of the scattering path operators follows that for a binary (bulk) system $A_c B_{1-c}$ (simple lattice, one atom per unit cell) the restricted averages $\langle \underline{\tau}^{ii}(\varepsilon) \rangle_{(i=\alpha)}$, $\alpha \in \{A, B\}$, have to meet the condition,

$$c\langle \underline{\tau}^{ii}(\varepsilon) \rangle_{(i=A)} + (1 - c)\langle \underline{\tau}^{ii}(\varepsilon) \rangle_{(i=B)} = \langle \underline{\tau}^{ii}(\varepsilon) \rangle \tag{166}$$

Because (166) is valid for all site indices $i \in I(J)$, it is sufficient to restrict this equation to $i = 0$ ($0 =$ origin of the underlying lattice) for a bulk, or—more general—to $i = p0$ ($p0$, origin of the $p$th layer, $\forall\, p$) for layered systems.

For a given set of $n$ atomic layers, containing also disordered layers, the so-called coherent scattering path operator $\tau_c(\varepsilon)$ is given by the following two-dimensional Brillouin zone integral,

$$\underline{\tau}_c^{pi,\,qj}(\varepsilon) = \frac{1}{\Omega_{BZ}} \int_{BZ} e^{-i\mathbf{k}_\parallel \cdot (\mathbf{T}_i - \mathbf{T}_j)} \hat{\underline{\tau}}_c^{pq}(\varepsilon, \mathbf{k}_\parallel)\, d^2 k_\parallel \tag{167}$$

where $pi$ and $qj$ denote site $i$ in layer $p$ and site $j$ in layer $q$, respectively. Moreover, $\hat{\underline{\tau}}_c^{pq}(\varepsilon, \mathbf{k}_\parallel)$ is the $(pq)$-th block of the supermatrix,

$$\hat{\underline{\tau}}_c(\varepsilon, \mathbf{k}_\parallel) = \left[\hat{\underline{t}}_c(\varepsilon)^{-1} - \hat{\underline{\underline{G}}}(\varepsilon, \mathbf{k}_\parallel)\right]^{-1} \tag{168}$$

Equation (167) implies two–dimensional translational invariance of the coherent medium for all layers under investigation, that is, that in each layer $p$ for the coherent single-site $t$-matrices the following translational invariance applies,

$$\underline{t}_c^{pi}(\varepsilon) = \underline{t}_c^{p0}(\varepsilon) = \hat{\underline{t}}_c^{p}(\varepsilon); \qquad \forall\, i \in I(L_2) \tag{169}$$

Using again supermatrices for a better visualization,

$$\hat{\underline{t}}_c(\varepsilon) = \begin{pmatrix} \hat{\underline{t}}_c^1(\varepsilon) & 0 & \cdots & \cdots & 0 \\ 0 & \ddots & & & \vdots \\ \vdots & & \hat{\underline{t}}_c^p(\varepsilon) & & \vdots \\ \vdots & & & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & \hat{\underline{t}}_c^n(\varepsilon) \end{pmatrix} \tag{170}$$

and

$$\hat{\underline{\tau}}_c(\varepsilon) = \begin{pmatrix} & \vdots & & \vdots & \\ \cdots & \hat{\underline{\tau}}_c^{pp}(\varepsilon) & \cdots & \hat{\underline{\tau}}_c^{pq}(\varepsilon) & \cdots \\ & \vdots & & \vdots & \\ \cdots & \hat{\underline{\tau}}_c^{qp}(\varepsilon) & \cdots & \hat{\underline{\tau}}_c^{qq}(\varepsilon) & \cdots \\ & \vdots & & \vdots & \end{pmatrix} \qquad p, q = 1, \ldots, n \tag{171}$$

quite clearly, a particular element of $\hat{\underline{\tau}}_c(\varepsilon)$,

$$\hat{\underline{\tau}}_c^{pq}(\varepsilon) = \underline{\tau}_c^{p0,\,q0}(\varepsilon) = \underline{\tau}_c^{p0,\,q0}(\varepsilon) = \frac{1}{\Omega_{BZ}} \int_{BZ} \hat{\underline{\tau}}_c^{pq}(\varepsilon, \mathbf{k}_\parallel)\, d^2 k_\parallel \tag{172}$$

refers then to the unit cells at the origin of $L_2$ in layers $p$ and $q$. Suppose now, in general, the concentration for constituents $A$ and $B$ in layer $p$ is denoted by $c_p^\alpha$ ($p = 1, \ldots, n$; $\alpha \in \{A, B\}$). By defining so-called impurity matrices, see also Eq. (159), that specify a single impurity of type $\alpha$ in the translational invariant coherent host formed by layer $p$, as

$$\hat{\underline{D}}_n^p(\varepsilon) = \underline{D}_n^{p\alpha}(\varepsilon) = \left[\underline{1} - \underline{\tau}_c^{p0,\,p0}(\varepsilon)\underline{m}_n^{p\alpha}(\varepsilon)\right]^{-1} \tag{173}$$

$$\hat{\underline{D}}_c^p(\varepsilon) = \tilde{\underline{D}}_c^{p\alpha}(\varepsilon) = \left[\underline{1} - \underline{m}_n^{p\alpha}(\varepsilon)\underline{\tau}_c^{p0,\,p0}(\varepsilon)\right]^{-1} \tag{174}$$

with

$$\underline{m}_\alpha^{p\alpha}(\varepsilon) = \underline{m}_\alpha^{p}(\varepsilon) = \widetilde{m}_\alpha^{p}(\varepsilon) = \underline{t}_\alpha^{p}(\varepsilon)^{-1} - \underline{t}_\alpha^{p}(\varepsilon)^{-1}, \quad \alpha \in \{A, B\} \tag{175}$$

where $\underline{t}_\alpha^{p}(\varepsilon)$ is the single-site $t$-matrix for constituent $\alpha$ in layer $p$, the coherent scattering path operator for the interface region, $\hat{\underline{\tau}}_L(\varepsilon)$ is obtained from the following inhomogeneous CPA condition,

$$\hat{\underline{\tau}}_c^{pp}(\varepsilon) = \sum_{\alpha=\{A,B\}} c_p^{\alpha}\langle\hat{\underline{\tau}}^{pp}(\varepsilon)\rangle_{p,\alpha} \tag{176}$$

$$\langle\hat{\underline{\tau}}^{pp}(\varepsilon)\rangle_{p,\alpha} = \hat{\underline{\tau}}_\alpha^{pp}(\varepsilon) = \hat{\underline{D}}_\alpha^{p}(\varepsilon)\hat{\underline{\tau}}_c^{pp}(\varepsilon) = \hat{\underline{\tau}}_c^{pp}(\varepsilon)\hat{\widetilde{\underline{D}}}_\alpha^{p}(\varepsilon) \tag{177}$$

that is, from a condition that implies solving *simultaneously* a layer-diagonal CPA condition for layers $p = 1, \ldots, n$. Once this condition is met then translational invariance in each layer under consideration is achieved,

$$\langle\hat{\underline{\tau}}^{pp}(\varepsilon)\rangle_{p,\alpha} \equiv \langle\underline{\tau}^{p0,p0}(\varepsilon)\rangle_{\{p0=\alpha\}} = \langle\underline{\tau}^{pi,pi}(\varepsilon)\rangle_{\{pi=\alpha\}},$$
$$\forall i \in I(L_2), \quad p = 1, \ldots, n \tag{178}$$

As before, restricted ensemble averages can be viewed as embedding an atom of type $\alpha$ into the two-dimensional translationally invariant coherent medium,

$$\langle\underline{\tau}^{p0,p0}(\varepsilon)\rangle_{\{p0=\alpha\}} = \hat{\underline{D}}_{it}^{p}(\varepsilon)\underline{\tau}_c^{p0,p0}(\varepsilon) = \underline{\tau}_c^{p0,p0}(\varepsilon)\hat{\widetilde{\underline{D}}}_{it}^{p}(\varepsilon) \tag{179}$$

Similarly, by specifying the occupation on two different sites the following restricted averages are obtained,

$$p \neq q: \langle\underline{\tau}^{pi,qj}(\varepsilon)\rangle_{\{pi=\alpha, qj=\beta\}} = \hat{\underline{D}}_{it}^{p}(\varepsilon)\underline{\tau}_c^{pi,qj}(\varepsilon)\hat{\widetilde{\underline{D}}}_{\beta}^{q}(\varepsilon) \tag{180}$$

$$p = q, \quad i \neq j: = \langle\underline{\tau}^{pi,pj}(\varepsilon)\rangle_{\{pi=n, pj=\beta\}} = \hat{\underline{D}}_{it}^{p}(\varepsilon)\underline{\tau}_c^{pi,pj}(\varepsilon)\hat{\widetilde{\underline{D}}}_{\beta}^{p}(\varepsilon) \tag{181}$$

where $\langle\underline{\tau}_c^{pi,qj}(\varepsilon)\rangle_{\{pi=\alpha, qj=\beta\}}$ has the meaning that site (subcell) $pi$ is occupied by species $\alpha$ and site (subcell) $qj$ by species $\beta$.

## 4. A PRACTICAL GREEN'S FUNCTION FORMULATION OF ELECTRIC TRANSPORT

### 4.1. Nonlocal Conductivity

A practical expression for the diagonal elements of the non-local conductivity tensor can be obtained by rewriting Eq. (99) in terms of Green's functions,

$$\sigma_{\mu\mu}^{pi,qj} = -\frac{\hbar}{4\pi V_{at}}\int_{\Omega_{pi}}d^3r_{pi}\int_{\Omega_{qj}}d^3r'_{qj}Tr(J_\mu[G^+(E_F; \mathbf{r}_{pi}, \mathbf{r}'_{qj}) - G^-(E_F; \mathbf{r}_{pi}, \mathbf{r}'_{qj})]$$

$$\times J_\mu[G^+(E_F; \mathbf{r}'_{qj}, \mathbf{r}_{pi}) - G^-(E_F; \mathbf{r}'_{qj}, \mathbf{r}_{pi})]) \tag{182}$$

where $\mu \in \{x, y, z\}$, $N_0$ is the total number of sites of a system of total volume $V = N_0V_{at}$ (assuming no lattice relaxation, thus $V_{at}$ is the same for all sites) with $G^\pm(E_F; \mathbf{r}_{pi}, \mathbf{r}'_{qj})$ referring to the up- and down-side limits of the Green's function. The integration is carried out over the $i$th unit cell in layer $p$, $\Omega_{pi}$, and the $j$-th unit cell in layer $q$, $\Omega_{qj}$ and $Tr$ denotes the trace over four-component spinors (relativistic formulation). Eq. (182), can be partitioned into four parts,

$$\sigma_{\mu\mu}^{pi,qj} = \lim_{\delta \to 0}\frac{1}{4}[\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^+, \varepsilon^+) + \tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^-, \varepsilon^-) - \tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^+, \varepsilon^-) - \tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^-, \varepsilon^+)] \tag{183}$$

each term thereof is easily expressed in terms of scattering path operators, namely

$$\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon_1, \varepsilon_2) = -\frac{\hbar}{\pi V_{oi}} tr\left[\underline{J}_{-\mu}^{pi}(\varepsilon_2, \varepsilon_1)\underline{\tau}_{clus}^{pi,qj}(\varepsilon_1)\underline{J}_{-\mu}^{qj}(\varepsilon_1, \varepsilon_2)\underline{\tau}_{clus}^{qj,pi}(\varepsilon_2)\right] \qquad (184)$$

where the underlined quantities refer to angular momentum representations and

$$\varepsilon_{1,2} = \varepsilon^{\pm} = E_F \pm i\delta \qquad (185)$$

In a relativistic formulation the current matrices are given by

$$\underline{J}_{\mu}^{pi}(\varepsilon_1, \varepsilon_2) = J_{\mu,QQ'}^{pi} = ec\int_{\Omega_{pi}} Z_Q^{pi}(\mathbf{r}_{pi}, \varepsilon_1)^+ \boldsymbol{\alpha}_\mu Z_{Q'}^{pi}(\mathbf{r}_{pi}, \varepsilon_2)d^3 r_{pi}, \qquad Q = (\kappa, \mu) \qquad (186)$$

with $\alpha_\mu$ denoting Dirac matrices, while in the non-relativistic case.

$$\underline{J}_{\mu}^{pi}(\varepsilon_1, \varepsilon_2) = J_{\mu,\Lambda\Lambda'}^{pi} = \frac{e}{m}\frac{\hbar}{i}\int_{\Omega_{pi}} Z_\Lambda^{pi}(\mathbf{r}_{pi}, \varepsilon_1)^+ \frac{\partial}{\partial r_{pi,\mu}}Z_{\Lambda'}^{pi}(\mathbf{r}_{pi}, \varepsilon_2)d^3 r_{pi}, \qquad \Lambda = (l, n) \qquad (187)$$

In the above equations, the $Z^{pi}(\mathbf{r}_{pi}, \varepsilon)$ are properly normalized regular scattering solutions of the radial Schrödinger or Dirac equation. It should be noted that in all examples shown further on exclusively relativistic current matrices have been used.

## 4.2. Nonlocal Conductivity in Disordered Systems

In order to describe substitutional binary alloys, configurational averages have to be performed in Eq. (184) [8, 19]. Omitting vertex corrections and using the single site approximation to the Coherent Potential Approximation (CPA) the site-diagonal terms are then defined as

$$\left\langle\tilde{\sigma}_{\mu\mu}^{pi,pi}(\varepsilon_1, \varepsilon_2)\right\rangle = \sum_\alpha c_\alpha tr\left[\underline{\tilde{D}}_\alpha^{pi}(\varepsilon_2)\underline{J}_{-\mu}^\alpha(\varepsilon_2, \varepsilon_1)\underline{D}_\alpha^{pi}(\varepsilon_1)\underline{\tau}_c^{pi,pi}(\varepsilon_1)\underline{J}_{-\mu}^\alpha(\varepsilon_1, \varepsilon_2)\underline{\tau}_c^{pi,pi}(\varepsilon_2)\right] \qquad (188)$$

or, by introducing the following quantity

$$\underline{\tilde{J}}_\mu^{pi,\alpha}(\varepsilon_1, \varepsilon_2) = \underline{\tilde{D}}_\alpha^{pi}(\varepsilon_1)\underline{J}_{-\mu}^\alpha(\varepsilon_1, \varepsilon_2)\underline{D}_\alpha^{pi}(\varepsilon_2) \qquad (189)$$

as

$$\left\langle\tilde{\sigma}_{\mu\mu}^{pi,pi}(\varepsilon_1, \varepsilon_2)\right\rangle = \sum_\alpha c_\alpha tr\left[\underline{\tilde{J}}_\mu^{pi,\alpha}(\varepsilon_2, \varepsilon_1)\underline{\tau}_c^{pi,pi}(\varepsilon_1)\underline{J}_{-\mu}^\alpha(\varepsilon_1, \varepsilon_2)\underline{\tau}_c^{pi,pi}(\varepsilon_2)\right] \qquad (190)$$

where $c_\alpha$ denotes the (homogeneous) concentration of the $\alpha$-th component. $\alpha \in \{A, B\}$, of a binary *bulk* alloy and the current matrix $\underline{J}_{-\mu}^\alpha$ refers to species $\alpha$.

For the off-site diagonal case, $(pi) \neq (qj)$, this kind of approach yields

$$\left\langle\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon_1, \varepsilon_2)\right\rangle = \sum_{\alpha,\beta} c_\alpha c_\beta tr\left[\underline{\tilde{D}}_\alpha^{pi}(\varepsilon_2)\underline{J}_{-\mu}^\alpha(\varepsilon_2, \varepsilon_1)\underline{D}_\alpha^{pi}(\varepsilon_1)\underline{\tau}_c^{pi,qj}(\varepsilon_1)\right.$$

$$\left. \times \underline{\tilde{D}}_\beta^{qj}(\varepsilon_1)\underline{J}_{-\mu}^\beta(\varepsilon_1, \varepsilon_2)\underline{D}_\beta^{qj}(\varepsilon_2)\underline{\tau}_c^{qj,pi}(\varepsilon_2)\right] \qquad (191)$$

or, in using Eq. (189),

$$\left\langle\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon_1, \varepsilon_2)\right\rangle = \sum_{\alpha,\beta} c_\alpha c_\beta tr\left[\underline{\tilde{J}}_\mu^{pi,\alpha}(\varepsilon_2, \varepsilon_1)\underline{\tau}_c^{pi,qj}(\varepsilon_1)\underline{\tilde{J}}_\mu^{qj,\beta}(\varepsilon_1, \varepsilon_2)\underline{\tau}_c^{qj,pi}(\varepsilon_2)\right] \qquad (192)$$

## 4.3. The "Large Cluster" Limit

If only unperturbed host atoms form the cluster then by increasing the size of the cluster the physical properties characteristic for the corresponding bulk or substrate have to be expected. A rigorous test for a *"real space formulation"* of the Kubo equation consists therefore in considering the following convergence procedure for the diagonal elements of the resistivity $\rho$

$$\rho_{\mu\mu} = \lim_{\delta \to 0} \rho_{\mu\mu}(r_0; \delta), \quad \rho_{\mu\mu}(r_0; \delta) = \lim_{r \to r_0} \rho_{\mu\mu}(r; \delta) \tag{193}$$

$$\rho_{\mu\mu}(r; \delta) = [\sigma^0_{\mu\mu}(r; \delta)]^{-1}, \quad \sigma^0_{\mu\mu}(r; \delta) = \sum_j \sigma^{0j}_{\mu\mu}(\delta) \tag{194}$$

where $r$ denotes the radius of a sphere with the origin in site $i = 0$ and $r_0$ is an arbitrarily large radius. The summation in Eq. (194) extends over all sites circumscribed by $r$; $\delta$ refers to the imaginary part of the Fermi energy. In performing the $\delta \to 0$ limit at the stage of Eq. (193) actually means that the side limits in Eq. (183) are taken at the last possible step.

Shrinking the sphere to a circle within the plane of a specific layer $p$ (e.g., surface layer), then in the $r \to \infty$ limit the following condition must be satisfied

$$\lim_{r \to \infty} \rho_{\mu\mu}(r, \delta) = \hat{\rho}^{pp}_{\mu\mu}(\delta), \quad \hat{\rho}^{pp}_{\mu\mu}(\delta) = \left[\hat{\sigma}^{pp}_{\mu\mu}(\delta)\right]^{-1} \tag{195}$$

where $\hat{\sigma}^{pp}_{\mu\mu}$ is the layer-diagonal conductivity of layer $p$,

$$\hat{\sigma}^{pp}_{\mu\mu}(\delta) = \frac{1}{4}\left[\tilde{\sigma}^{pp}_{\mu\mu}(\varepsilon^+, \varepsilon^+) + \tilde{\sigma}^{pp}_{\mu\mu}(\varepsilon^-, \varepsilon^-) - \tilde{\sigma}^{pp}_{\mu\mu}(\varepsilon^+, \varepsilon^-) - \tilde{\sigma}^{pp}_{\mu\mu}(\varepsilon^-, \varepsilon^+)\right] \tag{196}$$

for which each term on the rhs of the last equation can also be calculated directly using a two-dimensional lattice Fourier transformation,

$$\tilde{\sigma}^{pp}_{\mu\mu}(\varepsilon_1, \varepsilon_2) = -\frac{\hbar}{\pi V_{at}} \frac{1}{\Omega_{BZ}} \int_{BZ} d^2k_\parallel \, tr\left[\underline{J}^p_\mu(\varepsilon_2, \varepsilon_1)\underline{\tau}^{pp}(\varepsilon_1, \mathbf{k}_\parallel)\underline{J}^p_\mu(\varepsilon_1, \varepsilon_2)\underline{\tau}^{pp}(\varepsilon_2, \mathbf{k}_\parallel)\right] \tag{197}$$

Using a sphere of radius $r$ the summation must provide in the $r \to \infty$ limit the total conductivity of the bulk system,

$$\sigma^{total}_{\mu\mu} = \lim_{\delta \to 0} \sigma^{total}_{\mu\mu}(\delta) = \lim_{\delta \to 0}\left[\lim_{r \to \infty} \sigma^0_{\mu\mu}(r, \delta)\right] \tag{198}$$

Inverting $\sigma^{total}_{\mu\mu}$, the resistivity of a bulk system is obtained, which is zero for pure metals and finite for (disordered) alloy bulk systems (the so-called residual resistivity).

Quite clearly there are more efficient methods to evaluate resistivities for bulk or layered systems by making use of three- or two-dimensional lattice Fourier transformations, respectively. However, once it comes to determine, e.g., the electric properties of magnetic islands on surfaces, these methods are no longer applicable, and one has to rely on real space approaches as presented in here. It should be noted that the results in the "large cluster" limit presented later on are only illustrations of the reliability and applicability of the real space approach to the Kubo equation.

## 4.4. "Residual Resistivity" for Nanostructures

If no translational symmetry is present, then in principle one has to sum over all sites including the leads, contacts, and so forth; that is,

$$\tilde{\sigma}_{\mu\mu}(\varepsilon_1, \varepsilon_2) = \frac{1}{N_0}\sum_{i=1}^{N_0}\sum_{j=1}^{N_0} \tilde{\sigma}^{ij}_{\mu\mu}(\varepsilon_1, \varepsilon_2) \tag{199}$$

with $N_0 \approx 10^{23}$. Here, $i$ and $j$ denote sites without labelling layers explicitly. As such a procedure is numerically not accessible, the following quantity can be defined,

$$\tilde{\sigma}_{\mu\mu}(\varepsilon_1, \varepsilon_2; n) = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \tilde{\sigma}_{\mu\mu}^{ij}(\varepsilon_1, \varepsilon_2) \tag{200}$$

with $n$ being the number of sites in a chosen region ("cluster"). This implies, however, that the convergence properties of $\tilde{\sigma}_{\mu\mu}(\varepsilon_1, \varepsilon_2; n)$ with respect to $n$ have to be investigated. Because clearly enough a summation over all sites including the semi-infinite substrate would yield only the resistivity of the substrate, namely zero in the case of a pure metal, a kind of "residual resistivity" for finite clusters has to be defined,

$$\rho_{\mu\mu}^{n}(r) = \left[ \frac{1}{n} \sum_{i \in \text{chain}} \sum_{j=1}^{N(r)} \sigma_{\mu\mu}^{ij} \right]^{-1} \tag{201}$$

where, for example, in the case of a finite chain embedded in the surface of a suitable substrate $n$ denotes the number of atoms in the chain of type $\alpha$ and $N(r)$ is the number of atoms involved in the cluster (chain + substrate neighbourhood up to a chosen value $r$ for the circumscribing sphere).

## 4.5. Conductances

Linear response theory applies to an arbitrary choice for the perturbating electric field because the response function is obtained in the limit of a vanishing perturbation. Consider that a constant electric field, $E_z^q$, pointing along the $z$ axis, that is, normal to the planes, is applied in all cells of layer $q$. Denoting the $z$ component of the current density averaged over cell $i$ in layer $p$ by $j_z^{pi}$, the microscopic Ohm's law reads as

$$j_z^{pi} = \frac{1}{V_{at}} \sum_j \sigma_{zz}^{pi,qj} E_z^q \tag{202}$$

where $V_{at}$ is the volume of the unit cell in layer $p$. Note, that in neglecting lattice relaxations, $V_{at}$ is uniform for the whole system. According to the Kubo-Greenwood equation at zero temperature, see Eq. (99), the $zz$ component of the non-local conductivity tensor $\sigma_{zz}^{pi,qj}$ can be written as

$$\sigma_{zz}^{pi,qj} = -\frac{\hbar}{4\pi} \int_{\Omega_{pi}} d^3 r_{pi} \int_{\Omega_{qj}} d^3 r_{qj}' \, Tr\{J_z[G^+(E_F; \mathbf{r}_{pi}, \mathbf{r}_{qj}') - G^-(E_F; \mathbf{r}_{pi}, \mathbf{r}_{qj}')]$$

$$\times J_z[G^+(E_F; \mathbf{r}_{qj}', \mathbf{r}_{pi}) - G^-(E_F; \mathbf{r}_{qj}', \mathbf{r}_{pi})]\} \tag{203}$$

However, the total current flowing through layer $p$ can also be written as

$$I_{tot} = A_{\parallel} \sum_i j_z^{pi} = gU \tag{204}$$

where the summation has to be carried out for all sites in layer $p$ and the applied voltage $U$ is given by

$$U = E_z^q d \tag{205}$$

with $A_{\parallel}$ and $d$ denoting the area of the two-dimensional unit cell and the interlayer spacing, respectively ($V_{at} = A_{\parallel} d$). Combining Eqs. (202), (204) and (205) results in an expression for the conductance,

$$g = \frac{1}{d^2} \sum_i \sum_j \sigma_{zz}^{pi,qj} \tag{206}$$

where the summations should, in principle, be carried out over all cells in layers $p$ and $q$.

An alternative choice for the nonlocal conductivity tensor was given in Eq. (102), which is more practical in calculating the CPP conductance of a layered system than the nonlocal conductivity in the Kubo-Greenwood approach because, as shown by Baranger and Stone [9] for free electron leads, the second term appearing in Eq. (102) becomes identically zero when integrated over layers, $p \neq q$. This means also that the terms $\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^-,\varepsilon^-)$ and $\tilde{\sigma}_{\mu\mu}^{pi,qj}(\varepsilon^-,\varepsilon^-)$ should vanish after integration. It should be noted that very recently Mavropoulos et al. [20] rederived this result by assuming Bloch boundary conditions for the leads. According to these theoretical results, at zero temperature the diagonal elements of the nonlocal conductivity tensor between site $i$ in layer $p$ and site $j$ in layer $q$ can be written as

$$\sigma_{zz}^{pi,qj} = -\frac{1}{2}\tilde{\sigma}_{zz}^{pi,qj}(\varepsilon^-,\varepsilon^-) = \frac{\hbar}{2\pi}\int_{\Omega_{pi}} d^3r_{pi}\int_{\Omega_{qj}} d^3r'_{qj} Tr[J_z G^+(E_F;\mathbf{r}_{pi},\mathbf{r}'_{qj})J_z G^-(E_F;\mathbf{r}'_{qj},\mathbf{r}_{pi})]$$

$$= \frac{\hbar}{2\pi}tr[\underline{J}_z^{pi}(\varepsilon^-,\varepsilon^+)\underline{\tau}_{z,\text{clus}}^{pi,qj}(\varepsilon^+)\underline{J}_z^{qj}(\varepsilon^+,\varepsilon^-)\underline{\tau}_{z,\text{clus}}^{qj,pi}(\varepsilon^-)] \tag{207}$$

such that the expression for the conductance reduces to

$$g = \frac{\hbar}{2\pi d_\perp^2}\sum_i\sum_j\int_{\Omega_{pi}} d^3r_{pi}\int_{\Omega_{qj}} d^3r'_{qj} Tr[J_z G^+(E_F;\mathbf{r}_{pi},\mathbf{r}'_{qj})J_z G^-(E_F;\mathbf{r}'_{qj},\mathbf{r}_{pi})] \tag{208}$$

or, in terms of scattering path operators to

$$g = \frac{\hbar}{2\pi d_\perp^2}\lim_{\delta\to 0}\sum_i\sum_j tr[\underline{J}_z^{pi}(\varepsilon^-,\varepsilon^+)\underline{\tau}_{z,\text{clus}}^{pi,qj}(\varepsilon^+)\underline{J}_z^{qj}(\varepsilon^+,\varepsilon^-)\underline{\tau}_{z,\text{clus}}^{qj,pi}(\varepsilon^-)] \tag{209}$$

It has to be emphasized that because of the use of linear response theory and current conservation, the choice of layers $p$ and $q$ is arbitrary in Eqs. (208) and (209). If the layers $p$ and $q$ are asymptotically far away from each other, the above expressions naturally recover [20] the Landauer-Büttiker approach [3, 4], see section 2.2.

# 5. THE "LARGE CLUSTER" LIMIT

A real space version of Eq. (99) allows to study the interesting transition of electric transport properties from nanoscaled to macroscopic (mesoscopic) systems, simply by increasing the number of atomic sites included in the summation over sites. However, by doing so, such a procedure can also be used to document the numerical accuracy that can be achieved with a real space approach, since in the limit of two- or three-dimensional translational invariance corresponding theoretical results are available, obtained using appropriate lattice Fourier transformations. The following few sections serve exactly this purpose, namely to illustrate the convergence to semi-infinite and infinite systems.

## 5.1. Surface Layer of Ag(001)

The system studied is sketched in Fig. 2. The underlying parent lattice is an fcc structure corresponding to the experimental lattice spacing of fcc-Ag: $a_{3D} = 7.789$ a.u. and $a_{2D} = 5.508$ a.u.



Figure 2. Geometrical setup of a semi-infinite Ag(001) system.

$\sigma_{xx}^{0j}[(m\Omega cm)^{-1}]$

$\sigma_{zz}^{0j}[(m\Omega cm)^{-1}]$

**Figure 3.** Nonlocal conductivity $\sigma_{xx}^{0j}(x_j, y_j)$ (left) and $\sigma_{zz}^{0j}(x_j, y_j)$ (right) corresponding to the surface layer of Ag. $\delta = 1$ mRy.

The nonlocal conductivities in the surface layer were calculated according to Eqs. (183) and (184) with the real space scattering path operators being obtained by using 1830 $k_\parallel$ points in the two-dimensional irreducible wedge of the surface Brillouin zone. In Fig. 3, the $xx$ and $zz$ components of the nonlocal conductivity tensor $\sigma_{\mu\nu}^{0j}(x_j, y_j)$ are shown, where site 0 is fixed to the origin (0,0) of the surface layer, while the position of sites $j$ is varied in the (001)-oriented surface plane. As can be seen, for the out of plane conductivity ($zz$), only scatterers are important which are not too far away from the origin, while in the in-plane case ($xx$) also scatterers at farther distances do add non-negligible contributions to the corresponding component of the conductivity. Moreover, it should be noted that the $yy$ component is not shown because it is of similar form as the $xx$ component: the diagram of $\sigma_{xx}^{0j}(x_j, y_j)$ simply has to be rotated by 90°.

The shape of the non-local conductivities suggests that by performing the summation in Eq. (194) only for sites within the chosen plane of atoms, it should converge according to Eq. (195). In order to show that, two different square-shaped planar clusters were investigated, see Fig. 4, both having $C_{4v}$ symmetry which implies that $\rho_{\mu\mu}(r; \delta)$ has two independent components, namely

$$\rho_{xx} = \rho_{yy} \quad \text{and} \quad \rho_{zz} \tag{210}$$

The characteristic size ($r$) of the investigated clusters is given by the distance between the origin (0) and the farthermost atom from the origin, that is, can be viewed with respect to increasing sizes and fixed shape, namely in terms of $r_n = n \cdot a_{2D}$ for type 1 and $r_n = n\sqrt{2} \cdot a_{2D}$ for type 2, see Fig. 4. The number of atoms within a particular cluster is given by $N(n) = (2n^2 + 2n + 1)$ for type 1 and $N(n) = (4n^2 + 4n + 1)$ for type 2. Obviously, the clusters shown

square type 1

square type 2

**Figure 4.** Square- and diamond-like shapes of a cluster in the surface plane of an fcc(001) substrate.

**Figure 5.** Convergence study in the surface layer of an Ag(001) semi-infinite system. The in-plane ($xx$) and perpendicular to the plane ($zz$) resistivity components for two different cluster shapes are shown versus the characteristic size of the cluster ($r$). The horizontal line refers to the layer-diagonal resistivity calculated by Eqs. (195)–(197). Diamonds correspond to type 1 in Fig. 4, squares to type 2, $\delta = 1$ mRy. Reprinted with permission from [33], K. Palotás et al., *Phys. Rev. B* 67, 174404 (2003). © 2003, American Physical Society.

in Fig. 4 refer to $n = 3$. It has to be emphasized that this procedure is to show the validity of Eq. (195): as can be seen from Fig. 5, for both types of clusters a reliable convergence of the resistivity is achieved for $r > 15\ a_{2D}$.

## 5.2. Bulk Resistivities

The studied systems are summarized in Fig. 6. The nonlocal conductivities were calculated according to Eq. (183) and the side limits in Eq. (184) for Ag and in Eqs. (190), (192) for CuPt alloys with the real space scattering path operators being obtained by using 630 $k_\parallel$ points in the two-dimensional irreducible wedge of the surface Brillouin zone. In the following, three-dimensional clusters are assumed; the real space summation of the non-local conductivity tensor was performed according to Eq. (194). In addition by increasing the size of the clusters the convergence of Eq. (198) was studied and the obtained results were compared to known bulk resistivities, see Refs. [21, 22]. Clearly for large clusters the resistivity has to approach to the corresponding bulk value, namely to zero for pure metals and to the residual resistivity for (disordered) alloys. The clusters were chosen to be contained by a sphere of increasing radius; the origin of the spheres refers to the site denoted by 0 in Eq. (194). Table 1 shows the number of atoms ($N$) involved within a sphere with respect to $n$; the corresponding sphere radius $r_n$ is defined by

$$r_n = \frac{n}{\sqrt{2}} \cdot a_{3D}$$

Assuming the following behavior of the elements of the resistivity tensor with respect to the size of the cluster ($r$),

$$\rho_{\mu\mu}(r; \delta) = \rho_0(\delta) + \frac{\rho_1(\delta)}{r} \tag{211}$$



**Figure 6.** Geometrical setup of the Ag and CuPt bulk systems.

Table 1. The number of sites ($N$) in clusters of spherical shape.

| $n$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Cluster | | | | |
| $N(n)$ | 1 | 13 | 55 | 177 | 381 | 767 | 1289 | 2093 | 3055 | 4321 | 5979 | 7935 |

$\rho_0$ and $\rho_1$ being constants, it is obvious that

$$r\rho_{\mu\mu}(r; \delta) = r\rho_0(\delta) + \rho_1(\delta) \tag{212}$$

which means that the residual resistivity, $\rho_0(\delta)$ can be obtained by a linear fit of $r\rho_{\mu\mu}(r; \delta)$ with respect to $r$. In the case of substitutional alloys, the slope $(\rho_0(\delta), \delta \to 0)$ corresponds then to the residual resistivity, while for a pure bulk it should be zero. It should be noted that Eq. (211) is more or less an empirical finding which was used also quite a bit also in the experimental recording of resistivities.



Figure 7. Convergence study in bulk systems. The characteristic size of the cluster ($r$) times the resistiviy is shown versus the size of the cluster for three different imaginary parts ($\delta$) of the Fermi energy in order to evaluate the slope (residual resistivity), see Eq. (212). The $zz$ component of the resistivity is shown for fcc bulk Ag (top), and $xx$ and $zz$ components for fcc bulk CuPt alloys (bottom). Reprinted with permission from [33], K. Palotas et al., *Phys. Rev. B* 67, 174404 (2003). © 2003, American Physical Society.

## 5.2.1. Ag Bulk

The fcc bulk Ag structure has the same lattice constants as mentioned in Section 5.1. In principle it is sufficient to evaluate only one component of the resistivity because the system and also the clusters have cubic symmetry, which means that by choosing the coordinate system properly, the resistivity tensor has only one independent element, that is, the diagonal components must be identical ($\rho_{xx} = \rho_{yy} = \rho_{zz}$). Deviations from this behavior can be used to estimate numerical errors inherent to the calculational scheme and the fitting procedure. The actual fitting, see Eq. (212), was performed for each calculated value of $\delta$ ($\delta = 1, 2, 3$ mRy) considering the last three points of $r\rho_{zz}(r; \delta)$, see top part of Fig. 7. These points have been chosen because they refer to the biggest clusters considered, see Table 1. In order to obtain the real physical residual resistivity an extrapolation to $\delta = 0$ is needed, see Eq. (193). This extrapolation for Ag bulk structure is illustrated in the top part of Fig. 8 and demonstrates that an absolute error of roughly $0.05$ $\mu\Omega$cm was made in the applied fitting procedure.

## 5.2.2. $Cu_cPt_{1-c}$ Bulk

More interesting than pure bulk metals are disordered bulk alloys because the accuracy of the current approach can be directly compared with experimental data and results of previous calculations using three-dimensional periodic boundary conditions. For this reason, fcc $Cu_{0.50}Pt_{0.50}$ and $Cu_{0.75}Pt_{0.25}$ have been chosen with lattice constants $a_{3D}^{Cu_{0.50}Pt_{0.50}} = 7.140$ a.u. and $a_{3D}^{Cu_{0.75}Pt_{0.25}} = 6.995$ a.u. in order to test the reliability of the present approach. Again the fitting to a linear form to the last three points of $r\rho_{\mu\mu}(r)$ has been applied, see Eq. (212), as a function of $\delta$, see the bottom part of Fig. 7. As can be seen, the extrapolation can easily be performed because in the region $0 < \delta < 3$ mRy the resistivity depends linearly on $\delta$. In comparing the present results with previous calculations and available experimental data, see in particular Ref. [21], good quantitative agreement for both concentrations of CuPt is found: the results of Dulca et al. [22], for example, are 80.2 and 31.5 $\mu\Omega$cm for $Cu_{0.50}Pt_{0.50}$ and $Cu_{0.75}Pt_{0.25}$, respectively.

As already stated the numerical errors of the present approach can be judged best by determined by evaluating the difference between the in-plane and the perpendicular to the



Figure 8. Extrapolation to $\delta = 0$ for the investigated bulk systems. Open circles are obtained from the fitting procedure in Eq. (212), while full circles refer to the extrapolated values. Squares denote experimental results measured at room temperature [21]. Reprinted with permission from [33], K. Palotas et al., *Phys. Rev. B* 67, 174404 (2003). © 2003, American Physical Society.

**Figure 9.** Difference between the residual resistivity for the in-plane ($xx$) and the perpendicular to the plane ($zz$) component *versus* the imaginary part ($\delta$) of the Fermi energy for $Cu_{0.50}Pt_{0.50}$ and $Cu_{0.75}Pt_{0.25}$. Reprinted with permission from [33], K. Palotas et al., *Phys. Rev. B* 67, 174404 (2003). © 2003, American Physical Society.

plane elements of the residual resistivity, since the residual resistivities, $\rho_{xx}$ and $\rho_{zz}$ must be identical in cubic bulk systems. It can be seen from Fig. 9 that this difference is more or less independent of $\delta$ and is of order of a few tenth of a $\mu\Omega$cm.

## 6. MAGNETIC FINITE CHAINS IN THE SURFACE OF Ag(001)

In this section, single impurities and finite chains (length of 2–10 atoms) of Fe and Co embedded along the (110) direction ($x$) in the surface layer of Ag(001) are investigated, see also Fig. 2. In here, for matters of simplicity, a simple notation for the embedded chains is used, namely for example Co$_4$ for a Co chain of four atoms. For such a chain of our atoms the $y = 0$ plane-section of the system is shown in Fig. 10.

### 6.1. Nonlocal Conductivities

The influence of the chains to the in-plane transport in the surface layer was investigated by assuming a CIP geometry. The nonlocal conductivities were calculated according to Eq. (183), the scattering path operators of a specific cluster have been obtained n terms of the embedding equation, see Eq. (157). The real space host scattering path operators were calculated by using 210 $k_{\parallel}$ points in the two-dimensional irreducible wedge of the surface Brillouin zone.

Let 0 denote the origin in the surface plane ($x = 0$, $y = 0$), to which in the impurity case a single impurity is fixed. The $xx$-component of the nonlocal conductivity tensor between this site (0) and all the other atoms in the surface plane is shown in Fig. 11. As can be seen from this figure, the site-diagonal conductivity component of this site is larger for Co than for Fe, causing in turn of a higher resistivity of Fe chains after performing the summation in Eq. (213).

For chains the atom at the edge of a chain serves as origin, that is. it is located ir the origin of the surface. The $xx$-component of the nonlocal conductivity tensor between ths fixed site



**Figure 10.** Chain of four atoms in the surface layer of Ag(001), $y = 0$ plane-section.

Figure 11. Nonlocal conductivities $\sigma_{xx}^{0j}(x_j, y_j)$ in the surface layer of Ag in presence of Co or Fe impurity at site 0 or without any impurities (pure Ag surface). $\delta = 1$ mRy, $\tilde{M} = 3$.

and all other atoms in the surface plane is shown in Fig. 12 for Co and in Fig. 13 for Fe. It can be seen that in the impurity case, the shape of the conductivity is symmetric to the $x = 0$ plane, whereas in the case of finite chains the tensor-elements along the $+x$ direction (where the chain lies) are much larger than in other directions, causing thus an asymmetry. This shape also implies that by summing up the nonlocal conductivity $\sigma_{xx}^{0j}$ over sites $j$ in a three-dimensional cluster around the chain, a significant contribution arises from the magnetic atoms. For a Co chain with length of six atoms, for example, the contribution from the magnetic atoms amounts to about 63%. Furthermore, it can also be seen that the magnitude of the site-diagonal conductivities decrease for atoms forming a chain as compared to the corresponding single impurity.

## 6.2. "Residual Resistivities"

In Section 4.4, a "residual resistivity" for finite clusters was defined as

$$\rho_{\mu\mu}^n(r) = \left[ \frac{1}{n} \sum_{i \in \text{chain}}^{} \sum_{j=1}^{N(r)} \sigma_{\mu\mu}^{ij} \right]^{-1} \tag{213}$$

where $n$ denotes the number of atoms in the chain of type $\alpha$ (Fe or Co), and $N(r)$ is the number of atoms involved in the cluster (chain + environmental atoms up to the furthermost distance of $r$). It should be noted that for evaluating Eq. (213) three-dimensional clusters have to be used. Obviously, the convergence properties of $\rho_{\mu\mu}^\alpha(r)$ with respect to $r$ can be investigated by increasing size of the cluster. This is shown in Fig. 14. As can be seen in this figure $\rho_{xx}^\alpha(r)$ decreases for all chain lengths ($n$) monotonously and can in principle be extrapolated to large values of $N(r)$, see Eq. (212), while the difference, $\rho_{xx}^{Fe}(r) - \rho_{xx}^{Co}(r)$ remains finite and varies only slowly with respect to the cluster size. Furthermore, chains with length of three or five atoms differ distinctly from the rest, namely there is almost no difference whether Fe or Co atoms form the chain, i.e., the difference, $\rho_{xx}^{Fe}(r) - \rho_{xx}^{Co}(r)$ nearly vanishes for all cluster size considered.

The "residual resistivity" of finite clusters defined in Eq. (213) is a practical tool to study the influence of in-plane transport properties with respect to the orientation of magnetization

Figure 12. Nonlocal conductivities $\sigma_{xx}^{0j}(r, r')$ in the surface layer of Ag with Co atoms in a $Co_n$ chain in positions $(0,0), \ldots, (n-1, 0)$. $\delta = 1$ mRy. $\widetilde{M} = \hat{z}$.

($\widetilde{M}$). The calculated results of the $xx$-component of the resistivity are listed in Table 2. As can be seen, $\widetilde{M} = \hat{x}$, i.e., $\widetilde{M}$ parallel to the orientation of the chains provides the smallest resistivity for all Fe chains and for the most Co chains. There are two exceptions where this does not apply, namely $Co_2$ and $Co_3$. For these chains the smallest resistivity

**Figure 13.** Nonlocal conductivities $\sigma_{xx}^{0j}(x_j, y_j)$ in the surface layer of Ag with Fe atoms in a $Fe_n$ chain in positions $(0,0), \ldots, (n-1, 0)$. $\delta = 1$ mRy. $\hat{\mathbf{M}} = \hat{z}$.

is obtained for $\hat{\mathbf{M}} = \hat{y}$. This behavior is quite surprising in view of the resistivities for the other chains. In most cases the direction of magnetization $\hat{\mathbf{M}} = \hat{y}$ seems to yield the highest resistivity, however, the orientation of magnetization perpendicular to the chain ($\hat{y}$ and $\hat{z}$) results in minor differences in the resistivity. Moreover, in the impurity case, $\rho_{xx}^{Fe}(\hat{y})$ is

**Figure 14.** "Residual resistivities" of Fe (circles) and Co (triangles) chains. Open squares refer to $\rho_{1,1}^{Fe}(r) - \rho_{1,1}^{Co}(r)$. The length of the chains ($n$) is explicitly shown. $\delta = 0$ mRy (extrapolated). $M = \hat{z}$.

Table 2. "Residual resistivities" *versus* orientation of magnetization ($\hat{M}$). $\rho_{xx}(r = a_{3/)})[\mu\Omega\text{cm}]$ in Co and Fe chains. $\delta = 0$ mRy (extrapolated).

| Length | Co | | | Fe | | |
|--------|-------|-------|-------|-------|-------|-------|
|        | $\dot{x}$ | $\dot{y}$ | $\dot{z}$ | $\dot{x}$ | $\dot{y}$ | $\dot{z}$ |
| 1  | 120.1 | 126.9 | 123.5 | 198.3 | 225.3 | 219.6 |
| 2  | 162.6 | 162.2 | 165.4 | 200.5 | 213.1 | 213.2 |
| 3  | 151.4 | 140.1 | 143.9 | 142.0 | 153.8 | 152.0 |
| 4  | 109.0 | 113.7 | 113.2 | 166.7 | 176.4 | 176.2 |
| 5  | 122.6 | 126.6 | 128.8 | 122.2 | 130.1 | 129.7 |
| 6  | 94.9  | 100.1 | 97.0  | 132.3 | 138.6 | 137.6 |
| 7  | 85.6  | 89.7  | 88.3  | 119.5 | 125.9 | 127.3 |
| 8  | 86.1  | 89.4  | 87.8  | 108.0 | 111.9 | 111.7 |
| 9  | 74.0  | 78.6  | 77.6  | 110.2 | 114.7 | 117.1 |
| 10 | 73.4  | 77.2  | 74.9  | 91.3  | 93.4  | 94.9  |

by 13.6% larger than $\rho_{xx}^{Fe}(\dot{x})$, whereas $\rho_{xx}^{Co}(\dot{y})$ is only by 5.7% larger than $\rho_{xx}^{Co}(\dot{x})$, which means a higher sensitivity with respect to the orientation of the magnetization for the Fe impurity.

## 7. NANOCONTACTS

Nanocontacts made of gold are presumably the most studied atomic-sized conductors in the literature. A dominant peak very close to the conductance quantum, $1 \ G_0 = 2e^2/h$, has been reported for gold in the conductance histogram [23, 24] and attributed to the highly transmitting $sp$ channel across a linear monoatomic chain connecting the two electrodes. In this section, gold contacts are investigated in different geometries as well as the influence of transition metal impurities on the conductance is studied within the real-space approach described in Sections 4.1 and 4.5.

The host system for the embedding is shown in Fig. 15. It should be noted that all of the considered sites (Au, vacuum and impurities) refer to the positions of an underlying ideal fcc structure of gold with a lattice constant of $a_{3D} = 7.681$ a.u.

A schematic view of a typical contact is displayed in Fig. 16 with $N_1 = 5$ vacuum layers considered in the host system, see Fig. 15. As follows from the above, atomic sites refer to layers for which the following notation is used: $C$ denotes the *central layer*, $C - 1$ and $C + 1$ the layers below and above, and so forth. The contact consists of a central layer that contains 1 Au atom (the rest is built up from empty spheres); layers $C - 1$ and $C + 1$, see Fig. 17a, contain 4 Au atoms, layers $C - 2$ and $C + 2$ 9 Au atoms, and, though not shown, all other layers, namely $C - n$ and $C + n$ ($n \geq 3$), are completely filled with Au atoms, that is, denote full layers. The nonlocal conductivities were calculated according to Eq. (207), the scattering path operators of a specific cluster were obtained by the embedding equation, Eq. (157). The real space host scattering path operators were calculated by taking 210 $k_\parallel$ points in the two-dimensional irreducible wedge of the surface Brillouin zone.



Figure 15. Geometrical setup of Au(001) host system. A nanojunction between the two semi-infinite systems is modeled by embedding Au atoms into the vacuum region, see, for example, Figs. 16 and 17. The host characterized by $N_{Au}$ and $N_V$ sites that can be different for different contacts.

**Figure 16.** Schematic side view of a point contact between two semi-infinite leads embedded into the vacuum region (number of vacuum layers $N_i$ = 5). The layers are labeled by $C$, $C \pm 1$, and so forth. Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

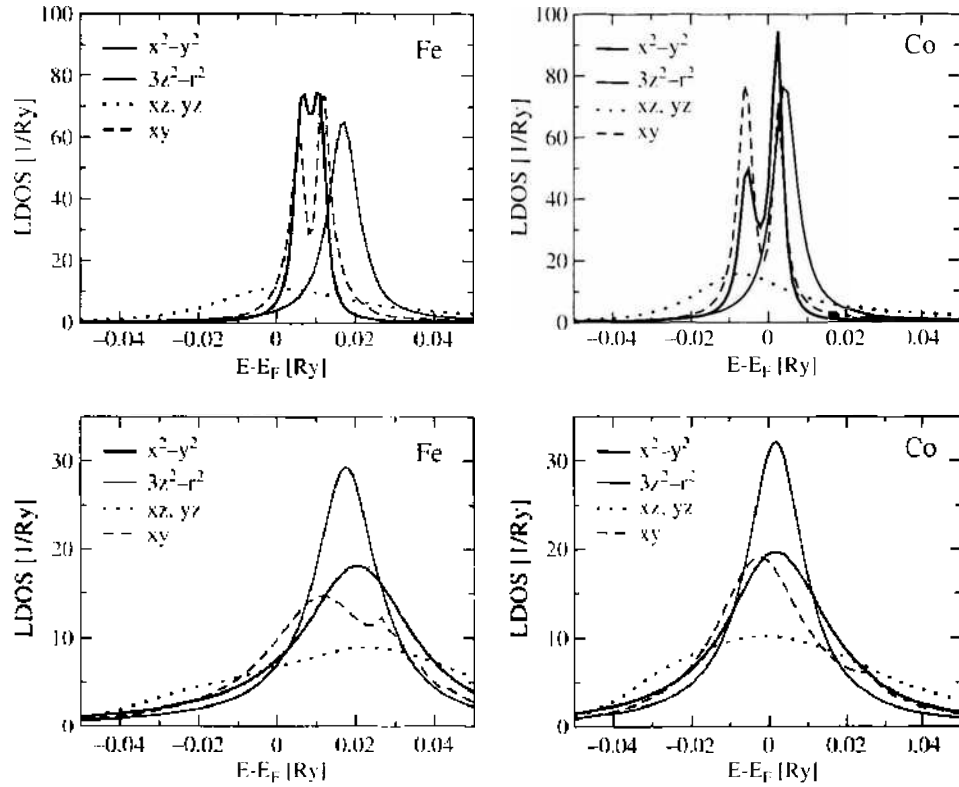## 7.1. Numerical Tests for Different Gold Contacts

As mentioned in Section 4.1, a finite Fermi level broadening, $\delta$, has to be used for the nonlocal conductivity, thus, also for conductance calculations. As an example, for the point contact depicted in Fig. 17a, the dependence of the conductance on $\delta$ is investigated. The summation in Eq. (209) was carried out up to convergence for the first two (symmetric) full layers ($p = C - 3$, $q = C + 3$). As can be seen from Fig. 18, the calculated conductances depend strongly but nearly linear on $\delta$. A straight line fitted for $\delta \geq 1.5$ mRy intersects the vertical axis at 2.38 $G_0$. Assuringly enough, a calculation with $\delta = 1$ $\mu$Ry resulted in $g = 2.40 G_0$. Although the nearly linear dependence of the conductance with respect to $\delta$ enables an easy extrapolation to $\delta = 0$, in the following all reported conductances refer to $\delta = 1$ $\mu$Ry.

For the same type of contact (Fig. 17a), the convergence of the summation in Eq. (209) over layers $p$ and $q$ was investigated, choosing different symmetric pairs of full layers.



**Figure 17.** Perspective view of some of the studied contacts between two fcc(001) semi-infinite leads. Only the partially filled layers are shown. (a) point contact (number of Au layers taken into account: $N_{lu}$=5, number of vacuum layers between the leads: $N_i$ = 5. (b) Slanted linear finite chain ($N_{au}$ = 7, $N_i$ = 7). (c) 2 × 2 finite chain ($N_{lu}$ = 6, $N_i$ = 9). Reprinted with permission from [34]. K. Palotas et al. *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

**Figure 18.** Calculated conductance as a function of the Fermi level broadening $\delta$ for the Au contact shown in Fig. 17a. The dashed straight line is a linear fit to the values for $\delta = 1.5, 2.0, 2.5$, and 3.0 mRy. Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

The convergence with respect to the number of atoms in the layers is shown in Fig. 19. Convergence for about 20 atoms can be obtained for the first two full layers ($p = C - 3$, $q = C + 3$), whereas the number of sites needed to get convergent sums gradually increases if one includes layers farther away from the contact atom. This kind of convergence property is qualitatively understandable, as the current flows from the contact within a cone of some opening angle that cuts out sheets of increasing area from the corresponding layers. As all the calculations were performed with $\delta = 1$ $\mu$Ry, current conservation has to be expected. Consequently, the calculated conductance ought to be independent with respect to the layers chosen for the summation in Eq. (209). As can be seen from Fig. 19 this is satisfied within a relative error of less than 10%. It should be noted, however, that for the pairs of layers, $p = C - n$, $q = C + n$, $n \geq 6$ convergence was not achieved within this accuracy: by taking more sites in the summations even a better coincidence of the calculated conductance values for different pairs of layers can be expected. Figure 19 also implies that an application of the Landauer-Büttiker approach to calculate the conductance of nanocontacts is numerically more tedious than the present one, since, in principle, two layers situated infinitely far from each other have to be taken in order to represent the leads.

Although only one Au atom is placed in the center of the point contact considered above, see Fig. 17a, the calculated conductance is more than twice as large as the conductance unit. This is easy to understand since the planes $C - 1$ and $C + 1$, each containing four Au atoms,



**Figure 19.** Conductance *versus* the number of sites included in the sum in Eq. (209) for the contact in Fig. 17a. The different curves show conductances as calculated between different pairs of layers. For a definition of the layer numbering see Fig. 16. Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

are relatively close to each other and, therefore, tunneling contributes quite a lot to the conductance through the contact. In order to obtain a conductance around 1 $G_0$, detected in the experiments, a linear chain has to be considered. The existence of such linear chains is obvious from the long plateau of the corresponding conductance trace with respect to the piezo voltage in the break-junction experiments. Because at present the computer code for the real space Kubo equation is restricted to geometrical structures confined to three-dimensional translational invariant simple bulk parent lattices, as the simplest model of such a contact a slanted linear chain was considered as shown in Fig. 17b. In there, the middle layer ($C$) and the adjacent layers ($C \pm 1$) contain only one Au atom, layers $C \pm 2$ and $C \pm 3$ four and nine Au atoms, respectively, whereas layers $C \pm 4$ refer to the first two full layers. The sum in Eq. (209) was carried out for two pairs of layers, namely for $p = C - 4$, $q = C + 4$ (full layers) and for $p = C - 2$, $q = C + 2$ (not full layers). The convergence with respect to the number of atoms in the chosen layers can be seen from Fig. 20. The respective converged values are 1.10 $G_0$ and 1.17 $G_0$. In the case of $p = C - 2$, $q = C + 2$ the contribution from the vacuum sites is nearly zero: considering only four Au atoms in the summation already gave a value for the conductance very close to the converged one. The small difference between the two calculated values, 0.07 $G_0$, most likely has to be attributed to the use of the atomic sphere approximation (ASA). Nevertheless, as expected, the calculated conductance is very close to the ideal value of 1 $G_0$.

Another interesting structure is the $2 \times 2$ chain described in Ref. [25], namely the structure depicted in Fig. 17c. The conductance for this structure was calculated by including to the summation 100 atoms from each of the first two full layers. As result a conductance of 2.58 $G_0$ was obtained. Papanikolaou et al. [25] got a conductance of 3 $G_0$ for an infinite Cu wire to be associated with three conducting channels within the Landauer approach. For an infinite wire the transmission probability is unity for all states, therefore, the conductance is just the number of bands crossing the Fermi level. For the present case of a finite chain, the transmission probability is less than unity for all the conducting states. This qualitatively explains the reduced conductance with respect to an infinite wire.

Finally, the dependence of the conductance on the thickness of the nanocontacts was studied. All the investigated structures have $C_{4v}$ symmetry and the central layer of the systems is a plane of reflection symmetry. The set-up of the structures is summarized in Table 3. Contact 0 refers to a broken contact which is embedded into a host with $N_{Au} = 7$ and $N_V = 7$ layers, see Fig. 15, while the others have different thicknesses from 1 up to 9 Au atoms in the central layer, and are embedded into a host characterized by $N_{Au} = 5$ and $N_V = 5$, see Fig. 15.

In Fig. 21, the calculated conductances are displayed as performed by including nearly 100 atoms from each of the first two full layers: $p = C - 4$, $q = C + 4$ for the broken contact and



Figure 20. Conductance versus the number of sites included in the sum in Eq. (209) for the slanted linear chain shown in Fig. 17b. Full circles are the results of summing in layers $p = C - 4$ and $q = C + 4$ (first full layers), and squares refer to a summation in layers $p = C - 2$ and $q = C + 2$ (layers containing four Au atoms). Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004. American Physical Society.

Table 3. Set-up of various nanocontacts.

| Layer position | Contact | | | | |
|---|---|---|---|---|---|
| | 0 | 1 | 4 | 5 | 9 |
| $C \pm 4$ | Full | Full | Full | Full | Full |
| $C \pm 3$ | 9 | Full | Full | Full | Full |
| $C \pm 2$ | 4 | 9 | 16 | 21 | 25 |
| $C \pm 1$ | 1 | 4 | 9 | 12 | 16 |
| $C$ | 0 | 1 | 4 | 5 | 9 |

Shown is the number of Au atoms in the layers as labeled by $C$, $C \pm 1$, etc., see Fig. 16. Contact 1 refers to Fig. 17a. Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

$p = C - 3$. $q = C + 3$ for all the other cases, see Table 3. It can be seen that the conductance is almost proportional to the number of Au atoms in the central layer. This finding can qualitatively be compared with the result of model calculations for the conductance of a three-dimensional electron gas through a connective neck as a function of its area in the limit of $\vartheta_0 = 90°$ for the opening angle [26]. In the case of the broken contact, the nonzero conductance can again be attributed to tunneling of electrons.

## 7.2. Gold Contact with an Impurity

In recent break junction experiments [27], remarkable changes of the conductance histograms of nanocontacts formed from AuPd alloys have been observed when varying the Pd concentration. Studying the effect of impurities placed into the nanocontact are, in that context, at least relevant for dilute alloys. The interesting question is whether the presence of impurities can be observed in the measured conductance. For that reason we investigated transition metal impurities such as Pd, Fe, and Co placed at various positions of the point contact as shown in Fig. 17a. For the notation of the impurity positions, see Fig. 22.

The calculated spin and orbital moments of the magnetic impurities are listed in Table 4. They were calculated with assuming the direction of magnetization to be parallel to the z axis ($\widehat{M} = \hat{z}$), that is, normal to the planes. Additional calculations of the magnetic anisotropy energy confirmed this choice. As usual for magnetic impurities with reduced coordination number [28], both for Fe and Co remarkably high spin moments, and in all positions of a Co impurity large orbital moments were obtained. In particular, the magnitude of the orbital moments is very sensitive to the position of the impurity. This is most obvious in the case of Fe, where at positions B and C the orbital moment is relatively small, but at position A a surprisingly high value of 0.47 $\mu_B$ was obtained.



**Figure 21.** Conductance *versus* the number of Au atoms in the central layer for the Au contacts described in Table 3. Reprinted with permission from [34], K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

**Figure 22.** Impurity positions (light gray spheres) in a Au point contact, see Fig. 17a. Reprinted with permission from [34]. K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004. American Physical Society.

The summation over 116 atoms from each of the first two full layers ($p = C - 3$, $q = C + 3$) in Eq. (209) has been carried out in order to evaluate the conductance. The calculated values are summarized in Table 5.

A Pd impurity (independent of position) reduces only little the conductance as compared to a pure Au point contact. This qualitatively can be understood from the local density of states (LDOS) of the Pd impurity as calculated for an imaginary part of the energy of $\delta = 1$ mRy, the real space scattering path operators by using 1830 $k_\parallel$ points in the 2D IBZ. It should be noted that the LDOS at site $i$ ($n_i$) is defined as follows

$$n_i(\varepsilon) = \mp \frac{1}{\pi} \int_{\Omega_i} d^3 r Im[G^\pm(\varepsilon, \mathbf{r}, \mathbf{r})] \tag{214}$$

where $\Omega_i$ denoted the volume of the $i$th unit cell. In Fig. 23, the corresponding LDOS at positions A and C is plotted. Clearly, the change of the coordination number (8 at position A and 12 at position C), that is, different hybridization between the Pd and Au $d$ bands, results into different widths for the Pd $d$-like LDOS. In both cases, however, the Pd $d$ states are completely filled and no remarkable change in the LDOS at Fermi level (conducting states) happens.

The case of magnetic impurities seems to be more interesting. As can be inferred from Table 5, impurities at position B change only very little the conductance. Being placed at position A, however, Fe and Co atoms increase the conductance by 11% and 24%, whereas at position C they decrease the conductance by 19% and 27%, respectively. In Ref. [25] it was found that single Fe, Co (and also Ni) defects in a 2 × 2 infinite Cu wire decreased the conductance. By analyzing the DOS, it was concluded that the observed reduction of the conductance is due to a depletion of the $s$-like states in the minority band. The above situation is very similar to the case of an Fe or Co impurity in position C of the point contact considered, even the calculated drop of the conductance ($\sim -20\%$ for Fe and $\sim -28\%$ for Co) agrees quantitatively well with our present result. Our result, namely, that Fe and Co impurities at position A increase the conductance, however, *cannot* be related to the results of Ref. [25]. In order to understand this feature, one carefully has to investigate the LDOS calculated for the point contact.

In Fig. 24, the minority $d$-like LDOS of the Fe and Co impurities in positions A and C are plotted as resolved according to the canonical orbitals $d_{z^2-x^2}$, $d_{xy}$, $d_{xz}$, $d_{yz}$ and $d_{x^2-y^2}$. It has

**Table 4.** Calculated spin and orbital moments of magnetic impurities placed at different positions in a Au point contact, see Fig. 22. $\hat{M} = \hat{z}$.

| Position | $S$ [$\mu_B$] | | $L$ [$\mu_B$] | |
|---|---|---|---|---|
| | Fe | Co | Fe | Co |
| A | 3.36 | 2.01 | 0.47 | 0.38 |
| B | 3.46 | 2.17 | 0.04 | 0.61 |
| C | 3.42 | 2.14 | 0.07 | 0.22 |

Reprinted with permission from [34]. K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004. American Physical Society.

Table 5. Calculated conductances of a Au point contact with impurities on different positions, see Fig. 22.

| Impurity position | Conductance $[G_0]$ | | |
|---|---|---|---|
| | Pd | Fe | Co |
| A | 2.22 | 2.67 | 2.97 |
| B | 2.24 | 2.40 | 2.26 |
| C | 2.36 | 1.95 | 1.75 |
| Pure Au | | 2.40 | |

Reprinted with permission from [34]. K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004. American Physical Society.

to be pointed out strongly that this kind of partial decomposition, usually referred to as the $(\ell, m, s)$ representation of the LDOS, is not unique within a relativistic formalism, since due to the spin-orbit interaction different spin- and orbital components are mixed. However, due to the large spin-splitting of Fe and Co the mixing of the majority and minority spin-states can be neglected.

As can be seen from Fig. 24, the LDOS of an impurity in position A is much narrower than in position C. This is an obvious consequence of the difference in the coordination numbers (8 for position A and 12 for positions C). Thus an impurity in position A hybridizes less with the neighboring Au atoms and, as implied by the LDOS, the corresponding $d$ states are fairly localized. Also to be seen is a spin-orbit induced splitting of about 8 mRy ($\sim$0.1 eV) in the very narrow $d_{x^2-y^2}$-$d_{xy}$ states of the impurities in position A. The difference of the band filling for the two kind of impurities shows up in a clear downward shift of the LDOS of Co with respect to that of Fe.

In Fig. 25, a comparison between a nonrelativistic and a relativistic calculation is displayed: the splitting in the $d_{x^2-y^2}$ and $d_{xy}$ states vanishes by turning off the spin-orbit coupling.

In order to explain the change in the conductance through the point contact caused by impurities in positions A and C, the $s$-like DOS at the center site, that is, at the narrowest section of the contact, is plotted in the top half of Fig. 26. As a comparison, the corresponding very flat $s$-like DOS is shown for a pure Au contact. For contacts with impurities this $s$-like DOS shows a very interesting shape, which can indeed be correlated with the corresponding $d_{3z^2-r^2}$-like DOS at the impurity site, see bottom half of Fig. 26. As clearly can be seen, the center positions and the widths of the $d_{3z^2-r^2}$-like DOS peaks and those of the respective (anti-)resonant $s$-like DOS shapes coincide well with each other. This kind of behavior in the DOS resembles the case studied by Fano for a continuum band and a discrete energy level in the presence of configuration interaction (hybridization) [29]. Apparently, by keeping this analogy, in the point contact the $s$-like states play the role of a continuum and



Figure 23. Local density of states of a Pd impurity in position A (solid line) and in position C (dashed line) of a Au point contact, see Fig. 22. Reprinted with permission from [34]. K. Palotas et al., *Phys. Rev. B* in press (2004). © 2004. American Physical Society.

**Figure 24.** Minority-spin orbital-resolved $d$-like local density of states of Fe and Co impurities in position A (upper panels) and in position C (lower panels) of a Au point contact, see Fig. 22. Reprinted with permission from [34], K. Palotás et al., *Phys. Rev. B* in press (2004). © 2004, American Physical Society.

the $d_{3z^2-r^2}$-like state of the impurity acts as the discrete energy level. Because these this two kinds of states share the same cylindrical symmetry, interactions between them can occur due to backscattering effects. It should be noted that similar resonant line-shapes in the STM I-V characteristics have been observed for Kondo impurities at surfaces [30, 31] and explained theoretically in Ref. [32].

Inspecting Fig. 26, the enhanced $s$-like DOS at the Fermi level at the center of the point contact provides a nice interpretation for the enhancement of the conductance when an Fe and Co impurity is placed at position A. As the peak position of the $d_{3z^2-r^2}$-like states of Fe is shifted upwards by more than 0.01 Ry with respect to that of Co, the corresponding resonance of the $s$-like states is also shifted and the $s$-like DOS at the Fermi level is decreased. This is also in agreement with the calculated conductances. In the case of impurities at position C, that is, in a position by $a_{3D} = 7.681$ a.u. away from the center of the contact,



**Figure 25.** Minority-spin orbital-resolved $d$-like local density of states of a Co impurity in position A, see Fig. 22. On the left nonrelativistic, on the right relativistic calculation is displayed.

**Figure 26.** Top left: minority-spin $s$-like local density of states at the center site of a Au point contact with an impurity at position A, see Fig. 22 (solid line: Co, dashed line: Fe). Top right: the same as before, but with an impurity at position C. As a comparison, in both figures the corresponding LDOS for the pure Au contact is plotted by dotted lines. The solid vertical lines highlight the position of the Fermi energy. Bottom: minority-spin $d_{3z^2-r^2}$ local density of states of the impurities (solid line: Co, dashed line: Fe) at positions A (left) and C (right). Vertical dashed lines mark the center positions of the $d_{3z^2-r^2}$-LDOS peaks.

the resonant line-shape of the $s$-like states is reversed in sign, therefore, one observes a decreased $s$-like DOS at the Fermi level, explaining in this case the decreased conductance, see Table 5. As, however, the $s$-like DOS for the case of a Co impurity is larger than for an Fe impurity, this simple picture cannot account correctly for the opposite relationship obtained for the corresponding conductances.

## 8. CONCLUSIONS

In the current paper, methods and approaches were introduced and discussed in order to describe the electric properties of "real space" nanostructures, that is, of systems with a finite number of atoms properly embedded in (metallic) substrates. It is indeed important to note that whenever structures nanoscaled in two dimensions (finite supported clusters) are considered, the influence of the substrate has to be taken into account. Furthermore, because in particular finite magnetic nanostructures (very small islands) are of technological interest, a relativistic approach has to be applied in order to describe adequately the orientation of the magnetization in these structures on all levels (electronic structure, electric transport). The last example shown, namely atom-sized contacts, refers to a topic that will be of increasing importance in many applications. in particular, since the conducting properties of such contacts can be modulated quite a bit by placing impurities in the very vicinity of the actual contact atoms.

## ACKNOWLEDGMENTS

## REFERENCES

1. P. Weinberger, *Physics Reports* 377, 281 (2003).
2. N. F. Mott, *Adv. Phys.* 13, 325 (1964).
3. R. Landauer, *IBM J. Res. Dev.* 1, 223 (1957); *Z. Phys. B* 21, 247 (1975); *Z. Phys. B* 68, 217 (1987); *IBM J. Res. Dev.* 32, 306 (1988); *J. Phys. Condens. Matter* 1, 8099 (1989).
4. M. Büttiker, *Phys. Rev. Lett.* 57, 1761 (1986); *IBM J. Res. Dev.* 32, 317 (1988).
5. R. Kubo, M. Toda, and N. Hashitsume, "Statistical Physics II: Non-equilibrium Statistical Mechanics." Springer, Berlin, 1985; R. Kubo, *J. Phys. Soc. Jpn.* 12, 570 (1957).
6. J. M. Luttinger, in "Mathematical Methods in Solid State and Superfluid Theory" (R. C. Clark and G. H. Derric, Eds.), Oliver and Boyd, Edinburgh, 1967, Chap. 4, p. 157.
7. D. A. Greenwood, *Proc. Phys. Soc. London* 71, 585 (1958).
8. W. H. Butler, *Phys. Rev. B* 31, 3260 (1985).
9. H. U. Baranger and A. D. Stone, *Phys. Rev. B* 40, 8169 (1989).
10. J. Zabloudil, R. Hammerling, L. Szunyogh, and P. Weinberger, "Electron Scattering in Solid Matter, a Theoretical and Computational Treatise." 1st Edn., Springer, 2004.
11. M. E. Rose, "Relativistic Electron Theory." Wiley, New York, 1961.
12. A. Messiah, "Quantum Mechanics." North Holland, Amsterdam, 1969.
13. S. L. Altmann and P. Herzig, "Point Group Theory Tables." Clarendon Press, Oxford, UK, 1994.
14. G. A. Korn and T. M. Korn "Mathematical Handbook for Scientists and Engineers." McGraw-Hill, New York, 1961.
15. B. L. Györffy and M. J. Stott, in "Band Structure Spectroscopy of Metals and Alloys" (D. J. Fabian and L. M. Watson, Eds.), Academic Press, London, 1973.
16. L. Szunyogh, B. Újfalussy, P. Weinberger, and J. Kollár, *Phys. Rev. B* 49, 2721 (1994).
17. R. Zeller, P. H. Dederichs, B. Újfalussy, L. Szunyogh, and P. Weinberger, *Phys. Rev. B* 52, 8807 (1995).
18. L. Szunyogh, B. Újfalussy, and P. Weinberger, *Phys. Rev. B* 51, 9552 (1995).
19. P. Weinberger, P. M. Levy, J. Banhart, L. Szunyogh, and B. Újfalussy, *J. Phys.: Condens. Matter* 8, 7677 (1996).
20. P. Mavropoulos, N. Papanikolaou, and P. H. Dederichs, *Phys. Rev. B* 69, 125104 (2004).
21. J. Banhart, H. Ebert, P. Weinberger, and J. Voitländer, *Phys. Rev. B* 50, 2104 (1994).
22. L. Dulca, J. Banhart, and G. Czycholl, *Phys. Rev. B* 61, 16502 (2000).
23. J. M. Krans, I. K. Yanson, Th. C. M. Govaert, R. Hesper, and J. M. van Ruitenbeck, *Phys. Rev. B* 48, 14721 (1993).
24. M. Brandbyge, J. Schiøtz, M. R. Sorensen, P. Stoltze, K. W. Jacobsen, J. K. Nørskov, L. Olesen, E. Laegsgaard, I. Stensgaard, and F. Besenbacher, *Phys. Rev. B* 52, 8499 (1995).
25. N. Papanikolaou, J. Opitz, P. Zahn, and I. Mertig, *Phys. Rev. B* 66, 165441 (2002).
26. J. A. Torres, J. I. Pascual, and J. J. Sáenz, *Phys. Rev. B* 49, 16581 (1994).
27. A. Enomoto, S. Kurokawa, and A. Sakai, *Phys. Rev. B* 65, 125410 (2002).
28. B. Lazarovits, L. Szunyogh, and P. Weinberger, *Phys. Rev. B* 65, 104441 (2002).
29. U. Fano, *Phys. Rev.* 124, 1866 (1961).
30. V. Madhavan, W. Chen, T. Jamneala, M. F. Crommie, and N. S. Wingreen, *Science* 280, 567 (1998).
31. H. C. Manoharan, C. P. Lutz, and D. M. Eigler, *Nature* 403, 512 (2000).
32. O. Újsághy, J. Kroha, L. Szunyogh, and A. Zawadowski, *Phys. Rev. Lett.* 85, 2557 (2000); O. Újsághy, G. Zaránd, and A. Zawadowski, *Solid State Commun.* 117, 167 (2001).
33. K. Palotas, B. Lazarovits, L. Szunyogh, and P. Weinberger, *Phys. Rev. B* 67, 174404 (2003).
34. K. Palotas, B. Lazarovits, L. Szunyogh, and P. Weinberger, *Phys. Rev. B* in press (2004).

# CHAPTER 7

# Transport Theory for Interacting Electrons Connected to Reservoirs

## Akira Oguri

Department of Material Science, Osaka City University,
Sumiyoshi-ku, Osaka, Japan

## CONTENTS

# 1. INTRODUCTION

Quantum transport through finite interacting-electron systems has been studied extensively in this decade. For instance, the Coulomb blockade and various effects named after Kondo, Aharanov-Bohm, Fano, Josephson, and so forth, in quantum dot systems have been a very active field of research. Furthermore, the realization of non-Fermi liquid systems such as the Tomonaga-Luttinger liquid in quantum wires and multichannel Kondo behavior in some novel systems have also been investigated by a number theorists.

To study transport properties of correlated electron systems, theoretical approaches that can treat correctly both the interaction and quantum interference effects are required. The Keldysh Green's function approach is one such method [1–9]. Specifically, the formulation for the nonlinear current–voltage profile by Caroli et al. has been applied widely to the quantum transport phenomena. In this report, we describe the outline of the Keldysh formalism in Section 2 and then apply it to a single Anderson impurity, which is a standard model of quantum dots in the Kondo regime in Section 3.

When a finite sample is connected to reservoirs that can be approximated by free-electron systems with continuous energy spectrums, the low-energy eigenstates of whole the system including the attached reservoirs are determined coherently. Thus, to understand the low-temperature properties, the information about the low-lying energy states of the whole system is required. The local Fermi-liquid theory [10, 11], which was originally introduced for the Kondo systems [12], is also applicable to the transport properties in wide classes of the interacting-electron systems at low temperatures. In Section 4, we reformulate the transport theory for the interacting systems connected to noninteracting leads based on the Kubo formalism. The dc conductance can be written in a Landauer-type form with a many-body transmission coefficient determined by a three-point correlation function. We also provide a brief introduction to Tomonaga-Luttinger model in Section 5 to take a quick look at the transport properties of a typical interacting-electron system in one dimension.

# 2. KELDYSH FORMALISM FOR QUANTUM TRANSPORT

## 2.1. Thermal Equilibrium

We start with a system that consists of three regions: a finite central region $(C)$ and two reservoirs on the left $(L)$ and the right $(R)$. The central region consists of $N$ resonant levels, and the interaction $U_{j_1 j_2 ; j_3 j_4}$ is switched on only for the electrons in this region. We assume that each of the reservoirs is infinitely large and has a continuous energy spectrum. The central region and reservoirs are connected with the mixing matrix elements $v_L$ and $v_R$, as illustrated in Figure 1. The complete Hamiltonian is given by

$$\mathcal{H}_{\text{tot}}^{\text{eq}} = \mathcal{H}_L + \mathcal{H}_R + \mathcal{H}_C^0 + \mathcal{H}_C^1 + \mathcal{H}_{\text{mix}} \tag{1}$$

$$\mathcal{H}_L = -\sum_{\substack{i,j \\ \sigma}} t_{ij}^L c_{i\sigma}^\dagger c_{j\sigma}, \qquad \mathcal{H}_R = -\sum_{\substack{i,j \\ \sigma}} t_{ij}^R c_{i\sigma}^\dagger c_{j\sigma} \tag{2}$$



Figure 1. Schematic picture of the system.

$$\mathcal{H}_C^0 = -\sum_{i,j} t_{ij}^l c_{i\sigma}^\dagger c_{j\sigma}, \qquad \mathcal{H}_C^i = \frac{1}{2} \sum_{i,j,k,l} U_{ij,kl} c_{i\sigma}^\dagger c_{k\sigma'}^\dagger c_{l\sigma'} c_{j\sigma} \tag{3}$$

$$\mathcal{H}_{mix} = -\sum_{\sigma} v_l \left[ c_{0\sigma}^\dagger c_{1\sigma} + H.c. \right] - \sum_{\sigma} v_R \left[ c_{N+1\sigma}^\dagger c_{N\sigma} + H.c. \right] \tag{4}$$

Here, $c_{j\sigma}^\dagger$ ($c_{j\sigma}$) creates (destroys) an electron with spin $\sigma$ at site $j$, and $\mu$ is the chemical potential. Also, $t_{ij}^l$, $t_{ij}^R$, and $t_{ij}^c$ are the intraregion hopping matrix elements in each of the three regions $L$, $R$, and $C$, respectively. The labels 1, 2, ..., $N$ are assigned to the sites in the central region. Specifically, the label 1 ($N$) is assigned to the site at the interface on the left (right), and the label 0 ($N + 1$) is assigned to the site at the reservoir-side of the left (right) interface. We will be using units $\hbar = 1$ unless otherwise noted.

The density matrix for the equilibrium state $\rho_{eq}$ is given by

$$\rho_{eq} = e^{-\beta\{\mathcal{H}_{tot}^{eq} - \mu(N_L + N_C + N_R)\}} / \mathrm{Tr}\, e^{-\beta\{\mathcal{H}_{tot}^{eq} - \mu(N_L + N_C + N_R)\}} \tag{5}$$

$$N_L = \sum_{i \in L, \sigma} c_{i\sigma}^\dagger c_{i\sigma}, \qquad N_C = \sum_{i \in C, \sigma} c_{i\sigma}^\dagger c_{i\sigma}, \qquad N_R = \sum_{i \in R, \sigma} c_{i\sigma}^\dagger c_{i\sigma} \tag{6}$$

Therefore, the Hamiltonian $\mathcal{H}_{tot}^{eq}$ and a single chemical potential $\mu$ determine the statistical weight in thermal equilibrium.

## 2.2. Statistical Weight for Nonequilibrium Steady States

When the voltage $V$ is applied, the contribution of the electrostatic potential has to be included into $\mathcal{H}_{tot}^{eq}$, as

$$\mathcal{H}_{tot} = \mathcal{H}_{tot}^{eq} + \mathcal{T}_{ext} \tag{7}$$

$$\mathcal{T}_{ext} = \Phi_L N_L + \Phi_R N_R + \sum_{i \in C, \sigma} \Phi_C(i) \tag{8}$$

Here, $\Phi_L$ and $\Phi_R$ are the potentials for the lead at $L$ and $R$, respectively, and the applied bias voltage corresponds to $eV = \Phi_L - \Phi_R$. To determine the potential profile in the central region $\Phi_C(i)$, the energy of the electric field should also be included into the Hamiltonian Eq. (2), and it should be determined self-consistently. However for simplicity, we assume that $\Phi_C(i)$ is a given function. In Fig. 2, two typical profiles are illustrated. For an insulating sample, there must be a finite electric field in the central region and the potential shows approximately a linear $i$-dependence as that in the panel (a). In an opposite case, for a metallic sample, the electric field vanishes in the central region, and the potential profile will become the one as shown in the panel (b). Realistic situations seem to be in between these two extreme cases.

In contrast to the thermal equilibrium, one cannot write down the density matrix that describes the nonequilibrium statistical weight simply by using a single chemical potential,

$$\rho \neq e^{-\beta\{\mathcal{H}_{tot} - \mu(N_L + N_C + N_R)\}} / \mathrm{Tr}\, e^{-\beta\{\mathcal{H}_{tot} - \mu(N_L + N_C + N_R)\}} \tag{9}$$

because this statistical weight describes the situation after the electrons have already been redistributed to gain the electrostatic potential energy. One possible statistical weight that describes a nonequilibrium steady state was introduced by Caroli et al. [4].



Figure 2. Examples for the profile of the electrostatic potential $\Phi_C(i)$ for (a) an insulating sample and (b) a metal sample, where $eV = \Phi_L - \Phi_R$.

In the formulation of Caroli et al., the coupling to the leads $\mathscr{H}_{\text{mix}}$ and the interaction $\mathscr{H}_C^U$ in the sample region are switched on adiabatically by separating the total Hamiltonian in the form

$$\mathscr{H}_{\text{tot}}(t) = \mathscr{H}_1 + \mathscr{H}_2(t) \tag{10}$$

$$\mathscr{H}_1 = \mathscr{H}_{1;L} + \mathscr{H}_{1;C} + \mathscr{H}_{1;R} \tag{11}$$

$$\mathscr{H}_{1;L} = \mathscr{H}_L + \Phi_L N_L, \quad \mathscr{H}_{1;R} = \mathscr{H}_R + \Phi_R N_R \tag{12}$$

$$\mathscr{H}_{1,C} = \mathscr{H}_C^0 + \sum_{i \in C, \sigma} \Phi_C(i) c_{i\sigma}^\dagger c_{i\sigma} \tag{13}$$

$$\mathscr{H}_2(t) = \left[ \mathscr{H}_{\text{mix}} + \mathscr{H}_C^U \right] e^{-\delta|t|} \tag{14}$$

Here, $\delta = 0^-$ is an positive infinitesimal. Because $\mathscr{H}_1$ has a quadratic form, it is possible to use the Wick theorem in the perturbation expansion with respect to $\mathscr{H}_2$. At $t = -\infty$ the two reservoirs and the impurity are isolated, so that the different chemical potentials $\mu_L$, $\mu_R$, and $\mu_C$ can be introduced into the three regions in the initial condition. The time evolution of the density matrix is determined by the equation

$$\frac{\partial}{\partial t} \rho(t) = -i\left[ \mathscr{H}_{\text{tot}}(t), \rho(t) \right] \tag{15}$$

The formal solution of this equation can be obtained, by using the interaction representation $\tilde{\rho}(t) = e^{i\mathscr{H}_1 t}\rho(t)e^{-i\mathscr{H}_1 t}$ and $\tilde{\mathscr{H}}_2(t) = e^{i\mathscr{H}_1 t}\mathscr{H}_2(t)e^{-i\mathscr{H}_1 t}$, as

$$\tilde{\rho}(t) = U(t, t_0)\tilde{\rho}(t_0)U(t_0, t) \tag{16}$$

$$U(t, t_0) = T\exp\left[ -i \int_{t_0}^{t} dt' \, \tilde{\mathscr{H}}_2(t') \right] \tag{17}$$

$$U(t_0, t) = \tilde{T}\exp\left[ i \int_{t_0}^{t} dt' \, \tilde{\mathscr{H}}_2(t') \right] \tag{18}$$

Here, T denotes the operator for the chronological time order, and $\tilde{T}$ is the anti-time-ordering operator. Note that Eq. (18) is the Hermite conjugate of Eq. (17). Caroli et al. have assumed that the initial condition at $t_0 \to -\infty$ is given by

$$\tilde{\rho}(-\infty) = \frac{e^{-\beta\{\mathscr{H}_{1;L} - \mu_L N_L\}} e^{-\beta\{\mathscr{H}_{1;C} - \mu_C N_C\}} e^{-\beta\{\mathscr{H}_{1;R} - \mu_R N_R\}}}{\text{Tr}\left[ e^{-\beta\{\mathscr{H}_{1;L} - \mu_L N_L\}} e^{-\beta\{\mathscr{H}_{1;C} - \mu_C N_C\}} e^{-\beta\{\mathscr{H}_{1;R} - \mu_R N_R\}} \right]} \tag{19}$$

Namely, at $t_0 \to -\infty$, each system is in a thermal equilibrium state with the chemical potential $\mu_{L,R,C}$. Here, $\mu_L - \mu_R = \Phi_L - \Phi_R = eV$.

## 2.3. Perturbation Expansion along the Keldysh Contour

The perturbed part $\tilde{\mathscr{H}}_2(t)$ is switched fully on at $t = 0$. Therefore, the expectation value of physical quantities are defined with respect to the density matrix at $t = 0$,

$$\langle \mathscr{O} \rangle \equiv \text{Tr}\left[ \rho(0) \mathscr{O}_s \right]$$
$$= \text{Tr}\left[ \tilde{\rho}(0) \mathscr{O}_s \right] = \text{Tr}\left[ \tilde{\rho}(-\infty) U(-\infty, 0) \mathscr{O}_s U(0, -\infty) \right] \tag{20}$$

where $\mathscr{O}_s$ is a Schrödinger operator, and $[\rho(0), \mathscr{H}_{\text{tot}}(0)] = 0$ for the stationary states. Equation (20) can be rewritten by using a property of the time-evolution operator. $U(-\infty, 0) = U(-\infty, +\infty)U(+\infty, 0)$, as

$$\langle \mathscr{O} \rangle = \langle U(-\infty, +\infty)U(+\infty, 0) \mathscr{O}_s U(0, -\infty) \rangle_0 \tag{21}$$

where $\langle \cdots \rangle_0 \equiv \text{Tr}[\tilde{\rho}(-\infty) \cdots]$. The stream of the time in Eq. (21) can be mapped onto a loop shown in Fig. 3, starting from $t = -\infty$, observing a quantity $\mathscr{O}$ at $t = 0$, then proceeding to $t = +\infty$, and then going back to $t = -\infty$. In the case of the usual $T = 0$ Green's function

**Figure 3.** The Keldysh contour for the time evolution.

with respect to the equilibrium ground state, the wavefunction at $t = +\infty$ is essentially the same with the one at $t = -\infty$ apart from a phase factor if the initial state has no degeneracy [13]. Thus, Eq. (21) can be decoupled at the time $t = +\infty$. However, this simplification does not take place in the nonequilibrium case of the initial condition Eq. (19), Therefore, one has to treat the time loop including the way back to $t \rightarrow -\infty$.

The time-dependent expectation value is defined by using the Heisenberg operator $\ell_H(t)$, and is written in the form

$$\langle \ell(t) \rangle \equiv \mathrm{Tr}[\rho(0)\ell_H(t)]$$

$$= \langle U(-\infty, +\infty)U(+\infty, t)\tilde{\ell}(t)U(t, -\infty)\rangle_0$$

$$= \langle U(-\infty, +\infty)\{\mathrm{T}\ U(+\infty, -\infty)\tilde{\ell}(t)\}\rangle_0$$

$$= \langle \mathrm{T}_C U_c \tilde{\ell}(t^-)\rangle_0 \tag{22}$$

Here, the relation among the Schrödinger $\ell_S$, interaction $\tilde{\ell}(t)$, and Heisenberg $\ell_H(t)$ representations are given by

$$\tilde{\ell}(t) = e^{iH_0 t}\ell_S e^{-iH_0 t}, \qquad \ell_H(t) = U(0, t)\tilde{\ell}(t)U(t, 0) \tag{23}$$

In Eq. (22), $\mathrm{T}_C$ and $U_c$ express the time order and time evolution along the Keldysh contour, respectively, and $t^-$ denotes the time in the $-$branch in Fig. 3.

The perturbation expansion with respect to $H_2$ can be carried out by substituting Eqs. (17) and (18), respectively, into $U(+\infty, -\infty)$ and $U(-\infty, +\infty)$ in Eq. (22), as

$$\langle \ell(t) \rangle = \sum_{n=0}^{\infty}\sum_{m=0}^{\infty}\frac{i^n}{n!}\frac{(-i)^m}{m!}\int_{-\infty}^{+\infty}dt_1' \cdots dt_n'\int_{-\infty}^{+\infty}dt_1 \cdots dt_m$$

$$\times \left\langle \left\{\mathrm{\bar{T}}\tilde{H}_2(t_1') \cdots \tilde{H}_2(t_n')\right\}\left\{\mathrm{T}\tilde{H}_2(t_1) \cdots \tilde{H}_2(t_m)\tilde{\ell}(t)\right\}\right\rangle_0 \tag{24}$$

The Wick's theorem is applicable to the average $\langle \cdots \rangle_0$ because $H_1$ has a bilinear form. However, because $U(+\infty, -\infty)$ gives a factor $(-i)^m$ for the $m$th order terms while $U(-\infty, +\infty)$ gives a factor $(+i)^n$ for the $n$th order terms, four types of the Green's functions are necessary to distinguish the contributions from these two branches.

## 2.4. Nonequilibrium Green's Function

We now introduce the four types of the Green's functions, which are required in the Feynman-diagrammatic approach to the perturbation expansion along the time loop,

$$G_{ij}^{--}(t_1, t_2) \equiv -i\langle \mathrm{T}c_{i\sigma}(t_1)c_{j\sigma}^{\dagger}(t_2)\rangle$$

$$= -i\langle \mathrm{T}_C U_c \tilde{c}_{i\sigma}(t_1^-)\tilde{c}_{j\sigma}^{\dagger}(t_2^-)\rangle_0 \tag{25}$$

$$G_{ij}^{++}(t_1, t_2) \equiv -i\langle \mathrm{\bar{T}}c_{i\sigma}(t_1)c_{j\sigma}^{\dagger}(t_2)\rangle$$

$$= -i\langle \mathrm{T}_C U_c \tilde{c}_{i\sigma}(t_1^+)\tilde{c}_{j\sigma}^{\dagger}(t_2^+)\rangle_0 \tag{26}$$

$$G_{ij}^{-+}(t_1, t_2) \equiv -i\langle c_{i\sigma}(t_1)c_{j\sigma}^{\dagger}(t_2)\rangle$$

$$= -i\langle \mathrm{T}_C U_c \tilde{c}_{i\sigma}(t_1^-)\tilde{c}_{j\sigma}^{\dagger}(t_2^+)\rangle_0 \tag{27}$$

$$G_{ij}^{+-}(t_1, t_2) \equiv -i\langle c_{j\sigma}^{\dagger}(t_2)c_{i\sigma}(t_1)\rangle$$

$$= -i\langle \mathrm{T}_C U_c \tilde{c}_{i\sigma}(t_1^+)\tilde{c}_{j\sigma}^{\dagger}(t_2^-)\rangle_0 \tag{28}$$

Here. $c_{i\sigma}(t_1)$ and $c^{\dagger}_{i\sigma}(t_2)$ are Heisenberg operators. $t^{\pm}$ denotes the time $+$ or $-$ branch in Fig. 3. For these Green's functions, there is another notation used widely in literatures, and the relation between that and the present one by Lifshitz-Pitaevskii [5] is summarized in Table 1.

Based on the Feynman-diagrammatic approach, the Dyson equation can be expressed in a $2 \times 2$ matrix form:

$$G_{ij}(\omega) = g_{ij}(\omega) + \sum_{lm} g_{il}(\omega)\Sigma_{lm}(\omega)G_{mj}(\omega)$$  (29)

$$G_{ij} = \begin{bmatrix} G_{ij}^{--} & G_{ij}^{-+} \\ G_{ij}^{+-} & G_{ij}^{++} \end{bmatrix}, \quad \Sigma_{lm} = \begin{bmatrix} \Sigma_{lm}^{--} & \Sigma_{lm}^{-+} \\ \Sigma_{lm}^{+-} & \Sigma_{lm}^{++} \end{bmatrix}$$  (30)

Here. $g_{ij}$ is the Green's function determined by the unperturbed Hamiltonian $\mathcal{H}_1$ and density matrix $\bar{\rho}(-\infty)$ for the initial isolated system. The Fourier transform has been carried out for stationary states,

$$G(t_1, t_2) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} G(\omega) e^{i\omega(t_1 - t_2)}$$  (31)

For instance, in the noninteracting case $\mathcal{H}_C^1 = 0$, the self-energy correction is caused only by the couplings between the sample and reservoirs $\mathcal{H}_{mix}$,

$$\Sigma_{ij}^0 = -v_L(\delta_{i,1}\delta_{j,0} + \delta_{i,0}\delta_{j,1})\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

$$-v_R(\delta_{i,N}\delta_{m,N+1} + \delta_{i,N+1}\delta_{m,N})\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$  (32)

Note that the four types Green's functions are not independent.

$$G^{--} + G^{++} = G^{-+} + G^{+-}, \quad \Sigma^{--} + \Sigma^{++} = -\Sigma^{-+} - \Sigma^{+-}$$  (33)

Thus, the Dyson equation (29) can be expressed in terms of three independent quantities by carrying out a Unitary transformation $P^{-1}GP$:

$$P = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$  (34)

$$\begin{bmatrix} 0 & G_{ij}^a \\ G_{ij}^r & F_{ij} \end{bmatrix} = \begin{bmatrix} 0 & g_{ij}^a \\ g_{ij}^r & F_{ij}^0 \end{bmatrix} + \sum_{lm}\begin{bmatrix} 0 & g_{il}^a \\ g_{il}^r & F_{il}^0 \end{bmatrix}\begin{bmatrix} \Omega_{lm} & \Sigma_{lm}^r \\ \Sigma_{lm}^a & 0 \end{bmatrix}\begin{bmatrix} 0 & G_{mj}^a \\ G_{mj}^r & F_{mj} \end{bmatrix}$$  (35)

Here. $G^r$ and $G^a$ are the retarded and advanced Green's functions, respectively.

$$G^r = G^{--} - G^{-+}, \quad G^a = G^{--} - G^{+-}, \quad F = G^{--} + G^{++}$$  (36)

$$\Sigma^r = \Sigma^{--} + \Sigma^{-+}, \quad \Sigma^a = \Sigma^{--} + \Sigma^{+-}, \quad \Omega = \Sigma^{--} + \Sigma^{++}$$  (37)

The function $F$ and $\Omega$ link closely to a nonequilibrium distribution. Alternatively, the original four Green's functions can be expressed with these three functions as

$$G^{--} = [F + (G^r + G^a)]/2, \quad G^{-+} = [F - (G^r + G^a)]/2$$  (38)

$$G^{+-} = [F - (G^r - G^a)]/2, \quad G^{++} = [F + (G^r - G^a)]/2$$  (39)

Table 1. Correspondence between two standard notations.

| Lifshitz-Pitaevskii | $G$ | $G$ | $G$ | $G$ |
|---|---|---|---|---|
| Alternative notation | $G$ | $G$ | $G$ | $G$ |

Similarly, the four self-energies are written in terms of $\Sigma^r$, $\Sigma^a$ and $\Omega$,

$$\Sigma \quad = [\Omega + (\Sigma^r + \Sigma^a)]/2, \quad \Sigma^{--} = [\Omega - (\Sigma^r + \Sigma^a)]/2 \tag{40}$$

$$\Sigma \quad = -[\Omega - (\Sigma^r - \Sigma^a)]/2, \quad \Sigma \quad = -[\Omega + (\Sigma^r - \Sigma^a)]/2 \tag{41}$$

The Dyson equation for the three functions are deduced from Eq. (35)

$$G^r = g^r + g^r \Sigma^r G^r, \quad G^a = g^a + g^a \Sigma^a G^a \tag{42}$$

$$F = F^0 + F^0 \Sigma^a G^a + g^r \Sigma^r F + g^r \Omega G^a \tag{43}$$

Here, we have suppressed the subscripts for simplicity, and these equations should be understood symbolically. Equation (43) can be solved formally by using Eq. (42) as

$$F = [1 - g^r \Sigma^r]^{-1} F^0 [1 + \Sigma^a G^a] + [1 - g^r \Sigma^r]^{-1} g^r \Omega G^a$$

$$= G^r \{g^r\}^{-1} F^0 \{g^a\}^{-1} G^a + G^r \Omega G^a \tag{44}$$

## 2.5. Green's Function for the Initial State

The unperturbed Green's function $g_{ij}$ is determined by $\mathcal{H}_1$ and the initial density matrix $\tilde{\rho}(-\infty)$. Initially at $t \to \infty$, the three regions are isolated and noninteracting. Therefore, $g_{ij;\nu}$ for $\nu = L, R, C$ is given by

$$g^r_{ij;\nu}(\omega) = \sum_n \frac{\phi_{n;\nu}(i)\phi^*_{n;\nu}(j)}{\omega - \epsilon_{n;\nu} + i\delta}, \quad g^a_{ij;\nu}(\omega) = \sum_n \frac{\phi_{n;\nu}(i)\phi^*_{n;\nu}(j)}{\omega - \epsilon_{n;\nu} - i\delta} \tag{45}$$

$$F^0_{ij;\nu}(\omega) = [1 - 2f_\nu(\omega)][g^r_{ij;\nu}(\omega) - g^a_{ij;\nu}(\omega)] \tag{46}$$

Here, $\epsilon_{n;\nu}$ and $\phi_{n;\nu}(i)$ are the one-particle eigenvalue and eigenstate of $\mathcal{H}_{1;\nu}$. The information about the statistical distribution is contained in the function $F^0_{ij;\nu}$ via $f_\nu(\omega) = f(\omega - \mu_\nu)$, where $f(\epsilon) = [e^{\beta\epsilon} + 1]^{-1}$. In the system we are considering, each of the reservoirs ($\nu = L, R$) has a continuous energy spectrum, and the isolated sample ($\nu = C$) has a discrete energy spectrum. Thus, the full Green's function becomes to depend only on $\mu_L$ and $\mu_R$ and does not depend on $\mu_C$ [14]. This is because the contribution of $F^0$ to the corresponding full one $F$ arises in a sandwiched form $\{g^r\}^{-1} F^0 \{g^a\}^{-1}$ as described in Eq. (44). Thus, the singular contributions of $\delta$ functions in $F^0 \propto [g^r - g^a]$ of the sample region are canceled out by the zero points of the inverse Green's functions in both sides, to yield $\{g^r\}^{-1} F^0 \{g^a\}^{-1} = 0$ for $\nu = C$.

## 2.6. Nonequilibrium Current for Noninteracting Electrons

The nonequilibrium average of the charge and current can be deduced from the Green's functions. For instance, by using Eqs. (28) and (31), an equal-time correlation function can be written in the form

$$\langle c^\dagger_{i\sigma} c_{j\sigma} \rangle = -iG^{-r}_{ji}(0, 0) = -i \int_{-\infty}^\infty \frac{d\omega}{2\pi} G^{-r}_{ji}(\omega) \tag{47}$$

Therefore, the current flowing from the left lead to the sample is given by

$$J_L = iev_L \sum_\sigma [c^\dagger_{1\sigma} c_{0\sigma} - c^\dagger_{0\sigma} c_{1\sigma}] \tag{48}$$

$$\langle J_L \rangle = 2ev_L \int_{-\infty}^\infty \frac{d\omega}{2\pi} [G^{-r}_{01}(\omega) - G^{-r}_{10}(\omega)] \tag{49}$$

The expectation value for the current flowing from the sample to right lead, $J_R$, can also be written in a similar form. In the noninteracting case $\mathcal{H}_C^I = 0$, Eq. (49) can be rewritten in

terms of retarded and advanced Green's functions which link the two different interfaces of the sample [4];

$$\langle J \rangle = \frac{2e}{h} \int_{-\infty}^{\infty} d\omega [f_L(\omega) - f_R(\omega)] \tilde{J}_{\shortmid\shortmid}(\omega) \tag{50}$$

$$\tilde{J}_{\shortmid\shortmid}(\omega) = 4\Gamma_L(\omega) G_{1N}^a(\omega) \Gamma_R(\omega) G_{N1}^r(\omega) \tag{51}$$

$$\Gamma_L(\omega) = -\mathrm{Im}\big[v_L^2 g_{00}'(\omega)\big], \quad \Gamma_R(\omega) = -\mathrm{Im}\big[v_R^2 g_{N+1N+1}'(\omega)\big] \tag{52}$$

Note that $\langle J_L \rangle = \langle J_R \rangle$ ($\equiv \langle J \rangle$) in steady states. The outline of the derivation are provided in Section 3.3 for a single Anderson impurity. Equation (51) implies that the current is determined by the electrons with the energy $\mu_R \lesssim \omega \lesssim \mu_L$ at low temperatures, where $\mu_L - \mu_R = eV$.

For interacting electron systems, the nonequilibrium current can not generally be written in the form of Eq. (50). It does only in a particular case where the connection of the two leads and the sample has a special symmetry described by a relation $\Gamma^L(\epsilon) = \lambda \Gamma^R(\epsilon)$ in the notation used in Ref. [15]. In the interacting case, the imaginary part of the self-energy caused by the inelastic scattering becomes finite. It links with the contributions of the vertex corrections, and the formulation becomes somewhat complicated. Nevertheless, in the linear-response regime, the dc conductance for interacting electrons can be expressed in a Landauer-type form quite generally [16–18] even in the case without the special symmetry mentioned above [19]. Specifically, at zero temperature $T = 0$, the imaginary part of the self-energy and vertex corrections for the current become zero at the Fermi energy $\omega = 0$, and the transmission probability can be written in the form of Eq. (51) with the interacting Green's functions. We discuss the details of these points in Section 4.

# 3. OUT-OF-EQUILIBRIUM ANDERSON MODEL

We now apply the Keldysh formalism to a single Anderson impurity connected to two leads as illustrated in Fig. 4. It corresponds to the $N = 1$ case of the Hamiltonian Eq. (1) and has been widely used as a model for quantum dots. For convenience, we change the label for the sites: the new one for the impurity site is given by $0 \Rightarrow d$, and that for the interfaces at the left and right leads are $0 \Rightarrow L$ and $N + 1 \Rightarrow R$, respectively.

## 3.1. Green's Function for the Anderson Impurity

The self-energy in the interacting case can be classified into two parts.

$$\mathbf{\Sigma}_{ij}(\omega) = \mathbf{\Sigma}_{ij}^0(\omega) + \mathbf{\Sigma}_U(\omega)\delta_{i,d}\delta_{j,d} \tag{53}$$

Here, $\mathbf{\Sigma}_{ij}^0$ corresponds the one defined in Eq. (32), which represents the effects purely due to the mixing with reservoirs and sample. The remaining part $\mathbf{\Sigma}_U$ contains the contributions of the onsite Coulomb interaction $U$. Substituting the self-energy Eq. (53) into the Dyson equation (29), we obtain a set of equations for the impurity Green's function,

$$G_{dd}(\omega) = g_{dd}(\omega) + g_{dd}(\omega)\mathbf{\Sigma}_U(\omega)G_{dd}(\omega)$$

$$- v_L g_{dd}(\omega)\tau_3 G_{Ld}(\omega) - v_R g_{dd}(\omega)\tau_3 G_{Rd}(\omega) \tag{54}$$

$$G_{Ld}(\omega) = -v_L g_L(\omega)\tau_3 G_{dd}(\omega) \tag{55}$$

$$G_{Rd}(\omega) = -v_R g_R(\omega)\tau_3 G_{dd}(\omega) \tag{56}$$



Figure 4. Anderson impurity connected to two leads.

where $g_L$ and $g_R$ are the Green's function at the interfaces at left and right, respectively. The explicit form of the unperturbed Green's function at the impurity site is given by $\{g_{dd}(\omega)\}^{-1} = (\omega - \epsilon_d)\tau_3$ with one of the Pauli matrices $\tau_3$. Substituting Eqs. (55) and (56) into Eq. (54), we obtain

$$G_{dd}(\omega) = g_{dd}(\omega) + g_{dd}(\omega)\left[\sigma(\omega) + \Sigma_L(\omega)\right]G_{dd}(\omega) \tag{57}$$

$$\sigma(\omega) = v_L^2\tau_3 g_L(\omega)\tau_3 + v_R^2\tau_3 g_R(\omega)\tau_3 \tag{58}$$

Therefore,

$$\{G_{dd}(\omega)\}^{-1} = \{G_{dd}^{(0)}(\omega)\}^{-1} - \Sigma_L(\omega) \tag{59}$$

$$\{G_{dd}^{(0)}(\omega)\}^{-1} = \{g_{dd}(\omega)\}^{-1} - \sigma(\omega) \tag{60}$$

Here, $G_{dd}^{(0)}(\omega)$ is the Green's function for the noninteracting case. Furthermore, an alternatively form of the Dyson equation $G = g + G\Sigma g$ yields

$$G_{dL}(\omega) = -v_L G_{dd}(\omega)\tau_3 g_L(\omega) \tag{61}$$

$$G_{dR}(\omega) = -v_R G_{dd}(\omega)\tau_3 g_R(\omega) \tag{62}$$

Thus, the intersite Green's functions can be deduced from $G_{dd}(\omega)$ by using Eqs. (55)–(56) and (61)–(62).

The voltage-dependence arises via the unperturbed Green's functions for the leads at $\nu = L$ and $R$,

$$g_\nu(\omega) = P\begin{bmatrix} 0 & g_\nu^a(\omega) \\ g_\nu^r(\omega) & F_\nu^0(\omega) \end{bmatrix}P^{-1} \tag{63}$$

$$F_\nu^0(\omega) = [1 - 2f_\nu(\omega)][g_\nu^r(\omega) - g_\nu^a(\omega)]$$
$$= -2i[1 - 2f_\nu(\omega)]\Gamma_\nu(\omega)/v_\nu^2 \tag{64}$$

where $\Gamma_\nu(\omega) \approx -v_\nu^2 \mathrm{Im}[g_\nu^r(\omega)]$. The pure mixing part of the self energy $\sigma(\omega)$ can be calculated from Eqs. (58) and (63),

$$\sigma(\omega) = P\begin{bmatrix} \Omega^{(0)}(\omega) & \sigma^r(\omega) \\ \sigma^a(\omega) & 0 \end{bmatrix}P^{-1} \tag{65}$$

$$\Omega^{(0)}(\omega) = v_L^2 F_L^0(\omega) + v_R^2 F_R^0(\omega) \tag{66}$$

$$\sigma^r(\omega) = v_L^2 g_L^r(\omega) + v_R^2 g_R^r(\omega) \tag{67}$$

and $\sigma^a(\omega) = \{\sigma^r(\omega)\}^*$. Then the noninteracting Green's function $G_{dd}^{(0)}(\omega)$ can be determined via Eq. (60),

$$G_{dd}^{(0)--}(\omega) = [1 - f_{eff}(\omega)]G_{dd}^{(0)r}(\omega) + f_{eff}(\omega)G_{dd}^{(0)a}(\omega) \tag{68}$$

$$G_{dd}^{(0)++}(\omega) = -f_{eff}(\omega)G_{dd}^{(0)r}(\omega) - [1 - f_{eff}(\omega)]G_{dd}^{(0)a}(\omega) \tag{69}$$

$$G_{dd}^{(0)-+}(\omega) = -f_{eff}(\omega)\left[G_{dd}^{(0)r}(\omega) - G_{dd}^{(0)a}(\omega)\right] \tag{70}$$

$$G_{dd}^{(0)+-}(\omega) = [1 - f_{eff}(\omega)]\left[G_{dd}^{(0)r}(\omega) - G_{dd}^{(0)a}(\omega)\right] \tag{71}$$

where

$$f_{eff}(\omega) = \frac{f_L(\omega)\Gamma_L + f_R(\omega)\Gamma_R}{\Gamma_L + \Gamma_R} \tag{72}$$

$$G_{dd}^{(0)r}(\omega) = \frac{1}{\omega - \epsilon_d - \sigma^r(\omega)} \tag{73}$$

and $G_{dd}^{(1)a}(\omega) = \{G_{dd}^{(1)r}(\omega)\}^*$. Thus, the effects of the bias voltage arise through the distribution function $f_{\text{eff}}(\omega)$.

The full Green's function for the impurity site can be expressed, using Eqs. (42)–(44) and (64)–(66), as

$$G_{dd}^r(\omega) = \frac{1}{\omega - \epsilon_d - \sigma^r(\omega) - \Sigma_L^r(\omega)} \qquad (74)$$

$$F_{dd}(\omega) = G_{dd}^r(\omega)\left[\Omega^{(1)}(\omega) + \Omega_L(\omega)\right]G_{dd}^a(\omega) \qquad (75)$$

where $G_{dd}^a(\omega) = \{G_{dd}^r(\omega)\}^*$. Note that $\Omega_L(\omega) = -\Sigma_L^{-+}(\omega) - \Sigma_L^+(\omega)$ and $F_{dd}(\omega)$ are pure imaginary. The four elements of $G_{dd}(\omega)$ can be written, using Eqs. (38)–(39) and (74)–(75), as

$$G_{dd}^{--}(\omega) = \left[1 - \tilde{f}_{\text{eff}}(\omega)\right]G_{dd}^r(\omega) + \tilde{f}_{\text{eff}}(\omega)G_{dd}^a(\omega) \qquad (76)$$

$$G_{dd}^{-+}(\omega) = -\tilde{f}_{\text{eff}}(\omega)G_{dd}^r(\omega) - \left[1 - \tilde{f}_{\text{eff}}(\omega)\right]G_{dd}^a(\omega) \qquad (77)$$

$$G_{dd}^{--}(\omega) = -\tilde{f}_{\text{eff}}(\omega)\left[G_{dd}^r(\omega) - G_{dd}^a(\omega)\right] \qquad (78)$$

$$G_{dd}^{-+}(\omega) = \left[1 - \tilde{f}_{\text{eff}}(\omega)\right]\left[G_{dd}^r(\omega) - G_{dd}^a(\omega)\right] \qquad (79)$$

where $G_{dd}^{+-}(\omega) = -\{G_{dd}^{-+}(\omega)\}^*$, and $\tilde{f}_{\text{eff}}(\omega)$ is a correlated distribution defined by

$$\tilde{f}_{\text{eff}}(\omega) = \frac{f_L(\omega)\Gamma_L + f_R(\omega)\Gamma_R - (1/2i)\Sigma_L^{-+*}(\omega)}{\Gamma_L + \Gamma_R - \text{Im}\Sigma_L^r(\omega)} \qquad (80)$$

With this distribution function, the number of the electrons in the impurity site can be written in the form

$$\langle n_d \rangle = 2 \int d\omega \tilde{f}_{\text{eff}}(\omega)\left(-\frac{1}{\pi}\right)\text{Im } G_{dd}^r(\omega) \qquad (81)$$

In the equilibrium case, $\mu \equiv \mu_L = \mu_R$, both $\tilde{f}_{\text{eff}}(\omega)$ and $f_{\text{eff}}(\omega)$ coincide with the Fermi function $f(\omega)$, because of the property Eq. (84).

## 3.2. Properties of the Green's Functions at $eV = 0$

We summarize here the properties of the Keldysh Green's function in the limit of the zero-bias voltage $eV = 0$, at which $\mu_L = \mu_R$. In this case, the four self-energies can be written in the form

$$\Sigma_{L;\text{eq}}^{--}(\omega) = [1 - f(\omega)]\Sigma_{L;\text{eq}}^r(\omega) + f(\omega)\Sigma_{L;\text{eq}}^a(\omega) \qquad (82)$$

$$\Sigma_{L;\text{eq}}^{+-}(\omega) = -f(\omega)\Sigma_{L;\text{eq}}^r(\omega) - [1 - f(\omega)]\Sigma_{L;\text{eq}}^a(\omega) \qquad (83)$$

$$\Sigma_{L;\text{eq}}^{-+}(\omega) = f(\omega)\left[\Sigma_{L;\text{eq}}^r(\omega) - \Sigma_{L;\text{eq}}^a(\omega)\right] \qquad (84)$$

$$\Sigma_{L;\text{eq}}^{+-}(\omega) = -[1 - f(\omega)]\left[\Sigma_{L;\text{eq}}^r(\omega) - \Sigma_{L;\text{eq}}^a(\omega)\right] \qquad (85)$$

Furthermore, in equilibrium, $F_{dd;\text{eq}}(\omega)$ and $\Omega_{L;\text{eq}}(\omega)$ are determined by the retarded and advanced functions,

$$F_{dd;\text{eq}}(\omega) = [1 - 2f(\omega)]\left[G_{dd;\text{eq}}^r(\omega) - G_{dd;\text{eq}}^a(\omega)\right] \qquad (86)$$

$$\Omega_{L;\text{eq}}(\omega) = [1 - 2f(\omega)]\left[\Sigma_{L;\text{eq}}^r(\omega) - \Sigma_{L;\text{eq}}^a(\omega)\right] \qquad (87)$$

Specificality at zero temperature, $\Sigma_{L;\text{eq}}^{-+}(\omega)$ and $\Sigma_{L;\text{eq}}^{+-}(\omega)$ vanish, respectively, at $\omega > \mu$ and $\omega < \mu$, because of the Fermi function in Eqs. (84) and (85). Similarly, the Green's functions $G_{dd;\text{eq}}^{-+}(\omega)$ and $G_{dd;\text{eq}}^{+-}(\omega)$ also vanish at $\omega > \mu$ and $\omega < \mu$, respectively. Therefore, at the

equilibrium ground state, the usual $T = 0$ formalism, which yields a single-component Dyson equation,

$$G_{dd;eq}(\omega) = G_{dd;eq}^{(0)}(\omega) + G_{dd;eq}^{(0)}(\omega)\Sigma_{L;eq}(\omega)G_{dd;eq}(\omega) \tag{88}$$

becomes available.

## 3.3. Current Through the Anderson Impurity

The operators for the current flows from the left reservoir to the sample $J_L$ and that from the sample to the right reservoir $J_R$ are given by

$$J_L = ie\sum_\sigma v_L\left[d_\sigma^\dagger c_{L\sigma} - c_{L\sigma}^\dagger d_\sigma\right] \tag{89}$$

$$J_R = ie\sum_\sigma v_R\left[c_{R\sigma}^\dagger d_\sigma - d_\sigma^\dagger c_{R\sigma}\right] \tag{90}$$

As mentioned in Section 2.6, the expectation values can be expressed in the form

$$\langle J_L\rangle = 2e\int_{-\infty}^\infty \frac{d\omega}{2\pi}\, v_L\left[G_{Ld}^{-+}(\omega) - G_{dL}^{-+}(\omega)\right] \tag{91}$$

$$\langle J_R\rangle = 2e\int_{-\infty}^\infty \frac{d\omega}{2\pi}\, v_R\left[G_{dR}^{-+}(\omega) - G_{Rd}^{-+}(\omega)\right] \tag{92}$$

The intersite Green's functions in the right-hand side of Eq. (92) can be expressed in terms of $G_{dd}(\omega)$ using Eqs. (55)–(56), and (61)–(62), as

$$P^{-1}G_{Rd}P = \begin{bmatrix} 0 & G_{Rd}^a \\ G_{Rd}^r & F_{Rd} \end{bmatrix} = -v_R\begin{bmatrix} 0 & g_R^a G_{dd}^a \\ g_R^r G_{dd}^r & g_R^r F_{dd} + F_R^{tt}G_{dd}^a \end{bmatrix} \tag{93}$$

$$P^{-1}G_{dR}P = \begin{bmatrix} 0 & G_{dR}^a \\ G_{dR}^r & F_{dR} \end{bmatrix} = -v_R\begin{bmatrix} 0 & G_{dd}^a g_R^a \\ G_{dd}^r g_R^r & G_{dd}^r F_R^{tt} + F_{dd}g_R^a \end{bmatrix} \tag{94}$$

Thus, using Eqs. (39), (75), and (93)–(94), we obtain

$$v_R[G_{dR}^{-+} - G_{Rd}^{-+}] = \frac{v_R}{2}[F_{dR} - F_{Rd}]$$

$$= \frac{v_R^2}{2}\left[(g_R^r - g_R^a)F_{dd} - F_R^{tt}(G_{dd}^r - G_{dd}^a)\right]$$

$$= \Gamma_R G_{dd}^r G_{dd}^a\left[4\Gamma_L(f_L - f_R) - i\Omega_r - 2(1 - 2f_R)\mathrm{Im}\Sigma_U^r\right] \tag{95}$$

Similarly, for the right-hand side of Eq. (91), we obtain

$$v_L[G_{Ld}^{-+} - G_{dL}^{-+}] = \frac{v_L}{2}[F_{Ld} - F_{dL}]$$

$$= \frac{v_L^2}{2}\left[-(g_L^r - g_L^a)F_{dd} + F_L^{tt}(G_{dd}^r - G_{dd}^a)\right]$$

$$= \Gamma_L G_{dd}^r G_{dd}^a\left[4\Gamma_R(f_L - f_R) + i\Omega_U + 2(1 - 2f_L)\mathrm{Im}\Sigma_U^r\right] \tag{96}$$

If the two couplings between the Anderson impurity and leads have a property $\Gamma_L(\omega) = \lambda\Gamma_R(\omega)$ with $\lambda$ being a constant [15], the expression for current can be simplified by taking an average, $\Gamma_L \times$ (95) $+ \Gamma_R \times$ (96), as

$$\langle J\rangle = \frac{\Gamma_L\langle J_R\rangle + \Gamma_R\langle J_L\rangle}{\Gamma_R + \Gamma_L}$$

$$= \frac{2e}{h}\int_{-\infty}^\infty d\omega[f_L(\omega) - f_R(\omega)]\frac{4\Gamma_L\Gamma_R}{\Gamma_R + \Gamma_L}[-\mathrm{Im}G_{dd}^r(\omega)] \tag{97}$$

Note that $\langle J\rangle = \langle J_L\rangle = \langle J_R\rangle$ for stationary states.

## 3.4. Perturbation Expansion with Respect to $\mathcal{H}_C^U$

So far, we have discussed general properties of the nonequilibrium Green's functions. To study how the interelectron interactions affect the transport properties, the self-energy $\Sigma_{t^-}(\omega)$ must be calculated with reliable methods. Because the noninteracting Green's function, which includes the couplings between the leads and the Anderson impurity have already been obtained in Eqs. (68)–(70), the remaining task is calculating $\Sigma_U(\omega)$, for instance, by taking $G_{dd}^{(0)}$ to be the unperturbed Green's function. Then, the interacting Green's function $G_{dd}$ are deduced via the Dyson equation (59).

The self-energy $\Sigma_{t^-}(\omega)$ can be calculated with the perturbation expansion with respect to the interelectron interaction $\mathcal{H}_t^U$. For generating the perturbation series, it is convenent to introduce an effective action,

$$S(\eta^\dagger, \eta) = S_0(\eta^\dagger, \eta) + S_t(\eta^\dagger, \eta) + S_{ex}(\eta^\dagger, \eta) \tag{98}$$

$$S_0(\eta^\dagger, \eta) = \sum_\sigma \int dt\, dt'\; \eta_{t\sigma}^\dagger(t) K_{dd}^{(0)}(t, t') \eta_\sigma(t') \tag{99}$$

$$S_t(\eta^\dagger, \eta) = -U \int dt \Big[ \eta_\uparrow^\dagger (t)\eta_{\uparrow-}(t)\eta_\downarrow^\dagger (t)\eta_{\downarrow-}(t)$$
$$-\eta_{\uparrow+}^\dagger (t)\eta_{\uparrow+}(t)\eta_{\downarrow+}^\dagger(t)\eta_{\downarrow+}(t) \Big] \tag{100}$$

Here, $\eta_\sigma^\tau(t) = (\eta_{\sigma-}^\dagger(t), \eta_{\sigma+}^\dagger(t))$ is a two-component Grassmann number corresponding to the $-$ and $+$ branches of the Keldysh contour shown in Fig. 3. The Kernel $K_{dd}^{(0)}$ is determined by the noninteracting Greens function,

$$K_{dd}^{(0)}(\omega) \equiv \left\{ G_{dd}^{(0)}(\omega) \right\}^{-1} \tag{101}$$

$$K_{dd}^{(0)}(t, t') = \int \frac{d\omega}{2\pi} K_{dd}^{(0)}(\omega) e^{-i\omega(t-t')} \tag{102}$$

In Eq. (100), the sign for the interaction along the $-$branch and that for $+$branch are different. This correspond to the sign arises in Eq. (24), and it is determined from which of the time-evolution operators, $U(+\infty, -\infty)$ or $U(-\infty, +\infty)$, the perturbation terms arise. For $S_{ex}(\eta^\dagger, \eta)$ in Eq. (98), we introduce an external source of two anticomutating c-numbers $j_\sigma^\dagger(t) = (j_{\sigma-}^\dagger(t), j_{\sigma+}^\dagger(t))$ following along the standard procedure [20],

$$S_{ex}(\eta^\dagger, \eta) = -\sum_\sigma \int dt \big[ \eta_\sigma^\dagger(t) j_\sigma(t) + j_\sigma^\dagger(t)\eta_\sigma(t') \big] \tag{103}$$

In this formulation, the Green's functions are generated from a functional $Z[j]$, as

$$Z[j] = \int D\eta^\dagger D\eta\; e^{iS(\eta^\dagger, \eta)} \tag{104}$$

$$G_{dd,\sigma}^{rr}(t, t') = -i \frac{1}{Z[0]} \frac{\delta}{\delta j_{\sigma r}^\dagger(t)} \frac{\delta}{\delta j_{\sigma r'}(t')} Z[j]\Big|_{j=0} \tag{105}$$

In the noninteracting case, the functional integration can be calculated analytically

$$Z^{(0)}[j] \equiv \int D\eta^\dagger D\eta\; e^{i[S_0(\eta^\dagger, \eta) + S_{ex}(\eta^\dagger, \eta)]}$$
$$= Z^{(0)}[0] \exp\Big[ -i \sum_\sigma \int dt\, dt'\, j_\sigma^\dagger(t) G_{dd}^{(0)}(t, t') j_\sigma(t') \Big] \tag{106}$$

The generating functional $Z[j]$ can be rewritten in the form

$$Z[j] = e^{iS_t(-i\delta/\delta j, -i\delta/\delta j)} Z^{(0)}[j] \tag{107}$$

Here, $\eta$ and $\eta'$ in the action $S_t$ $(\eta^+, \eta)$ has been replaced by the functional derivatives

$$\eta^+_{\sigma r}(t) \Rightarrow -i\frac{\delta}{\delta j_{\sigma r}(t)}, \qquad \eta_{\sigma r}(t) \Rightarrow i\frac{\delta}{\delta j^+_{\sigma r}(t)} \tag{108}$$

The perturbation series can be obtained by substituting Eq. (106) into (107) and then expanding $e^{iS_t}$ in a power series of $S_t$.

## 3.5. Fermi-Liquid Behavior at Low Bias Voltages

The out-of-equilibrium Anderson model has been studied by a number of theoretical approaches. In this section, we discuss briefly the low-bias behavior of the Green's function and differential conductance, which have been deduced from the Ward identities for the Keldysh formalism [21]. In equilibrium and linear-response regime, the low-energy properties at $\max(\omega, T) \ll T_K$ can be described by the local Fermi liquid theory, where $T_K$ is the Kondo temperature [12]. The results deduced from the Ward identities show that the nonlinear properties at small bias-voltages $eV \ll T_K$ can also be described by the local Fermi-liquid theory.

The low-energy behavior of $\mathrm{Im}\Sigma^r_U(\omega)$ has been calculated exactly up to terms of order $\omega^2$, $(eV)^2$, and $T^2$,

$$\mathrm{Im}\Sigma^r_U(\omega) = -\frac{\pi}{2}\{A_{eq}(0)\}^3|\Gamma_{\uparrow\downarrow;\downarrow\uparrow}(0,0;0,0)|^2$$

$$\times\left[(\omega - \alpha eV)^2 + \frac{3\Gamma_L\Gamma_R}{(\Gamma_L + \Gamma_R)^2}(eV)^2 + (\pi T)^2\right] \tag{109}$$

where $\Gamma_{\sigma\sigma';\sigma'\sigma}(\omega, \omega'; \omega', \omega)$ is the vertex function for the causal Green's function in the zero-temperature formalism, and $A_{eq}(\omega) = -\mathrm{Im}\ G^r_{dd;eq}(\omega)/\pi$. The parameter $\alpha$ is defined by $\alpha \equiv (\alpha_L\Gamma_L - \alpha_R\Gamma_R)/(\Gamma_L + \Gamma_R)$, where $\alpha_L$ and $\alpha_R$ are constants which have been introduced to specify how the bias voltage is applied to the equilibrium state. Namely, $\mu_L \equiv \alpha_L eV$ and $\mu_R \equiv -\alpha_R eV$ with $\alpha_L + \alpha_R = 1$.

The real part of the self-energy is generally complicated. However, it is simplified in the electron–hole symmetric case for $\epsilon_d = -U/2$, $\Gamma_L = \Gamma_R$, and $\mu_L = -\mu_R = eV/2$. In this case the spectral wight at the Fermi energy becomes $A_{eq}(0) = 1/(\pi\Delta)$ with $\Delta = \Gamma_L + \Gamma_R$, and the low-energy behavior of the real part of the self-energy is given by

$$\mathrm{Re}\ \Sigma^r_U(\omega) = (1 - z^{-1})\omega + O(\omega^3) \tag{110}$$

$$z^{-1} \equiv 1 - \frac{\partial\Sigma^r_{U;eq}(\omega)}{\partial\omega}\bigg|_{\omega=0} \tag{111}$$

Here, the constant Hartree term $U/2$ is included into the unperturbed part, and it set the position of the Kondo peak on the Fermi energy $\epsilon_d + U/2 = 0$. Therefore, $G^r(\omega)$ can be deduced exactly up to terms of order $\omega^2$, $T^2$ and $(eV)^2$ from Eqs. (109) and (110),

$$G^r(\omega) \simeq \frac{z}{\omega + i\tilde{\Delta} + i(\tilde{U}^2/2\tilde{\Delta}(\pi\tilde{\Delta})^2)[\omega^2 + 3/4(eV)^2 + (\pi T)^2]} \tag{112}$$

where the renormalized parameters are defined by

$$\tilde{\Delta} \equiv z\Delta, \qquad \tilde{U} \equiv z^2\Gamma_{\uparrow\downarrow;\downarrow\uparrow}(0,0;0,0) \tag{113}$$

The order $U^2$ result [14] can be reproduced from Eq. (112) by replacing $\tilde{U}$ with the bare Coulomb interaction $U$ and using the perturbation result for the renormalization factor $z = 1 - (3 - \pi^2/4)u^2 + \cdots$, where $u = U/(\pi\Delta)$.

This result shows that in the symmetric case the low-voltage behavior is characterized by the two parameters $\tilde{\Delta}$ and $\tilde{U}$. These two parameters are defined with respect to the equilibrium ground state, and the exact Bethe ansatz results exist for these parameters.

The width of the Kondo resonance $\tilde{\Delta}$ decreases with increasing $U$, and becomes close to $\tilde{\Delta} \simeq (4/\pi)T_K$ for large $U$ with the Kondo temperature defined by

$$T_K = \pi\Delta\sqrt{u/(2\pi)}\exp[-\pi^2 u/8 + 1/(2u)] \tag{114}$$

The Wilson ratio is usually defined by $R \equiv \tilde{\chi}_s/\tilde{\gamma}$, where $\tilde{\gamma}$ and $\tilde{\chi}_s$ are the enhancement factors for the $T$-linear specific heat and spin susceptibility, respectively [11]. Alternatively, it corresponds to the ration of $\tilde{U}$ to $\tilde{\Delta}$.

$$R - 1 = \frac{\tilde{U}}{(\pi\tilde{\Delta})} \tag{115}$$

The Wilson ratio takes a value $R = 1$ for $U = 0$, and it reaches $R = 2$ in the strong-coupling limit $U \to \infty$.

The nonequilibrium current $\langle J \rangle$ is calculated by substituting Eq. (112) into Eq. 97). Then, the differential conductance $dJ/dV$ are determined exactly up to terms of order $T^2$ and $(eV)^2$,

$$\frac{dJ}{dV} = \frac{2e^2}{h}\left[1 - \frac{1 + 2(R-1)^2}{3}\left(\frac{\pi T}{\tilde{\Delta}}\right)^2 - \frac{1 + 5(R-1)^2}{4}\left(\frac{eV}{\tilde{\Delta}}\right)^2 + \cdots \right] \tag{116}$$

Therefore, the nonlinear $(eV)^2$ term is also scaled by the resonance width $\tilde{\Delta}$, and the coefficient is determined by the parameter $R - 1$, or $\tilde{U}/(\pi\tilde{\Delta})$.

# 4. TRANSPORT THEORY BASED ON KUBO FORMALISM

We have discussed in Section 2.6 that in the noninteracting case the nonequilibrium current can be written in a Landauer-type form, as Eq. (50). The similar expression has been derived for interacting electrons in a special case, when the couplings between the leads and sample satisfy the condition $\Gamma^L(\epsilon) = \lambda\Gamma^R(\epsilon)$ in the notation used in Ref. [15]. For this condition to be held, the interacting sites must be classified into the following two groups: one group consists of the sites that are connected directly to both of the two leads, and other group consists of the sites that have no direct links (hopping matrix elements) to the leads. This condition restricts the application of Eq. (50). For instance, if there is an interacting site that is connected to only one of the two leads, the condition is not satisfied. Therefore, Eq. (50) is not applicable to a series of quantum dots as illustrated in Fig. 6. Nevertheless, in the linear-response regime, the Landauer-type expression of the dc conductance, Eq (139), can be derived quite generally without the condition mentioned above [19].

In Section 4.1, based on the Kubo formalism, we describe the outline of the derivation of Eq. (139) for interacting electrons. Our proof uses the analytic properties of the vertex corrections following along the Éliashberg theory of a transport equation for correlated electrons [22, 23]. The many-body transmission probability $\mathcal{T}(\epsilon)$ is given by Eq. (140), and it is written in terms of a three-point correlation function. At zero temperature, the imaginary part of the self-energy due to the interaction and the vertex corrections for the current become zero at the Fermi energy $\epsilon = 0$. Due to this property, the transmission probability at $T = 0$ is determined by the single-particle Green's functions as shown in Ec. (141). In Section 4.2, the current conservation law for the correlation functions is described with the generalized Ward identity, which expresses the relation between the self-energy and current vertex. In Section 4.3, we provide the Lehmann representation of the three-point functions to carry out the analytic continuation formally. It can be also used for nonperturbative calculations of $\mathcal{T}(\epsilon)$. We apply this formulation to a finite Hubbard chain in Section 4.4 and show an example of the transmission probability $\mathcal{T}(\epsilon)$ for interacting electrons.

## 4.1. Many-Body Transmission Coefficient $\mathcal{T}(\epsilon)$

We now consider the Hamiltonian $\mathcal{H}_{tot}^{eq}$ defined in Eq. (1) again, which is also illustrated in Fig. 1. The dc conductance $g$ can be determined in the Kubo formalism, and

it corresponds to the $\omega$-linear imaginary part of a current–current correlation function $K_{\alpha\alpha'}(\omega + i0^+)$:

$$g = e^2 \lim_{\omega \to 0} \frac{K_{\alpha\alpha'}(\omega + i0^+) - K_{\alpha\alpha'}(i0^+)}{i\omega} \tag{117}$$

$$K_{\alpha\alpha'}(i\nu_l) = \int_0^\beta d\tau \langle T_\tau J_\alpha(\tau) J_{\alpha'}(0) \rangle e^{i\nu_l \tau} \tag{118}$$

where $\alpha = L$ or $R$. The retarded correlation can be calculated via the analytic continuation $K_{\alpha\alpha'}(\omega + i0^+) = K_{\alpha\alpha'}(i\nu_l)\big|_{i\nu_l \to \omega + i0^+}$, where $\nu_l = 2\pi l/\beta$ is the Matsubara frequency. The current operator $J_\alpha$ is defined by

$$J_L = i \sum_\sigma v_L \left( c_{1\sigma}^\dagger c_{0\sigma} - c_{0\sigma}^\dagger c_{1\sigma} \right) \tag{119}$$

$$J_R = i \sum_\sigma v_R \left( c_{N+1\sigma}^\dagger c_{N\sigma} - c_{N\sigma}^\dagger c_{N+1\sigma} \right) \tag{120}$$

Here, $J_L$ is the current flowing into the sample from the left lead, and $J_R$ is the current flowing out to the right lead from the sample. These currents and total charge in the sample $\rho_C$ satisfy the equation of continuity

$$\rho_C = \sum_{j \in C, \sigma} c_{j\sigma}^\dagger c_{j\sigma} \tag{121}$$

$$\frac{\partial \rho_C}{\partial t} + J_R - J_L = 0 \tag{122}$$

Owing to this property, the dc conductance $g$ defined in Eq. (117) does not depend on the choice of $\alpha$ and $\alpha'$ [18]. Note that $K_{\alpha'\alpha}(z) = K_{\alpha\alpha'}(z)$ owing to the time-reversal symmetry of $\mathcal{H}$.

To calculate the $\omega$-linear imaginary part of $K_{\alpha\alpha'}(\omega + i0^+)$, we introduce the three-point correlation functions of the charge and currents,

$$\Phi_{C;jj'}(\tau; \tau_1, \tau_2) = \left\langle T_\tau \delta\rho_C(\tau) c_{j\sigma}(\tau_1) c_{j'\sigma}^\dagger(\tau_2) \right\rangle \tag{123}$$

$$\Phi_{L;jj'}(\tau; \tau_1, \tau_2) = \left\langle T_\tau J_L(\tau) c_{j\sigma}(\tau_1) c_{j'\sigma}^\dagger(\tau_2) \right\rangle \tag{124}$$

$$\Phi_{R;jj'}(\tau; \tau_1, \tau_2) = \left\langle T_\tau J_R(\tau) c_{j\sigma}(\tau_1) c_{j'\sigma}^\dagger(\tau_2) \right\rangle \tag{125}$$

where $\delta\rho_C \equiv \rho_C - \langle \rho_C \rangle$. These three functions can be expressed as functions of two Matsubara frequencies $i\nu$ and $i\varepsilon$,

$$\Phi_{\gamma;jj'}(\tau; \tau_1, \tau_2) = \frac{1}{\beta^2} \sum_{i\varepsilon, i\nu} \Phi_{\gamma;jj'}(i\varepsilon, i\varepsilon + i\nu) e^{-i\varepsilon(\tau_1 - \tau)} e^{-i(\varepsilon+\nu)(\tau - \tau_2)} \tag{126}$$

for $\gamma = C, L, R$. We mainly consider the electrons in the central region assuming $jj' \in C$. In the right-hand side of Eqs. (124) and (125), there still exist the creation and annihilation operators with respect to the leads at $0$ and $N+1$ in the operators $J_L$ and $J_R$. The correlation functions that include these two sites as one of the external points can be related to those defined with respect to the adjacent sites $1$ and $N$, by using the properties of the Green's function at the two interfaces:

$$\begin{cases} G_{0,j}(z) = -g_L(z) v_L G_{1,j}(z) & \text{for} \quad 1 \le j \le N+1 \\ G_{j,N+1}(z) = -G_{j,N}(z) v_R g_R(z) & \text{for} \quad 0 \le j \le N \end{cases} \tag{127}$$

Here, $g_L(z)$ and $g_R(z)$ are the local Green's functions at the interfaces of the isolated leads, $0$ and $N+1$, respectively. Using these properties, the three-point correlation functions for

$jj' \in C$ can be expressed as

$$\Phi_{\gamma;jj'}(i\varepsilon, i\varepsilon + i\nu) = \sum_{j_4 j_1, t} G_{j_4}(i\varepsilon) \Lambda_{\gamma;j_4 j_1}(i\varepsilon, i\varepsilon + i\nu) G_{j_1 j'}(i\varepsilon + i\nu) \tag{128}$$

where $\Lambda_{\gamma;j_4 j_1}$ includes all the vertex corrections. The corresponding bare current vertices are given by

$$\Lambda^{(0)}_{C;j_4 j_1}(i\varepsilon, i\varepsilon + i\nu) = \delta_{j_4 j_1} \tag{129}$$

$$\Lambda^{(0)}_{L;j_4 j_1}(i\varepsilon, i\varepsilon + i\nu) = \lambda_L(i\varepsilon, i\varepsilon + i\nu)\delta_{i,j_4}\delta_{i,j_1} \tag{130}$$

$$\Lambda^{(0)}_{R;j_4 j_1}(i\varepsilon, i\varepsilon + i\nu) = \lambda_R(i\varepsilon, i\varepsilon + i\nu)\delta_{N,j_4}\delta_{N,j_1} \tag{131}$$

with

$$\lambda_L(i\varepsilon, i\varepsilon + i\nu) = -iv_l^2[g_l(i\varepsilon + i\nu) - g_l(i\varepsilon)] \tag{132}$$

$$\lambda_R(i\varepsilon, i\varepsilon + i\nu) = iv_R^2[g_R(i\varepsilon + i\nu) - g_R(i\varepsilon)] \tag{133}$$

We now calculate the $\omega$-linear part of $K_{\alpha\alpha'}(z)$ taking $\alpha$ and $\alpha'$ to be $R$ and $L$, respectively. Using the three-point correlation functions, the current–current correlation function $K_{RL}(i\nu)$ can be expressed as

$$K_{RL}(i\nu) = \frac{1}{\beta}\sum_{j\varepsilon}\sum_{\sigma}\lambda_L(i\varepsilon, i\varepsilon + i\nu)\Phi_{R;11}(i\varepsilon, i\varepsilon + i\nu) \tag{134}$$

Paying attention to the analytic properties of the Green's functions, the summation over the Matsubara frequency can be rewritten in a contour-integral form. Then, by carrying out the analytic continuation $i\nu \to \omega + i0^+$, we obtain

$$K_{RL}(\omega + i0^+) = \sum_{\sigma}\left\{ -\int_{\sim}^{\sim} \frac{d\epsilon}{2\pi i} f(\epsilon) \lambda_L^{[1]}(\epsilon, \epsilon + \omega)\Phi_{R;11}^{[1]}(\epsilon, \epsilon + \omega)\right.$$

$$-\int_{-\infty}^{\infty}\frac{d\epsilon}{2\pi i}\left[f(\epsilon + \omega) - f(\epsilon)\right]\lambda_L^{[2]}(\epsilon, \epsilon + \omega)\Phi_{R;11}^{[2]}(\epsilon, \epsilon + \omega)$$

$$\left. +\int_{-\infty}^{\infty}\frac{d\epsilon}{2\pi i} f(\epsilon + \omega)\lambda_L^{[3]}(\epsilon, \epsilon + \omega)\Phi_{R;11}^{[3]}(\epsilon, \epsilon + \omega)\right\} \tag{135}$$

where $f(\epsilon) = (e^{\beta\epsilon} + 1)^{-1}$. The superscript with the bracket, that is, $[k]$ for $k = 1, 2, 3$, is introduced specifies the three analytic region of $\Phi_{R;11}(z, z + w)$ and $\lambda_L(z, z + w)$ in the complex $z$-plane. These regions are separated by the two lines, $\mathrm{Im}(z) = 0$ and $\mathrm{Im}(z + w) = 0$, as shown in Fig. 5. In each of the three regions, $\Phi_{R;11}(z, z + w)$ corresponds to the analytic



Figure 5. Three analytic regions of $\Phi_{11}(z, z + w)$.

function given by

$$
\begin{cases}
\Phi_{R;11}^{[1]}(\epsilon, \epsilon + \omega) = \Phi_{R;11}(\epsilon + i0^+, \epsilon + \omega + i0^-) \\[2mm]
\Phi_{R;11}^{[2]}(\epsilon, \epsilon + \omega) = \Phi_{R;11}(\epsilon - i0^+, \epsilon + \omega + i0^+) \\[2mm]
\Phi_{R;11}^{[3]}(\epsilon, \epsilon + \omega) = \Phi_{R;11}(\epsilon - i0^-, \epsilon + \omega - i0^+)
\end{cases}
\tag{136}
$$

These analytic properties can be clarified explicitly in the Lehmann representation, Eq. (158), provided in the next subsection. Similarly, the analytic continuation of the bare current vertex $\lambda_\alpha(i\varepsilon, i\varepsilon + i\nu)$ for $\alpha = L, R$ is given by

$$
\begin{cases}
\lambda_\alpha^{[1]}(\epsilon, \epsilon + \omega) = is_\alpha v_\alpha^2[g_\alpha^+(\epsilon + \omega) - g_\alpha^+(\epsilon)] \\[2mm]
\lambda_\alpha^{[2]}(\epsilon, \epsilon + \omega) = is_\alpha v_\alpha^2[g_\alpha^+(\epsilon + \omega) - g_\alpha^-(\epsilon)] \\[2mm]
\lambda_\alpha^{[3]}(\epsilon, \epsilon + \omega) = is_\alpha v_\alpha^2[g_\alpha^-(\epsilon + \omega) - g_\alpha^-(\epsilon)]
\end{cases}
\tag{137}
$$

where the factor $s_\alpha$ is defined such that $s_L = -1$ and $s_R = +1$. In this section, we distinguish the retarded and advanced Green's functions by the label $+$ and $-$, respectively, in the superscript. In the limit of $\omega \to 0$, the bare vertices for $k = 1$ and 3 vanish as $\lambda_\alpha^{[k]}(\epsilon, \epsilon + \omega) \propto \omega$. In contrast, for $k = 2$, it tends to a finite constant $\lambda_\alpha^{[2]}(\epsilon, \epsilon) = 2s_\alpha\Gamma_\alpha(\epsilon)$ with $\Gamma_\alpha(\epsilon) = -v_\alpha^2\mathrm{Im}[g_\alpha^+(\epsilon)]$. Correspondingly, the asymptotic behavior of $\Phi_\alpha^{[k]}(\epsilon, \epsilon + \omega)$ for small $\omega$ has been investigated by using the Lehmann representation of a four-point vertex function [22, 23], and the result is [19],

$$
\Phi_\alpha^{[k]}(\epsilon, \epsilon + \omega) \propto
\begin{cases}
\omega & \text{for} \quad k = 1 \\
\text{finite} & \text{for} \quad k = 2 \\
\omega & \text{for} \quad k = 3
\end{cases}
\tag{138}
$$

for $\alpha = L, R$. Therefore, taking the $\omega \to 0$ limit in Eq. (117) by using Eq. (135) for $K_{RL}(\omega + i0^+)$, we obtain

$$
g = \frac{2e^2}{h} \int d\epsilon \left( -\frac{\partial f}{\partial \epsilon} \right) \mathcal{T}(\epsilon)
\tag{139}
$$

$$
\mathcal{T}(\epsilon) = 2\Gamma_L(\epsilon)\Phi_{R;11}^{[2]}(\epsilon, \epsilon)
\tag{140}
$$

Thus, the dc conductance is determined by the three-point function for the analytic region $k = 2$. The analytic continuation is performed formally by using the Lehmann representation in Section 4.3. The result shows that $\Phi_{R;11}^{[2]}(\epsilon, \epsilon)$ can be expressed as a Fourier transform, Eq. (162), of a real-time retarded product in Eq. (160).

Specifically, at $T = 0$ the conductance is determined by the value of the transmission probability at the Fermi $\epsilon = 0$, and it can be written in the form [24–27],

$$
\mathcal{T}(0) = 4 \, \Gamma_L(0)G_{1N}^-(0)\Gamma_R(0)G_{N1}^+(0)
\tag{141}
$$

This is due to the property that the vertex corrections for the current vanishes at $T = 0$ and $\epsilon = 0$, as shown in Eq. (154) in Section 4.2. Furthermore, the reflection probability is given by

$$
\mathcal{R}(0) = |1 - 2i\Gamma_L(0)G_{11}^+(0)|^2 = |1 - 2i\Gamma_R(0)G_{NN}^+(0)|^2
\tag{142}
$$

The current conservation $\mathcal{T}(0) + \mathcal{R}(0) = 1$ follows from the identity in Eq. (156). Similarly, at zero temperature, the Friedel sum rule for interacting electrons is given by [28],

$$
\Delta N_{\mathrm{tot}} = \frac{1}{\pi i} \log[\det S]
\tag{143}
$$

where the $S$-matrix is defined by

$$S = \begin{bmatrix} 1 - 2i\Gamma_L(0)G_{11}^+(0) & -2i\Gamma_L(0)G_{1N}^+(0) \\ -2i\Gamma_R(0)G_{N1}^+(0) & 1 - 2i\Gamma_R(0)G_{NN}^+(0) \end{bmatrix} \tag{144}$$

In Eq. (143), $\Delta N_{tot}$ is the displacement of the total charge defined by

$$\Delta N_{tot} = \sum_{i \in C} \sum_\sigma \langle c_{i\sigma}^\dagger c_{i\sigma} \rangle$$

$$+ \sum_{i \in L} \sum_\sigma [\langle c_{i\sigma}^\dagger c_{i\sigma} \rangle - \langle c_{i\sigma}^\dagger c_{i\sigma} \rangle_L] + \sum_{i \in R} \sum_\sigma [\langle c_{i\sigma}^\dagger c_{i\sigma} \rangle - \langle c_{i\sigma}^\dagger c_{i\sigma} \rangle_R] \tag{145}$$

where $\langle \cdots \rangle_L$ and $\langle \cdots \rangle_R$ denote the ground-state average of isolated leads determined by $\mathcal{H}_L$ and $\mathcal{H}_R$, respectively.

## 4.2. Current Conservation and Ward Identity

The interelectron interactions generally cause the damping of excitations. Therefore, theoretically, the self-energy and vertex corrections must be treated consistently with the approaches that conserve the current. In this subsection, we discuss the current conservation using a generalized Ward identity.

The generalized Ward identity can be derived from the equation of continuity in the Matsubara form $-(\partial/\partial\tau)\delta\rho_\ell + iJ_R - iJ_L = 0$ [29],

$$-\frac{\partial}{\partial\tau}\Phi_{\ell,jj'}(\tau; \tau_1, \tau_2) + i\Phi_{R,jj'}(\tau; \tau_1, \tau_2) - i\Phi_{L,jj'}(\tau; \tau_1, \tau_2)$$

$$= \delta(\tau - \tau_2)G_{jj'}(\tau_1, \tau) - \delta(\tau_1 - \tau)G_{jj'}(\tau, \tau_2) \tag{146}$$

It can be expressed by using a $N \times N$ matrix representation for $jj' \in C$ with the Matsubara frequencies,

$$i\nu\,\Phi_\ell(i\varepsilon, i\varepsilon + i\nu) + i\,\Phi_R(i\varepsilon, i\varepsilon + i\nu) - i\,\Phi_L(i\varepsilon, i\varepsilon + i\nu) = G(i\varepsilon) - G(i\varepsilon + i\nu) \tag{147}$$

Here, $G(z) = \{G_{jj'}(z)\}$ and $\Phi_\gamma(z, z+w) = \{\Phi_{\gamma;jj'}(z, z+w)\}$. The matrix version of Eq. (128) is given by

$$\Phi_\gamma(z, z+w) = G(z)\Lambda_\gamma(z, z+w)G(z+w) \tag{148}$$

Thus, the identity can also be expressed using $\Lambda_\gamma(z, z+w) = \{\Lambda_{\gamma;jj'}(z, z+w)\}$, as

$$i\nu\Lambda_\ell(i\varepsilon, i\varepsilon + i\nu) + i\Lambda_R(i\varepsilon, i\varepsilon + i\nu) - i\Lambda_L(i\varepsilon, i\varepsilon + i\nu)$$

$$= \{G(i\varepsilon + i\nu)\}^{-1} - \{G(i\varepsilon)\}^{-1} \tag{149}$$

Furthermore, the Dyson equation for the single-particle Green's function can be expressed as

$$\{G(z)\}^{-1} = z1 - \mathcal{H}_C^0 - \mathcal{V}_{mix}(z) - \Sigma(z) \tag{150}$$

$$\mathcal{H}_C^0 = \begin{bmatrix} -t_{11}^C - \mu & -t_{12}^C & \cdots & \\ -t_{21}^C & -t_{22}^C - \mu & & \\ \vdots & & \ddots & \\ & & & -t_{NN}^C - \mu \end{bmatrix} \tag{151}$$

$$\mathcal{V}_{mix}(z) = \begin{bmatrix} v_L^2 g_L(z) & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & v_R^2 g_R(z) \end{bmatrix} \tag{152}$$

and $\Sigma(z) = \{\Sigma_{ij}(z)\}$ is the self-energy due to the interelectron interactions. Therefore, Eq. (149) represents a relation between the self-energy and vertex functions, and this identity must be satisfied in the conserving approaches. Carrying out the analytic continuation of Eq. (149) in the region $k = 2$, $i\varepsilon + i\nu \to \epsilon + \omega + i0^+$ and $i\varepsilon \to \epsilon - i0^-$, and then taking the limit of $\omega \to 0$, we obtain

$$\Lambda_R^{[2]}(\epsilon, \epsilon) - \Lambda_L^{[2]}(\epsilon, \epsilon) = -2\mathrm{Im}\,\boldsymbol{V}_{\mathrm{mix}}^+(\epsilon) - 2\,\mathcal{I}m\,\boldsymbol{\Sigma}^+(\epsilon) \tag{153}$$

At $T = 0$, $\epsilon = 0$, the imaginary part of the self-energy vanishes $\mathrm{Im}\boldsymbol{\Sigma}^+(0) = 0$, and then the current vertices become equal to the bare ones,

$$\Lambda_{R:j_1 j_1}^{[2]}(0,0) = 2\Gamma_R(0)\delta_{Nj_1}\delta_{Nj_1} \tag{154}$$

$$\Lambda_{L:j_4 j_1}^{[2]}(0,0) = -2\Gamma_L(0)\delta_{1j_4}\delta_{1j_1} \tag{155}$$

Correspondingly, $\Phi_{R:11}^{[2]}(0,0) = G_{1N}^-(0)2\Gamma_R(0)G_{N1}^+(0)$ at $T = 0$, and then the transmission probability is given by Eq. (141). Alternatively, at $T = 0$, the analytic continuation of Eq. (146) in region $k = 2$ is written in the form

$$G^+(0) - G^-(0) = G^+(0)[\boldsymbol{V}_{\mathrm{mix}}^+(0) - \boldsymbol{V}_{\mathrm{mix}}^-(0)]G^-(0) \tag{156}$$

and the $(1, 1)$ and $(N,N)$ matrix elements represent the optical theorem for Eqs. (141) and (142).

Particularly, for the single Anderson impurity at $N = 1$, Eq. (147) becomes a single-component equation with respect to the impurity site. In the region $k = 2$, it becomes $\Phi_R^{[2]}(\epsilon, \epsilon) - \Phi_L^{[2]}(\epsilon, \epsilon) = G^-(\epsilon) - G^+(\epsilon)$ in the limit of $\omega \to 0$. Furthermore, if the mixing terms have a property $\Gamma_L(\epsilon) = \lambda\Gamma_R(\epsilon)$, an additional relation $\Phi_L^{[2]}(\epsilon, \epsilon) = -\lambda\Phi_R^{[2]}(\epsilon, \epsilon)$ follows. Thus, in this case the dc conductance can be written in the form [14, 15],

$$g_{\mathrm{single}} = \frac{2e^2}{h} \int_{-\infty}^{\infty} d\epsilon \left(-\frac{\partial f}{\partial \epsilon}\right) \frac{4\Gamma_L\Gamma_R}{\Gamma_R + \Gamma_L} [-\mathrm{Im}G^+(\epsilon)] \tag{157}$$

## 4.3. Lehmann Representation for $\mathcal{T}(\epsilon)$

We now show that the transmission probability $\mathcal{T}(\epsilon)$ can be expressed in terms of a real-time retarded product in Eq. (160) via the Fourier transform Eq. (162). It shows a direct link between the transmission probability and dynamic correlation functions. To prove it, we first of all derive the Lehmann representation for $\Phi_{R:11}(i\varepsilon, i\varepsilon + i\nu)$, and then carry out the analytical continuation.

Inserting a complete set of the eigenstates, $\mathcal{H}|n\rangle = E_n|n\rangle$, into Eq. (125) and using Eq. (126), we obtain

$$\Phi_{R:11}(i\varepsilon, i\varepsilon + i\nu)$$

$$= \frac{1}{Z} \sum_{lmn} \langle l|c_{1\sigma}^\dagger|m\rangle\langle m|J_R|n\rangle\langle n|c_{1\sigma}|l\rangle \left[ \frac{e^{-\beta E_m}}{(i\varepsilon + i\nu + E_m - E_l)(i\nu + E_m - E_n)} \right.$$

$$\left. - \frac{e^{-\beta E_l}}{(i\varepsilon + E_n - E_l)(i\varepsilon + i\nu + E_m - E_l)} - \frac{e^{-\beta E_n}}{(i\nu + E_m - E_n)(i\varepsilon + E_n - E_l)} \right]$$

$$+ \frac{1}{Z} \sum_{lmn} \langle l|c_{1\sigma}|n\rangle\langle n|J_R|m\rangle\langle m|c_{1\sigma}^\dagger|l\rangle \left[ \frac{e^{-\beta E_n}}{(i\varepsilon + E_l - E_n)(i\nu + E_n - E_m)} \right.$$

$$\left. + \frac{e^{-\beta E_l}}{(i\varepsilon + E_l - E_n)(i\varepsilon + i\nu + E_l - E_m)} - \frac{e^{-\beta E_m}}{(i\varepsilon + i\nu + E_l - E_m)(i\nu + E_n - E_m)} \right] \tag{158}$$

where $Z = \mathrm{Tr}\,e^{-\beta\mathcal{H}}$. From Eq. (158), the analytic continuation to obtain $\Phi_{R:11}^{[k]}(\epsilon, \epsilon + \omega)$ for $k = 1, 2, 3$ can be carried out by replacing the imaginary frequencies $i\varepsilon$ and $i\nu$

by the real ones $\epsilon$ and $\omega$, respectively, with the infinitesimal imaginary parts shwn in Eq. (136). Then the same expressions for $\Phi_{R;11}^{[k]}(\epsilon, \epsilon + \omega)$ can be derived from the real-time functions

$$\Phi_{R;11}^{[1]}(t; t_1, t_2) = \theta(t - t_1)\theta(t_1 - t_2)\langle\left[\left\{c_{1\sigma}(t_1), c_{1\sigma}^{\dagger}(t_2)\right\}, J_R(t)\right]\rangle$$

$$+\theta(t_1 - t)\theta(t - t_2)\langle\left|c_{1\sigma}(t_1), \left[c_{1\sigma}^{\dagger}(t_2), J_R(t)\right]\right|\rangle \qquad (159)$$

$$\Phi_{R;11}^{[2]}(t; t_1, t_2) = \theta(t - t_1)\theta(t_1 - t_2)\langle\left|c_{1\sigma}^{\dagger}(t_2), \left[c_{1\sigma}(t_1), J_R(t)\right]\right|\rangle$$

$$-\theta(t - t_2)\theta(t_2 - t_1)\langle\left|c_{1\sigma}(t_1), \left[c_{1\sigma}^{\dagger}(t_2), J_R(t)\right]\right|\rangle \qquad (160)$$

$$\Phi_{R;11}^{[3]}(t; t_1, t_2) = -\theta(t - t_2)\theta(t_2 - t_1)\langle\left[\left\{c_{1\sigma}(t_1), c_{1\sigma}^{\dagger}(t_2)\right\}, J_R(t)\right]\rangle$$

$$-\theta(t_2 - t)\theta(t - t_1)\langle\left|c_{1\sigma}^{\dagger}(t_2), \left[c_{1\sigma}(t_1), J_R(t)\right]\right|\rangle \qquad (161)$$

where $J_R(t) \equiv e^{iHt}J_R e^{-iHt}$, and $\theta(t)$ is the step function. The commutators are defined by $[A, B] \equiv AB - BA$, and $\{A, B\} \equiv AB + BA$, as usual. The Fourier transform into the real frequencies is given by

$$\int_{-\infty}^{\infty} dt\, dt_1\, dt_2\, e^{i\omega t}\, e^{i\epsilon t_1}\, e^{-i\epsilon' t_2}\,\Phi_{R;11}^{[k]}(t; t_1, t_2)$$

$$= 2\pi\delta(\epsilon + \omega - \epsilon')\Phi_{R;11}^{[k]}(\epsilon, \epsilon + \omega) \qquad (162)$$

For example, a time-ordered function

$$F(t; t_1, t_2) = \theta(t - t_1)\theta(t_1 - t_2)\langle J_R(t)c_{1\sigma}(t_1)c_{1\sigma}^{\dagger}(t_2)\rangle \qquad (163)$$

is transformed into

$$F(\epsilon, \epsilon + \omega) = \frac{-1}{Z}\sum_{lmn}\frac{e^{-\beta E_m}\langle l|c_{1\sigma}^{\dagger}|m\rangle\langle m|J_R|n\rangle\langle n|c_{1\sigma}|l\rangle}{(\epsilon + \omega + E_m - E_l + i0^+)(\omega + E_m - E_n + i0^+)} \qquad (164)$$

Among the three real-time functions Eqs. (159)–(161), the function for the regin $k = 2$, that is, $\Phi_{R;11}^{[2]}(t; t_1, t_2)$ in Eq. (160) determines the transmission probability $\tau(\epsilon) = 2\Gamma_L(\epsilon)\Phi_{R;11}^{[2]}(\epsilon, \epsilon)$. Because the analytic continuation has already been done, the real-time correlation links directly to the transport coefficient. This formulation can be used for numerical calculations.

## 4.4. Application to a Hubbard Chain Connected to Leads

In this subsection, we apply the linear-response formulation to a finite Hubbard chain attached to reservoirs, which can be considered as a model for a series of quantun dots or atomic wires of nanometer size. A schematic picture of the model is shown in Fig. 6. The Hamiltonian parameters defined in Eq. (1) are taken as follows. We take $t_{ij}^C$ tobe the nearest-neighbor hopping $t$, and $U_{u;v;h}$ to be an on-site repulsion $U$. Specifically, we consider the electron–hole symmetric case, at which $\mu = 0$ and $\epsilon_0 + U/2 = 0$ with $\epsilon_0$ beng the onsite energy. We also assume that the two couplings are symmetric $\Gamma_L = \Gamma_R (\equiv \Gamma)$, and the local density of states of the leads is a constant.



Figure 6. Schematic picture of a finite Hubbard chain.

**Figure 7.** The order $U^2$ terms of (a) self-energy and (b) vertex corrections.

To examine the effects of the Coulomb interaction, we calculate the self-energy and vertex corrections up to terms of order $U^2$, the Feynman diagrams for which are illustrated in Fig. 7 [19]. These contributions satisfy the generalized Ward identity Eq. (153) that corresponds to the current conservation law. In Fig. 8, the results of $\mathcal{T}(\epsilon)$ for $N = 3, 4$ are plotted versus $\epsilon/t$ for $\Gamma/t = 0.75$ for three values of $U/(2\pi t)$: (—) 0.0, (–o–) 0.5, and (–•–) 1.0. The temperature $T/t$ is taken to be (a) 0.0, (b) 0.2 for $N = 3$ in the upper panels, and (c) 0.0, (d) 0.2 for $N = 4$ in the lower panels.

At low temperatures, there are $N$ resonance peaks that have one-to-one correspondence to resonant states of the unperturbed system. In addition to these resonance peaks, two broad peaks of atomic character appear at $\epsilon \simeq \pm U/2$ for large $U$. The resonance peaks become sharper with increasing $U$ at low temperatures as seen in the panels (a) and (c). However, the height of the peaks decreases with increasing $U$. One exception, which happens for odd $N$, is the Kondo resonance at the Fermi level $\epsilon = 0$. At this peak, the transmission probability reaches the unitary-limit value 1.0 for any values of $U$, when the systems have the inversion symmetry $\Gamma_L = \Gamma_R$ together with the electron-hole symmetry [27, 30]. The width of the Kondo resonance $T_K$ must decrease with increasing $N$. For even $N$, the transmission probability $\mathcal{T}(\epsilon)$ shows a minimum at $\epsilon = 0$. The characteristic energy scale in this case is the width of the valley, which eventually becomes the Mott-Hubbard gap in the limit of large $N$. The high energy profile of $\mathcal{T}(\epsilon)$ at $|\epsilon| \gtrsim 2t$ in the case of $N = 3$ is similar to that for $N = 4$. Namely, the high-energy part shows no notable $N$ dependence. For $U/(2\pi t) = 0.5$, the upper and lower Hubbard levels at $\epsilon \simeq \pm U/2$ exist inside the energy region corresponding to the one-dimensional band of the width $2t$. The two Hubbard levels got outside of this energy region for $U/2 \gtrsim 2t$. At finite temperatures, the resonance peaks at $|\epsilon| \lesssim 2t$ become



Transmission probability of N=3, 4 Hubbard model

**Figure 8.** Many-body transmission coefficient for $N = 3$ (upper two panels) and 4 (lower two panels) is plotted versus $\epsilon/t$ for $\Gamma/t = 0.75$ for three values of $U/(2\pi t)$: (—) 0.0, (–o–) 0.5, and (–•–) 1.0. The temperature is taken to be $T/t = 0.0$ for (a) and (c), and $T/t = 0.2$ for (b) and (d).

broad. and the peak height decreases with increasing $\Gamma$. The structures of the resonance peaks vanish eventually at higher temperatures, and then the even-odd oscillatory behavior disappears [19].

## 5. TOMONAGA-LUTTINGER MODEL

Transport through interacting systems in one dimension has been studied extensively for quantum wires, organic conductors, carbon nanotube, etc. In this section, we provide a brief introduction to a Tomonaga-Luttinger model [31], to take a quick look at the transport properties of a typical interacting system in one dimension.

### 5.1. Spin-Less Fermions in One Dimension

We start with the spin-less fermions described by the Hamiltonian.

$$H_0 - \langle H_0 \rangle_0 = \sum_k (\epsilon_k - \mu)\left[c_k^\dagger c_k - \langle c_k^\dagger c_k \rangle_0\right] \tag{165}$$

$$H_I = \frac{1}{2L} \sum_{qkk'} V_q c_{k+q}^\dagger c_{k'-q}^\dagger c_{k'} c_k \tag{166}$$

In Eq. (165), the ground-state energy for the noninteracting electrons has been subtracted. At low energies, the excitations near the Fermi level play a dominant role, so that $\epsilon_k$ can be linearized at the two Fermi points $k = \pm k_F$,

$$H_0 = \sum_k v_F(k - k_F)\left[a_k^\dagger a_k - \langle a_k^\dagger a_k \rangle_0\right]$$

$$+ \sum_k v_F(-k - k_F)\left[b_k^\dagger b_k - \langle b_k^\dagger b_k \rangle_0\right] \tag{167}$$

Here, $a_k$ ($b_k$) is the operator for the right-moving (left-moving) particles. The summation over $k$ in Eq. (167) should be restricted in a range $|k| - k_F < k_c$ with the cutoff momentum $k_c$ of the order a band width $D \sim v_F k_c$ as illustrated in Fig. 9. However, we assume that $k_c \to \infty$, and will introduce the cut-off for the momentum transfer $q$, when it is required [31]. For low-energy properties, the interactions between the electrons near the Fermi level is important. Therefore, the interaction Hamiltonian Eq. (166) can be simplified by taking only the scattering processes in which all the four momentums are close to one of the two Fermi points, that is, $k + q \simeq \pm k_F$, $k' - q \simeq \pm k_F$, $k' \simeq \pm k_F$, and $k \simeq \pm k_F$, into account;

$$H_I \simeq \frac{1}{2L}\left(\sum_{q \sim 0} V_q + \sum_{q \sim \pm 2k_F} V_q\right)\left(\sum_{k \simeq k_F} + \sum_{k \simeq -k_F}\right)\left(\sum_{k' \simeq k_F} + \sum_{k' \simeq -k_F}\right)c_{k+q}^\dagger c_{k'-q}^\dagger c_{k'} c_k \tag{168}$$



Figure 9. Linearized dispersion.

For the scattering process with small momentum transfer $q \simeq 0$, there are two types of possibilities for the initial momentums $k \simeq k'$ and $k \simeq -k'$. In contrast, in the case of the back scattering $q \simeq \pm 2k_F$, the incident momentums must have the opposite sign $k \simeq -k'$. These scattering processes can be described by a simplified Hamiltonian,

$$H_I \Rightarrow \mathcal{H}_4 + \mathcal{H}_2 + \mathcal{H}_1 \tag{169}$$

$$\mathcal{H}_4 = \frac{g_4}{2L} \sum_{qkk'} a^\dagger_{k+q} a^\dagger_{k'-q} a_{k'} a_k + \frac{g_4}{2L} \sum_{qkk'} b^\dagger_{k+q} b^\dagger_{k'-q} b_{k'} b_k \tag{170}$$

$$\mathcal{H}_2 = \frac{g_2}{2L} \sum_{qkk'} a^\dagger_{k-q} b^\dagger_{k'-q} b_{k'} a_k + \frac{g_2}{2L} \sum_{qkk'} b^\dagger_{k+q} a^\dagger_{k'-q} a_{k'} b_k \tag{171}$$

$$\mathcal{H}_1 = \frac{g_1}{2L} \sum_{q'kk'} b^\dagger_{k+q'-2k_F} a^\dagger_{k'+q'+2k_F} b_{k'} a_k + \frac{g_1}{2L} \sum_{q'kk'} a^\dagger_{k+q'+2k_F} b^\dagger_{k'-q'-2k_F} a_{k'} b_k \tag{172}$$

In Eq. (172), $q'$ is a small momentum defined such that $q = q' \pm 2k_F$. The coupling constants should be taken as $g_2 \simeq g_4 \simeq V_0$, and $g_1 \simeq V_{2k_F}$. However, in the following we treat these three constants to be independent parameters. The momentum-transfer cutoff is introduced for the summation over $q$ and $q'$. The Tomonaga-Luttinger model is defined by

$$\mathcal{H}_{TL} = \mathcal{H}_0 + \mathcal{H}_4 + \mathcal{H}_2 \tag{173}$$

Note that for the spin-less model there should be no distinction between $\mathcal{H}_2$ and $\mathcal{H}_1$ [31].

The interactions $\mathcal{H}_4$ and $\mathcal{H}_2$ can be expressed in terms of the density operators $\rho_1(p)$ and $\rho_2(p)$ defined by

$$\rho_1(p) = \sum_k \left[ a^\dagger_{k+p} a_k - \delta_{p,0} \langle a^\dagger_k a_k \rangle_0 \right] \tag{174}$$

$$\rho_2(p) = \sum_k \left[ b^\dagger_{k+p} b_k - \delta_{p,0} \langle b^\dagger_k b_k \rangle_0 \right] \tag{175}$$

Here, $\langle a^\dagger_k a_k \rangle_0 = \theta(k_F - k)$ and $\langle b^\dagger_k b_k \rangle_0 = \theta(k_F + k)$, and these terms are required to define the deviation from the noninteracting value without an ambiguity caused by the occupation of the negative energy states. Equations (170) and (171) can be rewritten in the following forms apart from a renormalization of the chemical potential that can be absorbed in to $k_F$.

$$\mathcal{H}_4 = \frac{g_4}{L} \sum_{p>0} \rho_1(-p)\rho_1(p) + \frac{g_4}{L} \sum_{p>0} \rho_2(-p)\rho_2(p) \tag{176}$$

$$\mathcal{H}_2 = \frac{g_2}{L} \sum_{p>0} \rho_1(-p)\rho_2(p) + \frac{g_2}{L} \sum_{p>0} \rho_2(-p)\rho_1(p) \tag{177}$$

These two density operators satisfy the commutation relations

$$[\rho_1(p), \rho_1(-p')] = \frac{Lp}{2\pi} \delta_{pp'}, \qquad [\rho_2(-p), \rho_2(p')] = \frac{Lp}{2\pi} \delta_{pp'} \tag{178}$$

One notable feature is that these two commutation relations are equivalent to those of the bose operators,

$$C_p = \sqrt{\frac{2\pi}{Lp}} \, \rho_1(p), \qquad C^*_p = \sqrt{\frac{2\pi}{Lp}} \, \rho_1(-p) \tag{179}$$

$$C_{-p} = \sqrt{\frac{2\pi}{Lp}} \, \rho_2(-p), \qquad C^*_{-p} = \sqrt{\frac{2\pi}{Lp}} \, \rho_2(p) \tag{180}$$

$$[C_p, C^*_{p'}] = \delta_{pp'}, \qquad [C_{-p}, C^*_{-p'}] = \delta_{pp'} \tag{181}$$

where $p > 0$. The commutation relation of the density operators and Hamiltonian can be calculated by using Eqs. (167) and (174)-(177) as

$$\left[\rho_1(p), \mathcal{H}_0\right] = v_F p\rho_1(p), \qquad \left[\rho_2(-p), \mathcal{H}_0\right] = v_F p\rho_2(-p) \tag{182}$$

$$\left[\rho_1(p), \mathcal{H}_4\right] = \tilde{g}_4 v_F p\rho_1(p), \qquad \left[\rho_2(-p), \mathcal{H}_4\right] = \tilde{g}_4 v_F p\rho_2(-p) \tag{183}$$

$$\left[\rho_1(p), \mathcal{H}_2\right] = \tilde{g}_2 v_F p\rho_2(p), \qquad \left[\rho_2(-p), \mathcal{H}_2\right] = \tilde{g}_2 v_F p\rho_1(-p) \tag{184}$$

where $\tilde{g}_4 = g_4/(2\pi v_F)$, and $\tilde{g}_2 = g_2/(2\pi v_F)$.

## 5.2. Two Conservation Laws

The operator for the charge and current are defined by

$$\rho_c(p) = \rho_1(p) + \rho_2(p), \qquad \rho_J(p) = \rho_1(p) - \rho_2(p) \tag{185}$$

In the real space, the operators for the left and right movers $\nu = 1, 2$ are written in the form

$$\rho_\nu(x) = \frac{1}{L} \sum_{p>0} (\rho_\nu(p)e^{ipx} + \rho_\nu(-p)e^{-ipx}) \tag{186}$$

The equation of motion for $\rho_c(x)$ and $\rho_J(x)$ are derived from the Heisenberg equation using the commutation relations in Eqs. (182)-(184).

$$\frac{\partial}{\partial t}\rho_c(x, t) + v_J \frac{\partial}{\partial x}\rho_J(x, t) = 0, \qquad v_J = v_F(1 + \tilde{g}_4 - \tilde{g}_2) \tag{187}$$

$$\frac{\partial}{\partial t}\rho_J(x, t) + v_N \frac{\partial}{\partial x}\rho_c(x, t) = 0, \qquad v_N = v_F(1 + \tilde{g}_4 + \tilde{g}_2) \tag{188}$$

Because there are two independent equations for $\rho_c(x, t)$ and $\rho_J(x, t)$, the explicit form of these Heisenberg operators can be calculated analytically;

$$\left(\frac{\partial^2}{\partial t^2} - v_\rho^2 \frac{\partial^2}{\partial x^2}\right)\rho_c(x, t) = 0, \qquad \left(\frac{\partial^2}{\partial t^2} - v_\rho^2 \frac{\partial^2}{\partial x^2}\right)\rho_J(x, t) = 0, \tag{189}$$

$$v_\rho^2 = v_J v_N \tag{190}$$

The relation among the three velocities $v_J$, $v_N$, and $v_\rho$ can be summarized as

$$v_J = K_\rho v_\rho, \qquad v_N = \frac{v_\rho}{K_\rho}, \qquad K_\rho \equiv \sqrt{\frac{1 + \tilde{g}_4 - \tilde{g}_2}{1 + \tilde{g}_4 + \tilde{g}_2}} \tag{191}$$

## 5.3. Charge and Current Correlation Functions

Owing to the property shown in Eq. (189), the correlation functions for the density operators

$$\chi'_{\mu\nu}(p, t) = i\frac{1}{L}\theta(t)\langle[\rho_\mu(p, t), \rho_\nu(-p)]\rangle, \qquad \text{for } \mu, \nu = 1, 2 \tag{192}$$

can also be calculated exactly. The equation of motion for these correlations are given by

$$i\frac{\partial}{\partial t}\chi'_{\mu\nu}(p, t) = -\frac{1}{L}\delta(t)\langle[\rho_\mu(p), \rho_\nu(-p)]\rangle - \theta(t)\frac{1}{L}\left\langle\left[\frac{\partial\rho_\mu(p, t)}{\partial t}, \rho_\nu(-p)\right]\right\rangle$$

$$= -\frac{p}{2\pi}\tau_{z,\nu}\delta(t) + i\theta(t)\frac{1}{L}\langle[[\rho_\mu(p, t), \mathcal{H}_1], \rho_\nu(-p)]\rangle \tag{193}$$

where $\tau^3$ is a Pauli matrix:

$$\tau^1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \tau^2 = \begin{bmatrix} 0 & -i \\ i & -0 \end{bmatrix}, \quad \tau^3 = \begin{bmatrix} 1 & -0 \\ 0 & -1 \end{bmatrix}, \quad \mathbf{1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \tag{194}$$

The commutation relation $[\rho_\mu(p, t), \mathcal{H}_{TL}]$ in Eq. (193) can be calculated by using Eqs. (182)–(184). Then, by carrying out the Fourier transform with respect to $t$, we obtain

$$\{\omega \tau^3 - v_\mu p (\cosh \varphi \mathbf{1} + \sinh \varphi \tau^1)\} \chi'(p, \omega) = -\frac{p}{2\pi} \mathbf{1}, \tag{195}$$

$$\cosh \varphi \equiv \frac{1 + \tilde{g}_4}{\sqrt{(1 + \tilde{g}_4)^2 - \tilde{g}_2^2}} = \frac{1}{2}\left(\frac{1}{K_\mu} + K_\mu\right) \tag{196}$$

$$\sinh \varphi \equiv \frac{\tilde{g}_2}{\sqrt{(1 + \tilde{g}_4)^2 - \tilde{g}_2^2}} = \frac{1}{2}\left(\frac{1}{K_\mu} - K_\mu\right) \tag{197}$$

Note that $\cosh \varphi \mathbf{1} + \sinh \varphi \tau^1 = \exp(\varphi \tau^1)$. The Bogoliubov transformation given by $\exp(\varphi \tau^1/2)$ has a property,

$$\tau^3 = \exp\left(\frac{\varphi \tau^1}{2}\right) \tau^3 \exp\left(\frac{\varphi \tau^1}{2}\right) \tag{198}$$

Therefore, Eq. (195) can be diagonalized as

$$\exp\left(\frac{\varphi \tau^1}{2}\right) \{\omega \tau^3 - v_\mu p \mathbf{1}\} \exp\left(\frac{\varphi \tau^1}{2}\right) \chi^r(p, \omega) = -\frac{p}{2\pi} \mathbf{1} \tag{199}$$

With this transformation by $\exp(\varphi \tau^1/2)$, the operators $C_p$ and $C^\dagger_p$ are transformed into

$$\begin{bmatrix} \gamma_p \\ \gamma^\dagger_{-p} \end{bmatrix} = \begin{bmatrix} \cosh\left(\frac{\varphi}{2}\right) & \sinh\left(\frac{\varphi}{2}\right) \\ \sinh\left(\frac{\varphi}{2}\right) & \cosh\left(\frac{\varphi}{2}\right) \end{bmatrix} \begin{bmatrix} C_p \\ C^\dagger_p \end{bmatrix} \tag{200}$$

$$\cosh\left(\frac{\varphi}{2}\right) = \frac{1}{2}\left(\frac{1}{\sqrt{K_p}} + \sqrt{K_p}\right), \quad \sinh\left(\frac{\varphi}{2}\right) = \frac{1}{2}\left(\frac{1}{\sqrt{K_p}} - \sqrt{K_p}\right) \tag{201}$$

where the Bose statistics is preserved for the new operators, $[\gamma_p, \gamma^\dagger_{p'}] = \delta_{pp'}$. The explicit form of $\chi^r(p, \omega)$ is determined by Eq. (199),

$$\chi'(p, \omega) = -\frac{p}{2\pi} \exp\left(-\frac{\varphi \tau^1}{2}\right) \{D^r_+(p, \omega)\tau^3 + D^r_-(p, \omega)\mathbf{1}\} \exp\left(-\frac{\varphi \tau^1}{2}\right)$$

$$= -\frac{p}{2\pi} \{D^r_+(p, \omega)\tau^3 + D^r_-(p, \omega)(\cosh \varphi \mathbf{1} - \sinh \varphi \tau^1)\}, \tag{202}$$

$$D^r_\pm(p, \omega) = \frac{1}{2}\left(\frac{1}{\omega - v_\mu p + i\delta} \pm \frac{1}{\omega + v_\mu p + i\delta}\right) \tag{203}$$

The charge susceptibility $\chi^r_c(p, \omega)$, which corresponds to the $\rho_c$-$\rho_c$ correlation function, is given by

$$\chi^r_c(p, \omega) = \sum_{\mu\nu} \chi^r_{\mu\nu}(p, \omega) = -\frac{K_\mu}{\pi v_\mu} \frac{(v_\mu p)^2}{(\omega + i\delta)^2 - (v_\mu p)^2} \tag{204}$$

Then, the uniform charge susceptibility is given by $\lim_{p \to 0} \chi^r_c(p, 0) = K_\mu/(\pi v_\mu)$. It becomes twice as large for the spin $1/2$ fermions.

The current operator is determined by Eqs. (185)–(187) as

$$J = e v_J \rho_J \tag{205}$$

Therefore, the $J$-$J$ correlation function is given by

$$K^r(p, \omega) = -\frac{e^2 K_p v_p}{\pi} \frac{(v_p p)^2}{(\omega + i\delta)^2 - (v_p p)^2} \tag{206}$$

Then, the conductivity can be calculated with the Kubo formula,

$$\sigma(p, \omega) = \frac{K^r(p, \omega) - K^r(p, 0)}{i\omega}$$

$$= \frac{e^2 K_p v_p}{\pi} \frac{i\omega}{(\omega + i\delta)^2 - (v_p p)^2} \tag{207}$$

The conductivity $\sigma(p, \omega)$ for a uniform $p = 0$ and stationary $\omega = 0$ field depends on the order of taking the limits of $p \to 0$ and $\omega \to 0$. The Drude weight corresponds to the $p \to 0$ limit,

$$\text{Re } \sigma(0, \omega) = e^2 K_p v_p \delta(\omega) \tag{208}$$

In the real space, the conductivity takes the form

$$\sigma(x, \omega) = \int_{-\infty}^{\infty} \frac{dp}{2\pi} \, \sigma(p, \omega) e^{ipx}$$

$$= \int_{-\infty}^{\infty} \frac{dp}{2\pi} \frac{ie^2 K_p v_p}{2\pi} \left[ \frac{1}{\omega - v_p p + i\delta} + \frac{1}{\omega + v_p p + i\delta} \right] e^{ipx}$$

$$= \frac{e^2 K_p}{2\pi} e^{i\frac{\omega}{v_p} x} \tag{209}$$

The dc conductance corresponds to the $\omega \to 0$ limit,

$$\sigma(x, 0) = \frac{e^2}{2\pi\hbar} K_p = \frac{e^2}{h} K_p \tag{210}$$

where $\hbar$ has been reinserted.

## 5.4. Boson Representation of the Hamiltonian

We have seen in the above that the bosonic excitations play an important role on the transport properties of the Tomonaga-Luttinger model. Correspondingly, there is one notable feature in the commutation relations for the density operators in Eqs. (182) and (183): the two parts of the Hamiltonian $H_0$ and $H_4/g_4$ show the same commutation relations. Therefore, one can introduce an effective Hamiltonian $\tilde{H}_0$ defined by

$$\tilde{H}_0 = \frac{2\pi v_f}{L} \cdot \sum_{r>0} \rho_1(-p)\rho_1(p) + \frac{2\pi v_f}{L} \sum_{r<0} \rho_2(-p)\rho_2(p) \tag{211}$$

which reproduces the commutation relation Eq. (182). Thus, the correlations can be calculated exactly by using $\tilde{H}_0$ as a replacement for $H_0$. The effective Hamiltonian is written in a bilinear form with the boson operators,

$$\tilde{H}_{13} = \tilde{H}_0 + H_4 + H_2$$

$$= \sum_{p>0} [C_p \quad C_p] v_p p \,(\cosh\varphi \, 1 + \sinh\varphi \tau^1) \begin{bmatrix} C_p \\ C^*_p \end{bmatrix}$$

$$= \sum_{p>0} v_p p (\gamma_p^\dagger \gamma_p + \gamma_{-p}^\dagger \gamma_{-p}) + \text{const} \tag{212}$$

In this section, we have discussed only the two-particle correlation functions. The equation of motion for the single-particle Green's function can also be written in a closed form [32, 33], and the precise calculations have been reported in Refs. [34, 35].

# REFERENCES

1. L. V. Keldysh. *Sov. Phys. JETP* 20, 1018 (1965).
2. J. Schwinger. *J. Math. Phys.* 2, 407 (1961).
3. L. P. Kadanoff and G. Baym. "Quantum Statistical Mechanics." Benjamin, New York. 1962.
4. C. Caroli, R. Combescot, P. Nozieres, and D. Saint-James. *J. Phys. C* 4, 916 (1971).
5. E. M. Lifshitz and L. P. Pitaevskii. "Physical Kinetics." Pergamon, Oxford, 1981.
6. G. D. Mahan: "Many-Particle Physics." Plenum Press. New York. 1990.
7. K. Cho, Z. Su, B. Hao, and L. Yu, *Physics Reports* 118, 1 (1985).
8. S. Datta. "Electronic Transport in Mesoscopic Systems." Cambridge University Press. Cambridge, 1997.
9. H. Haug and A. P. Jauho. "Quantum Kinetics in Transport and Optics of Semiconductors." Springer. Heidelberg. 1996.
10. P. Nozières. *J. Low Temp. Phys.* 17, 31 (1974).
11. K. Yamada, *Prog. Theor. Phys.* 53, 970 (1975); *Prog. Theor. Phys.* 54, 316 (1975).
12. A. C. Hewson. "The Kondo Problem to Heavy Fermions." Cambridge University Press, Cambridge, 1993.
13. A. L. Fetter and J. D. Walecka. "Quantum Theory of Many-Particle Systems." McGraw-Hill, New York, 1971.
14. S. Hershfield, J. H. Davies, and J. W. Wilkins, *Phys. Rev. B* 46 7046, (1992).
15. Y. Meir and N. S. Wingreen. *Phys. Rev. Lett.* 68, 2512 (1992).
16. R. Landauer. *Philos. Mag.* 21, 863 (1970).
17. M. Büttiker, Y. Imry, R. Landauer, and S. Pinhas, *Phys. Rev. B* 31, 6207 (1985).
18. D. S. Fisher and P. A. Lee, *Phys. Rev. B* 23, 6851 (1981); P. A. Lee and D. S. Fisher, *Phys. Rev. Lett.* 47, 882 (1981).
19. A. Oguri, *J. Phys. Soc. Jpn.* 70, 2666 (2001); *J. Phys. Soc. Jpn.* 72, 3301 (2003).
20. J. W. Negele and H. Orland. "Quantum Many-Particle Systems." Addison-Wesley, Redwood City, 1987.
21. A. Oguri, *Phys. Rev. B* 64, 153305 (2001).
22. G. M. Éliashberg. *Soviet Phys. JETP* 14, 886 (1962).
23. A. A. Abrikosov, L. P. Gor'kov, and I. Y. Dzyaloshinskii, "Quantum Field Theoretical Methods in Statistical Physics." Pergamon, London, 1965.
24. A. Oguri, *J. Phys. Soc. Jpn.* 66, 1427 (1997).
25. A. Oguri, *Phys. Rev. B* 56, 13422 (1997) [Errata: 58, 1690 (1998)].
26. A. Oguri. *Phys. Rev. B* 59, 12240 (1999).
27. A. Oguri. *Phys. Rev. B* 63 115305 (2001) [Errata: 63, 249901 (2001)].
28. J. S. Langer and V. Ambegaokar, *Phys. Rev.* 121, 1090 (1961).
29. J. R. Schrieffer. "Theory of Superconductivity." Benjamin Reading, Massachusetts, 1964.
30. A. Oguri and A. C. Hewson, *J. Phys. Soc. Jpn.* 74, 988 (2005).
31. J. Sólyom, *Adv. Phys.* 28, 201 (1979).
32. Y. Dzyaloshinskii and A. I. Larkin, *Soviet Phys. JETP* 38, 202 (1974).
33. E. U. Everts and H. Shultz, *Solid State Commun.* 15, 1413 (1974).
34. J. Voit, *Phys. Rev. B* 47, 6740 (1993).
35. V. Meden and K. Schönhammer, *Phys. Rev. B* 46, 15753 (1992).

# CHAPTER 8

# Computational Nanotechnology: Computational Design and Analysis of Nanosize Electronic Components and Circuits

## Jerry A. Darsey

*Department of Chemistry, University of Arkansas at Little Rock, Little Rock, Arkansas, USA*

## Dan A. Buzatu

*Division of Chemistry, National Center for Toxicological Research, Jefferson, Arkansas, USA*

## CONTENTS

# 1. INTRODUCTION

The fundamental building blocks of materials are undergoing a radical change. The dimension of these fundamental units is at the nanometer scale. The prefix "nano" represents a scale of one-billionth. That is, we are talking about devices that are at the billionth of a meter in dimension. The technology for such devices is in its infancy, yet an Internet search on the word nanotechnology has almost 4,430,000 hits. A similar search was performed in 1994, when we first got into this field, and the result was a little less than 1000 hits In other words, this field has exploded by more than three orders of magnitude in less than a decade.

The applications that are being explored for nanodimensional devices are too many to list. However, some of the more publicized applications are in diverse areas such as biomedicine, electronics, high-strength materials, photovoltaic materials, computing, transportation, artificial sensors, to name just a few. The National Science Foundation (NSF) has estimated that in the next 10 to 15 years, $1 trillion in business will be generated as a direct result of this new technology. The 2004 U.S. federal budget provided $847 million for nanotechnology research through the National Nanotechnology Initiative (NNI). Agencies represented in NNI include NSF, Department of Defense (DOD), Department of Energy (DOE), NASA, National Institutes of Health (NIH), and the Department of Homeland Security.

This chapter will concentrate on the application of computational nanotechnology to nanoscale electronics. Similar to the field of nanotechnology itself, the area of nanoscale electronics has also exploded.

Any review about developments in computational nanotechnology cannot be complete unless it addresses the parallel development of computational methods and computing platforms. The complexity of computational methods demonstrated each decade in the scientific literature, starting with the 1960s to present day, was directly related to the computational hardware available at that time. In the 1960s, computers were very large, relied on inefficient input–output devices, and, memory was measured in kilobytes (K). The most sophisticated computational calculations performed at that time were empirical and based on the laws of classical mechanics. The 1970s saw a movement toward quantum mechanics through semiempirical techniques and early forms of *ab initio* calculations, as computers decreased in physical size and gained CRT (cathode ray tube) terminals and modest increases in memory, which could be measured in a few megabytes (MB). By the 1980s, both semiempirical and *ab initio* calculations gained sophistication and accuracy through the development of advanced semiempirical methods such as the Austin method 1 (AM1) and the implementation of Hartree-Fock methods. This occurred because of increased access to high speed multiuser supercomputers such as the Cray systems and the development and introduction of RISC (reduced instruction set computer) chip–based workstations, such as the IBM RS6000 series and the Sun Microsystems Workstations in the late 1980s to early 1990s. The mid to late 1990s witnessed the coming of age and evolution of personal computer (PC)-based computational power, which resulted in inexpensive, large memory arrays, fast processor PC workstations, dual processor PC workstations, multiple node PC cluster computers, and large PC-based parallel array computers. Memory sizes and processor speeds increased significantly throughout the 1990s and has continued to increase throughout the first few years of the new millennium. There was also a significant improvement in accuracy and molecular size handling capability of the *ab initio* methods. This included the implementation of density functional theory (DFT) methods and compound models such as the Gaussian-1 and Gaussian-2 theories and the complete basis set (CBS) methods, which use combined results from lower computational cost techniques to achieve high accuracy results. The PC computational power, however, is reaching a limit under the current 32 bit chipset technology and

is currently moving toward 64 bit PC technology. This should once again provide large gains in computational power.

This chapter will start with a review of computational methods across the past four decades, the computing platforms on which they were implemented, and follow with the development of molecular electronics.

## 2. EMPIRICAL METHODS

Empirical methods are perhaps the first methods developed by chemists to study the conformation characteristics of molecules [1, 2]. These methods are based on a technique called molecular mechanics [3]. This technique is a mathematical formalism that attempts to reproduce molecular geometries, energies, and other features by adjusting bond lengths, bond angles, and torsion angles to values that are dependent on the hybridization of an atom and its bonding scheme [4]. Rather than using quantum physics, the method relies on the laws of classical Newtonian physics and some experimentally derived parameters to calculate geometry as a function of steric energy. The form of the equation, referred to as a force field calculation, is

$$E_{pot} = \sum E_{bnd} + \sum E_{ang} + \sum E_{tor} + \sum E_{oop} + \sum E_{nb} + \sum E_{el} \qquad (1)$$

$E_{pot}$ is the total steric energy made up of several different terms. These terms are $E_{bnd}$, the energy resulting from deforming a bond length from its natural value, which is calculated using Hooke's equation for the deformation of a spring. $E_{ang}$ is the energy resulting from bending a bond angle from its natural value and is also calculated from Hooke's law. $E_{tor}$ is the energy that results from deforming the torsion or dihedral angle. $E_{oop}$ is the out-of-plane bending component of the steric energy. $E_{nb}$ is the energy arising from nonbonded interactions. $E_{el}$ is the energy arising from Coulombic forces. When the terms shown in the general form of the force field are expanded, the equation takes the form

$$E_{pot} = \sum 1/2 K_b (b - b_0)^2 + \sum 1/2 K_\theta (\theta - \theta_0)^2 + \sum 1/2 K_\phi (1 + \cos N\phi)^2$$
$$+ \sum 1/2 K_\chi (\chi - \chi_0)^2 + \sum [(B/r)^{12} - (A/r)^6] + \sum (qq/r) \qquad (2)$$

The manner in which these terms are used to build a model is referred to as the functional form of the force field. The force constants

$$k_b, k_\theta, k_\phi, k_\chi \qquad (3)$$

and equilibrium values

$$b_0, \theta_0, k_\phi, k_\chi \qquad (4)$$

are atomic parameters that are experimentally derived, mostly from, NMR spectroscopy, IR spectroscopy, microwave spectroscopy, X-ray diffraction, or Raman spectroscopy. The energy of the atoms in a molecule is calculated and minimized using a variety of directional derivative or gradient techniques.

In contrast to *ab initio* methods, molecular mechanics is used to compute molecular properties that do not depend on electronic effects. These include geometry, rotational barriers, vibrational spectra, heats of formation, and the relative stability of conformers [5, 6]. Because the calculations are fast and efficient, these empirical techniques can be used to examine systems containing thousands of atoms. Unlike *ab initio* methods, molecular mechanics relies on experimentally derived parameters, so that calculations on new molecular structures may be misleading.

The method described above is used to calculate the energy of a compound in a specific 3D orientation and to optimize the geometry as a function of energy. This is usually accomplished by adjusting the coordinates of each of the atoms and recomputing the energy of the molecule until a minimum energy is obtained. This technique is also used to simulate the time-dependent behavior of molecules and is generally referred to as molecular dynamics. The first generation of computational empirical method algorithms were developed on large card-reading machines such as the IBM 360.

Figure 1. IBM 360 console, printer, and tape drives.

## 2.1. IBM 360

Figure 1 is an illustration of the large, room-sized IBM 360 that appeared toward the end of the 1960s. There were many versions of the IBM 360, from the IBM 360/20, /30, and /40, to the IBM 360/67 and /75, all the way up to the IBM 360/91, /95, and /195, which eventually became the IBM 370 technology. Core memory sizes for these machines ranged from a few K for the early version numbers, to several hundred KB for the midrange versions, to 1 to 4 MB for the later models (/91 and /95). Processor clock cycle times ranged from 3.6 $\mu$s per byte for the 360/20 all the way up to 0.125 $\mu$s per 8 bytes for the 360/195. These machines primarily used punch cards for input/output operations, which eventually were replaced by paper terminals.

## 3. SEMIEMPIRICAL APPROXIMATIONS

Semiempirical methods are self-consistent field (SCF) methods that approximate molecular orbital integrals using parameters derived from experimental data. With the exception of methods later developed that include $d$ orbitals, they all use a restricted basis set of one $s$ orbital and three $p$ orbitals ($p_x$, $p_y$, and $p_z$) per atom. All orbital overlap integrals are ignored in the Roothan-Hall secular equation, which has the form:

$$c_i|F - \varepsilon_i S|c_i = 0 \tag{5}$$

Pople and Segal introduced the first semiempirical computational method, complete neglect of differential overlap (CNDO), in the mid 1960s [7]. This early method did not include any contribution to the electronic wavefunction from molecular atomic overlap and used only spherical symmetrical orbitals. Other semiempirical algorithms that appeared in the late 1960s were based on the neglect of diatomic differential orbital overlap (NDDO) method, which includes the directionality of orbitals on the same atom for repulsion integrals. These included the MINDO and MINDO/2 (modified intermediate neglect of differential overlap) methods. MINDO and MINDO/2 were developed by Dewar and others in 1969 and 1970 [8, 9]. These techniques use empirical data to parameterize the repulsion integrals rather than analytical solutions.

In response to the quantum mechanical issues that arose as semiempirical calculations grew in size and complexity, principally that approximations neglected electron repulsion integrals involving a one-center overlap. Dewar and Thiel developed MNDO (modified neglect of overlap integrals) in 1977 [10]. MNDO theory achieved better determinations of

multicenter repulsion integrals by treating the total energy of the molecule, $E_{tot}^{mol}$, as the sum of the electronic energy, $E_{el}$, and the repulsion, $E_{AB}^{core}$, between the atoms A and B.

The mathematics involved in calculating the $E_{AB}^{core}$ value are estimated from experimental results. This fact allows researchers the freedom to manipulate the parameter set to fit their experimental designs. If a trend is suspected in a family of molecules, experimental data can be input for the calculations and the math, being based upon experimental results, will show conclusive results. Those molecules that fit the proposed pattern give reasonable results, while those that might appear in the same category, but are actually in a different class of molecules, show dramatically different results. Certain small deviations in bonding energies, bond lengths, or heats of formation occur because of the experimental results fusing with theoretical expectations.

The elements H, C, N, and O were parameterized by 1985. The AM1 and PM3 methods were developed by Dewar, Zoebisch, Healy, and Stewart in 1985–1986 to include these parameterizations and further improvements of the MNDO technique [11, 12].

## 4. HARTREE-FOCK CALCULATIONS

The Schrödinger equation introduced the idea of describing an electron using a wavefunction, treating the molecule as a collection of particles represented by another wavefunction, which is a function of time and the coordinates of the particles. Because the time-dependent Schrödinger equation is too complex to solve for a molecule, the wavefunction is separated into a time function and a spatial function using separation of variables. Ignoring the time function, the time-independent Schrödinger equation is formed.

$$H\Psi(r) = E\Psi(r) \tag{6}$$

This is a function of the positions of the nuclei and electrons of the molecule. The Born-Oppenheimer approximation allows further simplification by treating the nuclei as fixed entities that are orbited by the electrons. This approximation is made based on the fact that a nucleus is thousands of times more massive than the electron [13]. Using this approximation, the nuclear kinetic energy term can be removed, and the molecular electronic Hamiltonian is obtained:

$$H^{elec} = -\frac{1}{2} \sum_i^{electrons} \nabla^2 - \sum_i^{electrons} \sum_I^{nuclei} \left( \frac{Z_i}{|\vec{R}_J - \vec{r}_i|} \right)$$

$$+ \sum_i^{electrons} \sum_{j<i} \left( \frac{1}{\vec{r}_i} - \vec{r}_j \right) + \sum_I^{nuclei} \sum_{J<I} \left( \frac{Z_I Z_J \vec{R}_I}{\vec{R}_J} \right) \tag{7}$$

The preceding equation and the equations that are to follow are all in atomic units in order to keep from rewriting the fundamental constants (Planck's constant, mass of electron, etc.). There are restrictions that apply to the wavefunction for a particle. A major restriction is that if the wavefunction is integrated over all space, the probability that a particle is somewhere in that space is 1. This requires a normalization constant

$$\int_{-\infty}^{\infty} |c\Psi|^2 dv = 1 \tag{8}$$

According to molecular orbital theory, the approximation is made that the wavefunction is a combination of a normalized, orthogonal set of molecular orbitals $\psi_1, \psi_2, \psi_3$. The Hartree product is a simple way of combining these orbitals to form a wavefunction:

$$\psi(\vec{r}) = \psi_1(\vec{r}_1)\psi_2(\vec{r}_2)\ldots \tag{9}$$

However, this wavefunction does not take into consideration an electron's spin, which can be $+1/2$ or $-1/2$. If two of the orbitals are swapped, there is no sign change in the function; that

is, the product is not antisymmetric. In order to make the product antisymmetric. the orbitals are multiplied by one of two spin functions:

$$\alpha\left(+\frac{1}{2}\right) = 1 \quad \alpha\left(-\frac{1}{2}\right) = 0$$

$$\beta\left(+\frac{1}{2}\right) = 0 \quad \beta\left(-\frac{1}{2}\right) = 1 \tag{10}$$

When the orbital is multiplied by $\alpha$, the electron is spin up and when multiplied by $\beta$, the electron is spin down. The molecular orbitals are now called spin orbitals and now represent position and spin. Using spin orbitals, a determinant is built. Electrons are assigned to this determinant, in orbitals, in pairs of opposite spin. The rows of the determinant represent all the assignments of each electron to all possible spin orbital combinations [14]. There are two electrons assigned to each molecular orbital (MO), and the whole determinant is made up of $n = 2$ MOs. The whole determinant is multiplied by $1 = \sqrt{n!}$, which is the normalization constant, and is known as the Slater determinant. In order to best describe molecular orbitals mathematically, linear combinations of one electron function. called basis functions, are used. A type of function frequently used to form basis functions in most modern quantum chemistry computational packages such as Gaussian are Gaussian-type atomic functions. These have the form:

$$g(\alpha, \vec{r}) = c x^n y^m z^l \exp -(\alpha r^2) \tag{11}$$

where $\alpha =$ constant, which determines the radial size of the function; $n, m, l$ determine the type of orbital; and $c =$ constant for normalization. Linear combinations of these functions (which are centered on the same atomic nucleus) are formed, and these functions are called contracted Gaussian-type functions (CTGF) [15]:

$$\chi_\mu = \sum_p d_{\mu p} g_p \tag{12}$$

where $d_{\mu p}s$ are fixed constants in the basis set. The contracted Gaussian function represents an atomic orbital and is an approximation. The accuracy of the approximation improves with the number of primitive Gaussian functions used in the linear combination. The more primitive Gaussian functions used, the better the description of the electron density. The molecular orbital is described by the following equation

$$\phi_i = \sum_\mu c_{\mu i} \chi_\mu = \sum_\mu c_{\mu i} (d_{\mu p} g_p) \tag{13}$$

where $c_{\mu i}$ is molecular orbital expansion coefficients.

In order to solve for the expansion coefficients, the variational principle is used, which says that the energy of the exact wavefunction will always be lower than that of the approximate wavefunction. This turns into a minimization problem. where the Hartree-Fock variational energy is minimized. Its form is

$$E_{HF} = \langle \phi \hat{H}_d + V_{NN} \phi \rangle \tag{14}$$

where $V_{NN}$ is potential energy of nuclear–nuclear repulsion term. For a closed shell system with $1/2n$ doubly occupied orbitals, the energy is given by

$$E = \sum_{i=1}^{1/2n} 2h_{ii} + \sum_{i=1}^{1/2n} \sum_{j=1}^{1/2n} (2J_{ii} - K_{ii}) \tag{15}$$

where $h_{ii}$ is contribution of one electron term to the energy and is given by

$$h_{ii} = \int \phi_i(1)\left(-\frac{1}{2}\nabla_i^2 - \sum_{n=1}^{N} \frac{Z_n}{r_{1n}}\right)\phi_i(1)d\tau_1 \tag{16}$$

where $J_{ij}$ is Coulomb integrals, which are the electrostatic interaction between double occupied orbitals $\phi_i$ and $\phi_j$.

$$J_{ij} = \int \phi_i(1)\phi_j(2)\frac{1}{r_{12}}\phi_i(1)\phi_j(2)d\tau_{12} \tag{17}$$

The integration is over the coordinates of electron 1 and 2. $K_{ij}$ is the exchange integral and is given by

$$K_{ij} = \int \phi_i(1)\phi_j(2)\frac{1}{r_{12}}\phi_i(2)\phi_j(1)d\tau_{12} \tag{18}$$

The Hartree-Fock equations allow minimization of the energy while keeping the orbitals mutually orthogonal.

$$\hat{F}_i(1)\phi_i(1) = \varepsilon_i\phi_i(1) \tag{19}$$

where $\varepsilon_i$ is orbital energy for orbital $i$, and $\hat{F}_i$ is the Fock operator and is given by

$$\hat{F}_i(1) = \hat{h}(1) + \sum_j [2\hat{J}_j(1) - \hat{K}_j(1)] \tag{20}$$

where $\hat{h}(1)$ is one electron operator and is given by

$$\hat{h}(1) = -\frac{1}{2}\nabla_1^2 - \sum_{\alpha=1}^{N} \frac{Z_\alpha}{r_{1\alpha}} \tag{21}$$

and $\hat{J}_j$ is Coulomb operator and $\hat{K}_j$ is exchange operator, and their equations are

$$\hat{J}_j\phi_i = \phi_i \int |\phi_j(2)|^2 \frac{1}{r_{12}} dv_2 \tag{22}$$

$$\hat{K}_j\phi_i(1) = \phi_j \int \phi_j^*(2) \tag{23}$$

The SCF [16] procedure starts with a guess set of orbitals (wavefunction) that construct the Fock operator. The Fock operator yields a new set of orbitals after the solution of the Hartree-Fock equations. The resulting orbitals are used to build the Fock operator again, and another solution to the equations is obtained, once again leading to a new set of orbitals. The second term in the Fock operator equation [i.e., $\sum(2\hat{J}_j - K_j)$] represents the effect of the field of all other electrons in the molecule on one electron, and the iterative procedure is finished when the field stays unchanged (beyond a given criterion) [17]. Roothaan and Hall developed a system for the Hartree-Fock equations based on matrices [18]. Their equations are based on the orbital coefficients previously described and are:

$$FC = SC\epsilon \tag{24}$$

where $F$ is Fock matrix, $S$ is overlap matrix (overlap between orbitals), and $\epsilon$ is diagonal matrix of the orbital energies (where the terms $\epsilon_i$ = the energy of orbital $\chi_i$). The Fock matrix represents the effects of all the electrons on each orbital and is given by

$$F_{\mu\nu} = H_{\mu\nu}^{core} + \sum_{\alpha=1}^{N}\sum_{\sigma=1}^{N} P_{\lambda\sigma}\left[(\mu\nu|\lambda\sigma) - \frac{1}{2}(\mu\nu|\nu\sigma)\right] \tag{25}$$

where $H_{\mu\nu}^{core}$ is matrix of energy of $1e^-$ in a field of bare nuclei; and $(\mu\nu|\lambda\sigma)$ is $2e^-$ repulsion integral (each electron sees the rest of the electrons as an average field). There is no direct instantaneous electron–electron correlation. $P$ is the density matrix, which is given by

$$P_{\alpha\sigma} = 2 \sum_{i=1}^{occupied} c_{\lambda i}^* c_{\sigma i} \tag{26}$$

This is a closed shell method known as the restricted Hartree-Fock method, where the coefficients are summed over only the occupied MOs, each which holds two electrons (Pauli exclusion principle). Each of the terms in Eq. (24) is based on the MO expansion coefficients. The minimization of the equation yields the approximate solution to the exact wavefunction. Because the equation is not linear, the process is carried out using the iterative scheme that was previously described. When the scheme converges, the solution produces a complete set of occupied MOs ($\phi_{i,j}; \ldots$) and unoccupied MOs ($\phi_{a,b} \ldots$). The density matrix $P$ represents the probability distribution of the electrons in the molecule. The molecule's electrostatic potential can be calculated from the Hartree-Fock method using

$$V_i = \sum_A \frac{Z_A}{|r_i - R_A|} - \sum_{\mu, \nu} P_{\mu, \nu} \int \frac{\phi_\mu \phi_\nu}{|r_i - r'|} dr \tag{27}$$

where $P_{\mu, \nu}$ is an element of the density matrix, and $\phi_\mu$ and $\phi_\nu$ are basis functions. The electrostatic potential $E_i$, obtained from the atomic valence population $q_i$, is given by

$$E_i = \sum_{i=1}^{n} \frac{q_i}{r_{ij}} \tag{28}$$

The electrostatic potential is an exact one electron property that can be calculated at any point in space based on a molecular wavefunction.

## 4.1. Supercomputers and High-Speed Workstations

High-speed multiuser parallel processing supercomputers helped to advance the accuracy of *ab initio* and semiempirical calculations and allow them to model larger and more complex systems. In fact, supercomputers and high-speed workstations were the tools that allowed all of computational modeling and simulation science to advance as a whole. The first Cray 1 supercomputer was installed at Los Alamos National Laboratory in 1976. It performed 160 million floating-point operations per second (160 megaflops) and had an 8 megabyte main memory (Fig. 2).

In 1988, Cray Research introduced the Cray Y-MP. It was the world's first supercomputer to sustain over 1 gigaflop on many applications. It used multiple 333 megaflops processors to power the system to a record sustained speed of 2.3 gigaflops. The formation of Cray supercomputer centers in the 1980s and early 1990s at academic and government scientific



Figure 2. Cray 1 supercomputer. Image courtesy of Cray, Inc.

institutions provided researchers access to high-speed, powerful computing resources, which in turn advanced computational modeling in all fields.

RISC processor–based workstations also became tools of the trade for researchers, which led to many advances in computational modeling science. IBM developed a 32 bit RISC chip in the 1970s after engineers realized that complex instruction set chips (CISC) were not necessary, as 80% of the CPU's time was spent on basic instructions. This later developed into the IBM RS6000 series of computers, which appeared in the late 1980s to early 1990s, in a line of workstation computers using a 64 bit RISC architecture. The 64 bit chip could perform double precision floating-point operations in half the time a 32 bit chip could perform them. Digital Equipment Corporation (DEC) also introduced a series of high-speed 64 bit RISC processor-based workstations in 1992 called Alpha, as did Sun Microsystems in 1995, with the introduction of its 64 bit Ultra-SPARC (scalable performance Architecture)-based workstations. These machines became the benchtop standard for scientific computational modeling until the PC revolution of the late 1990s.

# 5. DENSITY FUNCTIONAL THEORY

## 5.1. Introduction

Density functional theory has proved a very successful approach for the description of ground state properties of metals, semiconductors, and insulators. The success of density functional theory (DFT) not only encompasses standard bulk materials but also complex materials such as proteins and carbon nanotubes. The whole theory is based on functionals of the electron density, which therefore plays the central role. However, the key functional, which describes the total energy of the electrons as a functional of their density, is not known exactly: the part of it that describes electronic exchange and correlation has to be approximated in practical calculations.

The main idea of DFT is to describe an interacting system of charged particles (such as electrons) through the system's electronic density and not its many-body wavefunction. For any system containing $N$ electrons, the basic calculations involve the spatial coordinates $x$, $y$, and $z$—rather than the $3 \times N$ degrees of freedom. Although DFT in principle gives a good description of ground state properties, practical applications of DFT are based on approximations for the so-called exchange-correlation potential. The exchange-correlation potential describes the effects of the Pauli principle and the Coulomb potential beyond a pure electrostatic interaction of the electrons. Possessing the exact exchange-correlation potential means that the many-body problem is solved exactly, and this is clearly not feasible in large molecular or solid-state systems.

A common approximation often made is the local density approximation (LDA), which locally substitutes the exchange-correlation energy density of an inhomogeneous system by that of an electron gas evaluated at the local density. Although many ground state properties are well described in the LDA, the dielectric constant may be overestimated by as much as 10% to 40% compared to experiment. This overestimation stems in part from the neglect of a polarization-dependent exchange correlation field in LDA compared to DFT.

In 1986, Schlüter et al. [19] calculated an accurate exchange-correlation potential for silicon using many-body perturbation theory. They showed that the "band-gap problem" (the observation that the electronic band gap of semiconductors in density functional theory calculations is only about 50% of the experimental band gap) was present even with a more accurate exchange-correlation potential and therefore corresponded to a non-analyticity in the functional, rather than an inadequacy of the local-density approximation normally used in density functional theory.

Related work includes an investigation of the DFT band-gap problem as a semiconductor is compressed to a metallic state [20], a study of the behavior of the exact exchange-correlation potential at semiconductor interfaces [21, 22], an investigation of exact DFT for a model semiconducting wire using Monte Carlo methods [23], a study of exact DFT in the presence of a macroscopic electric field [24] (which for an infinite solid requires DFT to be augmented to become a density-polarization functional theory), and an investigation of

DFT for a Hubbard model [25]. A similar technique has also been developed in which a new type of density functional theory [26], based on a simplified self-energy approach, has been demonstrated to outperform conventional Kohn-Sham DFT.

## 5.2. Background

The following derivation and discussion is based on that given by Ohno et al [27] and is by no means a comprehensive treatment. For more details of the derivation and a deeper discussion of its significance, please consult that text. Density functional theory is one of many *ab initio* techniques that attempt to solve the many-body Schrödinger equation:

$$H\Psi_i(1, 2, \ldots N) = E_i\Psi_i(1, 2, \ldots N)$$

where $H$ is the Hamiltonian of a quantum mechanical system composed of $N$ particles, $\Psi$ is its $i$th wavefunction, and $E_i$ is the energy eigenvalue of the $i$th state. The particle coordinates are usually associated with a spin and a position coordinate. For electronic systems with nonrelativistic velocities, the Hamiltonian for an $N$-electron system is:

$$H = -\frac{1}{2}\sum_{i=1}^{N}\nabla_i^2 + \sum_{i<j}\frac{1}{|r_i - r_j|} + \sum_{i=1}^{N}v(r_i)$$

where the first term of this equation represents the electron kinetic energy, the second term the electron–electron Coulombic interactions, and the third term is the Coulombic potential generated by the nuclei. This equation also assumes that the nuclei are effectively stationary with respect to electron motion (Born-Oppenheimer approximation).

Initial approaches to this problem attempted to transform the full $N$-body equation into $N$-single-body equations by using the Hartree-Fock (HF) approximation. This approximation basically affects the accuracy with which the "exchange-correlation" contribution to the total energy is calculated. The exchange contribution is a direct consequence of Pauli's exclusion principle, which prohibits two fermions from occupying the same quantum state. This reduces the probability of one electron being near another electron of the same spin. The correlation contribution is due to the reduction in probability of an electron being near another electron due to strong electron–electron Coulomb repulsion. HF only includes the correlation contribution for similar spin electrons, but neglects entirely the contribution for opposite spin electrons. Traditionally (and confusingly), the part of exchange-correlation included in HF is known as "exchange" and that neglected is known as "correlation."

In contrast to these methods, which try to determine approximations of the electron density or many-electron wavefunction, DFT can "exactly" calculate any ground-state property from the electron density. In 1964, Hohenberg and Kohn [28] considered the ground state of the electron–gas system in an external potential $v(r_i)$, and proved that the following density functional theorem holds exactly: There is a universal functional $F[\rho(r)]$ of the electronic charge density $\rho(r)$ that defines the total energy of the electronic system as:

$$E = \int v(r)\rho(r)dr + F[\rho(r)]$$

The energy of the system can be minimized to find the "true" electronic charge density in the external potential. However, this procedure is exact only for a nondegenerate ground state. Unfortunately, as yet an exact general form of the functional has not been found, so approximations must continue to be used.

## 5.3. The PC Revolution

Until about the mid to late 1990s, serious scientific computational work was performed on 64 bit RISC processor-based desktop workstations and multiuser supercomputers. In the spring of 1993, Intel introduced the first generation of Pentium processors: 60 and 66 MHz 32 bit chips, composed of 3.1 million transistors, capable of 100 million instructions per second (MIPS). Pipelining and larger L1 caches were used to significantly increase performance.

However, even though these processors represented a significant technological advancement for the home PC market, they still could not compete with the 64 bit RISC-based machines of the time. At that time, 64 bit RISC processors ranged from 66 to 100 MHz and could perform 2 floating-point operations in the same amount of time a 32 bit processor could perform one.

In 1995, Intel introduced the Pentium Pro, 32 bit 200 MHz processor composed of 5.5 million transistors and capable of 440 MIPS. This innovation, coupled with a larger L1 (onboard) cache, and the development of advanced motherboard technologies including the introduction of an L2 cache, larger memory sizes, fast serial buses, and dual processor platforms, allowed PCs to enter the arena of scientific and engineering computation. However, most computational researchers at the time were still using 64 bit RISC-based systems.

The true computational PC revolution occurred during the next 5 years; the direct result of a number of technological advancements based on the 32 bit microchip architecture. In 1997, Intel released the second generation of Pentium processors, the Pentium II. These came in faster 233, 266, and 300 MHz versions and included a much larger on-board cache. In the spring of 1998, Intel released the 333 MHz Pentium II processor, which used a new 0.25 micron manufacturing process to run faster and generate less heat than before. Intel's competitor Advanced Micro Devices (AMD) introduced the AMD K6-III 400 MHz version in 1999, which contained approximately 23 million transistors and was based on a 100 MHz super socket 7 motherboard. Both Intel and AMD released 1 GHz processors in 2000 containing significantly larger onboard caches and wider and faster serial buses. In 2003, an Intel Pentium 4 3.40 GHz processor became available, developed with 0.09–0.13 micron printed circuit technology, using hyperthreading technology and an 800 MHz system bus. Toward the end of 2003, both AMD and Intel introduced 64 bit processors, the AMD Athlon 64 and the Intel Itanium. With these innovations and more to come, PCs have truly become state of the art computational science platforms.

# 6. MOLECULAR ELECTRONICS

## 6.1. Introduction

The field of molecular electronics is fast becoming one of the fastest growing areas of nanotechnogy, and nanotechnology is developing into the fastest growing area of science. It essentially uses individual molecules to perform functions in electronic circuitry. There was a time when this was considered impossible. Currently, in most solid-state circuits, this is performed primarily by semiconductor devices. Molecules are hundreds to thousands of times smaller than the smallest semiconductor device. Electronic devices constructed from molecules will necessarily be hundreds of times smaller than their semiconductor-based counterparts. The classic talk that defined and probably gave birth to nanotechnology was given by Richard Feynman on December 29, 1959, at the annual meeting of the American Physical Society at the California Institute of Technology (Caltech). It was first published in the February 1960 issue of Caltech's *Engineering and Science*, which owns the copyright. It has been made available on the Web at Zyvex Corporation [29]. For an account of the talk and how people reacted to it, see Chapter 4 of *Nano!* by Ed Regis (Little/Brown, 1995).

The field of molecular electronics probably was given its first practical molecule by Mark A. Ratner and Arieh Aviram at Northwestern University, when they published their historical paper "Molecular Rectifiers," which appeared in the journal *Chemical Physics Letters* in November 1974 [30]. Since then, the research in this area has increased exponentially. But it was only recently, in 1997, that two separate groups demonstrated practical molecular rectification. One group was led by Metzger [31] at the University of Alabama and the other led by Reed [32] at Yale University.

Individual molecules can easily be made exactly the same by the billions and trillions. The dramatic reduction in size, and the sheer enormity of numbers in manufacture, are the principle benefits offered by the field of molecular electronics.

To illustrate the order of magnitude we are dealing with, let's look at how the ability to store information has changed. Memory is the primary means for storing information.

For example, a color photo requires about $10^5$ bytes. A typical modern personal computer has between $10^8$ to $10^9$ bytes. The human brain is estimated to have about $10^{13}$ bytes. To store all of the information in the Library of Congress, we would need an estimated $10^{20}$ bytes. However, if we are able to take advantage of the size of molecules, we would have available a potential storage capacity approaching $6.022 \times 10^{23}$ bytes of information. This number is ofcourse what chemists call a "mole." These numbers are cited in James M. Tour's book, *Molecular Electronics* [33]. As Professor Tour points out, the big question is "how to access such a vast memory in any usable timeframe" [33].

There are many areas that are considered part of the molecular electronics domain. As one of many possible examples that could be cited, a relatively recent conference, 'Towards Molecular Electronics," held on 23–28 June 2003 in Srem," Poland, can illustrate the breath of topics included in molecular electronics [34]. This conference was part of the Central European Conference on Advanced Materials and Nanotechnology. Topics included were molecular conducting and superconducting materials, molecular magnets; molecular opto-electronic materials; molecular display materials and devices; fullerenes and nanotubes; supramolecular systems; self-assembled and Langmuir-Blodgett (LB) films; nano-techniques and molecular manipulations; molecular electronic devices; molecular machines; molecular sensors; proton-conducting systems; molecular electronics and living organisms; and biocomputing and molecular neural networks. This conference was actually the fourth in a series of international meetings planned as an interdisciplinary forum for discussion on all aspects concerning molecular electronics, specifically materials (molecular conductors, molecular magnets, electrochromic organic materials), synthesis, structure and properties, applications of molecular materials in electronics and optoelectronics, supramolecular systems, and LB films and self-assembled monomolecular layers. To demonstrate one aspect of this enormous potential, let's look closer at LB films. A Langmuir-Blodgett film is a set of monolayers, or layers of organic material one molecule thick, deposited on a solid substrate. An LB film can consist of a single layer or many, up to a depth of several visible-light wavelengths. The term Langmuir-Blodgett comes from the names of a research scientist and his assistant, Irving Langmuir and Katherine Blodgett, who discovered unique properties of thin films in the early 1900s. Langmuir's original work involved the transfer of monolayers from liquid to solid substrates. Several years later, Blodgett expanded on Langmuir's research to include the deposition of multilayer films on solid substrates.

By transferring monolayers of organic material from a liquid to a solid substrate, the structure of a film can be controlled at the molecular level. Such films exhibit various electrochemical and photochemical properties [35]. This has led some researchers to pursue LB films as a possible structure for integrated circuits (ICs). Ultimately, it might be possible to construct an LB-film memory chip in which each data bit is represented by a single molecule. Complex switching networks might be fabricated onto multilayer LB-film chips [35].

As an additional illustration on how rapidly this field of molecular electronics is growing, we will summarize recent developments in just the past couple of years, some of which were written by Stu Borman [36] in the December 16, 2002, issue of *C&E News*. For example, the field of molecular electronics made some very significant developments in several key areas. One was the fabrication of nanoscale wires. These wires contained segments of varying chemical or dopant composition. The wires were the first to contain more than one junction within an individual nanowire or nanotube. The work was done by groups at Harvard University [37], Lund University, Sweden [38], and U.C. Berkeley [39, 40].

The patterned self-assembly of integrated semiconductor devices such as light-emitting diodes (LEDs) on surfaces was demonstrated by George M. Whitesides and coworkers at Harvard [41]. Patterned solder-coated areas placed on the substrates doubled as self-assembly receptors and as electrical connections.

Another recent development deals with a technique for transforming coated films into polymer-covered liquid-crystal layers, which could lead to thinner, cheaper, and more flexible liquid-crystal displays. These particular liquid-crystal displays were developed by Dirk J. Broer and coworkers [42, 43].

Phaedon Avouris and his group at IBM and colleagues found that carbon nanotube-based transistors can outperform silicon transistors—suggesting that it might be feasible to replace

silicon with carbon nanotubes in electronic devices when the size of silicon-based circuits can no longer be reduced [44]. By experimenting with different device structures, the IBM researchers were able to achieve the highest transconductance (the measure of the current carrying capability) of any carbon nanotube transistor to date. High transconductance means transistors can run faster, leading to more powerful integrated circuits. The researchers also discovered that carbon nanotube transistors produced more than twice the transconductance per unit width of top-performing silicon transistor prototypes. "Proving that carbon nanotubes outperform silicon transistors opens the door for more research related to the commercial viability of nanotubes," said Phaedon Avouris, manager of nanoscale science, IBM Research. "Carbon nanotubes are already the top candidate to replace silicon when current chip features just can't be made any smaller, a physical barrier expected to occur in about 10 to 15 years" [44, 45].

Paul L. McEuen and Daniel C. Ralph of Cornell and coworkers and (independently) a group led by Hongkun Park of Harvard created transistors in which a single molecule of a transition-metal organic complex bridges a nanometer-scale gap between electrodes [46]. In addition, a team led by Benjamin Mattes of Santa Fe Science & Technology showed that dopable polymers can be electrochemically cycled for up to 1 million cycles in ionic liquids without failure—suggesting that ionic liquids could be useful for fabricating and operating polymer electrochemical devices [47].

Groups led by Chang-Beom Eom of the University of Wisconsin, Madison, and Xiaoxing Xi at Pennsylvania State University independently developed oriented thin films of magnesium diboride that are potentially useful for making superconducting devices [48]. These would require less cooling than current niobium-based superconductor circuits.

A team led by Yet-Ming Chiang of MIT found that doping lithium iron phosphate with metal ions boosts its electrical conductivity by an astonishing eight orders of magnitude. The title of the article describing this nanomaterial is "Electronically Conductive Phospho-Olivines as Lithium Storage Electrodes." Their key limitation has been extremely low electronic conductivity, until now believed to be intrinsic to this family of compounds. It has been shown that controlled cation nonstoichiometry combined with solid-solution doping by metals supervalent to $Li^+$ increases the electronic conductivity of $LiFePO_4$ by a factor of $\sim 10^8$ [49]. This material is considered a potentially inexpensive electrode material for high-power-density lithium batteries [49].

A simple procedure for converting an insulating calcium–aluminum oxide to a transparent electrical conductor by heating the material and then exposing it to UV light was developed by Hideo Hosono and coworkers [50, 51] of Tokyo Institute of Technology. Materials that are good electrical conductors are not in general optically transparent, yet a combination of high conductivity and transparency is desirable for many. However, transparent conductors have been found to be very useful in optoelectronic devices.

A group led by Carlo D. Montemagno of UCLA devised a switch based on mutant $F_1$-ATP synthase that can turn a biomolecular nanomotor off and on [52]. The biophysical and biochemical properties of motor proteins have been well studied, but these motors also show promise as mechanical components in hybrid nano-engineered systems. The cytoplasmic $F_1$ fragment of the adenosine triphosphate synthase ($F_1$-ATPase) has the ability to function as an ATP-fueled rotary motor. It also has been integrated into self-assembled nanomechanical systems as a mechanical actuator. The results of this molecular system demonstrate the ability to engineer chemical regulation into a biomolecular motor and represent a crucial step toward controlling integrated nanomechanical devices at the single-molecule level.

Harry L. Anderson of Oxford University, Franco Cacialli of University College London, and coworkers prepared polyrotaxanes—polymer wires sheathed by cyclo-dextrin rings. They demonstrated that the coated wires act as semiconductors and used them to prepare blue and green LEDs [53]. Control of intermolecular interactions was found to be crucial to the exploitation of molecular semiconductors for both organic electronics and the viable manipulation and incorporation of single molecules into nano-engineered devices. Conjugated macromolecules, such as poly(para-phenylene), poly(4,4'-diphenylene vinylene) or polyfluorene, were investigated. The approach employed preserves the fundamental semi-conducting properties of the conjugated wires and is also found to be effective at both

increasing the photoluminescence efficiency and blue-shifting the emission of the conjugated cores, in the solid state, while still allowing charge-transport. These polymers were used to prepare single-layer light-emitting diodes with Ca and Al cathodes, and blue and green emission was observed. "The reduced tendency for polymer chains to aggregate allows solution-processing of individual polyrotaxane wires onto substrates, as revealed by scanning force microscopy" [53].

John T. Fourkas of Boston College and coworkers used multiphoton absorption to encode and read-out 3D data in molecular glasses and highly cross-linked polymers. "The materials are inexpensive, easy to process, and can store and read data robustly with an unamplified laser" [54]. Significant interest is expressed in three-dimensional laser-based optical data storage techniques, which can potentially provide efficient storage at densities significantly higher than those that are currently available from magnetic media. The development of inexpensive, efficient, and robust media has been a major obstacle in optical data storage. However, Fourkas's group has discovered a class of materials that become highly fluorescent on multiphoton absorption of pulses of 800-nm light from a Ti:sapphire oscillator. This makes them an excellent potential storage media. These materials are also inexpensive, have high optical quality, can be quickly processed, and can take a number of useful forms, including molecular glasses and highly cross-linked polymers. In addition, the potential for three-dimensional data storage at high densities makes these materials very valuable as storage media.

Another key molecular electronics advance in 2000 was the first demonstration of single-molecule electroluminescence by Robert M. Dickson and coworkers at the Georgia Institute of Technology [55]. Readily formed at nanoscale break junctions, arrays of individual spatially isolated, strongly electroluminescent $Ag_2-Ag_8$ nanoclusters perform complex logic operations within individual two-terminal nanoscale optoelectronics devices. Simultaneous electrical excitation of discrete room-temperature nanocluster energy levels directly yields AND, OR, NOT, XOR, and even full addition logic operations with either individual nanoclusters or nanocluster pairs as the active medium between only two electrodes. Imaged in parallel, noncontact electroluminescent readout obviates the need for electrically isolating individual features. This gated, pulsed, two-terminal device operation will likely drive future nano and molecular electronics advances without complicated nanofabrication [55, 56].

## 6.2. Conductive Polymers

Perhaps the first thought-of building circuits using molecules was realized when it was determined that polymers could be made conductive. In the early 1970s, a Japanese graduate student was trying to repeat the synthesis of polyacetylene, by linking together the molecules of ordinary acetylene gas used in welding. After the chemical reaction was completed, instead of the expected black powder, the student found a film coating on the inside of his glass reaction vessel. This material looked much like aluminum foil. He later realized that he had inadvertently added much more than the recommended amount of catalyst to cause the acetylene molecules to link together. Before this accident occurred, all carbon-based polymers were regarded as insulators. News about the foil-like film reached Alan MacDiarmid of the University of Pennsylvania. At that time, he was interested in nonmetallic electrical conductors. Because polyacetylene in its newly discovered form looked so much like a metal, MacDiarmid hypothesized that it could have the ability to conduct electricity like a metal. MacDiarmid invited the student's instructor to join his team in the United States, and this collaboration soon led to further findings. The University of Pennsylvania investigators confirmed that polyacetylene exhibited surprisingly high electrical conductivity.

The idea that plastic materials could conduct electricity was considered absurd. Indeed, plastics have extensively been used by the electronics industry because they did not exhibit this very property. They are used as inactive packaging and insulating material. This very narrow perspective is rapidly changing as a new class of polymers known as intrinsically conductive polymers or electroactive polymers are being discovered. Although this class of polymers is very young, the potential use of these is quite significant.

Conductive polymers have applications in many different areas: from conductive coatings and paints used on airplane fuselages, to fibers used in astronaut suits and clothing, just to

name a few. Being able to predict the conductivity of a polymer before it is even synthesized would be of great practical and financial importance. Conducting polymers have opened up many new possibilities for devices combining unique optical, electrical, and mechanical properties. The literature on conducting polymers, as reflected in its top-cited papers, shows the diversity of materials that can be used, optical effects achieved, and underlying physical processes. A partial list of publications on conductive polymers that may provide an historical perspective, can be seen in [57–64] and references therein. The molecules cited in these publications include organic polymers, copolymers, and conjugated polymers, such as polyacetylene, polyaniline, poly(para-phenylene), and poly(p-phenylenevinylene), to name a few. Some examples include polyphenylene and polythiophene, which showed electrical conductivities of up to 0.1 S cm$^{-1}$ [63]. These values are some of the highest values measured for a conductive polymer, but still less that that of most metals, yet these polymers were the first actually capable of conducting electricity. In addition, many of these polymers have high flexibility, and the devices made from them include light-emitting diodes (LEDs) and lasers, and the color of the light emitted can be chemically tuned. The main physical process involved is electroluminescence.

Another significant milestone was made when the polymer polyacetylene was exposed to traces of iodine or bromine vapor. This thin polymer film exhibited still higher electrical conductivity. The researchers discovered that by adding various selected impurities to polyacetylene, its electrical conductivity could be made to range widely—behaving as an insulator, like glass, to a conductor, like metal. The key breakthrough leading to practical application as batteries occurred in 1979 when one of Professor MacDiarmid's graduate students was investigating alternative ways for doping polyacetylene [62–64]. He placed two strips of polyacetylene in a solution containing the doping ions and passed an electric current from strip to strip. As expected, the positive ions migrated to one strip and the negative ions to the other. But when the current source was removed, the charge remained stored in the polyacetylene polymer. This stored charge could then be discharged if an electrical load was connected between the two strips, just as in a conventional battery [62].

Conducting polymers can also be induced to transfer electrons to other materials such as buckminsterfullerene. It is this last application, along with many others, which makes conductive polymers important in the field of nanotechnology and molecular electronics. The relative conductivities of some of the polymers synthesized are shown below, along with a comparison with copper metal and liquid mercury. Figure 3 below shows the relative conductivity of several conductive polymers compared to that of copper metal and liquid metallic mercury [57, 65]. Further details can be found on an excellent Web page by Colin Pratt [57].

Since then, it has been found that a little more that a dozen different polymers and polymer derivatives can undergo a transition to a good conducting material after being doped with a weak oxidation agent or reducing agent. These polymers are all various types of conjugated polymers. Early work has led to an understanding of the mechanisms of charge storage and charge transfer in these system. All have a highly conjugated electronic state. This also causes the main problems with the use of these systems, that of processibility and stability. Most early conjugated polymers were unstable in air and were not capable of being processed. The most recent research in this has been the development of highly conducting polymers with good stability and acceptable processing attributes.

## 6.3. Conductivity and Charge Storage

When it was decided that this phenomenon of conductivity in a polymer film was real and its significance was understood, an explanation was needed to help explain the mechanism of conductivity. One of the earliest explanations of conducting polymers borrowed from the theory of solid-state physics and used band theory as a method of explaining conductivity. This theory used the fact that a half-filled valence band would be formed from a continuous delocalized $\pi$-system. This would be an ideal condition for conduction of electricity. However, it turns out that the polymer can more efficiently lower its energy by bond alteration (alternating single and double bonds), which introduces a band width of about 1.5 eV,

Logarithmic conductivity ladder locating some metals and conducting polymers

**Figure 3.** Relative conductivity of several conductive polymers compared to that of copper metal and iquid metallic mercury. Image courtesy of Colin Pratt.

making it a high energy gap semiconductor. The polymer is transformed into a conductor by doping it with either an electron donator or an electron acceptor. This is reminiscent of doping of silicon-based semiconductors, where silicon is doped with either arsenic or boron. However, while the doping of silicon produces a donor energy level close to the conduction band or an acceptor level close to the valence band, this is not the case with conducting polymers. The evidence for this is that the resulting polymers do not have a high enough concentration of free spins, as determined by electron spin spectroscopy.

Initially, the free spin concentration increases with concentration of dopant. At larger concentrations, however, the concentration of free spins levels off at a maximum. To understand this, it is necessary to examine the way in which charge is stored along the polymer chain and its effect.

The polymer has the ability to store charge in at least two ways. In an oxidation process, it could either lose an electron from one of the bands or it could localize the charge over a small section of the chain. Localizing the charge causes a local distortion due to a change in geometry, which costs the polymer some energy. However, the generation of this local geometry decreases the ionization energy of the polymer chain and increases its electron affinity, making it more able to accommodate the newly formed charges. This method increases the energy of the polymer less than it would have if the charge was delocalized and, hence, takes place in preference of charge delocalization. This is consistent with an increase in disorder detected after doping by Raman spectroscopy. A similar scenario occurs for a reductive process.

Typical oxidizing dopants used include elements such as iodine, arsenic pentachloride, iron(III) chloride, and NOPF₆[47]. A typical reductive dopant could be sodium naphthalide. The main criteria is its ability to oxidize or reduce the polymer without lowering its stability or whether or not they are capable of initiating side reactions that inhibit the polymer's ability to conduct electricity. An example of the latter is the doping of a conjugated polymer with bromine. Bromine it too powerful an oxidant and adds across the double bonds to form $sp^3$ carbons. The same problem may also occur with NOPF₆ if left reacting too ong.

The oxidative doping of polypyrrole proceeds in the following way. An electron is removed from the p-system of the backbone, producing a free radical and a spinless positive charge. The radical and cation are coupled to each other via local resonance of the charge and the

radical. In this case, a sequence of quinoid-like rings is used. The distortion produced by this is of higher energy than the remaining portion of the chain. The creation and separation of these defects costs a considerable amount of energy. This limits the number of quinoid-like rings that can link these two bound species together. In the case of polypyrrole, it is believed that the lattice distortion extends over four pyrrole rings. This combination of a charge site and a radical is called a polaron. This could be either a radical cation or radical anion. This creates new localized electronic states in the gap, with the lower energy states being occupied by a single unpaired electrons. The polaron state of polypyrrole is symmetrically located about 0.5 eV from the band edges.

Upon further oxidation, the free radical of the polaron is removed, creating a new spinless defect called a bipolaron. This is of lower energy than the creation of two distinct polarons. At higher doping levels, it becomes possible that two polarons combine to form a bipolaron. Thus, at higher doping levels, the polarons are replaced with bipolarons. The bipolarons are located symmetrically with a band gap of 0.75 eV for polypyrrole. This eventually, with continued doping, forms into continuous bipolaron bands. Their band gap also increases as newly formed bipolarons are made at the expense of the band edges. For a very heavily doped polymer, it is conceivable that the upper and the lower bipolaron bands will merge with the conduction and the valence bands, respectively, to produce partially filled bands and metallic-like conductivity. This is shown in Fig. 4. Conjugated polymers with a degenerate ground state have a slightly different mechanism. As with polypryyole, polarons and



Figure 4. Illustration of several mechanisms for the conductivity of conductive polymers, which also shows the energy levels for the polaron and bipolaron mechanism of conductivity for various conductive polymers. Image courtesy of Colin Pratt.

bipolarons are produced upon oxidation. However, because the ground state structure of such polymers are twofold degenerate, the charged cations are not bound to each other by a higher energy bonding configuration and can freely separate along the chain. The effect of this is that the charged defects are independent of one another and can form domain walls that separate two phases of opposite orientation and identical energy. Figure 4 illustrates several mechanisms for the conductivity of conductive polymers and also shows the energy levels for the polaron and bipolaron mechanism of conductivity for various conductive polymers [47, 66].

These states are called solitons and can be charged or neutral. Solitons produced in polyacetylene are believed to be delocalized over about 12 C–H units with the maximum charge density next to the dopant counterion. The bonds closer to the defect show less amount of bond alternation than the bonds away from the center. Soliton formation results in the creation of new localized electronic states that appear in the middle of the energy gap. At high doping levels, the charged solitons interact with each other to form a soliton band that can eventually merge with the band edges to create true metallic conductivity. This is shown below in Fig. 5 [57, 67].

## 6.4. Charge Transport Process

Although solitons and bipolarons are known to be the main source of charge carriers, the precise mechanism is not yet fully understood. The problem lies in attempting to trace the path of the charge carriers through the polymer. All of these polymers are highly disordered, containing a mixture of crystalline and amorphous regions. It is necessary to consider the transport along and between the polymer chains and also the complex boundaries established by the multiple number of phases. This has been studied by examining the effect of doping, of temperature, of magnetism, and the frequency of the current used. These tests show that a variety of conduction mechanisms are used. The main mechanism used is by movement of charge carriers between localized sites or between soliton, polaron, or bipolaron states. Alternatively, where inhomogeneous doping produces metallic island dispersed in an insulating matrix, conduction is by movement of charge carriers between highly conducting domains. It is also known that charge transfer between these conducting domains also occurs by thermally activated hopping or tunneling. This is known to be consistent with conductivity being proportional to temperature [57].

## 6.5. Stability

There are two distinct types of stability. Extrinsic stability is related to the vulnerability to external environmental agents such as oxygen, water, peroxides, and so forth. This is determined by the susceptibility of the polymer's charged sites to attack by nucleophiles,



**Figure 5.** Illustration of energy levels for the polaron and soliton mechanism of conductivity for conductive polymers along with their corresponding energy bands. Image courtesy of Colin Pratt.

electrophiles, and free radicals. If a conducting polymer is extrinsic unstable, then it must be protected by a stable coating.

Many conducting polymers, however, degrade over time even in dry, oxygen-free environments. This intrinsic instability is thermodynamic in origin. It is most likely caused by irreversible chemical reactions between charged sites of polymers and either the dopant counterion or the p-system of an adjacent neutral chain, which produces an $sp^3$ carbon, breaking the conjugation. Intrinsic instability can also come from a thermally driven mechanism that causes the polymer to lose its dopant. This happens when the charge sites become unstable due to conformational changes in the polymer backbone. This has been observed in alkyl- substituted polythiophenes.

## 6.6. Possibilities?

Conjugated polymers may be made by a variety of techniques, including cationic, anionic, radical chain growth, coordination polymerization, step growth polymerization, or electrochemical polymerization. Electrochemical polymerization occurs by suitable monomers that are electrochemically oxidized to create an active monomeric and dimeric species, which react to form a conjugated polymer backbone. The main problem with electrically conductive plastic stems from the very property that gives it its conductivity, namely the conjugated backbone. This causes many such polymers to be intractable, insoluble films or powders that cannot melt. There are two main strategies to overcoming these problems. They are either to modify the polymer so that it may be more easily processed or to manufacture the polymer in its desired shape and form. There are, at this time, four main methods used to achieve these aims.

The first method is to manufacture a malleable polymer that can easily be converted into a conjugated polymer. This is done when the initial polymer is in the desired form and then, after conversion, is treated so that it becomes a conductor. The treatment used is most often thermal treatment. The precursor polymer used is often made to produce a highly aligned polymer chain, which is retained upon conversion. This is used for highly orientated thin films and fibers. Such films and fibers are highly anisotropic, with maximum conductivity along the stretch direction.

The second method is the synthesis of copolymers or derivatives of a parent conjugated polymer with more desirable properties. This method is the more traditional one for making improvements to a polymer. What is done is to try to modify the structure of the polymer to increase its processibility without compromising its conductivity or its optical properties. All attempts to do this on polyacetylene have failed, as they always significantly reduced its conductivity. However, such attempts on polythiophenes and polypyrroles proved more fruitful. The hydrogen on carbon 3 on the thiophene or the pyrrole ring was replaced with an alkyl group with at least four carbon atoms in it. The resulting polymer, when doped, has a comparable conductivity to its parent polymer while being able to melt, and it is soluble. A water-soluble version of these polymers has been produced by placing carboxylic acid group or sulfonic acid group on the alkyl chains. If sulfonic acid group, groups are used along with built-in ionizable groups. Then, such a system can maintain charge neutrality in its oxidized state, and so they effectively dope themselves. Such polymers are referred to as "self-doped" polymers. One of the most highly conductive derivatives of polythiophene is made by replacing the hydrogen on carbon 3 with a $-CH_2-O-CH_2CH_2-O-CH_2CH_2-O-CH_3$. This is soluble and reaches a conductivity of about 1000 S cm$^{-1}$ upon doping.

The third method is to grow the polymer into its desired shape and form. An insulating polymer impregnated with a catalyst system is fabricated into its desired form. This is then exposed to the monomer, usually a gas or a vapor. The monomer then polymerizes on the surface of the insulating plastic, producing a thin film or a fiber. This is then doped in the usual manner. A variation of this technique is electrochemical polymerization with the conducting polymer being deposited on an electrode either at the polymerization stage or before the electrochemical polymerization. This cast may be used for further processing of the conducting polymer. For instance, by stretching aligned bends of polyacetylene/polybutadiene, the conductivity increases 10-fold, due to the higher state of order produced by this deformation.

The final method used was the use of Langmuir-Blodgett trough to manipulate the surface active molecules into highly ordered thin films whose structure and thickness are controllable at the molecular layer. Amphiphilic molecules with hydrophilic and hydrophobc groups produces monolayers at the air–water surface interface of a Langmuir-Blodgett trough. This is then transferred to a substrate, creating a multilayer structure comprised of molecular stacks, which are typically about 2.5-nm thick. This is a development from the creation of insulating films by the same technique. The main advantage of this technique is ts unique ability to allow control over the molecular architecture of the conducting films produced. It can be used to create complex multilayer structures of functionally different molecular layers as determined by the chemist. By using alternating layers of conductor and insulator, it is possible to produce highly anisotropic film that is conducting within the plane of the film, but insulating across it. Table 1 shows various conductive polymers, their conductivity, stability, and processability [57].

## 6.7. Some Applications of Conductive Polymers

The extended $\pi$-systems of conjugated polymer are highly susceptible to chemical and electrochemical oxidation or reduction. These alter the electrical and optical properties of these polymers by controlling this oxidation and reduction. It is actually possible to precisely control these properties using the possible electrochemical states of the molecule. In addition, reactions are often reversible, therefore it is possible to systematically control the electrical and optical properties with a great deal of precision. It is even possible to switch from a conducting state to an insulating state and back to a conducting state many times.

There are two main groups of applications for these polymers. The first group uses their conductivity as its main property. The second group uses their electroactivity. They are shown in Table 2 [47].

Conductive polymer demand in the United States is expected to grow about 6.5% annually through 2006 [68]. Advances will result from more intensive use in many products including electronics in motor vehicles, appliances, and numerous portable consumer devices, such as portable radios, video and audio recorders, cell phones, cameras, to name only a ew.

It is estimated that the United States has a $1 billion conductive polymer industry. For an in-depth analysis, the reader is referred to [68] and references therein. This book presents historical (1992, 1996, and 2001) data and forecasts to 2006 by technology (e.g., carbon black powder–filled, metallized, paint-coated, fiber-filled); by resin (e.g., ABS, PVC, polyphenylene-based, PE, PP, TPE, nylon, polystyrene, inherently conductive); by function; by application (e.g., product components, antistatic packaging, materials handling, worksurface and flooring); and by market. The study also examines the market environment, details industry structure and market share, and profiles 51 key companies including GE Plastics, Noveon, DSM Engineering Plastics, LNP Engineering Plastics, PolyOne, RTP, and Illinois Tool Works.

## 6.8. Molecular Wires

A molecular wire can be defined as consisting of a molecule connected between two reservoirs of electrons. The molecular orbitals of the molecules must provide a favorable pathway

Table 1. Stability and processing attributes of several conducting polymers.

| Polymer | Conductivity ($\Omega^{-1}$ cm$^{-1}$) | Stability (doped state) | Processing possibilities |
|---|---|---|---|
| Polyacetylene | $10^{5}$–$10^{7}$ | Poor | Limited |
| Polyphenylene | 1000 | Poor | Limited |
| PPS | 100 | Poor | Excellent |
| PPV | 1000 | Poor | Limited |
| Polypyrroles | 100 | Good | Good |
| Polythiophenes | 100 | Good | Excellent |
| Polyaniline | 10 | Good | Good |

Table 2. Two groups of conductive polymers illustrating conductivity as a main property (group 1) and electroactivity as the main property (group 2).

| Group 1 | Group 2 |
| --- | --- |
| Electrostatic materials | Molecular electronics |
| Conducting adhesives | Electrical displays |
| Electromagnetic shielding | Chemical and biochemical sensors |
| Printed circuit boards | Rechargeable batteries and solid electrolytes |
| Artificial nerves | Drug-release systems |
| Antistatic clothing | Optical computers |
| Thermal sensors | Ion-exchange membranes |
| Piezoceramics | Electromechanical actuators |
| Active electronics | "Smart" structures |
| Molecular switches | |

for electrons when they connect between the leads. Such a system was first suggested in the early 1970s by Aviram and Ratner to have the ability to rectify current [30]. Experimentally, research on molecular wires has increased substantially over the past several years, with significant developments in rectification and other conductive phenomena being made [69–79]. There has also been an increase in the theoretical modeling of these systems [80–85]. For an overview of the current status of the molecular electronics field, see [86] and references therein. Additional references on molecular wires and their synthesis can be found in [87–104] and references therein.

Theoretical studies of the electronic conductance of a molecular wires bring together many different methods from chemistry and physics. Quantum theory is used to model the energetics of molecules. It is also incorporated into the study of the coupling between the molecule and the metallic reservoirs. Once these issues have been addressed, it is possible to proceed to the electron transport problem. Currently, Landauer theory [57; for a comprehensive review of Landauer theory, see 105] is used, which relates the conductance to the electron transmission probability. One of the molecules of current experimental interest as a molecular wire is 1,4-benzene-dithiolate (BDT) [106]. It consists of a benzene with two sulfur atoms attached, one on either end of the benzene ring. The sulfurs can bond effectively to gold nanocontacts, and the conjugated $\pi$ ring provides the delocalized electrons necessary for conduction. Two major unknowns of the experimental system are the geometry of the gold contacts and the nature of the bond between the molecule and these contacts.

For mesoscopic systems with discrete energy levels connected to continuum reservoirs, the transmission probability displays resonance peaks. Another possible transport phenomenon that has been predicted is the appearance of antiresonances [84, 107]. These occur when the transmission probability is zero and correspond to the incident electrons being perfectly reflected by the molecule. There is a simple mechanism controlling where the antiresonances occur in the transmission spectrum. This may be analyzed by applying a simple formula to the case of a molecular wire consisting of an "active" molecular segment connected to two metal contacts using a pair of finite $\pi$ conjugated chains. This calculation can show how an antiresonance can be generated near the Fermi energy of the metallic leads. This antiresonance is characterized by a noticeable drop in conductance. It is found that for this calculation, the analytic theory of antiresonances has predictive power [107].

A very interesting paper [108] involves the carbon-assisted synthesis of silicon nanowires. There has been a large research activity in the synthesis of inorganic nanowires and nanotubes in the past few years [30, 109, 110]. According to this manuscript, the use of carbon-assisted synthesis of silicon nanowires has been accomplished with silicon powders as well as solid substrates. The technique for synthesis involves heating an intimate mixture of silicon powder and activated carbon or a carbon-coated solid substrate in argon at 1200–1350°C and yields abundant quantities of crystalline nanowires. Besides being simple, the method eliminates the use of metal catalysts.

Silicon nanowires have received considerable attention, and several methods have been employed for their synthesis. These include thermal evaporation of Si powder [111], vapor–liquid–solid method involving liquid metal solvents with low solubility for Si [112, 113],

laser ablation [114, 115], and the use of silicon oxide in mixture with Si [116, 117 $SiO_2$-sheathed crystalline SiNWs have been obtained by heating Si–$SiO_2$ mixtures [118]. A excellent paper by Sharma and Sunkara reported that enhanced yields of silicon nanowires are obtained by heating a Si substrate coated with carbon nanoparticles at 1050°C under acuum [119]. The role of carbon is considered to be as it is in other carbothermal methodsof synthesizing nanowires of oxides, nitrides, and other materials, that is involving a vapor–solid mechanism, wherein the carbon reacts with the oxide, probably producing a suboxide-type species. The formation of silicon nanowires in the presence of carbon may be explaned by considering silicon to generally be covered by an oxide layer. This oxide layer is rediced by carbon into silicon monoxide by the reaction. The paper by Sharma and Sunkara also shows that crystalline silicon nucleates and grows perpendicular to the (111) direction to frm the nanowires [119]. Similar reactions have been proposed for the oxide-assisted syntiesis of silicon nanowires [115], although the monoxide-type species is generated by other nzans.

Researchers from Nanosys, Inc., have found a way to assemble large arrays of naiowires constructed from silicon or other semiconductor material into a densely packed thn film, then process the assembly to produce relatively efficient transistors on a variety of sirfaces. The technology may eventually be used to construct very large flat-panel displays, disiosable computing and storage electronics, and tiny radio-frequency identification devices. Tle work appears in the September 18, 2003, issue of *Nature* [120]. Researchers at Nanosys, Inc., also demonstrated that cadmium sulfide nanoribbons can also be used to produce tiin film transistors [121].

Researchers from the Hahn-Meitner Institute in Germany developed a new way b make flexible transistors. The method causes forests of vertical semiconductor nanowires o grow inside a plastic film. These thin films contain as many as 100 million nanowires persquare centimeter. The semiconductor nanowires serve as transistor channels, a metal laye serves as the gate electrode, and source and drain electrodes were added to the top and botom of the stack to complete the transistors. The gate electrode blocks or allows electricityto flow through a transistor channel [122].

An example of the molecular structure of a molecular wire made from a nanotubecan be seen in Fig. 6.

## 6.9. Molecular Switches

At the center of the microelectronics industry is the semiconductor switch. Becaus semiconductor switches can be manufactured at very small scales, and in combinationcan be made to perform all desired computational functions, the microelectronic switch has iecome



Figure 6. Typical single-walled nanotube used as a molecular wire.

the fundamental device in virtually all modern electronics. More than 25 years ago, Intel cofounder Gordon Moore observed that the number of transistors on a given piece of silicon would double every couple of years—a profound insight that was dubbed "Moore's law," and a law that still holds today [123]. In the future, however, the laws of physics will begin to play a role. Processors now come with dimensions measuring about 130 nanometers. Next year, it is estimated that chips will have features measuring around 100 nanometers. When these features are shrunk to 30 nanometers—a goal that is estimated to be achieved in 7 to 8 years—designers will begin to hit a design wall. Alternatives to current chips are seen in such new technologies as nanotechnology. However, Gordon Moore says, "crafting single transistors is one thing, but housing a billion of them on a chip is another" [124].

The switch pictured in Fig. 7 is a single molecule that exhibits classical switching properties [125]. It was developed by the California Molecular Electronics Corporation (CALMEC).

According to CALMEC scientists, the Chiropticene switch is a device that goes beyond the semiconductor switch in size reduction and cost. This switch is a single molecule that exhibits classical switching properties. Being only a molecule in size, it is hundreds of times smaller than even the smallest semiconductor switch. Chiropticene molecules are switchable between two distinct states that are spatial mirror images of each other. These mirror images are electronically and optically distinct, enabling sharp and stable switching properties. According to the company's Web site, some of the advantages of such a nanomolecular switch and molecules with similar properties are:

Stability: Two equal but opposite energy states in these molecules affords stability while assuring complete reversibility.

Speed: Electrical field switching will potentially provide femtosecond computational switching times. Optical switching will potentially provide nanosecond computational switching times.

Nanoassembly: The molecules will lend themselves to the new techniques of nanotechnology self-assembly, enabling the assembly of supramolecular device architectures.

Photonic advantage: The Chiropticene molecule capitalizes in novel ways on the unique properties of light in data manipulation: high bandwidth, frequency domain modulation, diffraction, refraction, reflection, superposition, and parallelism.

Threshold protection: An intervening neutral state prevents optical switching without electrical stimulation.

Molecular engineering: By the judicious selection of constituents, the Chiropticene molecule can be tuned to respond to selected laser frequencies and can be engineered to meet specific performance requirements.

Nondestructive readouts: The molecule can be interrogated by optical rotation without energy absorption.

Commercial attractiveness: Chiropticene-based devices are expected to have no moving parts and operate at room temperature. They are also expected to be enormously cost-effective to produce. For additional details on this and similar molecules, the reader is referred to the material found in California Molecular Electronics Corporation Web page [125].

Further advances in molecular switches have come from transistor design modifications. Recently, a group from IBM reported that their so-called top-gated nanotube field-effect



Figure 7. Chiropticene switch molecule. Image courtesy of CALMEC.

transistors (FETs) outperform state-of-the-art silicon FETs in terms of switching rate and the amount of current they can carry per width of conductor. One key difference between the latest design and earlier designs is that the silicon wafers that support the FETs do not function as gates. Instead, the gate is fabricated above the nanotube—allowing all FETs in contact with the silicon wafer to be switched independently. In addition, the new FET design benefits from switching voltages that are an order of magnitude lower than those needed to switch older FETs [126].

Just recently, the IBM team developed a catalyst-free procedure for preparing single-walled carbon nanotubes. Conventional methods for preparing single-walled tubes, such as discharge and vapor deposition techniques, rely on particles of nickel and cobalt or other catalytic transition metals. A problem with those procedures is that the metal particles left behind in the products disturb the electrical properties of the nanotubes, forcing scientists to take complex purification measures.

Avouris and coworkers have shown that single-walled nanotubes can be prepared from SiC in the absence of metal catalysts at temperatures above 1500°C [127]. The method produces nanotubes that are 1.2 to 1.6 nm wide and are aligned along the face of the support. Other researchers have reported techniques for growing vertically aligned products, which they dub nanotube forests. But nanotubes that grow horizontally are amenable to standard vapor deposition methods and may be connected in parallel to lower their electrical resistance.

Despite the recent advances, Avouris asserts that nanotube FETs remain "far from optimized." Improvements could be made by using thinner insulators with higher dielectric constants, he suggests. "But what's really needed is a better understanding of the mechanism of electrical switching in nanotube FETs" [127].

Carbon nanotubes aren't the only useful nanotechnology. Nanowires made of semiconductors such as silicon, gallium arsenide, and indium phosphide have been investigated as candidate materials for nanoelectronics. The field is advancing rapidly as researchers are making fast progress in synthesis, device fabrication, and testing.

Harvard University chemistry professor Charles M. Lieber and colleagues have published widely on carbon nanotubes and semiconducting nanowires [128–134]. These tubes are versatile building blocks that can be used in a bottom-up approach to constructing nanoelectronic circuits. Because the size, structure, and functional properties of nanowires can be controlled, they are readily produced via known synthetic procedures.

In 1998, Lieber and coworkers described a vapor–liquid–solid synthesis method in which laser light is used to ablate nanometer-sized metal clusters that serve as nucleation centers and catalysts for nanowire growth. The Harvard researchers used the technique to prepare uniform single-crystalline nanowires of silicon and germanium with diameters as small as 3 nm and lengths up to 30 mm. Since that time, the procedure has been used to prepare nanowires with even smaller diameters, and it has been extended to a wide variety of materials. Examples include III–V semiconductors such as GaAs and II–VI semiconductors such as ZnSe.

Recently, a number of research groups boosted the complexity of materials that can be prepared by the cluster-nucleation method. The teams prepared modulated structures—nanowires composed of dissimilar segments. The two-tone materials are made by turning the supplies of reactants on and off during synthesis with pulsed lasers or by other methods. These "heterostructured" products open the door to new sophisticated applications, such as terahertz-frequency photon emitters.

Lieber, Gudiksen, and coworkers used modulation methods to prepare nanowires with 21 alternating segments of GaAs and GaP. And they prepared Si and InP nanowires with modulated doping, such that the wires were endowed with alternating $p$- and $n$-type regions. Using similar methods, Peidong Yang and collaborators prepared Si–SiGe nanowires [135–139]. In addition, Lars Samuelson and colleagues synthesized InAs/GaAs nanowires [140]. These nanowires were grown by chemical beam epitaxy using gold nanoparticles as a catalyst. Photoluminescence measurements showed spectra consisting of sharp lines with energies and excitation power dependency behavior very similar to that observed for Stranski-Krastanow–grown InAs/GaAs [140] and GaN quantum dots [141]. By reducing the excitation power density, they were able to obtain a quantum dot spectrum consisting of only one single sharp

line—the exciton line. An excellent summary of these later developments can be found in an article in *C&E News* by Mitch Jacoby [142].

Table 3 list journals published by Elsevier Publishing Co. and sent monthly to anyone interested, via e-mail. The address to subscribe is e-alert.nano@elsevier.com.

## 6.10. Computational Design and Analysis of Nanomolecular Circuits

The computational techniques that have been discussed in some detail at the beginning of this chapter are powerful techniques for proposing, modeling, and simulating possible nanomolecular-based electronic circuits and components while strictly adhering to the laws of physics. Accurate simulation and approximation of large molecular systems is possible on desktop computers, which these days are primarily PCs. This is due to the development of modern algorithms, combined quantum methods, bigger and better basis sets and functionals, and inexpensive powerful computer platforms. We have previously used molecular mechanics, semiempirical, Hartree-Fock, and DFT computational methods to demonstrate how conductive polymer molecules could be used to design and analyze nanoscale molecular transistors and logic circuits with electronic properties analogous to macroscale electronic components [143, 144]. These molecules are composed of fragments of known conductive polymers, partially insulating polymers, as well as molecules that exhibit semiconducting properties. The properties of the polymer molecules are chosen to achieve conductivity, resistivity, rectification, and amplification. This can be demonstrated by considering the bulk properties of some of these polymers. The emeraldine base form of polyaniline, which has a conductivity of $10^1$ to $10^2$ S/cm, can be used as a molecular wire, and the trans isomer

Table 3. List of journals by Elsevier Publishing Co. containing recent articles on nanotechnology.

*Acta Materialia*
*Applied Catalysis A*
*Applied Catalysis B*
*Applied Surface Science*
*Biomaterials*
*Carbon*
*Chemical Physics Letters*
*Composites Part A*
*Composites Science and Technology*
*Computational Materials Science*
*Electrochemistry Communications*
*Electrochimica Acta*
*Journal of Colloid and Interface Science*
*Journal of Electroanalytical Chemistry*
*Journal of Magnetism and Magnetic Materials*
*Journal of Membrane Science*
*Journal of Non-Crystalline Solids*
*Journal of Nuclear Materials*
*Journal of Photochemistry and Photobiology A*
*Journal of Physics and Chemistry of Solids*
*Journal of the American Society for Mass Spectrometry*
*Journal of the European Ceramics Society*
*Materials Letters*
*Materials Science and Engineering A*
*Microelectronic Engineering*
*Microporous and Mesoporous Materials*
*Physics Letters A*
*Polymer*
*Solar Energy Materials and Solar Cells*
*Solid State Electronics*
*Solid State Ionics*
*Surface and Coatings Technology*
*Surface Science*

of undoped polyacetylene, which has a conductivity of $10^{-5}$ S/cm, can be used as a molecular resistor. We have previously proposed a molecular transistor, based on the molecular recitifier molecule proposed by Aviram and Ratner (A&R) in 1974 [30]. A&R proposed using a tetracyanoquinodimethane (tcnq) molecule connected to a tetrathiafulvalene (ttf) molecule by a triple methylene bridge as a molecular rectifier. Their calculations showed that the proposed molecule would allow a certain current flow in one direction through the circuit, but not in the reverse direction, capable of the behavior of a $p-n$ junction rectifier diode in microelectronic circuits. The mechanism they proposed was based on separation of the $p$ and $n$ molecular fragments with a barrier potential molecule, which requires electron tunneling for transmission. It also relied on the fact that a 1 to 2 eV potential bias would alter the molecular orbital energy levels enough to allow electron transmission from ttf to tcnq, but that conduction in the reverse direction required a 9 eV potential, thus achieving current rectification. A slightly modified version of this $p-n$ rectifier molecule concept was experimentally proven by Metzger et al. in 1997 [145]. We took the idea and further extended the molecule by adding to their molecule another triple methylene bridge and one more tcnq molecule. This extends their proposed $p-n$ junction, making it a $p-n-p$ junction. We believe this molecule to have $p-n-p$ junction transistor properties. Figure 8 is a wire frame representation of the proposed molecule.

Methods analogous to those used to analyze solid-state devices were used to analyze the current–voltage characteristics of this molecule. A number of assumptions were used in order to make the calculations possible. Conductive polyaniline (pani) polymer molecules were used for emmiter and collecter contacts at both tcnq ends of the molecule and in the middle of the molecule (at ttf) for the base contact. Each $n$ or $p$ fragment of the molecule (i.e., each ttf and tcnq) was treated as a separate molecular entity by assuming the methylene bridges (∩) would effectively separate the electronic $\pi$ states of the ttf and tcnq fragments, as there is no $\pi$ conjugation across the triple methylene bridge. According to the calculations, a threshold forward bias >0.4 eV from base (ttf) to collector (tcnq–pani) is necessary to raise the energy of the occupied levels of ttf to that of the unoccupied levels of tcnq, causing a threshold forward bias >1.27 eV across the whole molecule (i.e., pani → tcnq → ∩ → ttf → ∩ → tcnq → pani) to be amplified [143]. This transistor molecule was then used as a component of a simulated logic AND gate molecule, composed of two $p-n-p$ transistor molecules, polyaniline polymer molecules as molecular wires, polyacetylene molecules as resistors, and benzene rings as conductive junctions for the polyaniline. Figure 9 is a wire frame representation of the large AND gate circuit molecule that was simulated.

The dashed lines in Fig. 9 are used to denote the locations of the $p-n-p$ transistor molecules. The A and B inputs are through the polyacetylene resistors, which are attached to the $p-n-p$ transistor next to the ttf molecule. The circuit molecule would operate in the following manner. If a background forward bias >1.27 eV is applied across the whole molecule, both transistors would have to be activated from inputs A and B with a voltage >4.32 eV in order to provide an output, which would be at the benzene ring (the junction point of the EB-pani leads coming from the $p-n-p$ transistors and the resistor to the ground). Because the $p-n-p$ transistor molecules are linked together to complete the loop to the output, if one of the transistors is inactivated (i.e., A or B input <4.32 eV), there will be no current flow to the output of the circuit [143]. Thus, the circuit would behave as an AND logic gate.



Figure 8. Proposed $p-n-p$ transistor based on A&R's 1974 molecular rectifier.

**Figure 9.** Molecular AND circuit based on $p$-$n$-$p$ transistor molecules.

At the time we proposed these organic molecule transistors and polymer molecule circuits, we became aware that single-wall carbon nanotubes could be used to achieve the same or higher goals in molecular electronics. We not only recognized the potential of using single-wall nanotubes as molecular wires in nanosized electronic circuits, but also imagined that a single-wall carbon nanotube's electronic properties could be altered through the intercalation of transition metal atoms into the cavity of the nanotube. Our hypothesis was that a transition metal atom's outer $d$ electrons would interact with the conductive valence "band" of the nanotube, causing an altered hybrid electronic state, which would depend on the type of metal atom used. We further theorized that it might be possible to obtain $p$- or $n$-type semiconductor characteristics in the carbon nanotube depending on which metal is used to "dope" the tube.

We recognized from the start that this was a lofty goal, even when considering the faster generation of PCs that started to be available for calculations 2–3 years ago. The highest occupied molecular orbital (HOMO) to lowest unoccupied molecular orbital (LUMO) energy difference [a molecular orbital energy (MOE) gap analogous to a band gap] was used as an estimate of the conductive characteristics of the metal "doped" nanotube. This required accurate approximations of molecular energy levels for these systems. Small (hydrogen terminated) single-wall nanotubes composed of 50 carbon atoms, approximately 3.3 Å in diameter, and substituted with different atoms as high as period 6 on the periodic table were simulated. Substituting with increasingly heavy metal atoms was no small task, as relativistic effects start to affect the calculation in these increasingly heavy nuclei. For this reason, the Los Alamos National Laboratory Double Zeta (LANL2DZ) basis set was used within the Gaussian94 [146] package of programs. The LANL2DZ basis set uses an effective core potential (ECP) to represent the innermost (core) electrons and incorporates parameters that account for the relativistic effects that become more prominent in these heavier atoms and are mostly restricted to the core electrons. DFT-SCF single-point energy calculations were performed on structurally optimized nanotubes using the B3PW91 functional. This procedure was performed for most of the period 4, 5, and 6 transition metals and several group

**Figure 10.** Plot of the HOMO/LUMO orbital energy gap of small single-wall carbon nanotubes substituted with different transition metal atoms.

IIIA, IVA, and VA nonmetallic elements, each time substituting one atom into the small nanotube, optimizing the structure, and calculating the molecular orbital energy scheme. Figure 10 is a plot of the molecular orbital energy gap of these "doped" systems *versus* the atomic number of the element used for doping.

The general trend in Fig. 10 for the transition metals is a decreasing molecular orbital energy gap as period number increases. There are several factors that may contribute to this effect. As the diameter of the intercalated metal atom increases (increasing atomic number), there is probably an enhanced overlap between the conductive $\pi$ orbitals of the carbon nanotube and the outer $d$ orbitals of the metal atom. This suggests the existence of a hybrid orbital electronic state between the nanotube and metal atom that favors enhanced conductivity. However, there are drastic differences between the energy gaps of the odd and even electron systems. The odd electron systems have energy gaps 0.5 to 1.0 eV smaller than their even electron counterparts, suggesting that paired electron systems are much more stable, requiring considerably more energy to mobilize their electrons, and the odd electron systems have unpaired electrons that are much more unstable and thus mobile.

# 7. CONCLUSION

A lot has been discussed in this chapter that spans over four decades of science and technology. Much has happened in that span of time. Computers have evolved from large, few-transistor mastodons of yesteryear, capable of relatively few operations per second, to sleek microsized machines that can fit into one's palm for the most part and perform billions of operations in the same amount of time. Computational science has also advanced, taking advantage of the computer innovations each decade has offered, and advancing scientific modeling to a high level of precision. It was over 45 years ago that Richard Feynman proposed nanotechnology as the science of the future, and computational scientists have been working on it ever since. Today, it is clear that we are standing on the doorstep of that next scientific revolution. But it is sad to note that even though theoretical science achieved highly accurate modeling methods on powerful computers more than a decade ago and proposed novel uses for molecules and nanosized particulates for even longer, the scientific community has been very slow to recognize this potential. It is easy to say today that this revolution was coming, but in reality, 10 years ago most computational studies that proposed truly novel uses for molecules by exploiting their properties at the nanoscale were not taken seriously

and were hard to publish. It is only in the past 5 years that this message has finally been understood by the rest of the scientific community and the technology has been embraced. Computational scientists have predicted this technology all along.

In addition, there have also been many hopeful extrapolations of nanotechnology to date. One needs only to read the numerous articles found in the popular press, television, and the movies to see how much nanotechnology has been hyped. Predicted applications of nano-technology have run the gamut from killer nanorobots to replacing surgeons to protecting astronauts on deep-space missions. Nanotechnology also has its skeptics. They have pointed out that nanotechnology is based more on wishful thinking than science, and they "cast a jaundiced eye" on most of the claims and call them practically impossible.

The reality of nanotechnology lies somewhere between either side of this debate. Nano-technology has the potential to deliver amazing products, high value-added materials that can be used in commercial applications. It will also be at the core of most future research activity. This technology is expanding and will continue to do so for the foreseeable future. We hope this article will help, even if in some small way, to advance this cause.

# REFERENCES

1. G. Némethy and H. A. Scheraga, *Biopolymers* 3, 155 (1965).
2. N. L. Allinger and J. T. Sprague, *J. Am. Chem. Soc.* 95, 3893 (1973).
3. J. T. Sprague, J. C. Tai, Y. Yuh, and N. L. Allinger, *J. Comput. Chem.* 8, 581 (1987).
4. T. Liljefors, J. C. Tai, S. Li, and N. L. Allinger, *J. Comput. Chem.* 8, 1051 (1987).
5. J. C. Tai, and N. L. Allinger, *J. Am. Chem. Soc.* 110, 2050 (1989).
6. J. C. Tai, J.-H. Lii, and N. L. Allinger, *J. Comput. Chem.* 10, 635 (1989).
7. J. A. Pople and G. A. Segal, *J. Chem. Phys.* 43 (Suppl), S136 (1965).
8. N. C. Baird and M. J. S. Dewar, *J. Chem. Phys.* 50, 1262 (1969).
9. M. J. S. Dewar and E. Haselbach, *J. Am. Chem. Soc.* 92, 590 (1970).
10. M. J. S. Dewar and W. Thiel, *J. Am. Chem. Soc.* 99, 4899 (1977).
11. M. J. S. Dewar, E. G. Zoebisch, E. F. Healy, and J. J. P. Stewart, *J. Am. Chem. Soc.* 107, 3902 (1985).
12. J. J. P. Stewart, *J. Comp. Chem.* 10, 209 (1989).
13. M. Born and J. R. Oppenheimer, *Ann. Physik* 84, 457 (1927).
14. I. Levine, in "Quantum Chemistry" 4th Edition. p. 455. Prentice Hall, Englewood Cliffs, NJ.
15. E. Clementi and C. Roetti, *At. Data Nucl. Data Tables* 14, 177 (1974).
16. S. M. Blinder, *Am. J. Phys.* 33, 431 (1965).
17. D. M. Hirst, in "A Computational Approach to Chemistry." Blackwell Scientific Publications.
18. C. C. J. Roothaan, *Rev. Mod. Phys.* 23, 69 (1951); G. G. Hall, *Proc. Roy. Soc.* A205, 541 (1951).
19. R. W. Godby, M. Schlüter, and L. J. Sham, *Phys. Rev. Lett.* 56, 2415 (1986).
20. R. W. Godby and R. J. Needs, *Phys. Rev. Lett.* 62, 1169 (1989).
21. R. W. Godby, L. J. Sham, and M. Schlüter, *Phys. Rev. Lett.* 65, 2083 (1990).
22. R. W. Godby and L. J. Sham, *Phys. Rev. B* 49, 1849 (1994).
23. W. Knorr and R. W. Godby, Investigating exact density-functional theory of a model semiconductor, *Phys. Rev. Lett.* 68, 639 (1992).
24. X. Gonze, P. Ghosez, and R. W. Godby, *Phys. Rev. Lett.* 78, 294 (1997).
25. A. Schindlmayr and R. W. Godby, *Phys. Rev. B* 51, 10427 (1995).
26. P. Sanchez-Friera and R. W. Goldby, *Phys. Rev. Lett.* 85, 5611 (2000).
27. K. Ohno, K. Esfarjani, and Y. Kawazoe, "Computational Materials Science." Springer, Berlin, 1999.
28. P. Hohenberg and W. Kohn, *Phys. Rev.* 136, B864 (1964).
29. Available at http://www.zyvex.com/nanotech/feynman.html.
30. A. Aviram and M. A. Ratner, *Chem. Phys. Lett.* 29, 277 (1974).
31. R. M. Metzger, *J. Am. Chem. Soc.* 119, 10455 (1997).
32. M. A. Reed, *Proceedings of the IEEE*. 87, 652 (1999).
33. J. M. Tour, "Molecular Electronics: Chemical Insights, Chemistry, Devices, Architecture and Programming." World Scientific, New Jersey, 2003.
34. Available at http://main.amu.edu.pl/~tme01/TME03.htm.
35. Available at http://whatis.techtarget.com/definition/0,,sid9_gci517746.00.html.
36. Stu Borman, *C & E News* 80, 46 (2002).
37. M. S. Gudiksen, L. J. Lauhon, J. Wang, D. C. Smith, and C. M. Lieber, *Nature* 415, 617 (2002).
38. M. T. Bjork, B. J. Ohlsson, T. Sass, A. I. Persson, C. Thelander, M. Magnusson, K. Deppert, L. R. Wallenberg, and L. Samuelson, *Nano Lett.* 2, 87 (2002).
39. Y. Wu and P. Yang, *Chem. Mater.* 12, 605 (2000).
40. Y. Wu, R. Fan, and P. Yang, *Nano Lett.* 2, 83 (2002).
41. H. O. Jacobs, A. R. Tao, A. Schwartz, D. H. Gracias, and G. M. Whitesides, *Science* 296, 323 (2002).
42. R. Penterman, S. Klink, H. de Koning, G. Nisato, and D. J. Broer, *Nature* 417, 55 (2002).

*43.* K. Robbie, D. J. Broer, and M. J. Brett, *Nature* 399, 764 (1999).

*44.* S. J. Wind, J. Appenzeller, R. Martel, V. Derycke, and Ph. Avouris, *Appl. Phys. Lett.* 80, 3817 (2002

*45.* Available at http://www.nanoelectronicsplanet.com/nanochannels/research/article/0, 10497_1141291,0.html.

*46.* J. Park, A. N. Pasupathy, J. I. Goldsmith, C. Chang, Y. Yaish, J. R. Petta, M. Rinkoski, J. P. Setha, H. D. Abruna, P. L. McEuen, and D. C. Ralph, *Nature* 417, 722 (2002).

*47.* W. Lu, A. G. Fadeev, B. Qi, E. Smela, B. R. Mattes, J. Ding, G. M. Spinks, J. Mazurkiewicz, D. Zhu, G. G. Wallace, D. R. MacFarlane, S. A. Forsyth, and M. Forsyth, *Science* 297, 983 (2002).

*48.* S. D. Bu, D. M. Kim, J. H. Choi, J. Giencke, E. E. Hellstrom, D. C. Larbalestier, S. Patnaik, L. Cooey, C. B. Eom, J. Lettieri, D. G. Schlom, W. Tian, and X. Q. Pan, *Appl. Phys. Lett.* 81, 1851 (2002).

*49.* S. Chung, J. Bloking, and Y. Chiang, *Nat. Mater.* 1, 123 (2002).

*50.* K. Hayashi, S. Matsuishi, T. Kamiya, M. Hirano, and H. Hosono, *Nature* 419, 462 (2002).

*51.* H. Kawazoe, M. Yasukawa, H. Hyodo, M. Kurita, H. Yanagi, and H. Hosono, *Nature* 389, 939 (1997.

*52.* H. Liu, J. J. Schmidt, G. D. Bachand, S. S. Rizk, L. L. Looger, H. W. Hellinga, and C. D. Montemano, *Nat. Mater.* 1, 173 (2002).

*53.* F. Cacialli, J. S. Wilson, J. J. Michels, C. Daniel, C. Silva, R. H. Friend, N. Severin, P. Samori, J. P. Roe, M. J. O'Connell, P. N. Taylor, and H. L. Anderson, *Nat. Mater.* 1, 160 (2002).

*54.* C. E. Olson, M. J. R. Previte, and J. T. Fourkas, *Nat. Mater.* 1, 225 (2002).

*55.* T.-H. Lee, J. I. Gonzalez, and R. M. Dickson, *Proc. Natl. Acad. Sci. U.S.A* 99, 10272 (2002).

*56.* T.-H. Lee and R. M. Dickson, *Proc. Natl. Acad. Sci. U.S.A* 100, 3043 (2003).

*57.* Colin Pratt: Available at http://homepage.dtn.ntl.com/colin.pratt/cpoly.htm.

*58.* M. F. Rubner, "Molecular Electronics," Research Studies Press, 1992.

*59.* L. Alcacer, "Conducting Polymers," D. Reidel Publishing Company, 1987.

*60.* H. Naarmann, "Polymers to the Year 200 and Beyond," John Wiley & Sons, New York, 1993.

*61.* W. R. Salaneck, D. T. Clark, and E. J. Samuelsen, "Science and Application of Conducting Polymrs." IOP Publishing, 1991.

*62.* R. B. Kaner and A. G. MacDiarmid, *Sci. Amc.* 60 (1988).

*63.* C. K. Chiang, C. R. Fincher, Jr., Y. W. Park, A. J. Heeger, H. Shirakawa, E. J. Louis, S. C. Gau, ad A. G. MacDiarmid, *Phys. Rev. Lett.* 39, 1098 (1977).

*64.* C. R. Fincher, Jr., Y. W. Park, and A. J. Heeger, *Phys. Rev. Lett.* 39, 1098 (1977).

*65.* I. M. Cambell, "Introduction to Synthetic Polymers," p. 196, Oxford Science Publications, 1994.

*66.* M. F. Rubner, "Molecular Electronics," (G. J. Ashwell, Ed.), p. 79, Research Studies Press, 1992.

*67.* M. F. Rubner, "Molecular Electronics," (G. J. Ashwell, Ed.), p. 81, Research Studies Press, 1992.

*68.* "Conductive Polymers to 2006" (Study No. 1563), Freedonia Group, Inc., 2002.

*69.* C. Zhou, M. R. Deshpande, and M. A. Reed, *Appl. Phys. Lett.* 71, 611 (1997).

*70.* M. A. Reed, C. Zhou, C. J. Muller, T. P. Burgin, and J. M. Tour, *Science* 278, 252 (1997).

*71.* S. Frank, P Poncharal, Z. L. Wang, and W. A. de Heer, *Science* 280, 1744 (1998).

*72.* S. J. Tans, M. H. Devoret, and C. Dekker, *Nature* 386, 474 (1997).

*73.* R. P. Andres, J. D. Bielefeld, J. I. Henderson, D. B. Janes, V. R. Kolagunta, C. P. Kubiak, W. J. Mahney, and R. G. Osifchin, *Science* 273, 1690 (1996).

*74.* C. A. Mirkin, R. L. Letsinger, R. C. Mucic, and J. J. Storhoff, *Nature* 382, 607 (1996).

*75.* L. A. Bumm, J. J. Arnold, M. T. Cygan, T. D. Dunbar, T. P. Burgin, L. Jones II, D. L. Allara, J. M. our, and P. S. Weiss, *Science* 271, 1705 (1996).

*76.* S. Datta, W. Tian, S. Hong, R. Reifenberger, J. I. Henderson, and C. P. Kubiak, *Phys. Rev. Lett.* 79, 2530 (1997).

*77.* E. Emberly and G. Kirczenow, *Phys. Rev. B* (1998).

*78.* M. Magoga and C. Joachim, *Phys. Rev. B* 56, 4722 (1997).

*79.* S. Datta and W. Tian, *Phys. Rev. B* 55, R1914 (1997).

*80.* E. Emberly and G. Kirczenow, *Phys. Rev. B* (1998).

*81.* M. Magoga and C. Joachim, *Phys. Rev. B* 56, 4722 (1997).

*82.* S. Datta and W. Tian, *Phys. Rev. B* 55, R1914 (1997).

*83.* M. P. Samanta, W. Tian, S. Datta, J. I. Henderson, and C. P. Kubiak, *Phys. Rev. B* 53, R7626 (1996)

*84.* M. Kemp, A. Roitberg, V. Mujica, T. Wanta, and M. A. Ratner, *J. Phys. Chem.* 100, 8349 (1996).

*85.* V. Mujica, M. Kemp, A. Roitberg, and M. Ratner, *J. Chem. Phys.* 104, 7296 (1996).

*86.* A. Aviram. Ed., "Proceedings of the Conference on Molecular Electronics: Science and Tecnology," New York Academy of Science, 1998.

*87.* R. Landauer, *IBM J. Res. Dev.* 1, 223 (1957); R. Landauer, *Phys. Lett.* 85A, 91 (1981).

*88.* W. E. Jones, Jr., L. H. Hermans, B. Jiang, in "Molecular and Supramolecular Photochemistry" (V. Raramurthy and K. S. Schanze. Eds.), Vol. 2, Marcel Dekker, New York, 1999.

*89.* W. Doug Bates, F. Chen, D. M. Dattelbaum, W. E. Jones, Jr., and T. J. Meyer, *J. Phys. Chem. A* 03, 5227 (1999).

*90.* M. Gallagher, P. Dougherty, P. S. Tanner, D. C. Barbini, J. Schulte, and W. E. Jones, Jr., *Inorg. Chem* 38, 2953 (1999).

*91.* L. A. Worl, W. E. Jones, Jr., G. F. Strouse, J. N. Younathan, F. Danielson, K. A. Maxwell, M. Syora, and T. J. Meyer, *Inorg. Chem.* 38, 2705 (1999).

*92.* B. Jiang, Y. Zhang, S. Sahay, S. Chatterjee, and W. E. Jones, Jr., *SPIE Proceedings* 212 (1999).

*93.* L. Hermans and W. F. Jones, Jr., *Polymer Preprints*, 409 (2000).

94.  A. G. MacDiarmid, I. D. Norris, W. E. Jones, Jr., M. A. El-Sherif, J. Yuan, B. Han, and F. K. Ko, *Polymer Preprints*, 544 (2000).

95.  J. Yuan, M. A. El-Sherif, A. G. MacDiarmid, and W. E. Jones, Jr., *Proc. SPIE-Int. Soc. Opt. Eng.: Advanced Environmental and Chemical Sensing Technologies*, 4205, 170 (2000).

96.  A. G. MacDiarmid, W. E. Jones, Jr., I. D. Norris, J. Gao, M. Llaguno, A. T. Johnson, N. J. Pinto, F. K. Ko, and H. Okusaki, *Synthetic Metals* 119, 27 (2000).

97.  Y. Zhang, C. B. Murphy, S. Chatterjee, and W. E. Jones, Jr., *Polym. Mater. Sci. Eng.* 527 (2000).

98.  L. Huang, K. J. Seward, B. P. Sullivan, W. E. Jones, Jr., J. J. Mecholsky, and W. J. Dressick, *Inorg. Chim. Acta* 310, 227 (2000).

99.  J. Yuan, M. A. El-Sherif, A. G. MacDiarmid, and W. E. Jones, Jr., *Proc. SPIE-Int. Soc. Opt. Eng.* 482 (2000).

100.  Y. Zhang, C. B. Murphy, S. Chatterjee, and W. E. Jones, Jr., *Polym. Mater. Sci. Eng.* 527 (2000).

101.  D. M. Sarno, L. J. Matienzo, and W. E. Jones, Jr., *Mater. Res. Soc. Proc.* 6281 (2001).

102.  A. G. MacDiarmid, W. E. Jones, Jr., I. D. Norris, J. Gao, M. Llaguno, A. T. Johnson, N. J. Pinto, F. K. Ko, and H. Okusaki, *Synthetic Metals* 119, 27 (2001).

103.  D. M. Sarno, L. J. Matienzo, and W. E. Jones, Jr., *Inorg. Chem.* 40, 6308 (2001).

104.  D. Hohnholz, A. G. MacDiarmid, D. M. Sarno, and W. E. Jones, Jr., *J. Chem. Soc. Chem. Comm.* 2444 (2001).

105.  S. Datta, "Electronic Transport in Mesoscopic Systems," Cambridge University Press, Cambridge, 1995.

106.  M. A. Reed, C. Zhou, C. J. Muller, T. P. Burgin, and J. M. Tour, *Science* 278, 252 (1997).

107.  E. Emberly and G. Kirezenow, unpublished (1998).

108.  G. Gundiah, F. L. Deepak, A. Govindaraj, and C. N. R. Rao, *Chem. Phys. Lett.* 381, 579 (2003).

109.  C. Zhou, M. R. Deshpande, and M. A. Reed, *Appl. Phys. Lett.* 71, 611 (1997).

110.  M. A. Reed, C. Zhou, C. J. Muller, T. P. Burgin, and J. M. Tour, *Science* 278, 252 (1997).

111.  Y. H. Tang, Y. F. Zhang, N. Wang, C. S. Lee, X. D. Han, I. Bello, and S. T. Lee, *J. Appl. Phys.* 85, 7981 (1999).

112.  S. Frank, P. Poncharal, Z. L. Wang, and W. A. de Heer, *Science* 280, 1744 (1998).

113.  S. J. Tans, M. H. Devoret, and C. Dekker, *Nature* 386, 474 (1997).

114.  R. P. Andres, J. D. Bielefeld, J. I. Henderson, D. B. Janes, V. R. Kolagunta, C. P. Kubiak, W. J. Mahoney, and R. G. Osifchin, *Science* 273, 1690 (1996).

115.  C. A. Mirkin, R. L. Letsinger, R. C. Mucic, and J. J. Storhoff, *Nature* 382, 607 (1996).

116.  L. A. Bumm, J. J. Arnold, M. T. Cygan, T. D. Dunbar, T. P. Burgin, L. Jones II, D. L. Allara, J. M. Tour, and P. S. Weiss, *Science* 271, 1705 (1996).

117.  S. Datta, W. Tian, S. Hong, R. Reifenberger, J. I. Henderson, and C. P. Kubiak, *Phys. Rev. Lett.* 79, 2530 (1997).

118.  S.-Min Choe, J.-Ah Ahn, and O. Kim, *IEEE Electron Device Lett.* 22 (2001).

119.  S. Sharma and M. K. Sunkura, *Nanotechnology* 15, 130 (2004).

120.  P. Yang, *Nature* 425, 243 (2003).

121.  E. Smalley, *Technology Research News:* http://www.trnmag.com/Stories/2003/102203/Nanowires_make_flexible_circuits_102203.html (2003).

122.  E. Smalley, *Technology Research News:* http://www.trnmag.com/Stories/2003/102203/Nanowires_boost_plastic_circuits_Brief_102203.html (2003).

123.  G. E. Moore, *Electronics* 38 (1965).

124.  M. Kanellos, Available at http://news.com.com/2100-1001-942671.html.

125.  Available at http://www.calmec.com/chiropticene.htm.

126.  S. J. Wind, J. Appenzeller, R. Martel, V. Derycke, and Ph. Avouris, *Appl. Phys. Lett.* 80, 3817 (2002).

127.  V. Derycke, R. Martel, M. Radosavljević, F. M. Ross, and Ph. Avouris, *Nano Lett.* (online), Available at http://dx.doi.org/10.1021.nl0256309 (2002).

128.  S. S. Wong, J. D. Harper, P. T. Lansbury, and C. M. Lieber, *J. Am. Chem. Soc.* 120, 603 (1998).

129.  S. S. Wong, E. Joselevich, A. T. Woolley, C. L. Cheung, and C. M. Lieber, *Nature* 394, 52 (1998).

130.  S. S. Wong, A. T. Woolley, E. Joselevich, C. L. Cheung, and C. M. Lieber, *J. Am. Chem. Soc.* 120, 8557 (1998).

131.  J. H. Hafner, C. L. Chueng, and C. M. Lieber, *Nature* 398, 761 (1999).

132.  J. H. Hafner, C. L. Cheung, and C. M. Lieber, *J. Amer. Chem. Soc.* 121, 9750 (1999).

133.  C. L. Cheung, J. H. Hafner, and C. M. Lieber, *Proc. Natl. Acad. Sci. U.S.A.* 97, 3809 (2000).

134.  C. L. Cheung, J. H. Hafner, T. W. Odom, K. Kim, and C. M. Lieber, *Appl. Phys. Lett.* 76, 3136 (2000).

135.  B. Messer, J. H. Song, and P. Yang, *J. Am. Chem. Soc.* 122, 10232 (2000).

136.  Y. Wu and P. Yang, *Appl. Phys. Lett.* 77, 43 (2000).

137.  M. Huang, A. Choudrey, and P. Yang, *Chem. Commun.* 12, 1603 (2000).

138.  B. Messer, J. H. Song, M. Huang, Y. Wu, F. Kim, and P. Yang, *Adv. Mater.* 12, 1526 (2000).

139.  J. S. Lettow, Y. J. Han, P. Schmidt-Winkel, P. Yang, D. Zhao, A. Butler, G. D. Stucky, and J. Y. Ying, *Langmuir* 16, 8291 (2000).

140.  N. Panev, A. I. Persson, N. Sköld, and L. Samuelson, *Appl. Phys. Lett.* 83, 2238 (2003).

141.  N. Gogneau, D. Jalabert, E. Monroy, T. Shibata, M. Tanaka, and B. Daudin, *J. Appl. Phys.* 94, 2254 (2003).

142.  M. Jacoby, *C&E News* 80, 7 (2002).

143.  D. A. Buzatu and J. A. Darsey, in "Clusters and Nanostructure Interfaces" (P. Jena, S. Khanna, and B. K. Rao, Eds.), p. 831, World Scientific Publishing, Singapore, 2000.

144.  K. K. Taylor, D. A. Buzatu, and J. A. Darsey, in "Computational Studies, Nanotechnology, and Solution Thermodynamics of Polymer Systems" (Dadmun et al., Eds.), p. 159, Kluwer Academic/Plenum Publishers, New York, 2000.

*145.* R. M. Metzger, B. Chen, U. Holpfner, M. V. Lakshmikantham, D. Vuillaume, T. Kawai, X. Wu H. Tachibana, T. V. Hughes, H. Sakurai, J. W. Baldwin, C. Hosch, M. P. Cava, L. Brehmer, and G. J. Ashwel. *J. Am. Chem. Soc.* 119, 10455 (1997).

*146.* M. J. Frisch, G. W. Trucks, H. B. Schlegel, P. M. W. Gill, B. G. Johnson, M. A. Robb, J. R. Cheesman, T. Keith, G. A. Peterson, J. A. Montgomery, K. Raghavachari, M. A. Al-Laham, V. G. Zakrzewski, J. V. Ortiz, J. B. Foresman, J. Cioslowski, B. B. Stefanov, A. Nanayakkara, M. Challacombe, C. Y. Peng, P. Y. Ayala, W. Chen, M. W. Wong, J. L. Andres, E. S. Replogle, R. Gomperts, R. L. Martin, D. J. Fox, J. S. Binkley D. J. Defrees, J. Baker, J. P. Stewart, M. Head-Gordon, C. Gonzales, and J. A. Pople, Gaussian 94, Revision E. 1. Gaussian, Inc., Pittsburgh, 1995.

# CHAPTER 9

# Tunneling Models for Semiconductor Device Simulation

## Andreas Gehring, Siegfried Selberherr

*Institute for Microelectronics, Technical University Vienna, Vienna, Austria*

## CONTENTS

## 1. INTRODUCTION

The increasing demand for higher computing power, smaller dimensions, and lower power consumption of electronic devices leads to a pressing need to downscale semiconductor components. This process has already led to length scales where the electrical device characteristics are dominated by quantum-mechanical effects. One of the most interesting of these effects is the quantum-mechanical tunneling of charge carriers through classically forbidden regions.

It is therefore necessary to account for tunneling effects in the design of semiconductor devices. Several models of varying complexity and accuracy can be derived to describe

the tunneling current density in semiconductor devices. The models depend on two central quantities; namely, the supply function, which describes the supply of available electrons, and the transmission coefficient, which describes the probability that an electron can tunnel through the barrier. The supply function is determined by the energy distribution of the electrons. In equilibrium, this distribution can be approximated by a Maxwellian distribution. However, the electric field in miniaturized devices is so high that non-Maxwellian models have to be considered to describe accurately the shape of the distribution function and especially the shape of the high-energy tail of the distribution.

To calculate the transmission coefficient of a dielectric layer, Schrödinger's equation must be solved. One of the most frequently used methods is the Wentzel–Kramers–Brillouin (WKB) approximation, which, however, does not reproduce transmission coefficient oscillations as observed in thin gate dielectrics. To describe accurately tunneling through dielectric stacks, it is necessary to resolve the effects of wave function interference. This can be achieved using the transfer-matrix method with either constant or linear potential segments. However, this method is prone to numerical instabilities. A more promising approach is the quantum transmitting boundary method, which allows a stable and reliable evaluation of the transmission coefficient.

Unlike what is assumed in idealized models, dielectric layers are not ideal insulators. Caused by electric stress or processing conditions, defects arise in the dielectric that give rise to trap-assisted tunneling. This results in increased tunneling current at low bias, which is referred to as SILC (stress-induced leakage current). The trap-assisted tunneling process is caused by inelastic transitions of carriers supported by the emission of phonons As this is a transient process, it is necessary to account for the creation and annihilation of traps in the dielectric based on the rate equation of the traps.

All these effects are discussed in Section 2, which treats the theory of tunneling in semiconductors. This comprises modeling of the supply function, the transmission coefficient, and trap-assisted tunneling. In Section 3, several applications are presented. The general-purpose device simulator Minimos-NT is used for the simulation of gate leakage currents in metal-oxide-semiconductor (MOS) capacitors and MOSFETs (MOS field-effect transistors). Emphasis is put on modeling of the different tunneling paths in MOS transistors and on the evaluation of alternative high-$\kappa$ dielectric materials. Furthermore, several NVM (nonvolatile memory) devices such as electrically erasable programmable read-only memory EEPROM) devices, trap-rich dielectric, or multi-barrier tunneling based devices are investigated.

## 2. THEORY OF TUNNELING

This section outlines the theory of quantum-mechanical tunneling in semiconductor devices. Different tunneling mechanisms, such as direct-, Fowler–Nordheim, and trap-assisted tunneling are covered. As a first step, the Tsu–Esaki model is derived. The derivation of the supply function and the transmission coefficient is described in detail. Tunneling from quasi-bound states and compact tunneling models is covered as well. The section continues with the description of trap-assisted tunneling and discusses some of the most frequently used models.

### 2.1. Tunneling Mechanisms

In the silicon-dielectric-silicon structure sketched in Fig. 1, a variety of tunneling processes can be identified. Considering the shape of the energy barrier alone, Fowler–Nordheim (FN) tunneling and direct tunneling can be distinguished. However, a more rigorous classification distinguishes between ECB (electrons from the conduction band), EVB (electrons from the valence band), HVB (holes from the valence band), and TAT (trap-assisted tunneling) processes. The EVB process is caused by electrons tunneling from the valence band to the conduction band. It thus creates free carriers at both sides of the dielectric, which, for MOS transistors, gives rise to increased substrate current. The TAT process can either be elastic, which means that the energy of the carrier is conserved, or inelastic, where the carrier loses energy due to the emission of phonons. Furthermore, in dielectrics with a very high defect density, hopping conduction via multiple defects may occur.

**Figure 1.** Schematic of the tunneling processes in a silicon-dielectric-silicon structure. The different tunneling processes are indicated by arrows and are described in the text. The abbreviations EED and HED denote the electron and hole energy distribution function.

## 2.2. The Tsu–Esaki Model

The processes ECB and HVB shown in Fig. 1 can be investigated considering an energy barrier as shown in Fig. 2. Two semiconductor or metal regions are separated by an energy barrier with barrier height $q\Phi_B$ measured from the Fermi energy to the conduction band edge of the insulating layer. Electrons tunnel from Electrode 1 to Electrode 2. The distribution functions at both sides of the barrier are indicated in the figure.

In the derivation, the following assumptions are made:

- Effective-mass approximation: The different masses corresponding to the band structure of the considered material are lumped into a single value for the effective mass. This is denoted by $m_{eff}$ in the electrodes and $m_{diel}$ in the dielectric layer.
- Parabolic bands: The dispersion relation in semiconductors is approximated by

$$\varepsilon = \frac{\hbar^2 k^2}{2m_{eff}} = \frac{\hbar^2 (k_x^2 + k_y^2 + k_z^2)}{2m_{eff}} \tag{1}$$

with the wave vector $\mathbf{k} = k_x \mathbf{e}_x + k_y \mathbf{e}_y + k_z \mathbf{e}_z$.

- Conservation of parallel momentum: Only transitions in the $x$-direction are considered; the parallel wave vector $\mathbf{k}_\parallel = (k_y \mathbf{e}_y + k_z \mathbf{e}_z)$ is not altered by the tunneling process.

Figure 2. Schematic of an energy barrier with two electrodes that can be used to describe the ECB or HVB processes.

The net tunneling current density from Electrode 1 to Electrode 2 can be written as the net difference between current flowing from Side 1 to Side 2 and vice versa [1, 2]

$$J = J_{1\to2} - J_{2\to1} \tag{2}$$

The current density through the two interfaces depends on the perpendicular component of the wave vector $k_x$, the transmission coefficient $TC$, the perpendicular velocity $v_x$, the density of states $g$, and the distribution function at both sides of the barrier:

$$dJ_{1\to2} = qTC(k_x)v_x g_1(k_x)f_1(\varepsilon)[1 - f_2(\varepsilon)]dk_x$$
$$dJ_{2\to1} = qTC(k_x)v_x g_2(k_x)f_2(\varepsilon)[1 - f_1(\varepsilon)]dk_x \tag{3}$$

In this expression, it is assumed that the transmission coefficient only depends on the momentum perpendicular to the interface. The density of $k_x$ states $g(k_x)$ is

$$g(k_x) = \int_0^\infty \int_0^\infty g(k_y, k_y, k_z)dk_y dk_z \tag{4}$$

where $g(k_x, k_y, k_z)$ denotes the three-dimensional density of states in the momentum space. Considering the quantized wave vector components within a cube of side length $L$

$$\Delta k_x = \frac{2\pi}{L} \quad \Delta k_y = \frac{2\pi}{L} \quad \Delta k_z = \frac{2\pi}{L} \tag{5}$$

yields for the density of states within the cube

$$g(k_x, k_y, k_z) = 2\frac{1}{\Delta k_x \Delta k_y \Delta k_z}\frac{1}{L^3} = \frac{1}{4\pi^3} \tag{6}$$

where the factor 2 stems from spin degeneracy. For the parabolic dispersion relation (1), the velocity and energy components in tunneling direction obey

$$v_x = \frac{1}{\hbar}\frac{\partial\varepsilon}{\partial k_x} = \frac{\hbar k_x}{m_{eff}} \quad \varepsilon_x = \frac{\hbar^2 k_x^2}{2m_{eff}} \quad v_x dk_x = \frac{1}{\hbar}d\varepsilon_x \tag{7}$$

Hence, expressions (3) become

$$dJ_{1\to2} = \frac{q}{4\pi^3\hbar} TC(\mathscr{E}_x)d\mathscr{E}_x \int_0^\infty \int_0^\infty f_1(\mathscr{E})[1 - f_2(\mathscr{E})]dk_y\,dk_z$$

$$dJ_{2\to1} = \frac{q}{4\pi^3\hbar} TC(\mathscr{E}_x)d\mathscr{E}_x \int_0^\infty \int_0^\infty f_2(\mathscr{E})[1 - f_1(\mathscr{E})]dk_y\,dk_z \tag{8}$$

Using polar coordinates for the parallel wave vector components

$$k_\rho = \sqrt{k_y^2 + k_z^2} \qquad k_y = k_\rho \cos(\gamma)$$

$$\gamma = \arctan\left(\frac{k_z}{k_y}\right) \qquad k_z = k_\rho \sin(\gamma) \tag{9}$$

the current density evaluates to

$$J_{1\to2} = \frac{4\pi m_{\text{eff}} q}{h^3} \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} TC(\mathscr{E}_x)d\mathscr{E}_x \int_0^\infty f_1(\mathscr{E})[1 - f_2(\mathscr{E})]d\mathscr{E}_\rho$$

$$J_{2\to1} = \frac{4\pi m_{\text{eff}} q}{h^3} \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} TC(\mathscr{E}_x)d\mathscr{E}_x \int_0^\infty f_2(\mathscr{E})[1 - f_1(\mathscr{E})]d\mathscr{E}_\rho \tag{10}$$

In these expressions, the total energy $\mathscr{E}$ has been split into a longitudinal part $\mathscr{E}_\rho$ and a transversal part $\mathscr{E}_x$

$$\mathscr{E}_\rho = \frac{\hbar^2(k_y^2 + k_z^2)}{2m_{\text{eff}}} = \frac{\hbar^2 k_\rho^2}{2m_{\text{eff}}} \qquad \mathscr{E}_x = \frac{\hbar^2 k_x^2}{2m_{\text{eff}}} \tag{11}$$

Evaluating the difference $J = J_{1\to2} - J_{2\to1}$, the net current through the interface equals

$$J = \frac{4\pi m_{\text{eff}} q}{h^3} \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} TC(\mathscr{E}_x)d\mathscr{E}_x \int_0^\infty (f_1(\mathscr{E}) - f_2(\mathscr{E}))d\mathscr{E}_\rho \tag{12}$$

This expression is usually written as an integral over the product of two independent parts, which only depend on the energy perpendicular to the interface: the transmission coefficient $TC(\mathscr{E}_x)$ and the supply function $N(\mathscr{E}_x)$:

$$J = \frac{4\pi m_{\text{eff}} q}{h^3} \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} TC(\mathscr{E}_x)N(\mathscr{E}_x)d\mathscr{E}_x \tag{13}$$

which is the expression known as Tsu–Esaki formula. This model has been proposed by Duke [3] and was used by Tsu and Esaki for the modeling of tunneling current in resonant tunneling devices [4]. The values of $\mathscr{E}_{min}$ and $\mathscr{E}_{max}$ depend on the considered tunneling process:

- Electrons tunneling from the conduction band (ECB): $\mathscr{E}_{min}$ is the highest conduction band edge of the two electrodes; $\mathscr{E}_{max}$ is the highest conduction band edge of the dielectric.
- Holes tunneling from the valence band (HVB): $\mathscr{E}_{min}$ is the absolute value of the lowest valence band edge of the electrodes; $\mathscr{E}_{max}$ is the absolute value of the lowest valence band edge of the dielectric. The sign of the integration must be changed.
- Electrons tunneling from the valence band (EVB): $\mathscr{E}_{min}$ is the lowest conduction band edge of the two electrodes; $\mathscr{E}_{max}$ the highest valence band edge of the two electrodes. It must be checked if $\mathscr{E}_{min} < \mathscr{E}_{max}$.

The next sections concentrate on the calculation of the supply function and the transmission coefficient.

## 2.3. Supply Function Modeling

The supply function describes the difference in the supply of carriers at the interfaces of the dielectric layer. Following (12), it is given as

$$N(\mathcal{E}_x) = \int_0^\infty (f_1(\mathcal{E}) - f_2(\mathcal{E})) d\mathcal{E}_p \qquad (14)$$

where $f_1$ and $f_2$ denote the energy distribution functions near the interfaces. Because the exact shape of these distributions is usually not known, approximative shapes are commonly used. Furthermore, it is assumed that the distributions are isotropic.

### 2.3.1. Fermi–Dirac Distribution

In equilibrium, the energy distribution function of electrons or holes is given by the Fermi–Dirac statistics

$$f(\mathcal{E}) = \frac{1}{1 + \exp(\frac{\mathcal{E} - \mathcal{E}_F}{k_B T})} \qquad (15)$$

which can be derived from statistical thermodynamics [5]. Separating the longitudinal and transversal energy components $\mathcal{E} = \mathcal{E}_x + \mathcal{E}_p$ and splitting the integral in (14) $N(\mathcal{E}_x) = \xi_1(\mathcal{E}_x) - \xi_2(\mathcal{E}_x)$, the values of $\xi_1$ and $\xi_2$ become

$$\xi_i = \int_0^\infty f_i(\mathcal{E}) d\mathcal{E}_p = \int_0^\infty \frac{1}{1 + \exp(\frac{\mathcal{E}_x + \mathcal{E}_p - \mathcal{E}_{F,i}}{k_B T})} d\mathcal{E}_p \qquad i = 1, 2 \qquad (16)$$

This expression can be integrated analytically using

$$\int \frac{dx}{1 + \exp(x)} = \ln\left(\frac{1}{1 + \exp(-x)}\right) + C \qquad (17)$$

so expression (16) evaluates to

$$\xi_i = k_B T \ln\left[1 + \exp\left(-\frac{\mathcal{E}_x - \mathcal{E}_{F,i}}{k_B T}\right)\right] \qquad i = 1, 2 \qquad (18)$$

and the total supply function (14) becomes

$$N(\mathcal{E}_x) = k_B T \ln\left(\frac{1 + \exp(-\frac{\mathcal{E}_x - \mathcal{E}_{F,1}}{k_B T})}{1 + \exp(-\frac{\mathcal{E}_x - \mathcal{E}_{F,2}}{k_B T})}\right) \qquad (19)$$

### 2.3.2. Maxwell–Boltzmann Distribution

For nondegenerate semiconductors, the Fermi energy is located below the conduction band edge. Therefore, $\mathcal{E}_{min} - \mathcal{E}_F \gg k_B T$ holds in expression (13), and the Fermi–Dirac distribution (15) can be approximated by a Maxwell–Boltzmann (or Maxwellian) distribution

$$f(\mathcal{E}) = \exp\left(\frac{\mathcal{E}_F - \mathcal{E}}{k_B T}\right) \qquad (20)$$

Using this expression, $\xi$ in (14) becomes

$$\xi_i = \int_0^\infty f_i(\mathcal{E}) d\mathcal{E}_p = \int_0^\infty \exp\left(-\frac{\mathcal{E}_x + \mathcal{E}_p - \mathcal{E}_{F,i}}{k_B T}\right) d\mathcal{E}_p \qquad i = 1, 2 \qquad (21)$$

which evaluates to

$$\xi_i = k_B T \exp\left(-\frac{\mathcal{E}_x - \mathcal{E}_{F,i}}{k_B T}\right) \qquad i = 1, 2 \qquad (22)$$

and yields a supply function of

$$N(\mathcal{E}_x) = k_B T \left[\exp\left(-\frac{\mathcal{E}_x - \mathcal{E}_{F,1}}{k_B T}\right) - \exp\left(-\frac{\mathcal{E}_x - \mathcal{E}_{F,2}}{k_B T}\right)\right] \qquad (23)$$

### 2.3.3. Non-Maxwellian Distributions

The Fermi–Dirac or Maxwell–Boltzmann distribution functions are frequently used to describe the distribution of carriers in equilibrium, because they are the solution of Boltzmann's transport equation for the case of zero electric field. In the channel region of a MOSFET, however, the energy distribution deviates from the ideal shape implied by expressions (15) or (20). Carriers gain energy by the electric field in the channel, and they experience scattering events. Models to describe the distribution function of such hot carriers have been studied by numerous authors [6–8]. One possibility to describe the distribution of hot carriers is to use a heated Maxwellian distribution function

$$f(\varepsilon) = A \exp\left(-\frac{\varepsilon}{k_B T_n}\right) \tag{24}$$

where $T_n$ denotes the electron temperature and $A$ is a normalization constant. The validity of this approach, however, is limited. Figure 3 shows in the left part the contour lines of the heated Maxwellian distribution function at the Si–SiO$_2$ interface in comparison to Monte Carlo results. A Monte Carlo simulator employing analytical nonparabolic bands was used for this simulation for a MOSFET with a gate length of $L_g = 180$ nm and a thickness of the gate dielectric of 1.8 nm at a bias of $V_{DS} = V_{GS} = 1V$. It is evident that the heated Maxwellian distribution (full lines) yields only poor agreement with the Monte Carlo results (dashed lines). The distribution function at two points near the middle of the channel (point A) and near the drain contact (point B) are shown in the right part of this figure. Particularly, the high-energy tail in the middle of the channel is heavily overestimated by the heated Maxwellian model. This is unsatisfactory, because a correct description of the high-energy tail is crucial for the evaluation of hot-carrier injection at the drain side used for programming and erasing of EEPROM devices.

To obtain a better prediction of hot-carrier effects, Cassi and Riccó presented an expression to account for the non-Maxwellian shape of the electron energy distribution function [6]

$$f(\varepsilon) = A \exp\left(-\frac{\chi \varepsilon^3}{E^{1.5}}\right) \tag{25}$$

with $\chi$ as fitting parameter and $E$ being the local electric field in the channel. This local-field dependence was soon questioned by other authors such as Fiegna et al. [9], who replaced the electric field with an effective field calculated from the average electron energy to model the EEPROM writing process. Hasnat et al. used a similar form for the distribution function [10]

$$f(\varepsilon) = A \exp\left(-\frac{\varepsilon^\varsigma}{\eta(k_B T_n)^r}\right) \tag{26}$$



**Figure 3.** Comparison of the heated Maxwellian distribution (full lines) with the results from a Monte Carlo simulation (dotted lines) in a turned-on 180-nm MOSFET [150]. Neighboring lines differ by a factor of 10. The distributions at point A and B are compared with a cold Maxwellian distribution in the right figure.

They obtained values of $\xi = 1.3$, $\eta = 0.265$, and $\nu = 0.75$ by fitting simulation results to measured gate currents. However, these values fail to describe the shape of the distribution function along the channel when compared to Monte Carlo results [11]. A quite generalized approach to describe the shape of the electron energy distribution (EED) has been proposed by Grasser et al.

$$f(\epsilon) = A \exp\left[-\left(\frac{\epsilon}{\epsilon_{\text{ref}}}\right)^{h}\right]$$ (27)

In this expression, the values of $\epsilon_{\text{ref}}$ and $h$ are mapped to the solution variables $T_n$ and $\beta_n$ of a six moments transport model [12]. Expression (27) has been shown to reproduce appropriately Monte Carlo results in the source and the middle region of the channel of a turned-on MOSFET. However, this model is still not able to reproduce the high-energy tail of the distribution function near the drain side of the channel because it does not account for the population of cold carriers coming from the drain. This was already visible in the right part of Fig. 3 near the drain side of the channel: The distribution consists of a cold Maxwellian, a high-energy tail, and a second cold Maxwellian at higher energies. Expression (27) cannot reproduce the low-energy Maxwellian. A distribution function accounting for the cold carrier population near the drain contact was proposed by Sonoda et al. [8], and an improved model has been suggested by Grasser et al. [11]:

$$f(\epsilon) = A\left\{\exp\left[-\left(\frac{\epsilon}{\epsilon_{\text{ref}}}\right)^{h}\right] + c\exp\left(-\frac{\epsilon}{k_B T_L}\right)\right\}$$ (28)

Here, the pool of cold carriers in the drain region is correctly modeled by an additional cold Maxwellian subpopulation. The values of $\epsilon_{\text{ref}}$, $h$, and $c$ are again derived from the solution variables of a six moments transport model [11]. Figure 4 shows again the results from Monte Carlo simulations in comparison to the analytical model. A good match between this non-Maxwellian distribution and the Monte Carlo results can be seen.

This model for the distribution function, however, requires calculation of the third even moment of the distribution function: the kurtosis $\beta_n$. As an approximation, $\beta_n$ can be calculated by an expression obtained for a bulk semiconductor where a fixed relationship between $\beta_n$, $T_n$, and the lattice temperature $T_L$ exists:

$$\beta_{\text{Bulk}}(T_n) = \frac{T_L^2}{T_n^2} + 2\frac{\tau_\beta}{\tau_\epsilon}\frac{\mu_S}{\mu_n}\left(1 - \frac{T_L}{T_n}\right)$$ (29)

In this expression, $\tau_\epsilon$, $\tau_\beta$, $\mu_n$, and $\mu_S$ are the energy relaxation time, the kurtosis relaxation time, the electron mobility, and the energy flux mobility, respectively. The value of



Figure 4. Comparison of the non-Maxwellian distribution (full lines) with the results from a Monte Carlo simulation (dotted lines) in a turned-on 180-nm MOSFET [150]. Neighboring lines differ by a factor of 10. The distributions at point A and B are compared with a cold Maxwellian distribution in the right figure

$\tau_\beta \mu_\nu / \tau \cdot \mu_n$ can be approximated by a fit to Monte Carlo data [11]. Estimating the kurtosis from (29), the distribution (27) can be used within the energy-transport or hydrodynamic model. For a parabolic band structure, the expressions

$$T_n = \frac{2}{3}\frac{\Gamma(\frac{5}{2b})}{\Gamma(\frac{3}{2b})}\frac{\mathcal{E}_{\text{rel}}}{k_B} \tag{30}$$

$$\beta_n = \frac{3}{5}\frac{\Gamma(\frac{3}{2b})\Gamma(\frac{7}{2b})}{\Gamma(\frac{5}{2b})^2} \tag{31}$$

are found [12], where $\Gamma(x)$ denotes the Gamma function

$$\Gamma(x) = \int_0^\infty \exp(-\alpha)\alpha^{x-1}\,d\alpha \tag{32}$$

Though (30) can easily be inverted to obtain $\mathcal{E}_{\text{rel}}(T_n)$, the inversion of (31) to find $b(T_n)$ at $\beta_n(b) = \beta_{\text{Bulk}}(T_n)$ cannot be given in a closed form. Instead, a fit expression

$$b(T_n) = 1 + b_0\left(1 - \frac{T_L}{T_n}\right)^{b_1} + b_2\left(1 - \frac{T_L}{T_n}\right)^{b_3} \tag{33}$$

with the parameters $b_0 = 38.82$, $b_1 = 101.11$, $b_2 = 3.40$, and $b_3 = 12.93$ can be used. Using $\mathcal{E}_{\text{rel}}(T_n)$ and $b(T_n)$, the Monte Carlo distribution can be approximated without knowledge of $\beta_n$. Figure 5 shows simulation results for a 500-nm MOSFET using the heated Maxwellian distribution (24), the non-Maxwellian distribution (28), and the non-Maxwellian distribution (27) using (30) and (33) to calculate the values of $\mathcal{E}_{\text{rel}}$ and $b$. It can be seen that the fit to the results from Monte Carlo simulations is good. However, the emerging population of cold carriers near the drain end of the channel leads to a significant error in the shape of the distribution at low energy. This is important for certain processes, whereas in the case of tunneling, the high-energy tail is more crucial.

With expression (27) for the distribution function and the assumption of a Fermi–Dirac distribution in the polysilicon gate, the supply function (14) becomes

$$N(\mathcal{E}) = A_1 \frac{\mathcal{E}_{\text{rel}}}{b}\Gamma_i\left[\frac{1}{b}, \left(\frac{\mathcal{E}}{\mathcal{E}_{\text{rel}}}\right)^b\right] - A_2 k_B T_L \ln\left[1 + \exp\left(-\frac{\mathcal{E} + \Delta\mathcal{E}_c}{k_B T_L}\right)\right] \tag{34}$$

where $\Gamma_i(\alpha, \beta)$ denotes the incomplete gamma function

$$\Gamma_i(x, y) = \int_y^\infty \exp(-\alpha)\alpha^{x-1}\,d\alpha \tag{35}$$

In (34), the explicit value of the Fermi energy was replaced by the shift of the two conduction band edges $\Delta\mathcal{E}_c$. Assuming a Maxwellian distribution in the polysilicon gate, the supply function can be further simplified to

$$N(\mathcal{E}) = A_1 \frac{\mathcal{E}_{\text{rel}}}{b}\Gamma_i\left[\frac{1}{b}, \left(\frac{\mathcal{E}}{\mathcal{E}_{\text{rel}}}\right)^b\right] - A_2 k_B T_L \exp\left(-\frac{\mathcal{E} + \Delta\mathcal{E}_c}{k_B T_L}\right) \tag{36}$$

Using the accurate shape of the distribution (28), the expressions for the supply function become

$$N(\mathcal{E}) = A_1 \frac{\mathcal{E}_{\text{rel}}}{b}\Gamma_i\left[\frac{1}{b}, \left(\frac{\mathcal{E}}{\mathcal{E}_{\text{rel}}}\right)^b\right] + A_1 c k_B T_2 \exp\left(-\frac{\mathcal{E}}{k_B T_L}\right) - A_2 k_B T_L \ln\left[1 + \exp\left(-\frac{\mathcal{E} + \Delta\mathcal{E}_c}{k_B T_L}\right)\right] \tag{37}$$

for a Fermi–Dirac distribution, and

$$N(\mathcal{E}) = A_1 \frac{\mathcal{E}_{\text{rel}}}{b}\Gamma_i\left[\frac{1}{b}, \left(\frac{\mathcal{E}}{\mathcal{E}_{\text{rel}}}\right)^b\right] + A_1 c k_B T_2 \exp\left(-\frac{\mathcal{E}}{k_B T_L}\right) - A_2 k_B T_L \exp\left(-\frac{\mathcal{E} + \Delta\mathcal{E}_c}{k_B T_1}\right) \tag{38}$$

assuming a Maxwellian distribution in the polysilicon gate.

$$f(\ell) = A \exp\left(-\frac{\ell}{k_B T_n}\right)$$

$$f(\ell) = A\left(\exp\left(-\left(\frac{\ell}{\ell_{ref}}\right)^b\right)\right.$$
$$\left. + c \exp\left(-\frac{\ell}{k_B T_L}\right)\right)$$

$\ell_{ref}$, $b$, and $c$ derived from
$n$, $T_n$, and $\beta_n$.

$$f(\ell) = A \exp\left(-\left(\frac{\ell}{\ell_{ref}}\right)^b\right)$$

$\ell_{ref}$ and $b$ derived from
$n$ and $T_n$.

Figure 5. Different expressions for the energy distribution function in a 500-nm MOSFET compared with Monte Carlo results [152].

## 2.3.4. Normalization

When implementing the analytical expressions for the distribution function and the supply function into a device simulator, it is necessary to assure consistency: the carrier concentration defined by the analytical distribution function must match the carrier concentration

from the transport model used. Therefore, the normalization prefactor $A$ has to be evaluated from

$$n = \langle 1 \rangle = \frac{1}{4\pi^3} \int f(\mathbf{k}) d^3 k \tag{39}$$

This equation can be transformed to spherical coordinates using $k = (k_x^2 + k_y^2 + k_z^2)^{1/2}$

$$n = \frac{1}{4\pi^3} \int_{-\pi}^{\pi} d\alpha \int_0^\pi \sin\theta d\theta \int_0^\infty f(k)k^2 dk \tag{40}$$

For a parabolic dispersion relation we have $dk = m_{\text{eff}}/k\hbar^2 d\ell$, which finally leads to

$$n = \int_0^\infty f(\ell) \frac{4\pi\sqrt{2m_{\text{eff}}^3}}{h^3} \sqrt{\ell} \, d\ell \tag{41}$$

where the integration is performed from the conduction band edge $\ell_c = 0$. For a Maxwellian or heated Maxwellian distribution [expressions (20) or (24)], the normalization constant evaluates to

$$A = \frac{nh^3}{4\pi(k_B T_r)^{3/2}\Gamma(\frac{3}{2})\sqrt{2m_{\text{eff}}^3}} \tag{42}$$

where $T_r$ is either the lattice temperature (for the assumption of a Maxwellian distribution) or the carrier temperature (for the assumption of a heated Maxwellian distribution). Using the non-Maxwellian distribution (27), the normalization constant evaluates to

$$A = \frac{nh^3 b}{4\pi \ell_{\text{ref}}^{3/2}\Gamma(\frac{3}{2b})\sqrt{2m_{\text{eff}}^3}} \tag{43}$$

whereas for expression (28), it is

$$A = \frac{nh^3}{4\pi\left[\frac{\ell_{\text{ref}}^{3/2}}{b}\Gamma(\frac{3}{2b}) + c(k_B T_L)^{3/2}\Gamma(\frac{3}{2})\right]\sqrt{2m_{\text{eff}}^3}} \tag{44}$$

## 2.4. The Energy Barrier

For the calculation of the transmission coefficient, it is necessary to take the shape of the energy barrier into account. Electrons tunnel from a semiconductor or metal segment through a dielectric layer to another semiconductor or metal segment. Thus, the band diagram of a MOS capacitor has to be investigated. Furthermore, the image force, which leads to a reduction of both the electron and hole energy barrier for thin dielectrics, will be described in this section.

### 2.4.1. The Metal-Oxide-Semiconductor Capacitor

Figure 6 shows the band diagram and the electrostatic potential in a metal-oxide-semiconductor structure for different voltages at the metal contact [13–15]. A central quantity is the work function, which is defined as the energy required to extract an electron from the Fermi energy to the vacuum level. The work function of the semiconductor is

$$q\Phi_S = q\chi_S + \ell_g - \ell_i + \ell_v + q\Phi_f \tag{45}$$

where $\chi_s$ denotes the electron affinity of the semiconductor. The work function difference between the work function in the metal $q\Phi_M$ and the work function in the semiconductor $q\Phi_S$ is

$$q\Phi_{MS} = q\Phi_M - q\Phi_S \tag{46}$$

| Accumulation | Flatband | No bias | Inversion |
| --- | --- | --- | --- |
| $V_G = \Phi_{MS} - \phi_{surf} - \phi_{diel} < 0$ | $V_G = \Phi_{MS} < 0$ | $V_G = 0$ | $V_G = \Phi_{MS} + \phi_{surf} + \phi_{diel} > 0$ |

**Figure 6.** Band diagram and electrostatic potential in an $n$MOS structure (negative work function difference) in accumulation, under flatband condition, without bias, and under inversion condition.

The values of $\Phi_M$ and $\chi_S$ depend on the material, as shown in Table 1 [5, 16, 17]. However, the actual value of the work function of a metal deposited on $SiO_2$ is not exactly the same as that of the metal in vacuum [17].

As long as Boltzmann statistics can be applied, the Fermi potential $\Phi_f$ depends on the doping concentration of the semiconductor in the following way:

$$p\text{-type:} \quad \Phi_f = \frac{k_B T}{q} \ln\left(\frac{N_A}{n_i}\right) > 0, \tag{47}$$

$$n\text{-type:} \quad \Phi_f = -\frac{k_B T}{q} \ln\left(\frac{N_D}{n_i}\right) < 0 \tag{48}$$

**Table 1.** Electron affinity of various semiconductors (left), work function and the radius of the Fermi sphere of various metals (right) [18, 197].

| Semiconductor | $\chi_s$ (V) | Metal | $q\Phi_M$ (eV) | $k_i$ (nm$^{-1}$) |
| --- | --- | --- | --- | --- |
| Si | 4.05 | Al | 4.28 | 17.52 |
| Ge | 4.00 | Pt | 5.65 | |
| GaAs | 4.07 | W | 4.63 | |
| GaP | 3.80 | Mg | 3.66 | 13.74 |
| GaSb | 4.06 | Ag | 4.30 | 12.04 |
| InAs | 4.90 | Au | 4.80 | 12.06 |
| InP | 4.38 | Cu | 4.25 | 13.61 |
| InSb | 4.59 | Cr | 4.50 | |

where $N_A$, $N_D$, and $n_i$ denote the acceptor, donor, and intrinsic concentrations, respectively. The concentration-independent part of (46) is labeled $\Phi'_{MS}$:

$$q\Phi'_{MS} = q\Phi_M - q\chi_S - \mathscr{E}_g + \mathscr{E}_i - \mathscr{E}_v \tag{49}$$

The voltage that has to be applied to achieve flat bands is denoted the flatband voltage. If we deviate from this voltage, a space charge region forms near the interface between the dielectric and the semiconductor. The total potential drop across this space charge region is the surface potential $\phi_{surf}$. Due to this potential, all energy levels in the conduction and valence bands are shifted by a constant amount, therefore

$$\mathscr{E}_c(x) = \mathscr{E}_{c,0} - q\phi(x)$$
$$\mathscr{E}_v(x) = \mathscr{E}_{v,0} - q\phi(x) \tag{50}$$

where $\mathscr{E}_{c,0}$ and $\mathscr{E}_{v,0}$ are the conduction and valence bands in the flatband case. Note that in the flatband case $\phi(x) = 0$ in the whole structure.

In metals, the Fermi energy is located at a higher energy level than the conduction band. The difference between the conduction band edge in the metal and the Fermi energy in the metal can be calculated considering the free-electron theory of metals, which assumes that the metal electrons are unaffected by their metallic ions. The sphere of radius $k_f$ (the Fermi wave vector) contains all occupied levels and determines the electron concentration

$$k_f = \sqrt[3]{3\pi^2 n} \tag{51}$$

The values of the metal work function and $k_f$ for various metals are summarized in the right part of Table 1 [18]. The value of $\mathscr{E}_f - \mathscr{E}_c$ can then be calculated from the carrier concentration assuming a parabolic dispersion relation.

At the semiconductor side, the height of the energy barrier is given by $q\Phi_e$ for electrons and $q\Phi_h$ for holes. Note that in the derivation of the Tsu–Esaki formula, the barrier height $q\Phi_B$, which denotes the energetic difference between the Fermi energy and the band edge in the dielectric, is used. Depending on the considered tunneling process, $q\Phi_B$ must be calculated from $q\Phi_e$ or $q\Phi_h$.

### 2.4.2. Image Force Correction

When an electron approaches a dielectric layer, it induces a positive charge on the interface that acts like an image charge within the layer. This effect leads to a reduction of the barrier height for both electrons and holes [19–21]: The conduction band bends downward and the valence band bends upward, respectively. To account for this effect, the band edge energies (50) must be modified

$$\mathscr{E}_c(x) = \mathscr{E}_{c,0} - q\phi(x) + \mathscr{E}_{image}(x)$$
$$\mathscr{E}_v(x) = \mathscr{E}_{v,0} - q\phi(x) + \mathscr{E}_{image}(x) \tag{52}$$

where the image force correction in the dielectric with thickness $t_{diel}$ can be calculated as [22]

$$\mathscr{E}_{image}(x) = -\frac{q^2}{16\pi\kappa_{diel}} \sum_j (k_1 k_2)^j \left( \frac{k_1}{|x| + j t_{diel}} + \frac{k_2}{(j+1)t_{diel} - |x|} + \frac{2k_1 k_2}{(j+1)t_{diel}} \right) \tag{53}$$

where $x = 0$ is at the interface to the dielectric. The symbols $k_1$ and $k_2$ are calculated from the dielectric permittivities in the neighboring materials

$$k_1 = \frac{\kappa_{diel} - \kappa_{si}}{\kappa_{diel} + \kappa_{si}} \quad k_2 = \frac{\kappa_{diel} - \kappa_{metal}}{\kappa_{diel} + \kappa_{metal}} = -1 \tag{54}$$

Here, $k_2$ accounts for the interface between the insulator and the metal and evaluates to $-1$.

In the semiconductor, the band edge energies are also altered

$$F_{image}(x) = -\frac{q^2}{16\pi\kappa_{si}} \sum_j (k_1 k_2)^j \left( \frac{-k_1}{|x| + jt_{diel}} + \frac{k_2}{(j+1)t_{diel} + |x|} \right) \tag{55}$$

In practice, it is sufficient to evaluate the sums in (53) and (55) up to $j = 11$ [23]. Figure 7 shows the band edge energies in an MOS structure for a dielectric layer with a thickness of 2 nm and different dielectric permittivities for an applied bias of 0 V (left) and 2 V (right). A lower dielectric permittivity leads to a stronger band bending due to the image force and therefore strongly influences the transmission coefficient.

However, there is still some uncertainty if the image force has to be considered for tunneling calculations. Though it is used in some works [23–26], others neglect it or report only minor influence on the results [27–31]. For rigorous investigations, however, it is necessary to include it in the simulations. This, however, raises the need for a high spatial resolution along the dielectric. Simple models like the analytical Wentzel–Kramers–Brillouin (WKB) formula or the Gundlach formula are not valid for this case, as described in the following sections. It may therefore be justified to account for the image force barrier lowering by correction factors.

## 2.5. Transmission Coefficient Modeling

Now that the shape of the energy barrier has been treated, the calculation of the quantum-mechanical transmission coefficient can be investigated. The transmission coefficient $TC$ is defined as the ratio of the quantum-mechanical current density

$$J(r) = \frac{i\hbar q}{2m}(\Psi \nabla \Psi^* - \Psi^* \nabla \Psi) \tag{56}$$

due to an incident wave in Region 1 and a transmitted wave in Region $N$; see Fig. 8. The assumption of plane waves in both regions

$$\Psi_1(x) = A_1 \exp(i k_1 x)$$
$$\Psi_N(x) = A_N \exp(i k_N x) \tag{57}$$

leads to the transmission coefficient

$$TC = \frac{J_N}{J_1} = \frac{k_1 m_1}{k_N m_N} \frac{|A_N|^2}{|A_1|^2} \tag{58}$$



Figure 7. Effect of the image force in an nMOS device with a dielectric thickness of 2 nm at a gate bias of 0 V (left) and 2 V (right).

**Figure 8.** Schematic of an energy barrier of a single-layer dielectric. The potential energy $W(x)$ may either be the conduction band or the valence band energy, depending on the tunneling process. The linear and constant potential approximations refer to the transfer-matrix method described in Section 2.5.3.

Note that the quantum-mechanical current density (56) is equal in Region 1 and Region $N$. Considering only the incident wave in Region 1 and the transmitted wave in Region $N$ allows definition of a transmission coefficient $TC \leq 1$. The wave function amplitudes $A_1$ and $A_N$ can be found by solving the stationary Schrödinger equation [32]

$$\left[ -\frac{\hbar^2}{2m}\nabla^2 + W(\mathbf{r}) \right]\Psi(\mathbf{r}) = \varepsilon\,\Psi(\mathbf{r})$$
(59)

where $W(\mathbf{r})$ is an external potential energy, in the barrier region. This can be achieved by various methods. The Wentzel–Kramers–Brillouin approximation can be applied either analytically for a linear barrier or numerically for arbitrary barriers. Gundlach's method can be used for a single linear energy barrier, whereas the transfer-matrix and quantum transmitting boundary methods are applicable for arbitrary-shaped barriers. The transfer-matrix method can be applied using either constant or linear potential segments as shown in Fig. 8. The different methods will be described in this section, and a brief comparison at the end summarizes their advantages and shortcomings.

### 2.5.1. The Wentzel–Kramers–Brillouin Approximation

The Wentzel–Kramers–Brillouin approximation is one of the most frequently applied approximations to solve Schrödinger's equation [31, 33, 34]. Starting from the time-independent Schrödinger equation (59), the one-dimensional case reads

$$\left[ -\frac{\hbar^2}{2m}\frac{d^2}{dx^2} + W(x) - \varepsilon \right]\Psi(x) = 0$$
(60)

If the following *Ansatz* is used for the wave function

$$\Psi(x) = R(x)\exp\left( i\frac{S(x)}{\hbar} \right)$$
(61)

the equations

$$\frac{d^2 R}{dx^2} - \frac{R}{\hbar^2}\left(\frac{dS}{dx}\right)^2 + \frac{2m[\varepsilon - W(x)]}{\hbar^2} R = 0$$
(62)

and

$$R\frac{d^2S}{dx^2} + 2\frac{dR}{dx}\frac{dS}{dx} = 0 \tag{63}$$

for the real and imaginary part of (60) can be found. Equation (63) can be solved by

$$\frac{dS}{dx} = \frac{C}{R^2} \tag{64}$$

where $C$ is a constant. With (64), Eq. (62) becomes

$$\frac{1}{R}\frac{d^2R}{dx^2} - \frac{1}{\hbar^2}\left(\frac{dS}{dx}\right)^2 + \frac{2m[\mathscr{E} - W(x)]}{\hbar^2} = 0 \tag{65}$$

With the approximation

$$\frac{1}{R}\frac{d^2R}{dx^2} \ll \frac{1}{\hbar^2}\left(\frac{dS}{dx}\right)^2 \tag{66}$$

we can write

$$S(x) \approx \int \sqrt{2m[\mathscr{E} - W(x)]}\,dx \tag{67}$$

and the wave function $\Psi(x)$ becomes

$$\Psi(x) = R(x)\exp\left(\frac{i}{\hbar}\int\sqrt{2m[\mathscr{E} - W(x)]}\,dx\right) \tag{68}$$

Now we consider an energy barrier between the classical turning points $x_1$ and $x_2$ with an incoming wave $\Psi_1$ and a transmitted wave $\Psi_2$, and $x_2 > x_1$

$$\Psi_1(x \le x_1) \sim \exp\left(\frac{i}{\hbar}\int_{-\infty}^{x_1}\sqrt{2m(\mathscr{E} - W(x'))}\,dx'\right)$$
$$\Psi_2(x \ge x_2) \sim \exp\left(\frac{i}{\hbar}\int_{-\infty}^{x_2}\sqrt{2m(\mathscr{E} - W(x'))}\,dx'\right) \tag{69}$$

The transmission probability $TC(\mathscr{E})$ is proportional to $|\Psi_2(x_2)/\Psi_1(x_1)|^2$:

$$TC = \left|\frac{\exp(\frac{i}{\hbar}\int_{-\infty}^{x_2}\sqrt{2m[\mathscr{E} - W(x')]}\,dx')}{\exp(\frac{i}{\hbar}\int_{-\infty}^{x_1}\sqrt{2m[\mathscr{E} - W(x')]}\,dx')}\right|^2 = \left|\exp\left(\frac{i}{\hbar}\int_{x_1}^{x_2}\sqrt{2m[\mathscr{E} - W(x')]}\,dx'\right)\right|^2$$
$$= \exp\left(-\frac{2}{\hbar}\int_{x_1}^{x_2}\sqrt{2m[W(x') - \mathscr{E}]}\,dx'\right) \tag{70}$$

This expression can be evaluated for arbitrary barriers. In Ref. [33], however, it is shown that the WKB approximation is only valid for

$$m\hbar\frac{dW(x)}{dx} \ll \sqrt{|2m[W(x) - \mathscr{E}]|^3} \tag{71}$$

This inequality is fulfilled for points where the variation of the energy barier is small. The WKB approximation is therefore not valid in the close vicinity of the classical turning points.

The WKB approximation is often used for tunneling simulations and has been implemented in device simulators [1, 35, 36]. For a linear energy barrier, the numerical calculation of the integral in (70) can be avoided. Still, it is necessary to distinguish between regions where direct or Fowler–Nordheim tunneling takes place. For the direct tunneling regime, $\mathscr{E} < q\Phi_0$ holds (see Fig. 8). Therefore, the transmission coefficient

$$TC(\mathscr{E}) = \exp\left(-\frac{2}{\hbar}\int_0^{x_{diel}}\sqrt{2m_{diel}(q\Phi - qE_{diel}x - \mathscr{E})}\,dx\right) \tag{72}$$

evaluates to

$$TC(\ell) = \exp\left\{-4\frac{\sqrt{2m_{diel}}}{3\hbar q E_{diel}}[(q\Phi - \ell)^{3/2} - (q\Phi_0 - \ell)^{3/2}]\right\}$$  (73)

with $E_{diel}$ being the electric field defined as $V_{diel}/t_{diel}$ and $m_{diel}$ the electron mass in the dielectric. The symbols $\Phi$ and $\Phi_0$ denote the upper and lower barrier heights as shown in Fig. 8. The value of $\Phi_0$ is calculated assuming a linear potential in the barrier

$$\Phi_0 = \Phi - E_{diel}t_{diel}$$  (74)

For the Fowler–Nordheim tunneling regime it holds $\ell > q\Phi_0$, and therefore with $x_1$ defined by $q\Phi - qE_{diel}x_1 = \ell$, the transmission coefficient

$$TC(\ell) = \exp\left(-\frac{2}{\hbar}\int_0^{x_1} \sqrt{2m_{diel}(q\Phi - qE_{diel}x - \ell)}\, dx\right)$$  (75)

evaluates to

$$TC(\ell) = \exp\left[-4\frac{\sqrt{2m_{diel}}}{3\hbar q E_{diel}}(q\Phi - \ell)^{3/2}\right]$$  (76)

The WKB tunneling coefficient is frequently multiplied by an oscillating prefactor to repro-
duce Fowler–Nordheim-induced oscillations [37–41]. However, because no wave function
interference is taken into account, the general validity of this method is questionable.

### 2.5.2. The Gundlach Method

The Gundlach method [42] provides an analytical solution of Schrödinger's equation for a
linear energy barrier. The one-dimensional time-independent Schrödinger equation in this
case reads

$$\frac{d^2}{dx^2}\Psi(x) + \frac{2m}{\hbar^2}[\ell - W(x)]\Psi(x) = 0$$  (77)

with the linear potential energy $W(x)$ between the points $x_0$ and $x_1$, $W_0 = W(x_0)$, and
$W_1 = W(x_1)$,

$$W(x) = W_0 + (x - x_0)\frac{W_1 - W_0}{x_1 - x_0}$$  (78)

for $x_0 < x < x_1$ Using the abbreviations

$$l = -\left(\frac{\hbar^2}{2m}\frac{x_1 - x_0}{W_1 - W_0}\right)^{1/3}$$

$$\lambda = -\left(\frac{2m}{\hbar^2}\right)^{1/3}\left(\frac{x_1 - x_0}{W_1 - W_0}\right)^{2/3}\left(\ell - W_0 + x_0\frac{W_1 - W_0}{x_1 - x_0}\right)$$  (79)

and $u(x) = \lambda - x/l$, expression (77) turns into

$$\frac{d^2}{dx^2}\Psi(x) - \frac{1}{l^2}u(x)\Psi(x) = 0$$  (80)

With

$$\frac{d^2}{dx^2}\Psi(x) = \frac{d}{du}\frac{du}{dx}\left\{\frac{d}{du}\frac{du}{dx}\Psi[u(x)]\right\} = \frac{1}{l^2}\frac{d^2}{du^2}\Psi[u(x)]$$  (81)

Schrödinger's equation evolves into the Airy differential equation

$$\frac{d^2}{du^2}\Psi[u(x)] - u(x)\Psi[u(x)] = 0$$  (82)

The solutions of this differential equation are the Airy functions $Ai[u(x)]$ and $Bi[u(x)]$ [43], which are depicted in Fig. 9 together with their derivatives. The wave functions consist of linear superpositions of these Airy functions

$$\Psi(x) = A Ai[u(x)] + B Bi[u(x)]$$  (83)

where the function $u(x)$ is given as

$$u(x) = -\left(\frac{2m}{\hbar^2}\right)^{1/3}\left(\frac{x_1 - x_0}{W_1 - W_0}\right)^{2/3}[\ell - W(x)]$$  (84)

Assuming a constant electron mass in the dielectric, Gundlach derives an expression for the transmission coefficient [42]

$$TC = \frac{k_n}{k_1}\frac{4}{\pi^2}\left[\left(\frac{z'}{k_1}A + \frac{k_n}{z'}B\right)^2 + \left(\frac{k_n}{k_1}C + D\right)^2\right]^{-1}$$  (85)

where the abbreviations

$$A = Ai'(z_0)Bi'(z_s) - Ai'(z_s)Bi'(z_0)$$  (86)

$$B = Ai(z_0)Bi(z_s) - Ai(z_s)Bi(z_0)$$  (87)

$$C = Ai(z_s)Bi'(z_0) - Ai'(z_0)Bi(z_s)$$  (88)

$$D = Ai(z_0)Bi'(z_s) - Ai'(z_s)Bi(z_0)$$  (89)

have been used, and the symbols $z_0$, $z_s$, and $z'$ are given by

$$z_0 = (q\Phi_0 - \ell)\left(\frac{at_{diel}}{2q(\Phi - \Phi_0)}\right)^{2/3} \quad z_s = (q\Phi - \ell)\left(\frac{at_{diel}}{2q(\Phi - \Phi_0)}\right)^{2/3}$$  (90)

and

$$z' = -\left(\frac{a^2}{4}\frac{q\Phi - q\Phi_0}{t_{diel}}\right)^{1/3} \quad a = \frac{2}{\hbar}\sqrt{2m_{diel}}$$  (91)

The symbols $q\Phi$ and $q\Phi_0$ denote the two edges of the energy barrier as shown in Fig. 8. The Gundlach method is frequently used in the literature [25, 44] and implemented in device simulators. Numerical problems may occur for flat barriers ($\Phi \approx \Phi_0$) due to the exponential increase of the Airy functions $Bi$ and $Bi'$ for positive arguments. In practical implementations, the values of $z_0$ and $z_s$ have been bounded to values below $\approx 200$ to avoid floating point overflow.



Figure 9. The Airy functions $Ai$ and $Bi$ and their derivatives.

## 2.5.3. The Transfer-Matrix Method

The use of the transfer-matrix (TM) method for the calculation of the transmission coefficient of energy barriers is based on the work of Tsu and Esaki on electron tunneling through one-dimensional super lattices [4]. It has been used by numerous authors to describe tunneling processes in semiconductor devices [45–49]. The basic principle of the transfer-matrix method is the approximation of an arbitrary-shaped energy barrier by a series of piece-wise constant or piece-wise linear functions. Because the wave function in such barriers can easily be calculated, the total transfer matrix can be derived by a number of subsequent matrix computations. From the transfer matrix, the transmission coefficient can easily be derived.

### 2.5.3.1. Piecewise-Constant Potential

If an arbitrary potential barrier is segmented into N regions with constant potentials (see Fig. 8), the wave function in each region can be written as the sum of an incident and a reflected wave [50] $\Psi_j(x) = A_j \exp(ik_j x) + B_j \exp(-ik_j x)$ with the wave number $k_j = \sqrt{2m_j(\mathcal{E} - W_j)}/\hbar$. The wave amplitudes $A_j$, $B_j$, the carrier mass $m_j$, and the potential energy $W_j$ are assumed constant for each region $j$. With the interface conditions for energy and momentum conservation

$$\Psi_j(x^-) = \Psi_{j+1}(x^+) \tag{92}$$

$$\frac{1}{m_j}\frac{d\Psi_j(x^-)}{dx} = \frac{1}{m_{j+1}}\frac{d\Psi_{j+1}(x^+)}{dx} \tag{93}$$

the outgoing wave of a layer relates to the incident wave by a complex transfer matrix:

$$\begin{pmatrix} A_j \\ B_j \end{pmatrix} = \underline{T}_j \begin{pmatrix} A_{j-1} \\ B_{j-1} \end{pmatrix} \quad 2 \leq j \leq N \tag{94}$$

The transfer matrices are of the form

$$\underline{T}_j = \frac{1}{2}\begin{pmatrix} \left(1 + \frac{k_{j-1}}{k_j}\right)\gamma^{-k_j} & \left(1 - \frac{k_{j-1}}{k_j}\right)\gamma^{-k_j} \\ \left(1 - \frac{k_{j-1}}{k_j}\right)\gamma^{k_j} & \left(1 + \frac{k_{j-1}}{k_j}\right)\gamma^{k_j} \end{pmatrix}\begin{pmatrix} \gamma^{k_{j-1}} & 0 \\ 0 & \gamma^{-k_{j-1}} \end{pmatrix} \quad 2 \leq j \leq N \tag{95}$$

with the phase factor $\gamma = \exp[i\Delta(j-2)]$. The transmitted wave in Region N can then be calculated from the incident wave by subsequent multiplication of transfer matrices:

$$\begin{pmatrix} A_N \\ B_N \end{pmatrix} = \prod_{j=2\ldots N} \underline{T}_j \begin{pmatrix} A_1 \\ B_1 \end{pmatrix} \tag{96}$$

If it is assumed that there is no reflected wave in Region N and the amplitude of the incident wave is unity, (96) simplifies to

$$\begin{pmatrix} A_N \\ 0 \end{pmatrix} = \begin{pmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{pmatrix}\begin{pmatrix} 1 \\ B_1 \end{pmatrix} \tag{97}$$

and the transmission coefficient can be calculated from (58). The transfer-matrix method based on constant potential segments has the obvious shortcoming that, for practical barriers, the accuracy of the resulting matrix strongly depends on the chosen resolution. A more rigorous approach is to use linear potential segments.

**2.5.3.2. Piecewise-Linear Potential** A general barrier may consist of several segments with linear potential sandwiched between contact segments where the potential is constant, as depicted in Fig. 10. The wave functions within these four regions can be written as [confer (83) and (84) for a linear potential]

$$\Psi_1(x) = A_1 \exp(ik_1 x) + B_1 \exp(-ik_1 x) \tag{98}$$

$$\Psi_2(x) = A_2 \mathrm{Ai}[u_2(x)] + B_2 \mathrm{Bi}[u_2(x)] \tag{99}$$

$$\Psi_3(x) = A_3 \mathrm{Ai}[u_3(x)] + B_3 \mathrm{Bi}[u_3(x)] \tag{100}$$

$$\Psi_4(x) = A_4 \exp(ik_4 x) + B_4 \exp(-ik_4 x) \tag{101}$$

with $u(x)$ from (84) and the $x$-independent derivative

$$u' = \frac{du(x)}{dx} = -\left(\frac{2m}{\hbar^2}\right)^{1/3}\left(\frac{W_2 - W_1}{x_2 - x_1}\right)^{1/3} \tag{102}$$

The conditions for continuity of the wave functions and their derivatives yield the following equation system, where abbreviations for the left and right value of $u(x)$ in a layer $\overleftarrow{u}_j = u_j(l_{j-2})$, $\overrightarrow{u}_j = u_j(l_{j-1})$, and their derivatives $u'_j$ for $2 \leq j \leq N - 1$ have been used.

$$
\begin{aligned}
A_1 \exp(ik_1 l_0) + B_1 \exp(-ik_1 l_0) &= A_2 \mathrm{Ai}(\overrightarrow{u}_2) + B_2 \mathrm{Bi}(\overrightarrow{u}_2) \\[4pt]
A_1 ik_1 \exp(ik_1 l_0) - B_1 ik_1 \exp(-ik_1 l_0) &= A_2 \mathrm{Ai}'(\overrightarrow{u}_2)u'_2 + B_2 \mathrm{Bi}'(\overrightarrow{u}_2)u'_2 \\[4pt]
A_2 \mathrm{Ai}(\overleftarrow{u}_2) + B_2 \mathrm{Bi}(\overleftarrow{u}_2) &= A_3 \mathrm{Ai}(\overrightarrow{u}_3) + B_3 \mathrm{Bi}(\overrightarrow{u}_3) \\[4pt]
A_2 \mathrm{Ai}'(\overleftarrow{u}_2)u'_2 + B_2 \mathrm{Bi}'(\overleftarrow{u}_2)u'_2 &= A_3 \mathrm{Ai}'(\overrightarrow{u}_3)u'_3 + B_3 \mathrm{Bi}'(\overrightarrow{u}_3)u'_3 \\[4pt]
A_3 \mathrm{Ai}(\overleftarrow{u}_3) + B_3 \mathrm{Bi}(\overleftarrow{u}_3) &= A_4 \exp(il_2 k_4) + B_4 \exp(-il_2 k_4) \\[4pt]
A_3 \mathrm{Ai}'(\overleftarrow{u}_3)u'_3 + B_3 \mathrm{Bi}'(\overleftarrow{u}_3)u'_3 &= A_4 ik_4 \exp(il_2 k_4) - B_4 ik_4 \exp(-il_2 k_4)
\end{aligned}
\tag{103}
$$

The transfer matrices between adjacent layers are again calculated from (94). Using the first two equations of (103) and the Wronskian [43]

$$\mathrm{Wr}\{\mathrm{Ai}(z), \mathrm{Bi}(z)\} = \mathrm{Ai}(z)\mathrm{Bi}'(z) - \mathrm{Ai}'(z)\mathrm{Bi}(z) = \pi^{-1} \tag{104}$$

the matrix $\underline{T}_1$ can be simplified to

$$
\underline{T}_1 = \pi
\begin{pmatrix}
\exp(ik_1 l_0)\left(\mathrm{Bi}'(\overrightarrow{u}_2) - \mathrm{Bi}(\overrightarrow{u}_2)\dfrac{ik_1}{u'_2}\right) & \exp(-ik_1 l_0)\left(\mathrm{Bi}'(\overrightarrow{u}_2) + \mathrm{Bi}(\overrightarrow{u}_2)\dfrac{ik_1}{u'_2}\right) \\[14pt]
\exp(ik_1 l_0)\left(-\mathrm{Ai}'(\overrightarrow{u}_2) + \mathrm{Ai}(\overrightarrow{u}_2)\dfrac{ik_1}{u'_2}\right) & \exp(-ik_1 l_0)\left(-\mathrm{Ai}'(\overrightarrow{u}_2) - \mathrm{Ai}(\overrightarrow{u}_2)\dfrac{ik_1}{u'_2}\right)
\end{pmatrix}
$$

Using the next two lines of (103) yields

$$T_2 = \pi \begin{pmatrix} \mathrm{Ai}(\bar{u}_2)\mathrm{Bi}'(\bar{u}_3) - \dfrac{u_2'}{u_3'}\mathrm{Bi}(\bar{u}_3)\mathrm{Ai}'(\bar{u}_2) & \mathrm{Bi}(\bar{u}_2)\mathrm{Bi}'(\bar{u}_3) - \dfrac{u_2'}{u_3'}\mathrm{Bi}(\bar{u}_3)\mathrm{Bi}'(\bar{u}_2) \\[2ex] \dfrac{u_2'}{u_3'}\mathrm{Ai}(\bar{u}_3)\mathrm{Ai}'(\bar{u}_2) - \mathrm{Ai}(\bar{u}_2)\mathrm{Ai}'(\bar{u}_3) & \dfrac{u_2'}{u_3'}\mathrm{Ai}(\bar{u}_3)\mathrm{Bi}'(\bar{u}_2) - \mathrm{Bi}(\bar{u}_2)\mathrm{Ai}'(\bar{u}_3) \end{pmatrix}$$

and the last two equations yield with the phase factor $\gamma = \exp(il_2 k_4)$

$$T_3 = \frac{1}{2} \begin{pmatrix} \mathrm{Ai}(\bar{u}_3)\gamma^{-1} + \dfrac{u_3'}{ik_4}\mathrm{Ai}'(\bar{u}_3)\gamma^{-1} & \mathrm{Bi}(\bar{u}_3)\gamma^{-1} + \dfrac{u_3'}{ik_4}\mathrm{Bi}'(\bar{u}_3)\gamma^{-1} \\[2ex] \mathrm{Ai}(\bar{u}_3)\gamma - \dfrac{u_3'}{ik_4}\mathrm{Ai}'(\bar{u}_3)\gamma & \mathrm{Bi}(\bar{u}_3)\gamma - \dfrac{u_3'}{ik_4}\mathrm{Bi}'(\bar{u}_3)\gamma \end{pmatrix}$$

Though being more accurate than the constant potential approach, this method is computationally more expensive. This drawback, however, is offset by the fact that a lower resolution and thus fewer matrix multiplications are necessary to resolve an energy barrier consisting of linear potential segments.

Simulations using the transfer-matrix method have been reported by several authors [51–54]. Others compared the constant and linear potential approaches and found the constant potential method more feasible for device simulation [55]. The main advantage of the linear-potential transfer-matrix method is that for linear potential segments, the accuracy does not depend on the resolution as it does for the constant-potential transfer-matrix method. However, the evaluation of the Airy functions must be carefully implemented to avoid overflow.

Although the transfer-matrix method for constant or linear potential segments is intuitively easy to understand and implement, the main shortcoming of the method is that it becomes numerically instable for thick barriers. This has been observed by several authors [55–59]. The reason for the numerical problems is that during the matrix multiplications, exponentially growing and decaying states have to be multiplied, leading to rounding errors that eventually exceed the amplitude of the wave function itself for thick barriers.

These problems have been overcome by a further segmentation of the barrier into slices with more accurate transfer matrices [56], the use of scattering matrices instead of transfer matrices [57], iterative methods [58], or by simply setting the transfer matrix entries to zero if the decay factor $\sum k_j x_j$ exceeds a certain value of about 20 [55]. In the next section, a method will be presented that avoids this problem and allows a fast and reliable transmission coefficient estimation.

### 2.5.4. The Quantum Transmitting Boundary Method

An alternative method to solve the Schrödinger equation has been proposed by Frensley and Einspruch [60], which is based on the tight-binding quantum transmitting boundary method (QTBM) introduced by Lent [61]. It has been used to simulate electron transport in resonant tunneling diodes [59]. The method is based on the finite-difference approximation of the stationary one-dimensional Schrödinger equation (77) on an equidistant grid with an effective mass $m_j$ and a grid spacing $\Delta$

$$\underline{H}\Psi_j = -s_{j-1}\Psi_{j-1} + d_j\Psi_j - s_{j+1}\Psi_{j+1} = \mathscr{E}\Psi_j \tag{105}$$

where $s_j = \hbar^2/(2m_j\Delta^2)$ and $d_j = \hbar^2/(m_j\Delta^2) + W_j$. For the evaluation of the transmission coefficient, it is necessary to assume open boundary conditions. They are introduced by writing the wave functions at the boundaries of the simulation domain as

$$\Psi_1 = a_1 + b_1 \tag{106}$$

$$\Psi_N = a_N + b_N \tag{107}$$

and relate them to the wave functions outside of the simulation domain by

$$\Psi_0 = a_1 \exp(-ik_1\Delta) + b_1 \exp(ik_1\Delta),$$    (108)

$$\Psi_{N+1} = a_N \exp(-ik_N\Delta) + b_N \exp(ik_N\Delta)$$    (109)

This introduces four unknowns and two equations into the system. Setting

$$a_1 = \zeta_1\Psi_0 + \xi_1\Psi_1$$    (110)

$$a_N = \zeta_N\Psi_{N+1} + \xi_N\Psi_N$$    (111)

eliminates the unknown values of $b_1$ and $b_N$ and gives a linear system for the $N + 2$ complex values $\Psi_j$

$$\begin{pmatrix} \zeta_1 & \xi_1 & & & & \\ -s_1 & d_1 - \prime & -s_2 & & & \\ & -s_2 & d_1 - \prime & -s_3 & & \\ & & & \cdots & & \\ & & & -s_N & d_N - \prime & -s_N \\ & & & & \xi_N & \zeta_N \end{pmatrix} \begin{pmatrix} \Psi_0 \\ \Psi_1 \\ \Psi_2 \\ \cdots \\ \Psi_N \\ \Psi_{N+1} \end{pmatrix} = \begin{pmatrix} a_1 \\ 0 \\ 0 \\ \cdots \\ 0 \\ a_N \end{pmatrix}$$    (112)

Setting $a_1 = 1$ and $a_N = 0$ yields the values of the wave function in the whole simulation domain for an incident wave from the left side as in the transfer-matrix method. The method is easy to implement, fast, and more robust than the transfer-matrix method. A further advantage of this method is its suitability for two- and three-dimensional problems. It directly yields the values of the open-boundary wave functions, which can be used to estimate the carrier concentration in the dielectric. Note that the QTBM is closely linked with the nonequilibrium Green's function formalism (NEGF): The matrix in expression (112) is the inverse of the retarded Green's function for an open system without scattering. However, the values of $\zeta$ and $\xi$ are complex, so the matrix admits complex eigenvalues and complex solving routines are necessary.

## 2.5.5. Comparison

Figure 11 shows the transmission coefficient for the described methods for a triangular energy barrier (left) and a two-step nonlinear energy barrier (right). The inset shows the energy barrier and the values of $|\Psi|^2$ for an energy of 2.8 eV on a logarithmic scale. The dotted lines refer to the constant-potential transfer-matrix method. In the left figure, the numerical instability of the transfer-matrix method leads to an increasing transmission coefficient for energies below 1 eV. These numerical problems occur for both the constant-potential and the linear-potential approaches.

The Gundlach and analytical WKB methods deliver similar results for the triangular barrier. For the stacked dielectric shown in the right figure, the analytical WKB and Gundlach methods cannot be used. The numerical WKB, transfer-matrix, and QTB methods deliver similar results; however, the WKB method does not resolve oscillations in the transmission coefficient.

It can be concluded that for a single-layer dielectric, the analytical WKB method yields reasonable accuracy as compared to the other, computationally more expensive methods. For stacked dielectrics, however, only the numerical WKB, transfer-matrix, or QTB methods can be used in the first place. Because transfer-matrix-based methods exhibit problems regarding numerical stability, only the QTBM and the numerical WKB methods remain. Because the numerical WKB method also needs a numerical integration, its advantage in terms of computational effort is not high enough to rule out the QTBM. Furthermore, if resonance effects—such as in dielectrics with quantum wells—have to be taken into account, the QTBM remains as the method of choice for a reliable transmission coefficient estimation.

**Figure 11.** The transmission coefficient using different methods for a dielectric consisting of a single layer (left) and for a dielectric consisting of two layers (right) [140]. The shape of the energy barrier and the wave function at 2.8 eV is shown in the inset.

## 2.6. Bound and Quasi-Bound States

Up to now, it has been assumed that all energetic states in the substrate contribute to the tunneling current. However, the high doping and the high electric field in the channel leads to a quantum-mechanical quantization of carriers [62,63]. If it is assumed that the wave function does not penetrate into the gate, discrete energy levels can be identified. However, it cannot be assumed that electrons tunnel from these energies, as for the derivation of the levels, it was assumed that there is no wave function penetration into the dielectric. This leads to the *paradox* that was addressed by Magnus and Schoenmaker [64]: How can a bound state, which has vanishing current density, lead to tunneling current?

The answer is that it cannot. Taking a closer look at the conduction band edge of a MOSFET in inversion reveals that, depending on the boundary conditions, different types of quantized energy levels must be distinguished [65]: Bound states are formed at energies for which the wave function decays to zero at both sides of the dielectric. Quasi-bound states (QBS) have closed boundary conditions at one side and open boundary conditions at the other side. Free states, finally, are states that do not decay at any side of the dielectric layer. This is shown schematically in Fig. 12. The total tunnel current density therefore consists of



**Figure 12.** Free, bound, and quasi-bound states in a typical MOS inversion layer.

current from the QBS and from the free states:

$$J = q \sum_i \frac{n_\nu(\mathscr{E}_i)}{\tau_q(\mathscr{E}_i)} + \frac{4\pi m_{\text{eff}} q}{h^3} \int_{\mathscr{E}_{\text{min}}}^{\mathscr{E}_{\text{max}}} TC(\mathscr{E}) N(\mathscr{E}) d\mathscr{E} \tag{113}$$

where the symbol $n_\nu(\mathscr{E}_i)$ denotes the two-dimensional carrier concentration [66]

$$n_\nu = g_\nu \frac{m k_B T}{\pi \hbar^2} \ln\left[1 + \exp\left(\frac{\mathscr{E}_F - \mathscr{E}_i}{k_B T}\right)\right] \tag{114}$$

the symbol $g_\nu$ is the valley degeneracy, and $\tau_q$ is the lifetime of the quasi-bound state $\mathscr{E}_i$. The lifetime is based on Gamow's theory of nuclear decay [67] and denotes the time constant with which an electron leaks through the energy barrier. Because bound and quasi-bound states are closely related, the computation of bound states will be described first.

### 2.6.1. Eigenvalues of a Triangular Energy Well

To first order the conduction band edge in a MOSFET inversion layer can be approximated by a linear potential (this is actually done by various authors, see Refs. [68–71]). The solution of Schrödinger's equation for a linear potential has been derived in Section 2.5.2 and consists of a linear superposition of Airy functions. If the triangular energy well is defined as

$$W(x) = W_0 + \frac{W_1 - W_0}{x_1 - x_0} x \tag{115}$$

and no wave function penetration for $x <= x_0$ is taken into account, the wave function for $x > 0$ can be written as [62]

$$\Psi(x) = A\text{Ai}[u(x)] \tag{116}$$

$$\Psi(x_0) = A\text{Ai}[u(x_0)] = 0 \tag{117}$$

Therefore, $u(x_0)$ must equal one of the zeros of the Airy function $z_i$:

$$u(x_0) = z_i < 0 \tag{118}$$

With $u(x)$ from expression (84), the energy eigenvalues are found as

$$\mathscr{E}_i = W_0 - z_i \left(\frac{\hbar^2}{2m}\right)^{1/3} \left(\frac{W_1 - W_0}{x_1 - x_0}\right)^{2/3} \tag{119}$$

The first five zeros of the Airy function are $-2.34$, $-4.09$, $-5.52$, $-6.79$, and $-7.94$. These values are often used to approximate the quantized carrier concentration in the channel of MOS devices.

For the assumption of a triangular energy well, the wave function is approximately given as (see Section 2.6.1)

$$\Psi(x) = A\text{Ai}[u(x)] \tag{120}$$

with

$$u(x) = -\left(\frac{2m}{\hbar^2}\right)^{1/3} \left(\frac{x_1 - x_0}{W_1 - W_0}\right)^{2/3} (\mathscr{E} - W(x)) \tag{121}$$

The square of the wave function is a probability, therefore the normalization can be written as [62]

$$\int_0^\infty |\Psi[u(x)]|^2 \, dx = 1 \tag{122}$$

$$\int_0^\infty |A \mathrm{Ai}[u(x)]|^2 \, dx = 1 \tag{123}$$

$$\int_0^\infty A \mathrm{i}^2 \left\{ -\left(\frac{2m}{\hbar^2}\right)^{1/3} \left(\frac{x_1 - x_0}{W_1 - W_0}\right)^{2/3} [f_1 - W(x)] \right\} dx = \frac{1}{A^2} \tag{124}$$

where an infinite barrier is assumed for $x < 0$. With $x_0 = 0$, $W_0 = 0$, and the electric field

$$E = \frac{W_1}{q x_1} \tag{125}$$

the integral becomes

$$\int_0^\infty A \mathrm{i}^2 \left[ \left(\frac{2mqE}{\hbar^2}\right)^{1/3} \left(x - \frac{f_1}{qE}\right)\right] dx = \frac{1}{A^2}. \tag{126}$$

Substituting

$$\lambda(x) = \left(\frac{2mqE}{\hbar^2}\right)^{1/3} \left(x - \frac{f_1}{qE}\right) \tag{127}$$

$$d\lambda(x) = \left(\frac{2mqE}{\hbar^2}\right)^{1/3} dx \tag{128}$$

yields

$$\left(\frac{\hbar^2}{2mqE}\right)^{1/3} \int_{\lambda(0)}^\infty A \mathrm{i}^2 [\lambda(x)] \, d\lambda(x) = \frac{1}{A^2} \tag{129}$$

Using the expression [63]

$$\int_z^\infty A \mathrm{i}^2(x)\, dx = -z A \mathrm{i}^2(z) + A \mathrm{i}'^2(z) \tag{130}$$

and $\lambda(0) = \lambda_0$, the normalization constant becomes

$$A = \left(\frac{(\frac{2mqE}{\hbar^2})^{1/3}}{A \mathrm{i}'^2(\lambda_0) - \lambda_0 A \mathrm{i}^2(\lambda_0)}\right)^{1/2} \tag{131}$$

This method can be used to get an estimate of the first few eigenvalues of the system or to find initial values for the calculation of the eigenvalues described in the next section.

### 2.6.2. Eigenvalues of Arbitrary Energy Wells

To calculate the eigenvalues of an arbitrary energy well, it is necessary to solve Schrödinger's equation. This can be done using the method of finite differences. It is based on a discretization of the Hamiltonian on a spatial grid and given by (105), which is repeated here for convenience

$$\underline{H} \Psi_j = -s_j \Psi_{j-1} + d_j \Psi_j - s_{j+1} \Psi_{j+1} = \mathscr{E} \Psi_j$$

Though in Section 2.5.4, a constant value of the electron mass in the simulated region was used, a discretization that allows for a position-dependent carrier mass reads

$$d_j = \frac{\hbar^2}{4\Delta^2} \left(\frac{1}{m_{j-1}} + \frac{2}{m_j} + \frac{1}{m_{j+1}}\right) + W_j \tag{132}$$

and

$$x_i = \frac{\hbar^2}{4\Delta^2}\left(\frac{1}{m_{i-1}} + \frac{1}{m_j}\right)$$ (133)

The system Hamiltonian is tridiagonal and, for a six-point example, can be written similar to (112) but without the entries for $\zeta$ and $\xi$:

$$\begin{pmatrix} d_1 & -s_2 & & \\ -s_2 & d_2 & -s_3 & \\ & -s_3 & d_3 & -s_4 \\ & & -s_4 & d_4 \end{pmatrix}\begin{pmatrix} \Psi_1 \\ \Psi_2 \\ \Psi_3 \\ \Psi_4 \end{pmatrix} = \mathcal{E}\begin{pmatrix} \Psi_1 \\ \Psi_2 \\ \Psi_3 \\ \Psi_4 \end{pmatrix}$$ (134)

The values $\Psi_0$ and $\Psi_5$ must be 0 in this case; that is, closed boundary conditions are assumed. The system Hamiltonian is real and symmetric, therefore all eigenvalues are real. Though this matrix equation looks similar to (112), there are important differences. Here it is necessary to solve the eigenvalue equation to get a value for $\mathcal{E}_i$ and $\Psi_i$. In (112), any value of $\mathcal{E}$ leads to a valid solution for $\Psi_i$, and the solution is obtained by solving a complex equation system.

## 2.6.3. The Lifetime of Quasi-Bound States

The tunneling current from quasi-bound states in (113) depends on their quantum-mechanical lifetime $\tau_q$: In contrast to electrons in bound states, which have an infinite lifetime, electrons in quasi-bound states have a nonzero probability to tunnel through the energy barrier, thus their lifetime is finite [72–74]. This can be seen if the time evolution of the states is considered [75]

$$\Psi(t) = \Psi_0 \exp\left(-i\frac{\mathcal{E}_i}{\hbar}t\right)$$ (135)

where $\Psi_0$ is the initial wave function and the complex eigenenergy is

$$\mathcal{E}_i = \mathcal{E}_{re} - i\mathcal{E}_{im}$$ (136)

The time-dependent probability becomes

$$P(t) = \Psi^*(t)\Psi(t) = \Psi_0^2 \exp\left(-\frac{2\mathcal{E}_{im}}{\hbar}t\right) = \Psi_0^2 \exp\left(-\frac{t}{\tau_q}\right)$$ (137)

Thus, the imaginary component of the eigenenergy $\mathcal{E}$ is related to the decay time constant by

$$\tau_q = \frac{\hbar}{2\mathcal{E}_{im}}$$ (138)

The QBS are frequently used for tunneling current calculations [76–81]. Three methods are established to compute the lifetime of a quasi-bound state in MOS inversion layers: computing the full width half-maximum (FWHM) of the reflection coefficient resonances, using the quasi-classical formula based on the Wentzel–Kramers–Brillouin method, or from the complex eigenvalues of the non-Hermitian Hamiltonian. These methods will be described in the following.

**2.6.3.1. The Reflection Coefficient Resonances**   A quasi-bound state forms if one of the system boundary conditions is open ($\neq 0$) and the other one is closed ($=0$). The carrier wave function is reflected at the interface: there is no transmitted wave. Using the transfer-matrix method described in Section 2.5.3, the system can be described by

$$\begin{pmatrix} A_N \\ B_N \end{pmatrix} = \begin{pmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{pmatrix} \begin{pmatrix} A_1 \\ B_1 \end{pmatrix} \tag{139}$$

where the wave functions are plane waves

$$\Psi_j(x) = A_j \exp(\imath k_j x) + B_j \exp(-\imath k_j x) \tag{140}$$

However, no transmission coefficient can be defined for a quasi-bound state: The transmitted wave amplitude $A_N$ must vanish to fulfill the assumption of closed boundary conditions. Instead, a reflection coefficient can be defined, which is

$$RC(\epsilon) = \frac{B_1}{A_1} = -\frac{T_{21}}{T_{22}} \tag{141}$$

For free states, which is the kind of application investigated in Section 2.5.3, the transfer matrix is Hermitian:

$$T_{11} = T_{22}^* \tag{142}$$

$$T_{12} = T_{21}^* \tag{143}$$

It is shown in Ref. [72] that for a quasi-bound state, the transfer matrix is not Hermitian and its elements obey

$$T_{11} = T_{12}^*, \tag{144}$$

$$T_{21} = T_{22}^* \tag{145}$$

Therefore, the reflection coefficient $RC(\epsilon)$ can be written as

$$RC(\epsilon) = \exp[\imath \Theta(\epsilon)] \tag{146}$$

The phase $\Theta(\epsilon)$ varies only weakly at energies away from the resonance energy of the QBS, whereas near the QBS the phase changes strongly. Near the complex energy levels $\epsilon_j$, the derivative of the phase factor $\Theta(\epsilon)$ follows a Lorentzian distribution

$$\frac{d\Theta}{d\epsilon} = \frac{2\epsilon_j}{(\epsilon - \epsilon_{re})^2 + \epsilon_{im}^2} \tag{147}$$

where $2\epsilon_{im}$ is the FWHM value of $d\Theta/d\epsilon$. Thus, by calculating the phase of the reflection coefficient as a function of energy, the lifetimes can be determined. This method has been studied intensely by Cassan et al. [66, 82]. They reported numerical difficulties in the calculation of the value of $d\Theta/d\epsilon$, which is prone to numerical noise. Similar problems have been reported by other groups [83].

An alternative approach has been presented by Clerc et al., who noted that the lifetimes can also be extracted directly from the transfer matrix [49]. For a free state, $B_N = 0$ in (139), and the transmission coefficient becomes

$$TC = \left| \frac{A_N}{A_1} \right|^2 = \frac{1}{|T_{11}|^2} \tag{148}$$

For a quasi-bound state, $A_N = 0$. Therefore,

$$A_1 = T_{12} B_N \tag{149}$$

but, because $T_{11} = T_{12}^*$, the value of $|T_{11}|^{-2}$ may be evaluated as well—even if it cannot be interpreted as a transmission coefficient. The lifetime of the QBS is proportional o the resonance peak of the Lorentzian around the real component of the eigenenergy $\mathcal{E}_{re}$

$$\frac{1}{|T_{11}|^2} \propto \frac{1}{(\mathcal{E} - \mathcal{E}_{re})^2 + \frac{\hbar^2}{4\tau_q^2}}$$  (150)

but no derivative must be calculated this time. As an example of this method, the left part of Fig. 13 shows the shape of the conduction band edge of a MOS structure in the substrate, dielectric, and polysilicon gate. In the substrate, a triangular quantum well forms. Considering closed boundaries, eigenvalues and wave functions can be calculated. The corresponding wave functions are shown in the figure, where closed boundary conditions have been used at the boundaries of the simulation domain. Note the wave function penetration into the classically forbidden region of the dielectric layer. The eigenvalues of the quasi-bound states are located at 0.27, 0.47, 0.63, 0.76, 0.86, and 0.95 eV. The same information can be found when the value of $|T_{11}|^{-2}$ is investigated, as shown in right part of Fig. 13: Every quasi-bound state in the inversion layer manifests as a peak in the value of $|T_{11}|^{-2}$. The width of each peak is directly related to its lifetime.

**2.6.3.2. The Quasi-Classical Formula** The calculation of the lifetimes using the approaches shown so far is cumbersome and error-prone, as a precise value for the FWHM in regions where different QBS overlap is difficult to obtain. As an approximation, the lifetime of a QBS can be computed from the quasi-classical formula [83]

$$\tau_q = \frac{1}{TC(\mathcal{E}_i)} \int_0^{x_t} \sqrt{\frac{2m_i}{\mathcal{E}_i - \mathcal{E}_c(x)}}\, dx$$  (151)

where $\mathcal{E}_i$ is the resonance energy of the respective bound state and $x_t$ the classical turning point for this energy. The transmission coefficient $TC(\mathcal{E}_i)$ can be calculated by the transfer-matrix method or any other method that solves Schrödinger's equation.

**2.6.3.3. The Eigenvalues of the Non-Hermitian Hamiltonian** For open-boundary conditions, the system is described by a Hamiltonian that is not Hermitian and admits complex eigenvalues. The most straightforward way to calculate the lifetimes is to find directly the complex eigenvalues of the system Hamiltonian. This, however, is not easily possible because the eigenvalue problem is nonlinear [84]: The values of the matrix elements $\zeta$ and $\xi$ depend on the eigenvalue $\mathcal{E}$.



**Figure 13.** Wave function of quasi-bound states. Note the wave function penetration into classically forbidden regions (left). The respective value of $T_{11}^{-2}$ as a function of energy is shown in the right plot. The energy broadening around the poles is clearly visible.

Sophisicated methods have been developed to allow an easy solution of this matrix so that the lifetimes can be calculated [85–88]. First, the closed-boundary Hamiltonian is constructed, and the eigenvalues are calculated. In the one-dimensional case, the matrix is tridiagonal. I is shown in Ref. [89] that in this case, the LU algorithm is advantageous for the calculation of eigenvalues compared to the commonly used QR algorithm, which transforms the matrix into an upper Hessenberg matrix [90].

Then, the eigenvalues are filtered so that only the values remain that are located in the considered energy range. These values are then used as initial values for a Newton search around the closed-boundary eigenvalue [86, 88]. This is motivated by the fact that for $\mathscr{E}_i$ being an eigenvalue of $\underline{H}$, the determinant

$$m(\mathscr{E}_i) = \det(\underline{H} - \mathscr{E}_i \underline{I}) = 0 \tag{152}$$

must be zero. To find the roots of this equation, a Newton search around the closed-boundary eigenvalues $\mathscr{E}_i$ is used

$$\mathscr{E}_{i,j+1} = \mathscr{E}_{i,j} - \frac{m(\mathscr{E}_{i,j})}{m'(\mathscr{E}_{i,j})} \tag{153}$$

where $m(\mathscr{E})$ denotes the derivative of the determinant

$$m'(\mathscr{E}) = \frac{dm(\mathscr{E})}{d\mathscr{E}} \tag{154}$$

For a tridiagonal matrix, it is possible to find an analytical expression for $m'(\mathscr{E})$ [91, 92]. For general situations, however, the derivative can only be found numerically by

$$m'(\mathscr{E}_i) \approx \frac{m(\mathscr{E}_i + \Delta\mathscr{E}/2) - m(\mathscr{E}_i - \Delta\mathscr{E}/2)}{\Delta\mathscr{E}} \tag{155}$$

This has the advantage that it is not limited to one-dimensional problems but can be applied to any shape of the Hamiltonian.

The complex eigenvalues have been used to calculate the lifetimes of the structure shown in the left part of Fig. 13. The complex energies and lifetimes found are shown in Table 2 and agree with the values found using the method based on the evaluation of the reflection-coefficient.

## 2.7. Compact Tunneling Models

The above-presented models for the calculation of tunneling currents require a considerable computational effort. However, for practical device simulation, it is desirable to use compact models that do not require large computational resources. That may be necessary for a quick estimation of the dielectric thickness from IV data or to predict the impact of gate leakage on the performance of CMOS circuits [93–98]. The most frequently used model to describe tunneling is the Fowler–Nordheim formula [99]. The Tsu–Esaki expression (12) for the tunnel current density reads

$$J = \frac{4\pi q m_{\mathrm{eff}}}{h^3} \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} TC(\mathscr{E}_x) d\mathscr{E}_x \int_0^\infty [f_1(\mathscr{E}) - f_2(\mathscr{E})] d\mathscr{E}_\rho \tag{156}$$

Table 2. Eigenvalues found by using a resonance-finding algorithm based on the determinant of the open-boundary Hamiltonian.

| $\mathscr{E}_i$ | $\mathscr{E}_{\mathrm{R}}$ (eV) | $\mathscr{E}_{\mathrm{Im}}$ (eV) | $\tau_q$ (s) |
|---|---|---|---|
| 1 | 0.2695 | $1.503 \times 10^{-5}$ | $4.376 \times 10^4$ |
| 2 | 0.4695 | $1.830 \times 10^{-5}$ | $3.594 \times 10^3$ |
| 3 | 0.6256 | $5.285 \times 10^{-5}$ | $1.244 \times 10^1$ |
| 4 | 0.7549 | $2.794 \times 10^{-11}$ | $2.354 \times 10^{-3}$ |
| 5 | 0.8629 | $4.231 \times 10^{-5}$ | $1.555 \times 10^{-5}$ |
| 6 | 0.9503 | $2.005 \times 10^{-5}$ | $3.281 \times 10^{-11}$ |

where the total energy is split into a longitudinal and a transversal energy

$$\mathscr{E} = \mathscr{E}_\lambda + \mathscr{E}_\mu \tag{157}$$

The goal is to find a simple approximation of (156) that avoids numerical integration. As a first approximation, $T \to 0$ is assumed [1]. This allows replacement of the Fermi function $f(x)$ by the step function

$$f_1(\mathscr{E}) = f(\mathscr{E} - \mathscr{E}_{F,1}) = \begin{cases} 1 & \text{for } \mathscr{E} \le \mathscr{E}_{F,1} \\ 0 & \text{for } \mathscr{E} > \mathscr{E}_{F,1} \end{cases}$$

$$f_2(\mathscr{E}) = f(\mathscr{E} - \mathscr{E}_{F,2}) = \begin{cases} 1 & \text{for } \mathscr{E} \le \mathscr{E}_{F,2} \\ 0 & \text{for } \mathscr{E} > \mathscr{E}_{F,2} \end{cases} \tag{158}$$

Without loss of generality, it can be assumed that $\mathscr{E}_{F,1} > \mathscr{E}_{F,2}$ (see Fig. 14). The innermost integral can then be evaluated analytically for three distinct regions

$$\int_0^\lambda [f(\mathscr{E} - \mathscr{E}_{F,1}) - f(\mathscr{E} - \mathscr{E}_{F,2})] d\mathscr{E}_\mu = \mathscr{E}_{F,1} - \mathscr{E}_{F,2} \quad \text{for } \mathscr{E}_\lambda \le \mathscr{E}_{F,2}$$

$$= \mathscr{E}_{F,1} - \mathscr{E}_\lambda \quad \text{for } \mathscr{E}_{F,2} \le \mathscr{E}_\lambda \le \mathscr{E}_{F,1} \tag{159}$$

$$= 0 \quad \text{for } \mathscr{E}_\lambda > \mathscr{E}_{F,1}$$

This leads to the following expression for the current density:

$$J = \frac{4\pi q m_{eff}}{h^3} \left( \underbrace{\int_{-\infty}^{\mathscr{E}_{F,2}} TC(\mathscr{E}_x)(\mathscr{E}_{F,1} - \mathscr{E}_{F,2}) d\mathscr{E}_x}_{\approx 0} + \int_{\mathscr{E}_{F,2}}^{\mathscr{E}_{F,1}} TC(\mathscr{E}_x)(\mathscr{E}_{F,1} - \mathscr{E}_x) d\mathscr{E}_x \right) \tag{160}$$

The left integral represents tunneling current from electron states that are low in energy and face a high energy barrier. Hence, as a second approximation, the left integral is neglected. Still, it is necessary to insert an expression for the transmission coefficient in the right integral. For a single-layer dielectric, two shapes are possible: triangular and trapezoidal. First, the formula will be derived assuming a triangular shape.



Figure 14. Schematic of an energy barrier in the Fowler-Nordheim tunneling (left) and direct tunneling (right) regime.

### 2.7.1. Original Fowler–Nordheim Formula

The original Fowler–Nordheim formula assumes a triangular shape of the energy barrier. This is motivated by the fact that only tunneling at strong electric fields was studied. The WKB approximation (70) for the transmission coefficient reads

$$TC(\mathscr{E}_x) = \exp\left(-\frac{2}{\hbar}\int_0^{x_1}\sqrt{2m_{\text{diel}}(\mathscr{E}_c - \mathscr{E}_x)}\,dx\right) \tag{161}$$

The classical turning point $x_1$ is (see the left part of Fig. 14)

$$x_1 = \frac{\mathscr{E}_{\text{F,1}} + q\Phi_1 - \mathscr{E}_x}{qE_{\text{diel}}} \tag{162}$$

and the dielectric conduction band edge for a triangular barrier

$$\mathscr{E}_c(x) = \mathscr{E}_{\text{F,1}} + q\Phi_1 - qE_{\text{diel}}x \tag{163}$$

where the electric field in the dielectric $E_{\text{diel}}$ is caused by the different Fermi levels and the work function difference $\Delta\Phi_W$:

$$E_{\text{diel}} = \frac{\mathscr{E}_{\text{F,1}} - \mathscr{E}_{\text{F,2}} + q\Delta\Phi_W}{qt_{\text{diel}}} \tag{164}$$

The third approximation is to assume equal materials for both electrodes, so that $\Delta\Phi_W = 0$. The WKB-based transmission coefficient can then be applied and yields

$$TC(\mathscr{E}_x) = \exp\left(-2\frac{\sqrt{2m_{\text{diel}}}}{\hbar}\int_0^{x_1}\sqrt{\mathscr{E}_{\text{F,1}} + q\Phi_1 - qE_{\text{diel}}x - \mathscr{E}_x}\,dx\right) \tag{165}$$

$$= \exp\left[4\frac{\sqrt{2m_{\text{diel}}}}{3\hbar qE_{\text{diel}}}(\mathscr{E}_{\text{F,1}} + q\Phi_1 - qE_{\text{diel}}x - \mathscr{E}_x)^{3/2}\Big|_0^{x_1}\right] \tag{166}$$

$$= \exp\left[4\frac{\sqrt{2m_{\text{diel}}}}{3\hbar qE_{\text{diel}}}(-\mathscr{E}_{\text{F,1}} - q\Phi_1 + \mathscr{E}_x)^{3/2}\right] \tag{167}$$

$$= \exp\left\{-4\frac{\sqrt{2m_{\text{diel}}}}{3\hbar qE_{\text{diel}}}[q\Phi_1 - (\mathscr{E}_x - \mathscr{E}_{\text{F,1}})]^{3/2}\right\} \tag{168}$$

Using this expression in (160), the current density becomes

$$J = \frac{4\pi qm_{\text{eff}}}{h^3}\int_{\mathscr{E}_{\text{F,2}}}^{\infty}\exp\left\{-\frac{4\sqrt{2m_{\text{diel}}}}{3\hbar qE_{\text{diel}}}[q\Phi_1 - (\mathscr{E}_x - \mathscr{E}_{\text{F,1}})]^{3/2}\right\}(\mathscr{E}_{\text{F,1}} - \mathscr{E}_x)\,d\mathscr{E}_x \tag{169}$$

This integral cannot be solved analytically. Hence, the fourth approximation is to expand the square root into a first-order Taylor series around $q\Phi_1$:

$$[q\Phi_1 - (\mathscr{E}_x - \mathscr{E}_{\text{F,1}})]^{3/2} \approx (q\Phi_1)^{3/2} + \frac{3}{2}(\mathscr{E}_x - \mathscr{E}_{\text{F,1}})(q\Phi_1)^{1/2} \tag{170}$$

Inserting this expression into (169) and setting $\epsilon = \mathscr{E}_x - \mathscr{E}_{\text{F,1}}$ yields

$$J = \frac{4\pi qm_{\text{eff}}}{h^3}\exp\left[-\frac{4\sqrt{2m_{\text{diel}}}}{3\hbar qE_{\text{diel}}}(q\Phi_1)^{3/2}\right]\int_0^{\infty}\exp\left[\frac{2\sqrt{2m_{\text{diel}}}}{\hbar qE_{\text{diel}}}(q\Phi_1)^{1/2}\epsilon\right]\epsilon\,d\epsilon \tag{171}$$

With

$$\int\epsilon\exp(\lambda\epsilon)\,d\epsilon = \frac{1}{\lambda^2}\exp(\lambda\epsilon)(\lambda\epsilon - 1) \tag{172}$$

and

$$a = -\frac{4\sqrt{2m_{diel}}}{3\hbar q E_{diel}}(q\Phi_1)^{3/2}, \quad \lambda = \frac{2\sqrt{2m_{diel}}}{\hbar q E_{diel}}(q\Phi_1)^{1/2}$$  (173)

the current density becomes

$$J = \frac{4\pi q m_{eff}}{h^3}\exp(a)\int_{\epsilon_{F,2}-\epsilon_{F,1}}^{0}\exp(\lambda\epsilon)\epsilon\,d\epsilon$$  (174)

$$= \frac{4\pi q m_{eff}}{h^3}\exp(a)\frac{1}{\lambda^2}\exp[\lambda(\epsilon_{F,2}-\epsilon_{F,1})][\lambda(\epsilon_{F,2}-\epsilon_{F,1})-1]$$  (175)

The fifth assumption is now that $\epsilon_{F,1} \gg \epsilon_{F,2}$, leading to

$$J = \frac{4\pi q m_{eff}}{h^3}\exp(a)\frac{1}{\lambda^2}$$  (176)

or

$$J = \frac{q^3 m_{eff}}{8\pi m_{diel}\hbar q\Phi_1}E_{diel}^2\exp\left(-\frac{4\sqrt{2m_{diel}(q\Phi_1)^3}}{3\hbar q E_{diel}}\right)$$  (177)

which is the equation commonly known as the Fowler–Nordheim formula [100]. Note that there is a difference between the effective electron mass in the electrode ($m_{eff}$) and the effective electron mass in the dielectric ($m_{diel}$).

## 2.7.2. Correction for Direct Tunneling

The equation derived above is only valid for triangular barriers; that is, the case of high applied voltages. In Ref. [101], Schuegraf proposed a correction to the Fowler–Nordheim formula to account for tunneling in the direct tunneling regime. In this case, the transmission coefficient is

$$TC(\epsilon) = \exp\left(-\frac{2}{\hbar}\int_0^{t_{diel}}\sqrt{2m_{diel}(\epsilon_c-\epsilon_x)}\,dx\right)$$  (178)

where $t_{diel}$ is the dielectric thickness. The conduction band edge is again approximated by a linear shape

$$\epsilon_c(x) = \epsilon_{F,1} + q\Phi_1 - qE_{diel}x$$  (179)

The band edges $q\Phi$ and $q\Phi_0$ are given by (see the right part of Fig. 14)

$$q\Phi = \epsilon_{F,1} + q\Phi_1$$  (180)

$$q\Phi_0 = \epsilon_{F,1} + q\Phi_1 - qE_{diel}t_{diel}$$  (181)

As for the triangular energy barrier, it is assumed that the electrodes have equal work functions: $\Delta\Phi_W = 0$. Using these expressions, the transmission coefficient becomes

$$TC(\epsilon_x) = \exp\left(-2\frac{\sqrt{2m_{diel}}}{\hbar}\int_0^{t_{diel}}\sqrt{q\Phi-qE_{diel}x-\epsilon_x}\,dx\right)$$  (182)

$$= \exp\left[4\frac{\sqrt{2m_{diel}}}{3\hbar q E_{diel}}(q\Phi-qE_{diel}x-\epsilon_x)^{3/2}\Big|_0^{t_{diel}}\right]$$  (183)

$$= \exp\left\{-4\frac{\sqrt{2m_{diel}}}{3\hbar q E_{diel}}[(q\Phi-\epsilon_x)^{3/2}-(q\Phi_0-\epsilon_x)^{3/2}]\right\}$$  (184)

The exponent can be approximated using a first-order Taylor series expansion around $q\Phi_1$ and $q\Phi_1 - qE_{diel}t_{diel}$, respectively:

$$(q\Phi - \mathcal{E}_x)^{3/2} = (\mathcal{E}_{F,1} + q\Phi_1 - \mathcal{E}_x)^{3/2} \tag{185}$$

$$= [q\Phi_1 - (\mathcal{E}_x - \mathcal{E}_{F,1})]^{3/2} \tag{186}$$

$$\approx (q\Phi_1)^{3/2} + \frac{3}{2}(\mathcal{E}_x - \mathcal{E}_{F,1})(q\Phi_1)^{1/2} \tag{187}$$

$$(q\Phi_{11} - \mathcal{E}_x)^{3/2} = (\mathcal{E}_{F,1} + q\Phi_1 - qE_{diel}t_{diel} - \mathcal{E}_x)^{3/2} \tag{188}$$

$$= [(q\Phi_1 - qE_{diel}t_{diel}) - (\mathcal{E}_x - \mathcal{E}_{F,1})]^{3/2} \tag{189}$$

$$\approx (q\Phi_1 - qE_{diel}t_{diel})^{3/2} + \frac{3}{2}(\mathcal{E}_x - \mathcal{E}_{F,1})(q\Phi_1 - qE_{diel}t_{diel})^{1/2} \tag{190}$$

With the temporary variable $\eta$

$$\eta = (q\Phi - \mathcal{E}_x)^{3/2} - (q\Phi_0 - \mathcal{E}_x)^{3/2}$$

$$\approx -(q\Phi_1 - qE_{diel}t_{diel})^{3/2} + (q\Phi_1)^{3/2} - \frac{3}{2}(\mathcal{E}_x - \mathcal{E}_{F,1})[(q\Phi_1)^{1/2} - q\Phi_1 - qE_{diel}t_{diel})^{1/2}] \tag{191}$$

the tunnel current density becomes

$$J = \frac{4\pi q m_{eff}}{h^3} \int_{\mathcal{E}_{1,2}}^{\mathcal{E}_{F,1}} TC(\mathcal{E}_x)(\mathcal{E}_{F,1} - \mathcal{E}_x)d\mathcal{E}_x \tag{192}$$

$$\approx \frac{4\pi q m_{eff}}{h^3} \int_{\mathcal{E}_{1,2}}^{\mathcal{E}_{F,1}} \exp\left(-4\frac{\sqrt{2m_{diel}}}{3\hbar q E_{diel}}\eta\right)(\mathcal{E}_{F,1} - \mathcal{E}_x)d\mathcal{E}_x \tag{193}$$

With the abbreviations

$$a = \frac{4\pi q m_{eff}}{h^3} \tag{194}$$

$$b = -\frac{4\sqrt{2m_{diel}}}{3\hbar q E_{diel}}[(q\Phi_1)^{3/2} - (q\Phi_1 - qE_{diel}t_{diel})^{3/2}] \tag{195}$$

$$c = -\frac{2\sqrt{2m_{diel}}}{\hbar q E_{diel}}[(q\Phi_1)^{1/2} - (q\Phi_1 - qE_{diel}t_{diel})^{1/2}] \tag{196}$$

the tunnel current density can be written as

$$J = a\exp(b)\int_{\mathcal{E}_{1,1}}^{\mathcal{E}_{F,2}} \exp[c(\mathcal{E}_x - \mathcal{E}_{F,1})](\mathcal{E}_{F,1} - \mathcal{E}_x)d\mathcal{E}_x \tag{197}$$

With $\epsilon = \mathcal{E}_x - \mathcal{E}_{F,1}$ this yields

$$J = -a\exp(b)\int_{\mathcal{E}_{1,2}-\mathcal{E}_{F,1}}^{0} \exp(c\epsilon)\epsilon\, d\epsilon \tag{198}$$

Using (172), this integral becomes

$$J = \frac{a\exp(b)}{c^2}\{1 - \exp[-c(\mathcal{E}_{F,1} - \mathcal{E}_{F,2})][1 + c(\mathcal{E}_{F,1} - \mathcal{E}_{F,2})]\} \tag{199}$$

which, for $\mathcal{E}_{F,1} \gg \mathcal{E}_{F,2}$, simplifies to

$$J = \frac{a\exp(b)}{c^2} \tag{200}$$

or, inserting the expressions for $a$, $b$, and $c$

$$J = \frac{q^3 m_{\text{eff}}}{8\pi h m_{\text{diel}}[(q\Phi_1)^{1/2} - (q\Phi_1 - qV_{\text{diel}})^{1/2}]^2} E_{\text{diel}}^2$$

$$\times \exp\left\{-\frac{4\sqrt{2m_{\text{diel}}}}{3\hbar q E_{\text{diel}}}[(q\Phi_1)^{3/2} - (q\Phi_1 - qV_{\text{diel}})^{3/2}]\right\} \tag{20.1}$$

which is the equation used in Refs. [101, 102]. In some publications, the equation is rewrtten to make it more similar to the Fowler–Nordheim formula:

$$J = \frac{q^3 m_{\text{eff}}}{8\pi m_{\text{diel}} h q\Phi_1 B_1} E_{\text{diel}}^2 \exp\left(-\frac{4\sqrt{2m_{\text{diel}}(q\Phi_1)^3} B_2}{3\hbar q E_{\text{diel}}}\right) \tag{20.2}$$

with the additional correction terms $B_1$, $B_2$ given as

$$B_1 = \left[1 - \left(1 - \frac{qE_{\text{diel}}t_{\text{diel}}}{q\Phi_1}\right)^{1/2}\right]^2$$

$$B_2 = \left[1 - \left(1 - \frac{qE_{\text{diel}}t_{\text{diel}}}{q\Phi_1}\right)^{3/2}\right] \tag{20.3}$$

For a triangular barrier, the correction factors become $B_1 = B_2 = 1$, and the expression simplifies to (177). Note that using these equations, the minimum tunneling current occurs for $E_{\text{diel}} = 0$ V/m, which, for a work function difference $\neq 0$, does not occur at the minimum applied bias.

## 2.8. Trap-Assisted Tunneling

Besides direct or Fowler–Nordheim tunneling, which are one-step tunneling processes, defects in the dielectric layer give rise to tunneling processes based on two or more steps. This tunneling component is mainly observed after writing-erasing cycles in electrically erasable programmable read-only-memories (EEPROMs). It is therefore assumed that traps arise in the dielectric layer due to the repeated high-voltage stress. The increased tunneling current at low bias is called stress-induced leakage current (SILC) and is mainly responsible for the degradation of the retention time of nonvolatile memory devices [103]. It is now generally accepted that it is caused by inelastic trap-assisted tunnel transitions and that the traps are created by the electric high-field stress during the writing and erasing processes [103–108]. SILC has widely been studied and modeled in MOS capacitors [109–111] and EEPROM devices [112].

This section gives a brief overview of trap-assisted tunneling models, describes two frequently encountered models (Chang's and Iclmini's model), and elaborates on a sophisticated model originally proposed by Jimenez et al. The adaption of this model to allow its inclusion in device simulators is described in some detail.

### 2.8.1. Model Overview

Numerous models have been presented to describe trap-assisted tunneling in the gate dielectric of MOS devices. These models usually share the equation for the current density, which is given by an integration along the gate dielectric [113]:

$$J = q \int_0^{t_{\text{diel}}} \frac{N_T(x)}{\tau_c(x) + \tau_e(x)} \, dx \tag{20.4}$$

In this expression, $N_T$ denotes the trap concentration, and $\tau_c$ and $\tau_e$ denote the capture and emission times of the considered trap. Because both processes—capture and emission—must happen in sequence, they both determine the current density. However, differences exist in how the capture and emission times are calculated. Some models use constant capture and emission cross sections to calculate the respective times. Another important point is the

distribution in space, where the traps are usually assumed to follow a Gaussian distribution. The distribution in energy is also crucial. Commonly, it is either assumed that traps have a Gaussian distribution in energy or that they are located at a certain energy level below the dielectric conduction band. The assumption of a discrete energy level for specific trap types is backed by spectroscopic analyses [114]. Additionally, the tunneling process can either be elastic, where the energy of the tunneling electron is conserved, or inelastic, where the energy of the tunneling electron changes. Recent studies and experiments have shown strong evidence for the tunneling process being inelastic [115–117].

### 2.8.1.1. Chang's Model

A frequently used model is the generalized trap-assisted tunneling model presented by Chang et al. [118, 119]. The current density reads

$$J = q \int_0^{t_{diel}} A N_T(x) \frac{P_1(x) P_2(x)}{P_1(x) + P_2(x)} dx \tag{205}$$

where $A$ denotes a fitting constant, $N_T(x)$ the spatial trap concentration, and $P_1$ and $P_2$ the transmission coefficients of electrons captured and emitted by traps. Using $\tau_c \sim P_1/P_2$ and $\tau_e \sim P_2/P_1$, this expression reduces to (204). A similar model was used by Ghetti et al. [76]

$$J = \int_0^{t_{diel}} C_T N_T(x) \frac{J_{in} J_{out}}{J_{in} + J_{out}} dx \tag{206}$$

who assumed a constant capture cross section $C_T$ for the traps. The symbols $J_{in}$ and $J_{out}$ denote the capture and emission currents. Essentially the same formula was used by other authors as well [116, 120].

### 2.8.1.2. Ielmini's Model

Considerable research has been done by Ielmini et al. [121–124], who describe inelastic TAT and also take hopping conduction into account [125, 126]. They derive the trap-assisted current by an integration along the dielectric thickness and energy

$$J = \int_0^{t_{diel}} dx \int_{\mathscr{E}_{min}}^{\mathscr{E}_{max}} \tilde{J}(\mathscr{E}_T, x) d\mathscr{E} \tag{207}$$

where $\tilde{J}$ denotes the net current flowing through the dielectric, given as the difference between capture and emission currents through either side of the dielectric

$$\tilde{J}(\mathscr{E}_T, x) = J_{el} - J_{cl} = J_{cr} - J_{er} = q N_T' W_c \left( 1 - \frac{f_T(\mathscr{E}_T, x)}{f_l(\mathscr{E}_T, x)} \right) \tag{208}$$

where $f_T$ is the trap occupancy, $\mathscr{E}_T$ the trap energy, $W_c$ the capture rate, and $f_l$ the energy distribution function at the left interface. The symbol $N_T'$ denotes the trap concentration in space and energy. Ielmini further develops the model to include transient effects, and notes that in this case, the net difference between current from the left and right interfaces equals the change in the trap occupancy multiplied by the trap charge

$$(J_{cl} - J_{el}) + (J_{er} - J_{cr}) = q N_T \frac{\partial f_T}{\partial t} \tag{209}$$

an observation that will be revisited in Section 2.8.2.4. The model assumes a constant capture cross section.

### 2.8.1.3. Compact Trap-Assisted Tunneling Models

For application in circuit simulators or to catch a quick glimpse at the effects of trap-assisted tunneling, compact models are required. A frequently used expression is based on the work of Ricco et al. [109]. They describe the trapping- and detrapping processes by

$$J_{TAT} = JC_1 TC_1(N_T - n_T) = q\nu n_T TC_2 \tag{210}$$

where $J$ is the supply current density at the interface, $C_T$ the capture cross section, $TC_1$ and $TC_2$ the transmission coefficients from the left and right sides of the dielectric to the trap, $n_T$ the concentration of trapped electrons that is smaller than or equal to the trap concentration $N_T$, and $\nu$ their escape frequency. The highest contribution comes from traps that have $TC_1 \approx TC_2$; therefore, the trap-assisted tunnel current becomes

$$J_{TAT} = q\nu n_T TC = q\nu C_T N_1 \frac{J}{JC_T + q\nu} TC \tag{211}$$

A modified version of this expression was used by Ghetti et al. [111, 127]. Other more or less empirical trap-assisted tunneling models based on SILC measurements are presented in Ref. [128]. These comprise hopping conduction

$$J = C_1 E_{diel} \exp\left(-\frac{q\Phi_a}{k_B T}\right) \tag{212}$$

where $\Phi_a$ is an activation potential, and the frequently applied Poole–Frenkel tunneling formula [128–134]. This model describes the emission of trapped electrons and reads

$$J = AE_{diel} \exp\left(-\frac{\ell_T}{k_B T}\right) \exp\left(\frac{q}{k_B T} \sqrt{\frac{qE_{diel}}{\pi \kappa_0 r^2}}\right) \tag{213}$$

where $r$ is the refractive index of the dielectric, $\ell_T$ is the difference between the conduction band in the dielectric and the trap energy, and the coefficient $A$ depends on the trap concentration. The main motivation to use this expression is that the trap-assisted gate current density was found to be a linear function of the square root of the dielectric field, in contrast to the Fowler–Nordheim tunneling current, which is a linear function of the dielectric field. Note, however, that no trapping-detrapping considerations enter this equation.

### 2.8.2. The Model of Jiménez et al.

A model for trap-assisted inelastic tunneling has been developed by Jiménez et al. [135]. Their model is based on the theory of nonradiative capture and emission of electrons by multiphonon processes [136]. The main difference to the models described before is that it does not require constant capture cross sections as fitting parameters but calculates them for each trap based on the trap energy level and the shape of the energy barrier.

#### 2.8.2.1. Capture and Emission Probabilities

The tunneling model is based on a two-step tunneling process via traps in the dielectric that incorporates energy loss by phonon emission [135]. Figure 15 shows the basic two-step process of an electron tunneling from a region with higher Fermi energy (the cathode) to a region with lower Fermi energy (the anode). To avoid integration in energy, the initial electron energy is assumed to be located at the average kinetic energy, which, for the parabolic dispersion relation (1) and the Maxwellian distribution (20), is

$$\frac{\langle \ell \rangle}{\langle 1 \rangle} = \frac{\int_0^\infty \ell f(\ell) g(\ell) d\ell}{\int_0^\infty f(\ell) g(\ell) d\ell} = \frac{\int_0^\infty \ell^{3/2} \exp(-\frac{\ell}{k_B T}) d\ell}{\int_0^\infty \ell^{1/2} \exp(-\frac{\ell}{k_B T}) d\ell} = \frac{3}{2} k_B T \tag{214}$$

During the capture process $(W_c)$, the difference in total energy between the initial and final state is released by means of phonon emission $(\hbar\omega)$. An electron captured by a trap can then be emitted into the anode $(W_e)$.

**Figure 15.** The trap-assisted tunneling process [135].

The rate with which an electron with energy $\ell$ is captured by a trap located at position $x$ and energy $\varepsilon'$ is given by [137]

$$W_c(x, \varepsilon', \ell) = \frac{\pi}{\hbar^2 \omega} |V_c|^2 S \left( 1 - \frac{P}{S} \right)^2 I_p(\xi) \exp\left[ -(2f_p + 1)S + \frac{\Delta \ell}{2k_B T} \right] \qquad (215)$$

Here, $S$ is the Huang–Rhys factor, which characterizes the electron–phonon interaction [138]; $\hbar\omega$ is the energy of the phonons involved in the transitions, $\Delta \ell = \ell - \ell'$; and $P = \Delta \ell / \hbar\omega$ is the number of phonons emitted due to this energy difference. In the simulations, the value of $S\hbar\omega$ was used as fitting parameter.

The population of phonons $f_p$ is given by the Bose–Einstein statistics

$$f_p = \left[ \exp\left( \frac{\hbar\omega}{k_B T} \right) - 1 \right]^{-1} \qquad (216)$$

The function $I_p(\xi)$ is the modified Bessel function of order $P$, with

$$\xi = 2S\sqrt{f_p(f_p + 1)} \qquad (217)$$

The term $|V_c|^2$ in (215) denotes the transition matrix element, which is calculated by an integration over the trap cube [136]

$$|V_c|^2 = 5\pi S(\hbar\omega)^2 \frac{\hbar^2}{2m_{diel} \ell_T} \int_{x_0 - x_T/2}^{x_0 + x_T/2} |\Psi(x)|^2 dx \qquad (218)$$

In this expression, $x_T$ denotes the side length of the trap cube, estimated as

$$x_T = \frac{\hbar}{\sqrt{2m_{diel} \ell_T}} \left( \frac{4\pi}{3} \right)^{1/3} \qquad (219)$$

The symbol $\ell_T$ denotes the energy difference between the trap energy and the barrier conduction band edge as shown in Fig. 15. For the emission of electrons from the trap to the anode, elastic tunneling is assumed. Hence, the probability of emission to the anode is equal to the probability of capture from the anode, which is calculated from (215).

The numerical evaluation of (218) requires the calculation of the wave functions in the dielectric layer, which, however, degrades the computational efficiency of a multipurpose device simulator where simulation speed is crucial. To avoid this, the barriers have been transformed to take advantage of the well-known solutions for constant potentials. Two cases must be distinguished; namely, the case of a trapezoidal barrier and the case of a triangular barrier. The two cases are depicted in Fig. 16.

For capture processes and for emission processes where the electron faces a trapezoidal barrier, the barrier is transformed into a step function of height equal to the potential at the middle point between $x = 0$ and $x = x_0$ ($\notin_m$ in the left part of Fig. 16), $x_0$ being the position of the trap inside the dielectric. Assuming

$$\Psi(x \le 0) = A \sin(k_1 x + \alpha)$$

$$\Psi(x > 0) = B \exp(-k_2 x) \tag{220}$$

the wave function at the position of the trap becomes

$$\Psi(x) = A \sin\left[\arctan\left(\frac{m_{\text{diel}}}{m_{\text{eff}}} \frac{k_1}{k_2}\right)\right] \exp(-k_2 x) \tag{221}$$

where $m_{\text{diel}}$ and $m_{\text{eff}}$ are the electron masses in the dielectric and the neighboring electrode, respectively. The wave numbers are given by

$$k_1 = \frac{1}{\hbar}\sqrt{2m_{\text{eff}}(\notin - \notin_c)}$$

$$k_2 = \frac{1}{\hbar}\sqrt{2m_{\text{diel}}(\notin_m - \notin)} \tag{222}$$

For emission processes in which the barrier is triangular (the electron energy is above the dielectric conduction band at some point between the trap and the anode), two regions in the dielectric must be distinguished. The first one, between the interface at $x = 0$ and the point $x = x_{\text{FN}}$ (see the right part of Fig. 16) has the height $\notin_{\text{FN}}$. The height of the approximated barrier in the other region is then the value of the barrier, $\notin_m$, in the middle point between $x = x_{\text{FN}}$ and the position of the trap $x = x_0$. With this new barrier and the assumptions for the wave functions in the three regions

$$\Psi(x \le 0) = A \sin(k_1 x + \alpha_1) \tag{223}$$

$$\Psi(0 < x \le x_{\text{FN}}) = B \sin(k_2 x + \alpha_2) \tag{224}$$

$$\Psi(x_{\text{FN}} < x \le x_0) = C \exp[-k_3(x - x_{\text{FN}})] \tag{225}$$



Figure 16. The approximate shape of the barrier in the direct (left) and Fowler-Nordheim regime (right) [171].

the wave function at the position of the trap becomes

$$\Psi(x) = A\frac{\sin\alpha_1}{\sin\alpha_2}\sin(k_2 x_{\mathrm{TN}} + \alpha_2)\exp[-k_3(x - x_{\mathrm{TN}})]$$  (226)

with the symbols

$$\alpha_1 = \arctan\left(\frac{k_1}{k_2}\tan\alpha_2\right)$$

$$\alpha_2 = \arctan\left(\frac{k_2}{k_3}\right) - k_2 x_{\mathrm{TN}}$$  (227)

The corresponding wave numbers are given as

$$k_1 = \frac{1}{\hbar}\sqrt{2m_{\mathrm{eff}}(\mathscr{E} - \mathscr{E}_c)}$$

$$k_2 = \frac{1}{\hbar}\sqrt{2m_{\mathrm{diel}}(\mathscr{E} - \mathscr{E}_{\mathrm{TN}})}$$  (228)

$$k_3 = \frac{1}{\hbar}\sqrt{2m_{\mathrm{diel}}(\mathscr{E}_m - \mathscr{E})}$$

Using expression (221) and (226), the integration in (218) can be performed analytically, which allows the capture and emission probabilities to be calculated without the need for numerical integration.

**2.8.2.2. Capture and Emission Times** Once the capture and emission probabilities have been obtained, the corresponding times can be calculated. The inverse of the capture time is given by [135, 139]

$$\tau_c^{-1}(x) = \int_{\mathscr{E}}^{\infty} W_c(x, \mathscr{E}', \mathscr{E})g_c(\mathscr{E})f_c(\mathscr{E})d\mathscr{E}$$  (229)

where $g_c(\mathscr{E})$ denotes the two-dimensional density of states and $f_c(\mathscr{E})$ the electron energy distribution function in the cathode. For the above-stated assumption that all electrons are captured from the same energy level $\mathscr{E}_c + 3/2k_B T$ in the cathode, this expression can be approximated by

$$\tau_c^{-1}(x) \approx W_c\left(x, \mathscr{E}', \mathscr{E}_c + \frac{3}{2}k_B T\right)n_c$$  (230)

where $n_c$ is the sheet carrier concentration in the cathode, which is determined by the transport model used in the device simulator. The inverse of the emission time is [135]

$$\tau_e^{-1}(x) = \int_{-\infty}^{\mathscr{E}'} W_e(x, \mathscr{E}', \mathscr{E})g_a(\mathscr{E})[1 - f_a(\mathscr{E})]d\mathscr{E}$$  (231)

Assuming $f_a(\mathscr{E}) \approx 0$ in the anode and elastic tunneling for the emission process ($\mathscr{E} = \mathscr{E}'$), the emission time becomes

$$\tau_e^{-1}(x) \approx W_e(x, \mathscr{E}', \mathscr{E}')g_a(\mathscr{E}')\hbar\omega$$  (232)

where the energy loss is restricted to values less than $\hbar\omega$. To check the validity of the approximations for the wave functions, the resulting capture and emission times have been compared to results using a Schrödinger-Poisson solver for a MOS capacitor with the parameters $\mathscr{E}_T = 2.8$ eV, $S\hbar\omega = 1.6$ eV, and a trap concentration of $N_T = 10^{19}$ cm$^{-3}$. As can be seen in Fig. 17, the analytical and the numerical results are very close. Electrons are captured from the right and emitted to the left in this figure. Thus, for traps near the right side of the barrier, the capture time is very low and the emission time is very high. The oscillations in the emission time for high bias are due to the fact that in this regime, the energy barrier has a triangular shape, which gives rise to an oscillating wave function, in contrast to the decaying wave function for a trapezoidal barrier.

**Figure 17.** Comparison of the analytic solution with a numerical solution for the capture and emission times at a gate bias of 3 V (left) and 7 V (right) [171].

### 2.8.2.3. Steady-State Current

The total steady-state tunneling current is derived as the sum of the trap-assisted tunneling current (204) and the direct tunneling current computed from the Tsu–Esaki formula (13)

$$J = J_{\text{TAT}} + J_{\text{Tsu-Esaki}} \tag{233}$$

Figure 18 shows the dependence of the gate current density on the model parameters $E_T$ (trap energy level) and $S\hbar\omega$ for a fixed phonon energy of $\hbar\omega = 10$ meV in an MOS capacitor. For a low trap energy level, traps are located near the conduction band edge in the dielectric, and direct tunneling prevails. With increasing trap energy level, the trap-assisted component becomes stronger and exceeds the direct tunneling current for low bias. The current density shows a peak at low bias, which is due to the alignment of the trap energy level with the cathode conduction band edge. The Huang–Rhys factor has only a minor influence on the results, as shown in the right part of Fig. 18.

### 2.8.2.4. Transient Current

Models of trap-assisted transitions are commonly employed to calculate steady-state SILC in MOS capacitors, whereas transient SILC has hardly been studied [110, 121]. However, transient tunneling current becomes important at high switching speed where the transients of the trap charging and discharging processes may degrade



**Figure 18.** Dependence of the tunneling current on the trap energy level (left) and on the Huang–Rhys factor for a fixed phonon energy of 10 meV (right) [171].

signal integrity. For the calculation of transient SILC, it is necessary to calculate capture and emission times at each time step. Considering a spatial trap distribution $N_T(x)$ across the dielectric layer, the rate equation for the concentration of occupied traps at position $x$ reads

$$N_T(x)\frac{df_T(x,t)}{dt} = N_T(x)[1 - f_T(x,t)]\tau_c^{-1}(x,t) - N_T(x)f_T(x,t)\tau_e^{-1}(x,t) \qquad (234)$$

where $f_T(x,t)$ is the trap occupancy function, and $\tau_c(x,t)$ and $\tau_e(x,t)$ are the inverse capture and emission times of electrons by a trap placed at position $x$. In the static case, capture and emission processes are in equilibrium and $df_T(x,t)/dt = 0$. In the transient case, however, capture and emission times include transitions from the cathode and the anode

$$\tau_c^{-1}(x,t) = \tau_{ca}^{-1}(x,t) + \tau_{cc}^{-1}(x,t)$$
$$\tau_e^{-1}(x,t) = \tau_{ea}^{-1}(x,t) + \tau_{ec}^{-1}(x,t) \qquad (235)$$

where $\tau_{ca}$ and $\tau_{cc}$ are the capture times to the anode and to the cathode and $\tau_{ea}$ and $\tau_{ec}$ the corresponding emission times. To calculate the local trap occupancy, the differential equation (234) must be solved. If the capture and emission times $\tau_c^{-1}$ and $\tau_e^{-1}$ are constant over time, like in a discharging process with a constant potential distribution, the solution of (234) can be given in a closed form

$$f_T(x,t) = f_T(x,0)\exp\left(-\frac{t}{\tau_m(x,t)}\right) + \frac{\tau_m(x,t)}{\tau_c(x,t)}\left[1 - \exp\left(-\frac{t}{\tau_m(x,t)}\right)\right] \qquad (236)$$

with $\tau_m^{-1} = \tau_c^{-1} + \tau_e^{-1}$.

A more general approach is to look at the change of the trap distribution at discrete time steps. Integration of (234) in time between $t_i$ and $t_{i+1}$ and changing to discrete time steps yields

$$f_T(x,t_i) - f_T(x,t_{i-1}) \approx \tau_c^{-1}(x,t_{i-1})\Delta t_i - \tau_m^{-1}(x,t_{i-1})\bar{f}_i\Delta t_i$$

where the abbreviations $\Delta t_i = t_i - t_{i-1}$ and $\bar{f}_i = [f_T(x,t_i) + f_T(x,t_{i-1})]/2$ have been used. Thus, it is possible to write the trap distribution over time in the following recursive manner:

$$f_T(x,t_i) = A_i + B_i f_T(x,t_{i-1}) \qquad (237)$$

where the symbols $A_i$, $B_i$, and $C_i$ are calculated from

$$A_i = \frac{\tau_c^{-1}(x,t_i)\Delta t_i}{1 + C_i}$$
$$B_i = \frac{1 - C_i}{1 + C_i} \qquad (238)$$
$$C_i = \frac{\tau_m^{-1}(x,t_i)\Delta t_i}{2}$$

Once the time-dependent occupancy function in the dielectric is known, the tunnel current through each of the interfaces is

$$J_{TAT,Anode}(t) = q\int_0^{d_{el}} N_T(x)\{\tau_{ea}^{-1}(x,t) - f_T(x,t)[\tau_{ca}^{-1}(x,t) + \tau_{ea}^{-1}(x,t)]\}dx \qquad (239)$$

$$J_{TAT,Cathode}(t) = q\int_0^{d_{el}} N_T(x)\{\tau_{ec}^{-1}(x,t) - f_T(x,t)[\tau_{cc}^{-1}(x,t) + \tau_{ec}^{-1}(x,t)]\}dx \qquad (240)$$

## 2.9. Model Comparison

This section outlined a number of tunneling models useful for the simulation of tunneling in semiconductor devices. For practical device simulation, however, it is often not clear which model to select for the application at hand. Therefore, Table 3 summarizes the main model

Table 3. A hierarchy of tunneling models and their properties.

| | Fowler-Nordheim Model | Schuegraf Model | Tsu-Esaki Analytic WKB | Tsu-Esaki Gundlach | Tsu-Esaki Numeric WKB | Tsu-Esaki Transfer-Matrix | Tsu-Esaki QTBM | Inelastic TAT | Frenkel-Poole |
|---|---|---|---|---|---|---|---|---|---|
| FN tunneling | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| Direct tunneling | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| EVB tunneling process | | | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| QM current oscillations | | | | ✓ | | ✓ | ✓ | | |
| Dielectric stacks | | | | ✓ | ✓ | ✓ | ✓ | | |
| Numerical stability | | | | | | — | | | |
| Trap-assisted tunneling | | | | | | | | ✓ | ✓ |
| Trap occupancy modeling | | | | | | | | ✓ | |
| Transient TAT | | | | | | | | ✓ | |
| Computational effort | Low | Low | | | High | High | High | | Low |

Note: WKB, Wentzel–Kramers–Brillouin; QTBM, quantum transmitting boundary method; TAT, trap-assisted tunneling; FN, Fowler–Nordheim; EVB, electrons from valence band; QM, quantum mechanical.

features and also gives the approximate computational effort. The following points can be concluded [140]:

- Especially the Fowler–Nordheim, Schuegraf, and Frenkel–Poole models have a very low computational effort because they are compact models. However, they do not correctly reproduce the device physics and can only be used after careful calibration.
- The Tsu–Esaki formula with the analytical WKB or Gundlach method for the transmission coefficient combines moderate computational effort with reasonable accuracy. This approach can be used for the simulation of tunneling in devices with single-layer dielectrics.
- The inelastic TAT model allows simulation of all effects related with traps in the dielectric and, due to the analytical calculation of the overlap integral, poses only moderate computational effort. This model can be used for the simulation of leakage in EEPROMs or trap-rich dielectric devices (see Section 3.2.2.1).
- The Tsu–Esaki model with the numerical WKB, transfer-matrix, or QTB method to calculate the transmission coefficient represents the most accurate method usable for the simulation of tunneling through dielectric stacks, however, with high computational effort. The transfer-matrix method should be used with care due to its poor numerical stability.

## 3. APPLICATIONS

Gate leakage is one of the most important issues for contemporary complementary metal-oxide semiconductor (CMOS) devices. Based on the tunneling models outlined so far, two different application areas will be investigated in this section. First, gate leakage in contemporary MOS transistors will be studied and compared to measurements. Emphasis is put on the distinction between the different sources of the tunneling current; namely, the region below the gate and the region near the drain and source extensions.

Device engineers commonly rely on gate leakage measurements of turned-off devices to evaluate the power consumption of CMOS circuits. This may lead to erroneous results because for turned-on devices, hot-carrier tunneling prevails that may exceed the turned-off tunneling current. Models that are based on simplified assumptions of the carrier energy distribution function fail to predict gate leakage in such cases.

Advanced CMOS devices will use alternative dielectric materials as gate dielectrics. However, a pronounced trade-off between the height of the energy barrier and the dielectric permittivity exists. This makes the use of optimization necessary to find the optimum layer composition. Furthermore, alternative dielectrics are not ideal insulators but contain defects that give rise to trap-assisted tunneling. As a state-of-the-art example, tunneling in $ZrO_2$-based MOS capacitors will be studied and compared to measurements.

As a second important application area, nonvolatile memories will be studied. Unlike MOS transistors, nonvolatile memory devices represent an application where tunneling is not a spurious effect, but crucial for the device functionality. After a short review of nonvolatile memory technology, the tunneling current of conventional EEPROMs and advanced structures will be studied. In contrast to these devices, SONOS (silicon-oxide-nitride-oxide-silicon) EEPROM devices store the charge not on an isolated contact, but in a layer of trap-rich dielectric.

Recent efforts to reduce the charging time of nonvolatile memory devices resulted in multibarrier tunneling devices and EEPROMs with asymmetrically layered tunnel dielectrics. The operation of these devices will briefly be described at the end of this chapter. All simulations are performed using the device simulator Minimos-NT [141].

### 3.1. Tunneling in MOS Transistors

The gate leakage current in contemporary MOS transistors poses a major problem for further device scaling. This section describes simulation results of MOS transistors, outlines the effect of various device parameters, shows how to account for hot-carrier tunneling in turned-on devices, and elaborates on the use of alternative dielectric materials to replace $SiO_2$ as a gate dielectric. First, however, the tunneling paths in MOS transistor structures will be reviewed.

### 3.1.1. Tunneling Paths in MOS Transistors

Tunneling in an MOS transistor, as shown in the left part of Fig. 19, basically can be separated into a path between the gate and the channel and a path between the gate and the source and drain extension areas [142]. Tunneling in the source and drain extension areas can exceed tunneling in the channel by orders of magnitude. This is related to two effects: First, instead of $n$-$p$ or $p$-$n$ tunneling, $n$-$n$ or $p$-$p$ tunneling prevails. Second, the potential difference and thus the bending of the energy barrier is high. This increased tunneling current in the source and drain extension areas can be a serious problem if measurements are performed on long-channel MOSFETs to characterize their short-channel pendants, because the edge tunneling currents exceed the channel tunneling current by orders of magnitude. Furthermore, there is a fundamental difference between tunneling in MOS transistors and MOS capacitors [1, 143]. In contrast to MOS transistors, MOS capacitors, which are biased in strong inversion, cannot supply the amount of carriers as predicted by the tunneling model. This effect is termed *substrate-limited* tunneling, because the tunneling current is limited by the generation rate in the substrate. In the channel of an inverted MOS transistor, on the other hand, carriers can always be supplied by the source and drain contacts. This effect is depicted in the right part of Fig. 19.

### 3.1.2. Channel Tunneling

In this section, the effects of various device parameters on the gate leakage of MOS capacitors are studied. This is equivalent to tunneling in MOS transistors if only channel tunneling ($n$-$p$ or $p$-$n$) is considered, and the source, drain, and bulk contacts are grounded. The parameters investigated are

- the doping of the polysilicon gate contact,
- the doping of the substrate,
- the thickness of the dielectric layer,
- the barrier height of the dielectric,
- the carrier mass in the dielectric,
- the dielectric permittivity, and
- the lattice temperature.

The typical shape of the gate current density in turned-off $n$MOS and $p$MOS devices is depicted in Fig. 20. A $SiO_2$ gate dielectric thickness of 2 nm and an acceptor or donor doping of $5 \times 10^{17}$ cm$^{-3}$ and polysilicon gates was chosen. In the $n$MOS device, the majority electron tunneling current always exceeds the hole tunneling current due to the lower electron mass and barrier height (3.2 eV instead of 4.65 eV for holes). In the $p$MOS capacitor, however, the majority hole tunneling exceeds electron tunneling only for negative and low positive bias. For positive bias, the conduction band electron current again dominates due to its much lower barrier height [144].



**Figure 19.** The different tunneling paths (channel tunneling, source and drain extension tunneling) in a MOS transistor (left). In a MOS transistor biased in inversion (right), tunneling electrons are supplied from the source and drain reservoirs, which is not possible in a MOS capacitor.

**Figure 20.** Channel tunneling regions in an $n$MOS (left) and a $p$MOS (right). The insets show the approximate shape of the band edge energies.

### 3.1.2.1. Effect of the Polysilicon Gate Doping on the Channel Tunneling

Highly doped polysilicon is used as material for the gate contact to allow adjustable work functions and realize CMOS circuits. Figure 21 shows the electron and hole tunneling current density for different doping of the polysilicon gate contact. In the $n$MOS, gate leakage generally



**Figure 21.** Electron (left) and hole (right) current density in an $n$MOS (top) and a $p$MOS (bottom) with different doping of the polysilicon gate. Substrate doping is $10^{18}$ cm$^{-3}$; dielectric thickness is 2 nm.

increases with increasing doping of the polysilicon gate because tunneling current is dominated by electrons. In the $p$MOS. a higher polysilicon doping leads to reduced electron tunneling current and increased hole tunneling current. The effect on the overall leakage depends on the doping and the gate bias.

**3.1.2.2. Effect of the Substrate Doping on the Channel Tunneling**   Figure 22 shows the electron and hole tunneling current density for different doping of the substrate. With increasing substrate doping, the majority tunneling component (electrons in the $n$MOS, holes in the $p$MOS) is reduced in both the $n$MOS and $p$MOS devices, whereas the minority component increases.

**3.1.2.3. Effect of the Dielectric Thickness on the Channel Tunneling**   The physical thickness of the dielectric has the largest impact on the gate current density, as shown in Fig. 23. Increasing the gate dielectric thickness by 0.4 nm leads to a decrease of all tunneling current components by several orders of magnitude.

**3.1.2.4. Effect of the Barrier Height on the Channel Tunneling**   The main parameter, besides the thickness of the dielectric. influencing tunneling current is the height of the energy barrier. The influence of this parameter is depicted in Fig. 24. Different dielectric materials strongly differ in their work function difference to silicon. It must be distinguished between the barrier height for electrons and for holes. The most frequently used dielectric material $SiO_2$ has an electron barrier height of about 3.2 eV and a hole barrier height of



**Figure 22.** Electron (left) and hole (right) current density in an $n$MOS (top) and a $p$MOS (bottom) with different doping of the substrate. Gate polysilicon doping is $5 \times 10^{19}$ cm⁻³; dielectric thickness is 2 nm.

**Figure 23.** Electron (left) and hole (right) current density in an *n*MOS (top) and a *p*MOS (bottom) with different thickness of the dielectric layer. Gate polysilicon doping is $5 \times 10^{20}$ cm $^{-3}$; substrate doping is $5 \times 10^{18}$ cm $^{-3}$.

approximately 4.6 eV. The measurement of these material parameters is difficult, and values in the available literature vary widely (see Section 3.1.5).

### 3.1.2.5. Effect of the Carrier Mass in the Dielectric on the Channel Tunneling

Being the parameter with the highest uncertainty, the electron and hole mass in the dielectric is commonly used as a fitting parameter to reproduce measurements. Its influence on the gate current density is shown in Fig. 25. An increase in the carrier mass by 0.1 $m_0$ leads to a reduction in the gate current density by about a factor of 10. It must, of course, be held in mind that with the approaches described so far, tunneling is described by a single value for the carrier mass. Its use as a fitting parameter may thus well be justified. Recent investigations, however, report an increase of the electron mass with reducing thickness of the dielectric layer, which is backed by measurements and tight-binding band structure calculations [145–147].

### 3.1.2.6. Effect of the Dielectric Permittivity on the Channel Tunneling

The permittivity of the dielectric layer influences the tunneling current density in two ways: First, the shape of the energy barrier—and thus the transmission coefficient—changes. Second, the inversion charge—and thus the band edge energy—in the channel is affected. The effect of varying dielectric permittivity is shown in Fig. 26. Especially in the low-bias regime, a higher permittivity strongly increases the gate current density.

### 3.1.2.7. Effect of the Lattice Temperature on the Channel Tunneling

The lattice temperature enters the gate tunneling current via the electron energy distribution functions in the polysilicon gate and in the channel. The transmission coefficient, being based on

**Figure 24.** Effect of the electron and hole barrier height on electron tunneling current (left) and hole tunneling current (right) in an nMOS (top) and a pMOS (bottom) with 2-nm dielectric thickness. $10^{20}$ cm$^{-3}$ polysilicon, and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.

quantum-mechanical reasoning alone. is not affected by the lattice temperature. However, the supply function depends on the lattice temperature. The impact on the gate current density is shown in Fig. 27. Rising temperature increases the tunneling current density in all cases.

**3.1.2.8. Comparison to Measurements** Because almost all available measurements of gate leakage in MOS devices are performed on turned-off MOS transistors, a comparison with measurements will be given before turned-on devices are investigated in Section 3.1.4. The Tsu–Esaki model with an analytical WKB transmission coefficient is in good agreement with recently reported data for devices with different gate lengths and bulk doping [1, 142] as shown in Fig. 28 for nMOS (left) and pMOS devices (right) [148]. It can be seen that the gate current density can be reproduced over a wide range of dielectric thicknesses with a single set of physical parameters. Additional measurements have been performed on MOSFETs with a gate dielectric thickness of 1.5 nm (see the lower part of Fig. 28) and compared with the results of other simulators (UTQUANT [149] and MEDICI [150]). Under inversion condition the fit is not perfect, whereas under accumulation the measurements can be reproduced well. Note that with UTQUANT, the low-bias tunneling current cannot be reproduced, and MEDICI completely failed for the pMOS device.

**3.1.2.9. Validity of Compact Models** Because the computational effort for the numerical integration in Tsu–Esaki's formula or the evaluation of the quasi-bound states is numerically expensive, it is reasonable to ask if compact models can describe tunneling, at least for

**Figure 25.** Effect of the carrier mass on electron tunneling current (left) and hole tunneling current (right) in an *n*MOS (top) and a *p*MOS (bottom) with 2-nm dielectric thickness, $10^{20}$ cm$^{-3}$ polysilicon. and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.

single-layer dielectrics. The compact tunneling models outlined in Section 2.7 are compared in Fig. 29 for a symmetrical metal-dielectric-metal structure (left) and for an *n*MOS structure with 3-nm dielectric thickness (right). For the metal-dielectric-metal structure, Schuegraf's model yields almost the same results as the computationally much more expensive Tsu–Esaki model. The Fowler–Nordheim model delivers correct values only for high bias. It is thus only applicable to describe high-field transport through gate dielectrics, like program and erase cycles in EEPROM devices. For the MOS structure in the right part of Fig. 29, the Schuegraf model fails to describe the tunneling current density at low bias. For high bias, however, it may be used to provide an estimation of the gate current. The Fowler–Nordheim model totally fails for this application. Furthermore, the Fowler–Nordheim model shows the minimum gate current at minimum electric field in the dielectric, and not for the minimum gate bias.

### 3.1.3. Source/Drain Extension Tunneling

In the following examples, the same devices as in Section 3.1.1 are investigated, but this time only the tunneling current in the source and drain extension areas (*n-n* or *p-p*) is taken into account. Because the barrier height, carrier mass, and dielectric thickness shows the same impact on the gate current density as for the case of channel tunneling, the corresponding figures are omitted.

#### 3.1.3.1. Effect of the Polysilicon Gate Doping on the Source and Drain Extension Tunneling   Figure 30 shows the effect of the doping concentration in the polysilicon gate on the extension region gate current density. Increasing the polysilicon doping leads to a

**Figure 26.** Effect of the dielectric permittivity $\kappa/\kappa_0$ on electron tunneling current (left) and hole tunneling current (right) in an $n$MOS (top) and a $p$MOS (bottom) with 2-nm dielectric thickness. $10^{20}$ cm$^{-3}$ polysilicon, and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.

slight increase of the main tunneling component and to a strong decrease of the minority tunneling component in both nMOS and pMOS devices.

**3.1.3.2. Effect of the Substrate Doping Concentration on the Source and Drain Extension Tunneling** Figure 31 shows the effect of the substrate doping concentration on the extension region gate current density. Similar to the polysilicon gate doping, a higher substrate doping leads to increased majority and decreased minority tunneling current.

**3.1.3.3. Effect of the Dielectric Permittivity on the Source and Drain Extension Tunneling** Figure 32 shows the effect of the dielectric permittivity on the extension region gate current density. In contrast to the channel-tunneling case, the low-bias regime is not influenced by the permittivity. Furthermore, the influence on the majority tunneling current component depends on the bias: The electron tunneling component in the $n$MOS decreases for negative bias and increases for positive bias. The hole tunneling component in the $p$MOS shows exactly the inverse trend.

**3.1.3.4. Effect of the Lattice Temperature on the Source and Drain Extension Tunneling** Figure 33 shows the effect of the temperature on the extension region gate current density. Especially the minority carriers (holes in the $n$MOS, electrons in the $p$MOS) show strongly increased tunneling current with higher temperature. Unlike in the channel tunneling case, the majority tunneling component is hardly influenced by the temperature.

**Figure 27.** Effect of the lattice temperature on electron tunneling current (left) and hole tunneling current (right) in an $n$MOS (top) and a $p$MOS (bottom) with 2-nm dielectric thickness, $10^{20}$ cm$^{-3}$ polysilicon, and $5 \times 10^{18}$ cm$^{-4}$ substrate doping.

## 3.1.4. Hot-Carrier Tunneling in MOS Transistors

It has been shown in Section 2.3 that the distribution function in the channel of a turned-on MOS transistor heavily deviates from the shape implied by a Fermi–Dirac or Maxwellian distribution. A model for the non-Maxwellian shape of the distribution function was presented that accurately reproduced the carrier energy distribution along the channel.

To check the impact of this wrong high-energy behavior, the integrand of the Tsu–Esaki formula, namely the expression $TC(\ell)N(\ell)$, has been evaluated for a standard device, as shown in the left part of Fig. 34, and compared to Monte Carlo results [151, 152]. The simulated device had a gate length of 100 nm and a gate dielectric thickness of 3 nm. Though at low energies, the difference between the non-Maxwellian distribution function (28) and the heated Maxwellian distribution (24) seems to be negligible, the amount of overestimation of the incremental gate current density for the heated Maxwellian distribution reaches several orders of magnitude at 1 eV and peaks when the electron energy exceeds the barrier height. This spurious effect is clearly more pronounced for points at the drain end of the channel where the electron temperature is high. The non-Maxwellian shape of the distribution function, indicated by the full line, reproduces the Monte Carlo results very well.

The region of high electron temperature is confined to only a small area near the drain contact, as shown in the right part of Fig. 34, where the gate current density along the channel is compared to Monte Carlo results. At the point of the peak electron temperature, which is located at approximately $x = 0.8L_g$, the heated Maxwellian approximation overestimates

**Figure 28.** Comparison of simulations using different simulators with measurements of $n$MOS (left) and $p$MOS (right) devices [1, 142, 148].



**Figure 29.** Compact models for a metal-dielectric-metal structure (left) and an $n$MOS structure (right; literature values from [142]).

**Figure 30.** Effect of the polysilicon doping on the electron tunneling current (left) and the hole tunneling current (right) in the source and drain extension region of an $n$MOS (top) and a $p$MOS (bottom) with 2-nm dielectric thickness and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.

the gate current density by a factor of almost $10^6$. It will therefore have a large impact on the total gate current density. The cold Maxwellian approximation underestimates the gate current density in this region, whereas the non-Maxwellian distribution correctly reproduces the Monte Carlo results.

The non-Maxwellian shape yields excellent agreement, whereas the heated Maxwellian approximation substantially overestimates the gate current density especially near the drain region. Instead of the heated Maxwellian distribution, it appears to be better to use a cold Maxwellian distribution in that regime because it leads to a comparably low underestimation of the gate current density.

The effect of hot-carrier tunneling on the total gate current of the devices is shown in Fig. 35. In the left part of this figure, the gate current density for a 0.5-$\mu$m turned-on MOSFET with a dielectric thickness of 4 nm is shown as a function of the gate bias. Results from Monte Carlo simulations are also shown in this figure. For low gate voltages ($V_{GS} < V_{DS}$), the peak electric field in the channel increases with increasing gate bias. The electron temperature is high, and the heated Maxwellian approximation massively overestimates the total gate current. If the gate bias exceeds the drain-source voltage, however, the peak electric field in the channel is reduced [153]. Therefore, for $V_{GS} > V_{DS}$, the electron temperature reduces with increasing gate bias, and the heated Maxwellian approximation delivers correct results. The non-Maxwellian model (28) delivers correct results for all gate voltages.

Tunneling Models for Semiconductor Device Simulation



Figure 31. Effect of the substrate doping on the electron tunneling current (left) and the hole tunneling current (right) in the source and drain extension region of an $n$MOS (top) and a $p$MOS (bottom) with 2-nm dielectric thickness and $5 \times 10^{20}$ cm$^{-3}$ polysilicon doping.

The question remains if the hot-carrier tunneling current strongly depends on the gate length of the device. In the right part of Fig. 35, the gate current is given as a function of the gate length for different gate dielectric thicknesses (2.2 nm–3.0 nm). Again, Monte Carlo simulation results are used as reference. It can be seen that the heated Maxwellian distribution delivers correct results only for large gate lengths, whereas it totally fails for smaller devices. The use of a cold Maxwellian distribution, on the other hand, underestimates the gate current only slightly and seems to be the better choice if accurate modeling of the device physics is not that important or only a quick estimation is asked for. The non-Maxwellian model correctly reproduces the Monte Carlo results for all gate lengths and gate dielectric thicknesses.

### 3.1.5. Alternative Dielectrics for MOS Transistors

The further reduction of device dimensions makes the introduction of alternative dielectric materials necessary. Because none of the possible materials forms a native oxide on silicon, a thin interfacial layer of SiO$_2$ cannot be avoided. Thus, a two-layer band edge diagram is commonly assumed as depicted in Fig. 36 [154]. A wide variety of high-$\kappa$ materials can be considered as alternative dielectrics. However, several points must be considered when evaluating these materials:

(1) The dielectric permittivity $\kappa$.
(2) The barrier height for electrons q$\Phi_e$ and holes q$\Phi_h$ on silicon. These values are equivalent to the band edge offsets $\Delta E_c$ and $\Delta E_v$.

**Figure 32.** Effect of the dielectric permittivity $\kappa/\kappa_0$ on the electron tunneling current (left) and the hole tunneling current (right) in the source and drain extension region of an $n$MOS (top) and a $p$MOS (bottom) with 2-nm dielectric thickness and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.

(3) The thermodynamic stability of the dielectric material on silicon: The material must withstand all following processing steps.

(4) The quality of the interfaces: High interface roughness may cause increased scattering, which reduces the carrier mobility in the channel.

(5) The trap concentration, which leads to trap-assisted tunneling.

(6) The feasibility and integrability of the deposition method in the fabrication process.

Only the permittivity, the trap concentration, and the barrier heights influence the tunneling current. When looking at the barrier height and permittivity of various dielectrics in Table 4, one notices a strong trade-off between the barrier height and the dielectric permittivity: dielectrics with a high energy barrier have a low permittivity and vice versa; see Figs. 37 and 38. Hence, optimization becomes necessary to find the optimum material.

Choosing the material parameters from Table 4, the gate current density can be computed as a function of the gate bias [155]. It is commonly assumed that an underlying layer of $SiO_2$ cannot be avoided—or is even deliberately introduced to achieve a lower trap density at the interface to silicon. Thus, an underlying $SiO_2$ layer with a thickness of 0.5 nm was assumed. The thickness of the high-$\kappa$ layer was adjusted so that the effective oxide thickness (EOT) remains unchanged at 1 nm. The gate current density is shown in the left part of Fig. 39 as a function of the gate bias for different material combinations. The commonly assumed limit of 1 Acm$^{-2}$ gate leakage is also indicated. Both $SiO_2$ and $Si_3N_4$ show a much too high leakage, whereas $Ta_2O_5$, $ZrO_2$, and $HfO_2$ stay below 1 Acm$^{-2}$ at $V_{GS} = 1$ V. Due to the low

**Figure 33.** Effect of the lattice temperature on the electron tunneling current (left) and the hole tunneling current (right) in the source and drain extension region of an *n*MOS (top) and a *p*MOS (bottom) with 2-nm dielectric thickness and $5 \times 10^{18}$ cm$^{-3}$ substrate doping.



**Figure 34.** Integrand of Tsu–Esaki's equation (left) and gate current density along the channel (right) of a MOSFET with 100-nm gate length and 3-nm gate dielectric thickness [151, 152].

**Figure 35.** Gate current for different values of the gate bias (left). Dependency of the total gate current on the gate length (right) [151, 152].

conduction band offset, $TiO_2$ shows an especially pronounced current increase for positive gate bias.

To assess the material parameters necessary to stay below a specific maximum gate current density, the gate current has been calculated as a function of the conduction band offset and dielectric permittivity as shown in the right part of Fig. 39. Because it is often not possible to vary the thickness of the underlying $SiO_2$ layer, it was again fixed at 0.5 nm and the high-$\kappa$ thickness was adjusted to reach an EOT of 1.5 nm. The gate current density was evaluated at a fixed bias point of $V_{GS} = 1.5$ V and $V_{DS} = 0$ V. The current density decreases strongly with increasing conduction band offset. Increasing the value of the dielectric permittivity $\kappa$ also strongly reduces the leakage current due to the higher physical stack thickness. However, materials with a conduction band offset below 1 eV never reach acceptable gate current densities.

It may be asked which thickness of the high-$\kappa$ layer is necessary to achieve a certain gate current density. In the left part of Fig. 40, the gate current density is shown for an effective oxide thickness ranging from 0.5 nm to 2.0 nm as a function of the high-$\kappa$ layer thickness. Again, the stack consists of an underlying 0.5 nm layer of $SiO_2$ and the simulations are



**Figure 36.** Schematic of a band energy diagram of a stacked dielectric consisting of a thin underlying interface layer and a thick layer of a high $\kappa$ material with higher dielectric permittivity, but lower barrier height.

**Table 4.** Band gap energy and conduction band offset of various dielectric materials.

| | $\kappa/\kappa_0$ (1) | Band Gap $\mathscr{E}_g$ (eV) | Conduction Band Offset $\Delta\mathscr{E}_c$ (eV) | Valence Band Offset $\Delta\mathscr{E}_v$ (eV) | Reference |
|---|---|---|---|---|---|
| SiO$_2$ | 3.9 | 9.00 | 3.00 | 4.90 | [198] |
| | 3.9 | 9.00 | 3.50 | 4.40 | [193] |
| | 3.9 | 9.00 | 3.15 | 4.75 | [199] |
| | 3.9 | 8.90 | 3.20 | 4.60 | [200] |
| | | 9.00 | 3.50 | 4.40 | [201, 202] |
| | 3.9 | 9.00 | 3.00 | 4.90 | [41] |
| Si$_3$N$_4$ | 7.5 | 5.00 | 2.00 | 1.90 | [198] |
| | 7.6 | 5.00–5.30 | 2.40 | 1.50–1.80 | [193] |
| | 7.9 | 5.30 | 2.40 | 1.80 | [199] |
| | 7.0 | 5.10 | 2.00 | 2.00 | [200] |
| | | 5.30 | 2.40 | 1.80 | [201, 202] |
| | 7.5 | 5.00 | 2.00 | 1.90 | [41] |
| Ta$_2$O$_5$ | 25.0 | 4.40 | 1.40 | 1.90 | [41, 198] |
| | 23.0–25.0 | 4.40 | 0.30 | 3.00 | [193] |
| | 25.0 | 4.40 | 0.36 | 2.94 | [199, 202] |
| | 26.0 | 4.50 | 1.00–1.50 | 1.90–2.40 | [200] |
| | | 4.40 | 0.36 | 2.94 | [201] |
| TiO$_2$ | 40.0 | 3.50 | 1.10 | 1.30 | [41, 198] |
| | 39.0–110.0 | 3.00–3.27 | 0.00 | 1.90–1.97 | [193] |
| | 80.0–170.0 | 3.05 | 0.00 | 1.95 | [199] |
| | 80.0 | 3.50 | 1.20 | 1.20 | [200] |
| | | 3.05 | 0.00 | 1.95 | [201] |
| Al$_2$O$_3$ | 9.0 | 8.70 | 2.80 | 4.80 | [200] |
| | 8.0–9.0 | 8.8–9.00 | 2.78–2.80 | 4.92–5.10 | [193] |
| | 9.5–12.0 | 8.8 | 2.80 | 4.90 | [199] |
| | | 8.80 | 2.80 | 4.90 | [201] |
| | 10.0 | 8.80 | 2.80 | 4.90 | [202] |
| ZrO$_2$ | 23.0 | 5.80 | 1.40 | 3.30 | [202] |
| | 25.0 | 7.80 | 1.40 | 5.30 | [198, 200] |
| | 22.0–25.0 | 5.00–5.80 | 1.40 | 2.50–3.30 | [193] |
| | 12.0–16.0 | 5.70–5.80 | 1.40–1.50 | 3.10–3.30 | [199] |
| | | 5.80 | 2.50 | 2.20 | [201] |
| HfO$_2$ | 25.0 | 5.70 | 1.50 | 3.10 | [198, 200] |
| | 22.0–40.0 | 6.00 | 1.50 | 3.50 | [193] |
| | 16.0–30.0 | 4.50–6.00 | 1.50 | 1.90–3.40 | [199] |
| | | 6.00 | 1.50 | 3.40 | [201] |
| | 20.0 | 6.00 | 1.50 | 3.40 | [202] |
| Y$_2$O$_3$ | 15.0 | 5.60 | 2.30 | 2.20 | [200] |
| | 11.3–18.0 | 5.50–6.00 | 1.30 | 3.10–3.60 | [193] |
| | 4.4 | 6.00 | 1.30 | 3.60 | [201] |
| | 15.0 | 6.00 | 2.30 | 2.60 | [202] |
| ZrSiO$_4$ | 12.6 | 6.00 | 1.50 | 3.40 | [193] |
| | | 4.50 | 0.70 | 2.70 | [199] |
| | 3.8 | 6.00 | 1.50 | 3.40 | [201] |
| | | 6.00 | 1.50 | 3.40 | [202] |

performed at a fixed bias point of $V_{GS} = 1.5$ V and $V_{DS} = 0$ V. In this plot, the curves are only drawn for an EOT of 0.5 nm–2.0 nm, and conduction band offsets of $q\Phi_c = 1$ eV to $q\Phi_c = 3$ eV have been considered. For a conduction band offset of 1 eV, large high-$\kappa$ thicknesses are necessary to reduce the leakage. Such large stacks may pose problems due to fringing fields from the drain contact, which reduce the threshold voltage of the device.

The trade-off between the dielectric permittivity and the conduction band offset gives rise to further effects as shown in the right part of Fig. 40. If the EOT has to be held at a fixed value, an increase of the SiO$_2$ layer thickness causes a reduced thickness of the high-$\kappa$ layer. This is shown for different values of the permittivity ($\kappa = 8.0$ to $\kappa = 24.0$). So, the total stack thickness may be larger than 8 nm for $\kappa = 24$, or as small as 1.5 nm

**Figure 37.** Trade-off between electron barrier height (left) or hole barrier height (right) and the permittivity of various dielectric materials [41, 193, 198–202].

if only $SiO_2$ is used. Such a reduction of the total stack thickness, however, has no clear effect on the leakage. It may cause the gate current density at a specific bias point to stay constant, increase, or even decrease depending on the material parameters. For example, the gate leakage for a material with $\kappa = 24$ and a conduction band offset of 1 eV shows the maximum leakage at a $SiO_2$ layer thickness of approximately 0.8 nm. Therefore, a clear statement about the optimum thickness of the interface layer obviously depends on the material parameters.



**Figure 38.** Conduction and valence band edges of various dielectric materials compared to silicon [41, 193, 198–202].

**Table 5.** Layer thicknesses and effective oxide thickness of metal organic chemical vapor deposition-deposited $ZrO_2$ layers in nanometers, after Harasek [156].

| Layer Thickness | $t_{int}$ | $t_{high-k}$ | EOT |
|---|---|---|---|
| 6.9 | 0.75-2.0 | 6.15-4.9 | 2.0 |
| 12.7 | 0.3-1.0 | 12.4-11.7 | 3.0 |

the figure also shows the gate current for a 2-nm and a 3-nm $SiO_2$ layer (dotted lines). As expected, the measured current density is lower than for the $SiO_2$ counterparts. However, the Tsu–Esaki model cannot reproduce the measurements as it yields tunneling currents orders of magnitude lower than the measurements. This indicates the presence of strong trap-assisted tunneling due to a high trap concentration in the dielectric layer. By assuming a Frenkel–Poole-like conduction through the dielectric layer, the measurements could be reproduced (full lines). Note that in previous studies [156], tunneling through $ZrO_2$ layers fabricated by magnetron sputtering could be reproduced without considering trap-assisted tunneling. That indicates the presence of a high trap concentration due to the MOCVD process, in contrast to the sputtering process.

To clarify the trap energy level and concentration, the step response of the MOS capacitors has been measured as shown in the right part of Fig. 41 for the 12.7-nm $ZrO_2$ layer annealed in reducing conditions (forming gas) and the 6.9-nm layer annealed under oxidizing conditions [158]. The gate voltage is turned off after being fixed at a value of 2.5 V, and the resulting gate current is measured over time. The transient gate current exceeds the static gate current by orders of magnitude and decays very slowly. This behavior can be explained assuming defects in the dielectric layer [159]. Using the trap-assisted tunneling model outlined in Section 2.8.2, a trap energy level of 1.3 eV below the $ZrO_2$ conduction band edge, a trap concentration of $4.5 \times 10^{18}$ cm$^{-3}$, and an energy loss of 1.5 eV have been found. For the dielectric layer annealed under oxidizing conditions, a trap concentration of $4 \times 10^{17}$ cm$^{-3}$ was found.

To predict the performance of devices based on $ZrO_2$ dielectrics, a well-tempered MOS-FET as described in Ref. [160] with an effective channel length of 50 nm has been simulated. EOT thicknesses of 2-nm and 3-nm $SiO_2$ and respective $ZrO_2$ layers have been considered. The left part of Fig. 42 depicts the conduction band edge in the channel for different gate-source voltages. It can be seen that the barrier is slightly lower for the $ZrO_2$ layer at $V_{GS} = 1.2$ V, whereas it is strongly reduced at $V_{GS} = 0.1$ V, which is due to the pronounced fringing fields from the drain contact.



**Figure 41.** Stationary (left) and transient (right) gate current measurements of the $ZrO_2$ layers performed by Harasek [157] compared with simulations [158].

**Figure 42.** Well-tempered MOSFET conduction band edge along the channel for SiO₂ and ZrO₂ dielectrics (left). Influence of the dielectric trap concentration on the MOSFET threshold voltage (right) [158].

An additional topic of interest for high-$\kappa$ dielectrics is the influence of trapped charges in the high-$\kappa$ layer on the threshold voltage of the device. The trap concentration in the $ZrO_2$ layer was increased from $10^{15}$ cm$^{-3}$ to $10^{19}$ cm$^{-3}$ with full trap occupancy in the dielectric layer. It can be seen in the right part of Fig. 42 that the threshold voltage strongly increases with rising trap concentration. This effect is therefore contrary to the decrease of the threshold voltage due to fringing fields described above.

## 3.2. Tunneling in Nonvolatile Memory Devices

Tunneling effects are crucial not only for MOS transistors but also for nonvolatile semiconductor memory devices. In contrast to volatile memory devices, they retain the stored information without external power supply. Nonvolatile memory (NVM) devices can be read and programmed like random-access memory (RAM) devices, have a low power consumption, are mechanically robust, and offer the possibility of large-scale integration. They constitute about 10% of the total semiconductor memory market [161]. However, simulation of such devices is often carried out using simplified compact models [162–167]. For the case of stacked gate dielectrics or hot electron injection, such models do not capture the device physics and can reproduce measured data only on a fit-formula level. In this section, some examples of conventional EEPROM and alternative devices will be studied using the tunneling models described above.

### 3.2.1. Conventional EEPROM Devices

The basic operating principle of an EEPROM was presented by Kahng and Sze in 1967 at Bell Laboratories [168]. The device consists of a control gate and a floating gate on top of a conventional MOS transistor. A thin tunnel dielectric separates the floating gate from the channel. It must be thick enough to allow up to $10^5$ writing and erasing cycles without breakdown—common thicknesses are 6–8 nm. Applying a high positive voltage (about 8–12 V) on the control gate raises the potential of the floating gate by capacitive coupling. The high electric field in the tunnel dielectric ($\approx 10^9$ V/m) leads to Fowler–Nordheim tunneling of electrons from the substrate to the floating gate. The charge on the floating gate changes the threshold voltage of the underlying MOS transistor and is retained even if the control gate voltage is removed. A retention time of 10 years is required for consumer applications such as memory cards. Though EEPROM cells offer random access for writing and erasing of individual bits, Flash cells can be programmed selectively but erased only at once. This has the advantage of lower cell size. Due to the high electric field in the dielectric, degradation or even breakdown of the dielectric is a major concern. A comprehensive survey of NVM technology is given in Refs. [169] and [170].

### 3.2.1.1. Static SILC in EEPROMs

The speed of the programming and erasing process is one of the main figures of merit of an EEPROM cell. Therefore, strong electric fields are applied at the control gate to allow Fowler–Nordheim tunneling of carriers during programming and erasing cycles. However, due to this repeated high-field stress, trap centers in the dielectric are formed, which allow trap-assisted tunneling at low fields and thus reduce the retention time of the devices. This additional current at low bias is known as stress-induced leakage current (SILC) and represents one of the major reliability concerns in contemporary EEPROM devices [112, 135]. In the left part of Fig. 43, measured SILC after different stress times for a MOS capacitor with a dielectric thickness of 5.5 nm is shown [105]. The trap-assisted tunneling model outlined in Section 2.8.2 yields excellent agreement with the measured data if the trap concentration is used as a fitting parameter dependent on the stressing time (the model parameters are stated in the figure caption). The transition from the region of mainly trap-assisted tunneling for $V_{GS} < 5$ V to the region of Fowler–Nordheim tunneling for $V_{GS} > 5$ V is clearly visible. The right part of Fig. 43 shows the trap occupancy $f_T$ across the gate dielectric of a MOS capacitor using the gate voltage as parameter. The regions near the gate (right) and near the substrate (left) are only sparsely occupied. Near the gate, the emission time is much smaller than the capture time, and near the substrate, the trap energy lies above the electron energy in the cathode. Some of the trapped electrons face a triangular barrier for the emission process, giving rise to an additional peak in the trap occupancy near the gate side (the anode) of the dielectric. This is due to the wave function interference in the Fowler–Nordheim region (the oscillations are also observed in the emission time of the traps shown in Fig. 17).

### 3.2.1.2. Transient SILC in EEPROMs

It has been shown that the transient trap-assisted tunneling current can be described by a rate equation that gives rise to an exponential behavior of the tunneling current over time; see Section 2.8.2.4. The left part of Fig. 44 shows measurements of the gate current density of MOS capacitors as a function of time with dielectric thicknesses of 8.5 nm and 13.0 nm compared to simulations [104, 171]. Initially, the traps are empty, which can be achieved by applying flat band conditions. At $t = 0$ s, the gate voltage is turned on ($-5.8$ V and $-8.3$ V for the thinner and the thicker dielectric, respectively) and the traps are filled according to their specific capture and emission time constants. This charging current consists of an emission and a capture current, which may exceed the steady-state current by orders of magnitude. A good fit to the measured data can be achieved using the trap parameters indicated in the figure caption.



Figure 43. Comparison of simulations with measurements of a MOS capacitor with a dielectric thickness of 5.5 nm [105, 171] is shown on the left. The trap energy is 2.7 eV, the phonon energy 130 meV, and the Huang–Rhys factor is 10. The trap concentration was set to $9 \times 10^{17}$ cm$^{-3}$, $10^{17}$ cm$^{-3}$, $3 \times 10^{16}$ cm$^{-3}$, and $3 \times 10^{15}$ cm$^{-3}$ to fit the measurements (from top to bottom). The trap occupancy across the gate dielectric at different gate voltages is shown on the right.

The right part of Fig. 44 shows the gate current of a MOS capacitor for an applied rectangular pulse with a frequency of 100 kHz assuming initial flat band conditions. It can be seen that the time constants of the trap filling and emptying processes are not equal but depend on the applied voltage, as different voltages lead to different capture and emission times. The spikes in this figure are due to the sudden voltage change, whereas the trap concentration remains constant: In the transition from 3.0 V to 3.5 V, the barrier shape changes suddenly, and traps are rapidly emptied. Traps near the cathode are filled, and it takes several microseconds until the new steady-state is reached. Thus, dielectric materials that have such a high trap concentration may lead to considerable problems for high-frequency applications.

For EEPROM devices, the charging and discharging characteristics are crucial: Programming and erasing should happen as fast as possible; therefore, high voltages are applied. The discharging current over time, on the other hand, determines the retention time and must be very low. Furthermore, the programming and erasing pulses must be carefully optimized to avoid over-erase, as the tunnel current density for positive and negative voltages on the floating gate is not equal. This is frequently addressed in the literature [172, 173].

### 3.2.2. Alternative Nonvolatile Memory Devices

Strong efforts are undertaken to improve the standard floating-gate EEPROM cell in terms of integration density, endurance, reliability, program time, erase time, and retention time. EEPROM devices with a tunnel window near the drain contact have been introduced to reduce the charge loss from the floating gate and thus reach higher retention time. However, due to the small area of the tunnel window, high voltages have to be used at the drain contact, which again reduces cell reliability.

Recently, Caywood et al. proposed a device structure where nonselected cells are isolated from the drain and source contacts by two additional side gates [174]. In this device, electrons tunnel from the inverted channel to the floating gate. The large area reduces programming and erasing time. Furthermore, the capacitive coupling between the control gate and the floating gate is higher than in the standard EEPROM cell, which allows use of lower programming and erasing voltages. No drain-source bias is applied for charging, thus the power consumption is low and the injected electrons are less likely to cause degradation of the dielectric. The control gate functions as a select transistor that isolates unselected cells from the high voltages at the shared source and drain contacts during read and write access of neighboring cells.

In contrast to the reduction of the cell footprint, integration density can also be increased by storing more than one bit on a standard EEPROM cell. This can be achieved by tailoring the programming and erasing pulses in such a way that the threshold voltage falls into one of 4, 8. or 16 voltage ranges. The different threshold voltages can be distinguished by the sensing circuits, resulting in two, three, or four bits that can be stored in the cell. However. charge loss must be extremely low over time, and the threshold voltages have to be detected very precisely.

Single-poly devices have been proposed to integrate NVM devices in standard CMOS logic processes, thus enabling an embedded memory. The control gate lies next to the floating gate, and capacitive coupling is achieved by a layer of highly doped silicon. Though such devices can readily be integrated into existing CMOS process flows, they come at the cost of a large footprint.

A different approach to store more than one bit in a single memory cell is to split the floating gate into two separate segments. If a nonuniform doping in the source and drain side of the channel is used, different amounts of charge can be stored in each floating gate. Such device structures are either achieved using separate metallic floating gates [97] or using a layer of trap-rich dielectric [175].

In the following sections, three of the most promising alternative EEPROM devices will be studied in detail. These are

- Quantum dot and trap-rich dielectric based devices: In these devices, charging and discharging is achieved by tunneling of electrons to and from localized trapping centers in the dielectric.
- Multibarrier tunneling devices consist of a floating gate—or memory node—which is separated from the control gate by several thin dielectric layers. By the use of a side gate, the tunneling current through these barriers can be controlled selectively. In contrast to EEPROMs, the tunneling current flows from the floating gate to the control gate and not to the channel. Extremely high $I_{on}/I_{off}$ ratios can be achieved because the tunneling current is controlled by a separate side gate contact.
- Devices where the tunnel dielectric consists of stacked dielectrics that are engineered in such a way that they block tunneling in the off-state but allow strong tunneling in the on-state.

### 3.2.2.1. Nonvolatile Memory Devices Based on Trap-Rich Dielectrics

A SONOS (silicon-oxide-nitride-oxide-silicon) device is a nonvolatile memory where the charge is stored in a layer of trap-rich dielectric material instead of a floating gate as in an EEPROM. Figure 45 shows an example where a layer of $Si_3N_4$ is sandwiched between two layers of $SiO_2$. Electrons tunneling from the substrate are trapped and redistribute themselves in separate trapping centers. This has the advantage that the charge is stored independently in the traps. A leaky path in the tunnel dielectric cannot lead to full charge loss, as it is the case in conventional EEPROM devices. Therefore, reliability and retention time is increased [75, 125, 176–185].

The band diagram along the dielectric of such a device is shown in Fig. 45 for the programming, storing, and erasing processes. By applying a positive voltage at the gate contact, electrons tunnel through the tunnel dielectric into the trap region. The traps are filled with



Figure 45. Conduction band edge in a SONOS device for the programming, storing, and erasing process.

electrons and become negatively charged. Because of the tunnel dielectric, this charge is stored even if the bias is removed. To erase the memory cell, a negative voltage is applied on the gate contact, leading to a reduced potential barrier and a high tunneling current of electrons out of the traps. Important device parameters are the charging and discharging current through the dielectric, the drain current in the on- and off-state, and the retention time.

The trap-assisted tunneling model can be applied to simulate device characteristics of this device, where three layers of $SiO_2$ have been used and the trap concentration and trap energy level in the middle layer was chosen to resemble a layer of silicon nitride. The transient trap occupancy for a discharging process starting from an initial condition of 2 V at the gate contact is shown in Fig. 46. Initially, the traps are filled. Over time, the electrons leak through the lower dielectric into the channel. After $10^9$ s, almost no more charge is stored in the trap-rich dielectric.

### 3.2.2.2. Multibarrier Tunneling Devices

One of the main shortcomings of conventional EEPROM devices is that the current in the on-state and off-state—the programming and leakage currents—flow through the same tunnel dielectric and face the same energy barrier. They cannot be optimized independently: Increasing the thickness of the tunnel dielectric reduces the leakage, but also reduces the on-state current and thus increases the programming time. Multibarrier tunneling devices offer a solution to this problem. Planar localized-electron device memory (PLEDM) cells have been presented by Nakazato et al. in Ref. [186], and promising results have been reported [187–190]. The principle of a PLEDM is to put a PLED transistor (PLEDTR) on top of the gate of a conventional MOSFET, as shown in Fig. 47. The charge on the memory node, which acts as a floating gate, is provided by tunneling of carriers through the PLED transistor, which consists of a stack of $Si_3N_4$ barriers sandwiched between layers of intrinsic silicon. Upper and lower barriers prevent diffusion from the polysilicon contacts, whereas the middle barrier—the central shutter barrier (CSB)—blocks the tunneling current in the off-state. The PLED transistor has two side



Figure 46. Transient trap occupancy in the trap-rich dielectric layer of a SONOS device that is discharged from $t = 0$ s to $t = 10^9$ s.

**Figure 47.** Conduction band edge energy in the PLEDM device.

gates that are separated by a thin dielectric layer. In the on-state, the energy barriers are heavily reduced by the voltage on the side gates, causing a strong tunneling current to flow at the interface to the side gate dielectric. In the off-state, however, the side gates are turned off, and the energy barrier blocks the leakage current. As in a conventional EEPROM, the charge on the memory node is used to control the underlying MOS transistor. Only a small amount of charge has to be added to or removed from the memory node to change the state of the memory cell.

For the simulation of such devices, measurement results for a single $Si_3N_4$ barrier diode [190] have been used to calibrate the model, as shown in Fig. 48 [191]. For calibration, the carrier mass in the dielectric was used as a fit parameter. Electron and hole masses of $0.5~m_0$ and $0.8~m_0$ were found to reproduce the data. The $Si_3N_4$ barrier was modeled with a barrier height of 5 eV and a conduction band offset of 2 eV to the silicon conduction band edge with the relative dielectric permittivity being 7.5.

The effect of the position and size of the central shutter barrier as well as the effect of shrinking the stack width have been investigated. Two cell states have been assumed: an on-state with 3 V applied on the top contact and the side gate, and an off-state with 0.8 V applied on the memory node and 0 V on the side gate. In both states, the charging and discharging current was extracted. The PLEDTR had a stack width of 180 nm and a stack height of 100 nm. The thickness of the upper and lower barriers was set to 2 nm. The left part of Fig. 49 shows the effect of different CSB thicknesses on the on- and off-current of the device. Though the on-current is hardly influenced by the different thicknesses, the off-current is very sensitive to it. Also, the position of the CSB is crucial, because for a CSB located near the memory node, the energy barrier will be reduced in the off-state by the charge on the memory node. If, on the other hand, the CSB is placed near the top contact, the energy barrier is not suppressed in the off-state, and the off-current is much lower. The on-current is also reduced by this effect, but the amount of reduction is much lower as compared to the off-current, due to the fact that the on-current mainly depends on the voltage of the side gate. Thus, the $I_{on}/I_{off}$ ratio increases with the thickness of the central shutter barrier and is highest for a CSB located near the top contact. Such an asymmetry in

Figure 48. PLEDM calibration of the tunneling current density for a single Si₃N₄ layer with 1.5-nm and 2-nm thickness [191]. The measured values are taken from Ref. [190].

the IV characteristics depending on the position of the central shutter barrier has already been observed experimentally [190].

In Ref. [192], the feasibility of very narrow silicon-insulator stacks is shown. This encourages the assumption that a reduction of the stack width is possible. Figure 49 shows the on- and off-currents of the device with a CSB thickness of 10 nm for a stack width of 140 nm down to 20 nm. It can be seen that a reduction of the stack width leads to increasing on-currents and decreasing off-currents. The reason is that the current in the on-state, which mainly flows as a surface current near the side gate, is not reduced by the decreased width of the stack. It even increases for very low stack widths, which may be due to the fact that the energy barriers at the side of the stack merge for very low stack widths. The off-current, on the other hand, is directly proportional to the stack area and can thus be directly downscaled by shrinking the stack width. For a stack width of 20 nm, $I_{on}/I_{off}$ ratios of more than $10^{32}$ can be reached.

### 3.2.2.3. Nonvolatile Memory Devices Based on Crested Barriers
One of the most important figures of merit of a nonvolatile memory cell is its $I_{on}/I_{off}$ ratio: A high on-current leads to low programming and erasing times, and a low off-current increases the retention



Figure 49. On-current density and off-current density as a function of the thickness of the central shutter barrier (left) and the stack width (right) [191].

time of the device. This ratio can be increased if, for a given device, the tunneling current in the on-state (the charging/discharging current) is increased or, in the off-state (during the retention time), decreased. With a single-layer dielectric it is not possible to tune on- and off-current independently. However, if the tunnel dielectric is replaced by a dielectric stack of varying barrier height as shown in Fig. 50, it becomes possible. In this figure, the device structure and the conduction band edge in the on- and off-state are shown. The device consists of a standard EEPROM structure, where the tunnel dielectric is composed of three layers. The middle layer has a higher energy barrier than the inner and outer layers. The flat-band case is indicated by the dotted lines.

In the on-state, a high voltage is applied on the top contact. The middle energy barrier is strongly reduced and gives rise to a high tunneling current. If the dielectric would consist of a single layer, the peak of the energy barrier would not be reduced. Thus, the on-current is much higher for the layered dielectric. In the off-state, a low negative voltage—due to charge stored on the memory node—is applied. The middle barrier is only slightly suppressed and blocks tunneling. The off-current is only slightly lower than for a single-layer dielectric. This behavior results in a high $I_{on}/I_{off}$ ratio. A high suppression of the middle barrier in the on-state requires a low permittivity of the outer layers so that the potential drop in the outer layers is high [193]. This device design was first proposed by Capasso et al. in 1988 [194] based on AlGaAs–GaAs devices and later used by several authors [195, 196], where it became popular as *crested-barrier* memory or VARIOT (*varying oxide thickness device*).

The gate current density of the device depicted in Fig. 50 is shown as a function of the gate bias in the left part of Fig. 51. A stack thickness of 5 nm was chosen. Because the middle layers must have a high band gap, only few material combinations are possible. For the simulations, middle layers of $Al_2O_3$ and $SiO_2$ have been chosen, with outer layers of $Y_2O_3$, $Si_3N_4$, and $ZrO_2$. For comparison, full $SiO_2$ and $Si_3N_4$ stacks have also been simulated (the dotted and dash-dotted lines). Though $Y_2O_3$ shows a very high off-current, stacks with outer layers of $Si_3N_4$ or $ZrO_2$ and $Al_2O_3$ as middle layer show good ratios between the on-state (positive gate bias) current density and the off-state (negative gate bias) current density.

The important figure of merit, however, is the $I_{on}/I_{off}$ ratio. In the right part of Fig. 51, the $I_{on}/I_{off}$ ratio is shown for $Si_3N_4$ and $ZrO_2$ stacks with $SiO_2$ and $Al_2O_3$ middle layers as a function of the thickness of the middle layer. Also shown is the ratio for a layer of $SiO_2$ and $Si_3N_4$ alone. It is obvious that the ratio strongly depends on the thickness of the middle layer, and both minima and maxima can be observed. Only outer layers of $Si_3N_4$ lead to a



Figure 50. Device structure and operating principle of a nonvolatile memory based on crested barriers.

Figure 51. Gate current density as a function of the gate bias for different materials of the middle layer compared to full SiO$_2$ and Si$_3$N$_4$ layers (left). Ratio between the on-current and the off-current as a function of the middle layer thickness for different materials of the outer layers (Si$_3$N$_4$ and ZrO$_2$) and middle layers (Al$_2$O$_3$ and SiO$_2$) compared to the resulting current density using full layers of SiO$_2$ and Si$_3$N$_4$ (right).

significantly increased performance as compared to full layers of SiO$_2$ or Si$_3$N$_4$. A middle layer thickness around 1–2 nm for the assumed 6-nm stack gives optimum performance.

## 4. CONCLUSIONS

Tunneling effects in semiconductor devices were investigated. A hierarchy of tunneling models was outlined. Three main properties were identified to influence the tunneling process: The carrier energy distribution function, the transmission coefficient, and the presence of traps in the dielectric layer.

The energetic distribution of carriers was investigated using different approximations, such as the frequently applied Fermi–Dirac or Maxwell–Boltzmann statistics. However, these approximations are only valid near equilibrium. Comparisons with the results from Monte Carlo simulations showed that in turned-on devices, the distribution function strongly deviates from the ideal shape. Some non-Maxwellian models were reviewed, and it was found that a model that is based on the solution variables of a six-moments transport model accurately reproduces the Monte Carlo results.

The quantum-mechanical transmission coefficient can be computed from the solution of the stationary Schrödinger equation. Several approximations and analytical formulae were outlined. For a single-layer dielectric, the analytical WKB approximation or Gundlach's formula can be used. For arbitrary-shaped energy barriers, the numerical WKB, the transfer-matrix, or the quantum transmitting boundary method can be applied. It was found that the transfer-matrix method is prone to numerical problems due to the repeated matrix multiplications. The quantum transmitting boundary method turned out to be more robust.

Defects in the dielectric layer give rise to trap-assisted tunneling, which leads to an additional tunneling current at low bias. After reviewing several models from the literature, a recently presented inelastic trap-assisted tunneling model was adapted to avoid the numerical calculation of the overlap integral in the dielectric layer. This yielded a fully analytical model that was further developed to include transient trap charging and discharging effects.

Several examples were studied where a general distinction between tunneling in MOS transistors, where it is a parasitic effect, and tunneling in nonvolatile memory devices, where it is crucial for the device functionality, was made. Tunneling in MOS transistors was investigated, where special attention was paid to the investigation of the different tunneling paths from the gate to the channel and from the gate to the source and drain extension regions.

Furthermore, the importance of the carrier distribution functions for modeling of gate leakage in turned-on devices was shown. If a heated Maxwellian approximation was used for

the description of hot-carrier tunneling, the gate current density was heavily overestimated. This effect was found to be especially pronounced for devices with short gate lengths.

In future CMOS devices, the use of alternative dielectric materials instead of $SiO_2$ will make the reduction of the effective oxide thickness possible. Several candidate materials were studied, and it was found that they show a pronounced correlation between the barrier height and the permittivity. This makes optimization necessary to find the optimum layer composition. Furthermore, the investigation of a MOS capacitor with a $ZrO_2$ dielectric showed that the strong defect density makes the use of trap-assisted tunneling models a *sine qua non* for these materials.

In addition to MOS transistors, nonvolatile memory devices were studied. A general overview of nonvolatile memory technology was followed by an investigation of three selected device structures: devices where the floating gate contact is replaced by a layer of trap-rich dielectric, multibarrier tunneling devices, and devices that are based on crested barriers. Especially the multibarrier tunneling devices allow an extremely high $I_{on}/I_{off}$ ratio. The trap-rich dielectric devices, on the other hand, are easier to fabricate and have a smaller footprint. Devices that are based on crested barriers allow tuning of the on- and off-current density independently. However, the $I_{on}/I_{off}$ ratio heavily depends on the thicknesses of the dielectric layers, and simulation is necessary to find the optimum values. The investigated nonvolatile memory applications are expected to show high performance; however, the bad quality of the interface between the dielectric layers may offset the advantage in the $I_{on}/I_{off}$ ratio.

## ACKNOWLEDGMENTS

## REFERENCES

1.  J. P. Shiely, Ph.D. thesis, Duke University, 1999.
2.  R. Clerc, Ph.D. thesis, Institut National Polytechnique de Grenoble, 2001.
3.  C. B. Duke, "Tunneling in Solids," Academic Press, New York, 1969.
4.  R. Tsu and L. Esaki, *Appl. Phys. Lett.* 22, 562 (1973).
5.  N. Ashcroft and N. Mermin, "Solid State Physics," Harcourt College Publishers, Fort Worth, TX, 1976.
6.  D. Cassi and B. Riccò, *IEEE Trans. Electron. Devices* 37, 1514 (1990).
7.  A. Abramo and C. Fiegna, *J. Appl. Phys.* 80, 889 (1996).
8.  K.-I. Sonoda, M. Yamaji, K. Taniguchi, C. Hamaguchi, and S. T. Dunham, *J. Appl. Phys.* 80, 5444 (1996).
9.  C. Fiegna, F. Venturi, M. Melanotte, E. Sangiorgi, and B. Riccò, *IEEE Trans. Electron. Devices* 38, 603 (1991).
10. K. Hasnat, C.-F. Yeap, S. Jallepalli, S. A. Hareland, W.-K. Shih, V. M. Agostinelli, A. F. Tasch, and C. M. Maziar, *IEEE Trans. Electron. Devices* 44, 129 (1997).
11. T. Grasser, H. Kosina, C. Heitzinger, and S. Selberherr, *J. Appl. Phys.* 91, 3869 (2002).
12. T. Grasser, H. Kosina, and S. Selberherr, *J. Appl. Phys.* 90, 6165 (2001).
13. E. Nicollian and J. Brews, "MOS (Metal Oxide Semiconductor) Physics and Technology," Wiley, New York, 1982.
14. Y. Tsividis, "Operation and Modeling of the MOS Transistor," McGraw-Hill, New York, 1987.
15. S. M. Sze, "Physics of Semiconductor Devices," 2nd ed. Wiley, New York, 1981.
16. M. Levinshtein, S. Rumyantsev, and M. Shur, "Handbook Series on Semiconductor Parameters," Vol. 1, World Scientific, Singapore, 1996.
17. Y.-C. Yeo, T.-J. King, and C. Hu, *J. Appl. Phys.* 92, 7266 (2002).
18. W. Harrison, "Electronic Structure and the Properties of Solids," Dover Publications, New York, 1989.
19. E. H. Rhoderick and R. H. Williams, "Metal-Semiconductor Contacts," Oxford University Press, Oxford, UK, 1988.
20. W. Franz, in "Handbuch der Physik" (S. Flügger, Ed.), Vol. XVII, p. 155, Springer, Berlin, 1956.
21. M. V. Fischetti, S. E. Laux, and E. Crabbé, *J. Appl. Phys.* 78, 1058 (1995).
22. M. Kleefstra and G. C. Herman, *J. Appl. Phys.* 51, 4923 (1980).
23. F. Jiménez-Molinos, F. Gámiz, A. Palma, P. Cartujo, and J. A. Lopez-Villanueva, *J. Appl. Phys.* 91, 5116 (2002).
24. G. Yang, K. Chin, and R. Marcus, *IEEE Trans. Electron. Devices* 38, 2373 (1991).
25. A. Schenk and G. Heiser, *J. Appl. Phys.* 81, 7900 (1997).
26. A. Schenk, "Advanced Physical Models for Silicon Device Simulation," Springer, Wien, 1998.
27. C. Fiegna, E. Sangiorgi, and L. Selmi, *IEEE Trans. Electron. Devices* 40, 2018 (1993).
28. Z. A. Weinberg, *J. Appl. Phys.* 53, 5052 (1982).

29. W.-Y. Quan, D. M. Kim, and M. K. Cho, *J. Appl. Phys.* 92, 3724 (2002).
30. L. Larcher, A. Paccagnella, and G. Ghidini, *IEEE Trans. Electron. Devices* 48, 271 (2001).
31. B. Majkusiak, *IEEE Trans. Electron. Devices* 37, 1087 (1990).
32. E. Schrödinger, *Annalen der Physik* 79, 361 (1926).
33. A. Messiah, "Quantum Mechanics," Dover, New York, 2000.
34. S. Gasiorowicz, "Quantum Physics," John Wiley & Sons, New York, 1995.
35. A. Hadjadj, G. Salace, and C. Petit, *J. Appl. Phys.* 89, 7994 (2001).
36. S. Nagano, M. Tsukiji, E. Hasegawa, and A. Ishitani, *J. Appl. Phys.* 75, 3530 (1994).
37. L. F. Register, E. Rosenbaum, and K. Yang, *Appl. Phys. Lett.* 74, 457 (1999).
38. H. Y. Yang, H. Niimi, and G. Lucovsky, *J. Appl. Phys.* 83, 2327 (1998).
39. N. Yang, W. K. Henson, J. R. Hauser, and J. J. Wortman, *IEEE Trans. Electron. Devices* 46, 1464 (1999).
40. M. I. Vexler, N. Asli, A. F. Shulekin, B. Meinerzhagen, and P. Seegebrecht, *Microelectronic Engineering* 59, 161 (2001).
41. J. Zhang, J. S. Yuan, Y. Ma, and A. S. Oates, *Solid-State Electron.* 44, 2165 (2000).
42. K. H. Gundlach, *Solid-State Electron.* 9, 949 (1966).
43. M. Abramowitz and I. A. Stegun, "Handbook of Mathematical Functions," Dover, New York, 1972.
44. A. Shanware, J. P. Shiely, and H. Z. Massoud, in "Proceeding of the International Electron Devices Meeting," pp. 815–818, IEEE Press, Piscataway, NJ, 1999.
45. M. O. Vassell, J. Lee, and H. F. Lockwood, *J. Appl. Phys.* 54, 5208 (1983).
46. Y. Ando and T. Itoh, *J. Appl. Phys.* 61, 1497 (1987).
47. B. Zimmermann, E. Marclay, M. Ilegems, and P. Gueret, *J. Appl. Phys.* 64, 3581 (1988).
48. G. Yong, *Phys. Rev. B* 50, 17249 (1994).
49. R. Clerc, A. Spinelli, G. Ghibaudo, and G. Pananakakis, *J. Appl. Phys.* 91, 1400 (2002).
50. D. K. Ferry and S. M. Goodnick, "Transport in Nanostructures," Cambridge University Press, Cambridge, UK, 1997.
51. W. W. Lui and M. Fukuma, *J. Appl. Phys.* 60, 1555 (1986).
52. K. F. Brennan, *J. Appl. Phys.* 62, 2392 (1987).
53. D. C. Hutchings, *Appl. Phys. Lett.* 55, 1082 (1989).
54. J.-G. S. Demers and R. Maciejko, *J. Appl. Phys.* 90, 6120 (2001).
55. B. A. Biegel, Ph.D. thesis, Stanford University, 1997.
56. J. N. Schulman and Y.-C. Chang, *Phys. Rev. B* 27, 2346 (1983).
57. D. Y. K. Ko and J. C. Inkson, *Phys. Rev. B* 38, 9945 (1988).
58. T. Usuki, M. Saito, M. Takatsu, R. A. Kiehl, and N. Yokoyama, *Phys. Rev. B* 52, 8244 (1995).
59. D. Z. Y. Ting, E. T. Yu, and T. C. McGill, *Phys. Rev. B* 45, 3583 (1992).
60. W. R. Frensley and N. G. Einspruch, Eds., "Heterostructures and Quantum Devices, VLSI Electronics: Microstructure Science." Academic Press, New York, 1994.
61. C. S. Lent and D. J. Kirkner, *J. Appl. Phys.* 67, 6353 (1990).
62. A. P. Gnädinger and H. E. Talley, *Solid-State Electron.* 13, 1301 (1970).
63. F. Stern, *Phys. Rev. B* 5, 4891 (1972).
64. W. Magnus and W. Schoenmaker, *Microelectronics Reliability* 41, 31 (2001).
65. E. Anemogiannis, E. N. Glytsis, and T. K. Gaylord, *IEEE J. Quant. Electron.* 29, 2731 (1993).
66. E. Cassan, *J. Appl. Phys.* 87, 7931 (2000).
67. A. Schenk, in "Proceedings of the European Solid-State Device Research Conference" (H. Ryssel, G. Wachutka, and H. Grünbacher, Eds.), pp. 9–16, Frontier Group, 2001.
68. A. T. M. Fairus and V. K. Arora, *Microelectronics Journal* 32, 679 (2000).
69. N. Matsuo, Y. Takami, and Y. Kitagawa, *Solid-State Electron.* 46, 577 (2002).
70. S. Padmanabhan and A. Rothwarf, *IEEE Trans. Electron. Devices* 36, 2557 (1989).
71. M. J. van Dort, P. H. Woerlee, and A. J. Walker, *Solid-State Electron.* 37, 411 (1994).
72. G. Gildenblatt, B. Gelmont, and S. Vatannia, *J. Appl. Phys.* 77, 6327 (1995).
73. P. J. Price, *Phys. Rev. B* 45, 9042 (1992).
74. P. J. Price, *Appl. Phys. Lett.* 82, 2080 (2003).
75. A. Thean and J. P. Leburton, *IEEE Electron. Device Lett.* 22, 148 (2001).
76. A. Ghetti, A. Hamad, P. J. Silverman, H. Vaidya, and N. Zhao, in "Proceedings of the Simulation of Semiconductor Processes and Devices," pp. 239–242, IEEE Press, Piscataway, NJ, 1999.
77. S. H. Lo, D. A. Buchanan, Y. Taur, and W. Wang, *IEEE Trans. Electron. Devices* 18, 209 (1997).
78. S. Mudanai, Y. Fan, Q. Ouyang, A. F. Tasch, and S. K. Banerjee, *IEEE Trans. Electron. Devices* 47, 1851 (2000).
79. S. Mudanai, L. F. Register, A. F. Tasch, and S. K. Banerjee, *IEEE Electron. Device Lett.* 22, 145 (2001).
80. F. Rana, S. Tiwari, and D. A. Buchanan, *Appl. Phys. Lett.* 69, 1104 (1996).
81. W.-K. Shih, E. X. Wang, S. Jallepalli, F. Leon, C. M. Maziar, and A. F. Tasch, Jr., *Solid-State Electron.* 42, 997 (1998).
82. E. Cassan, P. Dollfus, S. Galdin, and P. Hesto, *IEEE Trans. Electron. Devices* 48, 715 (2001).
83. A. Dalla Serra, A. Abramo, P. Palestri, L. Selmi, and F. Widdershoven, *IEEE Trans. Electron. Devices* 48, 1811 (2001).
84. Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst, Eds., "Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide," SIAM, Philadelphia, 2000.
85. R. Lake, G. Klimeck, R. C. Bowen, and D. Jovanovic, *J. Appl. Phys.* 81, 7845 (1997).

86. R. C. Bowen, W. R. Frensley, G. Klimeck, and R. K. Lake, *Phys. Rev. B* 52, 2754 (1995).
87. C. Bowen, C. L. Fernando, G. Klimeck, A. Chatterjee, D. Blanks, R. Lake, J. Hu, J. Davis, M. Kulkarni, S. Hattangady, and I.-C. Chen, in "Proceedings of the International Electron Devices Meeting." pp. 35.1.1–35.1.4. IEEE Press, Piscataway, NJ, 1997.
88. G. Klimeck, R. Lake, C. Bowen, W. R. Frensley, and T. S. Moise, *Appl. Phys. Lett.* 67, 2539 (1995).
89. C. L. Fernando and W. R. Frensley, *J. Appl. Phys.* 76, 2881 (1994).
90. W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, "Numerical Recipes in C." Cambridge University Press, Cambridge, UK, 1997.
91. W. R. Frensley. *Superlattices & Microstructures* 11, 347 (1992).
92. "Nanotechnology Engineering Modeling Program (NEMO) Version 3.0." Raytheon TI Systems, 1996.
93. N. Arora, "MOSFET Models for VLSI Circuit Simulation." Springer, Berlin, 1993.
94. C.-H. Choi, K.-H. Oh, J.-S. Goo, Z. Yu, and R. W. Dutton, in "Proceedings of the International Electron Devices Meeting," pp. 30.6.1–30.6.4. IEEE Press, Piscataway, NJ, 1999.
95. C.-H. Choi, K.-Y. Nam, Z. Yu, and R. W. Dutton, *IEEE Trans. Electron. Devices* 48, 2823 (2001).
96. Y.-S. Lin, H.-T. Huang, C.-C. Wu, Y.-K. Leung, H.-Y. Pan, T.-E. Chang, W.-M. Chen, J.-J. Liaw, and C. H. Diaz, *IEEE Trans. Electron. Devices* 49, 442 (2002).
97. H. Lin, J. T.-Y. Chen, and J.-H. Chang, *Solid-State Electron.* 46, 1145 (2002).
98. S. Schwantes and W. Krautschneider, in "Proceedings of the European Solid-State Device Research Conference" (H. Ryssel, G. Wachutka, and H. Grünbacher, Eds.), pp. 471–474. Frontier Group, 2001.
99. R. H. Fowler and L. Nordheim, *Proc. Roy. Soc. A* 119, 173 (1928).
100. M. Lenzlinger and E. H. Snow, *J. Appl. Phys.* 40, 278 (1969).
101. K. F. Schuegraf and C. Hu, *IEEE Trans. Electron. Devices* 41, 761 (1994).
102. K. F. Schuegraf, C. C. King, and C. Hu, in "Proceedings of the Symposium on VLSI Technology," 1992, pp. 18–19.
103. S. Aritome, R. Shirota, G. Hemink, T. Endoh, and F. Masuoka, *Proc. IEEE* 81, 776 (1993).
104. R. Moazzami and C. Hu, in "Proceedings of the International Electron Devices Meeting," pp. 139–142. IEEE Press, Piscataway, NJ, 1992.
105. E. Rosenbaum and L. F. Register, *IEEE Trans. Electron. Devices* 44, 317 (1997).
106. S.-I. Takagi, N. Yasuda, and A. Toriumi, *IEEE J. Solid-State Circuits* 46, 348 (1999).
107. R. Rofan and C. Hu, *IEEE Electron. Device Lett.* 12, 632 (1991).
108. J. Wu, L. F. Register, and E. Rosenbaum, in "Proceedings of the International Reliability Physics Symposium," 1999, pp. 389–395.
109. B. Riccò, G. Gozzi, and M. Lanzoni, *IEEE Trans. Electron Devices* 45, 1554 (1998).
110. K. Sakakibara, N. Ajika, K. Eikyu, K. Ishikawa, and H. Miyoshi, *IEEE Trans. Electron. Devices* 44, 1002 (1997).
111. A. Ghetti, E. Sangiorgi, J. Bude, T. W. Sorsch, and G. Weber, *IEEE Trans. Electron. Devices* 47, 2358 (2000).
112. C.-M. Yih, Z.-H. Ho, M.-S. Liang, and S. S. Chung, *IEEE Trans. Electron. Devices* 48, 300 (2001).
113. A. I. Chou, K. Lai, K. Kumar, P. Chowdhury, and J. C. Lee, *Appl. Phys. Lett.* 70, 3407 (1997).
114. K. Komiya and Y. Omura, *Microelectronic Engineering* 59, 61 (2001).
115. T.-K. Kang, M.-J. Chen, C.-H. Liu, Y. J. Chang, and S.-K. Fan, *IEEE Trans. Electron. Devices* 48, 2317 (2001).
116. M. Lenski, T. Endoh, and F. Masuoka, *J. Appl. Phys.* 88, 5238 (2000).
117. S.-I. Takagi, N. Yasuda, and A. Toriumi, *IEEE Trans. Electron. Devices* 46, 335 (1999).
118. W. J. Chang, M. P. Houng, and Y. H. Wang, *J. Appl. Phys.* 89, 6285 (2001).
119. W. J. Chang, M. P. Houng, and Y. H. Wang, *J. Appl. Phys.* 90, 5171 (2001).
120. L. Larcher, A. Paccagnella, and G. Ghidini, *IEEE Trans. Electron. Devices* 48, 285 (2001).
121. D. Ielmini, A. S. Spinelli, M. A. Rigamonti, and A. L. Lacaita, *IEEE Trans. Electron. Devices* 47, 1258 (2000).
122. D. Ielmini, A. S. Spinelli, M. A. Rigamonti, and A. L. Lacaita, *IEEE Trans. Electron. Devices* 47, 1266 (2000).
123. D. Ielmini, A. S. Spinelli, A. L. Lacaita, A. Martinelli, and G. Ghidini, *Solid-State Electron.* 45, 1361 (2001).
124. D. Ielmini, A. S. Spinelli, A. L. Lacaita, and G. Ghidini, *Solid-State Electron.* 46, 417 (2002).
125. D. Ielmini, A. S. Spinelli, A. L. Lacaita, and A. Modelli, *Microelectronic Engineering* 59, 189 (2001).
126. D. Ielmini, A. S. Spinelli, A. L. Lacaita, and A. Modelli, *Solid-State Electron.* 46, 1749 (2002).
127. A. Ghetti, *Microelectronic Engineering* 59, 127 (2001).
128. C. Chaneliere, J. L. Autran, and R. A. B. Devine, *J. Appl. Phys.* 86, 480 (1999).
129. M. Houssa, R. Degraeve, P. W. Mertens, M. M. Heyns, J. S. Leon, A. Halliyal, and B. Ogle, *J. Appl. Phys.* 86, 6462 (1999).
130. M. Houssa, M. Tuominen, M. Naili, V. Afanas'ev, A. Stesmans, S. Haukka, and M. M. Heyns, *J. Appl. Phys.* 87, 8615 (2000).
131. D. Caputo, F. Irrera, S. Salerno, S. Spiga, and M. Fanciulli, in "Proceedings of the 4th European Workshop on Ultimate Integration of Silicon" (E. Sangiorgi and L. Selmi, Eds.), pp. 89–92. University of Udine, Udine, Italy, 2003.
132. B. DeSalvo, G. Ghibaudo, G. Pananakakis, B. Guillaumot, and G. Reimbold, *Solid-State Electron.* 44, 895 (2000).
133. E. Kameda, T. Matsuda, Y. Emura, and T. Ohzone, *Solid-State Electron.* 42, 2105 (1998).
134. J. A. López-Villanueva, J. A. Jiménez-Tejada, P. Cartujo, J. Bausells, and J. E. Carceller, *J. Appl. Phys.* 70, 3712 (1991).
135. F. Jiménez-Molinos, A. Palma, F. Gámiz, J. Banqueri, and J. A. Lopez-Villanueva, *J. Appl. Phys.* 90, 3396 (2001).

136. A. Palma, A. Godoy, J. A. Jimenez-Tejada, J. E. Carceller, and J. A. Lopez-Villanueva, *Phys. Rev. B* 56, 9565 (1997).

137. J. H. Zheng, H. S. Tan, and S. C. Ng, *J. Phys.: Condensed Matter* 6, 1695 (1994).

138. W. B. Fowler, J. K. Rudra, M. E. Zvanut, and F. J. Feigl, *Phys. Rev. B* 41, 8313 (1990).

139. M. Herrmann and A. Schenk, *J. Appl. Phys.* 77, 4522 (1995).

140. A. Gehring, H. Kosina, T. Grasser, and S. Selberherr, in "Proceedings of the 4th European Workshop on Ultimate Integration of Silicon" (E. Sangiorgi and L. Selmi, Eds.), pp. 131–134. University of Udine, Udine, Italy, 2003.

141. Minimos-NT 2.1 User's Guide, Institute for Microelectronics, Wien, 2004.

142. J. Cai and C.-T. Sah, *J. Appl. Phys.* 89, 2272 (2001).

143. H. Z. Massoud and J. P. Shiely, *Microelectronic Engineering* 36, 263 (1997).

144. Y. Shi, T. P. Ma, S. Prasad, and S. Dhanda, *IEEE Trans. Electron. Devices* 45, 2355 (1998).

145. M. Städele, B. Tuttle, B. Fischer, and K. Hess, *J. Comput. Electron.* 1, 153 (2002).

146. M. Städele, F. Sacconi, A. D. Carlo, and P. Lugli, *J. Appl. Phys.* 93, 2681 (2003).

147. F. Sacconi, M. Povolotskyi, A. D. Carlo, P. Lugli, and M. Städele, in "Proceeding of the 4th European Workshop on Ultimate Integration of Silicon" (E. Sangiorgi and L. Selmi, Eds.), pp. 125–128. University of Udine, Udine, Italy, 2003.

148. S. H. Lo, D. A. Buchanan, and Y. Taur, *IBM J. Res. Dev.* 43, 327 (1999).

149. S. Jallepalli, J. Bude, W.-K. Shih, M. R. Pinto, C. M. Maziar, and A. F. Tasch, Jr., *IEEE Trans. Electron. Devices* 44, 297 (1997).

150. MEDICI User's Manual, Synopsys, Mountain View, CA, 2003.

151. A. Gehring, T. Grasser, H. Kosina, and S. Selberherr, *J. Appl. Phys.* 92, 6019 (2002).

152. A. Gehring, T. Grasser, H. Kosina, and S. Selberherr, *Electron. Lett.* 39, 691 (2003).

153. L. Selmi, A. Ghetti, R. Bez, and E. Sangiorgi, *Microelectronic Engineering* 36, 293 (1997).

154. E. M. Vogel, K. Z. Ahmed, B. Hornung, K. Henson, P. K. McLarty, G. Lucovsky, J. R. Hauser, and J. J. Wortman, *IEEE Trans. Electron. Devices* 45, 1350 (1998).

155. A. Gehring, H. Kosina, and S. Selberherr, *J. Comput. Electron.* 2, 219 (2003).

156. Y.-Y. Fan, R. E. Nieh, J. C. Lee, G. Lucovsky, G. A. Brown, L. F. Register, and S. K. Banerjee, *IEEE Trans. Electron. Devices* 49, 1969 (2002).

157. S. Harasek, Ph.D. thesis, Technische Universität Wien, 2003.

158. A. Gehring, S. Harasek, E. Bertagnolli, and S. Selberherr, in "Proceedings of European Solid-State Device Research Conference" (J. Franca and R. Freitas, Eds.), pp. 473–476. Frontier Group, 2003.

159. P. Tanner, S. Dimitrijev, and H. B. Harrison, *Electron. Lett.* 31, 1880 (1995).

160. D. A. Antoniadis, I. J. Djomehri, K. M. Jackson, and S. Miller, "Well-Tempered" Bulk-Si NMOSFET Device Home Page, available at http://www-mtl.mit.edu/Well/.

161. W. D. Brown and J. Brewer, "Nonvolatile Semiconductor Memory Technology," IEEE Press, Piscataway, NJ, 1998.

162. A. Concannon, S. Keeney, A. Mathewson, R. Bez, and C. Lombardi, *IEEE Trans. Electron. Devices* 40, 1258 (1993).

163. S. Keeney, R. Bez, D. Cantarelli, F. Piccinin, A. Mathewson, L. Ravazzi, and C. Lombardi, *IEEE Trans. Electron. Devices* 39, 2750 (1992).

164. A. Kolodny, S. T. K. Nieh, B. Eitan, and J. Shappir, *IEEE Trans. Electron. Devices* 33, 835 (1986).

165. K. T. San, C. Kaya, D. K. Y. Liu, T.-P. Ma, and P. Shah, *IEEE Electron. Device Lett.* 13, 328 (1992).

166. R. Bouchakour, N. Harabech, P. Canet, P. Boivin, and J. M. Mirabel, in "Proceedings of the International Symposium on Circuits & Systems," 2001, pp. 822–825.

167. R. Duane, A. Concannon, P. O'Sullivan, M. O'Shea, and A. Mathewson, *Solid-State Electron.* 45, 235 (2001).

168. D. Kahng and S. M. Sze, *Bell Syst. Tech. J.* 46, 1288 (1967).

169. P. Pavan, R. Bez, P. Olivo, and E. Zanoni, *Proc. IEEE* 86, 1248 (1997).

170. P. Cappelletti, C. Golla, P. Olivo, and E. Zanoni, "Flash Memories," Kluwer Academic Publishers, Boston, 2000.

171. A. Gehring, F. Jiménez-Molinos, H. Kosina, A. Palma, F. Gámiz, and S. Selberherr, *Microelectron. Reliab.* 43, 1495 (2003).

172. P. Canet, R. Bouchakour, N. Harabech, P. Boivin, J. M. Mirabel, and C. Plossu, in "Proceedings of the 43rd IEEE Midwest Symposium on Circuits and Systems," 2000, pp. 1144–1147.

173. M. K. Cho and D. M. Kim, *IEEE Electron. Device Lett.* 21, 399 (2000).

174. J. M. Caywood, C. J. Huang, and Y. J. Chang, *IEEE Trans. Electron. Devices* 49, 802 (2002).

175. B. Eitan, P. Pavan, I. Bloom, E. Aloni, A. Frommer, and D. Finzi, *IEEE Electron. Device Lett.* 21, 543 (2000).

176. M. H. White, D. A. Adams, and J. Bu, *IEEE Circuits Devices* 22 (2000).

177. K.-T. Chang, W.-M. Chen, C. Swift, J. M. Higman, W. M. Paulson, and K.-M. Chang, *IEEE Electron. Device Lett.* 19, 253 (1998).

178. G. Iannaccone and P. Coli, *Appl. Phys. Lett.* 78, 2046 (2001).

179. A. Thean and J. P. Leburton, *IEEE Electron. Device Lett.* 20, 286 (1999).

180. J. J. Welser, S. Tiwari, S. Rishton, K. Y. Lee, and Y. Lee, *IEEE Electron. Device Lett.* 18, 278 (1997).

181. B. DeSalvo, G. Ghibaudo, G. Pananakakis, P. Masson, T. Baron, N. Buffet, A. Fernandes, and B. Guillaumot, *IEEE Trans. Electron. Devices* 48, 1789 (2001).

182. K. Han, I. Kim, and H. Shin, *IEEE Trans. Electron. Devices* 48, 874 (2001).

183. H. I. Hanafi, S. Tiwari, and I. Khan, *IEEE Trans. Electron. Devices* 43, 1553 (1996).

184. Y.-C. King, T.-J. King, and C. Hu, *IEEE Trans. Electron. Devices* 20, 409 (1999).
185. X. Tang, X. Baie, J.-P. Colinge, C. Gusting, and V. Bayot, *IEEE Trans. Electron. Devices* 49, 1420 (2002).
186. K. Nakazato, P. J. A. Piotrowicz, D. G. Hasko, H. Ahmed, and K. Itoh, in "Proceedings of the International Electronic Devices Meeting," pp. 179–182. IEEE Press, Piscataway, NJ, 1997.
187. H. Mizuta, K. Nakazato, P. J. A. Piotrowicz, K. Itoh, T. Teshima, K. Yamaguchi, and T. Shimada, in "Proceedings of the Symposium on VLSI Technology," 1998, pp. 128–129.
188. N. Nakazato, K. Itoh, H. Mizuta, and H. Ahmed, *Electron. Lett.* 35, 848 (1999).
189. K. Nakazato, K. Itoh, H. Ahmed, H. Mizuta, T. Kisu, M. Kato, and T. Sakata, in "Proceedings of the International Solid-State Circuits Conference," 2000, p. TA 7.4.
190. H. Mizuta, M. Wagner, and K. Nakazato, *IEEE Trans. Electron. Devices* 48, 1103 (2001).
191. A. Gehring, T. Grasser, B.-H. Cheong, and S. Selberherr, *Solid-State Electron.* 46, 1545 (2002).
192. H. Fukuda, J. L. Hoyt, M. A. McCord, and R. F. W. Pease, *Appl. Phys. Lett.* 70, 333 (1997).
193. J. D. Casperson, L. D. Bell, and H. A. Atwater, *J. Appl. Phys.* 92, 261 (2002).
194. F. Capasso, F. Beltram, R. J. Malik, and J. F. Walker, *IEEE Electron. Device Lett.* 9, 377 (1988).
195. K. K. Likharev, *Appl. Phys. Lett.* 73, 2137 (1998).
196. B. Govoreanu, P. Blomme, M. Rosmeulen, J. V. Houdt, and K. D. Meyer, *IEEE Electron. Device Lett.* 24, 99 (2003).
197. W. Harrison, "Solid State Theory." Dover, New York, 1979.
198. M. LeRoy, E. Lheurette, O. Vanbesien, and D. Lippens, *J. Appl. Phys.* 93, 2966 (2003).
199. C. M. Osburn, I. Kim, S. K. Han, I. De, K. F. Yee, S. Gannavaram, S. J. Lee, C.-H. Lee, Z. J. Luo, W. Zhu, J. R. Hauser, D.-L. Kwong, G. Lucovsky, T. P. Ma, and M. C. Öztürk, *IBM J. Res. Dev.* 46, 299 (2002).
200. G. D. Wilk, R. M. Wallace, and J. M. Anthony, *J. Appl. Phys.* 89, 5243 (2001).
201. J. Robertson, *J. Vac. Sci. Technol.* 18, 1785 (2000).
202. H.-S. P. Wong, *IBM J. Res. Dev.* 46, 133 (2002).

# CHAPTER 10

# Electronic Structure of Quantum Dots

## J. B. Wang, C. Hines, R. D. Muhandiramge

*School of Physics, The University of Western Australia, Crawley, Australia*

## CONTENTS

## 1. INTRODUCTION

With the increased sophistication of semiconductor technology, it is now possible to create man-made objects that display many of the characteristic properties normally associated with atoms. In semiconductors, all electrons are tightly bound to the nuclei except for a very small fraction of mobile electrons. These mobile electrons can be trapped between two layers of semiconductors, confined effectively to two dimensions. They are then restricted further by

some lateral confining potential to create a quantum dot, often referred as to an artificial atom. The ability of experimentalists to manipulate the size and shape of these artificial atoms and to use different materials in their construction has opened up a wide range of possibilities and areas for examination [1–10].

For example, it is possible to vary the exact number of mobile electrons in the dot by simply changing the voltage applied to the gate electrode, allowing one to scan through the periodic table of artificial atoms with ease. Another interesting aspect of quantum dots is that, unlike atomic systems, they are not limited to being spherically or circularly symmetric—we can have elliptic dots, rectangular dots, triangular dots, and even dots without any symmetry. Transitions never observed in the spectra of natural atoms can be obtained from the artificial ones. In the future, quantum dots may be used to build more efficient and precisely controlled lasers with otherwise inaccessible wavelengths, and also as vital components of nanoelectronic devices [3, 11–13]. It is also hoped that quantum dots may one day be able to help realize the dream of quantum computing [14–16].

In atomic systems, electrons are confined by the attractive Coulombic potential of the positively charged nucleus. In quantum dots the confinement of the electrons is instead the result of an artificially created potential, formed by electrodes connected to layers of semiconductor, as shown in Fig. 1a. From quantum theory we know that if the electrons are confined, then they are only allowed to possess certain discrete energies. In this sense, quantum dots are just like the natural atoms with well-defined energy level structures. What does change is that the energy structure associated with these artificial atoms now reflects the two-dimensional nature of the system. We get a different periodic table of these two-dimensional "elements," such as that shown in Fig. 1b.

The atomic-like structure of quantum dots has been demonstrated by the experimental work of Tarucha et al. [17], in which electrons are removed from a quantum dot one at a time by varying the gate voltage. A current will flow only if the number of electrons in the dot changes. It was found that only at particular voltages can an electron be removed because of the discrete nature of its energy structure, giving rise to well-separated current peaks, as shown in Fig. 2. The energy difference between consecutive current peaks is a measure of the addition energy, which is the energy required to add an extra electron to the dot and is shown in the inset of Fig. 2. Such measurements reflect the two-dimensional shell structure of the quantum dot. Note, for example, that the gap between the adjacent current peaks is significantly greater when the dot undergoes changes from two electrons to three, from six to seven, and from twelve to thirteen. This corresponds to the fact that a two-dimensional artificial atom has completely filled shells at electron numbers of $N = 2, 6, 12$. This differs from a three-dimensional natural atom, the shells of which are completely filled at atomic numbers of $N = 2, 10, 18, 36, \ldots$ instead. The data of Tarucha et al. [17] also exhibit secondary maxima in the addition energy, which correspond to half-filled shells at electron numbers $N = 4, 9, 16$. Theoretical analysis shows that such shell structures correspond to circularly symmetric confinement potentials, and deviation from this geometry can change the electronic structure in a significant way [18].

(a)      ~0.5 μm      (b)



Figure 1. Schematic representation of a typical quantum dot. Reprinted with permission from [1]. L. P. Kouwenhoven et al., *Rep. Prog. Phys.* 64, 701 (2001). © 2001, Institute of Physics.

The detailed electronic structure of quantum dot systems depends on the material, size, and geometry of the quantum dots, requiring *ab initio* calculations with full quantum mechanical treatment for their accurate description. This is the challenge imposed on theoreticians. In this chapter, we review and describe the theoretical models, basic numerical techniques, and several computational schemes developed to solve the corresponding Schrödinger's equation and to evaluate the atomic properties of quantum dots.

## 2. THEORETICAL MODEL

### 2.1. Two-Dimensional Mobile Electrons

Quantum dot systems are typically made up of many layers of semiconductor material with metallic gates used to form the patterned structures. Though themselves consisting of roughly $10^3$–$10^9$ atoms and a corresponding number or electrons, most of these electrons are tightly bound to the nuclei of the semiconductor material, and only a small proportion of the electrons are mobile. It is these mobile electrons that have the greatest effect on the electronic properties of the quantum dot [19].

As an example, the quantum dot shown in Fig. 1a is a double barrier heterostructure in which $In_{0.05}Ga_{0.95}As$ forms the central well that is sandwiched between two $Al_{0.22}Ga_{0.78}As$ barriers. The AlGaAs barriers confine the electrons motion in the vertical ($z$) direction, while the edge of the dot and the gate potential provide a variable lateral confinement. By regulating the gate potential, one can vary the number of electrons in the dot. The conducting source and drain are Si doped n-GaAs. The diameter of the dot is on the order of hundreds of nanometers, whereas the height is approximately 10 nm.

The confining potential can be separated into a vertical ($z$) and a lateral ($x$, $y$) component. The confining potential in the vertical direction is a very narrow and effectively infinitely high well [20], whereas the lateral confining potential $V(x, y)$ has a bowl-like shape. The energy level of the first excited state in the $z$ direction is generally hundreds of times greater than many of the low-energy states in the $x$-$y$ plane. This property allows us to model electron motion in a quantum dot as two dimensional, as the electrons are tightly confined in the $z$ direction, as they only occupy the ground state in this direction.

### 2.2. Effective Mass Approximation

At the length scales involved in quantum dots in which the de Broglie wavelength of the electron is on the order of the confinement region, electrons behavior is dominated by their quantum mechanical properties, most notably a quantisation of their possible energies. These properties are described by the electron's wavefunctions, the evolution of which is governed by the Schrödinger equation

$$i\hbar \frac{\partial \Psi(\mathbf{r}, \sigma, t)}{\partial t} = \hat{H}\Psi(\mathbf{r}, \sigma, t) \tag{1}$$

where the system Hamiltonian

$$\hat{H} = \frac{1}{2m}\left[-i\hbar\nabla + \frac{e}{c}\hat{A}(\mathbf{r})\right]^2 + V(\mathbf{r}) + g\mu_B B S_z + T \tag{2}$$

where $\psi(\mathbf{r}, \sigma, t)$ is the system wavefunction, $\mathbf{r}$ represents collectively the spatial coordinates of all electrons in the system, $\sigma$ represents the spin coordinates of all electrons, $\nabla$ is the gradient operator acting on the spatial coordinates of all electrons, $A$ is the external magnetic vector potential, the magnetic field $\hat{B} = \nabla \times \hat{A}$. $V$ is the electronic potential including the Coulomb interaction between electron pairs, $g$ is the Landé $g$ factor, $S_z$ is the $z$-component of the total spin, $g\mu_B B S_z$ is the Zeeman energy associated with the interaction of spin with the magnetic field, and the last term $T$ represents the spin and relativistic interactions.

The above equation provides both spatial and phase information of the system at all times, and thus a complete knowledge of all possible observables of the system under study. However, solving this equation for all electrons and particles in the semiconductor is not at all practical. Certain approximations need to be made by taking into account physically the most important aspects of the system. Guidance is taken from previous and current experimental and theoretical studies as to which effects can be neglected and which are important, with attempts made to understand the effects on the final solution.

Firstly, only mobile electrons will be considered, while all other electrons and nuclei provide a lattice-like potential that can be rather complex because of the large number of atoms in the material. However, it has been shown in theory and by experimentation [21–23] that such lattice effects and other local effects can be approximated by an electron with an effective mass $m^*$. The effective permittivity $\varepsilon^*$ is also different from the vacuum permittivity $\varepsilon_0$ because of a screening effect in the semiconductor. For gallium arsenide, the effective mass $m^*$ is approximately $0.0667m_e$, and the effective permittivity $\varepsilon^*$ is approximately $12.9\varepsilon_0$ [23–25]. Note that a number of assumptions are made here about the material: it is homogeneous and uniform, it has perfect crystalline structure and high purity, and it is nondeforming because of screening (i.e., the presence of the conduction electron does not alter the shape of the lattice potential). Second, the Landé $g$ factor is generally small for semiconductor materials (using GaAs as an example, [$g = -0.44$]). In this case, the Zeeman energy $g\mu_B B S_z$ is small and can thus be neglected [26, 27]. Other relativistic effects, such as spin-orbit coupling, are also negligibly small in comparison with the Coulomb energy, the exchange interaction, and the confinement potentials.

These assumptions lead to a much simplified Hamiltonian for the quantum dot system,

$$\hat{H} = \sum_{i=1}^{N}\left(\frac{1}{2m^*}(-i\hbar\nabla_i + e\hat{A}_i)^2 + V(x_i, y_i)\right) + \frac{1}{\varepsilon^*}\sum_{i<j}^{N}\frac{1}{\sqrt{|r_i - r_j|^2}}$$

$$= \sum_{i=1}^{N}\hat{H}_i + \frac{1}{\varepsilon^*}\sum_{i<j}^{N}\frac{1}{r_{ij}}$$

$$= \hat{H}^{\mathrm{dot}} + \hat{H}^{\mathrm{int}} \tag{3}$$

where $N$ is the number of mobile electrons confined inside the quantum dot, $\hat{H}_i$ is the single-electron Hamiltonian acting on the $i$th electron, $\hat{H}^{\mathrm{dot}}$ is the sum of all single-electron Hamiltonian terms, and $\hat{H}^{\mathrm{int}}$ is the sum of all two-electron Hamiltonian terms describing the Coulomb interactions between each pair of electrons. Furthermore, in this review chapter, we are mainly concerned with the electronic structure of quantum dots, which is independent of time. In this case, we can separate out the time-dependent factor from the Schrödinger equation and solve only the time-independent Schrödinger equation; that is,

$$\hat{H}\Psi(\mathbf{r}_1\cdots\mathbf{r}_N, \sigma_1\cdots\sigma_N) = E\Psi(\mathbf{r}_1\cdots\mathbf{r}_N, \sigma_1\cdots\sigma_N) \tag{4}$$

Although the Hamiltonian given by Eq. (3) does not contain any spin-related terms, the spin coordinates are still very important in the description of multielectron wavefunctions.

This is because electrons are fermions that must also obey the Pauli's exclusion principle. This is stated mathematically by the antisymmetry principle; that is, if both the spin and spatial coordinates of two electrons are interchanged, the wavefunction must change sign. In other words, two electrons with the same spin can not be at the same place at the same time, as both electrons would have the same coordinates, making the wavefunction vanish. A valid wavefunction would satisfy both the Schrödinger equation and the Pauli's exclusion principle.

## 3. BASIC ANALYTICAL AND NUMERICAL TECHNIQUES

This section details basic analytical and numerical techniques common to the various computational schemes developed and used to study the electronic structure of quantum dots and other nanosystems.

### 3.1. Single-Electron Quantum Dots

#### 3.1.1. Analytic Solutions for Parabolic Confinement Potential

If the quantum dot is formed by a circularly symmetric parabolic confinement potential and an uniform external magnetic field $B$ is applied perpendicular to the two dimensional $x$-$y$ plane, we have analytical solutions for such two-dimensional single-electron quantum dot systems. This was first established by Fock [28] and later, independently, by Darwin [29] and Landau [30]. Following Fock's work, we write the system Hamiltonian as the following:

$$\hat{H} = \frac{1}{2m^*}\left(-i\hbar\nabla - \frac{e}{c}\hat{A}\right)^2 + \frac{1}{2}m^*\omega_0^2 r^2$$

Because the magnetic field is applied perpendicularly to the $x$-$y$ plane, $\hat{A} = (+\frac{1}{2}B\hat{y}, -\frac{1}{2}B\hat{x}, 0)$ in the symmetric gauge. In this case, we have

$$\hat{H}\psi = \frac{1}{2m^*}\left(i\hbar\frac{\partial}{\partial x} + i\hbar\frac{1}{2}\frac{e}{c}B\hat{y}, \frac{\partial}{\partial y} - \frac{1}{2}\frac{e}{c}B\hat{x}, 0\right)^2\psi + \frac{1}{2}m^*\omega_0^2 r^2\psi \qquad (5)$$

where $\psi$ represents a single-electron wavefunction. In polar coordinates $(r, \theta)$; that is,

$$\hat{H}\psi = -\frac{\hbar^2}{2m^*}\left(\frac{\partial^2\psi}{\partial r^2} + \frac{1}{r}\frac{\partial\psi}{\partial r} + \frac{1}{r^2}\frac{\partial^2\psi}{\partial\theta^2}\right) + \frac{i\omega_c\hbar}{2}\frac{\partial\psi}{\partial\theta} + \left(\frac{m^*\omega_c^2 r^2}{8} + \frac{1}{2}m^*\omega_0^2 r^2\right)\psi \qquad (6)$$

where $\omega_c = \frac{eB}{m^*c}$ is the cyclotron frequency, and we henceforth refer to the strength of magnetic field in terms of $\omega_c$.

By separation of variables and the substitution of $\psi = \frac{1}{\sqrt{2}}f(r)e^{im\theta}$ into the single-electron Schrödinger equation $\hat{H}\psi = E\psi$, we have

$$\left(-\frac{1}{2r}\frac{\partial}{\partial r}\left(r\frac{\partial}{\partial r}\right) + \frac{m^2}{2r^2} + \frac{\Omega^2 m^{*2}}{2\hbar^2}r^2 - \frac{Em^*}{\hbar^2} - \frac{mm^*\omega_c}{2\hbar}\right)f(r) = 0 \qquad (7)$$

where $\Omega^2 = (\frac{1}{4}\omega_c^2 + \omega_0^2)$.

As $r \to 0$, the above differential equation simplifies to

$$\left(-\frac{1}{2r}\frac{\partial}{\partial r}\left(r\frac{\partial}{\partial r}\right) + \frac{m^2}{2r^2}\right)f(r) = 0 \qquad (8)$$

Substituting $f(r) = r^p$, we obtain

$$\frac{1}{2}(m^2 - p^2)r^{p-2} = 0$$

which has solution $p = \pm m$. As the solution should be finite at the origin, we need $p = |m|$.

As $r \to \infty$, the differential equation becomes

$$\left(-\frac{1}{2r}\frac{\partial}{\partial r}\left(r\frac{\partial}{\partial r}\right) + \frac{k^2}{2}r^2\right)f(r) = 0 \tag{9}$$

the solution of which is

$$f(r) = d_1 I_0\left(\frac{kr^2}{2}\right) + d_2 K_0\left(\frac{kr^2}{2}\right)$$

where $d_1$ and $d_2$ are constants, $k = m^*\Omega/\hbar$, and $I_0$ and $J_0$ are the modified Bessel functions. As $r \to \infty$, $I_0(kr^2/2)$ diverges. Thus, the physical solution is $K_0(kr^2/2)$, which is of the form $e^{-kr^2/2}$ for large $r$.

If we assume a trial wavefunction of the form

$$f(r) = r^m e^{-kr^2/2} g(r)$$

and let

$$\frac{\gamma^2}{2} = \frac{Em^*}{\hbar^2} + \frac{mm^*\omega_c}{2\hbar}$$

then we get a differential equation for $g(r)$,

$$[\gamma^2 - 2k(|m| + 1)]g(r) + (-2kr^2 + 2|m| + 1)g'(r) + rg''(r) = 0$$

The above equation has the general solution

$$g(r) = e_{1 \, 1}F_1\left[\frac{1}{4}\left(-\frac{\gamma^2}{k} - 2|m| + 2\right); 1 - |m|; kr^2\right] + e_{2 \, 1}F_1\left[\frac{1}{4}\left(-\frac{\gamma^2}{k} + 2|m| + 2\right); |m| + 1; kr^2\right] \tag{10}$$

where $e_1$ and $e_2$ are constants and $_1F_1$ is the hypergeometric function. To be able to normalize the solution, the hypergeometric function must be finite. This implies that $_1F_1(-n, b, z)$, where $n$ is an integer and $b \geq 1$. This condition is satisfied by the second hypergeometric function in Eq. (10), as $b = |m| + 1 \geq 1$, but not the first hypergeometric function. This gives us the condition

$$\frac{1}{4}\left(-\frac{\gamma^2}{k} + 2|m| + 2\right) = -n \tag{11}$$

that is,

$$E = (2n + |m| + 1)\hbar\Omega - \frac{1}{2}m\hbar\omega_c \tag{12}$$

The generalized Laguerre functions are defined by

$$L_n^m(kr^2) = \binom{n + m}{n}{}_1F_1(-n, m + 1; kr^2)$$

so the solution to Eq. (7) is

$$f(r) = N_{n,m} r^m e^{-\frac{1}{2}kr^2} L_n^m(kr^2)$$

where the $N_{n,|m|}$ is a normalization constant. After normalizing using the orthogonality relations of the Laguerre polynomials, the final result is

$$\psi_{nm}(r, \theta) = \left[2k^{|m|+1}\frac{n!}{(n + |m|)!}\right]^{1/2} r^{|m|} e^{-\frac{1}{2}kr^2} L_n^{|m|}(kr^2)\frac{1}{\sqrt{2\pi}}e^{im\theta} \tag{13}$$

These quantized energy levels are known as the Fock-Darwin states.

When there is no applied magnetic field, Eq. (12) becomes

$$E_{nm} = \frac{\hbar\omega_0}{2}(2n + |m| + 1)$$

and we can deduce that the $n$th highest energy level is $n$-fold degenerate. This degeneracy is broken by applying a magnetic field, as shown in Fig. 3. In the limit of high magnetic fields, the energy level relation becomes

$$E_{nm} = \frac{\hbar\omega_c}{2}(2n + |m| - m + 1)$$

and we can see Landau bands forming with energies $\frac{\hbar\omega_c}{2}$, $\frac{3\hbar\omega_c}{2}$, $\frac{5\hbar\omega_c}{2}$, ....

### 3.1.2. Finite-Difference Shooting Method

For other forms of confinement potentials, we do not have analytical solutions in general. However, if the confinement potential is circularly symmetric, that is, $V(r, \theta) = V(r)$, by using separation of variables we can reduce the two-dimensional Schrödinger's equation to a one-dimensional radial equation,

$$\left[ -\frac{1}{2r}\frac{d}{dr}\left(r\frac{d}{dr}\right) + \frac{m^2}{2r^2} + V(r)\right]R(r) = ER(r) \tag{14}$$

while the electron wavefunction is

$$\psi(r, \theta) = \frac{1}{\sqrt{2\pi}}e^{im\theta}R(r)$$

By a simple substitution $R(r) \rightarrow \Re(r)/\sqrt{r}$, the above radial equation can be written in the following general form:

$$\left( -\frac{1}{2}\frac{d^2}{dr^2} + V_{eff}(r)\right)\Re(r) = E\Re(r) \tag{15}$$

where the effective potential

$$V_{eff}(r) = V(r) - \frac{1}{8r^2} + \frac{m^2}{2r^2}$$

Equation (15) can be readily solved numerically by the finite-difference shooting method [31, 32], in which the second derivative is approximated by

$$\frac{d^2\Re(r)}{dr^2} \approx \frac{\Re(r_{n+1}) + \Re(r_{n-1}) - 2\Re(r_n)}{\Delta^2} \tag{16}$$



Figure 3. The evolution of single-electron energy levels as the external magnetic field increases.

where $\Re(r_n)$ are the discretized radial wavefunction and $\Delta$ is the size of the numerical grid interval. This leads to the following finite-difference equation:

$$\frac{1}{2}\frac{\Re(r_{n+1}) + \Re(r_{n-1}) - 2\Re(r_n)}{\Delta^2} + V_{\text{eff}}(r_n)\Re(r_n) \approx E\Re(r_n) \tag{17}$$

where the truncation error is on the order $O(\Delta^2)$. Equation (17) can be rearranged so that we evaluate $\Re(r_n + 1)$ using the values of $\Re(r_n)$ and $\Re(r_n - 1)$; that is,

$$\Re(r_{n+1}) \approx 2\Re(r_n) - \Re(r_{n-1}) + 2\Delta^2[V_{\text{eff}}(r_n) - E]\Re(r_n) \tag{18}$$

The numerical procedure is straightforward, assuming two arbitrary initial values for $\Re(r_0)$ and $\Re(r_1)$, and then using Eq. (18) to obtain the other values recursively. The only complication is that we do not know the value of energy $E$. If the wavefunction is unbound, there is a solution for any given energy, which can be obtained iteratively following the above recurrence relation. However, if the wavefunction is bound, the energy $E$ can only take certain quantized values. In this case, the shooting method starts with an estimated $E$ and then systematically changes its value until the final wavefunction satisfies the boundary conditions imposed by the confinement potential.

A Mathematica implementation of the shooting method is given below. We begin by first defining a confinement potential and various parameters for the calculation. We will integrate out to a distance of $x$max with increment $\Delta$.

$$V(x\_) := \frac{k^2 x^2}{2} - \frac{1}{8x^2} + \frac{m^2}{2x^2};$$

$m = 1; \quad k = 1.35;$

$x\text{max} = 5.0; \quad \Delta = 0.05;$

$V\text{values} = \text{Table}[\{x, V(x)\}, \{x, 0.03, x\text{max}, \Delta\}];$

$V\text{Length} = \text{Length}[V\text{values}];$

$\text{LowerEnergy} = 0; \quad \text{UpperEnergy} = 5;$

Equation (18) is then used for the integration; that is,

$$\psi_{n\_,\epsilon\_} := \psi_{n,\epsilon} = 2\{-\epsilon\Delta^2 + V[(n-1)\Delta]\Delta^2 + 1\}\psi_{n-1,\epsilon} - \psi_{n-2,\epsilon}; \tag{19}$$

Next we have the iterator

$\text{Do}[\epsilon = (\text{LowerEnergy} + \text{UpperEnergy})/2;$

$\quad \psi_{1,\epsilon} = 0; \quad \psi_{2,\epsilon} = 1;$

$\quad \psi \; \text{Flatten}[\text{Prepend}[\text{Table}[\psi_{n,\epsilon}, \{n, 3, V\text{Length}\}], \{\psi_{1,\epsilon}, \psi_{2,\epsilon}\}]];$

$\quad \psi_s = \text{Table}[\{n\Delta, \psi_s[[n]]\}, \{n, 1, V\text{Length}\}]; \tag{20}$

$\quad \text{list} = \text{Join}[\text{list}, \{\{\epsilon, \psi_s\}\}];$

$\quad \text{If}[\text{Last}[\psi_s] < 0, \text{UpperEnergy} = \epsilon, \text{LowerEnergy} = \epsilon],$

$\quad \{\text{Counter}, 1, 25\}]$

This searches for the correct energy $\epsilon$ via a bisection method. The implementation here runs for 25 steps. Figure 4 is generated by plotting the elements of *list* every five time steps.

### 3.1.3. Numerov-Cooley Matching Method

The matching method involves integrating a trial wavefunction both inward and outward to some chosen midpoint and matching the value and derivative at that point. The wavefunction can simply be scaled to be continuous, but the derivative can only be matched for certain

**Figure 4.** Convergence of the radial wavefunction calculated by the finite difference shooting method.

values of energy if the wavefunction is bound. The basic algorithm varies the energy as a parameter until the correction to the estimated energy $E$ becomes sufficiently small.

The Numerov-Cooley method uses a higher-order approximation than the finite-difference scheme [33]. The basic formula of this method is the three-term recurrence relation:

$$Y(r_{i+1}) + Y(r_{i-1}) - 2Y(r_i) = 2\Delta^2(V_{eff}(r_i) - E)\Re(r_i),\tag{21}$$

where

$$Y(r_i) = \left(1 - \frac{\Delta^2}{6}(V_{eff}(r_i) - E)\right)\Re(r_i).\tag{22}$$

Again a trial energy is assumed first. The wavefunction is calculated by integration inward from the outer boundary until the absolute value of $Y(r_i)$ stops increasing with decreasing $i$, and this point is set as the matching point $r_m$. This is to make sure that the solution is significantly large at the matching point to guarantee fast convergence. The wavefunction is then integrated outward from the inner boundary to $r_m$ and normalized so that $Y^{out}(r_m) = Y^{in}(r_m)$. The discrepancy between the inward and outward integrations at $r_m$ is used to obtain the correction for the energy eigenvalue:

$$D(E) = \frac{[-Y(r_{m-1}) + Y(r_m) - Y(r_{m+1})/(2\Delta^2)] + [V(r_m) - E]\Re(r_m)}{\sum \Re(r_i)^2}.\tag{23}$$

This calculation is repeated several times until $D(E)$ is sufficiently small. Then the correct wavefunction is obtained. The node number of the wavefunction can be counted easily to

determine the quantum number of the obtained wavefunction. The truncation error for the Numerov-Cooley method is of the order $O(\Delta^4)$, providing a much-improved numerical accuracy. By using Eq. (23), the convergence rate is also dramatically increased in comparison with the standard shooting method.

We begin the Mathematica implimentation by first defining the confinement potential and the various parameters used for the calculations.

$$V[r\_] := \frac{k^2 r^2}{2} - \frac{1}{8r^2} + \frac{m^2}{2r^2}; \quad k = 1.35; \quad m = 1;$$

$$\text{num} = 1000; \quad r = 4; \quad \Delta = \frac{r}{\text{num}};$$

$$x\text{grid} = \text{Table}[q\,l, \{q, 1, \text{num}\}]; \quad V\text{grid} = \frac{V}{(a\,x\text{grid})}; \quad V_{i\_} := V\text{grid}[[i]];$$

$$\chi = \text{Table}[0, \{\text{num}\}]; \quad \chi[[1]] = 0; \quad \chi[[2]] = 10^{-5}; \quad \chi[[\text{num}]] = 10^{-5}; \quad \chi[[\text{num} - 1]] = 10^{-4};$$

Equation (21) is then used for the integration:

$$Y_{i\_} := \left[ 1 - \Delta^2 \frac{(V_i - \ell)}{6} \right] \chi[[i]];$$

with initial energy

$$\ell = 2;$$

The iterator is

list = {};

$$\text{Do}[\text{mid} = \text{Catch}[\text{Do}[\chi[[i]] = \frac{2(V_{i-1} - \ell)\chi[[i + 1]]\Delta^2 + 2Y_{i+1} - Y_{i-2}}{1 - \Delta^2(V_i - \ell)/6};$$

$$\text{If}[\chi[[i]] < \chi[[i + 1]], \text{Throw}(i)], \{i, \text{num} - 2, 1, -1\}] + 1;$$

$$\text{Do}[\chi[[i]] = \frac{\chi[[i]]}{\chi[[\text{mid}]]}, \{i, \text{num}, \text{mid}, -1\}];$$

$$\text{Do}[\chi[[i]] = \frac{2(V_{i-1} - \ell)\chi[[i - 1]]\Delta^2 + 2Y_{i-1} - Y_{i-2}}{1 - \Delta^2(V_i - \ell)/6}, \{i, 3, \text{mid}\}]; \qquad (24)$$

$$\text{Do}[\chi[[i]] = \frac{\chi[[i]]}{\chi[[\text{mid}]]}, \{i, 1, \text{mid}\};$$

$$\text{list} = \text{Join}[\text{list}, \{\{\ell, \chi\}\}];$$

$$\ell\!j = \frac{1}{\text{Plus}(a\,a\,\chi^2} \left( \frac{1}{2l^2}(-Y_{\text{mid}-1} + 2Y_{\text{mid}} - Y_{\text{mid}+1}) + (V_{\text{mid}} - \ell)\chi[[\text{mid}]] \right);$$

$$\ell = \ell + \ell\!j.$$

{counter, 1, 4}]

where Eq. (23) is used to correct the estimated energy at each iteration. Convergence can be achieved usually within a couple of iterations, and the optimized wavefunctions are shown in Fig. 5 for each step.

### 3.1.4. Direct Matrix Method

The finite-difference shooting method and the Numerov-Cooley matching method can be used to solve one-dimensional Schrödinger equations, but they are not suitable for two or more dimensions. In a one-dimension problem, the boundary condition usually concerns only two points, which can be handled readily by the shooting or matching methods. In higher dimensions, however, the boundary condition is imposed on lines or surfaces, making the shooting or matching methods inappropriate [34].

**Figure 5.** Convergence of the radial wavefunction calculated by the Numerov-Cooley matching method.

A direct matrix method can be used to solve general differential eigenvalue equations in two or higher dimensions, where the wavefunction as well as the Laplacian operator and the potential energy operator are mapped onto a multidimensional grid. For example, in two dimensions, the Schrödinger equation can be written as the following, with the second derivatives in the approximate finite difference form:

$$
-\frac{1}{2}\left[\frac{\partial^2\psi(x, y)}{\partial x^2} + \frac{\partial^2\psi(x, y)}{\partial y^2}\right] + V(x, y)\psi
$$

$$
\approx -\frac{1}{2}\left[\frac{\psi(i+1, j) + \psi(i-1, j) - 2\psi(i, j)}{(\Delta x)^2} + \frac{\psi(i, j+1) + \psi(i, j-1) - 2\psi(i, j)}{(\Delta y)^2}\right]
$$

$$
+ V(i, j)\psi(i, j)
$$

$$
\approx E\psi(i, j). \tag{25}
$$

Finding the solution at each point $(i, j)$ of the numerical lattice gives a system of algebraic equations, which can be written in matrix form as

$$
\begin{bmatrix} & \vdots & \\ 0 \cdots 0 -1\ 0 \cdots 0 -1\ 4 - \dfrac{2m}{\hbar^2}\Delta^2(V_{i,j} - E)\ -1\ 0 \cdots 0 -1\ 0 \cdots 0 \\ & \vdots & \end{bmatrix} \begin{bmatrix} \vdots \\ \psi_{i,j-1} \\ \vdots \\ \psi_{i-1,j} \\ \psi_{i,j} \\ \psi_{i+1,j} \\ \vdots \\ \psi_{i,j+1} \\ \vdots \end{bmatrix} = 0 \tag{26}
$$

where we assume $\Delta x = \Delta y = \Delta$.

The above matrix equation can be solved by many well-established methods (see, e.g., Refs. [35, 36]). Note that the matrix on the left side of the above equation is sparse, with

many zero-valued elements. Making use of this information can save on storage requirements on the computer and on overall time needed to find the eigenvalues and eigenvectors. In addition, we are frequently only interested in the lowest eigenfunctions of a system.

The algorithm we use for solving the eigenvalue problems encountered in this work is the implicitly restarted Arnoldi factorization algorithm, as implemented in the set of library routines ARPACK [37]. In this method, a set of Schur vectors is calculated that gives rise to approximate eigenvalues and eigenvectors of the original matrix. The library is written in such a way that when the main routine is called, it returns the user with a vector. The user is requested to multiply the vector by the matrix for which the eigenvalues are required and then re-call the same routine. This process continues until convergence is achieved. When the library determines that a reasonable set of Schur vectors has been calculated, a separate routine is used to calculate the eigenvalues and eigenvectors from the previous iterative refinement.

The ARPACK library is most suitable for use on "sparse" matrices, in which the definition of "sparse" is such that the multiplication of a vector by the matrix is an order $O(n)$ operation. That is, if the vector has $n$ elements, and the matrix has $n^2$ elements, most of the elements of the matrix are zero, so that the number of multiplications and additions required to multiply the vector by the matrix is only proportional to $n$ (multiplying a dense matrix with $n^2$ nonzero elements is an order $O[n^2]$ process). Thus, the time required to complete the calculation is proportional to the length of the eigenvectors and the number of eigenvectors required. Therefore, the ARPACK library is most advantageous when only the first few eigenvectors are required, as is the case, for example, when computing the first few eigenfunctions of a Hamiltonian.

This method was used to calculate the energy structures of single-electron quantum dot systems with five different geometries; namely, circular, elliptic, triangular, square, and annular ring [18]. It is worth noting that, in a circularly symmetric quantum dot, the angular momentum operator commutes with the system Hamiltonian, and thus the energy eigenstates are simultaneous eigenstates of the angular momentum. In this case, the energy and angular momentum quantum numbers $(n, m)$ can be used to uniquely identify each eigenstate of the system. However, when the circular symmetry is broken and the confinement potential is of an arbitrarily complex geometry, we do not really have a general equivalent set of quantum numbers that can be used to uniquely identify each eigenstate. For example, one can use the number of nodes $(nx, ny)$ in the $(x, y)$ direction to identify the first five eigen states in the elliptic dot as $(nx, ny) = (0, 0), (1, 0), (0, 1), (2, 0),$ and $(1, 1)$, respectively. However, such classification does not work for the triangular, the square, or the annual ring dot. For each individual quantum dot with a certain degree of symmetry, it is possible to find some operators that would commute with the system Hamiltonian, but this cannot be generalized.

Presented in Figs. 6, 7, 8, and 9 are the electron densities of the wavefunctions that are the eigenstates of the Hamiltonian operator only, which nevertheless reflect the symmetry of the quantum dots under study. The first panel of each of these figures shows the confinement potential of each dot. The other panels illustrate the charge density distributions of the first five eigenstates of the system, which reflect the symmetry of the quantum dot. The corresponding energy levels are listed in Table 1. As shown, the degeneracy of the levels in the circular dot has been lifted to a certain extent in the other four quantum dots due to their lower symmetries.

In principle, this matrix approach is completely general and can be applied to arbitrarily high dimensions. However, it is not practical for problems in three or higher dimensions, nor for a two-dimensional problem in which a large numerical grid is needed. This is simply because the matrix size grows rapidly as the problem dimension or the number of grid points in each dimension increases. For example, if we use 100 points in each direction of the numerical grid, the matrix size is $10,000 \times 10,000$ for a two-dimensional problem. In three dimensions, the matrix size would be $10^{12}$.

## 3.1.5. Rayleigh-Ritz Variational Method

The Rayleigh-Ritz variation principal can be used to solve general differential eigenvalue equations in two or higher dimensions by using a trial solution that is a linear combination

**Figure 6.** An elliptic quantum dot. First panel: confinement potential; other panels: charge density distributions of the first five eigen-states.

of basis functions. The problem is again converted to a matrix eigenvalues problem, which solves for the coefficients of the linear combination. However, the size of the matrix is determined by the number of basis functions required in the expansion, which is in general drastically smaller than that required in direct matrix method.

For completeness, we present the basic theory here. We begin with an equation of the form

$$\hat{H} u_n(\mathbf{x}) = \lambda_n u_n(\mathbf{x}) \tag{27}$$

in which $\hat{H}$ is the system Hamiltonian in the Schrödinger equation (or any other linear differential operator), $u_n(\mathbf{x})$ are its eigenfunctions, and $\lambda_n$ are the corresponding eignevalues. The problem can be formulated as a functional equation,

$$\lambda[\psi] = \frac{\int \psi^*(\mathbf{x})\hat{H}\psi(\mathbf{x})d\mathbf{x}}{\int \psi^*(\mathbf{x})\psi(\mathbf{x})d\mathbf{x}} \tag{28}$$

It is easy to see that if $\psi$ is an eigenfunction of $\hat{H}$, $\psi = u_n$, then $\lambda[\psi]$ will be the eigenvalue of $\lambda_n$. It can also be shown that if some function $\psi$ causes $\lambda[\psi]$ to become stationary, that is,

$$\frac{\delta\lambda[\psi]}{\delta\psi} = 0 \tag{29}$$

then $\psi$ is an eigenfunction of $\hat{H}$ [38].

**Figure 7.** A triangular quantum dot. Same as Fig. 6.

The linear variational method expands $\psi$ in terms of a complete set of basis functions $\nu_i(\mathbf{x})$ so that

$$\psi(x) = \sum_{i=1}^{\infty} c_i \nu_i(x) \equiv \mathbf{c}^{\mathsf{T}} \nu(\mathbf{x}) \tag{30}$$

where $\{\mathbf{c}\}_i = c_i$ and $\{\nu(\mathbf{x})\}_i = \nu_i(\mathbf{x})$. Eq. (28) becomes

$$\lambda[\psi] = \frac{\sum_{i,j} c_i^* c_j \int \nu_i^*(\mathbf{x}) \hat{H} \nu_j(\mathbf{x}) d\mathbf{x}}{\sum_{i,j} c_i^* c_j \int \nu_j^*(\mathbf{x}) \nu_i(\mathbf{x}) d\mathbf{x}} \tag{31}$$

Introducing the Ritz matrix $\mathcal{H}$,

$$\{\mathcal{H}\}_{i,j} = \int \nu_i^*(\mathbf{x}) \hat{H} \nu_j(\mathbf{x}) d\mathbf{x} \tag{32}$$

and the overlap matrix $\mathcal{C}$

$$\{\mathcal{C}\}_{i,j} = \int \nu_i^*(\mathbf{x}) \nu_j(\mathbf{x}) d\mathbf{x} \tag{33}$$

Eq. (31) reads

$$\lambda[\psi] = \frac{\mathbf{c}^{\mathsf{T}} \mathcal{H} \mathbf{c}}{\mathbf{c}^{\mathsf{T}} \mathcal{C} \mathbf{c}} \tag{34}$$

**Figure 8.** A square quantum dot. Same as Fig. 6.

A variation of Eq. (31) with respect to the expansion coefficients

$$\frac{\delta \lambda[\psi]}{\delta c_i} = 0 \tag{35}$$

leads to the generalized matrix eignevalue equation, known as the Ritz matrix equation

$$H \mathbf{c} = \lambda^{\ell} \mathbf{c} \tag{36}$$

In general, $H$ and $\ell$ are infinite dimensional matrices. Truncating $H$ and $\ell$ to $\Lambda \times \Gamma$ matrices, the resulting $\Lambda$-dimensional matrix eignevalue problem,

$$(\ell^{-1} H)\mathbf{c}_i = \lambda \mathbf{c}_i \tag{37}$$

can be solved for nonsingular $\ell$. As a result of this truncation, the eigenvectors and eigenvalues are now approximations of the actual solutions,

$$\psi(\mathbf{x}) \simeq \psi_i(\mathbf{x}) = \sum_{i=1} c_i \nu_i(\mathbf{x}) = \mathbf{c}^T \boldsymbol{\nu}_i(\mathbf{x}) \tag{38}$$

As we increase $\Lambda$, we converge to the exact eignevectors and eigenvalues, with the condition that

$$\lambda_{\text{exact}} \leq \lambda[\psi_i] \tag{39}$$

implying that $\lambda[\psi_i]$ is an upper bound for the exact eigenvalues of $\tilde{H}$ [38].

**Figure 9.** An annular ring quantum dot. Same as Fig. 6.

The Ritz matrix Eq. (36) for a single-electron system can be readily constructed and solved, for example, by the ARPACK package, as described in the previous section. For multielectron systems, various computational schemes have been developed to obtain their wavefunctions and energy structures, which will be the topic of Section 4. Central to these schemes is the construction of two-particle interaction integrals, which is discussed in the following section.

## 3.2. Two-Particle Interaction Integrals

To construct the Ritz matrix for a multielectron system, we need to calculate inner products of the following form

$$H_{\alpha\beta}^{\gamma\eta} \equiv \langle \psi_\alpha(x, y)\psi_\beta(x', y') | \hat{H}^{int} | \psi_\gamma(x, y)\psi_\eta(x', y') \rangle \tag{40}$$

which is an integration in four dimensions and can be very time consuming. If we perform this integration as a Reiman sum, where each $\psi_i$ is represented by a grid of $n_x \times n_y$ points, we would need to evaluate the integrand $n_x^4 \times n_y^4$ times for each possible combination of

**Table 1.** Single-electron energy levels (a.u.).

| $n$th level | Circular | Elliptic | Triangular | Square | Ring |
|---|---|---|---|---|---|
| 1st | 1.0 | 1.1 | 0.9933 | 0.9954 | 1.005 |
| 2nd | 2.0 | 2.1 | 1.9604 | 1.977 | 1.121 |
| 3rd | 2.0 | 2.3 | 1.9604 | 1.977 | 1.121 |
| 4th | 3.0 | 3.1 | 2.879 | 2.761 | 1.450 |
| 5th | 3.0 | 3.3 | 2.974 | 2.929 | 1.472 |
| 6th | 3.0 | 3.5 | 2.974 | 3.194 | 1.996 |
| 7th | 4.0 | 4.1 | 3.852 | 3.698 | 1.996 |
| 8th | 4.0 | 4.3 | 3.852 | 3.698 | 2.688 |
| 9th | 4.0 | 4.5 | 3.966 | 4.197 | 2.702 |
| 10th | 4.0 | 4.7 | 4.008 | 4.197 | 3.531 |

$\alpha, \beta, \gamma, \eta$. Furthermore, the integrand has a singularity at the origin, where $r = r'$. It is an integrable singularity, but this necessitates a large number of grid points, especially close to the singularity, to achieve reasonable numerical accuracy.

In the following text, we discuss two methods developed for accurate evaluation of the two-particle interaction integrals. The first method reduces the four-dimension integration into a sum of one-variable integrals that can be evaluated algebraically. The second method casts the problem in the form of a Poisson's equation whose solution is the integration result.

### 3.2.1. Analytical Formulas

If Eq. (13) is used as the basis functions, we can actually evaluate the two-particle interaction integrals analytically to arbitrary precision. In this case, the integration can be written explicitly, in polar coordinates, as

$$H_{mnpq}^{abcd} \equiv \langle \psi_{mn}(x, y)\psi_{pq}(x', y')|\hat{H}^{int}|\psi_{ab}(x, y)\psi_{cd}(x', y')\rangle$$

$$= \int_0^\infty \int_0^\infty \int_0^{2\pi} \int_0^{2\pi} rs \frac{\psi_{nm}^*(r, \theta)\psi_{pq}^*(s, \phi)\psi_{ab}(r, \theta)\psi_{cd}(s, \phi)}{\sqrt{r^2 + s^2 - 2rs\cos(\theta - \phi)}} d\theta\, d\phi\, dr\, ds$$

$$= \int_0^\infty \int_0^\infty \int_0^{2\pi} \int_0^{2\pi} rs \frac{R_{nm}(r)R_{pq}(s)R_{ab}(r)R_{cd}(s)e^{i(b-m)\theta + i(d-q)\phi}}{4\pi^2\sqrt{r^2 + s^2 - 2rs\cos(\theta - \phi)}} d\theta\, d\phi\, dr\, ds \quad (41)$$

where $R_{nm}(r) = [2k^{|m|+1}n!/(n + |m|)!]^{1/2} r^{|m|} e^{-(k/2)r^2} L_n^{|m|}(kr^2)$.

Noting that the integrand is periodic in both $\theta$ and $\phi$ with period $2\pi$, we can apply the transformation $\phi \rightarrow \theta + \delta$, or equivalently, $\theta \rightarrow \phi + \delta$ with delta ranging from 0 to $2\pi$. These, respectively, give

$$H_{mnpq}^{abcd} = \frac{1}{4\pi^2} \int_0^{2\pi} e^{i(b-m+d-q)\theta} d\theta \int_0^\infty \int_0^\infty \int_0^{2\pi} rs \frac{R_{nm}(r)R_{pq}(s)R_{ab}(r)R_{cd}(s)e^{i(d-q)\delta}}{\sqrt{r^2 + s^2 - 2rs\cos(\delta)}} d\delta\, dr\, ds$$

$$\tag{42}$$

and

$$H_{mnpq}^{abcd} = \frac{1}{4\pi^2} \int_0^{2\pi} e^{i(b-m+d-q)\phi} d\phi \int_0^\infty \int_0^\infty \int_0^{2\pi} rs \frac{R_{nm}(r)R_{pq}(s)R_{ab}(r)R_{cd}(s)e^{-i(m-b)\delta}}{\sqrt{r^2 + s^2 - 2rs\cos(-\delta)}} d\delta\, dr\, ds$$

$$\tag{43}$$

Evaluating the single integral gives us

$$\int_0^{2\pi} e^{i(b-m+d-q)\theta} d\theta = \begin{cases} 0, & b - m + d - q \neq 0 \\ 2\pi, & b - m + d - q = 0 \end{cases} \tag{44}$$

This means the interaction integral is zero if the sum of the angular momentum quantum numbers is different for the two orbitals (i.e., $b + d \neq m + q$).

Using the fact that $\cos(-x) = \cos(x) \ \forall x \in \mathbb{R}$, and equating the two integrals (42) and (43), the complex part of the integral vanishes and we can replace $e^{i(d-q)\delta}$ with $\cos(|(d - q)|\delta)$. If we define a new constant $\alpha = |b - m| = |d - q|$, the integral (41) becomes

$$H_{mnpq}^{abcd} = \frac{1}{2\pi} \int_0^\infty \int_0^\infty \int_0^{2\pi} rs \frac{R_{nm}(r)R_{pq}(s)R_{ab}(r)R_{cd}(s)\cos(\alpha\delta)}{\sqrt{r^2 + s^2 - 2rs\cos(\delta)}} d\delta\, dr\, ds \quad (45)$$

We can break this integral down further by noting that

$$rsR_{nm}(r)R_{pq}(s)R_{ab}(r)R_{cd}(s) = 4\sqrt{k}^{|b|+|d|+|m|+|q|+4} e^{-k(r^2+s^2)} r^{|m|+|b|+1} s^{|q|+|d|+1} \frac{1}{|b|!|d|!|m|!|q|!}$$

$$\times \sqrt{\frac{(a + |b|)!}{a!}} \sqrt{\frac{(c + |d|)!}{c!}} \sqrt{\frac{(n + |m|)!}{n!}} \sqrt{\frac{(p + |q|)!}{p!}}$$

$$\times {}_1F_1(-a, 1 + |b|, kr^2) {}_1F_1(-c, 1 + |d|, ks^2)$$

$$\times {}_1F_1(-n, 1 + |m|, kr^2) {}_1F_1(-p, 1 + |q|, ks^2)$$

where $_1F_1(-n, 1 + |m|, kr^2) = \sum_{i=0}^{n}\{(-n)_i/[(1 + |m|)_i i!]\}(kr^2)^i$, and $(a)_i = a(a + 1)(a + 2)\cdots(a + i - 1)$ is the Pochhammer symbol or rising factorial. The expansion is finite, as $(-n)_i$ would be zero for $i \geq n + 1$.

It is clear then that the original integral (41) is simply a linear combination of integrals of the following form:

$$\int_0^\infty \int_0^\infty \int_0^{2\pi} \frac{\cos(\alpha\delta)e^{-k(r^2 + s^2)}k(\sqrt{k})^{a+b}r^a s^b}{\sqrt{r^2 + s^2 - 2rs\cos(\delta)}}\, dr\, ds\, d\delta. \tag{46}$$

Because any power of $r$ and $s$ coming from the expansion of the hypergeometric function will be even, we only have to worry about the powers of $r$ and $s$ coming from $r^{|m|+|b|+1}s^{|q|+|d|-1}$. Define the function $\sigma(n)$ by

$$\sigma(n) = \begin{cases} o, & n \text{ odd;} \\ e, & n \text{ even.} \end{cases} \tag{47}$$

This function has the property that $\sigma(n + m) = \sigma(\pm n \pm m)$. Remembering that $b - m + d - q = 0$, then $\sigma(b - m + d - q) = e \Rightarrow \sigma(|b| + |m| - (|d| + |q|)) = e$, so $(|b| + |m| + 1) - (|d| + |q| + 1)$ is even. This means $a - b$ will always be even.

We apply another transformation $r \to \rho\sin(\lambda)$ and $s \to \rho\cos(\lambda)$ with $\rho \in [0, \infty]$ and $\lambda \in [0, 2\pi]$. This gives us

$$\int_0^\infty \int_0^{\pi/2} \int_0^{2\pi} \rho \frac{\cos(\alpha\delta)e^{-k\rho^2}k(\sqrt{k})^{a+b}\rho^{a+b}\sin^a(\lambda)\cos^b(\lambda)}{\sqrt{\rho^2\sin^2(\lambda) + \rho^2\cos^2(\lambda) - 2\rho^2\sin(\lambda)\cos(\lambda)\cos(\delta)}}\, d\delta\, d\rho\, d\lambda$$

$$= \int_0^\infty e^{-k\rho^2}k(\sqrt{k})^{a+b}\rho^{a+b}\, d\rho \int_0^{\pi/2} \int_0^{2\pi} \frac{\cos(\alpha\delta)\sin^a(\lambda)\cos^b(\lambda)}{\sqrt{1 - \sin(2\lambda)\cos(\delta)}}\, d\delta\, d\lambda. \tag{48}$$

Now $\int_0^\infty e^{-k\rho^2}k(\sqrt{k})^{a+b}\rho^{a+b}\, d\rho = 1/2\sqrt{k}\,\Gamma((1 + a + b)/2)$, where $\Gamma$ is the Euler gamma function. Thus, our integral becomes:

$$\frac{1}{2}\sqrt{k}\,\Gamma\left(\frac{1 + a + b}{2}\right)\int_0^{\pi/2} \int_0^{2\pi} \frac{\cos(\alpha\delta)\sin^a(\lambda)\cos^b(\lambda)}{\sqrt{1 - \sin(2\lambda)\cos(\delta)}}\, d\delta\, d\lambda. \tag{49}$$

The important point to note here is that the integral is proportional to $\sqrt{k}$ for all the interaction terms, where $k$ is the effective harmonic well constant. This allows us to calculate the interaction integrals for $k = 1$ and then use these for any other value of $k$ by multiplying by an appropriate factor. This integral has a logarithmic singularity at $\lambda = \frac{\pi}{4}$, which results from the Coulomb interaction diverging when the two electrons are in the same location. However, this integral can be done in closed form, noting that the integrand is symmetric in $a$ and $b$, and furthermore, that $a$ and $b$ can only differ by an even number.

If we let $Z_0$ represent the part of the integrand that does not involve $\sin^a(\lambda)\cos^b(\lambda)$ and assume, without loss of generality, that $b \geq a$ and let $b - a = 2n$, we get for $n = 1$

$$\int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^{a+2}(\lambda)\,d\lambda = \int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^a(\lambda)(1 - \sin^2(\lambda))\,d\lambda$$

$$= \int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^a(\lambda)\,d\lambda - \int_0^{\pi/2} Z_0\sin^{a+2}(\lambda)\cos^a(\lambda)\,d\lambda \tag{50}$$

Because $\int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^{a+2}(\lambda)\,d\lambda = \int_0^{\pi/2} Z_0\sin^{a+2}(\lambda)\cos^a(\lambda)\,d\lambda$, we have

$$\int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^{a+2}(\lambda)\,d\lambda = \frac{1}{2}\int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^a(\lambda)\,d\lambda \tag{51}$$

Let $Z_n(a) = \int_0^{\pi/2} Z_0\sin^a(\lambda)\cos^{a+2n}(\lambda)\,d\lambda$, and we can use the same working as before to show that for $n \geq 2$,

$$Z_n(a) = Z_{n-1}(a) - Z_{n-1}(a + 2) \tag{52}$$

It can also be easily verified that

$$Z_0(a) = \frac{\pi^{3/2}\Gamma((1+\alpha+a)/2)}{2^{2\alpha}\,{}^a\Gamma((2+\alpha+a)/2)}\left(\frac{2\alpha}{\alpha}\right)\,{}_3F_2\left(\left\{\frac{1+2\alpha}{4},\frac{3+2\alpha}{4},\frac{1+\alpha+a}{2}\right\},\left\{1+\alpha,\frac{2+\alpha+a}{2}\right\},1\right)$$

(53)

and

$$Z_1(a) = \frac{1}{2}Z_0(a)$$

(54)

The above recursion relations (52–54) are sufficient to calculate the integral in terms of a sum of hypergeometric functions. *Mathematica* evaluates the hypergeometric functions exactly, which allows the exact calculation of interaction elements. This eliminates numerical problems that may arise in evaluating these integrals.

### 3.2.2. Poisson-Fourier Approach

For arbitrary basis functions, we do not have a general analytic formula for the interaction integral. However, we can cast the problem in the form of Poisson's equation, which can be solved by using Fourier transformations. In this way, the four-dimensional integration is reduced to two-dimensional Fourier transformations and a two-dimensional Reiman sum. More important, the singularity is removed from the integration.

In Cartesian coordinates, the interaction integral is

$$\langle\psi_\alpha(x,y)\psi_\beta(x',y')|\hat{H}^{int}|\psi_\gamma(x,y)\psi_\eta(x',y')\rangle$$

$$= \int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\frac{\psi_\alpha^*(x,y)\psi_\gamma(x,y)}{\sqrt{(x-x')^2+(y-y')^2}}\,dx\,dy\,\psi_\beta^*(x',y')\psi_\eta(x',y')\,dx'\,dy'$$

(55)

The essential step is to evaluate the following two-dimensional integral,

$$\int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\frac{\psi_\alpha^*(x,y)\psi_\gamma(x,y)}{\sqrt{(x-x')^2+(y-y')^2}}\,dx\,dy$$

(56)

for all possible values of $\alpha$ and $\gamma$, which contains the singularity. Directly using Reiman sums to evaluate the above integrals may give rise to inaccurate results unless a very large number of points are used close to the singularity. Alternatively, these integrals can be related to solutions to Poisson's equation in the following way.

The development of this technique is inspired by the fact that the three-dimensional integral

$$u(x,y,z) = \int\int\int\frac{\rho(x',y',z')}{\sqrt{(x-x')^2+(y-y')^2+(z-z')^2}}\,dx'\,dy'\,dz'$$

(57)

actually satisfies the Poisson's equation

$$\nabla^2 u(x,y,z) = \rho(x,y,z)$$

(58)

This equation can be solved accurately by standard numerical techniques such as the Fourier and cyclic reduction method described in Numerical Recipes [39].

To relate the two-dimensional integrals given by Eq. (56) to the Poisson's equation (58), we assume

$$\rho(x,y,z) = \rho_{\alpha\gamma}f(z)$$

(59)

where $\rho_{\alpha\gamma} = \psi_\alpha(x,y)\psi_\gamma(x,y)$ and $f(z)$ can be taken as the Dirac delta function (i.e., the electron is confined to an infinitely narrow sheet). Taking a discrete Fourier transformation

in $x$ and $y$ and a continuous Fourier transformation in $z$ of both $\rho(x, y, z)$ and $u(x, y, z)$, we have

$$\rho(x, y, z) = \sum_{k_x=0}^{n_x-1} \sum_{k_y=0}^{n_y-1} \rho_{k_x, k_y} e^{i2\pi k_x x/n_x} e^{i2\pi k_y y/n_y} \left[ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} F(k_z) e^{ik_z z} dk_z \right] \qquad (60)$$

and

$$u(x, y, z) = \sum_{k_x=0}^{n_x-1} \sum_{k_y=0}^{n_y-1} \left[ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} U(k_x, k_y, k_z) e^{ik_z z} dk_z \right] e^{i2\pi k_x x/n_x} e^{i2\pi k_y y/n_y} \qquad (61)$$

where $\rho_{k_x, k_y}$, $F(k_z)$, and $U(k_x, k_y, k_z)$ are the Fourier transforms of $\rho_{x, y}$, $f(z)$, and $u(x, y, z)$, respectively.

We can now apply the operator $\nabla^2$ to $u(x, y, z)$ and obtain

$$\nabla^2 u(x, y, z) = \sum_{k_x=0}^{n_x-1} \sum_{k_y=0}^{n_y-1} \left[ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (-k_x^2 - k_y^2 - k_z^2) U(k_x, k_y, k_z) e^{ik_z z} dk_z \right] e^{i2\pi k_x x/n_x} e^{i2\pi k_y y/n_y} \qquad (62)$$

Equating terms between Eqs. (60) and (62), we have

$$(-k_x^2 - k_y^2 - k_z^2) U_{k_x, k_y}(k_z) = \rho_{k_x, k_y} F(k_z) \qquad (63)$$

that is,

$$U(k_x, k_y, k_z) = \frac{\rho_{k_x, k_y} F(k_z)}{(-k_x^2 - k_y^2 - k_z^2)} \qquad (64)$$

Because we are only interested in the $z = 0$ plane and $F(k_z)$ is the Fourier transform of the Dirac delta function, by using the residue theorem, Eq. (61) becomes

$$u(x, y, z = 0) = \sum_{k_x=0}^{n_x-1} \sum_{k_y=0}^{n_y-1} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{\rho_{k_x, k_y} F(k_z)}{(-k_x^2 - k_y^2 - k_z^2)} dk_z e^{i2\pi k_x x/n_x} e^{i2\pi k_y y/n_y}$$

$$= \sum_{k_x=0}^{n_x-1} \sum_{k_y=0}^{n_y-1} \frac{\rho_{k_x, k_y}}{\sqrt{2\pi(k_x^2 + k_y^2)}} e^{i2\pi k_x x/n_x} e^{i2\pi k_y y/n_y} \qquad (65)$$

It is important to note that the discrete Fourier transformation implies periodic boundary conditions on the source term $\rho(x, y, z)$, which means the last point of the data set in each coordinate needs to be connected smoothly to the starting point in that coordinate. If $\rho(x, y, z)$ is close to zero at the edges of the numerical grid, such a boundary condition is then satisfied. It is therefore necessary to make sure the grid space is sufficiently large to have accurate solutions to the Poisson's equation.

Also note that the Poisson-Fourier approach removes the singularity at $x = x'$ and $y = y'$ in the interaction integral. However, we have a different singularity now at $k_x = k_y = 0$. This term corresponds to a zero frequency: that is, an addition of a constant component in $u(x, y, z)$. We can exclude this term in the summation, and then it is a simple matter to correct for this by evaluating the value of the integral via a Reiman sum at only a few points. The result of the Reiman sum can be compared with the results calculated using the Poisson-Fourier approach, and an appropriate shift is then applied. The use of a number of points (typically the four corners of the grid) means that any minor errors in the shift should average out.

In this way, the interaction integral can be evaluated by two Fourier transformations and a two-dimensional Reiman sum, instead of a four-dimensional Reiman sum, offering a significant time saving. The Poisson-Fourier approach is also much more accurate than the direct Reiman sums when using a moderate number of numerical points, as the singularity has been removed from the integration. Figure 10 uses a gaussian source potential to test

**Figure 10.** A cross section of the integrals $u(x, y, z)$ obtained by the direct Reiman sums and the Poisson-Fourier technique. Solid line: Poisson-Fourier 64 points; dashed line: Reiman sum 64 points; dashed-dotted line: Reiman sum 128 points; dotted line: Reiman sum 256 points.

the accuracy of the Poisson-Fourier method against evaluating the integral directly via a Reiman sum. As shown, a large number of points are required to achieve the same accuracy if the direct Reiman sum method is applied.

## 3.3. Numerical Differentiation Using Fourier Transformations

One of the tasks involved in quantum dot calculations is evaluating the derivatives of a function represented on a numerical grid. An example of this is solving the single-electron Schrödinger equation for a given potential. In this case, we need to evaluate the action of the Laplacian operator $\nabla^2$ on a grid of points accurately and efficiently.

The most common way to evaluate the derivative of a data set given on a numerical grid is the use of finite difference methods. For example, the 5-point finite difference method (FD5) for the second derivative at grid point $m$ has a formula

$$\left.\frac{d^2 f}{dx^2}\right|_m = \frac{1}{12\Delta x^2}(-f_{m-2} + 16f_{m-1} - 30f_m + 16f_{m+1} - f_{m+2}) \tag{66}$$

where $\Delta x$ is the grid spacing. If higher accuracy is required, the 7-point finite-difference method (FD7) can be applied; that is,

$$\left.\frac{d^2 f}{dx^2}\right|_m = \frac{1}{180\Delta x^2}(2f_{m-3} - 27f_{m-2} + 270f_{m-1} - 490f_m + 270f_{m+1} - 27f_{m+2} + 2f_{m+3}) \tag{67}$$

However, finite-difference methods are based on local approximations to the derivative operator, which brings about error. Also, the convergence with decreasing grid spacing is slow.

A more accurate and computationally efficient scheme [40], which is based on the discrete Fourier transform (DFT), uses the fact that if

$$f(x_1, x_2, \ldots, x_L) = \int \int \cdots \int F(k_1, k_2, \ldots, k_L) e^{2\pi i(k_1 x_1 + k_2 x_2 + \cdots + k_L x_L)} dk_1 dk_2 \cdots dk_L, \tag{68}$$

where $F(k_1, k_2, \ldots, k_L)$ is the Fourier transformation of $f(x_1, x_2, \ldots, x_L)$, then the partial derivaties are

$$f^{(n_1, n_2, \ldots, n_L)}(x_1, x_2, \ldots, x_L) = \int \int \cdots \int ((2\pi i k_1)^{n_1} (2\pi i k_2)^{n_2} \cdots (2\pi i k_L)^{n_L}) F(k_1, k_2, \ldots, k_L)$$

$$\times e^{2\pi i(k_1 x_1 + k_2 x_2 + \cdots + k_L x_L)} dk_1 dk_2 \ldots dk_L. \tag{69}$$

It follows that the derivaties can be evaluated by first calculating $F(k_1, k_2, \ldots, k_L)$, using the DFT, and second, taking the inverse DFT of the product:

$$(2\pi i k_1)^{n_1} (2\pi i k_2)^{n_2} \cdots (2\pi i k_L)^{n_L} F(k_1, k_2, \ldots, k_L). \tag{70}$$

| $\Delta x$ | FD5 | FD7 | DFT |
|---|---|---|---|
| 0.01 | $1.3 \times 10^{-5}$ | $5.0 \times 10^{-11}$ | $7.0 \times 10^{-11}$ |
| 0.02 | $2.1 \times 10^{-7}$ | $1.9 \times 10^{-9}$ | $1.2 \times 10^{-11}$ |
| 0.05 | $8.3 \times 10^{-6}$ | $4.7 \times 10^{-8}$ | $2.0 \times 10^{-12}$ |
| 0.1 | $1.3 \times 10^{-4}$ | $2.9 \times 10^{-6}$ | $5.9 \times 10^{-11}$ |
| 0.2 | $2.0 \times 10^{-3}$ | $1.7 \times 10^{-4}$ | $4.1 \times 10^{-14}$ |
| 0.3 | $9.0 \times 10^{-3}$ | $1.6 \times 10^{-3}$ | $8.1 \times 10^{-15}$ |

For $N$ data points of a one-dimensional function $f(x)$, separated by uniform spacing $\Delta x$, the continuous Fourier integral can be approximated using the DFF; that is,

$$F(k_j) = \sum_{m=0}^{N-1} f(x_m) e^{-2\pi i jm/N}, \tag{71}$$

where $x_m = m\Delta x$, $\Delta k = 1/(N\Delta x)$, and $k_j = j\Delta k$. The corresponding inverse DFT is

$$f(x_m) = \frac{1}{N} \sum_{j=0}^{N-1} F(k_j) e^{2\pi i jm/N}. \tag{72}$$

For any sampling interval $\Delta x$, the maximum frequency in its Fourier transform is $k_{max} = 1/(2\Delta x)$, which is the Nyquist critical frequency [41]. In other words, the valid frequency components are contained in $F(k_j)$ with $j \in [0, N/2]$, whereas $F(k_j) = F^*(k_{N-j})$ for $j \in [N/2 + 1, N - 1]$, where $*$ indicates complex conjugation. As a consequence, the inverse Fourier transform given by Eq. (69) is discretized to a good approximation as

$$f^{(n)}(x_m) = \frac{1}{N} \sum_{j=0}^{N-1} (2\pi i \bar{j}\Delta k)^n F(k_j) e^{2\pi i jm/N}, \tag{73}$$

where $\bar{j} = j$ for $j \in [0, N/2]$ and $\bar{j} = j - N$ for $j \in [N/2 + 1, N - 1]$. Similarly, for a two-dimensional data set of $N_1 \times N_2$ points, the discrete form of Eq. (69) is

$$f^{(n_1, n_2)}(x_{m_1}, x_{m_2}) = \sum_{j_1=0}^{N_1-1} \sum_{j_2=0}^{N_2-1} (2\pi i \bar{j}_1 \Delta k_1)^{n_1} (2\pi i \bar{j}_2 \Delta k_2)^{n_2} F(k_{j_1}, k_{j_2}) e^{2\pi i (m_1 j_1 + m_2 j_2)/(N_1 N_2)}. \tag{74}$$

The DFT scheme provides a global representation of the derivative operator and is highly accurate for functions with periodic boundary conditions.

Table 2 lists the maximum error in evaluating the second of $e^{-x^2}$, using finite difference method and the Fourier Transform method for various grid spacing.

# 4. COMPUTATIONAL SCHEMES FOR MULTIELECTRON QUANTUM DOTS

A number of computational schemes of varying sophistication can be used to obtain solutions of the multielectron Schrödinger equation

$$\hat{H}\psi(\mathbf{r}_1 \cdots \mathbf{r}_N, \sigma_1 \cdots \sigma_N) = E\psi(\mathbf{r}_1 \cdots \mathbf{r}_N, \sigma_1 \cdots \sigma_N) \tag{75}$$

where $\hat{H} = \hat{H}^{dot} + \hat{H}^{ee}$ is the system Hamiltonian defined by Eq. (3). The solutions to the above equation would provide a detailed theoretical description of multielectron quantum dots, and in particular their energy structures.

If the electrons were assumed to not interact with each other, the above equation could be reduced to a single-electron Schrödinger equation, which can be solved by the methods described in the previous section. These electrons would then sequentially fill the

single-electron energy levels, starting from the lowest state according to the Pauli's exclusion principle. The total energy of a multielectron quantum dot would be a simple sum of the energies of individual electrons in the dot.

However, the Coulomb interactions between the electrons are significant, especially when they are comparable with the confinement potential imposed by external electrods, and therefore cannot be ignored. Several computational schemes have been developed to deal with the interacting electrons in quantum mechanically confined systems, such as atoms and molecules. These methods have been extended to study the electronic structure of quantum dots and other nanosystems. Their strengths and limitations are reviewed in this section.

## 4.1. Constant Interaction Model

The simplest theoretical description for a multielectron quantum dot is the constant interaction model, in which the Coulomb interactions between an electron in the dot and all other electrons inside and outside of the dot are parameterized by a set of constant capacitances [1, 3, 42–44]. Under this approximation, the ground state energy of an $N$-electron quantum dot can be written as

$$U(N) = [e(N - N_0) - C_g V_g]^2 \frac{1}{2C} + \sum_{n=1}^{N} E_n(B) \qquad (76)$$

where $V_g$ is the gate voltage, $N_0$ is the number of electrons in the quantum dot when $V_g = 0$, $C_g$ is the capacitance between the dot and the gate, $C$ is the total capacitance between the dot and all electrodes (i.e., gate, source, and drain), and $E_n(B)$ is the energy of the $n$th single-electron state in an external magnetic field $B$.

The quantity measured by the experiments of Tarucha et al. [45] is the addition energy, which is the energy required to add an extra electron to the system. The chemical potential of the system is defined as

$$\mu(N) = U(N) - U(N - 1) \qquad (77)$$

and so we have

$$\mu(N) = \left(N - N_0 - \frac{1}{2}\right) \frac{e^2}{C} - eV_g \frac{C_g}{C} + E_N \qquad (78)$$

where $E_N$ is the highest filled single-electron state of the $N$-electron dot. The addition energy is given by

$$\Delta\mu(N) = \mu(N + 1) - \mu(N) = E_{N+1} - E_N + \frac{e^2}{C} = \Delta E + \frac{e^2}{C} \qquad (79)$$

Here the first term is the energy difference between the highest occupied and the lowest unoccupied single electron state, and the second term represents collectively the energy of electrostatic repulsion between the electrons. This assumes that the electrons in the dot will spread evenly over the dot area, regardless of the number of electrons in the dot, and the capacitance $C$ can be related to the physical dimensions of the quantum dot, using the parallel plate capacitor formula.

Despite the simplicity of the constant interaction model, it is able to describe the various properties of a quantum dot system. For instance, it provides a surprisingly good description of the addition energy for a quantum dot, explaining both the unusual peaks in the addition energy spectrum and the variations in the position of current peaks as a function of gate voltage in a varying magnetic field. It is also so simple that it can deal with a large number of electrons easily.

As an example we take the work of Tarucha et al. [17]. Figure 2a shows the measured current through a quantum dot as a function of gate voltage. The addition energies were obtained from the differences between the consecutive current peaks. Clearly seen are maxima values of the addition energy for 2, 6, and 12 electrons. Kouwenhoven et al. [1]

use the constant interaction model to explain the existence of large jumps in the addition energy at these "magic numbers." For multielectron quantum dots, there are degenerate or nearly degenerate single-electron energy levels. In particular, for a two-dimensional parabolic potential, the first energy level has two degenerate states, the second energy level has four, and the third level has six. In this case, the constant interaction model predicts that for 2, 6, and 12 electrons, $\Delta E$ will be nonzero, and there is thus an abnormally large addition energy.

The constant interaction model also provides an explanation for the variations in the measured chemical potential as a function of magnetic field. Continuing with the assumption of a parabolic potential, the single-electron solutions are the Fock-Darwin states, as plotted in Fig. 3. The crossing of the energy levels of the Fock-Darwin states with different angular momentum means that the ground state for a given number of electrons changes the value of its total angular momentum as the magnetic field is varied. Figure 11 shows the variation in the experimental chemical potential as the applied magnetic field increases, along with the chemical potential predicted by the constant interaction model.

Although the constant interaction model is capable of explaining some of the qualitative features of quantum dots, it is not very useful as a tool for predicting parameters from *ab initio* calculations. The model assumes that the value of $\Delta E$ derives from the energy level spacing for non-interacting electrons and is therefore most valid in situations in which the wavefunction and energy is not changed significantly by interaction with other electrons. This is the situation in which the electron is strongly confined by an external potential.

## 4.2. Hartree-Fock Method

In 1928, Hartree [46] established a mean-field model for quantum systems, in which each electron is assumed to experience an averaged repulsive potential caused by all the other electrons in the system. The multielectron Schrödinger equation is thus reduced to a single-electron Hartree equation

$$\left[ \hat{H}_i + \sum_{j=1}^{N} \int \psi_j^*(r_j) \frac{1}{r_{ij}} \psi_j(r_j)\, dr_j \right] \psi_i(r_i) = E_i \psi_i(r_i) \tag{80}$$

where $\hat{H}_i$ is the single-electron Hamiltonian acting only on the $i$th electron, as defined in Eq. (3); $\psi_i(r_i)$ is the single-electron wavefunction for the $i$th electron; and $E_i$ is the corresponding eigenenergy.

A self-consistent procedure was proposed to solve the Hartree equation. In this approach, an initial guess is made for the wavefunctions of all electrons except the $i$th electron. The above Hartree equation is then solved for the wavefunction of the $i$th electron. The solution is then used as an improved trial wavefunction in the calculation for other electrons.



Figure 11. Chemical potential versus applied magnetic field: (a) experimental data; (b) calculated results using the constant interaction model. Reprinted with permission from [17]. S. Tarucha et al., *Phys. Rev. Lett.* 77, 3613 (1996). © 1996, American Physical Society.

Self-consistency is obtained by repeating this process iteratively for one electron at a time until no further change to the wavefunctions is noticeable.

However, the Hartree method did not consider the antisymmetry requirement on the wavefunctions imposed by Pauli's exclusion principal. In 1930, Fock and Slater [47] made an extension known as the Hartree-Fock method, in which the total wavefunction of a multielectron system is described by a Slater determinant that is antisymmetric by definition. A derivation of the Hartree-Fock equations is given below that follows some of the steps given by Bransden and Jochain [48].

Assuming $E_0$ is the ground-state energy of the system and $\Phi$ is a trial wavefunction, the variational principle tells us that

$$E_0 \leq E[\Phi] \equiv \langle \Phi | \hat{H} | \Phi \rangle = \langle \Phi | \hat{H}^{\text{dot}} + \hat{H}^{\text{int}} | \Phi \rangle \tag{81}$$

In the Hartree-Fock method, the trial wavefunction is defined as a Slater determinant; that is,

$$\Phi_b(\mathbf{q}_1, \mathbf{q}_2, \ldots, \mathbf{q}_N) = \frac{1}{\sqrt{N!}} \begin{vmatrix} \psi_\alpha(\mathbf{q}_1) & \psi_\beta(\mathbf{q}_1) & \cdots & \psi_\nu(\mathbf{q}_1) \\ \psi_\alpha(\mathbf{q}_2) & \psi_\beta(\mathbf{q}_2) & \cdots & \psi_\nu(\mathbf{q}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \psi_\alpha(\mathbf{q}_N) & \psi_\beta(\mathbf{q}_N) & \cdots & \psi_\nu(\mathbf{q}_N) \end{vmatrix} \tag{82}$$

where $\mathbf{q}_i \equiv (\mathbf{r}_i, \sigma_i)$ represents both the spatial and the spin coordinate of the $i$th electron, $\psi_\lambda(\mathbf{q}_i) = u_\lambda(\mathbf{r}_i)\chi_\lambda$ is the spin-orbital of the $i$th electron with a collective quantum number $\lambda$, and $u$ and $\chi$ are respectively the spatial and spin wavefunction. This definition, in conjuction with the requirement of orthonormality, that is,

$$\langle \psi_\mu | \psi_\lambda \rangle = \delta_{\mu, \nu} \tag{83}$$

ensures that the total wavefunction is antisymmetric. The Slater determinant can also be written as

$$\Phi_b(\mathbf{q}_1, \mathbf{q}_2, \ldots, \mathbf{q}_N) = \frac{1}{\sqrt{N!}} \sum_P (-1)^P P \psi_\alpha(\mathbf{q}_1)\psi_\beta(\mathbf{q}_2) \cdots \psi_\nu(\mathbf{q}_N) = \sqrt{N!} \mathcal{A}\Phi \tag{84}$$

where $P$ is the permutation operator, $\mathcal{A} = \frac{1}{N!}\sum_P(-1)^P P$ is the anitsymmetrization operator, and $\Phi$ is a simple product of the individual spin-orbitals.

The antisymmetrization operator possesses two important properties; first, $\mathcal{A}^* = \mathcal{A}$ and $\mathcal{A}^2 = \mathcal{A}$, and second, it commutes with both parts of the Hamiltonian $[\hat{H}^{\text{dot}}, \mathcal{A}] = [\hat{H}^{\text{int}}, \mathcal{A}] = 0$. As a result, we have

$$\langle \Phi_b | \hat{H}^{\text{dot}} | \Phi_b \rangle = N! \langle \Phi | \mathcal{A} \hat{H}^{\text{dot}} \mathcal{A} | \Phi \rangle$$

$$= N! \langle \Phi | \hat{H}^{\text{dot}} \mathcal{A} | \Phi \rangle$$

$$= \sum_{i=1}^N \sum_P (-1)^P \langle \Phi | \hat{H}_i P | \Phi \rangle$$

$$= \sum_{i=1}^N \langle \Phi | \hat{H}_i | \Phi \rangle$$

$$= \sum_\lambda \langle \psi_\lambda(\mathbf{q}_i) | \hat{H}_i | \psi_\lambda(\mathbf{q}_i) \rangle \tag{85}$$

where $\lambda = 1 \ldots N$. In the second to last line, the sum over permutations is dropped because the orthogonality condition means that only the identity permutation will have a nonzero

value. In the last line, the sum over $i$ is replaced with a sum over all occupied orbitals $\lambda$, noting that the terms will evaluate identically for any given $i$. We also have

$$
\begin{aligned}
\langle \Phi_D | \hat{H}^{int} | \Phi_D \rangle &= N! \langle \Phi | \cdot \hat{H}^{int} \cdot | \Phi \rangle \\
&= N! \langle \Phi | \hat{H}^{int} \cdot | \Phi \rangle \\
&= \sum_{i,j} \sum_P (-1)^P \left\langle \Phi \left| \frac{1}{r_{ij}} P \right| \Phi \right\rangle \\
&= \sum_{i,j} \left( \left\langle \Phi \left| \frac{1}{r_{ij}} \right| \Phi \right\rangle - \left\langle \Phi \left| \frac{1}{r_{ij}} P_{i,j} \right| \Phi \right\rangle \right)
\end{aligned}
\tag{86}
$$

where the singular permutation operator $P_{ij}$ interchanges the coordinates of electrons $i$ and $j$. Taking into account the orthogonality of the individual spin-orbitals, the above formula can be simplified further as

$$
\langle \Phi_D | \hat{H}^{int} | \Phi_D \rangle = \frac{1}{2} \sum_\lambda \sum_\mu \left\langle \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \left| \frac{1}{r_{ij}} \right| \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \right\rangle - \left\langle \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \left| \frac{1}{r_{ij}} \right| \psi_\mu(\mathbf{q}_i) \psi_\lambda(\mathbf{q}_j) \right\rangle
\tag{87}
$$

Therefore, we have

$$
\begin{aligned}
E[\Phi_D] = \sum_\lambda \langle \psi_\lambda(\mathbf{q}_i) | \hat{H}_i | \psi_\lambda(\mathbf{q}_i) \rangle &+ \frac{1}{2} \sum_\lambda \sum_\mu \left\langle \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \left| \frac{1}{r_{ij}} \right| \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \right\rangle \\
&- \left\langle \psi_\lambda(\mathbf{q}_i) \psi_\mu(\mathbf{q}_j) \left| \frac{1}{r_{ij}} \right| \psi_\mu(\mathbf{q}_i) \psi_\lambda(\mathbf{q}_j) \right\rangle
\end{aligned}
\tag{88}
$$

The variational equation reads

$$
\delta E - \sum_{\lambda\mu} \epsilon_{\lambda\mu} \delta \langle \psi_\lambda | \psi_\mu \rangle = 0
\tag{89}
$$

where the Lagrange multipliers have the property $\epsilon_{\lambda\mu} = \epsilon_{\mu\lambda}^*$. Making a unitary transformation

$$
\psi_\lambda' = \sum_\mu U_{\mu\lambda} \psi_\mu
\tag{90}
$$

does not affect the functional $E[\Phi_D]$ at all and only affects $\Phi_D$ by a phase factor. However, this allows us to arrange for the matrix of elements $\epsilon_{\lambda\mu}$ to be diagonal; that is,

$$
\delta E - \sum_\lambda E_\lambda \delta \langle \psi_\lambda | \psi_\lambda \rangle = 0
\tag{91}
$$

Proceeding with the variation leads to the system of equations

$$
\hat{H}_i \psi_\lambda(\mathbf{q}_i) + \sum_{\mu=1} \int \psi_\mu^*(\mathbf{q}_j) \frac{1}{r_{ij}} \psi_\mu(\mathbf{q}_j) d\mathbf{q}_j \psi_\lambda(\mathbf{q}_i) - \sum_{\mu=1} \int \psi_\mu^*(\mathbf{q}_j) \frac{1}{r_{ij}} \psi_\lambda(\mathbf{q}_j) d\mathbf{q}_j \psi_\mu(\mathbf{q}_i) = E_\lambda \psi_\lambda(\mathbf{q}_i)
\tag{92}
$$

which are the Hartree-Fock equations.

To solve the Hartree-Fock equations, an iterative procedure was again adopted. At the $n$th iteration, we have an estimate for each spin-orbital denoted by $\psi_\lambda^n$. We can rewrite the Hartree-Fock equations as

$$
\hat{H}_i \psi_\lambda^{n-1}(\mathbf{q}_i) + \sum_\mu \int \psi_\mu^{n}(\mathbf{q}_j) \frac{1}{r_{ij}} \psi_\mu^{n}(\mathbf{q}_j) d\mathbf{q}_j \psi_\lambda^{n-1}(\mathbf{q}_i) - \sum_\mu \int \psi_\mu^{n}(\mathbf{q}_j) \frac{1}{r_{ij}} \psi_\lambda^{n-1}(\mathbf{q}_j) d\mathbf{q}_j \psi_\mu^{n}(\mathbf{q}_i)
$$
$$
= E_\lambda^{n-1} \psi_\lambda^{n-1}(\mathbf{q}_i)
\tag{93}
$$

Separating the spin part from the spatial part of the wavefunctions; that is, $\psi_\mu(\mathbf{q}_i) = u_\mu(\mathbf{r}_i)\chi_\mu$, we obtain

$$\hat{H}_i u_\lambda^{n-1}(\mathbf{r}_i) + \sum_{\mu=1} \int u_\mu^n(\mathbf{r}_j) \frac{1}{r_{ij}} u_\mu^n(\mathbf{r}_j)\langle\chi_\mu|\chi_\mu\rangle \, d\mathbf{q}_j u_\lambda^{n-1}(\mathbf{r}_i)$$

$$-\sum_{\mu=1} \int u_\mu^n(\mathbf{r}_j) \frac{1}{r_{ij}} u_\lambda^{n-1}(\mathbf{r}_j)\langle\chi_\mu|\chi_\lambda\rangle \, d\mathbf{q}_j u_\mu^n(\mathbf{r}_i) = E_\lambda^{n+1} u_\lambda^{n-1}(\mathbf{r}_i) \qquad (94)$$

Self-consistency is obtained by repeating this process iteratively until the difference between $\psi_\lambda^{n+1}$ and $\psi_\lambda^n$ is negligibly small.

One of the first theoretical studies of quantum dots came from Kumar et al. [49], who used the Hartree method and ignored exchange and correlation effects completely. In their work, the Hartree equation [Eq. (80)] was solved self-consistently in conjunction with the Poisson's equation

$$\epsilon^* \nabla^2 \varphi + \rho \qquad (95)$$

which provides the electrostatic potential $-e\varphi$. The boundary conditions were determined by voltages applied to the gates and contacts. The electron density $\rho$ in the quantum dot was related to the single-electron wavefunctions at each iteration. They obtained the energies for the lowest states of a quantum dot with seven electrons as a function of the applied magnetic field, as shown in Fig. 12. Their results provides a qualitative description of the main features of a multielectron quantum dot.

Soon after, Pfannkuche et al. [50] performed Hartree-Fock self-consistent calculations for a two-electron quantum dot (i.e., the artificial helium) and compared them with results obtained from the Hartree method and a direct numerical diagonalization of the two-particle Hamiltonian. The researchers obtained excellent agreement between the Hartree-Fock and the exact numerical calculations for the lowest triplet state ($S = 1$) but obtained markedly different results for the ground singlet state ($S = 0$), as shown in Fig. 13. This indicates that spin correlation is not fully accounted for in the Hartree-Fock formalism for such a system. Plotted in the same figure is the researchers' Hartree calculation, which does not distinguish the singlet and triplet states and significantly overestimates the ground state energies.



Figure 12. The lowest energy states of a quantum dot with 7 electrons versus applied magnetic field using the Hartree method. Reprinted with permission from [49], A. Kumar et al., *Phys. Rev. B* 42, 5166 (1990). © 1990, American Physical Society.

Figure 13. The lowest energy states of a quantum dot with 2 electrons versus applied magnetic field using the Hartree method. Hartree-Fock method and exact numerical calculation. Reprinted with permission from [50], D. Pfannkuche et al., *Phys. Rev. B* 47, 2244 (1993). © 1993, American Physical Society.

Nevertheless, many authors have subsequently used the Hartree-Fock method to study the electronic structure of quantum dot systems, as it explicitly takes into account the exchange effect, and at least partially the electron–electron corrections [18, 51–57]. The most commonly adapted approach is the so-called unrestricted Hartree-Fock method, in which the two electrons in the same shell, but with opposite spins, are not restricted to having the same spatial wavefunction. This generally provides more accurate energies in comparison with the restricted Hartree-Fock approach. An implementation of the unrestricted Hartree-Fock method in *Mathematica* was developed by McCarthy et al. [58].

Oaknin et al. [59] and Palacios et al. [51] carried out both Hartree and Hartree-Fock calculations to study the properties of a multielectron quantum dot under the influence of an external magnetic field. They also carried out a comparison between the unrestricted Hartree-Fock calculations and exact numerical diagonalization solutions for up to five electrons, as shown in Fig. 14, and found that the Hartree-Fock ground state energies are overestimated in general, especially for zero or low external magnetic fields. Also shown in this figure is the transition of the ground state from the paramagnetic state to the ferromagnetic state when the applied magnetic field is around $4T$ to $6T$.

Fujito et al. [52] applied the same approach to studying anisotropic parabolic quantum dots. The researchers obtained the energy and wavefunction of the ground states for large and small quantum dots containing up to 12 electrons. They found that the large dots had a ground state that is completely spin polarized, whereas small dots were not. They also studied the effect of the vertical extent of the dot, using a three-dimensional model potential

$$V(\mathbf{r}) = \frac{1}{2}\omega_{\perp}^2(x^2 + y^2) + \frac{1}{2}\omega_z^2 z^2 \tag{9b}$$

Figure 15 shows the capacitance versus the number of electrons in a quantum dot with the same $\omega_{\perp}$ but different vertical extend $\omega_z$. Bielińska-Wąż et al. [60] studied two electron systems confined in a three-dimensional well. which is described by a combination of a coulombic $(1/r)$ and a harmonic $(\omega^2 r^2)$ term, to see how changes to the potential affect the spectra of excited states.

Several researchers [53, 55, 56] used the Hartree-Fock method to examine Wigner crystallization in single and double quantum dots, which refers to the transition from a system

**Figure 14.** Chemical potential $\mu(N)$ calculated using the Hartree-Fock method (dashed lines) and the configuration interaction method (solid lines). Reprinted with permission from [51], J. J. Palacios et al., *Phys. Rev. B* 50, 5760 (1994). © 1994, American Physical Society.

of tightly bound electrons to a more loosely bound state, in which electrons behave in a more classical manner, forming isolated puddles in the electron charge distribution [61]. Müller et al. [53] demonstrated that the quantum dots undergo a gradual transition to a spin-polarized Wigner crystal with increasing magnetic field strength. Yannouleas et al. [55]



**Figure 15.** capacitance versus the number of electrons for various dot sizes with $l_s = 10.7$ nm: (a) $l_z = 0$ nm; (b) $l_z = 0.53$ nm; (c) $l_z = 2.13$ nm; and (d) $l_z = 7.11$ nm. Reprinted with permission from [52], M. Fujito et al., *Phys. Rev. B* 53, 9952 (1996). © 1996, American Physical Society.

discussed three types of Wigner crystallization, including pure spin density waves [62], spatial localization of charge density in a single quantum dot, and the separation of charge distribution onto each individual dot of a multidot system, even at zero external magnetic. The researchers predicted that Wigner crystallization occurs in quantum dots when the interelectron Coulomb repulsion and the harmonic confinement potential become comparable, as shown in Fig. 16. Reusch et al. [56] also observed Wigner crystallization in their calculations at fairly high electron density, as shown in Fig. 17. However, this appears to be a premature claim and is likely an artifact of broken spatial symmetry in the Hartree-Fork formalism, as more exact approaches seem to indicate [56, 63].

Bednarek's group carried out a series of restricted and unrestricted Hartree-Fock calculations (RHF and UHF) [54, 57, 64–66] to study spherical and cylindrical quantum dots made of Si and GaAs containing up to 20 electrons. They evaluated the critical values for the confinement potential parameters, which allow the quantum dot to contain 1–20 electrons. Poisson's equation was used to determine the shape of the lateral confining potential in a quantum dot, as shown in Fig. 18, based on measurements of the dimensions of an actual device. The Schrödinger equation was solved self-consistently with the Poisson's equation to obtain the wavefunctions and the energy levels of the system. In the particular case of two electrons confined in single isolated or double coupled quantum dots, they proposed an effective electron–electron potential to replace the Coulomb interaction, enabling the solution of this model problem exactly. Shown in Fig. 19 is a direct comparison between the exact solutions and RHF and UHF. It was found that the prediction for the triplet state energies from both RHF and UHF is accurate for all sizes of the quantum dot. For the singlet state energies, the UHF results are reasonably accurate, especially for small or extremely large dot size, whereas the RHF provides pretty poor predictions.

Wang and Hines [18] performed calculations for a circularly symmetric quantum dot and a triangular quantum dot with up to 13 electrons. The confinement potential is defined as $V(r) = m^*k^2(x^2 + y^2)/2$ and $V(r) = m^*k^2(x^2 + y^2)(1 + 2\cos(3\theta)/7)/2$, respectively. The energy levels of the two dots with a single electron are plotted in Fig. 20 to assist the assignment of ground state configurations. For the circular dot, the first shell is fully filled with two electrons, and the subsequent shells are fully filled with four, six, and eight electrons. For the triangular dot, the maximum number of electrons for each shell is different, being two, four, two, four, four, and two for the first to the sixth shells, respectively. The configurations listed in Tables 3 and 4 are found to provide the lowest energy and therefore correspond to the ground state of these systems. The shell filling order indeed obeys Hund's rule in both cases. As shown, each shell is filled with electrons of the same spin until half-full and is then filled with electrons of the opposite spin until full. Also listed in Tables 3 and 4 are their calculated results for the total energy of the quantum dots with $k = 0.2$ eV.

A good way to demonstrate the shell structure of quantum dots is to plot the addition energy, which is defined by Eq. (79). As the energy required to add an additional electron to the next shell is larger than to add an electron to the same shell, the peaks in the addition energy spectrum correspond to the number of electrons in a full-filled shell. Similarly, extra energy is required to add an electron of opposite spin because of exchange effect, and thus the secondary peaks in the addition energy spectrum correspond to the number of electrons in a half-filled shell.



Figure 16. Electron density distribution of a six electron quantum dot with   $\omega = 5$ meV and $B = 0$: (a) $R = 0.95$, (b) $R = 1.48$, (c) $R = 3.18$, where $R$ is the ratio between the interelectron Coulomb repulsion and the harmonic confinement potential. Reprinted with permission from [55], C. Yannouleas and U. Landman, *Phys. Rev. Lett.* 82, 5325 (1999). © 1999, American Physical Society.

**Figure 17.** Electron density distribution of a seven-electron quantum dot, forming a Wigner molecule as the Coulomb interaction strength between the electrons is increased. Reprinted with permission from [56], B. Reusch et al., *Phys. Rev. B* 63, 113313 (2001). © 2001, American Physical Society.



**Figure 18.** Potential energy calculated from the Poisson equation as a function of cylindrical coordinates. After reprinted with permission from [65], S. Bednarek et al., *Phys. Rev. B* 64, 195303 (2001). © 2001, American Physical Society.



**Figure 19.** Estimated errors of the UHF (solid curve) and RHF (dashed curve) methods for the singlet state and both UHF and RHF (dotted curve) for the triplet state of a two-electron quantum dot. Reprinted with permission from [65], S. Bednarek et al., *Phys. Rev. B* 64, 195303 (2001). © 2001, American Physical Society.

**Figure 20.** Energy levels (in a.u.) for a single electron in the circularly symmetric dot (solid lines) and the triangular dot (dashed lines). The labels on the left (right) side denote the shells for the circularly symmetric (triangular) dot. The numbers in the bracket represent the maximum number of electrons allowed in the corresponding shell.

Shown in Fig. 21 are their calculated addition energies for the circular and triangular dots. Clearly seen are the expected full-filling peaks at $N = 2, 6$, and 12 electrons for the circular dot and $N = 2, 6, 8$, and 12 electrons for the triangular dot. The researchers also obtain the secondary half-filling peaks at $N = 4$ and 9 for the circular dot and $N = 4$ for the triangular dot. In their previous work [58], a secondary peak at $N = 8$ was reported for the circular dot. This was found to be an artifact of the limited matrix size used in their earlier calculation. Note that, for the triangular dot, the energies for third and fourth shells are very close. As a result, the $N = 8$ full-filling peak is not as pronounced as the others, and the $N = 7$ and $N = 10$ half-filling peaks are not obvious from their calculations.

The spatial electron charge densities $\rho(x, y) = \sum_i |\psi_i(x, y)|^2$ for the circular dot and the triangular dot are plotted in Figs. 22 and 23. The total number of electrons in the dots ranges from $N = 2$ to $N = 10$, shown from left to right and from top to bottom. The number of spin-up and spin-down electrons for each dot is determined according to Tables 3 and 4, which correspond with the ground state configuration. Indeed, the densities reflect the symmetry of the quantum dots.

## 4.3. Density Functional Theory

Density-functional theory, like the Hartree-Fock method, is also a self-consistent mean-field model. However, the electron density distribution $n(r)$, rather than the multielectron wavefunction (namely, the Slater determinant), is used in its formulation. The main advantage of this method is that it can deal with a huge number of electrons in the system simultaneously, whereas the wavefunction approaches (such as the Hartree-Fock method and the configuration interaction method to be discussed in Section 4.4) are limited to quantum systems with only a small number of active electrons.

**Table 3.** Shell filling order for the circular dot ($k = 0.2$ eV). The full shells are labeled by '' and half-shells by '.

| $N$ | 1st Shell | 2nd Shell | 3rd Shell | E(ev) |
|---|---|---|---|---|
| 1 | ↑ | | | 0.200 |
| 2'' | ↑↓ | | | 0.452 |
| 3 | ↑↓ | ↑ | | 0.916 |
| 4' | ↑↓ | ↑ ↑ | | 1.409 |
| 5 | ↑↓ | ↑↓ ↑ | | 1.952 |
| 6'' | ↑↓ | ↑↓ ↑↓ | | 2.517 |
| 7 | ↑↓ | ↑↓ ↑↓ | ↑ | 3.287 |
| 8 | ↑↓ | ↑↓ ↑↓ | ↑ ↑ | 4.078 |
| 9' | ↑↓ | ↑↓ ↑↓ | ↑ ↑ ↑ | 4.889 |
| 10 | ↑↓ | ↑↓ ↑↓ | ↑↓ ↑ ↑ | 5.748 |
| 11 | ↑↓ | ↑↓ ↑↓ | ↑↓ ↑↓ ↑ | 6.633 |
| 12'' | ↑↓ | ↑↓ ↑↓ | ↑↓ ↑↓ ↑↓ | 7.522 |

Table 4. Shell filling order for the triangular dot ($k = 0.2$ eV). The full shells are labeled by **

| N | 1st Shell | 2nd Shell | 3rd Shell | 4th Shell | E(eV) |
|---|---|---|---|---|---|
| 1 | ↑ | | | | 0.199 |
| 2 ** | ↑↓ | | | | 0.449 |
| 3 | ↑↓ | ↑ | | | 0.905 |
| 4 | ↑↓ | ↑ | ↑ | | 1.386 |
| 5 | ↑↓ | ↑↓ | ↑ | | 1.918 |
| 6 ** | ↑↓ | ↑↓ | ↑↓ | | 2.472 |
| 7 | ↑↓ | ↑↓ | ↑↓ | ↑ | 3.208 |
| 8 ** | ↑↓ | ↑↓ | ↑↓ | ↑↓ | 3.973 |
| 9 | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑ | 4.780 |
| 10 | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑ ↑ | 5.617 |
| 11 | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ ↑ | 6.483 |
| 12 ** | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ ↑↓ | 7.381 |

At the heart of the density-functional theory is the self-consistent single-electron Kohn-Sham equation [67]:

$$\frac{-\hbar^2}{2m^*}\nabla^2\psi_i(\mathbf{r}) + [V_{ext} + V_e(\mathbf{r}) + V_{xc}(\mathbf{r})]\psi_i(\mathbf{r}) = \varepsilon_i\psi_i(\mathbf{r}) \qquad (97)$$

developed from the Hohenberg-Kohn theorems [68]. The term $V_{ext}$ represents the external electric potential imposed by, for example, external electrodes; $\psi_i$ is the wavefunction for the $i$th electron, which is solved from the Kohn-Sham equation to provide the electron density distribution $n(\mathbf{r})$, defined as

$$n(\mathbf{r}) = \sum_{i=1}^{N} |\psi_i(\mathbf{r})|^2 \qquad (98)$$



Figure 21. Addition energies for a circular dot (top) and a triangular dot (bottom) (♦ $k = 0.2$ eV; ★ $k = 0.1$ eV; ■ $k = 0.05$ eV).

**Figure 22.** Charge density distribution for a circular dot ($k = 0.2$ eV) with $N = 2$ to $N = 10$ electrons (from left to right and top to bottom).

The Coulomb potential is then given by

$$V_c(\mathbf{r}) = \frac{e^2}{4\pi\varepsilon} \int \frac{n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}'$$

whereas the exchange-correlation potential $V_{xc}(\mathbf{r})$ depends functionally on the electron density distribution $n(\mathbf{r})$. If the exact exchange-correlation functional $E_{xc}[n(\mathbf{r})]$ is used, the Kohn-Sham equation incorporates all many-particle effects. However, exchange effects come directly from the antisymmetrization of wavefunctions as required by the Pauli's exclusion principle. In the density-functional theory, this is a major problem because the mathematical object is the electron density distribution function, rather than the electron wavefunction, making evaluation of the exchange interaction intrinsically difficult. For many quantum systems, this functional cannot be exactly defined, and recent work has involved a considerable amount of empirical parameterization [69, 70].



**Figure 23.** Charge density distribution for a triangular dot ($k = 0.2$ eV) with $N = 2$ to $N = 10$ electrons (from left to right and top to bottom).

The simplest and the most widely used representation for $E_{xc}[n(\mathbf{r})]$ is the so-called local-density approximation (LDA): that is.

$$E_{xc}^{LDA} \equiv \int e_{xc}(\zeta)n(\mathbf{r})\,d\mathbf{r}.$$                  (99)

where $\zeta$ represents the spin polarization and $e_{xc}$ is the exchange-correlation energy. LDA is valid for homogeneous two-dimensional electrons and also for systems with small variation in electron density. The exchange-correlation energy can be parameterized as

$$e_{xc}(\zeta) = \frac{a_0(\zeta)(1 + a_1(\zeta)\sqrt{x})}{1 + a_1(\zeta)\sqrt{x} + a_2(\zeta)x + a_3(\zeta)\sqrt{x}}$$      (100)

where $x$ relates to the electron density and is defined as the radius of a sphere containing one electron. The coefficients $a_i(0)$ and $a_i(1)$ were determined by Tanatar and Ceperley [71] for the ground state of two-dimensional electron gas using the Green's function Monte Carlo method. For other values of $\zeta$, one can use [62, 72, 73]

$$e_{xc}(\zeta) = e_{xc}(0) + \frac{(1 + \zeta)^{3/2} + (1 - \zeta)^{3/2} - 2}{2^{3/2} - 2}[e_{xc}(1) - e_{xc}(0)].$$   (101)

The Kohn-Sham equations are solved iteratively. This is similar to the Hartree method described in Section 4.2. The wavefunction of each electron is solved, taking into account a potential field determined by the average position of all other electrons. After a solution is obtained, the potential field is recalculated, and the Kohn-Sham equation is solved for a new solution. The calculation is thus iterated until both the potential field and the solution cease to change.

Macucci et al. [74–76] used the LDA formalism and the parameterized exchange-correlation energy given by Eq. (100), ignoring the spin polarization effect. Their calculated results for the addition energies of up to 25 electrons are shown in Fig. 24a. Lee et al. [77] improved the functional approximation for $E_{xc}[n(\mathbf{r})]$ by using a generalized



Figure 24. Addition energies calculated by Macucci et al. Reprinted with permission from [76], M. Macucci et al., *Phys. Rev. B* 55, 4879 (1997). © 1997, American Physical Society. (top) with confinement potential of $\hbar\omega =$ 4, 3, 2.5 meV and Lee et al. Reprinted with permission from [77], I. H. Lee et al., *Phys. Rev. B* 57, 9035 (1998). © 1998, American Physical Society. (bottom) with $\hbar\omega =$ 20, 10, 4 meV.

gradient approximation developed by Perdew et al. [78, 79] to correct for some of the local effects. Their results are shown in Fig. 24b.

Koskinen et al. [62] performed density-functional calculations for quantum dots containing up to 46 electrons. Their calculations included spin explicitly (therefore called spin-density-functional theory) and established a rich variety of magnetic structures in the ground state even in the absence of an external magnetic field. However, some of their observed features, such as the so-called spin-density wave, were later found to be artifacts of broken spin symmetry in the density-functional formalism. Hirose and Wingreen [80] extended their work to 58 electrons and also studied elliptical external potentials. They found that Hund's rule determines the ground state spin configuration only for circular parabolic quantum dots with a small number of electrons (around 22 for the model potential used in their calculations). For elliptic confinement potentials, the quantum dot exhibits a more Pauli-like behavior.

Wensauer et al. [81] extended the spin-density functional theory further to study two laterally coupled quantum dots (often named quantum-dot molecules) with two, four, and eight electrons. For two electrons, the spin-density functional theory predicted a transition from a spin-unpolarized ground state to a spin-polarized ground state. This was found again to be an artifact of the density-functional formalism. However, for higher numbers of electrons, the results are expected to be more reliable as the concept of electron density becomes more applicable. This is, to some extent, demonstrated by the calculations of Wensauer et al. for four and eight electrons.

## 4.4. Configuration Interaction Method

Both the Hartree-Fock method and density-functional theory are self-consistent mean-field models, but they are fundamentally different in the way of treating exchange and correlation effects. In the Hartree-Fock method, exchange is considered exactly by the proper antisymmetrization of wavefunctions, whereas correlation in the motion of the electrons is neglected because only a time-averaged effective potential is used in its formulation. In density-functional theory, both exchange and correlation effects are included exactly in principle but only approximately in practice in most situations, and perhaps worse still, there is no clear route for the density-functional theory to provide convergent results. To fully understand the exchange and correlation effects and to establish an accurate description of the quantum dot systems, benchmark calculations with arbitrarily high numerical precision would be required.

As discussed in the previous section, the Hartree-Fock method uses a single Slater determinant that is an antisymmetric function and optimizes the single-electron wavefunctions used in the construction of that determinant. A more accurate formalism is to use a linear combination of all Slater determinants that can be formed from a given set of noninteracting spin-orbitals, each of which describes a different configuration. The number of spin-orbitals can be increased to improve the accuracy of the calculated results. This method is termed as the configuration interaction (CI) method, which provides a convergent route to obtain numerically exact solutions of multielectron systems.

Specifically, we have

$$\Psi = \sum_{i=1}^{N_{det}} d_i \Phi_D^i \tag{102}$$

where $\Phi_D^i$ is the Slater determinant as defined in Eq. (82) for each different configuration. The number of Slater determinants to be included in the expansion is given by

$$N_{det} = \binom{\aleph}{N^\uparrow} \binom{\aleph}{N^\downarrow} \tag{103}$$

where $\aleph$ is the number of available spin-orbitals to be used in the expansion; $N^\uparrow$ and $N^\downarrow$ are the number of electrons with up and down spin, respectively; and $\aleph > N^\uparrow + N^\downarrow$, which is the total number of electrons in the system. The linear Rayleigh-Ritz variation principal can be used to determine the expansion coefficients $d_i$; namely, by solving the eigenvalue problem of Hamitonian matrix with elements $\langle \Phi_D^i | \hat{H} | \Phi_D^j \rangle$, as described in detail in Section 3.1.5.

The main difficulty of the CI method lies in the rapid explosion of the number of Slater determinants to be included in the expansion as the number of spin-orbitals increases. For example, to represent a six-electron quantum dot with 10 available spin-orbitals, we need 44,100 Slater determinants in the expansion, whereas with 20 spin-orbitals, the number of Slater determinants to be included is 1,502,337,600.

As a result, a truncated CI approach is almost inevitable in practice. First, the set of spin-orbital bases needs to be of reasonable size, and so one is required to select the most appropriate basis functions to approximate the true wavefunction. Second, assuming that the lowest few energy states of the systems under study are of major concern, one can determine the truncation in the CI expansion on the basis of total noninteracting energy $E_I^{non} = \langle \Phi_D^i | \hat{H}^{don} | \Phi_D^i \rangle$. Determinants $\Phi_D^i$ giving rise to lower values of $E_I^{non}$ are generally more important to the lowest-energy states than those with higher values of $E_I^{non}$, and they are thus included preferentially. A cutoff value can be set so that only determinants with $E^{non}$ below this value are included in the expansion. Convergence is achieved if the solutions and energies remain the same when increasing the number of spin-orbitals and the cutoff noninteraction energy in the CI expansion.

The most computationally intensive aspect of the procedure is the evaluation of a large number of Hamiltonian matrix elements

$$\langle \Phi_D | \hat{H} | \Phi_D' \rangle = N! \langle \Phi | \mathcal{A} \hat{H} \mathcal{A} | \Phi' \rangle$$
$$= N! \langle \Phi | \hat{H} \mathcal{A} | \Phi' \rangle \tag{104}$$

where $\Phi_D$ is a Slater determinant and $\Phi$ is a simple product of the individual spin-orbitals. Slater and Condon [100–102] established a set of rules (the so-called Slater-Condon rules), which allow us to reduce the $N$-electron integral (104) to a sum of one- or two-electron integrals and, furthermore, to identify zero Hamiltonian matrix elements.

For completeness, we include the derivation of these rules below. Note that the standard Slater-Condon rules are only applicable if the two Slater determinants $\Phi_D$ and $\Phi_D'$ are lined up in maximum coincidence. For example, if we have $\Phi = \psi_1 \psi_2 \psi_3 \psi_4 \psi_6$, $\Phi' = \psi_1 \psi_3 \psi_4 \psi_5 \psi_6$, $\Phi'$ would need be aligned up by pairwise permutation to $\psi_1 \psi_5 \psi_3 \psi_4 \psi_6$. We found it is easier to simply use the ordered lists $\Phi$ and $\Phi'$ as they are (i.e., not necessarily aligned) but to introduce a simple phase factor in these rules. The formulas derived below require $\psi$ be an orthogonal set; it is, however, possible to derive similar formulas for a nonorthogonal basis [82].

THEOREM 4.1. *In the case that $\hat{H} = h_0$, which is independent of electron coordinates,* $\langle \Phi_D | \hat{H} | \Phi_D' \rangle = h_0$ *if $\Phi = \Phi'$ and otherwise $\langle \Phi_D | \hat{H} | \Phi_D' \rangle = 0$.*

PROOF. Following Eq. (84), we have

$$\langle \Phi_D | \hat{H} | \Phi_D' \rangle = h_0 N! \langle \Phi | \mathcal{A} | \Phi' \rangle$$
$$= h_0 \sum_{\hat{P}} (-1)^P \langle \Phi | \hat{P} \Phi' \rangle$$
$$= h_0 \sum_{\hat{P}} (-1)^P \langle \psi_1 | \psi_{P1}' \rangle \langle \psi_2 | \psi_{P2}' \rangle \cdots \langle \psi_N | \psi_{PN}' \rangle \tag{105}$$

Because of the orthogonality of the $\psi_i$, the above expression is 0 unless $\psi_i = \psi_{Pi}'$ $\forall i$. This is only possible if all elements of $\Phi$ are the same as the elements of $\Phi'$. Because each element of $\Phi$ is unique, clearly only a single permutation satisfies this condition. Because the set $\Phi$ has its elements arranged in the same order as $\Phi'$, this is the identity permutation, and thus $(-1)^P = 1$. By setting $h_0 = 1$, we find that the Slater determinants formed from a set of orthonormal orbitals are themselves an orthonormal set.

THEOREM 4.2. *In the case that $\hat{H} = \sum_{i=1}^N \hat{h}_i$, where $h_i$ is the one-electron operator involving only coordinates of the ith electron,*

a. $\langle \Phi_D | \hat{H} | \Phi_D' \rangle = 0$, *if $\Phi$ and $\Phi'$ differ by more than one orbital;*

b. $\langle \Phi_D | \hat{H} | \Phi_D' \rangle = (-1)^{l-m} \langle \psi_l | \hat{h}_l | \psi_m' \rangle$, *if $\Phi$ and $\Phi'$ differ by one orbital $\psi_l$ versus $\psi_m'$, where l is the position of $\psi_l$ in $\Phi$ and m is the position of $\psi_m'$ in $\Phi'$;*

c. $\langle \Phi_D | \hat{H} | \Phi_D' \rangle = \sum_{i=1}^N \langle \psi_i | \hat{h}_i | \psi_i \rangle$, *if $\Phi = \Phi'$.*

**Figure 25.** Energy levels obtained from the CI method for a square nano box. Reprinted with permission from [83]. G. W. Bryant. *Phys. Rev. Lett.* 59, 1140 (1987). © 1987. American Physical Society.

PROOF. The contribution of $\hat{h}_i$ to $\langle \Phi_p | \hat{H} | \Phi'_p \rangle$ is given by

$$\langle \Phi_p | \hat{h}_i | \Phi'_p \rangle = \sum_{\tilde{p}} (-1)^{\tilde{p}} \langle \psi_1 | \psi'_{p_1} \rangle \langle \psi_2 | \psi'_{p_2} \rangle \cdots \langle \psi_i | \hat{h}_i | \psi'_{p_i} \rangle \cdots \langle \psi_N | \psi'_{p_N} \rangle$$

$$= \sum_{\tilde{p}} (-1)^{\tilde{p}} \langle \psi_i | \hat{h}_i | \psi'_{p_i} \rangle \prod_{l \neq i} \langle \psi_l | \psi'_{p_l} \rangle \tag{106}$$

which is 0 unless $\psi_j = \psi'_{p_j}$ $\forall j \neq i$ because of the orthogonality of the basis function $\psi_i$.

Let us assume that there are at least two orbitals $\psi_l$ and $\psi_m$, which appear in $\Phi$ but not $\Phi'$. If this is the case, then there does not exist any value $i$ for which $\psi_j = \psi'_{p_j}$ $\forall j \neq i$, and therefore $\langle \Phi_p | \hat{h}_i | \Phi'_p \rangle = 0$ $\forall i$. Thus, $\langle \Phi_p | \hat{H} | \Phi'_p \rangle = 0$ when $\Phi$ and $\Phi'$ differ by more than one orbital.

Now assume there is only one orbital $\psi_l$, which is in $\Phi$ but not in $\Phi'$. Then there exists only one value $i = l$ for which the conditions $\psi_j = \psi'_{p_j}$ $\forall j \neq i$ can be satisfied.



**Figure 26.** Energy levels obtained from the CI method including the influence of an external magnetic field. Reprinted with permission from [84]. P. A. Maksym and T. Chakraborty. *Phys. Rev. Lett.* 65, 108 (1990). © 1990. American Physical Society.

In this case,

$$\langle \Phi_p | \hat{H} | \Phi'_p \rangle = \sum_{i=1} \langle \Phi_p | \hat{h}_i | \Phi'_p \rangle$$

$$= \sum_{i=1} \sum_{P} (-1)^P \langle \psi_i | \hat{h}_i | \psi'_{p_i} \rangle \prod_{j=i} \langle \psi_j | \psi'_{p_j} \rangle$$

$$= (-1)^P \langle \psi_i | \hat{h}_i | \psi'_{p_i} \rangle \qquad (107)$$

where $(-1)^P = (-1)^{l-m}$ as we need $|l - m|$ pair permutations to align $\Phi$ and $\Phi'$.
    As an example consider the Slater Determiants formed from

$$\Phi = \psi_1 \psi_2 \psi_3 \psi_4 \psi_6$$
$$\Phi' = \psi_1 \psi_3 \psi_4 \psi_5 \psi_6 \qquad (108)$$

Clearly we need to find $\hat{P}$ such that

$$P\Phi' = \psi_1 \psi_5 \psi_3 \psi_4 \psi_6 \qquad (109)$$

This can be achieved through two pairwise permutations

$$\Phi' = \psi_1 \psi_3 \psi_4 \psi_5 \psi_6$$
$$P_{3,4}\Phi' = \psi_1 \psi_3 \psi_5 \psi_4 \psi_6 \qquad (110)$$
$$P_{2,3}P_{3,4}\Phi' = \psi_1 \psi_5 \psi_3 \psi_4 \psi_6$$

and therefore $(-1)^P = (-1)^2 = 1$. This is equivalent to $(-1)^{|l-m|} = (-1)^2 = 1$, where $l = 2$ and $m = 4$.
    Finally we consider the case $\Phi = \Phi'$, where the conditions $\psi_j = \psi'_{p_j}$ $\forall j \neq i$ can be satisfied for any choice of $i$ and we have

$$\langle \Phi_p | \hat{H} | \Phi_p \rangle = \sum_{i=1} \langle \psi_i | \hat{h}_i | \psi_i \rangle. \qquad (111)$$



Figure 27. Chemical potentials for two to eight electrons in three different dots: (a) $\hbar\omega = 109$ meV, (b) $\hbar\omega = 54.8$ meV and (c) $\hbar\omega = 10.96$ meV. Solid lines are from the configuration interaction calculations, alternate dash and dot lines are from the Hartree-Fock calculations, and dashed lines are from the constant interaction model

THEOREM 4.3.   *In the case that $H = \sum_{i<j}^{N} \hat{g}_{i,j}$, which is the two-electron operator involving only coordinates of the $i$th and $j$th electron,*

a. $\langle \Phi_D | H | \Phi'_D \rangle = 0$, *if $\Phi$ and $\Phi'$ differ by more than two orbitals;*

b. $\langle \Phi_D | H | \Phi'_D \rangle = (-1)^{|l-m|-|s-t|}(\langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_m \psi'_t \rangle - \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_t \psi'_m \rangle)$, *if $\Phi$ and $\Phi'$ differ by two orbitals, $\psi_l$ and $\psi_s$ in $\Phi$, and $\psi'_m$ and $\psi'_t$ in $\Phi'$.*

c. $\langle \Phi_D | H | \Phi'_D \rangle = (-1)^{|l-m|} \sum_{i \neq l}^{N}(\langle \psi_l \psi_i | \hat{g}_{l,i} | \psi'_m \psi_i \rangle - \langle \psi_l \psi_i | \hat{g}_{l,i} | \psi_i \psi'_m \rangle)$, *if $\Phi$ differs by one orbital, $\psi_l$ in $\Phi$, and $\psi'_m$ in $\Phi'$.*

d. $\langle \Phi_D | H | \Phi'_D \rangle = \sum_{i<j}^{N}(\langle \psi_i \psi_j | \hat{g}_{i,j} | \psi_i \psi_j \rangle - \langle \psi_i \psi_j | \hat{g}_{i,j} | \psi_j \psi_i \rangle)$, *if $\Phi = \Phi'$.*

PROOF.   In general, we have

$$\langle \Phi_D | H | \Phi'_D \rangle = \sum_{i<j}^{N} \sum_{\bar{p}} (-1)^P \langle \psi_1 | \psi'_{P1} \rangle \langle \psi_2 | \psi'_{P2} \rangle \cdots \langle \psi_i \psi_j | \hat{g}_{i,j} | \psi'_{Pi} \psi'_{Pj} \rangle \cdots \langle \psi_N | \psi'_{PN} \rangle$$

$$= \sum_{i \cdot j}^{N} \sum_{\bar{p}} (-1)^P \langle \psi_i \psi_j | \hat{g}_{i,j} | \psi'_{Pi} \psi'_{Pj} \rangle \prod_{k \neq i,j} \langle \psi_k | \psi'_{Pk} \rangle. \tag{112}$$

which is $0$ unless

$$\psi_k = \psi'_{Pk} \quad \forall k \neq i, j \tag{113}$$

because of the orthogonality of the basis set $\psi$.



Figure 28. Energy levels for a parabolic quantum dot with $\hbar\omega = 109$ meV: configuration interaction calculations (solid lines). Hartree-Fock calculations (dashed lines).

Clearly, if $\Phi$ and $\Phi'$ differ by more than two orbitals, there is no permutation operator capable of satisfying the conditions given by Eq. (113), and therefore $\langle \Phi_D | H | \Phi'_D \rangle = 0$. If $\Phi$ and $\Phi'$ differ by two orbitals, $\psi_l \psi_s$ in $\Phi$ and $\psi'_m$, $\psi'_t$ in $\Phi'$, there are only two possible permutations $\hat{P}$ and $\hat{L} = \hat{P}_{l,s} \hat{P}$ satisfying these conditions, where

$$\psi_k = \psi'_{Pk} \quad \forall k \neq l, s.$$

$$\psi'_m = \psi'_{Pl},$$ (114)

$$\psi'_t = \psi'_{Ps},$$

and

$$\psi_k = \psi'_{Lk} \quad \forall k \neq l, s.$$

$$\psi'_m = \psi'_{Ls}.$$ (115)

$$\psi'_t = \psi'_{Lt}.$$

Then we have

$$\langle \Phi_D | \hat{H} | \Phi'_D \rangle = (-1)^P \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_{Pl} \psi'_{Ps} \rangle + (-1)^L \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_{Ll} \psi'_{Ls} \rangle$$

$$= (-1)^P \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_m \psi'_t \rangle + (-1)^L \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_t \psi'_m \rangle$$

$$= (-1)^P ( \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_m \psi'_t \rangle - \langle \psi_l \psi_s | \hat{g}_{l,s} | \psi'_t \psi'_m \rangle ).$$ (116)



Figure 29. Energy levels for a parabolic quantum dot with $\hbar\omega = 54.8$ meV: configuration interaction calculations (solid lines), Hartree-Fock calculations (dashed lines).

noting that $\hat{L} = \hat{P}_{l,s}\hat{P}$, which means $(-1)^{l'} = -(-1)^{l''}$. Also, as it takes $|l - m|$ permutations to align $\psi_l$ and $\psi'_m$ and $|s - t|$ permutations to align $\psi'_j$ and $\psi_s$, we have $(-1)^{l'} = (-1)^{|l-m|+s-t}$.

If $\Phi$ and $\Phi'$ differ by only one orbital, $\psi_l$ in $\Phi$ and $\psi'_m$ in $\Phi'$, the conditions given by Eq. (113) can be satisfied when $i = l$. But $j$ can take on any value allowed by the original definition of $H$. For any given value of $j$, there are two possible permutations that give nonzero results, so again we find

$$\langle \Phi_D | \hat{H} | \Phi'_D \rangle = \sum_{j \neq l}((-1)^P \langle \psi_l \psi_j | \hat{g}_{l,j} | \psi'_m \psi'_{l'j} \rangle + (-1)^l \langle \psi_l \psi_j | \hat{g}_{l,j} | \psi'_{l,j} \psi'_m \rangle)$$

$$= \sum_{j \neq l} \langle \Phi_D | \hat{g}_{l,j} | \Phi'_D \rangle = (-1)^{l-m}(\langle \psi_l \psi_j | \hat{g}_{l,j} | \psi'_m \psi_j \rangle - \langle \psi_l \psi_j | \hat{g}_{l,j} | \psi_j \psi'_m \rangle). \quad (117)$$

If $\Phi$ and $\Phi'$ are identical, no permutation is necessary to align $\Phi$ and $\Phi'$, and all $\hat{g}_{l,j}$ contribute to give

$$\langle \Phi_D | H | \Phi'_D \rangle = \sum_{i<j} \langle \Phi_D | \hat{g}_{l,j} | \Phi'_D \rangle$$

$$= \sum_{i<j}(\langle \psi_i \psi_j | \hat{g}_{l,j} | \psi'_i \psi'_j \rangle - \langle \psi_i \psi_j | \hat{g}_{l,j} | \psi'_j \psi'_i \rangle) \quad (118)$$



**Figure 30.** Energy levels for a parabolic quantum dot with $\hbar\omega = 10.96$ meV: configuration interaction calculations (solid lines), Hartree-Fock calculations (dashed lines).

The pioneering theoretical study on the electronic structure of quantum dots was carried out by Bryant, using the CI method for up to six electrons [83]. Bryant modeled the quantum dot system as an ultrasmall square box, and his calculated energy levels for noninteracting electrons and interacting electrons in various sizes of boxes are shown in Fig. 25. Despite the simplicity of the model, a very rich variety of electron properties, including unpaired electrons, weakly correlated states, and Wigner crystallization, were predicted. Maksym and Chakraborty [84] extended the CI calculations to include the effects of an external magnetic field, which was applied perpendicular to the plane of the dot. A more realistic confinement potential was also used in their calculations and was assumed to be parabolic in the form $V(r) = \frac{1}{2}m\omega_d^2 r^2$. This was believed to be a good approximation to most experimentally studied quantum dot systems, especially when the number of electrons is small. Figure 26 shows their calculated energy levels for three and four electrons at four different magnetic field strengths as a function of $J$, which is the sum of the single-electron angular momentum number $l$. Subsequently, many research groups have used the CI approach to study a wide variety of quantum dots and other nanosystems [20, 50, 51, 63, 85–90].

In this work, we make a direct comparison of the calculated chemical potential using the constant interaction model, the HF method, and the CI method, which is defined by Eq. (77) as the difference between the ground state energy for $N$ and $N - 1$ electrons. As demonstrated in Fig. 27, the HF and CI calculations are in excellent agreement, especially for tightly confined quantum dots (i.e., those with high $\omega$ values), but they become more different as the confinement becomes weaker and electron–electron correlation starts to dominate. The chemical potential provided by the simple formula of the constant interaction model is different from the HF and CI calculation, but the general trend is the same. For comparison purposes, we fit the capacitive term to the spacing between electrons being added to the same orbital as that given by the configuration interaction method at zero field.

Figures 28 to 30 show the lowest energy levels of the three quantum dots containing up to seven electrons. The general trend of increasing accuracy of the HF method with increasing confinement strength is clearly demonstrated. In addition to this, we note that Hartree-Fock calculations are more accurate for spin polarized states. This is because the unpolarized



Figure 31. Energy difference between the fully polarized ($S = 3$) and paramagnetic ($S = 0$) states for a quantum dot with six electrons as a function of $1/r_s = \sqrt{\pi n}$, where $r_s$ is the Wigner-Seitz radii and $n$ the electron density. Reprinted with permission from [63], S. Reimann et al., Phys. Rev. B 62, 8108 (2000). © 2000, American Physical Society.

**Figure 32.** Charge density of a quantum dot with six electrons in a harmonic well: exact result (solid line) and the DFT result (dashed line). Reprinted with permission from [63]. S. Reimann et al., *Phys. Rev. B* 62. 8108 (2000). © 2000, American Physical Society.

configuration has a higher degree of correlation between the motion of the electrons, and therefore the mean field model of Hartree-Fock is less accurate.

Reimann et al. compared the energy levels and electron charge densities obtained using the CI method and the density functional theory. As demonstrated in Figs. 31 and 32, for a quantum dot containing six electrons, the density-functional theory predicts the correct trend for both the energies and the densities. However, the differences are also clearly visible.

## 5. SPIN-ADAPTED CONFIGURATION INTERACTION APPROACH

The CI method takes into account the full interaction and correlation of the electrons in the system as long as the numerical results converge with an increasing number of basis functions. However, the standard CI method involves the calculation of a very large number of interaction integrals and the inversion of large matrixes, which can be prohibitively expensive in terms of computer resources. Reimann et al. had employed matrices of dimensions up to 108,375, with 67,521,121 nonzero elements for a six-electron quantum dot. Calculations for any higher number of electrons in the system were not considered numerically possible using the standard CI formalism [63, 88].

In this section, we develop a novel spin-adapted configuration interaction method to calculate the energy structure of an N-electron quantum dot. By isolating spin eigenstates of the system, which in turn restrict the possible form of the spatial wavefunctions, we are able to reduce significantly the size of the configuration interaction matrices in comparison with the standard CI method. This has allowed us to investigate quantum dot systems with electron numbers greater than six using PCs and Mathematica. A similar approach has been employed in quantum chemistry, and an in-depth analysis of the methodology can be found in Ref. [91]. However, it has not yet been used in mesoscopic physics.

### 5.1. An Orthonormal Spin-Adapted Basis

#### 5.1.1. Spin Eigenfunctions

The projection of the spin angular momentum operator for an $N$-electron system is given by

$$\hat{S}_z = \sum_{i=1}^{N} \hat{S}_z(i).$$                                          (119)

where $\tilde{S}_z(i)$ is the projection of the spin operator for the $i$th electron. The eigenfunctions of $\hat{S}_z$ are products of the single electron eigenfunctions; namely,

$$\theta_k = \sigma(1)\sigma(2)\cdots\sigma(N),\tag{120}$$

where $\sigma(i) = \alpha(i)$ or $\beta(i)$ represents spin up or spin down of the $i$th electron, respectively. We term $\theta_k$ the elementary spin eigenfunctions. The eigenvalue for a particular $\theta_k$ is simply the sum of the eigenvalues of the individual spin operators. If we define $\mu$ as the number of $\alpha$s and $\nu$ as the number of $\beta$s in $\theta_k$, then

$$\tilde{S}_z\theta_k = \frac{(\mu - \nu)\hbar}{2}\theta_k\tag{121}$$

where $\mu + \nu$ is equal to $N$, the number of electrons in the system.

For a given number of electrons, any linear combination of eigenfunctions $\theta_k$ with the same eigenvalues is also an eigenfunction of $\hat{S}_z$ with the same eigenvalue. It is easy to verify that a complete set of eigenfunctions for the particular eigenvalue $(\mu - \nu)\hbar/2$ is given by all possible permutations of the eigenfunction $\theta_1 = \alpha(1)\alpha(2)\cdots\alpha(\mu)\beta(\mu+1)\beta(\mu+2)\cdots\beta(\mu + \nu)$. There are $\binom{N}{\mu}$ of these functions, where $\binom{N}{\mu} = \frac{N!}{\mu!(1-\mu)!}$ is the binomial coefficient, and they form a complete orthonormal basis for eigenfunctions of $\hat{S}_z$ with that particular eigenvalue.

The elementary spin eigenfunctions are generally not eigenfunctions of $\hat{S}^2$, but linear combinations of elementary spin eigenfunctions can be constructed to be simultaneously eigenfunctions of $\hat{S}^2$. In the $N$-electron case, the $\hat{S}^2$ operator has the following form [91]:

$$\hat{S}^2 X = \sum_{i \neq j} P_{ij} X + \frac{N}{4}(4 - N)X\tag{122}$$

where $P_{ij}$ is a permutation operator and $X = \sum_k c_k \theta_k$ are simultaneous eigenfunctions of both $\hat{S}^2$ and $\hat{S}_z$.

Let $X = \sum_k c_k \theta_k(N, M)$ be a linear combination of all elementary spin eigenfunctions of $\hat{S}_z$ for a given number of electrons $N$ and the projection of total spin quantum number $M = (\mu - \nu)/2$. Let $\mathbf{v}$ be the corresponding vector of coefficients; that is, $\mathbf{v} = (c_1, c_2, \ldots, c_{\binom{N}{\mu}})$. Define the matrix $\{\hat{S}^2\}$ by $\{\hat{S}^2\}_{ij} = \langle\theta_i|\hat{S}^2|\theta_j\rangle$, whose elements can be calculated using Eq. (122). Then the action of $\hat{S}^2$ on $X$ can be represented by the matrix operation $\{\hat{S}^2\}\cdot\mathbf{v}$. The eigenfunctions of $\hat{S}^2$ are then found by solving the matrix eigenvalue problem.

As an example, the calculation of the spin eigenfunctions for three electrons with $M = 1/2$ is presented. In this case, the possible elementary spin eigenfunctions are $\theta_1 = \alpha\alpha\beta$, $\theta_2 = \alpha\beta\alpha$, and $\theta_3 = \beta\alpha\alpha$.

$$\hat{S}^2|\theta_1\rangle = \sum_{i \neq j} P_{ij}|\theta_1\rangle + \frac{N}{4}(4 - N)|\theta_1\rangle$$

$$= P_{12}|\alpha\alpha\beta\rangle + P_{13}|\alpha\alpha\beta\rangle + P_{23}|\alpha\alpha\beta\rangle + \frac{3}{4}(4 - 3)|\alpha\alpha\beta\rangle$$

$$= |\alpha\alpha\beta\rangle + |\beta\alpha\alpha\rangle + |\alpha\beta\alpha\rangle + \frac{3}{4}|\alpha\alpha\beta\rangle$$

$$= \frac{7}{4}|\theta_1\rangle + |\theta_2\rangle + |\theta_3\rangle\tag{123}$$

Similar calculations lead to the matrix elements of $\{\hat{S}^2\}$:

$$\{\hat{S}^2\} = \begin{pmatrix} \frac{7}{4} & 1 & 1 \\ 1 & \frac{7}{4} & 1 \\ 1 & 1 & \frac{7}{4} \end{pmatrix}\tag{124}$$

The eigenvalues of the above matrix are $\frac{15}{4} = \frac{3}{2}(\frac{3}{2} + 1)$ and $\frac{3}{4} = \frac{1}{2}(\frac{1}{2} + 1)$, which corresponds to the total angular momentum quantum numbers $S = \frac{3}{2}$ and $S = \frac{1}{2}$, respectively. The orthonormalized spin eigenfunctions are obtained by first solving for the eigenvectors and then applying the Gram-Schmidt orthonormalization procedure. These are

$$X_1\left(3, \frac{3}{2}, \frac{1}{2}\right) = \frac{1}{\sqrt{3}}(\alpha\alpha\beta + \alpha\beta\alpha + \beta\alpha\alpha) \tag{125}$$

$$X_1\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{1}{\sqrt{2}}(-\alpha\alpha\beta + \beta\alpha\alpha) \tag{126}$$

$$X_2\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{1}{\sqrt{6}}(-\alpha\alpha\beta + 2\alpha\beta\alpha - \beta\alpha\alpha) \tag{127}$$

Here $X_k(N, S, M)$ are termed as spin eigenfunctions as opposed to the elementary spin eigenfunctions, and $k$ is a positive integer denoting different eigenfunctions in a multidimensional spin eigenspace.

## 5.1.2. Spin Adapted Basis

Next we need to include the spatial wavefunctions to form an orthonormal and properly antisymmetrized basis. The basis elements have the form $C_0 A \Phi X_k(N, S, M)$, where $C_0$ is a normalization constant, $A = \frac{1}{\sqrt{N!}}(-1)^P P$ is the antisymmetrizer, $X_k(N, S, M)$ is the spin eigenfunction, $\Phi = \psi_1\psi_2\cdots\psi_N$ is a representative spatial wavefunction, and $\psi$ are single-electron orbitals (i.e., eigenfunctions of the single-electron Hamiltonian). We occasionally abbreviate $X_k(N, S, M)$ to $X_k$, when the values of $N$, $S$, and $M$ are fixed in a given calculation. Different elements of this basis have different spatial and spin combinations.

Note that both the spin and spatial variables are permuted in $C_0 A \Phi X_k$, as required by the antisymmetry principle. We also note that any product spatial wavefunction with three or more orbitals that are the same will vanish when multiplied by a spin eigenfunction and antisymmetrized. This implies that the greatest number of orbitals that can be the same in a representative spatial wavefunction is two, which we call a doubly occupied orbital.

In the following, we demonstrate that by using the spin eigenfunctions that are also eigenfunctions of the permutation operators $P_{2i-1,2i}$, for $i = 1, 2, \ldots$, where $i \leq \frac{N}{2}$ and spatial wavefunctions that only have doubly occupied orbitals sequentially at positions $(2i - 1, 2i)$ for $i = 1, 2, \ldots$; where $i \leq \frac{N}{2}$, we can construct a orthogonal basis in which vanishing elements can be easily identified and the normalization constant is of a simple form. This is similar to the Serber construction [92]. Following the nomenclature of Salmon and Ruedenberg [93], as also adopted by Pauncz [91], we will call a pair of numbers of the form $(2i, 2i - 1)$—a geminal pair. Correspondingly, we call the set of permutations mentioned in the previous paragraph the geminal permutations.

Because the $S^2$ operator commutes with any permutation and the geminal transposition operators pairwise commute, it is possible to construct a complete set of orthonormal spin eigenfunctions that commute with all of these operators. To compute this basis, we first obtain the orthonormal spin eigenfunctions, $X^0 = \{X_1^0, \ldots, X_f^0\}$, as in the previous section. We then use these to compute the representation matrix $U(P_{12})$. The representation matrix for any permutation $P$ is defined by

$$P X_k(N, S, M) = \sum_{i=1}^{f} U(P)_{ik} X_i(N, S, M) \tag{128}$$

It can be shown that for any two permutations $R$ and $P$,

$$U(RP) = U(R)U(P) \tag{129}$$

satisfying the condition for a representation of the symmetric group. This is independent of the particular spin basis chosen.

Note that $X''_k(N, S, M)$ are not, in general, eigenfunctions of $P_{12}$. However, the orthonormalized eigenvectors of $U(P_{12})$ will give the linear combination of vectors of $X''$ that form a new orthonormal basis, $X' = \{X'_1, \ldots, X'_j\}$, for the spin space that are also eigenfunctions of $P_{12}$.

For example, we can start with the three electron spin eigenfunctions, with $S = \frac{1}{2}$ and $M = \frac{1}{2}$ given by Eqs. (126) and (127), and obtain

$$P_{12}X''_1\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{1}{2}X''_1\left(3, \frac{1}{2}, \frac{1}{2}\right) + \frac{\sqrt{3}}{2}X''_2\left(3, \frac{1}{2}, \frac{1}{2}\right),$$

$$P_{12}X''_2\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{\sqrt{3}}{2}X''_1\left(3, \frac{1}{2}, \frac{1}{2}\right) - \frac{1}{2}X''_2\left(3, \frac{1}{2}, \frac{1}{2}\right)$$

(130)

Therefore

$$U(P_{12}) = \begin{pmatrix} \frac{1}{2} & \frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{-1}{2} \end{pmatrix}$$

(131)

The eigenvalues of the above matrix are $\pm 1$, corresponding to the eigenvectors $(\frac{\sqrt{3}}{2}, \frac{1}{2})$ and $(-\frac{1}{2}, \frac{\sqrt{3}}{2})$. This means that the new spin eigenfunction basis is

$$X'_1\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{\sqrt{3}}{2}X''_1\left(3, \frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2}X''_2\left(3, \frac{1}{2}, \frac{1}{2}\right)$$

$$= \frac{1}{\sqrt{6}}(-2\alpha\alpha\beta + \alpha\beta\alpha + \beta\alpha\alpha)$$

(132)

$$X'_2\left(3, \frac{1}{2}, \frac{1}{2}\right) = \frac{-1}{2}X''_1\left(3, \frac{1}{2}, \frac{1}{2}\right) + \frac{\sqrt{3}}{2}X''_2\left(3, \frac{1}{2}, \frac{1}{2}\right)$$

$$= \frac{1}{\sqrt{2}}(\alpha\beta\alpha - \beta\alpha\alpha)$$

(133)

The basis elements are now eigenfunctions of $P_{12}$. Note that the first spin eigenfunction is symmetric in the first two electrons, whereas the second is antisymmetric. This means that the first will vanish if we multiply it by a spatial wavefunction in which the first two electrons are in the same orbital and then antisymmetrize. This will not be the case for the second spin eigenfunction; thus, we are allowed to combine it with spatial wavefunctions in which the first two orbitals are the same.

The symmetry properties of the spin eigenfunctions dictate with which spatial orbitals it can be combined. This is a general feature that is developed in the next section. Specifically, we will see that the choice of spin eigenfunctions that are simultaneous eigenfunctions of geminal transpositions allows us to restrict the spatial representative wavefunctions we need to consider to those with the doubly occupied orbitals in the geminal positions.

If we have more than three electrons, we can use $X'$ to calculate the representation matrix $U(P_{34})$, which creates a new basis $X^2$. This process continues till we use all possible geminal transposition operators. The pairwise commutative property of the geminal transposition operators guarantees that subsequent bases will still be eigenfunctions of earlier geminal transpositions after each transformation. After this process, we will have produced a spin-adapted basis, $X = \{X_1, \ldots, X_j\}$, in which the spin eigenfunctions are simultaneous eigenfunctions of all the geminal transposition operators.

There are other construction methods, such as the Serber construction [92], which produce bases with the same properties. In the Serber construction, spin eigenfunctions are constructed by adding two electron singlet and triplet spin eigenfunctions using the angular momentum addition formula and the Clebsch-Gordan coefficients. Degenerate eigenfunctions can be traced using their Serber path symbol [91, 94]. For an odd number of electrons,

the last electron can be added using the one electron–addition formulas. This produces a orthonormal spin eigenspace, with the spin eigenfunctions also being eigenfunctions of the geminal transposition operators.

As the calculation of the spin eigenfunctions in our studies takes only a small proportion of the computational effort compared to the calculation of the Hamiltonian matrix, the above matrix eigenvalue method was chosen for simplicity of implementation. The Serber method may prove useful for systems with larger numbers of electrons.

## 5.2. Properties of the Spin-Adapted Basis

In this section, we discuss the general properties of the spin-adapted basis $X = \{X_1, \ldots, X_f\}$. Each $X_k$ is a simultaneous eigenfunction of $\tilde{S}^2$ and the geminal transposition operators. We also derive the formula for the normalization constant $C_\Phi$.

We are restricting the position of any doubly occupied orbitals in the spatial wavefunctions to be sequentially in the geminal positions. Also, two representative wavefunctions for different basis elements are never noninvariant permutations of each other. With this restriction, we show that we can construct an orthonormal basis.

The following two theorems were used in proving the properties of the spin adapted basis:

THEOREM 5.1.   Let $\Phi = \psi_1 \psi_2 \cdots \psi_n$ with $\psi_{2i-1} = \psi_{2i}$ (i.e., there is a doubly occupied orbital at this position). Then if $A\Phi X_k \neq 0$ we have $P_{2i-1,2i} X_k = -X_k$ and $U(P_{2i-1,2i})_{kk} = -1$.

The above theorem is simply a result of the antisymmetry principle. As the spatial wavefunction is symmetric under $P_{2i-1,2i}$, the spin eigenfunction must be antisymmetric to compensate. The condition $A\Phi X_k \neq 0$ is included because we can not have zero elements in the basis if we want a well-constructed matrix eigenvalue problem.

If the spatial wavefunction $\Phi$ has $d$ pairs of doubly occupied orbitals, then the invariance group $S_\Phi$ is given by all products of transpositions in the set $\{P_{1,2}, P_{3,4}, \ldots P_{2d-1,2d}\}$. As each of these elements is of order 2 (i.e., $P_{i,j}^2$ is the identity permutation), this set has $2^d$ elements given by $P = P_{1,2}^{n_1} * P_{3,4}^{n_2} * \cdots * P_{2d-1,2d}^{n_d}$ where $n_i = 0$ or $1$, $i = 1, \ldots, d$.

THEOREM 5.2.   Let $\Phi = \psi_1 \psi_2 \cdots \psi_n$ with $\psi_1 = \psi_2, \ldots, \psi_{2d-1} = \psi_{2d}$ where $d = II(\Phi)$ is the number of doubly occupied orbitals in $\Phi$. Then if $P \in S_\Phi = \{P \in S_N : P\Phi = \Phi\}$, with $A\Phi X_k \neq 0$ we have $U(P)_{kk} = (-1)^P$ where $(-1)^P$ is the parity of the permutation $P$.

PROOF.   Now if $P \in S_\Phi$ then $P = P_{1,2}^{n_1} * P_{3,4}^{n_2} * \cdots * P_{2d-1,2d}^{n_d}$ where $n_i = 0$ or $1$, $i = 1, \ldots, d$. Thus, using the multiplicative property of the representation matrices,

$$U(P) = U(P_{1,2}^{n_1})U(P_{3,4}^{n_2}) \ldots U(P_{2d-1,2d}^{n_d})$$

$$= U(P_{1,2})^{n_1} U(P_{3,4})^{n_2} \ldots U(P_{2d-1,2d})^{n_d}$$

where $U(P)^0$ is the identity matrix. The $U(P_{2i-1,2i})$ are diagonal, so

$$U(P)_{kk} = U(P_{1,2})_{kk}^{n_1} U(P_{3,4})_{kk}^{n_2} \ldots U(P_{2d-1,2d})_{kk}^{n_d}$$

$$= (-1)^{n_1} (-1)^{n_2} \ldots (-1)^{n_d}$$

$$= (-1)^{n_1 + n_2 + \cdots + n_d} \tag{134}$$

where we have used $U(P_{2i-1,2i})_{kk} = -1$ [Theorem 5.1]. The parity of $P$ is given by $(-1)^P = (-1)^{n_1 + n_2 + \cdots + n_d}$, so $U(P)_{kk} = (-1)^P$, which is the result we require.

As the antisymmetrizer is Hermitian and proportional to its square [91], we can simplify the inner product of basis elements.

THEOREM 5.3.   The basis functions $A\Phi X_j$ and $A\Psi X_k$ are orthogonal if $\Phi \neq P\Psi$ for all permutations $P$.

PROOF. As the antisymmetrizer is Hermitian and proportional to its square [91], we can simplify the inner product of basis elements:

$$\langle A\Phi X_j | A\Psi X_k \rangle = \langle \Phi X_j | A A^* \Psi X_k \rangle$$

$$= \langle \Phi X_j | A^2 \Psi X_k \rangle$$

$$= \langle \Phi X_j | \sqrt{N} A \Psi X_k \rangle$$

$$= \sum_P (-1)^P \langle \Phi | P\Psi \rangle \langle X_j | P X_k \rangle$$

$$= \sum_P (-1)^P \langle \Phi | P\Psi \rangle U(P)_{jk}, \qquad (135)$$

but $\langle \Phi | P\Psi \rangle = 0 \ \forall P \in S_N$. Thus,

$$\langle A\Phi X_j | A\Psi X_k \rangle = 0. \qquad (136)$$

This means that basis elements are orthogonal for different spatial wavefunctions. Note that the basis does not include spatial wavefunctions that are noninvariant permutations of each other. The only case in which two basis elements would have the same spatial wavefunction is when it is multiplied by a different spin eigenfunction. This case is dealt with by the next theorem.

THEOREM 5.4. *The basis elements* $A\Phi X_j$ *and* $A\Phi X_k$ *are orthogonal where* $A\Phi X_j$, $A\Phi X_k \neq 0$. *That is, for* $j \neq k$ *their inner product is* 0. *Furthermore, for* $j = k$, *their inner product is* $2^d$, *where* $d$ *is the number of pairs of doubly occupied orbitals in* $\Phi$.

PROOF. Because

$$\langle A\Phi X_j | A\Phi X_k \rangle = \sum_P (-1)^P \langle \Phi | P\Phi \rangle U(P)_{jk} \qquad (137)$$

and

$$\langle \Phi | P\Phi \rangle = \begin{cases} 1, & P \in S_\Phi \\ 0, & \text{otherwise}, \end{cases} \qquad (138)$$

Eq. (137) transforms to:

$$\langle A\Phi X_j | A\Phi X_k \rangle = \sum_{P \cdot S_\Phi} (-1)^P U(P)_{jk}. \qquad (139)$$

If $P \in S_\Phi$, then $P = P_{1,2}^{n_1} P_{3,4}^{n_2} \cdots P_{2d-1,2d}^{n_d}$, where $n_i = 0$ or 1, $i = 1, \ldots, d$ and $d = \Pi(\Phi)$ is the number of doubly occupied orbitals in $\Phi$.

If $j \neq k$, then

$$\langle A\Phi X_j | A\Phi X_k \rangle = 0 \qquad (140)$$

as $U(P)_{jk}$ is diagonal for $P \in S_\Phi$.

If $j = k$, then

$$\langle A\Phi X_k | A\Phi X_k \rangle = \sum_{P \in S_\Phi} (-1)^P U(P)_{kk}$$

$$= \sum_{P \in S_\Phi} (-1)^P (-1)^P$$

$$= |S_\Phi|$$

$$= 2^d, \qquad (141)$$

where we have used theorem 5.2, which states $U(P)_{kk} = (-1)^P$ when $A\Phi X_k \neq 0$ and $P \in S_\Phi$. Thus, $C_\Phi A\Phi X_k$ is a properly normalized basis function with $C_\Phi = 1/\sqrt{2^d}$.

THEOREM 5.5. *The basis function* $A(P^i\Phi)X_k$ *is a linear combination of the basis functions* $A\Phi X_i$ *for* $i = 1, \ldots, f$ *where* $f$ *is the dimension of the spin eigenfunction space.*

PROOF.  Let $P'$ be the part of permutation $P$ that acts on the spatial component of the basis function and $P''$ be the part that acts on the spin component.

$$A(P'\Phi)X_k = \frac{1}{\sqrt{N!}}\sum_I (-1)^q Q^r P^r \Phi Q'' X_k$$

$$= \frac{1}{\sqrt{N!}}\sum_I (-1)^q T^r \Phi T''(P^{-1})'' X_k$$

$$= \frac{1}{\sqrt{N!}}\sum_I (-1)^q T^r \Phi T'' \sum_{m=1}^I U(P^{-1})_{mk} X_m$$

$$= \sum_{m=1}^I U(P^{-1})_{mk} \frac{1}{\sqrt{N!}}\sum_I (-1)^t(-1)^{q-t} T^r \Phi T'' X_m$$

$$= \sum_{m=1}^I (-1)^r U(P^{-1})_{mk} \frac{1}{\sqrt{N!}}\sum_I (-1)^t T^r \Phi T'' X_m$$

$$= \sum_{m=1}^I (-1)^r U(P^{-1})_{mk} A\Phi X_m. \tag{142}$$

noting that $T = QP$. Thus, $A(P\Phi)X_k$ is a linear combination of the basis functions $A\Phi X_m$, with the coefficients given by $(-1)^r U(P^{-1})_{mk}$.

This means that the spatial wavefunctions from two different basis elements are never noninvariant permutations of each other. Therefore, if representative spatial wavefunctions have doubly occupied orbitals, we can choose to put the doubly occupied orbitals sequentially in the geminal positions. Thus, if $\Phi$ has $d$ doubly occupied orbitals, they will be in the positions $(1, 2), (3, 4)\cdots(2d - 1, 2d)$. For example, we would use $\psi_1\psi_1\psi_2\psi_2\psi_3$ but not $\psi_1\psi_2\psi_1\psi_3\psi_2$.

To construct a basis for the matrix eigenvalue problem for $N$ electron, we first calculate the spin-adapted basis. We must then choose a set of representative spatial wavefunctions that are products of the single-electron energy eigenfunctions, making sure any doubly occupied orbitals appear first and that none are a noninvariant permutation of another in the set. We then combine the spatial wavefunctions with each spin eigenfunction and antisymmetrize, eliminating any basis elements that are zero. Last, we multiply each basis element by the appropriate normalization constant. The accuracy of the final calculation will depend on which spatial wavefunctions are chosen but will improve by increasing basis size.

## 5.3. Simplifying the Hamiltonian Elements

We wish to use the basis described in the previous section to calculate the Hamiltonian matrix elements of the form $\langle C_q A\Phi X_j|\hat{H}|C_\psi A\Psi X_k\rangle$, where $\hat{H} = \hat{H}_0 + \hat{H}_{int}$; $\hat{H}_0 = \sum_{i=1}^N \hat{H}_{0i}$ is the single-electron component and $\hat{H}_{int} = \sum_{i,j}^N \hat{H}(i,j)$ is the interaction component that acts pairwise. This is the general form of a spin-free Hamiltonian in which spin–orbit and spin–spin interactions are neglected. This gives us

$$\langle C_q A\Phi X_j|\hat{H}|C_\psi A\Psi X_k\rangle = \langle C_q A\Phi X_j|\hat{H}_0 + \hat{H}_{int}|C_\psi A\Psi X_k\rangle$$

$$= \langle C_q A\Phi X_j|\hat{H}_0|C_\psi A\Psi X_k\rangle + \langle C_q A\Phi X_j|\hat{H}_{int}|C_\psi A\Psi X_k\rangle. \tag{143}$$

### 5.3.1. Single-Electron Integral

The single-electron integral is quite straightforward, noting that the single-electron orbitals are the eigenfunctions of $\hat{H}_{0i}$, that is, $\hat{H}_{0i}\Phi(r_i) = E_i\Phi(r_i)$, where $E_i$ is the energy eigenvalue. Therefore,

$$\langle C_q A\Phi X_j|\hat{H}_0|C_\psi A\Psi X_k\rangle = C_q C_\psi \sum_P (-1)^p \langle\Phi|\hat{H}_0|P\Psi\rangle U(P)_{jk}$$

$$= C_q C_\psi \sum_P \sum_i (-1)^p E_{P_{\psi_i}}\langle\Phi|P\Psi\rangle U(P)_{jk}. \tag{144}$$

where $\hat{H}_0(P\Psi) = E_{P(i)}P\Psi$, and $E_{P(i)}$ is the single-electron energy eigenvalue of the $i$th orbital in $P\Psi$. However, $\langle \Phi | P\Psi \rangle = 0$ $\forall P$ unless $\Phi = P\Psi$ for some $P \in S_{ii}$. In the basis construcion we only need to use representative wavefunctions that are not permutations of each other. Thus, if $\Phi = P\Psi$ for some $P$, then $P \in S_{ii}$; that is,

$$\langle C_\Phi A \Phi X_j | \hat{H}_0 | C_\Psi A \Psi X_k \rangle = 0 \quad \text{for} \quad \Phi \neq \Psi. \tag{145}$$

Otherwise,

$$\langle C_\Phi A \Phi X_j | \hat{H}_0 | C_\Phi A \Phi X_k \rangle = C_\Phi^2 \sum_{i=1}^{N} \sum_{P \in S_\Phi} (-1)^P E_{P(i)} \langle \Phi | \Phi \rangle U(P)_{jk}$$

$$= C_\Phi^2 \sum_{i=1}^{N} 2^{H(\Phi)} E_i \delta_{jk}$$

$$= \delta_{jk} \sum_{i=1}^{N} E_i \tag{146}$$

where $\delta_{ij}$ is the Kronecker delta. In other words, if the basis functions $A\Phi X_j$ and $A\Psi X_k$ are the same in both the spin and spatial components, that is, $\Phi = \Psi$ and $j = k$, then $\langle C_\Phi A \Phi X_j | \hat{H}_0 | C_\Phi A \Phi X_j \rangle$ is the sum of single-electron eigenvalues $\sum_{i=1}^{N} E_i$. Otherwise, $\langle C_\Phi A \Phi X_j | \hat{H}_0 | C_\Psi A \Psi X_k \rangle = 0$.

## 5.3.2. Reduction of the Sum Over Permutations in the Interaction Integral

Because the interaction Hamiltonian acts pairwise, that is, $H_{int} = \sum_{i > j} H(i, j)$, where $H(i, j)$ is the incraction term between the $i$th and $j$th electron,

$$\langle \Phi | \hat{H}(i, j) | \Psi \rangle = \langle \psi_1 \psi_2 \cdots \psi_N | \hat{H}(i, j) | \phi_1 \phi_2 \cdots \phi_N \rangle$$

$$= \langle \psi_i \psi_j | \hat{H}(i, j) | \phi_i \phi_j \rangle \prod_{k \neq i, j}^{N} \langle \psi_k | \phi_k \rangle \tag{147}$$

Thus $\langle \Phi | \hat{H}(i, j) | \Psi \rangle = 0$, unless $\Phi_k = \Psi_k$ $\forall k \neq i, j$.

Consequently,

$$\langle C_\Phi A \Phi X_k | \hat{H}_{int} | C_\Psi A \Psi X_l \rangle = C_\Phi C_\Psi \sum_P (-1)^P \langle \Phi | \hat{H}_{int} | P\Psi \rangle U(P)_{kl}$$

$$= C_\Phi C_\Psi \sum_{i > j} \sum_P (-1)^P \langle \Phi | \hat{H}(i, j) | P\Psi \rangle U(P)_{kl}$$

$$= C_\Phi C_\Psi \sum_{i > j} \sum_P \left( (-1)^P \langle \psi_i \psi_j | \hat{H}(i, j) | (P\Psi)_i (P\Psi)_j \rangle \right.$$

$$\left. \times \prod_{m \neq i, j} \langle \psi_m | (P\Psi)_m \rangle U(P)_{kl} \right) \tag{148}$$

where we have used the notation that $(P\Psi)_i$ is the $i$th orbital in $P\Psi$; for example, if $\Psi = \psi_1 \psi_2 \psi_3$ and $P\Psi = \psi_2 \psi_3 \psi_1$, then $(P\Psi)_2 = \psi_3$ as $\psi_3$ is the second orbital in $P\Psi$.

Let $L_{i,j}$ be a permutation operator that aligns $\Phi$ and $\Psi$, such that $\psi_k = (L_{i,j} \Psi)_k$ $\forall k \neq i, j$. If this permutation does not exist for a particular $\Phi$, $\Psi$, $i$, and $j$, then define $L_{ij}$ as 0 and $U[L_{ij}]_{kl} = 0$. For this permutation, $\prod_{k \neq i, j} \langle \Phi_k | (L_{i,j} \Psi)_k \rangle = 1$.

This permutation will not be unique as the set of permutations $L_{i,j} Q$ and $P_{i,j} L_{i,j} Q$, for $Q \in S_\Psi$ will also be line up permutations, where $P_{i,j}$ is the transposition of the electron $i$ and $j$. These sets are not necessarily distinct, as $L_{i,j} Q$ and $P_{i,j} L_{i,j} Q$ will produce the same integral if $(P\Psi)_i = (P\Psi)_j$.

Define $\omega_i$ as $(L_{i,j}\Psi)_i$ and $\omega_j$ as $(L_{i,j}\Psi)_j$. For all permutations not in $L_{i,j}Q$ or $P_{i,j}L_{i,j}Q$, we have $\prod_{k \neq i,j}\langle\Phi_k|(P\Psi)_k\rangle = 0$. Therefore,

$$\langle C_\Phi A\Phi X_j|\hat{H}_{int}|C_\Psi A\Psi X_k\rangle$$

$$= C_\Phi C_\Psi \sum_{i \cdot j}\sum_{Q \in S_\Psi} 2^{-\delta(\omega_i, \omega_j)}[(-1)^{l \cdot q}U(L_{i,j}Q)_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$+ (-1)^{p \cdot l \cdot q}U(P_{i,j}L_{i,j}Q)_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle]$$

$$= C_\Phi C_\Psi \sum_{i \cdot j}\sum_{Q \in S_\Psi} 2^{-\delta(\omega_i, \omega_j)}((-1)^{l}(-1)^q U(L_{i,j})_{km}U(Q)_{ml}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- (-1)^{l}(-1)^q U(P_{i,j}L_{i,j})_{km}U(Q)_{ml}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle)$$

$$= C_\Phi C_\Psi \sum_{i \cdot j}\sum_{Q \in S_\Psi} 2^{-\delta(\omega_i, \omega_j)}((-1)^{l}(-1)^q U(L_{i,j})_{kl}(-1)^q\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- (-1)^{l}(-1)^q U(P_{i,j}L_{i,j})_{kl}(-1)^q\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle)$$

$$= C_\Phi C_\Psi \sum_{i \cdot j}|S_\Psi|(-1)^{l} \cdot 2^{-\delta(\omega_i, \omega_j)}(U(L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- U(P_{i,j}L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle)$$

$$= \frac{C_\Phi}{C_\Psi}\sum_{i \cdot j}(-1)^{l} \cdot 2^{-\delta(\omega_i, \omega_j)}(U(L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- U(P_{i,j}L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle) \qquad (149)$$

which is a reduction of Eq. (148) using the line-up permutation. In particular, the $2^{-\delta(\omega_i, \omega_j)}$ factor makes sure integrals are not counted twice. The sum over all permutations is removed as permutations produce vanishing integrals because of the orthogonality of the orbitals.

## 5.3.3. Hamiltonian Elements in Special Cases

Let us first define the orbital difference between two spatial wavefunctions $\Phi$ and $\Psi$ as the number of orbitals that appear in $\Psi$ but do not appear in the corresponding number of times in $\Phi$. For example, the orbital difference between $\Phi = \psi_1\psi_1\psi_2\psi_3\psi_4$ and $\Psi = \psi_1\psi_2\psi_3\psi_4\psi_5$ is 1, because the orbital $\psi_5$ appears in $\Psi$ but not in $\Phi$, whereas the other four orbitals are in both functions.

We can now simplify Eq. (149) further, based on the orbital difference of the two spatial functions $\Phi$ and $\Psi$ in $\langle C_\Phi A\Phi X_j|\hat{H}_{int}|C_\Psi A\Psi X_k\rangle$. The following formulas are general in that they will work for any set of orthonormal single-electron orbitals. They are similar to the Slater-Condon rules for matrix elements between Slater determinates.

**Case 1. Orbital Difference Equals Zero**   With the basis we have chosen, if the orbital difference of two spatial wavefunctions is zero, the orbitals must be equal. This is because we are not including spatial wavefunctions that are noninvariant permutations of each other in the basis. This means the line-up permutation $L_{i,j}$ will be the same for any pair $(i, j)$; namely, the identity permutation $P_0$, the corresponding matrix of which is the identity matrix. Therefore, we have

$$\langle C_\Phi A\Phi X_j|\hat{H}_{int}|C_\Psi A\Psi X_k\rangle = \frac{C_\Phi}{C_\Psi}\sum_{i \cdot j}(-1)^{l} \cdot 2^{-\delta(\omega_i, \omega_j)}(U(L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- U(P_{i,j}L_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle)$$

$$= \frac{C_\Phi}{C_\Psi}\sum_{i \cdot j} 2^{-\delta(\omega_i, \omega_j)}(\delta_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_i\omega_j\rangle$$

$$- U(P_{i,j})_{kl}\langle\psi_i\psi_j|\hat{H}(i,j)|\omega_j\omega_i\rangle) \qquad (150)$$

**Case 2. Orbital Difference Equal to One**   In this case there is one orbital in $\Phi$ that does not appear in $\Psi$. We label this orbital $\psi_{\mathrm{dif}}$. The only pairs $(i, j)$ in Eq. (149) that give nonzero results are those containing the orbital $\psi_{\mathrm{dif}}$. If $\psi_{\mathrm{dif}}$ appears once in position $i$, then these pairs are

$$(1, i), (2, i), \ldots, (i - 1, i), (i, i + 1), \ldots, (i, N)$$

If $\psi_{\mathrm{dif}}$ appears twice in positions $i$ and $i + 1$, remembering that doubly occupied orbitals are next to each other in the basis, then these pairs are

$$(1, i), (2, i), \ldots, (i - 1, i), (i, i + 1), \ldots, (i, N) \quad \text{and}$$

$$(1, i + 1), (2, i + 1), \ldots, (i - 1, i + 1), (i + 1, i + 2), \ldots, (i + 1, N).$$

It is important to note that $\psi_i = \psi_{i+1} = \psi_{\mathrm{dif}}$, so the integrals in the second case are exactly the same as the integrals in the first case except the pair $(i, i + 1)$ only appears once. We also note that if $L_{i,j}$ is the line-up permutation for the pair $(i, j)$, then $P_{i,i+1}L_{i,j}$ is the line-up permutation for the pair $(i + 1, j)$. Using this information and Eq. (149), we obtain

$$\langle C_\Phi A\Phi X_k | \hat{H}_{\mathrm{int}} | C_\Psi A\Psi X_l \rangle$$

$$= \frac{C_\Phi}{C_\Psi} \sum_{j=1, j \neq i}^{N} 2^{-\delta(\omega_i, \omega_i)} (-1)^{l_{i,j}} \big( U(L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_i \omega_j \rangle$$

$$- U(P_{i,j}L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_j \omega_i \rangle \big) + (n_{\mathrm{dif}} - 1)$$

$$\times \frac{C_\Phi}{C_\Psi} \Bigg( \sum_{j=1, j \neq i+1}^{N} \big[ 2^{-\delta(\omega_{i+1}, \omega_i)} (-1)^{p_{i,i+1}+l_{i,j}} \big( U(P_{i,i+1}L_{i,j})_{kl} \langle \psi_{i+1} \psi_j | \hat{H}(i+1,j) | \omega_{i+1} \omega_j \rangle$$

$$- U(P_{i,i+1}P_{i,j}L_{i,j})_{kl} \langle \psi_{i+1} \psi_j | \hat{H}(i+1,j) | \omega_j \omega_{i+1} \rangle \big) \big]$$

$$- 2^{-\delta(\omega_{i+1}, \omega_i)} (-1)^{l_{i,i+1}} \big[ U(L_{i,i+1})_{kl} \langle \psi_i \psi_{i+1} | \hat{H}(i,i+1) | \omega_i \omega_{i+1} \rangle$$

$$- U(P_{i,i+1}L_{i,i+1})_{kl} \langle \psi_i \psi_{i+1} | \hat{H}(i,i+1) | \omega_{i+1} \omega_i \rangle \big] \Bigg), \quad (151)$$

where $n_{\mathrm{dif}}$ is the number of times the orbital $\psi_{\mathrm{dif}}$ appears in $\Phi$. The $(n_{\mathrm{dif}} - 1)$ term accounts for the cases in which $\psi_{\mathrm{dif}}$ occurs twice in $\Phi$.

Using the fact that $\psi_i = \psi_{i+1}$ and

$$U(P_{i,i+1}R)_{kl} = U(P_{i,i+1})_{kk} U(R)_{kl} = (-1)^{p_{i,i+1}} U(R)_{kl} = -U(R)_{kl},$$

we obtain

$$\langle C_\Phi A\Phi X_k | \hat{H}_{\mathrm{int}} | C_\Psi A\Psi X_l \rangle$$

$$= \frac{C_\Phi}{C_\Psi} \sum_{j=1, j \neq i}^{N} 2^{-\delta(\omega_i, \omega_i)} (-1)^{l_{i,j}} \big( U(L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_i \omega_j \rangle$$

$$- U(P_{i,j}L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_j \omega_i \rangle \big) + (n_{\mathrm{dif}} - 1)$$

$$\times \frac{C_\Phi}{C_\Psi} \Bigg( \sum_{j=1, j \neq i}^{N} 2^{-\delta(\omega_i, \omega_i)} (-1)(-1)^{l_{i,j}} \big( (-1) U(L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_i \omega_j \rangle$$

$$- (-1) U(P_{i,j}L_{i,j})_{kl} \langle \psi_i \psi_j | \hat{H}(i,j) | \omega_j \omega_i \rangle \big)$$

$$- 2^{-\delta(\omega_{i+1}, \omega_i)} (-1)^{l_{i,i+1}} \big[ U(L_{i,i+1})_{kl} \langle \psi_i \psi_{i+1} | \hat{H}(i,i+1) | \omega_i \omega_{i+1} \rangle$$

$$- U(P_{i,i+1}L_{i,i+1})_{kl} \langle \psi_i \psi_{i+1} | \hat{H}(i,i+1) | \omega_{i+1} \omega_i \rangle \big] \Bigg). \quad (152)$$

Grouping terms and simplifying further, we get

$$\langle C_\Phi A\Phi X_k | \hat{H}_{int} | C_\Psi A\Psi X_l \rangle$$

$$= \frac{C_\Phi}{C_\Psi} n_{diff} \sum_{j=1, j=i}^{N} 2^{-\delta(\omega_i,\omega_j)} (-1)^{l} \{ U(L_{i,j})_{kl} \langle \psi_i\psi_j | \hat{H}(i,j) | \omega_j\omega_i \rangle$$

$$- U(P_{i,j}L_{i,j})_{kl} \langle \psi_i\psi_j | \hat{H}(i,j) | \omega_j\omega_i \rangle \}$$

$$- \frac{C_\Phi}{C_\Psi} (n_{diff} - 1) \{ 2^{-\delta(\omega_i,\omega_{i+1})} (-1)^{l+1} \{ U(L_{i,i+1})_{kl} \langle \psi_i\psi_{i+1} | \hat{H}(i,i+1) | \omega_i\omega_{i+1} \rangle$$

$$- U(P_{i,i+1}L_{i,i-1})_{kl} \langle \psi_i\psi_{i-1} | \hat{H}(i,i+1) | \omega_{i-1}\omega_i \rangle \}\}$$

$$(153)$$

Note that we have used many of the mathematical properties developed earlier to obtain the above formula.

**Case 3. Orbital Difference Equal to Two**  If the orbital difference of the $\Phi$ and $\Psi$ is 2, then only one selection of orbitals $\psi_i$ and $\psi_j$ will give a valid line-up permutation; that is, the orbitals that appear in $\Phi$ but do not appear in $\Psi$. Although there is only one choice for the orbitals, this could correspond to more than one pair of $i$, $j$, as each of $\psi_i$ and $\psi_j$ could appear up to twice in $\Phi$.

Let $v_i$ and $v_j$ denote the orbitals in $\Phi$ that do not appear in $\Psi$. Let $n(v, \Phi)$ be the number of times the orbital $v$ appears in the function $\Phi$.

$$\langle C_\Phi A\Phi X_j | \hat{H}_{int} | C_\Psi A\Psi X_k \rangle = \frac{C_\Phi}{C_\Psi} 2^{-s(v_i,v_j,\Phi)} 2^{n(v_i,v_j,\Phi)} (-1)^{l} \cdot [ U(L_{i,j})_{kl} \langle v_i v_j | \hat{H}(i,j) | \omega_i\omega_j \rangle$$

$$- U(P_{i,j}L_{i,j})_{kl} \langle v_i v_j | \hat{H}(i,j) | \omega_i\omega_i \rangle ]$$

$$(154)$$

where $s(v_i, v_j, \Phi)$ is defined by

$$s(v_i, v_j, \Phi) = \begin{cases} 0, & v_i = v_j \\ n(v_i, \Phi) + n(v_j, \Phi) - 2, & v_i \neq v_j \end{cases}$$

The justification for this formula is similar to that for orbital difference of one.

**Case 4. Orbital Difference Greater than or Equal to Three**  In this case a line up permutation, $L_{ij}$, does not exist for any $i$ and $j$, no matter which permutation you choose. A a consequence,

$$\langle C_\Phi A\Phi X_j | \hat{H}_{int} | C_\Psi A\Psi X_k | \rangle = \frac{C_\Phi}{C_\Psi} \sum_{i,j} (-1)^{l} 2^{-\delta(\omega_i,\omega_j)} (U(L_{i,j})_{kl} \langle \psi_i\psi_j | \hat{H}(i,j) | \omega_i\omega_j \rangle$$

$$- U(P_{i,j}L_{i,j})_{kl} \langle \psi_i\psi_j | \hat{H}(i,j) | \omega_i\omega_j \rangle)$$

$$= 0 \qquad\qquad (155)$$

where the last step follows as we define $U(L_{ij})$ to be the zero matrix if $L_{ij}$ does not exist.

Note that in each of these cases only one line-up permutation needs to be calculated. We also know how to calculate the special spin basis in which the spin eigenfunctions are also eigenfunctions of the transpositions $P_1$ ., $P_{1,2}$ . $P_{1,n}$ and so on, and elements of the representation matrices $U(P)$.

## 5.4. Results and Discussion

The main advantage of the SACI approach is that it allows the calculation of energy levels and system wavefunctions to arbitrarily high precision, while using less computational

Figure 33. Convergence of the ground state energies with increasing basis size. The harmonic well parameter $\hbar\omega_0 = 10$ meV.

resources than the conventional CI methods. The results presented in this section can then act as benchmarks for other approximate methodologies developed to study more complex systems.

For illustration purpose, we consider again circularly symmetric quantum dots of the form $V(r) = \frac{1}{2}m^*\omega_0^2 r^2$, where $\omega_0$ relates to the steepness of the harmonic well and the external magnetic field is applied perpendicular to the plane of the dot. Figure 33 shows the convergence of the ground state energies for multielectron quantum dots at zero magnetic



Figure 34. The first 50 energy levels for a two electron quantum dot with $\hbar\omega_0 = 10$ meV.

field with increasing basis size. Here we take the harmonic well parameter $\hbar\omega_0$ to be 10 meV and the values of $m^* = 0.067$ and $\varepsilon = 13.1$ to correspond to the experiments using GaAs dots [45]. From the convergence calculation, we are confident in asserting that the ground state energies obtained from our calculations for two, three, four, five, six, and seven electrons are accurate to within 0.03%, 0.024%, 0.045%, 0.05%, 0.1%, 0.15%, and 0.2%, respectively.

From the same calculation, we also obtained accurate information on the excited states of these systems. Figure 34 shows the energy level structure for a two-electron quantum dot. In the triplet states, the two electrons have the same spin and are thus forbidden from being in the same orbital by Pauli's exclusion principle, whereas in the singlet states, the two electrons have different spin, and thus can be in the same orbital. As a result, the basis functions used in the singlet wavefunction expansion include both doubly and singly occupied orbitals, whereas the basis for the triplet system can only include singly occupied orbitals. This means the $1/r$ Coulomb interaction is less significant for the triplet states than the singlet states, and thus there is more degeneracy in the triplet energy levels of the system.

The charge densities for two and three electrons at zero magnetic field are shown in Figures 35 and 36. The most notable aspect of these plots is that the spatial extent of the electron density ranges from 50 to 100 nm. Thus, the electron density at the edge of a dot of diameter 500 nm is negligible, and therefore edge effects should not be important in the modeling. The fact that the electron density is concentrated close to the center of the dot also means the harmonic well approximation should be quite reasonable, as any circularly symmetric potential well can be approximated as a harmonic well near its center.

Figure 37 gives the calculated addition energy, which is defined by Eq. (79), at zero magnetic field at various values of $\hbar\omega_0$. The plot shows large maxima at two and six electrons, where the shells are filled, and a secondary maxima at four electrons because of the shell being half filled. From the graphs, we see that the peaks are more pronounced for larger values of $\hbar\omega_0$, which is a more tightly confined well. This is in agreement with the measurements of Tarucha et al. [45].



**Figure 35.** Charge densities for the lowest 6 states with $l = 0$ for a two-electron quantum dot (top left to bottom right).

**Figure 36.** Charge densities for the lowest 6 states with $L = 0$ for a three electron quantum dot (top left to bottom right).

The wavefunctions of the electrons in quantum dots are spread over a much larger area than those in natural atoms, and the excited state energies are much closer to the ground state energies. This makes the energy levels in the quantum dot more sensitive to changes in magnetic field than in natural atoms. As an example, Figure 38 shows the energies of various states as a function of magnetic field for a quantum dot with two electrons. In this case, the harmonic well parameter $\hbar\omega_0$ is set to be 5.5 meV to compare with the experimental results of Tarucha et al. [45]. The most noticeable feature of this figure is the singlet to triplet transition in the ground state that occurs at around $4T$, which is in agreement with the experimental data.

These calculations were carried out using *Mathematica* on a PC with the Pentium IV 2.4 GHz processor. This work can be readily extended by using a more efficient programming



**Figure 37.** Addition energies obtained using the SACI approach, where the energies for 8 electron dots were also required and evaluated with accuracy within $0.2\%$. The plots correspond to $\hbar\omega_0 = 4$ mev (solid lines), $\hbar\omega_0 = 10$ meV (alternate dashes and dots) and $\hbar\omega_0 = 20$ meV (dashed line).

**Figure 38.** The evolution of the energies of a 2-electron quantum dot with external magnetic field. The harmonic well parameter $\hbar\omega_0 = 5.5$ meV.

language such as Fortran or C++ and more powerful computers to study quantum dots containing more electrons or having more complex geometry.

## 6. CONCLUSIONS

Quantum dots have shown themselves to be tiny laboratories in which fundamental concepts in quantum mechanics can be tested and a new regime of physics can be learnt. This has led to a huge amount of recent activity investigating various aspects of these systems. This chapter focused on the theoretical models and computational schemes developed over the last decade to understand the electronic structure of few-electron quantum dot and other nanosystems, especially in which one can obtain numerically exact solutions.

In particular, we have presented the spin-adapted configuration interaction theory, which takes into account all perceivable interactions and full spin correlation effects, giving a level of accuracy not available in alternative approaches such as the density functional theory or the self-consistent Hartree-Fock method. It also has an advantage over the standard configuration interaction method using Slater determinants in that a much smaller basis is needed since the spin eigenfunctions already have the required symmetry. The results presented in the last section of this chapter provide a benchmark against which one can test other approximate methodologies required to study more complex systems.

For a broader overview of this rapidly developing field, we refer the reader to other review papers and books on related topics: for instance, Jacak et al. (the first comprehensive review on both experimental and theoretical work on quantum dots) [27], Chakraborty (on electronic and optical properties of quantum dots and antidots) [95], Maksym et al. (on molecular aspects of electron correlation in quantum dots) [96], Alhassid (on the statistical theory of quantum dots) [97], Kouwenhoven et al. (on electron transport experiments on few-electron quantum dot devices) [1]; Reimann and Manninen (on shell structures in artificial atoms) [88], and Aleiner et al. (on quantum transport and charge fluctuations) [98], which also contains many references to other works in this field).

The fabrication of quantum dots and other nano-structured systems is currently under way at many research institutes around the world. Much of the research undertaken thus far is concerned with the actual fabrication of these devices and observing fundamental quantum phenomena. Recent efforts have started to focus on using these structures to build usable computing devices such as single-electron devices, quantum bits and logic gates, and laser devices. Many studies have demonstrated nonintuitive behavior and promising potential for novel electronic and laser devices with eccentric functions. As the fabrication of these structures becomes more widely available and their properties understood, they will start to be increasingly important in the laser and electronic industry.

However, to study these quantum phenomena systematically through experiments is difficult and very costly because this would require a new device to be fabricated for each

geometric configuration. In this regard, computer simulations provide a very powerful way of providing detailed and often very accurate information about these systems. Through the use of computer code and an appropriate model description, potential problems and novel electronic devices may be identified and studied. Questions that previously could only be speculated on can now be addressed in great detail. With further development of theoretical models and new computational algorithms, considerable new information will be gained and added to the knowledge base of nano-electronic devices. This field of research is still in its infancy.

## ACKNOWLEDGMENTS

## REFERENCES

1. L. P. Kouwenhoven, D. G. Austing, and S. Tarucha, *Rep. Prog. Phys.* 64, 701 (2001).
2. L. P. Kouwenhoven and C. Marcus, *Physics World* 11, 35 (1998).
3. R. C. Ashoori, *Nature* 379, 413 (1996).
4. K. Ono, D. Austing, Y. Tokura, and S. Tarucha, *Science* 297, 1313 (2002).
5. T. Fujisawa, D. Austing, Y. Tokura, Y. Hirayama, and S. Tarucha, *Nature* 419, 278 (2002).
6. W. Lu, Z. Ji, L. Pfeiffer, K. West, and A. Rimberg, *Nature* 423, 422 (2003).
7. J. Petta, A. Johnson, C. Marcus, M. Hanson, and A. Gossard, *Phys. Rev. Lett.* 93, 186802 (2004).
8. J. Elzerman, R. Hanson, L. van Beveren, B. Witkamp, L. Vandersypen, and L. Kouwenhoven, *Nature* 430, 431 (2004).
9. J. Folk, R. Potok, C. Marcus, and V. Umansky, *Science* 299, 679 (2003).
10. N. Craig, J. Taylor, E. Lester, C. Marcus, M. Hanson, and A. Gossard, *Science* 304, 565 (2004).
11. M. A. Reed, *Scientific American* pp. 98–102 (1993).
12. D. Gammon, *Nature* 405, 899 (2000).
13. E. E. Vdovin, A. Levin, A. Patanè, L. Eaves, P. C. Main, Y. N. Kahnin, Y. V. Dubrovskii, M. Menini, and G. Hill, *Science* 290, 122 (2000).
14. G. Burkard and D. Loss, *Phys. Rev. B* 59, 2070 (1999).
15. D. P. DiVincenzo, D. Bacon, J. Kempe, G. Burkard, and K. B. Whaley, *Nature* 408, 339 (2000).
16. H. A. Engle, P. Recher, and D. Loss, *Solid State Communication* 119, 229 (2001).
17. S. Tarucha, D. G. Austing, T. Honda, R. J. van der Hage, and L. P. Kouwenhoven, *Phys. Rev. Lett.* 77, 3613 (1996).
18. J. B. Wang and C. Hines, *Journal of Computational Electronics* (2003).
19. L. Kouwenhoven, D. Austing, and S. Tarucha, *Reports on Progress is Physics* 64, 701 (2001).
20. T. Ezaki, N. Mori, and C. Hamaguchi, *Phys. Rev. B* 56, 6428 (1997).
21. K. Andres, R. N. Bhatt, P. Goalwin, T. M. Rice, and R. E. Walstedt, *Phys. Rev. B* 24, 244 (1981).
22. T. Pand and S. G. Louie, *Phys. Rev. Lett.* 65, 1635 (1990).
23. B. K. Tanner, "Introduction to the Physics of Electrons in Solids," University Press, Cambridge, 1995.
24. C. Kittel, "Introduction to Solid State Physics," 7th Edn. John Wiley and Sons Inc., 1996.
25. M. Macucci, K. Hess, and G. Jafrate, *J. Appl. Phys.* 77, 3267 (1995).
26. K. J. Thomas, J. T. Nicholls, M. Y. Simmons, M. Pepper, D. R. Mace, and D. A. Ritchie, *Phys. Rev. Lett.* 77, 135 (1996).
27. L. Jacak, P. Hawrylak, and A. Wòjs, "Quantum Dots." Springer, Berlin, 1998.
28. V. Fock, *Zeitschrift für Physik* 47, 446 (1928).
29. C. Darwin, in "Proc. Cambridge Philos. Soc." (1930), Vol. 27, pp. 86–90.
30. L. Landau, *Zeitschrift für Physik* 84, 629 (1930).
31. J. S. Bolemon, *Am. J. Phys.* 40, 1511 (1972).
32. W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, "Numerical Recipies in Fortran 77." 2nd Edn. Cambridge University Press, 1992, pp. 701–775.
33. J. W. Cooley, *Mathematica of Computation* 15 (1961).
34. N. J. Giordano, "Computational Physics." Prentice-Hall, 1997.
35. Y. Saad, "Numerical Methods of Large Eigenvalue Problems." Halsted Press, 1992.
36. R. B. Lehoucq and D. C. Sorensen, *SIAM J. Matrix Analysis and Applications* 17, 789 (1996).
37. R. B. Lehoucq, D. C. Sorensen, and Y. Yang, see www.camm.rice.edu/software/ARPACK.
38. B. Bransden and C. Joachain, "Physics of Atomsn and Molecules." Longman Scientific and Technical, 1983, p. 120.
39. W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, "Numerical Recipies in Fortran 77." 2nd Edn. Cambridge University Press, 1992, pp. 848–852.
40. R. Kosloff, *J. Phys. Chem.* 92, 3554 (1988).

*41.* S. A. Teukoslky. W. T. Vetterling, and B. P. Flannery. "Numereal Recipes in Fortran 90." Cambridge University Press, 1996.

*42.* McEuen, *Phys. Rev. Lett.* 66, 1926 (1991).

*43.* McEuen, *Phys. Rev. B* 45, 11419 (1992).

*44.* R. C. Ashoori. *Phys. Rev. Lett.* 71, 613 (1993).

*45.* S. Tarucha, T. Honda, D. G. Austing. Y. Tokura, K. Muraki, T. H. Oosterkamp. J. W. Janssen, and L. P. Kouwenhoven. *Physica E* 3, 112 (1998).

*46.* D. R. Hartree, *Proc. Cambridge Philos. Soc.* 24, 89 (1928).

*47.* V. Fock. *Zeitschrift für Physik* 61, 126 (1930).

*48.* B. Bransden and C. Joachain, "Physics of Atomsn and Molecules." Longman Scientific and Technical, 1983, pp. 320–327.

*49.* A. Kumar, S. E. Laux, and F. Stern, *Phys. Rev. B* 42, 5166 (1990).

*50.* D. Pfannkuche, V. Gudmundsson, and P. Maksym. *Phys. Rev. B* 47, 2244 (1993).

*51.* J. J. Palacios, L. Martin-Moreno, G. Chiappe, E. Louis, and C. Tejedor. *Phys. Rev. B* 50, 5760 (1994).

*52.* M. Fujito, A. Natori, and H. Yasunaga. *Phys. Rev. B* 53, 9952 (1996).

*53.* H. M. Müler and S. E. Koonin. *Phys. Rev. B* 54, 14532 (1996).

*54.* S. Bednarek, B. Szafran, and J. Adamowski. *Phys. Rev. B* 59, 13036 (1999).

*55.* C. Yannouleas and U. Landman. *Phys. Rev. Lett.* 82, 5325 (1999).

*56.* B. Reusch, W. Häusler, and H. Grabert. *Phys. Rev. B* 63, 113313 (2001).

*57.* S. Bednarek, T. Chwiej, J. Adamowski, and B. Szafran. *Phys. Rev. B* 67, 205316 (2003).

*58.* S. A. McCarthy, J. B. Wang, and P. C. Abbott. *Computer Physics Communications* 141, 175 (2001).

*59.* J. H. Oaknin, J. J. Palacios, L. Brey, and C. Tejedor. *Phys. Rev. B* 49 (1994).

*60.* D. Bielińska-Wąż, J. Karwowski, and G. H. F. Diercksen. *J. Phys. B* 34, 1987 (2001).

*61.* E. P. Wigner. *Phys. Rev.* 46, 1002 (1934).

*62.* M. Koskinen, M. Manninen, and S. M. Reimann. *Phys. Rev. Lett.* 79, 1389 (1997).

*63.* S. Reimann, M. Koskinen, and M. Manninen. *Phys. Rev. B* 62, 8108 (2000).

*64.* B. Szafran, J. Adamowski, and S. Bednarek. *Phys. Rev. B* 61, 1971 (2000).

*65.* S. Bednarek, B. Szafran, and J. Adamowski. *Phys. Rev. B* 64, 195303 (2001).

*66.* S. Bednarek, B. Szafran, and J. Adamowski. *Phys. Rev. B* 61, 4461 (2000).

*67.* W. Kohn and L. J. Sham. *Phys. Rev.* 140, A1133 (1965).

*68.* P. Hohenberg and W. Kohn. *Phys. Rev.* 136, B864 (1964).

*69.* W. Kohn, *Rev. Mod. Phys.* 71, 1253 (1999).

*70.* J. A. Pople, *Rev. Mod. Phys.* 71, 1267 (1999).

*71.* B. Tanatar and D. M. Ceperley. *Phys. Rev. B* 39, 5005 (1989).

*72.* U. von Barth and L. Hedin. *J. Phys. C* 5, 1629 (1972).

*73.* J. P. Perdew and A. Zunger, *Phys. Rev.* 23, 5048 (1981).

*74.* M. Macucci, K. Hess, and G. J. Iafrate. *Phys. Rev. B* 48, 17354 (1993).

*75.* M. Macucci, K. Hess, and G. J. Iafrate. *J. Appl. Phys.* 77, 3267 (1995).

*76.* M. Macucci, K. Hess, and G. J. Iafrate. *Phys. Rev. B* 55, 4879 (1997).

*77.* I. H. Lee, V. Rao, R. M. Martin, and J. P. Leburton. *Phys Rev. B* 57, 9035 (1998).

*78.* J. P. Perdew, K. Burke, and M. Ernzerhof. *Phys. Rev. Lett.* 77, 3865 (1996).

*79.* I.-H. Lee and R. M. Martin. *Phys. Rev. B* 56, 7197 (1997).

*80.* K. Hirose and N. S. Wingreen. *Phys. Rev. B* 59, 4604 (1999).

*81.* A. Wensauer, O. Steffens, M. Suhrke, and U. Rössler, *Phys. Rev. B* 62, 2605 (2000).

*82.* R. McWeeny and B. T. Sutcliffe, "Methods of Molecular Quantum Mechanics." Academic Press, London, 1980, p. 386.

*83.* G. W. Bryant, *Phys. Rev. Lett.* 59, 1140 (1987).

*84.* P. A. Maksym and T. Chakraborty. *Phys. Rev. Lett.* 65, 108 (1990).

*85.* T. Ezaki, Y. Sugimoto, N. Mori, and C. Hamaguchi. *Semicond. Sci. Technol.* 13, A1 (1998).

*86.* M. Eto, *J. Appl. Phys.* 36, 3924 (1997).

*87.* M. Eto, *J. Appl. Phys.* 38, 376 (1999).

*88.* S. Reimann and M. Manninen, *Rev. Mod. Phys.* 74, 1283 (2002).

*89.* D. Pfannkuche, R. R. Gerhardts, P. A. Maksym, and V. Gudmundsson. *Physica B* 189, 6 (1993).

*90.* S. Yang, A. H. MacDonald, and M. D. Johnson. *Phys. Rev. Lett.* 71, 3194 (1993).

*91.* R. Pauncz, "The Construction of Spin Eigenfunctions An Exercise Book." Kluwer Academic/Plenum Publishers, New York, 2000.

*92.* R. Serber, *Phys. Rev.* 45, 461 (1934).

*93.* W. Salmon and K. Ruedenberg. *J. Chem. Phys.* 57, 2776 (1972).

*94.* R. Pauncz, *International Journal of Quantum Chemistry* 12, 369 (1977).

*95.* T. Chakraborty, "Quantum dots: A survey of the properties of artificial atoms." Elsevier, Amsterdam, 1999.

*96.* P. A. Maksym, H. Imamura, G. P. Malion, and H. Aoki. *J. Phys.: Condens. Matter* 12, R299 (2000).

*97.* Y. Alhassid. *Rev. Mod. Phys.* 72, 895 (2000).

*98.* I. Aleiner, P. Brouwer, and L. Glazman, *Physics Reports* 358 309 (2002).

*99.* M. Macucci, K. Hess, and G. Jairate, *Phys. Rev. B* 55, R4879 (1997).

*100.* J. C. Slater. *Phys. Rev.* 38, 1109 (1931).

*101.* I. N. Levine, *Quantum Chemistry*, Prentice-Hall, 1991.

*102.* J. P. Dahl, *Introduction to the Quantum World of Atoms and Molecules*, World Scientific, 2001.

# CHAPTER 11

# Spatiotemporal Dynamics of Quantum-Dot Lasers

## Edeltraud Gehrig, Ortwin Hess

*Advanced Technology Institute, School of Electronics and Physical Sciences, University of Surrey, Guildford, Surrey, United Kingdom*

## CONTENTS

## 1. INTRODUCTION

The quantum-dot laser (QDL) is a complex nonlinear system in which the spatiotemporal dynamics of light fields propagating within a cavity is fundamentally linked with the physical properties of an ensemble of QDs. Due to its characteristic physical properties such as, for example, discrete energy levels, strong carrier localization, and low amplitude phase coupling (frequently expressed $\alpha$ factor) QDLs represent extremely promising laser sources for many applications (e.g., in optical communication networks) [1]. Moreover, recent technological progress in the field of quantum-dot lasers [2] has led to the concept and realization of innovative laser devices suitable for the generation of long-wavelength radiation with good spatial and spectral purity.

For theory and modeling, QDLs represent—due to the high complexity of carrier relaxation and light-matter coupling—many challenges. To unravel the complexity of this interplay, a profound theoretical analysis of this novel laser type is necessary, in particular for interpretation of recent experimental results (e.g., on beam quality, ultrafast time dynamics, and high-speed performance). Even more, clear insight is needed for technological design of quantum dot lasers with improved physical properties. In recent years, the impressive technological progress in the field of quantum-dot lasers has gone hand in hand with a development of various new theories that specifically focus on the physical properties of QDs and QDLs. Although the investigation of the electronic and optical properties of QDs represents a formidable task all by itself, for the QDL one has to set up a theoretical basis that combines the particular quantum optics of an ensemble of QDs with the special material properties of QDs. Our mesoscopic theory discussed here is based on a self-consistent space-dependent integration of material properties into spatiotemporally resolved equations for the light field and carrier dynamics. It thereby bridges theoretical descriptions of microscopic material properties of QDs with macroscopic phenomenological laser theories. This chapter is organized as follows: In Section 2 we derive and explain the QD Maxwell-Bloch equations. Section 3 shows results of our computational modeling of the spatiotemporal dynamics of QD lasers. Section 4 presents simulation results on ultrashort time dynamics, and Section 5 concludes this chapter.

## 2. THEORETICAL DESCRIPTION: QUANTUM-DOT MAXWELL-BLOCH EQUATIONS

In this section, we derive and discuss the quantum-dot Maxwell-Bloch equations (QDMBEs) [3]. The QDMBEs combine the semiconductor Bloch equations describing an ensemble of QDs (QDSBEs) with a suitable form of Maxwell's wave equation. The QDSBEs consider, in particular, a hierarchy of carrier relaxation processes including both intradot scattering as well as the interaction with carriers and phonons in the embedding medium. To represent the dynamic sub-wavelength variations in the light field dynamics, we will derive from Maxwell's equation a wave equation tailored for quantum-dot lasers. The coupled system of partial differential equations then constitutes the quantum dot Maxwell-Bloch equations that model on a mesoscopic basis the spatiotemporal light–matter interplay that characterizes a QDL.

Figures 1 and 2 illustrate the typical geometry of a section of a QDL. In Fig. 1, the active layer contains an ensemble of spatially distributed QDs that are embedded in the quantum well wetting layer (WL). Depending on the epitaxial growth process, the laser may consist of several layers defining vertical "QD stacks" (columns). Light propagates within the active layer in the resonator predominantly along the longitudinal ($z$) direction. This dynamics of the light fields is described by Maxwell's wave equations ("MWEs") considering the spatiotemporal changes of the light fields propagating in the forward ("+") and backward ("−") directions in the laser cavity (see Fig. 2). The layered vertical ($y$) structure is considered via effective material and device parameters. In particular, these are the effective refractive index and the guiding properties of the layer as well as the physical properties of the QD stack (vertically averaged energy levels, damping rates or QD size). The vertically averaged physical properties characterize an "effective" QD. The properties that enter the QDSBEs in a self-consistent way are the energy levels, the initial occupation of the levels (established,

**Figure 1.** Schematic of an idealized active layer of a quantum dot laser: columns of identical pyramidal quantum dots are aligned on a perfect grid.

e.g., via optical or electrical pumping), the dipole matrix elements coupling the individual electron and hole levels, as well as the size of the QDs.

Via the polarization (dipole density) of the active QD medium, the light fields are locally coupled to the dynamics of the carriers and to the interlevel dipole dynamics (described on the basis of the QDSBEs). In particular, their spatial and spectral characteristics are fully taken into account and include, for example, the localization of the dots in the medium, fluctuations in size and shape of the QDs, the spatially dependent light field propagation and diffraction, as well as spatially dependent scattering processes and carrier transport.

The time-dependent calculation of the carrier distributions and the light field dynamics allow an explicit consideration of the individual timescales of the various interaction processes. The relevant timescales range from the femtosecond regime (for the fast carrier scattering processes) up to the picosecond and nanosecond regime (for the dynamics of the propagating light fields and of the spatial carrier density).

## 2.1. Carrier Dynamics within a Quantum Dot

Starting from the single particle density matrices for the electrons, $n^e = \langle c^+ c \rangle$, and holes, $n^h = \langle d^+ d \rangle$, and for the interlevel polarization, $p = \langle d^+ c \rangle$, where $c$ and $d$ are the local annihilation operators for electrons and holes, respectively, one can derive semiconductor Bloch equations specifically for quantum dots. The resulting quantum dot semiconductor Bloch equations ("QDSBEs") mesoscopically describe the dynamic changes of the electron and hole distributions inside the dot (for each energy level) and the dynamics of the interlevel dipoles (for each combination of electron and hole energy levels). If one considers an ensemble of quantum dots as the active medium in a quantum dot laser, additional terms and effects are of relevance and have to be included in the model. These are



**Figure 2.** Schematic of the quantum dot laser model geometry. The counter-propagating light fields ($E^\pm$) spatio-temporally couple with carriers in the ensemble of quantum dots. Characteristic fluctuations in size and location of the quantum dots are effectively represented on a numerical grid with equally spaced grid points in the lateral ($x$) and propagation ($z$) direction (details see text).

contributions describing the electrical injection of carriers (pumping) $\Lambda^{c,h}$ (including Pauli-blocking), induced recombination (with generation rate $g^{c,h}$), spontaneous recombination of the carriers ($\Gamma_{sp}$), carrier–carrier and carrier–phonon scattering for the intradot relaxation $(\partial n^{c,h}/\partial t|_{QD}^{c-ph})$ and the interaction with the wetting layer $(\partial n^{c,h}/\partial t|_{QD-WL})$. The dynamics of the occupation of electrons (e, level index "i") and holes (h, level index "j"), $n^{c,h}$, and the dynamics of the interlevel polarizations $p^{\pm}$ (coupled to the forward (+) and backward (–) propagating optical fields) within a QD are then governed by the equations of motion

$$\frac{\partial n^c(i)}{\partial t} = \Lambda^c(i)\,[D^c(i) - n^c(i)] + g^c(i) - \gamma_{nr}n^c(i)$$

$$- \sum_j \Gamma_{sp}n^c(i)\cdot n^h(j) + \left.\frac{\partial n^c(i)}{\partial t}\right|_{QD}^{c-ph} + \left.\frac{\partial n^c(i)}{\partial t}\right|_{QD\ WL}$$

$$\frac{\partial n^h(j)}{\partial t} = \Lambda^h(j)\,[D^h(j) - n^h(j)] + g^h(j) - \gamma_{nr}n^c(i)$$

$$- \sum_j \Gamma_{sp}n^h(j)\cdot n^c(i) + \left.\frac{\partial n^c(i)}{\partial t}\right|_{QD}^{c\ ph} + \left.\frac{\partial n^c(i)}{\partial t}\right|_{QD-WL} \tag{1}$$

$$\frac{\partial p^{\pm}(j,i)}{\partial t} = -[i\bar{\omega}(j,i) + \gamma_p]p^{\pm}(j,i) - \tfrac{i}{\hbar}\left[n^c(i) + n^h(j)\right]\mathcal{U}^{\pm}$$

$$- \frac{i}{\hbar}\delta\mathcal{U}_{sl}^{\pm} + F_p q^p + \left.\frac{\partial p^{\mp}(j,i)}{\partial t}\right|_{QD}^{p-ph}$$

where $\gamma_{nr}$ represents the rate due to nonradiative recombination, and $\gamma_p$ denotes the dephasing rate of the interlevel dipole. The pump term

$$\Lambda^c(l) = \Gamma_{QDS}\frac{I\eta}{eh}\frac{n_{eq}^c(l)}{\sum_l n_{eq}^c(l)\,[D^c(l) - n^c(l)]} \tag{2}$$

mesoscopically represents the carrier injection and includes the pump-blocking effect (c = e, h and $l = i, j$ for electrons and holes, respectively). It depends on the absolute injection current, $I$, pump efficiency $\eta$, and the thickness of the active area, $h$. $D^c(l)$ denotes the degeneracy of an end energy level (i.e., the maximum occupation with carriers). $\Gamma_{QDS}$ describes the reduction of the pump efficiency resulting from the vertically arranged QDs, that is, the "spatial overlap" between carrier injection and a vertical stack of QDs in the medium.

The generation rates given by

$$g^c(i) = \text{Re}\left(\frac{i}{\hbar}\sum_j\{[\,\mathcal{U}^+p^{+*}(j,i) + \mathcal{U}^-p^{-*}(j,i)] - [\,\mathcal{U}^{+*}p^+(j,i) + \mathcal{U}^{-*}p^-(j,i)]\}\right)$$

$$g^h(j) = \text{Re}\left(\frac{i}{\hbar}\sum_i\{[\,\mathcal{U}^+p^{+*}(j,i) + \mathcal{U}^-p^{-*}(j,i)] - [\,\mathcal{U}^{+*}p^+(j,i) + \mathcal{U}^{-*}p^-(j,i)]\}\right) \tag{3}$$

depend on the interlevel polarization $p$ and on the optical field contributions of spontaneous and induced emission constituting the local field $\mathcal{U}^{\pm}$. The Langevin noise term $F_p q^p$ describes dipole fluctuations [4] with amplitude $F_p = (\Gamma\sqrt{2\hbar\epsilon_r})/(n_r^2 L\sqrt{\epsilon_0\omega_{il}})$. The local fields $\mathcal{U}^{\pm} = d(j,i)E^{\pm} + \delta\mathcal{U}^{\pm}$ are composed of the optical light field contributions $E^{\pm}$ as well as those induced by Coulomb screening in each quantum dot and by the Coulomb interactions between the carriers in the QD and the carriers in the wetting layer, $\delta\mathcal{U}$. $d(j,i)$ is the interlevel dipole matrix element. The interlevel polarization depends via $\bar{\omega}(j,i) = \hbar^{-1}(\varepsilon^c + \varepsilon^h) - \omega$ ($\omega$ is the frequency of the propagating light fields) on the carrier energies $\varepsilon^{c,h}$ that are given by

$$\varepsilon^c(l) = \epsilon^c(l) + \delta\varepsilon^c(l) \tag{4}$$

with the unperturbed level energies $\epsilon^{c,h}$ (i.e., neglecting the carrier dynamics). The characteristic level energies $\epsilon^{c,h}$ of the unperturbed QD are microscopically calculated [5] and are self-consistently included in the theory. The Coulomb-induced screening that leads to a renormalization of these energy levels and also results in additional local field contributions strongly depends on the specific QD design (size, shape). These respective corrections have been determined in detailed calculations (e.g., [6–8]) and are represented in the QDSBEs (1) in the form of spatially dependent energies ($\hbar\omega$ and local field contributions $\delta H^{-}$).

## 2.2. Carrier Relaxation Dynamics

In each quantum dot, the relaxation of the electrons and holes is determined by a variety of physical mechanisms. These are the intradot relaxation ($\partial n^{c,h}(l)/\partial t|_{QD}$) via acoustic and optical phonons or multiphonon processes as well as scattering ($\partial n^{c,h}(l)/\partial t|_{QD-WL}$) between the carriers in the QD and the carriers and phonons of the wetting layer. The physical properties of an individual QD (size, shape, energy levels) and the phonon distribution thereby determine the relevance of the various relaxation processes. Here we will use dynamic scattering rates for carrier–phonon relaxation processes on the basis of microscopically determined matrix elements for the respective interaction. The elastic scattering between the QD carriers and the carriers of the wetting layer will be considered on the basis of perturbation theory.

### 2.2.1. Intradot Relaxation

The scattering rates for carrier–phonon intradot relaxation generally include emission ("−") and absorption ("+") of longitudinal acoustical (LA) phonons, longitudinal optical (LO) phonons, and, in particular, multiphonon processes (± 2 LO, ± 2 LA, and ± LO ± LA). They are determined on the basis of microscopic calculations, allowing a self-consistent mesoscopic inclusion of all scattering processes that are relevant in QD lasers via a dependence of the rates on the spatially and temporally varying carrier and phonon distributions within the QD and the surrounding layers. Starting from the quantum kinetic equations of motion of the single particle density matrices with respect to the carrier–phonon Hamiltonian [9], factorizing the intraband and interlevel matrices into single-particle density matrices and using the Markov approximation (i.e., assuming slowly varying distributions), one obtains after adiabatically eliminating the dynamics of the density matrices the following equations:

$$\frac{\partial n^c(i)}{\partial t}\Big|_{QD}^{c\ ph} = \sum_{i_1 > i} 2\frac{|g^c|^2}{\hbar^2} J(i_1,i)\{(n_q+1)n^c(i_1)[D^c(i)-n^c(i)] - n_q n^c(i)[D^c(i_1)-n^c(i_1)]\}$$

$$-\sum_{i_1 < i} 2\frac{|g^c|^2}{\hbar^2} J(i,i_1)\{(n_q+1)n^c(i)[D^c(i_1)-n^c(i_1)] - n_q n^c(i_1)[D^c(i)-n^c(i)]\}$$

$$-\sum_{j}\Big\{\sum_{j_1 < j}\Big[2\frac{g^h g^{c*}}{\hbar^2} J(j,j_1)p(i,j_1)p^*(i,j) - 2\frac{g^{h*} g^c}{\hbar^2} J(j,j_1)p^*(i,j_1)p(i,j)\Big]$$

$$-\sum_{j_1 > j}\Big[2\frac{g^h g^{c*}}{\hbar^2} J(j_1,j)p(i,j_1)p^*(i,j) - 2\frac{g^{h*} g^c}{\hbar^2} J(j_1,j)p^*(i,j_1)p(i,j)\Big]\Big\}$$

(5)

$$\frac{\partial n^h(j)}{\partial t}\Big|_{QD}^{h\ ph} = \sum_{j_1 > j} 2\frac{|g^h|^2}{\hbar^2} J(j_1,j)\{(n_q+1)n^h(j_1)[D^h(j)-n^h(j)] - n_q n^h(j)[D^h(j_1)-n^h(j_1)]\}$$

$$-\sum_{j_1 < j} 2\frac{|g^h|^2}{\hbar^2} J(j,j_1)\{(n_q+1)n^h(j)[D^h(j_1)-n^h(j_1)] - n_q n^h(j_1)[D^h(j)-n^h(j)]\}$$

$$-\sum_{i}\Big\{\sum_{i_1 < i}\Big[2\frac{g^c g^{h*}}{\hbar^2} J(i,i_1)p(i_1,j)p^*(i,j) - 2\frac{g^{c*} g^h}{\hbar^2} J(i,i_1)p^*(i_1,j)p(i,j)\Big]$$

$$-\sum_{i_1 > i}\Big[2\frac{g^c g^{h*}}{\hbar^2} J(i_1,i)p(i_1,j)p^*(i,j) - 2\frac{g^{c*} g^h}{\hbar^2} J(i_1,i)p^*(i_1,j)p(i,j)\Big]\Big\}$$

Thereby, the phonon distributions have been approximated by their quasi-equilibrium distribution given by the respective Bose statistics, $n_q = 1/\{\exp[\hbar\omega_q/(kT)]-1\}$ with phonon frequency $\omega_q$. $g^{c-h}$ is the coupling constant of the respective carrier–phonon interaction [10]. The function $J$ describes the dependence of the carrier–phonon interaction on the contributing QD level energies and the energy of the respective phonon. It also contains the damping resulting from higher order contributions. For carriers interacting with an optical (LO) phonon, $J$ can be expressed by a Lorentzian line shape

$$J(l_1,l_2) = \frac{\tau_{LO}^{-1}}{\hbar^{-2}(\mathcal{E}_{l_1} - \mathcal{E}_{l_2} - \mathcal{E}_{LO})^2 + \tau_{LO}^{-2}} \tag{6}$$

with $l = i$ for electrons and $l = j$ for holes. The lifetime $\tau_{LO}$ with

$$\tau_{LO}^{-1} = \frac{2\pi}{\hbar^2}|V^{Anh}|^2 \sum_q \delta(2\omega_{LA}(q) - \omega_{LO})[2n_q + 1] \tag{7}$$

includes the decay of optical phonons into two acoustical phonons via the anharmonic interaction potential $V^{Anh}$. The lifetime of acoustical phonons is usually much longer than the lifetime of optical phonons. As a result, the function $J$ can be approximated by a $\delta$-function, $J(l_1,l_2) = \delta(\mathcal{E}_{l_1} - \mathcal{E}_{l_2} - \mathcal{E}_{ph})$, in the case of direct interaction of carriers with acoustical phonons.

In addition to the emission and absorption of one single phonon, the influence of multiphonon relaxation has to be considered. Among all multiphonon processes, the most relevant ones are the emission/absorption of an LO phonon accompanied by the absorption/emission of an acoustical phonon ($\pm$ LO $\pm$ LA) and the emission/absorption of two acoustical phonons ($\pm$ LA $\pm$ LA). The respective relaxation terms derived in analogy to the above read

$$\left.\frac{\partial n^c}{\partial t}\right|_{c\text{-}2\text{ph}}(l) = 2\sum_{l_1, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l_1,l_2) + g_{q_2}^c g_{q_1}^c J_{q_2}(l,l_2)\right|^2$$

$$\times \left\{(n_{q_2} n_{q_1} + 1)n^c(l_1)[D^c(l) - n^c(l)] - (n_{q_2} + 1)n_{q_1} n^c(l)[D^c(l_1) - n^c(l_1)]\right\}$$

$$+ 2\sum_{l_1, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l_2,l) + g_{q_1}^c g_{q_2}^c J_{q_2}(l_1,l_2)\right|^2$$

$$\times \left\{(n_{q_2} + 1)(n_{q_1} + 1)n^c(l_1)[D^c(l) - n^c(l)] - n_{q_2} n_{q_1} n^c(l)[D^c(l_1) - n^c(l_1)]\right\}$$

$$+ 2\sum_{l_1, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l_2,l) + g_{q_1}^c g_{q_2}^c J_{q_2}(l_2,l_1)\right|^2$$

$$\times \left\{n_{q_2}(n_{q_1} + 1)n^c(l_1)[D^c(l) - n^c(l)] - (n_{q_2} + 1)n_{q_1} n^c(l)[D^c(l_1) - n^c(l_1)]\right\}$$

$$+ 2\sum_{l, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l,l_2) + g_{q_1}^c g_{q_2}^c J_{q_2}(l_1,l_2)\right|^2$$

$$\times \left\{(n_{q_2} + 1)n_{q_1} n^c(l_1)[D^c(l) - n^c(l)] - (n_{q_2} + 1)n_{q_1} n^c(l)[D^c(l_1) - n^c(l_1)]\right\}$$

$$+ 2\sum_{l, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l,l_2) + g_{q_1}^c g_{q_2}^c J_{q_2}(l_2,l_1)\right|^2$$

$$\times \left\{n_{q_1} n_{q_2} n^c(l_1)[D^c(l) - n^c(l)] - (n_{q_1} + 1)(n_{q_2} + 1)n^c(l)[D^c(l_1) - n^c(l_1)]\right\}$$

$$+ 2\sum_{l, l_1} \sum_{l, l_2} \frac{1}{\hbar}\left|g_{q_1}^c g_{q_2}^c J_{q_1}(l_2,l_1) + g_{q_1}^c g_{q_2}^c J_{q_2}(l_2,l)\right|^2$$

$$\times \left\{(n_{q_2} + 1)n_{q_1} n^c(l_1)[D^c(l) - n^c(l)] - n_{q_2}(n_{q_1} + 1)n^c(l)[D^c(l_1) - n^c(l_1)]\right\} \tag{8}$$

where $J_{q_1}$ and $J_{q_2}$ are line shapes that depend on the phonon energy and damping rate. $g_{q_1}^c$, $g_{q_2}^c$ are the coupling constants of the respective carrier–phonon interaction.

### 2.2.2. Scattering Processes Between Quantum Dots and the Wetting Layer

In addition to the intradot relaxation, the dynamics of the quantum dot laser depends on the carrier–carrier and carrier–phonon scattering processes that occur between the QDs and the wetting layer in which they are embedded. Those are both inelastic emission and absorption of phonons as well as elastic collision processes. For the inelastic scattering processes, we will consider the inelastic Coulomb interaction between QD carriers and the 2D carrier plasma of the wetting layer via Auger recombination, the ionization of a QD via excitation of carriers by absorption of a phonon as well as carrier capture from the wetting layer in a (up to then unoccupied) state of the QD by emission of a phonon, that is,

$$\frac{\partial n^c(l)}{\partial t}\bigg|_{QD-WL} = \frac{\partial n^c(l)}{\partial t}\bigg|_{QD-WL}^{Aug} + \frac{\partial n^c(l)}{\partial t}\bigg|_{QD-WL}^{c-ph} \tag{9}$$

The relaxation rates

$$\frac{\partial n^c}{\partial t}\bigg|_{QD-WL}^{c-ph}(l) = \Big[ f_{QD-WL}^c N_{WL}^c \big[ D^c(l) - n^c(l) \big] (n_q + 1) - n^c n_q \big[ \Delta^c(l) \big] \Big] J(l,\infty) \tag{10}$$

describe the interactions between the discrete energy levels of the QDs and the states of the surrounding quantum well layer via emission and absorption of optical and acoustical phonons, where the level description $\infty$ in $J(l,\infty)$ refers to the respective energy of the valence and conduction bands of the wetting layer. In (10), the scaling factor $f_{QD-WL}^c$ represents the fraction of wetting layer states to which a single (effective) QD couples. It is determined by the dot density and by the epitaxial structure defining e.g. the potential barrier between the QD and the surrounding layers. $n_q[\Delta^c(l)]$ denotes the phonon number available at energy values higher than the energy given by the potential step between the respective QD level and the wetting layer. The Auger carrier capture kinetics may be attributed to the following processes: (1) A QD electron or hole in the wetting layer collides with a 2D electron and is captured by the QD. The final state of the second 2D electron is then a wetting layer state of higher energy. (2) A 2D hole is captured via Coulomb scattering with a QD electron by the dot while the electron is excited into a wetting layer state. We will represent the two processes by the rates [11]

$$\frac{\partial n^c}{\partial t}\bigg|_{QD-WL}^{Aug}(i) = -B_{he} N_{WL}^h \sum_i n^c(i)[D^h(j) - n^h(j)] + C_{ee}(N_{WL}^c)^2[D^c(i) - n^c(i)]$$
$$+ C_{eh} N_{WL}^c N_{WL}^h [D^c(i) - n^c(i)]$$
$$\frac{\partial n^h}{\partial t}\bigg|_{QD-WL}^{Aug}(j) = B_{he} N_{WL}^h \sum_i n^c(i)[D^h(j) - n^h(j)] + C_{hh}(N_{WL}^h)^2[D^h(j) - n^h(j)]$$
$$+ C_{he} N_{WL}^h N_{WL}^c [D^h(j) - n^h(j)]$$

(11)

In the wetting layer, the carriers are not as strongly localized as in the quantum dot islands and may therefore diffuse within the layer. With $N_{WL}^{c,h}$ denoting the local density of electrons (c = e) and holes (c = h), the dynamics of wetting layer carriers is represented by the diffusion equation

$$\frac{\partial N_{WL}^c}{\partial t} = \Delta D_l(\Delta N_{WL}^c) + \frac{J}{ed} + \frac{\partial N^c}{\partial t}\bigg|_{QD-WL} - \gamma_{sp} N_{WL}^c N_{WL}^h - \gamma_{WL}^{nr} N_{WL}^c \tag{12}$$

with a pump term describing carrier injection and a rate for nonradiative emission processes. In (12), the change in carrier density due to Auger relaxation can be expressed as [11]

$$\frac{\partial N^c}{\partial t}\Big|_{QD-WL} = \sum_{i,j} B_{he} N_{WL}^h n^c(i)[D^h(j)-n^h(j)]n_{QD} - \sum_i C_{ee}(N_{WL}^c)^2[D^c(i)-n^c(i)]n_{QD}$$

$$- \sum_i C_{ch} N_{WL}^c N_{WL}^h [D^c(i)-n^c(i)]n_{QD}$$

$$\frac{\partial N^h}{\partial t}\Big|_{QD-WL} = -\sum_{i,j} B_{he} N_{WL}^h n^c(i)[D^h(j)-n^h(j)]n_{QD} - \sum_i C_{hh}(N_{WL}^h)^2[D^h(j)-n^h(j)]n_{QD}$$ 

$$- \sum_i C_{he} N_{WL}^h N_{WL}^c [D^h(j)-n^h(j)]n_{QD}$$

(13)

In (11) and (13), $B_{he}$, $C_{ee}$, $C_{ch}$, and $C_{he}$ are the respective Auger capture coefficients, which we take from the detailed calculation in [11]. $n_{QD}$ is the dot density.

The elastic Coulomb scattering processes between the QDs and the wetting layer are treated on the basis of perturbation theory [12, 13]. Elastic collisions do not change the occupation of the levels. However, they may lead to significant changes in electronic energies and damping that result in a spectral shift and spectral broadening represented by an energy correction term $\delta^{\prime c,h}$ and the dipole damping $\tau_p$. These spatiotemporally varying quantities are self-consistently included in the ODSBEs. They lead to spatially dependent line shapes and frequency differences $\bar{\omega}$ between the frequency of the propagating light field and the eigen frequencies of the spatially localized QDs in the laser structure. The shift of the emission frequency and the carrier damping rate resulting from the elastic Coulomb scattering between QD and the surrounding layer are

$$\Delta\omega_{QD-WL}^{c,h} = \sqrt{\frac{2kT}{m^{c,h}}}\sigma_\omega^{c,h} N_{WL}^{c,h}$$

$$\Delta\gamma_{QD-WL}^{c,h} = \sqrt{\frac{2kT}{m^{c,h}}}\sigma_\gamma^{c,h} N_{WL}^{c,h}$$

(14)

where $\sigma_\gamma^{c,h}$ and $\sigma_\omega^{c,h}$ denote the intersection areas of the scattering processes given by [13]

$$\sigma_\gamma^{c,h} = \int_0^\infty 2db_{QD}\left\{1-\cos\left[\int_{-\infty}^\infty dt\Delta\omega_{QD}^{c,h}(t)\right]\right\}$$

$$\sigma_\omega^{c,h} = \int_0^\infty 2db_{QD}\sin\left[\int_{-\infty}^\infty dt\Delta\omega_{QD}^{c,h}(t)\right]$$

(15)

with

$$\Delta\omega_{QD}^{c,h} = -\frac{C_3^{c,h}}{r_{QD}^3} - \frac{C_4^{c,h}}{r_{QD}^4}$$

$$C_3^{c,h} = \frac{\pm e^2}{4\pi\varepsilon_0 n_{eff}^2\hbar}\frac{\beta h_{QD}^2}{2}$$

$$C_4^{c,h} = \frac{e^2}{16\pi^2\varepsilon_0^2 n_{eff}^4\hbar}\frac{k_1}{\hbar}$$

(16)

The sign $\pm$ in $C_3^{c,h}$ refers to the situation where the carriers in the QD and in the wetting layer that participate in the collision process have equal ($+$) or different ($-$) sign, respectively. $b_{QD}$ is the spatially dependent collision parameter. $h_{QD}$ and $r_{QD}$ are the height and the radius of the QD. $n_{eff}$ is the effective index of the material. The coefficient $\beta$ is a measure of the linear Stark effect (resulting from expressing the component of the dipole in the direction of the QD axis as $\beta h_{QD}$), and $k_1$ is a coefficient describing the quadratic Stark effect which can be estimated from the dipole moment and the eigen energies [14]. The collision-induced correction $\Delta\omega_{QD-WL}^{c,h}$ is added to the spatially dependent energy renormalization $\delta^{\prime\prime}$

and $\Delta\gamma_{QD}^{c,h}{}_{wI}$ contributes to the damping rates $\gamma^{c,h}$ and $\gamma^p$ in the density matrices. We note that the spatial dependence of $b_{QD}$, $r_{QD}$ and $h_{QD}$ modeled in the form ($X = b_{QD}$, $h_{QD}$, $r_{QD}$)

$$X_{QD} = X_{QD}^{in}(1 + X_{QD}^{fluc})$$  (17)

where $X_{QD}^{in}$ denotes the average value, and $X_{QD}^{fluc}$ is the spatially dependent fluctuation, which represents an arbitrarily distributed ensemble of quantum dots of varying size and shape. $b_{QD}^{fluc}$ considers for example the spatial fluctuation of the collision parameter resulting from a spatial localization of the QDs in the laser. The higher the amplitude of the fluctuations, the higher is, the degree of disorder in the spatial distribution of the effective QDs. $r_{QD}$ and $h_{QD}$ are the average radius and height of the QDs, respectively. This leads to spatially dependent energy corrections and damping rates.

The QDSBEs including the dynamic intradot scattering and the interactions with the wetting layer constitute a fundamental basis for a microscopic analysis of the relevant physical processes such as the influence of many-body interactions, spontaneous recombination, carrier relaxation, and carrier injection. Via the generation rate and the dipole dynamics at each location in the laser structure, the carrier dynamics within the QD (1) and the wetting layer (12) fundamentally linked to the light field dynamics that, in turn, is described by a suitable wave equation.

## 2.3. Optical Field Dynamics: Counterpropagation and Diffraction

Spatiotemporal light field dynamics plays a major role for relevant physical quantities such as the spatiospectral gain and induced index of the system that, in combination with the complex carrier dynamics, determine output quantities of the laser system (i.e., emission wavelength, spectral bandwidth, saturation properties and temporal emission characteristics). A realistic theoretical treatment consequently requires full consideration of the spatially and temporally varying optical fields (associated with spontaneous and induced emission processes) that are mesoscopically coupled to the dynamics of the electrons and holes in the QDs. The QD ensemble represents a strongly inhomogeneous gain medium with spatially distributed QDs with individual material properties (dielectric constant, refractive index, etc.). This spatial inhomogeneity is even more intensified by the spatiotemporal dynamics of the carrier distributions in the QDs and in the wetting layer as well as by the nonlinear interaction of both carrier systems with each other. One may immediately sense that these space and time-dependent variations lead to strong phase changes during the propagation of the light fields in the laser resonator. Consequently, the calculation of the light field dynamics has to include the temporal and spatial changes of the field amplitudes in an appropriate manner.

We start from Maxwell's equations for the optical field $E$ and the polarization $P$ and the material equations and derive the wave equation

$$\frac{1}{\varepsilon_0}\nabla\nabla\cdot P + \nabla^2 E - \frac{1}{c^2}\frac{\partial^2}{\partial t^2}E = \mu_0\frac{\partial^2}{\partial t^2}P$$  (18)

where $\varepsilon_0$ and $\mu_0$ are the permittivity and the permeability in vacuum, respectively, and $c$ is the velocity of light. Insertion of the ansatz

$$E = e^{i\beta z - i\omega t}(E_T + e_z E_z)$$
$$P = e^{i\beta z - i\omega t}(P_T + e_z P_z)$$  (19)

for the optical fields and the polarization leads to

$$\left(-\beta^2 + 2i\beta\frac{\partial}{\partial z} + \frac{\partial^2}{\partial z^2}\right)(E_T + e_z E_z) + \nabla_T^2(E_T + e_z E_z) - \nabla_T\nabla_T - \nabla_T E_T$$

$$-\nabla_T\left(i\beta + \frac{\partial}{\partial z}\right)e_z E_z - \left(i\beta + \frac{\partial}{\partial z}\right)\nabla_T E_T + \frac{1}{c^2}\left(\omega^2 + 2i\omega\frac{\partial}{\partial t} - \frac{\partial^2}{\partial t^2}\right)(E_T + e_z E_z)$$

$$= -\mu_0\left(\omega^2 + 2i\omega\frac{\partial}{\partial t} - \frac{\partial^2}{\partial t^2}\right)(P_T + e_z P_z)$$  (20)

with the propagation constant $\beta$ and frequency $\omega$. With the main field propagation in parallel and antiparallel to the resonator axis in (20), we may neglect the mixed derivatives, $\nabla_T(\partial/\partial z)$, $(\partial/\partial z)\nabla_T$. Similarly, $\nabla_T E_T \approx -i\beta E_z$ such that the deviates $\nabla_T \nabla_T E$ and $\nabla_T i\beta E_z$ can also be safely omitted. Disregarding the second-order derivate $(\partial^2/\partial t^2)P$ of the polarization (in analogy to the microscopic Bloch equations where one implicitly assumes a linear response function), we finally obtain the following effective wave equation for the counter-propagating $(+,-)$ optical fields in a QD laser:

$$\nabla_T^2 E^\pm \pm 2i\beta\frac{\partial}{\partial z}E^\pm + \frac{\partial^2}{\partial z^2}E^\pm + \frac{2i\omega}{c^2}\frac{\partial}{\partial t}E^\pm - \frac{1}{c^2}\frac{\partial^2}{\partial t^2}E^\pm = -\mu_0\omega^2 P^\pm - 2i\omega\frac{\partial}{\partial t}P^\pm + F_E q^E \qquad (21)$$

The Langevin noise term $F_E q^E$ that has been added to (21) has been derived from quantum Maxwell-Bloch equations [4]. It considers spontaneous light field fluctuations that depend via $F_E = (\sqrt{2\hbar\omega_0})/(\sqrt{\epsilon_r\epsilon_0})$ on the specific material parameters and on the emission wavelength of the device. $q^E(r,t)$ obeys the correlation relation

$$\langle q^E(r,t)q^E(r',t')\rangle = \kappa\delta(r-r',t-t') \qquad (22)$$

where $\kappa = 1/(2L)\ln[R_1 R_2]$ corresponds to the damping rate of the resonator. The polarization of the active semiconductor medium

$$P^\pm = V^{-1}\sum_{i,j}d(i,j)p^\pm(i,j) \qquad (23)$$

is the source of the optical fields ($V$ denotes the normalization volume of the crystal).

Our derivation of the "quantum-dot Maxwell-Bloch equations" (QDMBEs) reflects the spirit of describing the (spatiotemporal) dynamics of (spatially inhomogeneous) semiconductor lasers on the basis of Maxwell-Bloch equations [15]. The QDMBEs mesoscopically consider the dynamics of the carrier distributions in the dots and the interlevel dipoles together with the spatiotemporal dynamics of the optical fields (including spontaneous light fields, amplified spontaneous emission (ASE) and induced recombination).

Specific laser configurations of an actual device that is characterized by its geometry, mirror reflectivity, current injection, and so forth, are fundamentally included in our description. They enter the theory in the form of boundary conditions for the dynamically varying optical fields and the carriers, like the pump term of the QDMBEs. The laser cavity induces additional counterpropagation and wave-guiding effects which superimpose the carrier–field dynamics. The resulting complex dynamic spatiospectral interactions between the QDs, the optical fields, and the surrounding layers influence the emission properties (e.g., the temporal behavior of the optical fields, emission spectra). In the following Section 3, we will discuss selective results of numerical simulations that illustrate the interplay of light field and carrier dynamics in quantum dot lasers.

## 3. COUPLED SPATIOTEMPORAL LIGHT FIELD AND INTER-/INTRALEVEL CARRIER DYNAMICS IN QUANTUM-DOT LASERS

In the following, we will present selective results of numerical simulations based on QDMBEs. Specifically, we will consider spontaneous and induced light emission in the active dot medium and analyze the influence of spatial inhomogeneities (in quantum dot parameters such as dot size, level energies, dipole matrix elements) on the spatiotemporal light field and carrier dynamics. For specificity, the QDL structure is assumed to consist of three dot layers (InAs/GaAs [5]) with a dot density of $10^{10}$ cm$^{-2}$. The dots are assumed to be of pyramidal shape with base length 12 nm and with three energy and five hole levels. The length of the laser is 1 mm, its width (of the active zone) 10 $\mu$m. At every location in the medium, the QDSBEs are coupled to a diffusion equation describing the spatial distribution of the carriers in the surrounding layers via dynamic scattering terms. The general representation of physical properties and components of a QDL by the QDMBEs is sketched in Fig. 2.

In the simulation of the spatiotemporal dynamics, the spatial dependence of the carriers in the wetting layer (WL) and the propagating light fields ($E^\pm$) are considered via a numerical grid with equally spaced grid points in the lateral ($x$) and propagation ($z$) directions. The local distribution of QDs is defined by spatial coordinates with respect to this grid. Thereby, the spatial distance between the position of each QD and the center of the respective cell (with length $\Delta z$ and width $\Delta x$) is saved in a spatially dependent variable that is used for the collision rates between the QD and the wetting layer. Each mesh (size $\Delta x \cdot \Delta z$) contains the following information: number of (effective) QDs in the mesh, $N_{QD}$ (note that a "hole" in the spatially distributed QDs, i.e., $N_{QD} = 0$, is also possible), position of the QDs (i.e., their distance from the center of the area ($\Delta x \cdot \Delta z$)), and the individual material properties of the QDs. The specific laser configuration is defined by the size of the medium and the reflectivities of the facets ($R_1$, $R_2$) and enters the theoretical description as boundary conditions for the optical light fields.

A convenient way to visualize the spatiotemporal light dynamics of spontaneous and induced emission processes is in the form of snapshots of the spatial light fields and of carrier distributions. To additionally grasp the complex microscopic carrier relaxation dynamics, we will focus on the dynamics of the level occupations. On route, we will analyze the influence of spatially varying quantum dot properties (e.g., dot size, level energies, dipole matrix elements) on the spatiotemporal light field and carrier dynamics.

## 3.1. Spatiotemporal Light Field Dynamics: Interplay of Spontaneous and Induced Emission

One of the very characteristic properties of a laser is the buildup of coherence in the light field from the initial spontaneous emission. In the quantum dot laser whose active medium is fundamentally characterized by the spatially inhomogeneous ensemble distribution of active quantum dot sources, one would expect that this transition is determined by this very feature together with the dependence on the (electrical or optical) excitation of the system. For low electrical injection current, Fig. 3 shows typical snapshots of the spatial light field distribution in our model quantum dot laser structure. The time interval between the snapshots—showing



Figure 3. Dynamics of the luminescence pattern of spontaneous emission of an idealized quantum dot laser with perfect and uniform dot arrangement and identical parameters for dot size, level energies, and dipole matrix elements. Bright colors indicate high levels of intensity. The time between successive snapshots is 3 ps. Reprinted with permission from [3], E. Gehrig and O. Hess. *Phys. Rev. A* 65, 033804 (2002). ©2002. American Physical Society.

the speckle distribution that is a characteristic of the quantum dot medium—is 3 ps. In this example, we have assumed the very ideal case where the distribution of the dots in the structure is uniform (i.e., the dots are positioned with constant dot-to-dot distance). Furthermore for each dot within the structure, we have used an identical set of parameters for dot size, level energies, and dipole matrix element. The injection of carriers has been chosen such that the occupation of the energy levels of the dots is near transparency. In spite of the "ideal" conditions assumed for the laser structure, spatial fluctuations in light and carrier distribution arise. They are the result of spontaneous light fluctuations, microscopic carrier relaxation dynamics, and the nonlinear coupling between the light fields and the charge carrier plasma. The carrier dynamics within each dot is determined by processes such as carrier injection, spectral hole burning, intradot carrier relaxation via phonon emission, and absorption and carrier–carrier and carrier–phonon interaction with the wetting layers as well as screening. We will later focus on them in more detail. For now we can see that for the light field dynamics the underlying physical processes consist of both coherent (in the case of, e.g., induced recombination) and incoherent contributions (e.g., spontaneous emission, carrier relaxation). Consequently, they vary from dot to dot even when identical dot parameters and an ideal uniformity of the dot distribution in the layers are assumed.

The interplay of incoherent and coherent interactions yields a spatially varying number of electrons and holes in the energy levels of the quantum dots. Together, the spontaneous and induced light emitted by a quantum dot then contributes to the forward and backward propagating light fields and is thus transferred to the neighboring dots leading to complex spatiotemporally varying light–matter interactions. The propagating light fields, on the other hand, experience a spatially dependent modification via the interaction with the quantum dot ensemble. In combination with the diffraction of the light field this leads to a spatially varying light field dynamics (Figs. 3a–3c). The nonlinear and inhomogeneous light–matter interaction and the carrier dynamics affect the spatial charge carrier density at the same time. For the time frame of Fig. 3c, Fig. 4 shows the distribution of electrons as an example. The spatially varying level occupation and the formation of characteristic optical patterns are a direct consequence of spontaneous light fluctuations and scattering. The microscopic intradot scattering of the carriers within the dots via emission and absorption of phonons, the interaction of the "dot carriers" with the carriers and the phonons of the wetting layer, and the nonlinear coupling to the propagating light fields leads to a spatially varying occupation of the dots and subsequently to complex transverse carrier dynamics. It is important to note that the interplay of light with the carriers results in a spatiotemporally varying occupation, although we have assumed the "ideal" case of uniform carrier injection and regular matrix-like positioning of quantum dots that each have identical properties (size, level energies, matrix elements). The spatiotemporal light field dynamics changes if we



Figure 4. Snapshot of the spatial electron distribution corresponding to the luminescence pattern of Fig. 3(c). Reprinted with permission from [3]. E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002). ©2002, American Physical Society.

increase the excitation level (carrier injection) by increasing the respective pump term in the Bloch equations so that the dots are almost completely filled with carriers. In this case, the snapshots (again with time steps of 3 ps) of Figs. 5a–5c show the result of a significant inversion: light amplification by induced recombination occurs in addition to the spontaneous emission processes. The first intensity distribution is taken 100 ps after the initial excitation of the dots. In the longitudinal ($z$) direction, one can observe dynamically varying intensity modulations. These longitudinal structures are typical for the onset of laser oscillations of a device immediately after start-up. They are a measure of the characteristic internal coherence length scales that typically lie in the micrometer regime. In time, the structures lead to intensity spiking and relaxation oscillations in the light emission. In the lateral direction (i.e., parallel to the output facet), the intensity is rather uniform when compared to Fig. 3. This uniformity originates from induced emission processes which now play a major role in the overall behavior of the device: The initial filling of the dots establishes a carrier inversion and thus a high gain. Due to the increased influence of induced emission processes, a spatiotemporal coherence builds up that via the propagating light fields is transferred in both time and spatial dimensions. The coupling of the carriers in the dots with the propagating light fields in combination with the high gain characterizing the dot medium may then lead to a narrow-band stable laser output.

## 3.2. The Spatially Inhomogeneous Quantum-Dot Ensemble: Influence of Disorder

Contrary to the ideal situation assumed so far, slight dot-to-dot variations in size, energy levels, and material parameters exist in real quantum dot laser systems. In addition, the dots are not equally positioned on a grid within the layers. The respective variations in quantum-dot parameters and dot-to-dot distance depends on the material system and the epitaxial growth process of the particular quantum dot system. In the numerical simulations in the following, we will discuss some first steps in the analysis of disorder in quantum dot lasers and its influence on the spatiotemporal light fields and carrier dynamics and on emission spectra. In particular, we will focus on spatial fluctuations in the position of the quantum dots, their size, energy levels, dipole matrix elements, and scattering rates that vary from dot



Figure 5. Spatiotemporal dynamics of stimulated emission in quantum dot lasers pumped above threshold. The time interval between successive snapshots is 3 ps. Reprinted with permission from [3], E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002). ©2002, American Physical Society.

to dot in a real laser device. Thereby, we can systematically change the spatial fluctuation of the individual parameters. In this section, we will focus on the influence of the variation of dot size in the inhomogeneously distributed quantum-dot ensemble that constitutes the active laser material. A consequence of this will be a variation of energy levels and scattering matrix elements.

The dynamics of the light fields in quantum dot lasers is determined by induced and spontaneous recombination processes that, in turn, depend on the spatiospectral carrier distribution in the dot levels. The mutual influence of light and matter is particular strong during the start-up regime of the laser and leads to characteristic oscillations in the time domain. Figure 6 shows the first 10 ns (after start-up) of the dynamics of the optical nearfields (left) and carrier distribution (right) at the output facet of an InGaAs quantum-dot laser. Its width and cavity length are 50 $\mu$m and 500 $\mu$m, respectively, and the density of quantum dots



Figure 6. Dynamics of the optical nearfield emission (left) at the output facet and corresponding carrier density dynamics (right) of a quantum dot laser (InGaAs, width 50 $\mu$m, cavity length 500 $\mu$m, $j = 2.0 j_{tr}$, dot density $10^{11}$ cm$^{-2}$) with spatially varying dot parameters. From top to bottom, the amplitude of the (Gaussian-shaped) fluctuations is 1%, 4%, and 8%.

(self-organized growth) is $10^{11}$ cm$^{-2}$. In the top frames of Fig. 6, the QD parameters deviate only slightly from their average values (1% variance), whereas in the middle and lower frames the variance of the parameter values of the spatially distributed QDs (i.e., their size, dipole matrix elements, energy-levels) is Gaussian with variances of 4% and 8%, respectively. The nearfield intensity distributions show modulations on a picosecond timescale. They originate from dynamic interactions between the light fields and the dot carriers ranging from the femtosecond timescale (in the case of microscopic carrier scattering) up to the picosecond and nanosecond timescales (reflecting the resonator round trip time of the propagating light fields and the slow build-up and decay of the spatial carrier density). In combination, the light diffraction and the spatially dependent interaction of light with the carriers in the dots and in the wetting layer lead to formation of dynamic optical patterns. The timescales of the carrier dynamics thereby are transformed into characteristic interaction lengths like the coherence length via the propagation of the light fields. In combination with the diffraction of the light fields this leads to transverse modulations in the micrometer regime.

Dynamic carrier capture and escape, the complex intradot level dynamics and the diffraction of the propagating light fields lead to a transverse coupling of the field distribution arising at the individual dots. Via the light field propagation, these fluctuations are then transfered to respective longitudinal ($z$) and temporal ($t$) changes. Due to the coupling and interplay of spatial with temporal degrees of freedom, these material inhomogenieties thus affect both the transverse nearfield distribution as well as the dynamics of the light fields. With increasing fluctuation amplitude, a characteristic filament structure evolves. The spatial variation of dot size and level energies directly determines the spatial dependence of the nearfield. Via the term $i\omega + \gamma_p$ of the QD Bloch equations, the spatial fluctuations in the dot properties are transferred to spatiospectral changes in the interlevel polarization. This, in turn, affects the amplitude and phase of the propagating light fields. Thereby, the individual level energies lead to locally varying transition energies and frequencies that contribute to the spatial and spectral properties of the propagating light fields. Via the spatial light field polarization these are passed on to the propagating light fields as dynamic changes in both amplitude and phase. As Fig. 6 shows, an increase in energy fluctuations consequently leads to corresponding fluctuations in the light field dynamics. Via the generation rate, the dynamics of the carrier system depends on the light fields that in turn are spatiotemporally modified by the spatially varying dot parameters. In addition, the dynamics of the carrier relaxation processes (phonon emission or absorption and interaction with carriers and phonons of the wetting layers) depends on the energy differences of the levels involved and thus is also directly affected by the spatially varying dot parameters. As a consequence, dynamic characteristic filament structures evolve in both the light field and carrier distributions.

The spatially varying dot properties not only determine the emission dynamics, they also induce dynamic changes in the real and imaginary parts of the light fields. Our numerical simulation allows us to analyze the influence of the spatially varying dot properties on the spatiospectral emission characteristics. Figures 7a–7c show spatially resolved emission spectra corresponding to Fig. 6 with 1%, 4%, and 8% fluctuations. The vertical axis refers to the lateral position at the output facet of the quantum dot laser; the horizontal axis shows the frequency dependence in a range of 600 GHz. With a cavity length of the device of 500 $\mu$m, this corresponds to approximately 7 longitudinal modes. The spectral width of each longitudinal mode is determined by a large variety of physical effects. The characteristic times of induced and spontaneous recombination define a lower limit for a laser linewidth. However, the real spectral width is significantly broadened. First, in large-area lasers, the transverse degree of freedom leads to a characteristic transverse migration of the light fields that is determined by the dynamic interplay of light diffraction, dynamic self-focusing, as well as carrier scattering and relaxation. As a consequence, a group of transverse modes arises for each longitudinal mode. Second, the transverse dynamics is determined by inhomogeneous broadening resulting, for example, from the spatially varying quantum-dot parameters. They affect both the coherent light–matter coupling (via the carrier dependence of the generation rate) and the incoherent processes such as carrier–carrier and carrier–phonon interactions.

Figure 7 clearly demonstrates that an increase in fluctuation amplitude affects both the spectral as well as the spatial degree of freedom. The spatially varying transition energies

**Figure 7.** Emission spectra of a quantum dot laser with Gaussian fluctuations with amplitude (from top to bottom) 1%, 4% and 8%. The vertical axis refers to the lateral position at the output facet of the quantum dot laser; the horizontal axis shows the frequency dependence in a range of 600 GHz. With a cavity length of the device of 500 μm, this corresponds to approximately 7 longitudinal modes.

lead to significant broadening of the emission spectra. This originates from the variance in transition energy and (indirectly) the intradot and interdot scattering dynamics. The light propagation and diffraction in combination with carrier scattering and relaxation not only lead to a coupling of longitudinal and transverse degrees of freedom but also to a coupling of spectral and transverse dynamics. As a consequence, the spatially resolved emission spectra show inhomogenieties in both the transverse dimension (i.e., over the lateral extension of the quantum dot broad-area device) and the spectrum.

## 3.3. Dynamic Filamentation and Beam Quality of Quantum-Dot Lasers

In the following, we present a comparative study of numerical simulations and experiments on the spatiotemporal dynamics and emission characteristics of quantum-well and quantum-dot lasers of identical structure. The simulations show that, in the quantum-dot laser, the strong localization of carrier inversion and the small amplitude-phase coupling enable a significant improvement of beam quality compared to quantum-well lasers of identical geometry. Near-field profiles and beam quality ($M^2$) parameters calculated on the basis of time-dependent effective Maxwell-Bloch equations into which the physical properties of the active media are included via space-dependent material parameters, effective time constants, and matrix elements are fully confirmed by experimental measurements. Together they indicate that in the quantum-dot laser, the strong localization of carrier inversion and the small amplitude-phase coupling enable a significant improvement of beam quality compared with quantum-well lasers of identical geometry. Figure 8 shows in direct comparison the theoretical (a, b) and experimental (c) [16, 17] results of the nearfield characteristics of quantum-well (Fig. 8, left) and quantum-dot lasers (Fig. 8, right) of identical waveguide design. The width of the lasers is 6 μm. the cavity length was 1.3 mm. Figure 8a shows the calculated nearfield dynamics of a quantum-well (left) and quantum-dot laser (right). In the example, the injection current density was chosen such that the output power was 60 mW. For small output powers, the nearfield of a quantum-well laser is still rather uniform. However, if we operate the laser at higher power levels, a very different behavior can be observed. Physical processes such as carrier diffusion and scattering in combination with

**Figure 8.** Simulated spatiotemporal nearfield dynamics (a) and temporally averaged nearfields (b: theory, c: experiment) at the output facet of a quantum-well (left) and a quantum-dot (right) laser. Both lasers have the same geometry (width 10 µm, cavity length 1 mm). Reprinted with permission from [18], E. Gehrig et al., *Appl. Phys. Lett.* 84, 1650 (2004). © 2004. American Institute of Physics.

light diffraction lead to a complex transverse migration of the light fields (Fig. 8a, left). In contrast, the transverse dynamics of the light fields propagating in the quantum-dot laser (Fig. 8a, right) is still rather uniform (for the same output power). This complex light field dynamics is governed by the (space- and time-dependent) mutual interplay of carriers with spontaneous and induced emission processes. In combination with (counter-) propagation effects and transverse migration of the light fields, it affects and determines the nearfield and beam quality that can be measured and observed in an experimental investigation.

A comparison of the calculated and measured time-averaged nearfields of a quantum-well laser (left) and a quantum-dot laser (right) (Fig. 8b) demonstrates that the increased influence of the transverse degree of freedom leads in the quantum-well laser to the formation of filaments, whereas the quantum-dot laser shows a Gaussian-shaped uniform nearfield distribution. The theoretical results obtained with the Maxwell-Bloch equations are in good agreement with an experimental measurement of the nearfield distributions [17] (Fig. 8c). The side lobes next to the laser ridge that can be seen in Fig. 8c result from current-spreading in the cladding and waveguide layers. The suppressed transverse light field dynamics observed in experiment and simulation clearly demonstrate the promising device performance of quantum-dot lasers compared to large area lasers and laser-amplifiers, which show a strong tendency for filamentation formation [20–22].

By varying the carrier injection current and the laser width in the simulation, we have systematically analyzed the spatiotemporal light field dynamics of quantum-well and quantum-dot lasers. The theoretical and measured values of the beam quality factor $M^2$ as obtained by the spatiotemporal simulation are depicted in Fig. 9 in dependence on stripe width (a) and output power (b) (for the same output power of 20 mW). With increasing stripe width (Fig. 9a), the transverse degree effectively gets more important: Physical processes such as carrier diffusion and diffraction of the light fields lead to characteristic dynamic optical patterns that typically lie in the micrometer regime. In combination with the dynamic phase changes, this results in a deterioration of the beam quality, that is, the $M^2$-parameter of the

**Figure 9.** Calculated and measured [17] beam quality parameter ($M^2$) in dependence of stripe width (a) and output power (b). Reprinted with permission from [16], E. Gehrig et al., *Appl. Phys. Lett.* 84, 1650 (2004). © 2004, American Institute of Physics.

quantum-well and the quantum-dot laser increases with increasing stripe width. However, due to the strong localization of the carriers and the reduced $\alpha$-factor, the $M^2$-values of the quantum-dot laser are always smaller than the respective values of the quantum-well structure. In particular, the quantum-dot laser shows a characteristic threshold near 8 $\mu$m. In this intermediate stripe regime, the quantum-dot laser is still single mode, whereas $M^2 > 2$ for the quantum-well laser.

The dependence of $M^2$ on output power is shown in Fig. 9b (for a stripe width of 6 $\mu$m). In the quantum-well laser, an increase in the injection current density not only increases the output power but simultaneously leads to an increase in the $M^2$-parameter. This is a direct consequence of dynamic carrier diffusion and light diffraction affecting the light fields during their propagation in the laser. In contrast, the quantum-dot laser shows almost no dependence on output power. The dependence of $M^2$ on stripe width and output power can be confirmed by experimental measurements performed on the same devices [17].

Our numerical results clearly demonstrate that quantum-dot lasers have a much better beam quality compared to quantum-well lasers of same geometry. The strong localization of the carriers in the dots in combination with the reduced amplitude phase coupling thus guarantees a good spatial quality.

For an analysis of the spectral properties, Fig. 10 shows calculated (spatially resolved) emission spectra of the quantum-well (a) and the quantum-dot laser (b). In the figure, the vertical axis denotes the lateral coordinate of the laser; the horizontal axis refers to the frequency. Both structures show a set of longitudinal modes that coexist in the laser—according to the Fabry-Perot modes of the cavity. In the quantum-well laser (Fig. 10a), the width of the laser is larger than typical interaction length scales of the laser. As a consequence, a characteristic transverse spatiospectral coupling arises, leading to spectral broadening of the individual modes. In addition, each longitudinal mode is surrounded by a set of transverse modes. In the case of the quantum-dot laser (Fig. 10b), the strong carrier localization and the discrete energy levels lead to reduced carrier diffusion, a reduction of transverse dynamics, and to

quantum well laser                    quantum dot laser



Figure 10. Calculated emission spectrum of a quantum-well (left) and a quantum-dot laser (right). Reprinted with permission from [16], E. Gehrig et al., *Appl. Phys. Lett.* 84, 1650 (2004). © 2004, American Institute of Physics.

low amplitude-phase coupling (alpha-factor). As a consequence, the emission spectrum is of much higher spectral purity than in the situation of the quantum-well laser.

For characteristic sets of material parameters describing the active quantum-well or the quantum-dot media, the QD Maxwell-Bloch approach allows a realistic simulation of the spatiotemporal dynamics of quantum-dot and quantum-well lasers complementing experimental measurements of nearfield profiles and beam quality factors [17]. The self-consistent inclusion of all relevant geometrical parameters and material properties (e.g., refractive index and waveguide structure, dot density, and spatial dot distribution) provides a fundamental description of the underlying physical processes and guarantees a realistic modeling of the laser behavior. In particular, the simulations allow the systematic variation of the individual parameters and properties with respect to their influence on beam quality and power. Our results clearly indicate that the quantum-dot laser is a highly promising laser source for the generation of long-wavelength radiation with improved spatial and spectral purity. Experiment and modeling together may thus significantly contribute to the development of optimized, innovative quantum-dot devices.

# 4. ULTRASHORT TIME DYNAMICS: PULSE PROPAGATION IN A QUANTUM-DOT WAVEGUIDE

A typical means of probing the internal ultrashort time dynamics of a quantum-dot laser is to inject an optical pulse into the laser and analyze the ultrafast dynamics of the output signal. The dynamics of ultrashort pulses propagating in a quantum-dot amplifier is determined by a complex nonlinear coupling and dynamic interplay of light fields and carriers in the spatially inhomogeneous quantum-dot ensemble. Computational modeling shows that in spite of the large complexity, the strong localization of the carrier inversion and the low-amplitude phase coupling may allow the amplification and transmission of ultrashort light pulses with minimum deterioration of the pulse properties (e.g., pulse shape, duration). Simulation results of the nonlinear pulse propagation in quantum-dot optical amplifiers allow visualization and interpretation of fundamental nonlinear processes such as selective depletion and refilling of quantum-dot energy levels, leading to a complex gain and index dynamics that affect the amplitude and phase of a propagating light pulse. Computational modeling thus may lay the foundation for an optimization and tayloring of pulse properties.

## 4.1. Dynamic Shaping and Amplification of Ultrashort Pulses

For the simulation of a typical pulse propagation configuration, we consider a light pulse propagating in a laser waveguide (width of the structure 10 $\mu$m, length 1 mm) filled with an inverted quantum-dot ensemble. The snapshots displayed in Fig. 11 show the intensity (a–c) and the carrier density (d–f) in the active area of a QD laser during the propagation of a light pulse whose frequency corresponds to the transition energy of the QDs. The injection current density has been chosen such that the population of the dots in the layers of the QD waveguide are significantly above transparency. For the dot-to-dot fluctuation, a variance of 5% has been assumed. The time between successive plots is 3 ps. Figures 11a and 11d represent the spatial distributions of the intensity (a) and the carrier density (d) immediately after

Figure 11. Propagation of an ultrashort pulse (full width at half maximum of 500 fs) tuned to resonance of an inverted quantum-dot ensemble: (a–c): snapshots of the light field and (d–f): corresponding snapshots of the carrier density. The time between successive snapshots is 3 ps. Reprinted with permission from [3], E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002). © 2002, American Physical Society.

optical injection. It is important to note that initially the lateral spatial shape of the injected light field is Gaussian-shaped with a width (FWHM) of 6 $\mu$m, and the temporal profile of the pulse is chosen Gaussian as well with a full width at half maximum of 500 fs. Immediately after injection, the light pulse starts to interact with the ensemble of populated QDs. During its propagation through the laser, the pulse locally reduces the population in the dots established by the injection current by induced recombination. With continuing propagation. the light pulse is significantly amplified (Fig. 11b). Due to the nonlinear light–matter interaction between the pulse and the spatially distributed dots, a complex spatiotemporal behavior arises. It is directly reflected in the dynamic spatial structures in both intensity and carrier density. The pulse is laterally structured and temporally distorted via the interaction with the dots (Figs. 11b and 11c). At the same time, spatial hole burning effects can be observed in the carrier distribution Fig. 11e and 11f). The partial refilling of the dots—determined by carrier injection. carrier capture and thermalization via carrier relaxation—defines a finite "response time" of the QD medium. As a consequence. the spatial extension of the hole burnt by the light pulse (Fig. 11f) significantly exceeds the spatial area covered by the optical pulse (Fig. 11e). The microscopic response and relaxation of the excited charge carrier system typically occurs on ultrashort timescales (50 fs ... 500 fs). As a result, we expect the dynamic shaping and amplification of the light pulse to be strongly dependent on the duration of a light pulse propagating in a QD laser amplifier. In order to investigate the influence of pulse duration, we compare the dynamics of pulses with durations of 150 fs and 1.5 ps, respectively. Figure 12 shows snapshots of the propagating light pulses near the output facet of a QD laser amplifier (cavity length 500 $\mu$m). The figures refer to the small signal regime (Figs. 12a and 12b) and to the saturation regime (Figs. 12c and 12d) for a pulse with a duration (full width at half maximum) of 150 fs (Fig. 12a and 12c) and 1.5 ps (Fig. 12b and 12d). respectively.

Figure 12. Snapshot of the intensity in the active area of an optically injected quantum dot laser (width of active laser stripe: 10 μm, cavity length 500 μm) with a duration of (a, c) 150 fs and (b, d) 1.5 ps of the injected pulse. Parts (a) and (b) refer to the small signal regime, parts (c) and (d) to the saturation regime.

The characteristic spatial and temporal distortions the light pulse experiences during propagation in an inverted dot medium strongly depend on the duration of the light pulse: If the QD laser amplifier is operated in the small signal regime (e.g., small values of the injected light pulses), a light pulse more or less retains its shape during its propagation in an inverted QD laser amplifier, independent of the duration of the injected pulse. However, if we increase the input power level (approaching the saturation regime of the QD waveguide), the dynamic pulse shaping shows a strong dependence on the pulse duration. Figure 12 displays the profile of light pulses with a duration of 150 fs (a, c) and 1.5 ps (b, d) after the propagation in a QD laser amplifier (width of active laser stripe: 10 μm, cavity length 500 μm). Please note that only a part around the pulse and not the entire active area is plotted. The leading part of the pulse with duration of 1.5 ps (Fig. 12d) is significantly amplified by the nonlinear interaction with the charge carrier system. The trailing part "sees" the reduced carrier inversion and consequently experienced a smaller gain leading to an asymmetric pulse shape. In this case, it is, in particular, the spatial variation of inversion and gain that determines the amplification and shaping of the pulse. The situation is changed if the pulse duration is in the order of magnitude of the carrier relaxation (150 fs; Fig. 12c). In that case, the microscopic intradot carrier dynamics as well as the interaction with carriers and phonons of the embedding medium gain importance. In combination with light diffraction and propagation, this leads to a characteristic curvature of the pulse front and modulations in the trailing part. The amount of spatial and temporal distortions the light pulse experiences consequently strongly depends on both, spatial effects (such as dot density, uniformity of the dot distribution, spatial fluctuations) and microscopic "spectral" effects (determined by the characteristic relaxation times and microscopic dot properties). The dynamic amplification and shaping of ultrashort pulses are thus—via the coupling between wave equation and QD-Bloch equations—the result of a complex carrier dynamics within the dots.

## 4.2. Luminescence of Optically Excited Quantum-Dot Media

In Section 4.1, we have considered the propagation of a resonant light pulse in an inverted quantum-dot medium, that is, an electrically pumped QDL. Although this certainly represents the preferred mode of operation of quantum-dot lasers in most applications, in many current experimental setups, however, one investigates the luminescence of optically excited QD media by an (ultra-) short optical pump pulse. In the QD laser, this corresponds to an approximately δ-shaped excitation of carriers into one or more high-energy carrier reservoirs either in the dots themselves (direct optical pumping) or by carrier capture from the optically pumped wetting layer (indirect pumping). Because the dynamic interplay between the dots and the wetting layer is determined by a large variety of relaxation processes involving the dot-carriers and the carriers in the wetting layer, the excitation of the dots via the wetting layer represents a particularly interesting case. In the following, we will thus focus on the

case of an ensemble of (initially empty) dots that is dynamically filled from the wetting layer (the high-energy carrier reservoir). The dynamic coupling between the dot-carriers and the carriers of the wetting layer is determined by the density of quantum dots, the individual dot properties, and the epitaxial growth process. These factors contribute simultaneously and lead to a characteristic response of the quantum-dot system. Thus, the resulting luminescence is influenced by and directly reflects the multitude of physical quantities that are involved in the characteristic excitation and relaxation processes. These are characteristic (material) properties like the optical matrix elements for the various (intradot and dot-wetting layer) carrier and carrier–phonon interactions, the transition matrix elements of the dot levels involved, the energies of the dot levels and the wetting layer states, the dot density, and the spatial dot distribution. Because a detailed analysis and variation of all these parameters (some of which are presently not even known in detail) involves extensive simulations, in the following we will restrict ourselves to an investigation of the influence of the filling degree, the coupling strength and the energy levels. Thereby we implicitly assume the remaining parameters to be (spatially-dependent) constants.

Figures 13, 14, and 15 show the characteristic luminescence, that is, the laterally averaged intensity at the output facet of an optically pumped QDL structure. The degree of initial filling of the high-energy reservoir and the coupling strength between the dots and their environment determines the time constants for carrier capture into the dots and the degree of dot filling. In combination, this may lead to very different characteristic emission behavior discussed below.

### 4.2.1. Influence of Excitation Strength

The three curves displayed in Fig. 13 show the dependence of the luminescence on the filling degree of the carrier reservoir. The solid, gray, and dotted lines show the situation where all, half, and significantly less than half of the available wetting layer states are filled with carriers, respectively. For a moderate initial excitation, the loading of the QDs with carriers is comparatively slow, leading to delayed onset of light emission. An increase in the initial carrier filling of the reservoir provides a higher inversion in the dots, resulting in intense light emission. The variation of the reservoir not only determines the instant and intensity of the light emission, it also affects the shape of the curve: For high excitation, a selective saturation of individual transitions may occur. As a consequence, a second transition can be involved, leading to mode beating, temporal modulations, or a second peak in the emission curve.



Figure 13. Dependence of the QD luminescence on excitation level. Solid, gray, and dotted lines correspond to completely, half, and significantly less than half-filled wetting layer states, respectively. Reprinted with permission from [3]. E. Gehrig and O. Hess, Phys. Rev. 1 65, 033804 (2002). © 2002. American Physical Society.

**Figure 14.** Dependence of the QD luminescence on the coupling strength. High coupling strength (black) provides a fast filling and efficient refilling of the dots that are partially depleted via induced emission processes. Weak coupling (gray) leads to a gradual depletion of the carrier reservoir. Reprinted with permission from [3]. E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002). © 2002, American Physical Society.



**Figure 15.** Dependence of the QD luminescence on the transition energies. The black lines pertain to a QDL with the highest matrix elements separated by more than the LO phonon energy. The gray curves illustrate the case where the respective transition energies are very close to each other. The respective emission properties represented by the solid curve in Fig. 15b shows one intense peak belonging to the main transition. The dashed curves visualize the carrier occupation for the dot system with close transition energies. Reprinted with permission from [3]. E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002). © 2002, American Physical Society.

### 4.2.2. Influence of Coupling Strength

In order to analyze the influence of the coupling strength, we have calculated the intensity and level occupation in the dots in relation to the fraction of wetting layer states to which the dots couple (normalized to a unit cell $\Delta z \times \Delta x$). In a given laser structure, this value is determined by the dot density and by the potential step between the dots and their environment determined by the size and shape of the dots as well as the particular material systems and epitaxial growth processes. The resulting dynamic behavior of the emitted intensity is depicted in Fig. 14. First, a high coupling strength provides a faster filling of the (initially) empty dots. As a consequence, the dot occupation reaches the characteristic threshold value that is apparent in Fig. 14 at an earlier time step. Second, it enables efficient refilling of the dots that are partially depleted via induced emission processes. In combination, this leads to an intense peak in the emission curve (black). A weak coupling (gray), on the other hand, leads to a slow depletion of the carrier reservoir. As a result, it delays the onset of light emission and "stretches" the shape of the emission curve in time.

### 4.2.3. Influence of Quantum-Dot Size and Growth: Variation of Energy Levels

Variation in the size and epitaxial growth of a quantum dot has a direct consequence for its energy levels. These variations in eigen energies enter directly in the QD Bloch equations. In the following, we will consider two channels of transitions with the highest transition matrix elements for two different cases: (1) a QD system with close transition energies (i.e., separated by less than the LO phonon energy) and (2) a QD system where the carrier levels belonging to the two most dominant transitions differ by an energy much higher than the LO phonon energy. For these two examples, Fig. 15 shows the temporal behavior of the (electron) level occupations (Fig. 15a) and the resulting emission curve (Fig. 15b) after the initial excitation of the QDs. The black lines in Fig. 15a pertain to the QDL where the transitions with the highest matrix elements are separated by more than the LO phonon energy. The gray curves correspond to the situation where the respective transition energies are very close to each other. For the example with higher level separation, the carriers populating the two QD levels are mainly decoupled: the carrier recombination is mostly restricted to one level (belonging to the transition with the highest dipole matrix element), whereas the second level absorbs carriers from the reservoir. The respective emission properties represented by the solid curve in Fig. 15b show one intense peak belonging to the main transition. For the same excitation conditions, the dashed curves visualize the carrier occupation for the dot system with close transition energies. In this situation, the two main carrier levels interact and "interfere" via dynamic carrier and phonon scattering. The resulting emission curve (dashed line) shows two maxima. The specific shape of the temporal emission characteristics is thus a direct consequence of the dynamic interplay of competing transitions and spectral modes.

## 4.3. Inter- and Intralevel Carrier Dynamics During the Pulse Propagation

The spatiotemporal dynamics discussed in the past sections are the finger-prints of a complex carrier-dynamics in the charge carrier ensemble in the dots. In the QDMBEs, the inter- and intralevel dynamics are automatically represented and calculated. The level dynamics allows visualization of the microscopic interactions occurring within the dots. As an example, Fig. 16 visualizes the temporal dynamics of the level occupation in quantum dots during the passage of an ultrashort (150 fs) light pulse. Displayed are the level occupations of the hole levels (normalized to their initial value) at the center of the output facet. The respective electron level occupations show a qualitatively similar behavior. We will consider the gain (wavelength of the injected light pulse: 1270 nm) and transparency (wavelength of the injected light pulse: 1216 nm) regime, respectively. The pulse leads to an optical excitation of the carrier system and to a selective level-hole burning (e.g., Fig. 16a): Depending on the dipole matrix elements for the individual states and depending on the frequency detuning of the pulse with respect to the frequency of the respective electron and hole states, a reduction of the individual level occupation and a partial refilling via carrier injection and microscopic

Figure 16. Dynamics of the occupation of hole levels during the propagation of a light pulse: (a) amplification regime. (b) transparency regime.

scattering processes occurs. The microscopic scattering processes involved in this "level-burning" are determined by emission and absorption of phonons, multiphonon interactions, and the interaction with the carriers and phonons of the wetting layer. The magnitude of the various "channels" for relaxation mechanisms thereby depend on the QD energy levels, on the energy difference to the surrounding layers, and on the coupling of a dot to its next neighbors. The dot-to-medium and dot-to-dot interactions thereby are determined by the dot density and the light propagation that mesoscopically couples the QDs. In particular, the levels of a dot may "talk" to each other via carrier and phonon scattering. This dynamic carrier exchange and level dynamics then leads to characteristic modulations that can be seen near transparency (Fig. 16b). We note that this phonon interaction responsible for the fast intradot redistribution of carriers is limited to level separations approximately equal to the LO phonon energy: More deeply confined dots with level separations much larger than this value would not display this behavior but show decoupled level dynamics and selective level hole burning instead. In particular, the degree of coupling between the individual levels strongly effects the saturation behavior of a QD laser waveguide that has been optically injected with a light pulse: For the same characteristics of an injected light pulse (i.e., input power, duration, pulse shape). a QD medium with decoupled level energies (i.e., energy separation larger than the LO phonon energy) will show a stronger saturation behavior than a QD ensemble with a strong coupling of the level energies where a mutual exchange of level occupations via scattering and relaxation dynamics occurs.

The ultrashort carrier dynamics is directly reflected in the microscopic gain which can be derived from the (spatially averaged) microscopic distribution of the dipole density that is dynamically calculated within the framework of the QD Bloch equations. The computational results can then directly be compared to the results of an experimental pump-probe measurement. Figure 17 shows experimental [21] and theoretical results on the dynamics of the gain for three different wavelengths of the injected light. The longest wavelength (1270 nm) corresponds to the amplification regime, whereas the shortest wavelength (1216 nm) corresponds to the transparency regime. The intradot scattering (via emission and absorption of phonons) of the gain typically occurs on timescales of a few hundred femtoseconds up to a few picoseconds leading to a fast partial recovery of the gain. Near transparency the partial exchange of inversion via absorption and emission of phonons leads to characteristic modulations in the gain dynamics that are directly correlated to the corresponding increased level dynamics within the charge carrier ensemble of the QDs. The modulation depth and period thereby depends on the energy separation of the levels as well as on injection current and input power. The temporal regime after approximately 1 ps is characterized by a comparatively slow reestablishment of the inversion via the injection current and the interaction with the carriers of the embedding layers.

The calculation of the ultrashort time dynamics of the charge carriers clearly demonstrates that it is both operation conditions (e.g., injection current density and input power) and material properties (in particular, the energetic separation of the energy levels) that affects the dynamic shaping as well as the saturation characteristics of a propagating light pulse.



**Figure 17.** Measured (a) [21] and calculated (b) gain change during propagation of a short pulse (150 fs) at 1216 nm (solid line), 1250 nm (dashed line), and 1270 nm (dotted line).

## 4.4. Spatiospectral Gain and Index Dynamics

A light pulse propagating in a semiconductor laser directly couples to the carrier populations in the energy levels. Thereby, it induces highly nonequilibrium distributions. As a consequence, dynamic changes in both gain and refractive index arise that relax on femto- and picosecond timescales toward a new equilibrium distribution. Within this temporal regime, the complex dynamics of the charge carriers leads to strong distortions of both the temporal and the spectral profile of the pulse. In order to theoretically investigate the ultrafast carrier dynamics, we calculate the dynamics of the real- and the imaginary part of the polarization at the output facet of a semiconductor laser amplifier that has been optically excited by an ultrashort light pulse (duration 150 fs). This spatially varying polarization is directly correlated and composed of microscopic dipoles that are dynamically calculated within the framework of the Bloch equations. Thereby, every combination of electron and hole states is taken into account via the corresponding dipole matrix element. As a consequence, the polarization dynamics directly reflects the microscopic spectral gain (proportional to the negative imaginary part of the polarization) and the induced refractive index (proportional to the negative real part of the polarization). The lateral and temporal dependence of $P(x,t)$ at the output facet consequently will allow a fundamental analysis of the physical processes that are responsible for temporal and spectral properties of an amplified light pulse. The results of the simulations of gain and induced refractive index are summarized in Fig. 18 in a time window of 1 ps. The bright spot on top of the figures indicates the corresponding dynamics (duration and position with respect to the displayed time axis) of the light pulse. In the calculation, the injection current density has been set to 2.5 $I_{thr}$. The central pulse frequency is located within the amplifier gain bandwidth leading to an amplification of the pulse. During its propagation, the pulse reduces the carrier distribution in the dots, leading to dynamic level hole burning. The gain and induced refractive index are via the polarization of the Bloch equation correlated to the dynamics of electrons and holes in the dots. They thus directly reflect the highly nonequilibrium carrier dynamics in the dots determined by hole burning, carrier injection, as well as carrier–carrier and carrier–phonon scattering processes. The partial refilling of the hole via carrier relaxation and carrier injection typically occurs on timescales of a few hundred femtoseconds. These processes consequently determine the temporal and spectral shape of the "spatio-spectral" trench burnt by the pulse in the spatiospectral gain and index. The level hole burning and carrier heating induced by the pulse "shape" the gain during the passage of the pulse (negative values in Fig. 18a correspond to high gain), leading to amplifying and absorbing regions. At the same time, the reduction of the inversion in the dot levels induces a dynamic induced refractive index (negative values in Fig. 18b correspond an increase in the induced index). In particular, if the pulse duration is < 1 ps, the duration of the carrier relaxation may significantly exceed the pulse duration (i.e., the distributions are still excited after the passage of the pulse). As a result, the pulse envelope may not "see" the entire spectral index dispersion (that typically leads to significant spectral broadening of picosecond pulses) but select only a part of the index curve. The propagating light pulse may consequently experience either a red or a blue shift, depending



Figure 18. Dynamics of imaginary (a) and real (b) part of the polarization (calculated at the output facet of a QD laser amplifier) during the propagation of an ultrashort (150 fs) pulse.

on, for example, input power and injection current shaping the corresponding minima and maxima in the gain and index distributions. The spatiotemporally resolved calculation of the pulse-induced gain and index dynamics thus allows one, for a given laser geometry, to directly obtain and analyze all information on fundamental interaction processes that affects and shapes the nonlinear amplification process and spectral dynamics of ultrashort pulses. As an example, we will in the following vary the width of the active area and analyze the resulting gain and index dynamics. Figure 19 shows for three values of the width of the active area the lateral and temporal dependence of the imaginary and real part of the polarization reflecting gain and index at the output facet of a QD laser amplifier during the propagation of an ultrashort laser pulse. Dark areas reflect high gain and index, whereas light shading reflects absorption and a reduction of carrier-induced index, respectively. The top row shows the corresponding intensity plot. Due to the nonlinear light–matter interactions, the pulse continuously "shapes" its gain and index distribution during propagation. This leads to characteristic temporal modulations in both gain and index. Because carrier relaxation occurs on timescales of a few hundred femtoseconds, the dynamic gain and index distribution are still excited after the passage of the pulse. The temporal shift in gain and index visualize the pulse-induced gain (dark shading) within the pulse, as well as induced index dispersion resulting from the dynamics of the real and imaginary part of the interlevel polarization. In the QD laser amplifier with cavity width $w = 10$ $\mu$m, the distributions are rather uniform. As a consequence, the pulse will be characterized by a uniform lateral beam profile. An increase in lateral width leads to characteristic modulation in transverse direction. In particular, these structures show a characteristic curvature in lateral dimension. This is a direct consequence of the intensity-dependent carrier depletion in the levels induced by the spatial profile of the pulse. It is these optical patterns of dynamic gain and index that shape the pulse in the next time step.

## 4.5. Pulse Propagation in QD and QW Optical Amplifiers

The comparative study of numerical simulations and experiments on the spatiotemporal dynamics and emission characteristics of quantum-well and quantum-dot lasers of identical structure (Section 3.3) has demonstrated that, in the quantum-dot laser, the strong localization of carrier inversion and the small amplitude-phase coupling enable a significant improvement of beam quality compared to quantum-well lasers of identical geometry. It would thus be of interest to investigate whether a similar difference in beam quality can be observed in quantum well and quantum-dot amplifiers that have been optically injected by an ultrashort light pulse. For a comparison of the beam properties of an ultrashort light pulse (150 fs) propagating in a quantum-dot and quantum-well laser of same material and geometry, Fig. 20 shows temporal snapshots of the nearfield profile taken at the output facet during the passage of the pulse maximum. The distributions shown on the left of Fig. 20 were calculated for stripe widths of $w = 10$ $\mu$m. 30 $\mu$m. and 50 $\mu$m, respectively. The pulse



Figure 19. Dynamics of intensity (top row), imaginary (middle row), and real part (bottom row) of the polarization during the propagation of a 150 fs pulse in dependence on stripe width.

Figure 20. Comparison of the spatial beam profile at the output facet of quantum-dot and quantum-well amplifiers during the propagation of the pulse maximum for stripe widths of $w = 10$ $\mu$m. 30 $\mu$m, and 50 $\mu$m (left) and the dependence on injection current density (right).

propagating in the quantum-dot laser amplifier of 10-$\mu$m stripe width shows a laterally uniform profile, whereas the pulse propagating in the quantum-well laser of identical stripe width is characterized by a two-lobe structure. The strong carrier localization in the dots in combination with the low alpha factor lead to a reduced influence of the lateral degree of freedom and consequently to an improved beam quality compared to quantum-well laser of identical material and geometry. The nearfields of the quantum-dot amplifiers with stripe widths of 30 and 50 $\mu$m show small transverse modulation originating from the mutual interaction of light diffraction and dynamic intradot relaxation and scattering between the dots and their environment. However, the structures are still less pronounced than in the case of the quantum-well laser. In a quantum-dot laser, the transverse extension thus is—due to the reduced transverse light-matter coupling—not as crucial as in the bulk media, leading to less transverse beam distortions during the propagation and amplification of ultrashort light pulses. Furthermore, the pulse propagating in the quantum-dot laser amplifier is less affected by an increase in injection current density (Fig. 20, right).

## 5. CONCLUSIONS

In conclusion, we have set up a mesoscopic theory on the basis of a Maxwell-Bloch description. The resulting QDMBEs consist of coupled spatiotemporally resolved wave equations and QD Bloch equations for the electron and hole levels within each quantum-dot of a quantum-dot ensemble inside a quantum-dot laser.

We have presented results of numerical simulations that aim to mesoscopically represent realistic QD laser structures. The simulations include, in particular, microscopic QD properties, spatially dependent QD parameters and fluctuations, spatially inhomogeneous light propagation, and dynamic scattering. The carrier scattering processes are considered on a mesoscopic level and include both the intradot relaxation and the interactions between the QD carriers and the surrounding layers. The specific laser configuration of a model device is considered via the macroscopic boundary conditions and constraints. The QDMBEs allow calculation and visualization of spatial distributions of the light field intensity and carriers. The dynamic calculation of level occupations provides a detailed analysis of the various relaxation processes. Furthermore, computational results on the basis of the QD-Maxwell-Bloch equations allow the calculation of measurable quantities such as beam quality, nearfields, and optical spectra. Simulations of the ultrashort time dynamics in optically injected quantum-dot lasers allow a fundamental analysis and visualization of the light-field and carriers dynamics, as well as the calculation of spatiospectral gain and index dynamics affecting amplitude and phase of an ultrashort light pulse propagating in the spatially inhomogeneous QD ensemble.

Furthermore, the calculation of level occupations provides a detailed analysis of the various relaxation processes. For a specific set of parameters (injection current, pulse, and power of the injected light pulse) the quantum-dot Maxwell-Bloch theory allows a microscopically founded interpretation of spatiospectral saturation and dynamic pulse shaping. The mesoscopic theory and computational modeling discussed in this chapter may thus establish a basis for linking the microscopic analysis of QD material properties with the quantum electronics of modern quantum-dot laser systems.

## REFERENCES

*1.* D. Bimberg and N. N. Ledentsov, *J. Phys. Condens. Matter* 15, R1 (2003).

*2.* D. Bimberg, M. Grundmann, and N. N. Ledentsov, "Quantum Dot Heterostructures," John Wiley, Chichester, 1999.

*3.* E. Gehrig and O. Hess, *Phys. Rev. A* 65, 033804 (2002).

*4.* H. F. Hofmann and O. Hess, *Phys. Rev. A* 59, 2342 (1999).

*5.* O. Stier, M. Grundmann, and D. Bimberg, *Phys. Rev. B* 59, 5688 (1999).

*6.* M. Braskén, M. Lindberg, D. Sundholm, and J. Olsen, *Phys. Rev. B* 61, 7652 (2000).

*7.* K. Oshiro, K. Akai, and M. Matsuura, *Phys. Rev. B* 59, 10850 (1999).

*8.* P. G. Bolcatto and C. R. Proetto, *Phys. Rev. B* 59, 12487 (1999).

*9.* E. Gehrig and O. Hess, "Spatio-Temporal Dynamics and Quantum Fluctuations in Semiconductor Lasers," Springer, Heidelberg, 2003.

*10.* J. Schilp, T. Kuhn, and G. Mahler, *Phys. Rev. B* 50, 5435 (1994).

*11.* A. V. Uskov, J. McInnerney, F. Adler, H. Schweizer, and M. H. Pilkuhn, *Appl. Phys. Lett.* 72, 58 (1998).

*12.* J. L. Pan and P. L. Hagelstein, *Phys. Rev. B* 49, 2554 (1994).

*13.* A. V. Uskov, K. Nishi, and R. Lang, *Appl. Phys. Lett.* 74, 3081 (1999).

*14.* G. Bastard (Ed.), "Wave Mechanics Applied to Semiconductor Heterostructures," Halsted, New York, 1988.

*15.* O. Hess and T. Kuhn, *Phys. Rev. A* 54, 3347 (1996).

*16.* E. Gehrig, O. Hess, C. Ribbat, R. L. Sellin, and D. Bimberg, *Appl. Phys. Lett.* 84, 1650 (2004).

*17.* Ch. Ribbat, R. L. Sellin, and D. Bimberg, Technical Universität Berlin, Germany, 2004.

*18.* R. Lang, A. Hardy, R. Parke, D. Mehuys, and S. O'Brien, *IEEE J. Quant. Electr.* 30, 658 (1994).

*19.* J. R. Marciante and G. P. Agrawal, *IEEE J. Quant. Electr.* 32, 590 (1996).

*20.* E. Gehrig, O. Hess, and R. Wallenstein, *IEEE J. Quant. Electr.* 35, 320 (1999).

*21.* M. van der Poel, D. Birkedal, and J. Hvam, Research Center COM, Technical University of Denmark, Lyngby, Denmark, 2004.

# CHAPTER 12

# Theoretical Investigations of Silicon Quantum Dots

## Lin-Wang Wang

*Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, California, USA*

## CONTENTS

## 1. INTRODUCTION

Silicon quantum dots can be divided into two categories: the relatively large quantum dots made on Si substrate, and the nanometer crystallites. Current lithography techniques can be used to make quantum dots and wires as small as 20 nm [1, 2]. These quantum dots are

used for single-electron devices [1, 3, 4] taking advantage of their Coulomb blockade and quantum transport effects [5–7]. One or two monolayer Ge depositions on Si substrate can produce self-assembled Ge or GeSi alloy quantum dots [8–10]. This type of quantum dots has a base size larger than 50 nm [8–10]. The exciton radius in bulk silicon is about 5 nm; thus, to change the electron wavefunction and to see the quantum confinement effects, one needs to used the smaller, nanometer quantum dots. These are the quantum dots on which we will focus in this review.

There are many ways to synthesize nanometer Si crystallites. Because the covalent Si–Si bond is much stronger than the bonds in II–VI and III–V materials, a high temperature (e.g., 1000°C) is often required to form this bond, and hence to synthesize Si quantum dots. As a result, the wet chemistry methods [11–13] used to synthesize various types of II–VI and III–V quantum dots are difficult to apply to synthesizing IV–IV Si quantum dots. Under high temperature, the Si quantum dot can be synthesized either in gas phase or in an embedding matrix. Under these conditions, the synthesis parameters (temperature, density, etc.) cannot be controlled as well as they can in wet chemistry conditions. This has resulted in rather large size distributions in most synthesized Si quantum dots; which makes the small Si quantum dot study rather difficult and is one of the main reasons why there are still so many controversies after so many years of studies.

In the gas phase synthesis, Si sources are provided by organic gas molecules containing Si atoms (e.g., disilane) [14] or a melted and vaporized solid Si source [15, 16]. The aerosol reaction can happen in a chamber and is usually followed by reactions with passivating gases (e.g., H or O). Another common way to provide Si source is via laser ablation [17], where part of the Si crystal is evaporated under laser pulse. Similarly, sputtering caused by ionic bombardment can also be used to generate Si crystallites [18].

For all these techniques (gas phase, laser ablation, and sputtering), the Si crystallites are formed in a open space (air or vacuum) and then simply collected on a substrate. Another approach is to ion-implant Si in some substrate (e.g., $SiO_2$), and then form nanometer Si quantum dots through thermal annealing [19]. Recently, it has been demonstrated that this ion implantation method can synthesize Si nanocrystals in quartz or in $SiO_2$ on silicon substrate, and the resulting samples have achieved optical amplification—a major step toward Si-based laser [20].

One extremely important Si nanosystem is porous Si (p-Si). P-Si has been popular ever since Canham [21] reported strong photoluminescence (PL) of it 13 years ago. Bulk Si is an indirect material with almost no photoluminescence. This is a problem for building optical devices based on Si substrate and for integrating optics with Si-based microelectronics. P-Si can be produced by anodic etching (this was known back in 1956 by Uhlir [22]). In this etching process [23], Si wafer is used as the anode in an electrolyte solution containing HF. When the electric current is larger than a critical current density, p-Si will be formed on the surface of the Si wafer. After this etching process, the Si wafer is taken out of the solution and dried in the air with special care (to avoid cracking the skeleton of the porous structure). Because of the simplicity and low cost of this procedure, p-Si has become very attractive. It provides a potential way to integrate the optics on the Si wafer under the same technology of Si microelectronics. During the last 10 years, there has been almost one paper per day for p-Si. It is fair to say that p-Si has dominated the Si nanometer crystallite research. In some extent, porous Si is a synonym of Si nanometer crystallite.

Current p-Si can be made with porosity up to 95% [24]. The porousness is generated when the etching channels penetrate the Si wafer and start to branch [25–27]. After the etching, only a skeleton of the Si remains in the p-Si. The morphology is shown in Fig. 1 [28], where we can see the vertically penetrating etching channels and the remaining Si skeletons. However, it is not obvious why the remaining Si skeleton in the p-Si should be considered as Si quantum dots, and not wires or some more complicated structures. A transmission electron microscope (TEM) image of the p-Si is shown in Fig. 2 [29]. From that figure, we can see a nanometer-sized grainy structure. A schematic showing how the etched p-Si become isolated quantum dots is summarized in Fig. 3 [30]. As one can imagine, there could be many sizes and shapes of such quantum dots in a p-Si. This makes quantitative comparison between experiment and theory for p-Si more difficult than Si nanocrystallites synthesized by the other methods

Figure 1. Scanning electron microscope images of cleaved cross-section of porous Si. Reprinted with permission from [28]. J. P. Gonchond et al., in "Microscopy of Semiconducting Materials 1991" (A. G. Cullis and N. J. Long, Eds.), p. 235. Institute of Physics. Bristol. 1991. © 1991, Institute of Physics.

discussed earlier. One advantage of the p-Si is that all the remaining Si skeletons are still aligned roughly in the original crystal orientation as the Si wafer. This can be shown from the transmission electron diffraction pattern [29]. This makes it possible to do polarization experiments for photoemission [31]. The x-ray experiments [14] and direct TEM observations [32] have confirmed that the internals of the Si quantum dots are still in a bulk diamond structure with a lattice constant within 0.25% of the bulk Si value. This is true both for the p-Si and the nanocrystallites synthesized by the methods discussed above.

As complex as the p-Si structures shown in Fig. 1 and Fig. 2 are, it is no wonder that there have been many suggestions for the mechanism of the strong PL from p-Si. These suggestions range from crystalline Si to surface state recombinations. Six suggestions are summarized in Fig. 4 [33]; they are: crystalline silicon, hydrogenated amorphous silicon, surface hydrides, defect states, molecules, and surface state recombinations. Over the years, however, only the crystalline silicon model has survived. This is mainly because of the following observations: first, the PL energy changes with p-Si grainy size, having larger PL energy for smaller sizes, showing quantum confinement effects; second, phonon-assisted PL emissions with the phonon energies matching the ones in bulk phonon assisted transitions, as shown in Fig. 5 [34], and with the spectrum shape agrees with phonon replicas analysis, as shown in Fig. 6 [35]; third, good agreements between experimental measurements and theoretical



Figure 2. Transmission electron microscope images of porous Si samples: nanometer scale, columnar Si structures arrowed. Reprinted with permission from [29]. A. G. Cullis and L. T. Canham. *Nature* 353, 335 (1991). © 1991, the Nature Publishing Group.

Figure 3. Idealized schematic steps in the oxidization process of highly porous silicon. Reprinted with permission from [30], L. T. Canham. "Optical Properties of Low Dimensional Silicon Structures" (D. C. Benshael et al., Eds.), NATO ASI Series, Vol. 244, p. 81. Kluwer Academic Publishers, Dordrecht, 1993. © 1993, Kluwer Academic Publishers.

calculations based on Si quantum dot models. This includes PL energy confinement effects, exchange splitting, phonon-assisted transition, emission lifetimes, Auger process rate, and surface passivation effects. All these agreements are difficult to be achieved by the other models shown in Fig. 4.

The strong photoluminescence in Si nanocrystallite and p-Si was first explained by the breaking of the momentum conservation rule in nanostructures. As a result, the valence-band maximum (VBM) to conduction-band minimum (CBM) transition becomes pseudodirect, and dipoles are allowed. However, the PL lifetime calculated this way is larger than the experimentally measured lifetime (Fig. 7 [36]). Later, it was found both experimentally and theoretically that the phonon-assisted emission is still stronger than the zero phonon emission (Fig. 5). This means that the pseudodirect transition might not be the main reason why the PL intensity is strong in these nanosystems. A more relevant explanation is the exciton localization. In the Si crystallite, an excited exciton is physically confined within one quantum dot. As long as there is no surface nonradiative center, the exciton will eventually decay radiatively with an emission of a PL phonon. This explains why although the PL lifetime in Si quantum dot is very long (in microseconds), the PL quantum efficiency is still high (on the order of 50% for the best Si nanocrystals and 10% for p-Si). In the bulk Si, the exciton will move around and will usually be captured by some nonradiative center, or will aggregate around some neutral impurity. In the later case, although the impurity will not annihilate the exciton directly by emitting phonons, the Auger process among different excitons will kill the exciton nonradiatively. As a results of these processes, the PL quantum efficiency in the bulk Si is only 10⁻⁴% [37]. In Si nanocrystals, when the power of the optical pump is large enough, there will be more than one exciton in each quantum dot. Then, similar Auger processes will be activated to annihilate excitons nonradiatively. This is the reason of the nonlinear optics in Si nanocrystals. However, in the linear optical regime, there is only one exciton in one quantum dot; thus, the Auger process does not play any important roles, and the PL quantum efficiency is high.

a) CRYSTALLINE SILICON

b) HYDROGENATED AMORPHOUS SILICON

c) SURFACE HYDRIDES

d) DEFECTS

NBOHO

SD

SD

PB

e) MOLECULES

f) SURFACE STATES

KEY

O – silicon atoms

● – oxygen atoms

• – hydrogen atom

Figure 4. The six groups of models proposed to explain the Photoluminescence (PL) from porous Si. (a) The PL comes from the interior of the crystalline silicon. (b) The PL comes from hydrogenated amorphous silicon surrounding the crystalline silicon. (c) The radiation comes from the silicon hydride bonds. (d) The PL comes from various defects as radiative centers. (e) The PL comes from Si-based polymers with a structure of siloxene. (f) The PL comes from surface localization states. Reprinted with permission from [33], A. G. Cullis et al., *J. Appl. Phys.* 82, 909 (1997). © 1997. The American Institute of Physics.



Figure 5. Photoluminescence (PL) spectrum taken at 2 K with resonant excitation. Vertical dotted lines indicate the position of the no-phonon (NP) onset and its transverse optical (TO) and transverse acoustic (TA) phonon replicas. Reprinted with permission from [34], P. D. J. Calcott et al., *J. Lumin.* 57, 257 (1993). © 1993. Elsevier Science.

**Figure 6.** (a) The experimentally observed variation of the Photoluminescence (PL) line shape (at 2 K) with laser energy. For each spectrum, the value of the laser energy is indicated by a vertical dotted line. (b) The modeled variation of the PL line shape with laser energy, assuming that all the PL is from quantum confined crystalline Si. Reprinted with permission from [35], P. D. J. Calcott et al., in "Microcrystalline and Nanocrystalline Semiconductors" (R. W. Collins et al., Eds.), Materials Research Society, p. 465. Pittsburgh, PA. 1994. © 1994, Materials Research Society.

Although the exciton localization enhances the radiative decay of the exciton, it also causes problems for device applications. In device applications, not only PL is important; what is more important is the electroluminescence (EL). To realize EL, the electron and hole need to be provided electrically, which requires electron and hole conductivities. This is in direct conflict with the electron and hole localizations in these nanosystems. That is why the EL efficiency in the p-Si is about 10 times smaller than the PL efficiency. Although the PL efficiency in p-Si can be 10%, the best EL efficiency is still less than 1% [38]. Improving the EL efficiency is still a major challenge.

In the early days of Si nanocrystal and p-Si studies, different experiments often produced very different results. In addition, the calculated optical band gap was often found to be larger than the experimentally measured ones. Recently, it has been realized that



**Figure 7.** Calculated recombination rate of an excited electron-hole pair in silicon crystalline (crosses) with respect to the photon energy at 300 K. Continuous line plots the experimental dependence of decay rates on photon energy for three 65% porosity layers that differ by oxidation level. Reprinted with permission from [36], M. Lannoo and C. Delerue, in "Structural and Optical Properties of Porous Silicon Nanostructures" (G. Amato et al., Eds.), p. 187. Gordon and Breach, Amsterdam. 1997. © 1997, Gordon and Breach.

the surface passivation and oxidization are extremely important. In ambient air, oxygen will incorporate into a fresh-etched p-Si surface in a few minutes [39, 40]. Thus, unless extreme cares are taken, there will always be oxygen on the surface. In addition, to stablize the surface passivation, the p-Si samples are often intentionally oxidized [41]. This oxidization will introduce a red-shift in the PL energy, as shown by many recent theoretical calculations. This explains why previous calculations assuming a pure H passivation have produced PL energies larger than the experimental ones. After being exposed to ambient air for a long time, the surface of the p-Si will be oxidized aggressively and might eventually become pure $SiO_2$, as illustrated in Fig. 3. As a result, the interior Si nanocrystallite becomes smaller with time, causing larger quantum confinements and blue-shifts of the PL energy, as shown in Fig. 8 [30]. Thus, extreme care must be taken when comparing experiments with the theoretical calculations, or when comparing different experimental results. The surface passivation condition must be taken into account. These surface passivation conditions for porous Si are summarized in Table 1 [33].

In summary, Si nanocrystallites and p-Si are complex systems compared with other nanocrystallites like CdSe. Experimentally, their sizes and shapes are more difficult to control than the II–VI quantum dots, and their surface passivations are critical to the PL process. However, the motivations are high because of the potentials to integrate optics with Si-based microelectronics and to use the mature Si process technologies in making these devices. The recently realized optical amplification resulting from silicon nanocrystals on Si substrate [20] and the possibility of using optical communication as the interconnect in a Si chip [42] have stimulated new interest in these systems. Unlike the lithography-generated large Si quantum dots, where the interest is in using them as single electron devices under the Coulomb blockade, the interest in smaller Si nanocrystallites is mainly in its optical properties. This is the area we will focus on in this review. Physically, the bulk Si is an indirect band-gap material. This is different from the II–VI direct band-gap material. As a result, phenomena like the phonon-assisted transitions play important roles in Si quantum dots, which are not important in direct band-gap II–VI quantum dots. These new phenomena provide the physical incentives to investigate these systems.



Figure 8. Photoluminescence spectra from a porous Si layer (porosity 77%, thickness 11.6 $\mu$m) following storage in air for the indicated times. Reprinted with permission from [30], L. T. Canham, in "Optical Properties of Low Dimensional Silicon Structures" (D. C. Benshael et al., Eds.), NATO ASI Series, Vol. 244, p. 81. Kluwer Academic Publishers, Dordrecht, 1993. © 1993. Kluwer Academic Publishers.

Table 1. Surface chemical terminations for porous Si under various synthesis conditions.

| Condition of Porous Si | Chemical Termination of Skeleton Surface |
| --- | --- |
| *In-situ* in HF during and after formation | $SiF_xH_y$ |
| Freshly etched in interambient air | $SiH_x$ |
| Placed in ambient air for hours | $SiO_xH_y$ |
| Chemically or anodically oxidized | $SiO_xH_y$ |
| Aged in ambient air for months to years | $SiH_xC_z$ |
| Rapidly thermally oxidized at high temperatures | $SiO_xH_y$ |

*Source:* Reprinted with permission from [33]. A. G. Cullis et al., *J. Appl. Phys.* 82, 909 (1997). © 1997, The American Institute of Physics.

The poor experimental characterization of these systems, on the one hand, makes the comparison with theory difficult, but on the other hand, it also makes reliable theoretical investigations important and necessary. The indirect band gap means the widely used effective mass k.p method for the II–VI quantum dots can no longer be used effectively for the Si quantum dots. The sensitivity to surface passivation makes atomistic calculations necessary. Thus, new methodology development is essential. Like bulk Si and Si surfaces, Si quantum dot has also served as a test ground for many theoretical methods of computational nanoscience.

There are already many excellent review articles for p-Si and Si nanocrystallites [33, 43–45], both for experimental measurements and theoretical results. Not to repeat these previous reviews, but here we will focus on the theoretical methodologies, their developments, and their results in Si nanosystems. We will use Si nanocrystallites as a test ground for these different theoretical methods.

## 2. SINGLE-PARTICLE METHODOLOGIES

### 2.1. A Simple Procedure to Calculate the Optical Band Gap

In a few eminent papers [46–48] about 20 years ago, Brus proposed a simple effective mass model to calculate the ionization and optical energies for nanocrystals. In this model [46], the energy of an electron inside a quantum dot has been divided into the kinetic energy (which is described by the effective mass model) and an electrostatic potential energy caused by a surface polarization potential $P(r)$. This $P(r)$ is the classical electric–static potential caused by a point charge at position $r$ inside a nanocrystal of dielectric constant $\epsilon$. The idea is to compare the bulk and the quantum dot for their electrostatic interaction energies between the bare electron and the medium. The extra interaction energy for an electron at $r$ of a quantum dot is $P(r)$. If the nanocrystal is infinitely large, and $r$ is away from the surface, $P(r)$ will be zero. This is similar to the image potential when a point charge is near a dielectric medium or near a metal surface, but here, the electron is inside the medium, not outside. Thus, this potential $P(r)$ is also called the image potential of the charge at $r$, and it is induced by the surface of the nanocrystal. On quantization, $P(r)$ is treated as an on site potential in the particle's Schrodinger's equation. For an effective mass theory, we then have

$$\left[ -\frac{1}{2m}\nabla^2 + P(r) \right]\psi(r) = E\psi(r) \tag{1}$$

here $m$ is the effective mass and $E$ is the ionization energy (or say the electron affinity) of the small nanocrystal. Note that the above equation is only for the electron inside the nanocrystal. Outside the nanocrystal, the wavefunction is zero for an infinite barrier. In other words, we should also add another confinement potential $V_{conf}(r)$ in Eq. (1), where $V_{conf} = 0$ for $r$ inside the dot, and $V_{conf} = \infty$ for $r$ outside the dot. Equation 1 is the basic single-particle Schrodinger's equation, although the effective mass kinetic energy and the infinity confinement potential $V_{conf}(r)$ can be replaced by more realistic models.

When there are two particles, Brus [47] proposed an classical electrostatic interaction energy among these two particles as

$$V'(\mathbf{r}_1, \mathbf{r}_2) = \pm \frac{e^2}{\epsilon|\mathbf{r}_1 - \mathbf{r}_2|} \pm P_M(\mathbf{r}_1, \mathbf{r}_2) + P(\mathbf{r}_1) + P(\mathbf{r}_2) \tag{2}$$

here the first two terms are interactions between these two particles, whereas the last two terms are self-interactions, as in Eq. (1). The $P_M$ is the interaction between one particle's surface-induced polarization and the other particle's charge. Note that $P(\mathbf{r}) = P_M(\mathbf{r}, \mathbf{r})/2$. In Eq. (2), the plus is for two same charged particles (two electron or two hole), the minus is for an electron–hole system. Place $V'(\mathbf{r}_1, \mathbf{r}_2)$ in an effective mass quantum confinement Hamiltonian; we then have [47]:

$$H_{\text{exciton}} = -\frac{1}{2m_e}\nabla^2 - \frac{1}{2m_h}\nabla^2 - \frac{e^2}{\epsilon|\mathbf{r}_e - \mathbf{r}_h|} + \text{polarization} \tag{3}$$

Again, $m_e$ and $m_h$ are the electron and hole effective mass, respectively. In Eq. (3), polarization $= -P_M(\mathbf{r}_e, \mathbf{r}_h) + P(\mathbf{r}_e) + P(\mathbf{r}_h)$. In solving Eq. (3), the screened electron-hole interaction term $-\frac{e^2}{\epsilon|\mathbf{r}_e - \mathbf{r}_h|}$ + polarization can be treated iteratively. As a result, the effective mass solution for Eq. (3), under infinite barrier of spherical radius R can be written down analytically as

$$E^* \simeq E_g + \frac{\pi^2}{2R^2}\left[\frac{1}{m_e} + \frac{1}{m_h}\right] - \frac{1.8e^2}{\epsilon R} + \text{small-term} \tag{4}$$

here $E_g$ is the bulk semiconductor band gap. The small–term is caused by the polarization term in Eq. 3. The majority term of $-P_M(\mathbf{r}_e, \mathbf{r}_h)$ cancels that in $P(\mathbf{r}_e)$ and $P(\mathbf{r}_h)$ because $P(\mathbf{r}) = P_M(\mathbf{r}, \mathbf{r})/2$. However, because in $-P_M(\mathbf{r}_e, \mathbf{r}_h)$, $\mathbf{r}_e$ and $\mathbf{r}_h$ are not the same, this cancellation is not complete. Nevertheless, this remaining small-term is often small and can be ignored in practice for spherical quantum dots.

Equation (4) outlines a general strategy to calculate the exciton energy (or, say, the optical band gap) of a nanocrystal: first calculate the single-particle electron and hole eigen energies under quantum confinement [Eq. (1) without the $P(\mathbf{r})$ term], then calculate the electron and hole screened interaction inside the nanocrystal perturbatively. This is, however, only an *ansatz* derived from classical electrostatic considerations, not from quantum mechanical formalism, which treats the whole nanocrystal as a many-body quantum system. This approach especially ignores the dynamic screening of the system and the local field effects of the dielectric function. Thus, this approach is probably good only for quantum dots larger than a few atomic bond length. When Eq. (1) is solved without the $P(\mathbf{r})$ term, the definition of the single-particle eigen energy $E$ needs to be treated with care. Here, we will assume that the single-particle Hamiltonian $H$ [excluding the $P(\mathbf{r})$ term] is local in real space (i.e., the operation of $H$ is defined point by point), and the operation of $H$ at a given point $\mathbf{r}$ depends only on the local atomic structure and charge density around $\mathbf{r}$. Thus, if the interior of a quantum dot has a bulk crystal structure, then the $H$ inside the dot will be the same as the $H$ inside a bulk. $H$ will only be changed near the dot surface with a confinement potential. This gives an unique definition of the single-particle Hamiltonian and its eigen energy. As we will show later, the single-particle energy defined this way will be different from the quasi-particle energy defined in a usual GW formalism. However, the single-particle energy defined here is a natural extension of the traditional quantum confinement treatment in lower-dimension systems like superlattices and quantum wells. In the following text, we will discuss a few approaches to calculate the single-particle wavefunctions and eigen energies, using the single-particle Hamiltonians defined here.

## 2.2. Effective Mass and K.P Methods

Effective mass and k.p methods [49–52] have been used widely to study the quantum confinement effects in superlattices and quantum wells. They have also been used to study the energy levels of quantum dots, especially for direct materials [53]. For a single-band effective

mass theory, the electron wavefunction is described by its smooth envelope function $\phi(\mathbf{r})$. The equation for $\phi(\mathbf{r})$ under confinement potential $V_{conf}$ is

$$\left[ -\frac{1}{m^*} \nabla^2 + V_{conf}(\mathbf{r}) \right] \phi(\mathbf{r}) = E\phi(\mathbf{r}) \tag{5}$$

here, $m^*$ is the effective mass of the electron or hole. For a simple model, $V_{conf}$ is zero for $\mathbf{r}$ inside the nanocrystal, and infinite for $\mathbf{r}$ outside. Under this assumption, Eq. (5) can be solved analytically for simple geometric systems. For a spherical quantum dot with radius $R$, $E = \frac{\pi^2}{2m^*R^2}$. However, it was found that such calculated band gaps are much larger than the experimentally measured values [54, 55] for Si quantum dots. Effective mass theory often overestimates the quantum confinement effect.

Effective mass theory can be considered as a lowest-order Taylor expansion of the bulk band structure around a given high-symmetry k-point. If there are degenerated states at this k-point (e.g., at the $\Gamma$-point top of valence bands) single-band effective mass Hamiltonian can no longer describe the band structure properly. In this case, the wavefunction needs to be described as a sum of these few bands: $\sum_n u_n(\mathbf{r})\phi_n(\mathbf{r})$, where $u_n(\mathbf{r})$ is the Bloch state is of band $n$ at the symmetry k-point, and $\phi_n(\mathbf{r})$ is the envelope function of this Bloch state. Then the electron wavefunction is represented by the envelope function set $[\phi_n(\mathbf{r})]$, and the k.p Hamiltonian is used to describe the movements of $[\phi_n(\mathbf{r})]$ as

$$H_{k.p}[M \times M][\phi_n(\mathbf{r})] + V_{conf}(\mathbf{r})[\phi_n(\mathbf{r})] = E[\phi_n(\mathbf{r})] \tag{6}$$

here $H_{k.p}[M \times M]$ is a $M \times M$ matrix with matrix elements consisting of $\nabla$, $\nabla^2$ and k.p parameters, and $M$ is the number of bands included in the k.p model. The construction of $H_{k.p}[M \times M]$ is purely based on the symmetry of the bands at the high-symmetry k-point. For a given band symmetry, the lowest-order band structure Taylor expansion with $\mathbf{k}$ has a fixed form. After changing $k^2$ to $\nabla^2$, and $\mathbf{k}$ to $\nabla$, we get the Hamiltonian matrix $H_{k.p}[M \times M]$ in the above equation. A detailed description of the k.p Hamiltonian can be found in Refs. [50, 51, 56]. Note that the single-band effective mass Eq. (5) is just a special case of the k.p model with $M = 1$. For most semiconductors, to describe the top of valence band at the $\Gamma$-point, a six-band (including the spin) k.p model is often used. If the band gap is small, the coupling between the conduction band and the valence band is important. Then an eight-band k.p model with six bands at the top of valence band and two bands at the bottom of conduction band is often used.

Six-band and eight-band k.p models are used very successfully in the calculations of semiconductor superlattices and quantum wells [57]. As a natural extension, they have also been widely used to study quantum dots. Under spherical approximation of the k.p model, the six-band and eight-band k.p wavefunction [Eq. (6)] can be solved analytically [58, 59]. This has contributed to much of our understanding of electronic structures of many quantum dots, especially for direct-material systems. However, even for direct materials like CdSe and InP, problems still exist for the use of k.p Hamiltonian. Unlike many superlattices and quantum well systems, in which the size can be rather large, the colloidal quantum dots are often only a few nanometers in size. In reciprocal space, this might correspond to a k point 1/3 toward the Brillouin zone boundary. As a result, the lowest-order Taylor expansion might no longer be adequate. Indeed, it is found that the k.p result and a more accurate calculation might differ by 50% in their confinement energy [60], and sometimes the energy order of different states might be wrong in the k.p calculations [61]. In addition, the boundary condition is often problematic for k.p Hamiltonian, especially for the eight-band k.p model. Without care, spurious states might appear in k.p calculations [62].

Silicon is an indirect material. Thus, the application of the k.p model is often complicated. As mentioned above, in the early days, single-band effective mass theory was used for Si quantum dot band structure calculations [54]. However, this is apparently inadequate for the valence bands. Recently, Niquet et al. [63] have tested the six-band k.p model with the more accurate tight-binding calculation for Si quantum dots. It was found that the k.p method is not accurate even for quantum dot diameters as large as 8 nm. For the conduction band, the intermixing of the six bands near $X$ point-degenerated valleys makes the lowest-conduction

band state $a_1$, $e$, or $t_1$ symmetries, depending sensitively on the quantum dot sizes [64]. To adequately describe these state symmetries, a intermixing between these degenerated valleys has to be considered. This is beyond the commonly used k.p models.

## 2.3. Truncated Crystal Methods

As discussed above, one of the problems of effective mass theory is that the involved effective k points, which are inversely proportional to the size of the nanosystem, are far away from the Γ point. As a result, the band structure is no longer under the parabolic approximation, as assumed in the effective mass theory. To correct this problem, one approach is to directly use the bulk-band structure. This has been called the truncated crystal method (TCM). In this method, the energy $E$ of the one-electron state in a nanosystem is taken directly to be the bulk-band structure energy $\epsilon_n^{bulk}(k)$ ($n$ is a bulk-band index) at some appropriate k-point $k$. Thus, the nanostructure states are mapped into the bulk-band structure, and the task is to find the relationship between the nanostructure states and the k-point.

For thin films, the truncated crystal method often works well. Figure 9 [65] shows one mapping scheme between the directly calculated thin-film eigen energies and their corresponding bulk energies for Si (001) thin-film [65]. These k points along the (001) direction are taken as $\pi j/L$, with $L$ being the thin-film length and $j$ being the eigen state index 0, 1, 2 .... The error of this mapping is less than 50 meV. The bulk-band structure and the thin-film eigen states are calculated using empirical pseudopotential method (EPM), and the surfaces of the Si thin film are not passivated [65].

For an exact one dimensional system, with a Hamiltonian like $[-\frac{1}{2m}\frac{d^2}{dx^2} + v(x)]$ inside $[-L/2, L/2]$ [and $v(x+a) = v(x)$, $v(x) = -v(x)$], and infinite potential outside this region, it has been proved by Ren [66, 67] that the thin-film eigen energy is exactly $E_{n,j} = \epsilon_n^{bulk}(\pi j/L)$. This agrees well with the above Si thin-film example, although the Si thin film is not exactly a mathematical one-dimensional case, and the boundary does not abruptly go to infinity. Another interesting finding from both the Si thin-film calculation and Ren's derivation for the one-dimensional system is that there could be a "zero" confinement state, with its eigen energy independent of the film thickness. The wavefunction of this "zero" confinement state is uniformly distributed inside the thin film. This state is a result of a bulk-state nodal plane parallel to the surface. Consequently, the existence of this state depends sensitively on the boundary condition. For example, it is found that when the surface of the thin film



**Figure 9.** Mapping of the directly calculated film energy levels (solid dots) at the two-dimensional Brillouin zone center Γ onto the truncated crystal energy level for (a) a 12-layer Si (001) film and (b) an 11-layer Si(001) film. The vertical dotted lines indicate the quantized $k_z$ values. Thus, the intersections between the bulk dispersion and the vertical dotted lines give the truncated crystal energy levels. Reprinted with permission from [65], S. B. Zhang et al., *Phys. Rev. B* 48, 11204 (1993). © 1993, The American Physical Society.

is passivated by hydrogen, the energy of this "zero" confinement state will change with thin film thickness.

Rama Krishna and Friesner (RKF) [68] have applied the truncated crystal idea to a spherical quantum dot with a radius $R$. There the k is taken as $\pi j/R$ along (111) direction. Both conduction-band states and valence-band states are calculated this way. When compared with the results of experiments, these results are found to be better than the effective mass results. However, unlike in the thin-film cases, it was found later [69] that the eigen energies calculated from the RKF method and the direct calculation method do not agree well. The RKF method underestimated the quantum confinement energy, as shown in Fig. 10 [55, 68, 69]. There are two reasons for this discrepancy: first, there is no prior reason to argue that why the k point should be in (111) direction, but not in other directions. For a nonisotropic band structure like the Si top of the valence band, this ambiguity of direction can cause problems. Second, The truncated crystal approach does not allow band mixing. This is okay for one-dimensional confinement systems like the thin film, where different bands do not mix together. However, as can be seen from spherical quantum dot k.p calculations for the valence-band states, different bands will be mixed. This is also true in other shapes (e.g., the cubic shape) of quantum dot systems [69]. Actually, there is no known quantum dot shape with which different valence-band states will not mix. In contrary, the RKF formula always takes the highest bulk valence-band energy without mixing as the quantum dot eigen-state energy. As a result, the RKF valence-state energies are always higher than the directly calculated values, and the confinement energy is underestimated.

Attempts have been made to improve the truncated crystal method for nanosize quantum dots. For example, instead of one k point, multiple k points, or even k-point integrals can be used. To get the appropriate band mixing is still a challenge. One idea is to use the k.p method to figure out the band mixing and k-point coefficients [70], and then use the bulk-band structure to correct the k.p band dispersion error. Before such methods become reliable, however we have to rely on direct atomistic calculations.

## 2.4. Tight-Binding Methods

One problem of the continuum effective mass and k.p method is that it does not have the information of the intrinsic length scale of the atomic structure. As a result, it does not have important information like the Brillouin zone, and the Hamiltonian has the same form: one is either calculating a 100-nm system or a 2-nm system. To overcome this problem, atomistic calculation should be used. The tight-binding (TB) method is the simplest method to include the atomic structure in the calculation. In an empirical TB method, the atomic TB basis functions are not explicitly assumed. What do enter into the calculations are the parameters of onsite energies and overlap energies between neighboring atom basis sets.



Figure 10. Band gap versus the effective size $d$ for three prototype quantum dot shapes. ♦, sphere; ×, rectangular boxes; ■, cubic boxes. Also shown are effective mass (EMA) results [55] and truncated crystal results of Rama Krishna and Frisna [68]. Reprinted with permission from [69], L. W. Wang and A. Zunger, *J. Phys. Chem.* 98, 2158 (1994), © 1994, The American Chemical Society.

Ren and Dow [71] have first used the empirical TB method to calculate Si nanocrystals containing 3019 Si atoms. The quantum dot surface is passivated with hydrogen atoms. A $sp^3p^*$ TB basis and nearest-neighbor interaction were used in their calculation. This leads to a $16397 \times 16397$ matrix for its Hamiltonian. The $T_d$ symmetry of the spherical quantum dot is used to partially diagonalize the Hamiltonian into subblocks according to their irreducible representations. This reduces the dimension of the Hamiltonian to $\sim 1000$. The reduced Hamiltonian is then diagonalized directly. Density of states of the nanocrystals are calculated and their evolution is traced from five-atom small clusters to 3019-Si atom quantum dot. It was found that for the 3019-Si quantum dot, the density of the state is very close to that of a bulk state.

The similar TB method has been used by Hill and Whaley [72-74] to calculate even larger quantum dots. For large systems, instead of diagonalizing the Hamiltonian directly, they have used a time evolution method. In this method, an initial random wavefunction $\psi(0)$ is evolved with time based on the time-dependent Schrodinger's equation:

$$\psi(t) = e^{-iHt}\psi(0) \tag{7}$$

here $H$ is the TB Hamiltonian. Numerically, Eq. (7) is evaluated by separating the time $t$ into $N$ steps $dt$: $e^{-iHt} = [e^{-iHdt}]^N$, and expanding the $e^{-iHdt}$ using Troter expansion [75] and split operator form [76]. After $\psi(t)$ is obtained, the single-particle density of state $n(E)$ can be derived from the Fourier transform of $\psi(t)$. More specifically, we have

$$n(E) = \frac{1}{\pi}\mathrm{Re}\int_0^\infty e^{iEt}\langle\psi(0)|\psi(t)\rangle\,dt \tag{8}$$

For a finite system, and for a long time run, one can find peaks in $n(E)$ indicating individual states. If the eigen-state energy $E_j$ is known for a given state, then its wavefunction can be produced from the Fourier transform of $\psi(t)$

$$\psi(\mathbf{r}, E_j) = \frac{1}{2\pi}\int_{-\infty}^{+\infty} \psi(\mathbf{r}, t)e^{iE_jt}\,dt \tag{9}$$

This time-dependent approach can also be used to calculate two particle states of an exciton [77] and for optical absorptions [78]. There, the Hamiltonian is the two-particle Hamiltonian with a Coulomb interaction between the electron and the hole, and the wavefunction is a two-particle wavefunction with one basis index for the electron and one index for the hole. Quantum dots with 147 atoms have been calculated in this way [77, 78].

Both Ren and Dow [71] and Hill and Whaley's [74] work have used the TB parameters from Ref. [79]. Using this TB model, Hill and Whaley [74] have calculated the Si optical gap as a function of the quantum dot size; the results seemed agree well with the experiment. However, it was later pointed out that this agreement with the experiment is fortuitous [80]. The problem is that the $sp^3s^*$ basis with nearest-neighbor interaction is not flexible enough to describe accurately the Si bulk conduction band. Extreme care must be taken when using these empirical TB methods. Delerue, Allan, and Lannoo (DAL) [81] have used a TB model with third-nearest-neighbor interactions and an orthogonal basis set. This model describes the conduction band well, and the results agree with more reliable pseudopotential methods. As in the work of Ren, DAL [81] also used symmetry to reduce the dimension of the TB Hamiltonian matrix and diagonalize the reduced matrix directly. They have calculated Si nanocrystallites and different shapes of wires. Figure 11 [81] shows these results.

Many works have been published by the DAL group using the third-nearest-neighbor TB model, which is taken from Ref. [82]. However, recently, Niquet and DAL [63] pointed out that even this TB Hamiltonian has its problems. Its valence-band effective mass parameters $\gamma_2 = 1.233$ and $m_l^* = 0.567$ are far from the experimentally measured values of $\gamma_2 = 0.320$, $m_l^* = 0.916$. A more rigorous fitting procedure is proposed to fit the TB band structure to GW-calculated band structures throughout the whole Brillouin zone. In total, there are 21 parameters for this new third-nearest-neighbor TB model.

In the TB calculations, the surface of a Si nanocrystal is often passivated by hydrogen atoms. The TB nearest-neighbor matrix elements $V_{H-Si}$ between H and Si can be scaled from

**Figure 11.** Calculated optical band-gap energies for various silicon crystallites (+) or wires (100, ×; 111, *; 111, ○) with respect to their diameter $d$. The black dots and squares are the experimental result of Ref. [37]. The dashed line is the band-gap energy for the crystallites including the Coulomb interaction between the electron and the hole. Reprinted with permission from [81]. C. Delerue et al., *Phys. Rev. B* 48, 11024 (1993). © 1993, The American Physical Society.

the Si–Si matrix elements $V_{\text{Si-Si}}$ according to–Harrison's rule [83]: $V_{\text{II-Si}} = V_{\text{Si-Si}}(d_{\text{Si-Si}}/d_{\text{II-Si}})^2$; here, $d_{\text{Si-Si}}$ and $d_{\text{II Si}}$ are the bond distances [71]. Another way to treat the surface passivation is simply to remove the dangling bond states from the calculated results. These dangling bond states can also be depleted from the TB Hamiltonian before the matrix is diagonalized. This is done by removing the hybrid $sp^3$ dangling bond orbital from the TB Hamiltonian basis set (e.g., remove the Hamiltonian matrix columns and rows expanded by these $sp^3$ basis) [73]. This is an unique artificial passivation only applicable to TB calculations. The ability to describe the surface atomistically is one big advantage of the TB model compared with the continuum effective mass k.p model. Using TB Hamiltonian, nanocrystals with different surface passivations have been studied by Hill and Whaley [73].

## 2.5. Empirical Pseudopotential Method

One problem of the above TB method is its lack of explicit basis functions. The $s$, $p$, and $s^*$ orbital bases are never explicitly assumed. This causes problems to calculate physical properties like dipole transitions and Coulomb and exchange interactions. Although explicit basis functions can be added after the TB eigen states have been calculated, these basis functions are not an intrinsic part of the TB Hamiltonian and its fitting process; thus, their compatibility is a problem. One way to overcome this problem is to use EPM. In EPM, the wavefunction is explicitly expressed by planewave basis, and its formalism is the same as the modern ab initio calculations. This makes it easy to be improved by the ab initio calculations, as in the semiempirical pseudopotential approach and the charge-patching method, which will be discussed later.

The EPM method was first used in 1960s by Cohen et al. [84, 85] to fit the semiconductor band structures. Even today, for many semiconductor systems, it still gives the best available band structures. In the EPM method, the wavefunction is expanded by a planewave basis set

$$\psi_i(\mathbf{r}) = \sum_q C_i(\mathbf{q})e^{i\mathbf{q}\cdot\mathbf{r}} \tag{10}$$

Usually, the planewave reciprocal lattice vector $\mathbf{q}$ is chosen inside a sphere of cutoff energy $E_{\text{cut}}$. The EPM Schrodinger's equation is

$$H\psi_i(\mathbf{r}) = \left[-\frac{1}{2}\nabla^2 + V(\mathbf{r})\right]\psi_i(\mathbf{r}) = E_i\psi_i(\mathbf{r}) \tag{11}$$

and

$$V(\mathbf{r}) = \sum_{atom} v_{atom}(|\mathbf{r} - \mathbf{R}_{atom}|) \tag{12}$$

here, $\mathbf{R}_{atom}$ are the atomic positions and $v_{atom}(r)$s are the spherical EPM atomic potentials. To variationally change $C_i(\mathbf{q})$ of Eq. (10) when solving the Schrodinger's equation $H\psi_i = E_i\psi_i$, it is equivalent to diagonalize the matrix $\langle e^{-i\mathbf{q}_1 \cdot \mathbf{r}}|H|e^{i\mathbf{q}_2 \cdot \mathbf{r}}\rangle \equiv \langle \mathbf{q}_1|H|\mathbf{q}_2\rangle$. This requires the evalution of $V(\mathbf{q}_1 - \mathbf{q}_2)$, which is the Fourier transformation of $V(\mathbf{r})$ in Eq. (12). Note that $V(\mathbf{q}) = \sum_{atom} v_{atom}(|\mathbf{q}|)e^{i\mathbf{q}\cdot\mathbf{R}_{atom}}$. In the original EPM procedure for bulk crystals, typically only three reciprocal vector values of $v_{atom}(|\mathbf{q}|)$ are assumed to be nonzero, and all the higher $|\mathbf{q}|$ values of $v_{atom}(|\mathbf{q}|)$ are assumed to be zero. These few nonzero $v_{atom}(|\mathbf{q}|)$ values are used as fitting parameters to fit the experimental bulk band structure at various special k points.

This turned out to be an extremely successful procedure. Surprisely, most of the II–VI, III–V, and VI–VI semiconductor band structures can be fitted well using this simple procedure. This means that the single particle description is indeed a good picture for semiconductor band structure and a local potential $V(\mathbf{r})$ can be used as a mean field to represent the effects of the complex many-body electron-electron interactions in a crystal.

To extend this traditional EPM approach to nanostructures, we need two improvements.

The first improvement is to have a continuous $v_{atom}(q)$ curve. For a nanostructure, the supercell is very large, and hence the reciprocal lattice vector $\mathbf{q}$ is much more dense than in a crystal. As a result, we need a continuous $v_{atom}(q)$ curve, not just a few discrete points. We also need to fit the pseudopotential for surface passivating atoms like H. We will fit a continuous $v_{atom}(q)$ under a certain form to a series of experimental data and to detailed first-principles calculations on relevant prototype systems. This will include bulk-band structures, clean surface work function, and the density of states of chemisorbed surfaces. Unlike the case in TB approaches, we will be able to compare the ensuing potential $V(\mathbf{r})$ with screened first-principles local density approximation (LDA) results. This gives us the ability to improve $V(\mathbf{r})$ in the future.

The second improvement is to solve Eq. (11). Even for a relatively small kinetic energy cutoff $E_{cut}$ for the planewave basis set, there could be roughly 50 planewaves for each atom. That translates into $\sim 50,000$ basis for a $\sim 1000$-atom system. Very often, we want to solve systems as large as 10,000 atoms. As a result, the direct diagonalization method based on $\langle \mathbf{q}_1|H|\mathbf{q}_2\rangle$ can no longer be used. Even the conventional conjugate gradient method [86] that is often used in *ab initio* calculations cannot be used because it scales as $O(N^3)$. However, we may not need all the eigenstates. For example, to study, the threshold optical properties of semiconductor quantum structures, what do we need are the eigenvalues and eigenfunctions of the band edge states, the total and local electronic density of states, and the optical absorption spectra. With these three properties calculated, most of the optical properties of the system can be determined. The "folded spectrum method" (FSM) has been developed by Wang and Zunger [87] to calculate the band edge states, and the "generalized moments method" (GMM) has been developed by Wang [88] to calculate the density of states and optical absorption spectra. We will discuss the FSM in this section and the GMM in a later section.

### 2.5.1. Constructing the EPM Hamiltonian

The continuous Si local pseudopotential in reciprocal space is presented in the following form [65]

$$v_{Si}(q) = a_1(q^2 - a_2)/(a_3 e^{a_4 q^2} - 1) \tag{13}$$

The coefficients were fitted to the bulk-band structure at high-symmetry points [89–93], the effective masses [94, 95], and the surface work function [96]. The bulk-band structure was calculated in a plane wave basis with an energy cutoff of 4.5 Ry (the same cutoff is used in subsequent calculations) and a lattice constant of 5.43 Å. The fit gave $a_1 = 0.2685$, $a_2 = 2.19$,

**Table 2.** Comparison of the empirical pseudo potential method (EPM) bulk Si band structures and effective masses, and experimental results.

| Property | EPM | Experiment |
|---|---|---|
| $\Gamma_{25'v}$ | 0 | 0 |
| $\Gamma_{1c}$ | −12.57 | −12.5(6)[a] |
| $\Gamma_{15c}$ | 3.24 | 3.35(1)[i] |
| $\Gamma_{2'c}$ | 4.12 | 4.15(5)[b] |
| $L_{2'v}$ | −10.19 | −9.3(4)[c] |
| $L_{1v}$ | −7.25 | −6.8(2)[c] |
| $L_{3'v}$ | −1.28 | −1.2(2)[b] |
| $L_{1c}$ | 2.18 | 2.04(6)[c] |
| $L_{3c}$ | 4.02 | 3.9(1)[b] |
| $X_{4v}$ | −3.01 | −2.9[b] |
| $X_{1c}$ | 1.32 | 1.13(?)[a] |
| $\Sigma_{min}$ | 4.47 | −4.48[b] |
| $E_{gap}$ | 1.167 | 1.124[d] |
| $H$ | 4.96 | 4.9[h] |
| $m_l(e)$ | 0.928 | 0.916[f] |
| $m_t(e)$ | 0.199 | 0.19[f] |
| $m^{(2)}_{1,X}(h)$ | 0.272 | 0.34[g] |
| $m^{(1)}_{1,X}(h)$ | 0.168 | 0.15[g] |
| $m^{(2)}_{1,l}(h)$ | 0.669 | 0.69[g] |
| $m^{(1)}_{1,l}(h)$ | 0.098 | 0.11[g] |

*Note:* The numbers in the bracket of the experimental data indicate the estimated error in the last digit. Both $m^{(1)}_{1,X}(h)$ and $m^{(1)}_{1,l}(h)$ stand for the non-spin-coupled effective hole mass [defined as $(\hbar k)^2/2\Delta E$] in the $\Gamma$ $X$ and $\Gamma$ $L$ directions, where $i$ denotes the band degeneracy. The variable is the work function. Energies are in electron volts, and effective masses are in the unit of electron mass.

*Source:* Reprinted with permission from [97]. L. W. Wang et al. "Semiconductor Nanoclusters," Studies in Surface Science and Catalysis, Vol. 103, p.161, Elsevier, 1996 + 1996, Elsevier Science.
[a] from Ref. [89], [b] from Ref. [90], [c] from Ref. [91], [d] from Ref. [92], [e] from Ref. [93], [f] from Ref. [94], [g] from Ref. [95], [h] from Ref. [96].

$a_3 = 2.06$, and $a_4 = 0.487$ in atomic units (Hartree for energy, inverse Bohr for $q$). Table 2 [97] compares the fitted quantities with the experimentally measured ones [89-96]. We can see that the fitted-band energies are within 0.1 eV of the experimental data, which is similar to the experimental uncertainty.

Figure 12 [97] compares the current atomic Si pseudopotential $v_{Si}(q)$ with the Fourier transform of the (self consistently) screened local LDA pseudopotential [98]. The closeness



**Figure 12.** Comparison between the Si empirical pseudopotential $v_{atom}(q)$ and the first principle LDA local pseudopotential for bulk Si. The first principle potential $V(G)$ is decomposed into atomic potentials according to $\sum_R \exp(iR \cdot G)v_{atom}(G) = V(G)$, where $G$ is the bulk reciprocal lattice vector and $R$ is the Si atomic position. Reprinted with permission from [97]. L. W. Wang and A. Zunger, in "Semiconductor Nanoclusters" (P. V. Kamat and D. Meisel, Eds.), Studies in Surface Science and Catalysis, Vol. 103, p. 161, Elsevier Science, Amsterdam, 1996. © 1996, Elsevier Science.

of the EPM potential and the LDA potential ensures that the EPM potential is realistic. However, the small difference between the empirical pseudopotential and the LDA potential reflects the fact that the EPM potential corrects the LDA band-gap error [99, 100].

Before the fitting of the hydrogen potential, we must find way to determine the atomic positions and geometries of the surface hydrogen atoms. In the work of Wang and Zunger [97], the quantum dot surface atomic configuration is modeled following H-covered Si thin-film (100), (110), (111) surfaces. These surface reconstructions are well studied both experimentally and theoretically [101–104]. One example is given in Fig. 13 [104], showing the (001) surface with one surface Si atom being passivated by two hydrogen atoms. It turns out that all the surface atoms of a convex quantum dot can be considered as in one of those three flat surfaces.

After the surface atomic positions are obtained, the surface hydrogen pseudopotential is fitted to the surface local density of states (LDOS) of the above three H-covered flat Si surfaces [(100), (110), (111)]. The important point here is to get the correct energies of the surface Si–H bonds. There are ultraviolet photoemission spectroscopy [105, 106] (UPS) and angle-resolved electron-energy-loss spectroscopy [107] (AR-EELS) measurements for the LDOS of these H-passivated surfaces. These experiments indicate that the bonding Si–H states are located around $E_v - 5$ eV, where $E_v$ is the valence-band maximum. For the energy positions of the conduction band Si–H antibonding states, *ab initio* calculations results can be used. These calculations indicate that the unoccupied-state antibonding surface state should be around at $\geq E_c + 1$ eV. Figure 14 shows the fitted LDOS for the three surfaces. The resulting $v_H(q)$ is

$$v_H(q) = -0.1416 + 9.802 \times 10^{-3} q + 6.231 \times 10^{-2} q^2 - 1.895 \times 10^{-2} q^3 \quad \text{when } q \leq 2$$

$$= 2.898 \times 10^{-2}/q - 0.3877/q^2 + 0.9692/q^3 - 1.022/q^4 \quad \text{when } q > 2 \tag{14}$$

Figure 15 shows the contour plots of the real-space empirical pseudopotential potential $V(\mathbf{r})$ produced from Eq. (12) of a H-covered (100) Si film and the total screened potential from a selfconsistent LDA calculation. The two potentials are very close, showing the transferability of the empirical pseudopotential $V(\mathbf{r})$.

The above empirical pseudopotentials are fitted solely from band structures and density of states. The *ab initio* calculations are used for their band energies, but not for their screened potentials. Another approach is to use the LDA potential $V_{LDA}(\mathbf{r})$ directly and to try to use Eq. (12) to fit this LDA potential. Using Eq. (12), the $v_{atom}(q)$ for the discrete reciprocal lattice vector $\mathbf{q}$ can be solved point by point for different crystal structures if the $V_{LDA}(\mathbf{q})$ is not zero. As a result, the so-obtained $v_{atom}(q)$ can be plotted as a function of $q$, as shown in Fig. 16 [108]. Turns out that this $v_{atom}(q)$ obtained from different crystal structures falls nicely into a single curve and can be represented by a smooth form $v_{atom}^{LDA}(q)$. This $v_{atom}^{LDA}(q)$ is called spherical approximation of the LDA potential. It was found that [108] the spherical approximation can reproduce the original band structure within 0.1 eV. After this step, $v_{atom}^{LDA}(q)$ is modified slightly to correct the LDA band-gap error. The final result $v_{atom}^{SLPM}(q)$ is



**Figure 13.** The atomic structure of calculated canted dihydride (001) Si surface viewed from (110). Reprinted with permission from [104], J. E. Northrup, *Phys. Rev. B* 44, 1419 (1991) © 1991, The American Physical Society.

Figure 14. Surface local density of states of three primary H covered silicon surfaces calculated using empirical pseudopotentials of Si and H. The vertical arrows indicate the energies of the Si–H bonding and antibonding states. The dashed vertical lines show the bulk valance-band maximum and conduction-band minimum (CBM). Reprinted with permission from [97], L. W. Wang and A. Zunger, in "Semiconductor Nanoclusters" (P. V. Kamat and D. Meisel, Eds.), Studies in Surface Science and Catalysis, Vol. 103, p. 161. Elsevier Science, 1996. © 1996, Elsevier Science.

called semiempirical pseudopotential (SEPM) [108]. The advantage of this SEPM is that its potential is constructed to be close to the LDA results and, consequently, its wavefunction has a 99% overlap with the original LDA wavefunction. This SEPM approach has been applied to CdSe [108], InP [109], and Si [108] systems, representing II–VI, III–V, and IV–IV semiconductor systems, respectively.

### 2.5.2. Solving the EPM Eigenstates

After the EPM Hamiltonian H is known, the next step is to solve the single-particle wavefunction $\psi_i(r)$ from Schrodinger Eq. (11). In the early study, the methods in the total energy calculations are used to calculate the eigenstates. This method is used to study 100-atom quantum wire systems. One of the results for different quantum wire states is shown in Fig. 17 [110]. However, for large systems, the methods used in the total energy calculations scale as $O(N^3)$; faster methods have to be used.

The method in total energy calculation [86] is to minimize the energy $\langle \psi | \hat{H} | \psi \rangle$ by varying the expansion coefficients $C(q)$ of $\psi$ in Eq. (10). The first $\psi$ obtained this way is the lowest-energy eigenstate of $\hat{H}$. To find a higher state, one needs to orthogonalize $\psi$ to all previously converged energy eigenstates. The enforcement of this orthogonalization scales as $N^3$.

The key to avoiding this problem is to solve only a few states near the band gap by realizing that information of these band-edge states is enough to determine many optical properties of the nanosystems. To solve the interior eigenstates near-energy $E_{ref}$ without solving all the other states below $E_{ref}$, the eigen energy spectrum of $H$ is first folded into the spectrum of

**Figure 15.** The contour plots of the (110) cross section of the total potential of a H-covered (001) Si slab. The contour level interval is 0.25 Hartree. Reprinted with permission from [97], L. W. Wang and A. Zunger, in "Semiconductor Nanoclusters" (P. V. Kamat and D. Meisel, Eds.), Studies in Surface Science and Catalysis, Vol. 103, p. 161. Elsevier Science, 1996. © 1996, Elsevier Science.



**Figure 16.** The spherical local density approximation (SLDA) potential $v_{SLDA}(|G|)$ as obtained from self-consistent bulk LDA calculations of five crystal structures and cell volumes. Diamond symbols represent the results for individual $|G|$. Solid lines represent least square fits of all the diamond symbols. Dashed lines represent the empirically adjusted potential to fit the experimental excitations. Reprinted with permission from [108]), L. W. Wang and A. Zunger, Phys. Rev. B 51, 17398 (1995). © 1995, The American Physical Society.

**Figure 17.** (001)-direction averaged wave-function squares, energy separations, and lifetimes for the near band-gap states. The dots denote positions of outer Si atoms. CBM stands for conduction band minimum. VBM stands for valence band maximum. Reprinted with permission from [110], C. Y. Yeh et al., *Appl. Phys. Lett.* 63, 3455 (1993). © 1993, The American Institute of Physics.

$(H - E_{ref})^2$. This is shown schematically in Fig. 18. If the reference energy $E_{ref}$ is within the band gap, then either the CBM (conduction band minimum) or the VBM (valence band maximum) will be the first eigen state in the spectrum of $(H - E_{ref})^2$. Notice that the eigen state $\psi_i$ satisfying Eq. (11) also satisfies

$$(H - E_{ref})^2\psi_i = \left[-\frac{1}{2}\nabla^2 + V(\mathbf{r}) - E_{ref}\right]^2 = (E_i - E_{ref})^2\psi_i \qquad (15)$$



**Figure 18.** A schematic view of folding the spectrum $\{\epsilon_i\}$ to spectrum $\{(\epsilon_i - \epsilon_{ref})^2\}$. CBM, conduction-band minimum; VBM, valance-band maximum. Reprinted with permission from [97], L. W. Wang and A. Zunger, in "Semiconductor Nanoclusters" (P. V. Kamat, D. Meisel, Eds.). Studies in Surface Science and Catalysis, Vol. 103, p. 161, Elsevier, 1996. © 1996, Elsevier Science.

Because numerically, $H$ and $(H - E_{ref})^2$ has the same matrix dimension, that means they have the same number of eigenstates. Then, an eigenstate of $(H - E_{ref})^2$ must also be an eigenstate of $H$. Thus, we can solve the eigenstates of $(H - E_{ref})^2$ to get the interior eigenstates of $H$. To solve the lowest eigenstates of $(H - E_{ref})^2$, the traditional conjugate gradient method is used to minimize $F = \langle \psi|(H - \epsilon_{ref})^2|\psi\rangle$. However, comparing this calculation with the minimization of $\langle \psi|H|\psi\rangle$, the use of $(H - \epsilon_{ref})^2$ considerably slows down the convergence of the standard minimization methods. This problem is solved here by using a preconditioned conjugate gradient method with a large number of conjugate gradient steps. To calculate $F$, we apply twice $[-\frac{1}{2}\nabla^2 + V(\mathbf{r}) - \epsilon_{ref}]$ to $\psi(\mathbf{r}) = \sum_q C(\mathbf{q})e^{i\mathbf{q}\mathbf{r}}$. The term $-\frac{1}{2}\nabla^2\psi$ is computed in reciprocal space, and $V(\mathbf{r})\psi(\mathbf{r})$ is obtained by using the fast Fourier transformation (FFT) to transform $C(\mathbf{q})$ to real space $\psi(\mathbf{r})$, then applying $V(\mathbf{r})$ to $\psi(\mathbf{r})$ and transforming the product back to $\mathbf{q}$ space. The detailed of this procedure is described in Ref. [111]. A parallel code called PESCAN has been developed that can be routinely used to calculate systems for a few thousand atoms, or even nearly a million atom systems [112]. Because only a few wavefunctions are calculated, the computational effort scales linearly to the size $N$ of the system. Recently, the folded spectrum method has also been used to solve the TB Hamiltonian [113].

## 2.6. The Local Density Approximation Method

The above empirical pseudopotential is not solved self-consistently. This might be a problem for the quantum dot surface, especially when different surface passivations and surface states need to be calculated. It could also be a problem for impurity and semiconductor alloys. These problems can be solved by the self-consistent *ab initio* methods like the density-functional theory (DFT) [114, 115], and its LDA [115].

The general flow chart of a LDA calculation is described in Fig. 19. Basically, $M$ lowest-energy occupied states of Eq. (11) are solved using methods like the conjugate gradient method [86]. Then the total charge density $\rho(\mathbf{r})$ of the system is given by

$$\rho(\mathbf{r}) = \sum_{i=1}^{M} |\psi_i(\mathbf{r})|^2 \tag{16}$$



Figure 19. A flow-chart of the self-consistent local density approximation calculation.

and the potential $V(\mathbf{r})$ in Eq. 11 is then calculated as

$$V(\mathbf{r}) = \sum_{atom} \hat{v}_{bare}(\mathbf{r} - \mathbf{R}_{atom}) + \int \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3\mathbf{r}' + \mu_{xc}^{LDA}[\rho(\mathbf{r})] \tag{17}$$

here $\hat{v}_{bare}(\mathbf{r} - \mathbf{R}_{atom})$ is the bare, unscreened LDA pseudopotential, which might contain nonlocal parts [116], and $\mu_{xc}^{LDA}(\rho)$ is the LDA exchange-correlation function. In the self-consistent LDA calculation, the $V(\mathbf{r})$ calculated from Eq. (17) is taken back into Eq. (11), and the $\psi_i(\mathbf{r})$ are recalculated, until the input and output potential $V(\mathbf{r})$ become the same (self-consistent). The details of this procedure can be found in Ref. [117].

The biggest computational challenge of the self-consistent LDA method is its $O(N^3)$ scaling. However, using a limited linear combination of atomic orbital (LCAO) basis, Delley and Steigmeier [118] were able to calculate 700-atom Si quantum dots. They get a band-gap increase that scales like $1/R$. Another improvement came when Chelikowsky, Troubllier, and Saad [119] developed a real space finite difference method to calculate the wavefunction of Eq. (11). Compared to planewave basis expansion of Eq. (10), the real space calculation does not need the fast Fourier transformation, which is the most time-consuming part of the computation. Another advantage of the real space calculation is that it can be parallelized naturally, with each computer processor calculating one real space part of the grid. Using this method, Ogut, Chelikowsky, and Louie [120] have calculated 800 atom Si quantum dots using the LDA method. Although there is no direct comparison between the real space method and the more traditional and mature planewave method, it is believed that the real space method is probably more accurate than the limited-basis LCAO method.

Although the LCAO and the real space methods, combined with large-scale supercomputers, can be used to calculate thousand-atom systems, these methods still scale as $O(N^3)$. It is still difficult to calculate even larger systems or other materials with $d$ electron states. One well-studied approach is the self-consistent $O(N)$ method. A detailed review about this method is given by Goedecker [121]. Unfortunately, even after a decade of research, there is still no reliable $O(N)$ ab initio calculations and robust computer codes, except for the strongest covalent systems like graphite [122]. Thus, an alternative approach is needed to make LDA-type calculations for thousand-atom systems easy and practical.

One of such alternative approaches is the charge patching method (CPM) [123–125]. The basic assumption of the CPM is that the charge density at a given point depends only on the atomic arrangement around this point. This is true if there is no long-range external electric field and if there is a band gap in the material. On the basis of this assumption, the idea is to calculate the charge densities from some small prototype systems and then transfer the charge densities from these systems to get the charge density of a given large nanosystem.

More specifically, charge density motifs are calculated from the charge densities of the prototype systems as

$$m_{I_\alpha}(\mathbf{r} - \mathbf{R}_\alpha) = \rho_{LDA}(\mathbf{r}) \frac{w_\alpha(|\mathbf{r} - \mathbf{R}_\alpha|)}{\sum_{\mathbf{R}_\alpha} w_\alpha(|\mathbf{r} - \mathbf{R}_\alpha|)} \tag{18}$$

here $\mathbf{R}_\alpha$ is a atomic site of atom typed $\alpha$, and $m_{I_\alpha}(\mathbf{r} - \mathbf{R}_\alpha)$ is the charge density motif belonging to this atomic site, and $\rho_{LDA}(\mathbf{r})$ is the self-consistently-calculated charge density of a prototype system. The variable $w_\alpha(r)$ is an exponential decay function; thus, $w_\alpha(|\mathbf{r} - \mathbf{R}_\alpha|)/\sum_{\mathbf{R}_\alpha} w_\alpha(|\mathbf{r} - \mathbf{R}_\alpha|)$ is a distribution function that divides the space into regions belonging to each atom. Note that $m_{I_\alpha}(\mathbf{r} - \mathbf{R}_\alpha)$ is a localized function and, hence, can be stored in a fixed-size numerical array. The $I_\alpha$ in $m_{I_\alpha}$ is used to denote the atomic-bonding environment of the atom $\alpha$ at $\mathbf{R}_\alpha$.

After $m_{I_\alpha}$ for all the possible bonding environments $I_\alpha$ is obtained, the charge density of a given nanosystem can be generated by patching the charge motifs together

$$\rho_{patch}(\mathbf{r}) = \sum_{\mathbf{R}_\alpha} m_{I_\alpha}(\mathbf{r} - \mathbf{R}_\alpha) \tag{19}$$

here the atomic bonding environment $I_\alpha$ of atom $\alpha$ at $\mathbf{R}_\alpha$ should be the same as in Eq. (18).

This charge-patching method can be used to generate the charge densities of carbon fullerenes [124], semiconductor alloys [125], semiconductor impurities [123], and semiconductor quantum dots. The resulting patched charge density is typically within 1% of the self-consistently-calculated LDA charge density, and the resulting single-particle eigenstate energies are within 50 meV of the direct self-consistent calculation results. Consider that the typical numerical uncertainty (due to basis function truncations, different nonlocal pseudopotential treatments, and different nonlocal pseudopotentials) of a LDA calculation is about 50 meV, this charge-patching method can be considered as accurate as the direct *ab initio* calculations. After the charge density is generated, Eq. (17) can be used to generate the LDA potential $V(r)$, and Khon–Sham Eq. (11) can be solved by the folded spectrum method for a few band edge states. Figure 20 shows the VBM and CBM states of a 1000-atom Si quantum dot calculated by the charge-patching method. Such calculations take only 1 hour on a 64-processor IBM SP machine. If direct LDA calculation for the same system was attempted, it would take a few weeks.

One problem of the LDA calculation is that its band gap is severely underestimated [99, 100]. Strictly speaking, the density-functional theory is only valid for ground-state properties, and there is no physical meaning for the Kohn–Sham eigen energies [115]. One way to circumvent this conceptional difficulty is to use time-dependent DFT, which will be discussed later. In practices, LDA calculation does give reliable results for quantum dot confinement energies (the band-gap increase as a function of the quantum dot size). It was found that the Si LDA quantum confinement energy is practically the same as the confinement energy, as calculated from the empirical pseudopotential results and the third-nearest-neighbor tight-binding results, as shown in Fig. 21 [69, 126–129]. One practical way to correct the LDA band-gap error is to modify the LDA potential $V(r)$ slightly. One way to do this is the Kristenson method, which puts a small spherical potential at the crystal interstitial region. Another way is to modify the $s$, $p$, and $d$ nonlocal pseudopotentials [123].

## 3. MANY-BODY METHODOLOGIES

All the methods in Section 2 solve the single-particle eigenstates and energies. To calculate the optical transitions, usually simple zero-order electron-hole interactions are added on top of the single-particle eigen energies. However, there are cases in which the many-body effects might be important. In this section, we will discuss a few methods that either treat many-body effects explicitly (configuration interaction method), that are derived from many-body arguments and theories (e.g., time-dependent LDA [TDLDA], GW + Bethe–Salpeter equation), or that directly treat the nanosystem as a whole, using many-body descriptions (quantum Monte Carlo method).



(a)CBM                    (b)VBM

$Si_{1327}$ quantum dot

**Figure 20.** The conduction-band minimum and valance-band maximum of a 1000-Si-atom quantum dot passivated with H atoms. The charge density is generated with the charge-patching method. The conduction-band minimum and valance-band maximum states are calculated with the folded-spectrum method.

**Figure 21.** Energy gap versus $1/d$ for H-terminated Si dots, wires, and slabs. Local density approximation. (filled and empty dots) Ref. [126]. (+) Ref. [127]. ($\triangle$) Ref. [128]. (*) Ref. [129]. Empirical pseudopotential: (×) Ref. [69]. Reprinted with permission from [126]. B. Delley and E. F. Steigmeier, *Appl. Phys. Lett.* 67, 2370 (1995). © 1995, The American Institute of Physics.

## 3.1. Time-Dependent Density Functional Theory Methods

As discussed above, rigorously speaking, DFT only works for ground states. One way to study the excited state is to study the response of the system under time-dependent external electric field perturbation. Rigorous theory [130–132] can be derived similar to the time-independent DFT, which relates the true many-body time-dependent response to the Kohn–Sham noninteractive system response. Under the time-dependent DFT (TDDFT), one solves the time-dependent Kohn–Sham equation:

$$\left[-\frac{1}{2}\nabla^2 + V(\mathbf{r}, t)\right]\psi_i(\mathbf{r}, t) = i\frac{\partial}{\partial t}\psi_i(\mathbf{r}, t)$$ (20)

and

$$\rho(\mathbf{r}, t) = \sum_{i=1}^{M}|\psi_i(\mathbf{r}, t)|^2$$ (21)

Under rigorous TDDFT, $V(\mathbf{r}, t)$ at time $t$ is a functional of the charge-density function $\rho(\mathbf{r}, t')$ for all $t' < t$. However, under LDA and adiabatic approximation, $V(\mathbf{r}, t)$ depends only on $\rho(\mathbf{r}, t)$ and the relationship between $V(\mathbf{r}, t)$ and $\rho(\mathbf{r}, t)$ is the same as in Eq. (17), but replacing $V(\mathbf{r})$ with $V(\mathbf{r}, t)$, $\rho(\mathbf{r})$ with $\rho(\mathbf{r}, t)$ for the same time $t$ and adding an external perturbation term $V_{\text{ext}}(\mathbf{r}, t)$. The absorption spectrum can be calculated from the charge-density response $\rho(\mathbf{r})$. This approximation is called TDLDA.

One way to numerically calculate TDLDA is to integrate explicitly Eq. (20) in time [133]. To do that, a step function in time is used for the external potential $V_{\text{ext}}(\mathbf{r}, t)$. This is rather like doing Eq. (7). However, unlike in Eq. (7), where only one single-particle wavefunction is integrated, in Eq. (20), all $M$ occupied states are integrated. The absorption spectra for all the energies are calculated from the time Fourier transform of $\rho(\mathbf{r}, t)$ and $V_{\text{ext}}(\mathbf{r}, t)$.

The direct time integration method is good for relatively large systems. However, for a small system, a more efficient way is to solve Eq. (20) in frequency space. To do that, one has to assume that the perturbation is small, so linear response theory can be applied. Using a linear response theory, an exciton energy $\omega$ can be solved from the following equation

$$\sum_{jl}\left[(\epsilon_j - \epsilon_k)\delta_{ij}\delta_{kl} + (f_i - f_k)K_{ik,jl}(\omega)\right]C_{jl} = \omega C_{ik}$$ (22)

here $\epsilon_i$ and $\epsilon_k$ are the LDA ground-state Kohn–Sham eigen energies, and $f_i$ and $f_k$ are the occupation numbers of Kohn–Sham eigenstates $\psi_i$ and $\psi_k$, respectively. For an adiabatic TDLDA approximation, $K_{ik,il}(\omega)$ is independent of $\omega$ as

$$K_{ik,il} = \iint d\mathbf{r}_1 d\mathbf{r}_2 \psi_i(\mathbf{r}_1)\psi_k(\mathbf{r}_1)\left[\frac{1}{|\mathbf{r}_1 - \mathbf{r}_2|} + \delta(\mathbf{r}_1 - \mathbf{r}_2)\frac{\partial \mu_{xc}^{LDA}\{\rho(\mathbf{r}_1)\}}{\partial \rho(\mathbf{r}_1)}\right]\psi_l(\mathbf{r}_2)\psi_j(\mathbf{r}_2) \qquad (23)$$

Very often, Eq. (22) is rewritten as a quadratic form for numerical stability [134]

$$\sum_{il}\left[\omega_{ik}^2 \delta_{ij}\delta_{kl} + 2\sqrt{f_{ik}\omega_{ik}}K_{ik,il}(\omega)\sqrt{f_{jl}\omega_{jl}}\right]C_{jl} = \omega^2 C_{ik} \qquad (24)$$

where $\omega_{ik} = \epsilon_i - \epsilon_k$, $f_{ik} = f_i - f_k$. In Eq. (23), the first term $\frac{1}{|\mathbf{r}_1 - \mathbf{r}_2|}$ can be called the exchange interaction [because $\psi_i(\mathbf{r}_1)$ and $\psi_k(\mathbf{r}_1)$ belong to valence and conduction bands, respectively]. The second term can be called the screened Coulomb interaction (this comes from the comparison to configuration interaction and the GW + Bethe–Salpeter equation methods). The Coulomb interaction here is no longer an nonlocal integral between $\mathbf{r}_1$ and $\mathbf{r}_2$. This is because in the LDA formalism, the screened exchange interaction term is written as a local exchange-correlation functional, instead of being an explicit integral based on $1/r$ interaction.

The TDLDA has been extensively used for optical spectra for molecules and small clusters. For these small systems, the TDLDA results often agree well with the experimental measurement, as shown in Fig. 22 [135]. Compared to LDA results (i.e., using $\epsilon_i - \epsilon_k$ as the transition energies), the TDLDA is much improved. Usually, some major LDA peaks are blue-shifted. This shift is often caused by the unscreened repulsive exchange interaction in Eq. (23). For extremely small systems like molecules, this exchange interaction is very strong, especially for high-energy $\{i, k\}$ pairs. Recently, Benedict et al. [135] found out that the screened Coulomb interaction in Eq. (23) does not play any important roles for small clusters. In addition, for the lowest excitation state (exciton energy), the TDLDA result is roughly the same as the LDA eigen energy difference (CBM–VBM). For bulk system, it is known that [136] the TDLDA band edge will be the same as the LDA band edge. Thus, the result does not improve the LDA band-gap error. Strictly speaking, this is only a problem of the TDLDA. For TDDFT, the band gap should be correct, even though the exact DFT Kohn–Sham eigen energy band gap could still be different from the true band gap. The TDLDA also does not improve the bulk optical absorption spectrum from the LDA



Figure 22. Calculated and measured optical spectra of SiH$_4$. TDLDA, time-dependent local density approximation; BSE, Bethe–Salpeter equation; LDA, local density approximation. Reprinted with permission from [135], L. X. Benedict et al., *Phys. Rev. B* 68, 85310 (2003). © 2003, The American Physical Society.

results, as shown in Fig. 23. All these problems come from the screened Coulomb interaction term in Eq. (23). For larger systems, the Coulomb interaction becomes important. An on site interaction, as in Eq. (23) cannot describe the correct long-range behavior of this interaction. For not very small nanosystems, the Coulomb interaction becomes larger than the exchange interaction: It cannot be ignored. Thus, for larger nanosystems and bulks, not only the TDLDA inherits the LDA band gap error but it also lacks the accurate screened Coulomb interaction. Note that the diagonal screened Coulomb interaction is negative. The lack of an effective Coulomb interaction will make the TDLDA have difficulty capturing the bound exciton effect, which causes red-shift of the peak weights in the optical absorption spectra (Fig. 23).

TDDFT can also be implemented by density-functional approximations other than the LDA approximation. One often-used density-functional approximation is the B3LYP functional [137]. This functional is a hybrid of exact Hartree–Fock exchange with local and gradient-corrected exchange and correlation terms. More exactly, the exchange and correlation function is expressed as:

$$E_{xc}^{B3LYP} = (1 - a_0)E_x^{LSDA} + a_0 \Delta E_x^{HF} + a_x \Delta E_x^{B88} + a_c E_c^{LYP} + (1 - a_c)E_c^{VWN} \qquad (25)$$

here $E_x^{LSDA}$ is the local spin density approximation (LSDA) exchange local density functional, $E_x^{HF}$ is the Hartree–Fock exchange energy, $\Delta E_x^{B88}$ is the Becke's gradient correction [138] to the exchange functional, $E_c^{LYP}$ is the Lee–Yang–Parr correlation functional [139], and $E_c^{VWN}$ is the Vosko, Wilk, and Nusair local correlation functional [140]. The variables $a_0$, $a_x$ and $a_c$ are fitting parameters. They are used to fit the properties of many small molecules. The numbers used for B3LYP are $a_0 = 0.2$, $a_x = 0.72$, and $a_c = 0.81$. Strictly speaking, the B3LYP is no longer a density-functional theory because the Hatree–Fock exchange energy $E_x^{HF}$ is expressed as a functional of the single-particle wavefunctions, not as a functional of the total charge density, and the Schrodinger's equations for the wavefunctions are not the local potential Kohn–Sham equations. B3LYP can be viewed as a hybrid method between the LDA (or GGA [generalized gradient approximation]) and the Hatree–Fock method. The B3LYP method gives quite good band gaps for various bulk crystals [141]. For bulk Si, the agreement with the experiment is about 0.3 eV for various special k-points, which is only slightly larger than the results of the GW and quantum Monte Carlo methods. Because it also has the explicit exchange integral, it will give a long-range Coulomb interaction in Eq. (22) and (23). Thus, the TDDFT–B3LYP could mitigate the two problems of TDLDA discussed above. In practice, as the HF exchange integral is used, B3LYP is often solved using Gaussian atomic orbital basis. Recently, TDDFT–B3LYP has been used to calculate



Figure 23. Solid line: calculated time-dependent local density approximation (TDLDA) result for bulk Si absorption spectrum. Dashed line: experiment. Reprinted with permission from [135], L. X. Benedict et al., *Phys. Rev. B* 68, 085310 (2003). © 2003 The American Physical Society.

small Si clusters, and the result is found to agree well with the multireference second-order perturbation theory (MR-MP2) [142], although a more recent study [143] found that the basis set used in Ref. [142] might not be converged.

## 3.2. Configuration Interaction Method

Another approach to calculate the many-body interaction is via the configuration interaction approach. This is used in conjunction with the single-particle methods described in Section 2. After the single-particle eigenstates in Eq. (11) are obtained, we can use them to form Slater determinants to solve the many-body configuration interaction (CI) Hamiltonian. More specifically, one can construct the ground-state Slater determinant as

$$\Phi_0(\mathbf{r}_1, \ldots, \mathbf{r}_M) = \mathcal{A}[\psi_1(\mathbf{r}_1) \ldots \psi_v(\mathbf{r}_v) \ldots \psi_M(\mathbf{r}_M)] \tag{26}$$

here $\psi_i(\mathbf{r})$ are the single-particle wavefunctions from Eq. (11), and $\mathcal{A}$ is the antisymmetrizing operator. A single-excitation Slater determinant that replaces a valence state $\psi_v$ with a conduction-band state $\psi_c$ is

$$\Phi_{v,c}(\mathbf{r}_1, \ldots, \mathbf{r}_M) = \mathcal{A}[\psi_1(\mathbf{r}_1) \ldots \psi_c(\mathbf{r}_v) \ldots \psi_M(\mathbf{r}_M)] \tag{27}$$

Then, the exciton wavefunction $\Psi$ can be expressed as a linear combination of $\Phi_{v,c}$ as

$$\Psi = \sum_{v=1}^{N_v} \sum_{c=1}^{N_c} C_{v,c} \Phi_{v,c} \tag{28}$$

The coefficients $C_{v,c}$ are obtained by diagonalizing the Hamiltonian matrix

$$\sum_{v'c'} H_{vc,v'c'} C_{v'c'} = \sum_{v'c'} [(E_c - E_v)\delta_{v,v'}\delta_{c,c'} + K_{vc,v'c'} - J_{vc,v'c'}]C_{v'c'} = EC_{vc} \tag{29}$$

here $E_v$ and $E_c$ are the single-particle eigen energies in Eq. (11), and $E$ is the exciton energy. The $K_{vc,v'c'}$ and $J_{vc,v'c'}$ are the exchange and Coulomb interactions, respectively

$$K_{vc,v'c'} = \iint \frac{\psi_c(\mathbf{r}_1)\psi_{v'}^*(\mathbf{r}_1)\psi_v(\mathbf{r}_2)\psi_{c'}^*(\mathbf{r}_2)}{\bar{\epsilon}(\mathbf{r}_1, \mathbf{r}_2)|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \tag{30}$$

$$J_{vc,v'c'} = \iint \frac{\psi_v(\mathbf{r}_1)\psi_{v'}^*(\mathbf{r}_1)\psi_{c'}(\mathbf{r}_2)\psi_c^*(\mathbf{r}_2)}{\bar{\epsilon}(\mathbf{r}_1, \mathbf{r}_2)|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \tag{31}$$

Equations (29), (30), and (31) can be derived by applying Eq. (27) to the original many-body Hamiltonian. However, in Eq. (30) and (31), we have used the dielectric screening for both the exchange and Coulomb interactions. For a bulk exciton, the screening of the Coulomb interaction can be derived from many-body theory [144, 145]. An alternative derivation is given by Strinati [146], using the GW formalism. Formally, the effective dielectric screening can be rewritten as

$$\frac{1}{\bar{\epsilon}(\mathbf{r}_1, \mathbf{r}_2)|\mathbf{r}_1 - \mathbf{r}_2|} = \int \epsilon_{\text{bulk}}^{-1}(\mathbf{r}_1, \mathbf{r}) \frac{1}{|\mathbf{r} - \mathbf{r}_2|} d\mathbf{r} \tag{32}$$

here $\epsilon_{\text{bulk}}^{-1}(\mathbf{r}_1, \mathbf{r})$ is the bulk inversion of the dielectric function $\epsilon(\mathbf{r}_1, \mathbf{r})$. This is different from $\epsilon^{-1}(\mathbf{r}_1, \mathbf{r})$, as the direct whole-space inversion of the $\epsilon(\mathbf{r}_1, \mathbf{r})$ matrix for the quantum dot system. The $\epsilon^{-1}(\mathbf{r}_1, \mathbf{r})$ contains a surface polarization potential $P$, as discussed in Section 2.1. The bulk inversion $\epsilon_{\text{bulk}}^{-1}(\mathbf{r}_1, \mathbf{r})$, however, has no such surface polarization potential. However, the use of this bulk $\epsilon_{\text{bulk}}^{-1}(\mathbf{r}_1, \mathbf{r})$ in Eq. (29) and (31) corrabates well with the fact that the single-particle energy $E_v$ and $E_c$ are derived from Eq. (11), which contains no surface polarization term $P$ from Eq. (1). If the surface polarization term $P$ is included in the single-particle equation, as in Eq. (1), then the full whole-space inverse dielectric function $\epsilon^{-1}(\mathbf{r}_1, \mathbf{r})$ should be used. The surface polarization terms will cancel out at the end. This is the situation in the GW + Bethe–Salpeter formalism, as will be discussed in the next section.

It is often argued that the exchange interaction $K_{\mathrm{r}c,\mathrm{r}c}$ should not be screened. This can be viewed from the two-particle Green's function construction, where screening of the exchange term will cause double counting [147]. Nevertheless, in practice, it is found that the exchange consisted of a long-range and a short-range term [148]. Although the short range term should be unscreened, the long-range term should be screened by the bulk dielectric function [149–151]. The use of the effective dielectric function $\bar{\epsilon}(r_1, r_2)$ in Eq. (30) does satisfy this condition because $\bar{\epsilon}(r_1, r_2) \to 1$ for $|r_1 - r_2| \to 0$. The problem of the contradiction to the Green's function arguments can be resolved by realizing that if only a limited configuration space is used in Eq. (29), then the effect of other unused configurations can be included in the exchange screening term [147].

Notice that the CI in Eq. (29) has the same form as in Eq. (22) for the TDLDA, although the expression for the exchange and the Coulomb interaction is different. CI method is often used in conjunction with the non-selfconsistent single-particle methods discussed in Section 2. CI method has been used to calculate very large systems, including pyramidal quantum dots with nearly one million atoms [152]. The key here is to use a limited window for the configurations. One can never afford to use the full configuration space for such large systems. Another advantage for the CI approach is that it can also be used to study other many-body problems, such as two electrons, or multi-excitons, and Auger effects [153]. It is difficult to study these systems by the TDLDA and GW + Bethe–Salpeter equation approaches. However, cautious must be used when dealing with the screening issues for those multiparticle excitations.

## 3.3. GW and Bethe–Salpeter Equation Approach

GW and Bethe–Salpeter equation have been used to study optical absorption spectra for small molecules and clusters. This approach can be separated into two steps. In the first step, the GW quasiparticle eigen energies are calculated. This is a single-particle step, rather like the methods in Section 2. In the second step, the Bethe–Salpeter equation is solved for electron–hole pair excitons, based on quasiparticle energies and wavefunctions.

Quasiparticle is defined as the poles in frequency space in the single-particle Green's function

$$G(rt, r't') = -i\langle M | T\hat{\psi}(rt)\hat{\psi}^{\dagger}(r't') | M \rangle \tag{33}$$

here $\hat{\psi}(rt)$ is the particle creation operator, $|M\rangle$ is the $M$ particle-ground state, and $T$ is the time-ordering operator. It can be shown that [154, 155] the pole energy $E_i$ corresponds to the energy for adding one electron into the system to state $i$, $\epsilon_i = E(M + 1, i) - E(M)$ (when $\epsilon_i$ is larger than the Fermi energy), or for taking one electron out from the system, $\epsilon_i = E(M) - E(M - 1, i)$ (when $\epsilon_i$ is less than the Fermi energy). The quasiparticle energy $\epsilon_i$ can be calculated with the GW approximation [156, 157] as

$$\left[ -\frac{1}{2}\nabla^2 + \sum_{\mathrm{atom}} \hat{v}_{\mathrm{bare}}(r - R_{\mathrm{atom}}) + \int \frac{\rho(r')}{|r - r'|} d^3r' \right] \psi_i(r) + \int \Sigma(r, r': \epsilon_i)\psi_i(r')dr' = \epsilon_i\psi_i(r) \tag{34}$$

Note that, compared with the LDA potential in Eq. (17), the LDA exchange-correlation potential $\mu_{xc}^{LDA}[\rho(r)]$ has been replaced by a nonlocal self-energy potential $\Sigma(r, r': \epsilon_i)$. This self-energy potential has the following GW expression

$$\Sigma(r, r', \omega) = -\sum_k \psi_k(r)\psi_k(r)\left[ f_k W(r, r', \epsilon_k - \omega) + \frac{1}{\pi} \int \frac{\mathrm{Im}\,W(r, r', \omega')}{\omega - \epsilon_k - \omega' + i\delta} d\omega' \right] \tag{35}$$

here the dynamically screened interaction is

$$W(r, r', \omega) = \int \epsilon^{-1}(r, r_1, \omega)\frac{1}{|r_1 - r'|} dr_1 \tag{36}$$

Here $\epsilon^{-1}(r, r_1, \omega)$ is the inversion of the dielectric matrix $\epsilon(r, r_1, \omega)$.

Supposely Eqs. (34) and (35) should be solved self-consistently. However, in reality [100], one usually uses the LDA Kohn–Sham wavefunctions and eigen energies for $\psi_k$ and $\epsilon_k$ in Eq. (35), and takes an expectation value $\langle \psi_k | \Sigma | \psi_k \rangle$ for the self-energy term, instead of solving Eq. (34) exactly. This is called zero-order approximation of the GW procedure. Recently, some self-consistent calculations [158–163] of Eqs. (34) and (35) showed that the self-consistency actually makes the spectrum properties worse (e.g. for quasiparticle excitation energy, bandwidth, lifetimes, and plasmon satellites). Thus, in practice, one should just take the zero-order results of the GW equation.

As for the quasiparticle energy being the pole of the single-particle Green's function of Eq. (33), the information about the exciton energies is contained in the two-particle Green's function:

$$G(r_1 t_1, r_2 t_2, r'_1 t'_1, r'_2 t'_2) = -\langle M | T\psi(r_1 t_1)\psi(r_2 t_2)\psi^\dagger(r'_2, t'_2)\psi^\dagger(r'_1, t'_1) | M \rangle \qquad (37)$$

By taking $t_1 = t'_1 + 0^-$, $t_2 = t'_2 + 0^-$, the two-particle Green's function can be transformed into frequency space as $G_2(\omega)$, and the exciton energies are the poles in $G_2(\omega)$. The Dyson's equation for this two-particle Green's function is [146, 155] (written in a concise way):

$$G_2(\omega) = G_2^{(0)}(\omega) + G_2^{(0)}(\omega)K'(\omega)G_2(\omega) \qquad (38)$$

here $G_2^{(0)}(\omega) = G_1 G_1$ is the noninteracting two-particle Green's function, and $K'(\omega)$ is an electron-hole interaction kernal. Equation (38) for electron-hole pairs is also called the Bethe–Salpeter equation [146]. It can be solved by expanding the exciton wavefunction $|M, S\rangle$, using quasiparticle electron-hole pairs

$$|M, S\rangle = \sum_v \sum_c C_{vc} \hat{a}_v^\dagger \hat{b}_c^\dagger | M \rangle \qquad (39)$$

Here $\hat{a}_v^\dagger$ creates a quasihole, and $\hat{b}_c^\dagger$ creates a quasi-electron. On the basis of Eq. (39), the eigen state equation for $C_{vc}$ is

$$(\epsilon_c - \epsilon_v)C_{vc} + \sum_{v'c'}(K_{vc,v'c'} - J_{vc,v'c'})C_{v'c'} = \Omega_S C_{vc} \qquad (40)$$

here $\epsilon_v$, $\epsilon_m$ are the quasi-particle eigen energies of Eq. (34), and $\Omega_S$ is the exciton energy. The $K$ and $J$ are the exchange and screened Coulomb interactions. Although the $K_{vc,v'c'}$ is the same as in Eq. (30) without the screening $\bar{\epsilon}(r_1, r_2)$ [or, say, the first term in Eq. (23)], the screened Coulomb interaction is

$$J_{vc,v'c'} = \int dr dr' \psi_c^*(r)\psi_{c'}(r)\psi_v(r')\psi_{v'}^*(r')\frac{i}{2\pi}\int d\omega e^{-i\omega 0^+} W(r, r', \omega)$$

$$\times \left[ (\Omega_S - \omega - (\epsilon_{c'} - \epsilon_v) + i0^+)^{-1} + (\Omega_S + \omega - (\epsilon_c - \epsilon_{v'}) + i0^+)^{-1} \right] \qquad (41)$$

where $W(r, r', \omega)$ is the screened Coulomb interaction given by Eq. (36).

GW plus the Bethe–Salpeter equation is considered one of the most accurate methods to calculate the optical absorption spectra and excited-state electronic structures. It has been used to calculate molecules and bulk crystals. Figure 24 [164] shows the GW–Bethe–Salpeter equation results for the optical absorption spectra for bulk Si. The agreement with the experiment is excellent. It is the first time the lower energy peak emerged in the calculated results. This peak has long been expected to be caused by excitonic effects. Without the electron-hole interaction in the Bethe–Salpeter equation, this peak does not exist. The GW–Bethe–Salpeter equation results can be compared with the TDLDA results shown in Fig. 23. The TDLDA result is very similar to the original LDA result. The problem is that the Coulomb interaction in TDLDA is a local term [the second term in Eq. (23)]. This local approximation is inadequate to describe the excitonic bounding effect.

The Bethe–Salpeter equation of Eq. (40) can be compared with the TDLDA of Eq. (22) and the CI of Eq. (29). These three equations are exactly the same, except that the meanings

**Figure 24.** Calculated optical absorption spectrum of Si with (solid lines) and without (dashed lines) electron-hole interaction using GW–Bethe–Salpeter equation, with three valence bands, six conduction bands, 500 k points in the BZ, and an artificial broadening of 0.15 eV. The open and filled circles are experimental data. Reprinted with permission from [164], M. Rohlfing and S. G. Louie, *Phys. Rev. B* 62, 4927 (2000). © 2000, The American Physical Society.

of the single-particle eigen energies $\epsilon_i$, the exchange interaction $K$, and the screened Coulomb interaction $J$. In the GW–Bethe–Salpeter equation approach, the quasi-particle energy $\epsilon_i$ in Eq. (34) can be defined as the total energy difference $E(M+1, i) - E(M)$ for conduction bands, and $E(M) - E(M-1, i)$ for valence bands. Phenomenologically, this energy includes a surface polarization term, as described by $P(\mathbf{r})$ in Eq. (1). Another way to look at this is that the screened interaction $W(\mathbf{r}, \mathbf{r}', \omega)$ includes a whole space inversion $\epsilon^{-1}(\mathbf{r}, \mathbf{r}_1, \omega)$ of the dielectric matrix $\epsilon(\mathbf{r}, \mathbf{r}_1, \omega)$. Compared to the bulk inverse dielectric function used in Eq. (32), the nanocrystal inverse dielectric function $\epsilon^{-1}$ contains a surface polarization term. This surface polarization energy is unscreened and scales as $1/R$; thus, it is a rather large energy. In contrary, in the conventional implementation of the CI approach, the $P(\mathbf{r})$ term has been removed from Eq. (1), and the single-particle Schrodinger's equation becomes Eq. (11), with the potential $V(\mathbf{r})$ being the same as the bulk value at the interior of the quantum dot. As a result, the single-particle energy $E_i$ in the CI approach (in the methods of Section 2) is no longer $E(M+1, i) - E(M)$ or $E(M) - E(M-1, i)$. To relate to these ionization energies, a surface polarization term needs to be added [165]. Thus, the GW single-particle energy and the CI single-particle energy are fundamentally different. For the TDLDA of Eq. (22), the single-particle energy is just the Kohn–Sham eigen energy. This energy is closer to the CI single-particle energy. In the LDA Kohn–Sham equation, the potential $V(\mathbf{r})$ of Eq. (17) at the interior of a quantum dot will be the same as its bulk value. There will be no surface polarization term. Again, this is probably because the LDA exchange correlation interaction is local. There is no whole-space inverse dielectric function as in the GW formula. However, a difference between TDLDA and CI approach is that there is a band-gap error in LDA eigen energies. This will even be true for exact density-function theory. For the exact TDDFT, however, other terms will correct this band-gap error of the Kohn–Sham eigen energies. This is not true in the case of TDLDA.

The $1/R$ polarization term in the GW quasiparticle eigen energy cancels a term in the screened Coulomb interaction $J$ in Eq. (40) and (41). This is because the same surface polarization term also exists in the $W(\mathbf{r}, \mathbf{r}', \omega)$ of Eq. (41). This cancellation is completely analogous to the cancellation of the $P_M(\mathbf{r}_1, \mathbf{r}_2)$ and $P(\mathbf{r}_1)$, $P(\mathbf{r}_2)$ terms in Eq. (2) of the classical phenonemological analysis. Recently, using TB GW–Bethe–Salpeter equation calculations, Delerue, Lannoo and Allan [166] have shown numerically that the Coulomb correction term exactly cancels the polarization term in the self-energy of the quasi-particle eigen energy. Thus, at the end, the result of GW–Bethe–Salpeter equation should be similar to the results of CI, where the Coulomb interaction $J$ is screened by the bulk inverse dielectric function, not the nanostructure dielectric function whole-space inversion. For the TDLDA

approach, the Coulomb interaction is local; thus, it is also effectively screened by the bulk dielectric function, not the nanostructure one. This corroborates well with the fact that the single-particle equation in TDLDA Eq. (22) does not include the surface polarization effect. The different cancellation scheme can be seen by comparing the calculated absorption spectra from the many-body Eqs. (40), (22), and (29) with the spectra calculated from the single particle eigen energies alone. As shown in Fig. 25 [167], the Bethe–Salpeter equation absorption spectra for small $Si_nH_m$ clusters red-shifted from the results calculated using the GW eigen energies alone. This is mainly because of the negative surface polarization energies in the Coulomb interaction $J$. However, as shown in Fig. 22, the TDLDA optical absorption spectrum blue-shifted from the LDA results. This is because in the TDLDA calculation, the polarization energy does not exist in both the Coulomb interaction term and the LDA Kohn–Sham eigen energy. As a result, it is the exchange interaction that dominates the spectrum shift. Notice that, if the total energy differences $E(M + 1) - E(M)$ and $E(M) - E(M - 1)$ from LDA calculations, instead of the Kohn–Sham eigen energies, are used in the many-body equation [e.g. Eq. (22)], then the corresponding interactions with the surface polarization term must be used [120, 168]; one cannot use the Coulomb interaction of Eq. (31) and (32).

## 3.4. Quantum Monte Carlo Methods

All the above methods discussed so far are derived from solid-state many-particle approaches. They treat the excited system as an electron and hole two-particle system. Under the quantum Monte Carlo (QMC) approach [169], the whole system is described by a many-body wavefunction. There are variational quantum Monte Carlo methods [170, 171] (VMC) and diffusion quantum Monte Carlo methods [172, 173] (DMC). Under VMC,



Figure 25. Calculated optical absorption spectrum of $Si_nH_m$ clusters using GW–Bethe–Salpeter equation. The spectra include an artificial broadening of 0.04 eV. The dotted lines show the spectra without electron-hole interaction. Reprinted with permission from [167]. M. Rohlfing and S. G. Louie, *Phys. Rev. Lett.* 80, 3320 (1998). © 1998, The American Physical Society.

the many-body wavefunction $\Psi(X)$ is approximated by a Slater determinant multiplied by a Jastrow term

$$\Psi(X) = D^{\uparrow}(R)D^{\downarrow}(R)\exp\left[\sum_{i=1}^{M}\chi(r_i) - \sum_{i-j}u(|r_i - r_j|)\right]$$ (42)

here $X = \{r_i, s_i\}$ for $i = 1, M$. Usually, HF or LDA single-particle wavefunctions are used to construct the slater determinant $D$, and parametrized forms are used to express $\chi$ and $u$. Within VMC, the total energy of the system is obtained via the variational minimum of the expectation value $E = \langle\Psi|H|\Psi\rangle/\langle\Psi|\Psi\rangle$, here $H$ is the many-body Hamiltonian. This multidimensional integral is evaluated using the Monte Carlo method

$$E = \frac{\int \Psi(X)H\Psi(X)dX}{\int |\Psi(X)|^2 dX} = \frac{\int \left[\frac{H\Psi(X)}{\Psi(X)}\right]|\Psi(X)|^2 dX}{\int |\Psi(X)|^2 dX}$$ (43)

Under the metropolis Monte Carlo scheme, the above integral can be statistically averaged using a walker with its equilibrium probability distribution in the multidimensional space $X$ equals $|\Psi(X)|^2$. Under this scheme

$$E = \frac{1}{N_s}\sum_{s}\left[\frac{H\Psi(X_s)}{\Psi(X_s)}\right]$$ (44)

where $X_s$ is the position of the walker's $s$ step, and $N_s$ is the total number of the steps. The $|\Psi(X)|^2$ probability distribution will be reached if in the metropolis scheme, a new step $X' = \delta X + X$ from the current step position $X$ will be accepted when $\mu = |\Psi(X')|^2/|\Psi(X)|^2 > 1$, and accepted with probability $\mu$ when $\mu < 1$. This Monte Carlo technique using $|\Psi(X)|^2$ as the walker distribution is also called "important sampling." It samples the important region heavily, thus reduces the statistical fluctuations. Note that if $\Psi(X)$ is exact, satisfying the Schrodinger's equation $H\Psi = E\Psi$, then the sampling value in Eq. (44) is always $E$, and thus there is no statistical fluctuation at all. The main computational effort in the VMC is the evalution of $\Psi(X)$ and $H\Psi(X)$ for a given $X$, especially when one step $\delta X$ is taken, typically moving only one $r_i$ among the $M$ particle coordinations. It turns out [174] that the evalution of both $H\Psi(X)$ and $\Psi(X)$ depend heavily on the evalution of the Slater determinant in Eq. (42).

For the diffusion QMC method (DMC), the many-body imaginary-time Schrodinger's equation is treated as a classical diffusion equation [172, 173]. Following an analogy of a classical diffusion problem, the equilibrium distribution (the solution of the Schrodinger's equation) can be reached following the particle diffusion, which can be described by Monte Carlo walkers. However, for Fermion system, antisymmetry is required for the many-body wavefunction. This pauses a sign problem, which is usually approximately solved by a fixed nodal approximation, where an auxiliary wavefunction is used to give the fixed nodal for the DMC wavefunction. This auxiliary wavefunction also serves as a guiding wavefunction for important sampling. Usually, the VMC wavefunction of Eq. (42) is used as this auxiliary wavefunction.

With the use of pseudopotentials [174], both VMC and DMC methods have been used to solve systems of up to a dozen atoms. Williamson et al. [175] also showed that QMC methods can be used to solve exciton energies for excited states. The key is to replace one single-particle valence wavefunction with a conduction-band wavefunction in the Slater determinant $D$ in Eq. (42). The DMC under such a new nodal structure automatically takes care of the resulting correlation effects. The so-calculated Si band structure agrees well with the experimental values, as shown in Fig. 26. QMC is one of the most accurate methods for small system calculations. Recently, Grossman et al. [176] have compared the QMC and the GW–Bethe–Salpeter equation results for $SiH_4$ and $CH_4$; excellent agreements are found for these two methods, ranging from ionization potentials and electron affinities to various optical excitation energies.

The above-discussed QMC methods are limited to a dozen atoms. However, the development of a linear-scaling QMC method has changed that limit [177]. The idea is to use

**Figure 26.** The diffusion quantum Monte Carlo band structure for bulk Si (filled circles with error bars). The solid lines are the result of empirical pseudopotential fitting to experimental data. Reprinted with permission from [175], A. J. Williamson et al., *Phys. Rev. B* 57, 12140 (1998). © 1998. The American Physical Society.

localized Wannier functions in the Slater determinants. Notice that, under a unitary transformation among the single-particle wavefunctions, the Slater determinant $D$ will not change. However, under such unitary transformation, the single-particle eigenfunctions can be changed into localized Wannier functions. As a result, these Wannier functions can be truncated in space and stored on real space grid. This makes the Slater determinant $D$ sparse for any given point X; hence, it also makes its calculation proportional to the size $N$, instead of $N^3$ as in the old scheme. This approach allows the QMC calculation of a few hundred atoms and makes it possible to use the QMC method to calculate small quantum dots [178].

# 4. PHYSICAL PROPERTIES OF SILICON QUANTUM DOTS

In Sections 2 and 3, we discussed different methods to calculate the electronic structures and optical properties of nanostructures. In this section, we will discuss the results of these methods applied to nanometer Si quantum dots and will compare these results to experimental measurements.

## 4.1. Optical Band Gap

The single-particle eigen energies of H passivated Si quantum dots and wires calculated using the methods of Section 2 are shown in Fig. 21. Results for three types of methods are shown: empirical TB method, empirical pseudopotential method, and LDA methods. All these methods give basically the same curve. Notice that for the TB method, the results of Ren [71] and Hill–Whaley [72–74] are not included because it was found later [80] that their TB parameters do not reproduce the proper Si bulk-band structures. For the LDA method, a rigid shift is used to correct the bulk LDA band-gap error. From the good mutual agreements between these methods, one can say that the single-particle eigenstate energies for H-passivated Si quantum dots are well understood.

In addition to the three methods discussed above, the effective mass theory and the truncated crystal method do not produce accurate single-particle band-gap energies. This is shown in Fig. 10. The effective mass approximation [54, 55] severely overestimates the band-gap opening, whereas the truncated crystal method of RKF [68] underestimates the band-gap opening. Note that the truncated crystal method works well for two-dimensional

thin film, but not for wires and quantum dots. Figure 10 also shows that the band gap for different shapes of quantum dots aggregate into a single curve when plotted as a function of the total number of Si atoms in the quantum dot.

To be compared with the experimental optical band gap, the electron-hole interactions must be included. Thus, the methods of Section 3 must be used. The results of these methods are shown in Fig. 27 for H-passivated Si quantum dots [143]. Note that, the DMC method and GW-BSE method produce almost the same band gap for very small quantum dots. For slightly larger quantum dots up to a diameter of 1.6 nm, the DMC result is almost 1 eV above all the other results. It remains to be seen how accurate these DMC results are, for example, when compared with well-controlled experiments (perhaps for other material quantum dots such as CdSe). The TDLDA method is found to be almost the same as the LDA method, taking the CBM-VBM of the LDA single-particle eigen energy as the band gap. This means that both the exchange and Coulomb interaction in the TDLDA results are very small. Note that in Fig. 27, the TDLDA result is taken from the lowest possible transition in the absorption spectrum. This is different from Ref. [179], in which a cutoff criterion is used in the calculated absorption spectrum to determine the optical band gap. The band gap obtained that way is higher than the result shown in Fig. 27, and brings it to a better agreement with the GW-Bethe-Salpeter equation and QMC results. However, it was argued that [143] to compare to the exciton energy calculated in GW-Bethe-Salpeter equation and QMC method, the first possible transition in TDLDA absorption spectrum should be taken. In addition to TDLDA, TDDFT-B3LYP was performed in Ref. [143] and Ref. [142]. Although their results are slightly different because of the use of different basis sets, overall, the TDDFT-B3LYP band gap is below the DMC result, especially for relatively large quantum dots. However, in the work of Ref. [142], it was shown that for small molecules, the TDDFT-B3LYP result agrees, with the MR-MP2 quantum chemistry calculations. The results for the TB method and the EPM can be considered as the lowest-order results of the CI in Eq. (29) (i.e., only the zero-order screened Coulomb interactions between



Figure 27. Size dependence of optical gaps of silicon nanoclusters, calculated using diffusion quantum Monte Carlo (DMC), GW-Bethe-Salpeter equation (BSE), local density approximation (LDA), and time-dependent LDA (TDLDA), time-dependent density-functional theory (DFT) with B3LYP functional (TDDFT-B3LYP), semiempirical tight-binding, and semiempirical pseudopotential methods. Note the DMC and GW-BSE results are almost the same for the few small clusters. Reprinted from [143]. A. J. Williamson et al., *Phys. Rev. Lett.* 89, 196803 (2002). © 2002. The American Physical Society.

the VBM and CBM states are taken into account). The TB and EPM results agree well with each other. However, they are between the TDLDA result and TDDFT-B3LYP results.

Overall, among the theoretical calculations, the DMC result is above all the other methods for $d \leq 1.5$-nm Si quantum dots. The LDA and TDLDA have the lowest band gap, followed by the TB and EPM limited CI results. Slightly higher than the TB and EPM results are the TDDFT-B3LYP results. For very small quantum dots, the DMC results agree with the GW-Bethe–Salpeter equation results.

The comparisons with the experiments are much more complicated. This is for two reasons: the uncertainty of the quantum dot sizes and the complication of surface passivations. As discussed in Section 1, the experimental data are very scattered. The situation is illustrated in Fig. 28 without listing all the overwhelming experiment data exist in this field. From Fig. 28, we can see that, when the experimental band-gap is measured from the absorption spectrum, the measured band-gap opening agrees well with the EPM calculated openings. However, the measured experimental PL energy is systematically lower than the theoretically calculated one. Notice that the other theoretical results (DMC and TDDFT-B3LYP) are even larger than the EPM results shown in Fig. 28. Thus, this conclusion is true for all the theoretical results (although most experimental data are in a quantum dot size region larger than the theoretical calculations shown in Fig. 27, except for EPM and TB methods). This difference between the calculated exciton band gap and the experimentally measured PL energy will be reconciled in Section 4.5 by carefully considering the surface passivations.

In addition to the optical band gaps, x-ray emission has been used to measure the band shoulder shifts of the valence bands and conduction bands separately [180, 181]. When these shifts are plotted as functions of the Si quantum dot sizes, the measured shifts are still smaller than the theoretically calculated results. However, the sizes of these quantum dots are measured by the atomic force microscopy (AFM) via the heights of the quantum dots on a substrate. It is not clear whether these quantum dots have spherical shapes. In addition, quantum dot aggregations on the substrate might also reduce the effective quantum confinement [180]. Nevertheless, when the valence band shifts and conduction band shifts are plotted against each other, the size uncertainty issue can be eliminated. The calculated [182] $\Delta\epsilon_{VBM}/\Delta\epsilon_{CBM}$ curve agrees well with the experimental results, as shown in Fig. 29.

## 4.2. Screening Effects and Dielectric Functions

Dielectric screening in a small quantum dot is one of the most fundamental physical issues in nanoscience. It is also a rather unsettled issue at this point. As one can see in Section 3, many of the methods involve dielectric functions and inverse dielectric functions for the quantum dot systems. In the single-particle formalism of Section 2, the dielectric function



Figure 28. Calculated excitonic band gap [69] and measured band gap and photoluminescence (PL) energies as a function of nanocrystal sizes. The calculated and measured bandgaps (from absorption spectrum) are in good agreement, whereas the PL energy is consistently lower than the excitonic bandgap. Reprinted with permission from [44], P. M. Fauchet et al., *Phys. Stat. Sol.* (a) 165, 3 (1998). © 1998, Wiley InterScience.

Figure 29. Conduction-band edge shifts versus valance-band edge shifts. (1), and (2), and (3) are calculated results using different dielectric constants, and the experimental data are from Ref. [181]. Reprinted with permission from [182]. L. W. Wang and A. Zunger. *Phys. Rev. Lett.* 73. 1039 (1994). © 1994. The American Physical Society.

and its corresponding surface polarization also play essential roles. A better understanding of the dielectric function of a quantum dot will help us come up with a better model for the semiempirical methods discussed in Section 2 and with a better definition of the single-particle eigen energies and surface image potential interactions.

We discussed the classical screening model in Section 2 when we introduced the surface polarizations. In that model, the dielectric constant inside the quantum dot is $\epsilon$, which is taken as the same as in the bulk system. Tsu et al. [183] pointed out that this constant inside the dot might be different from the bulk. This is argued using a generalized Penn's model (GPM) [184]. Within this model, the dielectric constant is expressed as a function of the absorption spectrum peak position. Because this peak position will shift under quantum confinement, it is then argued that the dielectric constant might also change. The generalized Penn's model given by Tsu is

$$\epsilon_s(R) = 1 + \frac{\epsilon_b - 1}{1 + \alpha/R^l} \tag{45}$$

here $\epsilon_b = 11.4$ is the bulk dielectric constant, and $\alpha = 10.93$ Å and $l = 2$.

To test this phenomenological model, Wang and Zunger [182] have carried out a direct calculation based on the EPM method. The task is to calculate the imaginary part of the dielectric function (e.g., the optical absorption spectrum). Then, by a frequency integration, one can get the real part of the dielectric function. This calculation is done using a GMM (generalized moments method) [88]. Within this method, to calculate the imaginary part $\epsilon_2(E)$ of the dielectric function, we first calculate a two-dimensional spectral function

$$\tau(E_1, E_2) = \sum_{i,t} |\langle \psi_i|\hat{p}|\psi_j\rangle|^2 \delta(E_1 - E_i)\delta(E_2 - E_f) \tag{46}$$

After $\tau(E_1, E_2)$ is obtained, the $\epsilon_2(E)$ is calculated following the random phase approximation [185]

$$\epsilon_2(E) = \frac{A}{\Omega E^2} \int_{-\infty}^{E_s} dE_1 \int_{E_c}^{\infty} dE_2 \tau(E_1, E_2)\delta(E - E_2 + E_1) \tag{47}$$

where $A = 8\pi^2 e^2\hbar^2/3m^2$, and $\Omega$ is the volume of the system. To calculate $\tau(E_1, E_2)$ without knowing all the single-particle eigen states $\psi_i$ and $\psi_j$ in Eq. (46), one can first generate the two-dimensional Chebychev moments $\Gamma_{n,m}$ of the spectra $\tau(E_1, E_2)$. These moments $\Gamma_{n,m}$ can be generated in the following way

$$\Gamma_{n,m} = \langle\langle\psi_i|\hat{p}T_n(\hat{H}) \cdot \hat{p}T_m(\hat{H})|\psi_j\rangle\rangle \tag{48}$$

here $T_n$ and $T_m$ are Chebychev polynomials, $\psi_r$ is a random starting wavefunction, and the outer angle bracket denotes a statistical average over the random wavefunction $\psi_r$. Then, from $\Gamma_{n,m}$ to $\tau(E_1, E_2)$, one has the following transformation:

$$\tau(E_1, E_2) = \left(\frac{2}{\pi}\right)^2 (1 - E_1^2)^{-1/2}(1 - E_2^2)^{-1/2} \sum_{n,m} T_n(E_1)T_m(E_2)\Gamma_{n,m}b_n b_m \qquad (49)$$

here $b_n = 0.5$ for $n = 1$, and 1 for all the other $n$. This transformation can be carried out using FFT because $T_n(E)$ can be cast as a cosine function. Using this GMM method, the density of state (DOS) and joint DOS (JDOS) can also be calculated. The details of GMM are described in Ref. [88].

Using the GMM formalism, Wang and Zunger [182] calculated the optical absorption spectra for Si quantum dots, as shown in Fig. 30 [186]. Then the dielectric constants as functions of the quantum dots are calculated as an integral from the optical absorption spectrum. Indeed, the dielectric constant is found to reduce significantly from its bulk value, as shown in Fig. 31. For a quantum dot, different dielectric constant can be defined. One as the total polarization of the quantum dot, another as the screening inside the quantum dot. These two dielectric constants are different. Using the formula of Eq. (45), the EPM result for the total polarization can be fitted by $\alpha = 4.25$ Å and $l = 1.25$, and the internal screening dielectric constant can be fitted by $\alpha = 6.9$ Å and $l = 1.37$. The small exponential component $l$ compared to the generalized Penn's model is reminiscent of the soft scaling of the calculated EPM quantum confinement compared with effective mass result.

The dielectric constant was also calculated by using self-consistent calculations. In Ref. [187], Lannoo, Delerue, and Allan calculated a donor impurity at the center of a quantum dot with a LCAO basis and a Hartree Coulomb interaction. Self-consistent potentials and donor binding energies are calculated, and the results are fitted to dielectric screening



Figure 30. Calculated optical absorption spectra of Si systems with different sizes. The experimental data in (a) are from Ref. [186]. The vertical arrows denote band gap values. Reprinted with permission from [182], L. W. Wang and A. Zunger, *Phys. Rev. Lett.* 73, 1039 (1994). © 1994, The American Physical Society.

**Figure 31.** Calculated reduced dielectric constants in a Si quantum dot. The dashed line is for the generalized Penn's model. $\epsilon_s$ is for total polarization, and $\tilde{\epsilon}_s$ is for exciton screening. Reprinted with permission from [182], L. W. Wang and A. Zunger, *Phys. Rev. Lett.* 73, 1039 (1994). © 1994, The American Physical Society.

models. The so-calculated results can be fitted with $\alpha = 9.2$ Å and $l = 1.18$. Furthermore, the dielectric screening has been calculated using LDA methods [120]. The total polarization of the quantum dot is calculated under an external perturbation. Again, the result can be fitted by Eq. (45) with $\alpha = 9.7$ Å and $l = 1.3$. Experimentally, it is very difficult to measure the dielectric constant of a quantum dot directly. Krauss and Brus [188] have used a AFM tip and tip-substrate capacitance to measure the total polarization and charge effect of a single quantum dot when it is positioned between the AFM tip and the substrate. Their estimate of the dielectric constant for a CdSe quantum dot roughly agrees with the theoretical prediction [189].

One recent advance on the issue of dielectric screening is the work by Ogut, Burdick, Saad, and Chelikowsky [190]. In this work, the dielectric matrix $\epsilon(r_1, r_2)$ is calculated and the whole matrix is inverted, using some mathematical tricks. The inverse dielectric function $\epsilon^{-1}(r_1, r_2)$, written as an effective $\tilde{\epsilon}(r_1, r_2)$ in the style of Eq. (32) [replacing $\epsilon_{\text{bulk}}^{-1}$ with $\epsilon^{-1}$ in Eq. (32)], is shown in Fig. 32. This $\tilde{\epsilon}(r_1, r_2)$ approaches 1 before the quantum dot radius $R$. Unfortunately, the quantum dots are too small to derive any definite conclusions regarding issues like the surface polarization model.

## 4.3. Photoluminescence Lifetime and Phonon Coupling

Because of the finiteness of the quantum dot, the translational symmetry of the periodic crystal has been broken. As a result, there is no momentum conservation rule in the optical transition. The indirect Si transition in bulk has changed into a pseudodirect transition. If the



**Figure 32.** The averaged screening functions $\tilde{\epsilon}(r) = 1/\tilde{\epsilon}^{-1}(0, r)$ for the three quantum dots $SiH_4$, $Si_5H_{12}$, and $Si_{35}H_{36}$. Their effective radii are 3.2, 5.4, and 10.4 a.u., respectively. Reprinted with permission from [190], S. Ogut et al., *Phys. Rev. Lett.* 90, 127401 (2003). © 2003, The American Physical Society.

single-particle wavefunction $\psi_i$ has been calculated from Eq. (11), the optical transition lifetime $\tau$ between states $i$ and $f$ can be calculated as

$$\frac{1}{\tau} = \frac{4}{3}\frac{\alpha\omega n}{m^2 c^2}|\langle\psi_i|\hat{p}|\psi_f\rangle|^2$$                    (50)

where $n = 2.6$ is the effective refractive index of Si quantum dot, and $\omega$ is the photon angular frequency, $\alpha = e^2/\hbar c$. Equation (50) has been used to calculate the CBM to VBM recombination rate at room temperature under the EPM approach; the result is shown in Fig. 33 [191, 192]. The empirical TB calculation produce almost identical results [36] shown in Fig. 7. From Figs. 7 and 33, we see that the calculated PL life time is about 10 times longer than the experimentally measured value. This means that although the CBM to VBM transition is allowed and pseudodirect, this transition itself is not strong enough to explain the observed PL lifetime.

Experimentally, it is shown (Fig. 5) that the phonon-assisted transitions are strong, with their PL peak amplitudes stronger than the nonassisted transition. Thus, the key to get a better agreement with the experiment for the transition lifetime is to consider the phonon-assisted transitions. This was first done by Hybertsen [193]. Hybertsen used an effective mass model for the electron wavefunctions and bulk parameters for the electron–phonon couplings. He was able to show that the phonon-assisted transitions are indeed about an order of magnitude larger than the zero-phonon transition, as shown in Fig. 34. Recently, Delerue, Allan, and Lannoo [194] have calculated the phonon-assisted transitions using a TB model. In this calculation, the electron and phonon systems are treated as a whole. This allows a treatment of multiphonon processes. All the phonon modes are calculated using a valence force field model, and the coupling between all the phonon modes and the transition electronic states are calculated explicitly using Harrison's rule [83] for changes of TB parameters following the atomic displacements. Again strong phonon-assisted transitions are obtained. In particular the resonant PL spectrum shapes with the phonon replica peaks are calculated, and they compare well with the experimental results, as shown in Fig. 35. All the work has demonstrated that the PL lifetime and phonon-assisted transitions in Si quantum dots are relatively well understood.



Figure 33. The radiative recombination rate as a function of the luminescence photon energy. Experimental curves (1) and (2) are from Ref. [191] and [192], respectively. The symbols ◆, +, and ■ represent calculated results for spherical, rectangular, and cubic quantum dots and are for zero-phonon process. Reprinted with permission from [69], L. W. Wang and A. Zunger, *J. Phys. Chem.* 98, 2158 (1994). © 1994, The American Chemical Society.

**Figure 34.** Radiative recombination time as a function of the blue-shift of the photon energy from the bulk Si band edge: zero-phonon transitions (dots): transverse optical (TO) phonon-assisted transitions (line). The top scale indicates the equivalent cube size. Reprinted with permission from [193], M. S. Hybertsen, *Phys. Rev. Lett.* 72, 1514 (1994). © 1994. The American Physical Society.

## 4.4. Exchange Splitting

The PL lifetime discussed in previous sections is the lifetime at room temperature. At very low temperature, the PL lifetime can be much longer, and the fine structure of the electronic state can be deduced from the lifetime analysis. By studying the PL lifetime dependence on the temperature, Calcott et al. [195] have proposed a two-level model. In this model, the PL comes from two energy levels split by $\Delta$. The upper level has a short lifetime $\tau_U$, and the lower level has a long lifetime $\tau_L$, and this level is three fold degenerated. These two levels are occupied by thermal distribution under $kT$. As a result, the total PL life time $\tau$ is

$$\tau^{-1} = [3\tau_L^{-1} + \tau_U^{-1}\exp(-\Delta/kT)]/[3 + \exp(-\Delta/kT)] \tag{51}$$

When formula is used to fit the experimental data, they get $\tau_L$ and $\tau_U$, as shown in Fig. 36. When the temperature is high, the $\tau_L$ and $\tau_U$ change with $T$. That might indicate that other energy levels are involved, and Eq. (51) might no longer be a good description, but when the temperature is low, $\tau_L$ and $\tau_U$ are independent of $T$, showing the validity of this model. These two levels cannot be explained by phonon replicas. Instead, they have been assumed to be the result of exchange splitting [195]. In particular, this will explain the threefold degeneracy of the lower level state and its long radiation lifetime: it is a $S = 1$ spin forbidden dark state. The fitted energy splitting $\Delta$ as a function of the PL energy (and thus the quantum dot size) is shown in Fig. 37.



**Figure 35.** Experimental resonant photoluminescence spectrum (full line) compared to the theoretical result for a hydrogen-passivated cluster without (dashed line) or with (dotted line) phonon-assisted enhancement. Reprinted with permission from [194], C. Delerue et al., *Phys. Rev. B* 64, 193402 (2001). © 2001. The American Physical Society.

**Figure 36.** Temperature dependence of the Photoluminescence (PL) decay time. Each broken line is a fit of the data at a single PL energy to the model of a thermal-equilibrium two-level system, with the upper level lifetime $\tau_l$ and lower level lifetime $\tau_l$ and a fixed energy splitting $\Delta$. The ratio of the lower state to the upper state degeneracies is assumed to be three. Reprinted with permission from [195]. P. D. J. Calcott et al., *J. Phys. Condens. Matter.* 5, L91 (1993). © 1993, The Institute of Physics.

There is another way to measure this energy splitting directly; called onset energy in resonantly excited PL. It is found that the measured PL spectrum has a onset energy $\Delta_{on}$ below the exciting laser energy. Above this energy (i.e., in the energy range of $E_{laser} - \Delta_{on}$ and $E_{laser}$), there is almost no PL signal. The PL spectrum starts below $E_{laser} - \Delta_{on}$. This is like a Stokes shift, not in the PL peak, but in the starting point of the PL spectrum. This energy loss can be easily explained with the two-level model; thus, the onset energy $\Delta_{on}$ is a direct measurement of the two-level splitting $\Delta$. The result of this onset energy is also shown in Fig. 37. Note that this onset energy is systematically smaller than the $\Delta$ measured from PL life time analysis of Eq. 51.

Theoretically, Martin et al. [196] have calculated the exchange splitting using the empirical TB model. This can be viewed as a limited CI calculation using Eqs. 29 and 30. Here, only the VBM and CBM states are used in the CI configuration with their four possible spin combinations. These four degenerated states are split by the exchange interaction of



**Figure 37.** Exchange splitting between the two lowest excitonic levels in asymmetrical silicon crystallites with respect to their excitonic band gap. Crystallites have undulating ellipsoidal shapes with a longer axis in the (100) direction (open circles), (110) direction (open triangles), and (111) direction (+). Crosses correspond to the average of the splitting over all the orientations. Black squares are the first onsets measured by selectively excited photoluminescence, and black dots are the energy splitting derived from the fit of temperature dependence of the lifetime. Reprinted with permission from [196]. E. Martin et al., *Phys. Rev. B* 50, 18258 (1994). © 1994, The American Physical Society.

Eq. (30), producing a $S = 1$ lower energy state and $S = 0$ higher energy state. The calculated splitting for different shapes of quantum dots is shown in Fig. 37. Although of the same order of magnitude, the calculated results are systematically smaller than the experimentally measured values. Leung and Whaley [78] recalculated this exchange splitting using their TB model; their result is much larger than the result of Martin et al. [196], although both used empirical TB models. As shown in Fig. 38 [197], Leung's result agrees well with the $\Delta$ from lifetime measurement in Ref. [195]. This situation further changed when Reboredo et al. [198] recalculated this exchange splitting using the EPM. In the EPM, the exchange interaction in Eq. (30) can be calculated directly without making any assumptions for the atomic wavefunctions like in the TB models. As shown in Fig. 38, the resulting EPM splits agree well with the experimentally measured onset energy in Ref. [195]. There are other experimental $\Delta$ using the lifetime analysis. Some [197] of these results agree with the thermal analysis data in Ref. [195], and some [43] agree with the spectrum onset data in Ref. [195]. Interestingly, Reboredo et al. [64] have also proposed another model, in which the energy levels of a perfectly $T_d$ symmetry Si quantum dot can be split by the Coulomb interaction alone, without the help of the exchange interaction.

## 4.5. Effects of Surface Passivations

As we discussed earlier, the experimental PL energy depends sensitively on the surface passivation and oxidization. Table 1 lists the possible surface passivation atoms for different synthesis methods and postsynthesis conditions. Figure 8 shows the aging effect of Si quantum dots placed in the atmosphere as the surface is graduately oxidized away. Figure 39 shows the scatter of experimental PL energies produced by different experiments [199]. The higher-energy points in Fig. 39 might be the result of different PL peaks, which are related to bulk direct transition $\Gamma_{25} - \Gamma_{15}$. Even for the lower-energy points in the shaded area, however, the scattering is about 1 eV. Nevertheless, as shown in Fig. 28 all these scattered experimental PL energies are below the theoretically calculated ones, assuming full H passivations. Now most researchers believe that this difference is the result of the existence of surface O atoms. As shown in Table 1, O atoms exist in almost all the samples, unless the sample is under very careful control. A careful experiment has been carried out by Wolkin et al. [39] to address the question of surface O and the red-shift it produces. The result is



**Figure 38.** Exchange splitting between the two lowest excitonic levels and comparison to different methods, both theories and experiments. The calculated results are Reboredo for Ref. [64], and Leung for Ref. [78]. The experimental results are Calcott et al. (opt) (Therm) from Ref. [195], Kovalev et al. from Ref. [43], and Borngersma et al. from Ref. [197]. Reprinted with permission from [64]. F. A. Reboredo et al., *Phys. Rev. B* 61, 13073 (2000). © 2000. The American Physical Society.

Figure 39. Summary of different experimental data on peak PL energy versus Si nanocrystal size. Reprinted with permission from [199], J. P. Wilcoxon et al., *Phys. Rev. B* 60, 2704 (1999). © 1999, The American Physical Society.

shown in Fig. 40. It is shown that for Si quantum dots with diameter less than 2 nm, the introduction of O can cause a PL red-shift as large as 1.0 eV. However, before the red-shift, the fully H-passivated quantum dot has a PL energy that agrees excellently with the EPM and empirical TB calculated results.

The O passivation effect on Si quantum dot is currently under intense theoretical investigation. It has been calculated by many methods, including empirical TB, HF-CI, LDA, TDLDA, and QMC. In the same paper by Wolkin et al. [39], Allan and Delerue have calculated the O effect by considering a Si$=$O double bond (silinone). They found that the band gap is, indeed, reduced. They have classified three zones according to the quantum dot diameter $d$. When $d > 3$ nm, the surface O atom has no effect on the PL. When $1.6 < d < 3$ nm, the surface O atom will produce a localized conduction band state. When $d < 1.6$ nm, the surface O will also produce a localized valence-band state. The calculated overall band gap agrees well with the experimentally measured values, as shown in Fig. 40.



Figure 40. Comparison between experimental and theoretical PL energies as a function of crystallite size. The upper line is the calculated band gap for fully H passivated Si quantum dots, and the lower line is the calculated band gap in the presence of one Si$=$O bond. The ● and ○ are the measured peak PL energies for freshly synthesized and oxidized Si crystallites respectively. Reprinted with permission from [39], M. V. Wolkin et al., 82, 197 (1999). © 1999, The American Physical Society.

The effects of surface O atom have also been studied by Caldas [200], using a quantum chemistry semiempirical HF-CI method. In this study, it is found that the Si=O can be energetically costly compared with other O surface incorporations, like an Si–O–Si bridge, with one O replacing two H atoms bonded to two neighboring Si atoms. However, it has also been found that the Si–O–Si has a small effect on the optical transition energy. On the other hand, the Si=O bond will produce a large red-shift, as in the empirical TB calculation. Furthermore, it is found that a large Stokes shift exists in the excited state, corresponding to a large surface relaxation. This is true even in purely H bonded surfaces. Such excited-state Stokes shift might cause some H defect states, as proposed in Ref. [201].

The same O passivation effects have been studied using the LDA method. In the LDA calculations [202–204], as before, it was found that the the Si=O bond will introduce a large band-gap reduction, whereas the Si–O–Si bridge configuration will only introduce a small reduction. This is shown in Fig. 41 by the CBM–VBM band gap of the LDA calculation. In reality, many O atoms can be on the surface. Thus, the effects of multiple O atom have also been studied. In the work of Luppi et al. [203], a saturation level is found after a few O atom incorporations, and the final energy level is independent of the quantum dot sizes (although they are all small quantum dots up to 35 Si atoms). This is shown in Fig. 42. However, in the calculation of Puzder [202], for a 66-Si-atom quantum dot, the band gap continues to decrease almost down to the bulk band gap as the number of O increases. This is shown in Fig. 43. This result seems different from that in Fig. 42.

Similar systems have also been studied by Vasiliev et al. [205], using the TDLDA method. Again, red-shift has been found. However, they found almost equal amount of optical band-gap red-shifts for Si–O–Si and Si=O passivations, although the absorption spectra shapes are different for these two types of passivations. This conclusion is different from all the other calculations. This might be the result of the special criterion used in their study to determine the optical band gap. Especially in this case, as reported in Ref. [205], the direct optical transition between CBM and VBM for both the Si–O–Si and Si=O are forbidden.

Finally, the QMC method is used to calculate the O passivation [178]. The results are shown in Fig. 44 [206, 207]. The QMC results are similar to the LDA results, but they



Figure 41. Comparison of local density approximation (LDA)-predicted conduction-band minimum-valence-band maximum gaps for H-passivated silicon nanoclusters (■), with a single Si=O bond (♦), Si–O–Si bridge (▲), and interstitial oxygen site (♦). HOMO-LUMO, highest occupied molecular orbit–lowest unoccupied molecular orbit. Reprinted with permission from [202], A. Puzder et al., J. Chem. Phys. 117, 6721 (2002). © 2002, The American Institute of Physics.

**Figure 42.** The local density approximation–calculated conduction-band minimum–valance-band maximum energy gaps as a function of the number of Si=O bonds at the cluster surface. The different values for the same cluster at fixed number of Si=O bonds correspond to different relative locations of the Si=O bonds at the surface. Reprinted with permission from [203], M. Luppi and S. Ossicini, *J. Appl. Phys.* 94, 2130 (2003). © 2003. The American Institute of Physics.

are up-shifted, as shown in Fig. 27. However, the shift for the purely H-passivated case is more than the shift for the O-passivated case. This is probably because for the O-passivated case, both the CBM and VBM are localized near the Si=O. Consequently, the electron-hole Coulomb interaction is larger in this case than in the purely H-passivated case. Thus, if the Coulomb interactions of Eq. (31) were included in the LDA result, the LDA Si=O curve in Fig. 44 should be lowered down relative to the purely H-passivated curve. This would bring the LDA result closer to the QMC result for the crossover quantum dot size. Overall, the comparison with the experiment results is illustrated in Fig. 44. We see that using the O passivation, we can at least qualitatively explain the scattered experimental data.



**Figure 43.** The local density approximation–calculated conduction-band minimum–valance-band maximum energy gaps as functions of the number of Si=O bonds and Si—O—Si bridges. The system is a $Si_{66}$ cluster. Homo-, umo, highest occupied molecular orbit–lowest unoccupied molecular orbit. Reprinted with permission from [202], A. Puzder et al., *J. Chem. Phys.* 117, 6721 (2002). © 2002. The American Institute of Physics.

**Figure 44.** Energy gap versus diameter for H-passivated silicon nanoclusters with and without Si$=$O bonds. The upper panel (a) shows results of local density approximation (LDA) and quantum Monte Carlo (QMC) calculations. The calculated gaps of the fully hydrogenated clusters are marked by filled circles. The thick lines are exponential fits to these calculated gaps that approach the respective theoretical bulk gaps. The lower panel (b) compares the results of experiments (Schuppler, Ref. [206]; van Buuren, Ref. [180]; Dinh, Ref. [207]) with the (solid) curves fitted to the QMC gaps. Reprinted with permission from [178], A. Puzder et al., *Phys. Rev. Lett.* 88, 97401 (2002). © 2002, The American Physical Society.

## 5. CONCLUSIONS

Si quantum dots can be divided into two classes: one that is the lithographic or surface deposition that synthesized relatively large quantum dots, and the other is the small nanometer quantum dots synthesized via a variety of methods. For the large Si quantum dots (10–30 nm), one of their main applications is the single-electron device. Thus physical phenomena like Coulomb Blockade and quantum transport are important. For the small quantum dots (1–7 nm), their main application is for optical devices. Thus, their optical properties are essential. In this review, we have focused on the small quantum dots and their optical properties. One especially important type of the small Si quantum dot is the p-Si. Since the discovery of the strong luminescence of p-Si, tremendous effort has been devoted to its study. In the last decade, it has been one of the main topics in nanoscience.

Like bulk silicon and silicon surfaces, Si quantum dot has been a test ground for methodology developments and has played an important role in computational nanoscience. It is an interesting physical system because it is indirect, that provides many new phenomena that do not exist in a direct material system. However, it is also a difficult system because the

experimental samples are often inhomogeneous with large size and shape distributions. Its optical properties also depend sensitively on the surface passivations. In the early days, this has caused a lots of confusion, but after careful clarification, the sensitivity on surface passivation can also provide rich physical phenomena. Compared with the organic surface passivations in II–VI quantum dots, the Si surface passivation is actually simpler and amendable to theoretical investigations.

After a decade of work, a picture of the mechanism for optical transitions in Si quantum dot has emerged. It is now widely believed that the PL of p-Si and most other small Si quantum dots comes from the interior quantum confinement states. There are four major optical properties in Si quantum dots: quantum confinement, lifetime, phonon replica, and exchange splitting. For all these four properties, theoretical calculations have yield qualitative to quantitative agreements with experiments. For the blue-shift quantum confinement effect on PL, for very small quantum dots (<2 nm), O passivation is needed to explain the experimental results. This is still an ongoing investigation, and it has been studied intensively by various theoretical methods. It is still not known, for example, given multiple O passivations, what is the critical size $d_c$, beyond which the states involved in optical transition change from surface states to interior states of the quantum dot. Different calculations have not converged on this issue, and the agreement with the experiment is qualitative at best. If $d_c$ is large, then one has to rethink the results in lifetimes, phonon-assisted transitions, and exchange splittings. At present, the calculated values for these properties agree well with experimental results; but the calculations are done under the assumption that the PL transitions are from the interior eigenstates and that the surface is fully passivated by H. Cleaner samples with narrower size distributions are highly desirable for better theory-experiment comparisons.

On the issues of methodology developments, there are still a lot of work to do. Computational nanoscience is still in its infancy. New methodologies, especially embedding techniques like the charge-patching method, are highly desirable. O(N) scaling is essential for any methods to deal with thousands of atoms. One of such methods is the single-particle-CI approach. However, at present this approach is based on classical model derivations and physical intuitions. A more solid derivation based on many-body theory will be very helpful. Another very promising approach is the QMC method. It is now possible to calculate a few hundred atoms using this approach. The conceptual simplicity of this method is one of its main appealing qualities, but this method is still relatively new for excited-state and large-system calculations. More testing is needed for its accuracy. If there is no rigorous proof that the DMC many-body wavefunction based on a single excitation configuration Slater determinant (replacing one single-particle valence-band state with one conduction-band state) will be orthogonal to the ground-state DMC many-body wavefunction based on the ground-state Slater determinant. If the CI scheme is any guide, then the excited state many-body wavefunction should be a linear combination of many single-excitation configuration Slater determinants, not a single Slater determinant. A linear combination of Slater determinants probably has different nodal structures than the single Slater determinant. Without any proof for these issues, the best approach is to test this method with experiments, both for bulks and for quantum dots. For GW–Bethe–Salpeter equation, one challenge is to make it scalable to larger systems. At present, this method can only be applied to systems with dozens of atoms. For the GW method, one issue is to deal with the intrinsically nonlocal self-energy term. Some numerical approximation might be essential there. For the Bethe–Salpeter equation the diagonalization under the single-excitation configuration space is a bottleneck. Perhaps some dynamic time integration scheme that changes the Bethe–Salpeter equation to a time-dependent single-particle equation is the key. This scheme would be similar to the TDLDA scheme, which connects the time-dependent Kohn–Sham equation with the Bethe–Salpeter equation as in Eq. (22). It will be interesting to compare Eqs. (22), (29), and (40) more carefully for TDLDA, CI and GW–Bethe–Salpeter equation respectively. Such comparison will reveal the connection between these methods, and perhaps will give the empirical approaches some solid ground and give the high-end approaches some ways of approximations. At the center of this comparison is the understanding of the dielectric screening in nanostructures, including its local field effects and dynamic screening effects. This dielectric screening is used in all three methods, and the different approximations of it have contributed to the differences of these

three approaches. Finally, for large systems, another way to solve Eqs. (22), (29), and (40), is to write the two-particle wavefunction in a variational form [e.g., as in Eq. (42) with a Jastrow term]. It needs to be proven that the variational minima of such two-particle (electron and hole) wavefunctions with an appropriate two-particle Hamiltonians are the solutions of Eqs. (22), (29), and (40).

## ACKNOWLEDGMENT

## REFERENCES

1. L. Zhuang, L. Guo, and S. Y. Chou, *Appl. Phys. Lett.* 72, 1205 (1998).
2. A. C. Irvine, Z. A. K. Durrani, H. Ahmed, and S. Biesemans, *Appl. Phys. Lett.* 73, 1113 (1998).
3. E. Leobandung, L. Guo, Y. Wang, and S. Y. Chou, *J. Vac. Sci. Technol. B* 13, 2865 (1995).
4. H. Ishikuro and T. Hiramoto, *Appl. Phys. Lett.* 71, 3691 (1997).
5. R. A. Smith and H. Ahmed, *J. Appl. Phys.* 81, 2699 (1997).
6. R. Ashoori, H. L. Stormer, J. S. Weiner, L. N. Pfeiffer, K. W. Baldwin, K. W. West, *Phys. Rev. Lett.* 71, 613 (1993).
7. A. D. Yoffe, *Adv. Phys.* 42, 173 (1993).
8. P. Sutter, P. Zahl, and E. Sutter, *Appl. Phys. Lett.* 82, 3454 (2003).
9. U. Denker, M. Stoffel, and O. G. Schmidt, 90, 196102 (2003).
10. P. D. Miller, C. Liu, W. L. Henstrom, J. M. Gibson, Y. Huang, P. Zhang, T. I. Kamins, D. P. Basile, and R. S. Williams, *Appl. Phys. Lett.* 46, 46 (1999).
11. C. B. Murray, D. J. Norris, and M. G. Bawendi, *J. Am. Chem. Soc.* 115, 8706 (1993).
12. J. E. B. Katari, V. L. Colvin, and A. P. Alivisatos, *J. Phys. Chem.* 98, 4109 (1994).
13. O. I. Micic, C. J. Curtis, K. M. Jones, J. R. Sprague, and A. J. Nozik, *J. Phys. Chem.* 98, 4966 (1994).
14. K. A. Littau, P. J. Szajoski, A. J. Muller, A. R. Kortan, and L. E. Brus, *J. Phys. Chem.* 97, 1224 (1993).
15. L. N. Dinh, L. L. Chase, M. Balooch, W. J. Siekhaus, and F. Wooten, *Phys. Rev. B* 54, 5029 (1996).
16. T. van Buuren, L. N. Dinh, I. Jimenez, L. J. Terminello, M. Grush, T. A. Calcott, and J. A. Carlisle, *Mater. Res. Soc. Symp. Proc.* 452, 171 (1996).
17. T. Makimura, Y. Kunii, N. Ono, and K. Murakami, *Appl. Surf. Sci.* 127-129, 388 (1998).
18. R. K. Soni, L. F. Fonseca, O. Resto, M. Buzaianu, S. Z. Weisz, J. Lumin. and S. Lumin. 83-84, 187 (1999).
19. T. S. Iwayama, S. Nakao, K. Saitoh, *Appl. Phys. Lett.* 65, 1814 (1994).
20. L. Pavesi, L. Dal Negro, C. Mazzoleni, G. Franzo, and F. Priolo, *Nature* 408, 440 (2000).
21. L. T. Canham, *Appl. Phys. Lett.* 57, 1046 (1990).
22. A. Uhlir, *Bell Syst. Tech. J.* 35, 333 (1956).
23. A. Halimaoui, in "Properties of Porous Silicon" (L. T. Canham, Ed.), p. 12. The Institution of Electrical Engineers, London, 1998.
24. J. von Behren, E. H. Chimowitz, and P. M. Fauchet, *Adv. Mater.* 9, 921 (1997).
25. R. L. Smith and S. D. Collins, *J. Appl. Phys.* 71, R1 (1992).
26. M. I. J. Beale, J. D. Benjamin, M. J. Uren, N. G. Chew, and A. G. Cullis, *J. Cryst. Growth* 73, 622 (1985).
27. V. Lehmann, *J. Electrochem. Soc.* 140, 2258 (1993).
28. J. P. Gonchond, A. Hamimaoui, and K. Ogura, in "Microscopy of Semiconducting Materials 1991" (A. G. Cullis and N. J. Long, Eds.), p. 235. IOP, Bristol, 1991.
29. A. G. Cullis and L. T. Canham, *Nature* 353, 335 (1991).
30. L. T. Canham, in "Optical Properties of Low Dimensional Silicon Structures" (D. C. Benshael, L. T. Canham, and S. Ossicini, Eds.), NATO ASI Series, Vol. 244, p. 81. Kluwer Academic, Dordrecht, 1993.
31. D. Kovalev, M. Ben-Chorin, J. Diener, F. Koch, Al. L. Efros, M. Rosen, M. A. Gippius, and S. G. Tikhodeep, *Appl. Phys. Lett.* 67, 1585 (1995).
32. H. Takagi, H. Ogawa, Y. Yamazaki, A. Ishizaki, and T. Nakagiri, *Appl. Phys. Lett.* 56, 2379 (1990).
33. A. G. Cullis, L. T. Canham, and P. D. J. Calcott, *J. Appl. Phys.* 82, 909 (1997).
34. P. D. J. Calcott, K. J. Nash, L. T. Canham, M. J. Kane, and D. Brumhead, *J. Lumin.* 57, 257 (1993).
35. P. D. J. Calcott, K. J. Nash, L. T. Canham, and M. J. Kane, in "Microcrystalline and Nanocrystalline Semiconductors" (R. W. Collins, C. C. Tsai, M. Hirosi, F. Koch, and L. Brus, Eds.), Materials Research Society, p. 465. Pittsburgh, PA, 1994.
36. M. Lannoo, G. Allan, and C. Delerue, in "Structural and Optical Properties of Porous Silicon Nanostructures" (G. Amato, C. Delerue, and H.-J. von Bardeleben, Eds.), p. 187. Gordon and Breach, Amsterdam, 1997.
37. E. Yablonovitch, *Phys. Rev. Lett.* 57, 249 (1986).
38. A. Loni, A. J. Simons, T. I. Cox, P. D. J. Calcott, and L. T. Canham, *Electronics Lett.* 31, 1288 (1995).
39. M. V. Wolkin, J. Jorne, P. M. Fauchet, G. Allan, and C. Delerue, *Phys. Rev. Lett.* 82, 197 (1999).
40. L. T. Canham, M. R. Houlton, W. Y. Leong, C. Pickering, and J. M. Keen, *J. Appl. Phys.* 70, 422 (1991).
41. L. Tsybeskov, S. P. Duttagupta, and P. M. Fauchet, *Solid State Commun.* 95, 429 (1995).

42. L. C. Kimerling, *Appl. Surf. Sci.* 159, 8 (2000).
43. D. Kovalev, H. Heckler, G. Polisski, and F. Koch, *Phys. Stat. Sol. (b)* 215, 871 (1999).
44. P. M. Fauchet, J. von Behren, K. D. Hirschman, L. Tsybeskov, and S. P. Duttagupta, *Phys. Stat. Sol. (a)* 165, 3 (1998).
45. O. Bisi, S. Ossicini, and L. Pavesi, *Surf. Sci. Rep.* 38, 1 (2000).
46. L. E. Brus, *J. Chem. Phys.* 79, 5566 (1983).
47. L. E. Brus, *J. Chem. Phys.* 80, 4403 (1984).
48. L. E. Brus, *J. Phys. Chem.* 90, 2444 (1986).
49. G. Bastard, "Wave Mechanics Applied to Semiconductor Heterostructures," Les editions de physique. Less Ulis, 1988.
50. J. M. Luttinger and W. Kohn, *Phys. Rev.* 97, 869 (1955).
51. E. O. Kane, in "Semiconductor and Semimetals" (R. K. Willardson and A. C. Beer, Eds.), Vol. 1, p. 75. Academic Press, New York, 1966.
52. M. G. Burt, *J. Phys. Condens. Matter* 4, 6651 (1992).
53. A. L. Efros and A. V. Rodina, *Phys. Rev. B* 47, 10005 (1993).
54. S. Furukawa and T. Miyasato, *Phys. Rev. B* 38, 5726 (1988).
55. T. Takagahara and K. Takeda, *Phys. Rev. B* 46, 15578 (1992).
56. M. A. Cusack, P. R. Briddon, and M. Jaros, *Phys. Rev. B* 54, R2300 (1996).
57. G. A. Baraff and D. Gershoni, *Phys. Rev. B* 43, 4011 (1991).
58. P. C. Sercel and K. J. Vahala, *Phys. Rev. B* 42, 3690 (1990).
59. Al. L. Efros and M. Rosen, *Phys. Rev. B* 58, 7120 (1998).
60. H. Fu, L. W. Wang, and A. Zunger, *Phys. Rev. B* 57, 9971 (1998).
61. L. W. Wang and A. Zunger, *J. Phys. Chem. B* 102, 6449 (1998).
62. L. W. Wang, *Phys. Rev. B* 61, 7241 (2000).
63. Y. M. Niquet, C. Delerue, G. Allan, and M. Lannoo, *Phys. Rev. B* 62, 5109 (2000).
64. F. A. Reboredo, A. Franceschetti, and A. Zunger, *Phys. Rev. B* 61, 13073 (2000).
65. S. B. Zhang, C. Yeh, and A. Zunger, *Phys. Rev. B* 48, 11204 (1993).
66. S. Y. Ren, *Phys. Rev. B* 64, 35322 (2001).
67. S. Y. Ren, *Ann. Phys.* (San Diego), 301, 22 (2002).
68. M. V. Rama Krishna and R. A. Friesner, *Phys. Rev. Lett.* 67, 629 (1991).
69. L. W. Wang and A. Zunger, *J. Phys. Chem.* 98, 2158 (1994).
70. S. Y. Ren, *Phys. Rev. B* 55, 4665 (1997).
71. S. Y. Ren and J. D. Dow, *Phys. Rev. B* 45, 6492 (1992).
72. N. A. Hill and K. B. Whaley, *J. Chem. Phys.* 99, 3707 (1993).
73. N. A. Hill and K. B. Whaley, *J. Chem. Phys.* 100, 2831 (1994).
74. N. A. Hill and K. B. Whaley, *Phys. Rev. Lett.* 75, 1130 (1995).
75. H. F. Troter, *Proc. Am. Math. Soc.* 10, 545 (1959).
76. M. D. Feit, J. A. Fleck, and A. Steiger, *J. Comput. Phys.* 47, 413 (1982).
77. N. A. Hill and K. B. Whaley, *J. Chem. Phys.* 210, 117 (1996).
78. K. Leung and K. B. Whaley, *Phys. Rev. B* 56, 7455 (1997).
79. P. Vogl, H. P. Hjalmarson, and J. D. Dow, *J. Phys. Chem. Solids* 44, 365 (1983).
80. C. Delerue, M. Lannoo, and G. Allan, *Phys. Rev. Lett.* 76, 3038 (1996); N. A. Hill and K. B. Whaley, *ibid.* 76, 3039 (1996).
81. C. Delerue, G. Allan, and M. Lannoo, *Phys. Rev. B* 48, 11024 (1993).
82. C. Tserbak, H. M. Polatoglou, and G. Theodorou, *Phys. Rev. B* 47, 7104 (1993).
83. W. A. Harrison, "Electronic Structure and the Properties of Solids." Freeman, San Francisco, 1980.
84. M. L. Cohen and T. K. Bergstresser, *Phys. Rev.* 141, 789 (1966).
85. M. L. Cohen and J. R. Chelikowsky, Eds., "Electronic Structure and Optical Properties of Semiconductors." Springer, Berlin, 1988.
86. M. P. Teter, M. C. Payne, and D. C. Allan, *Phys. Rev. B* 40, 12255 (1989).
87. L. W. Wang and A. Zunger, *J. Chem. Phys.* 100, 2394 (1994).
88. L. W. Wang, *Phys. Rev. B* 49, 10154 (1994).
89. M. Welkowsky and R. Braunstein, *Phys. Rev. B* 5, 497 (1972).
90. W. E. Spicer and R. C. Enden, in "Proceedings of the Ninth International Conference of the Physics of Semiconductors" (S. M. Ryvkin, Ed.), Vol. 1, p. 65. Nauka, Moscow, 1968.
91. R. Hulthen and N. G. Nilsson, *Solid State Commun.* 18, 1341 (1976).
92. W. Bludau, A. Onton, and W. Heinke, *J. Appl. Phys.* 45, 1846 (1974).
93. L. Ley, S. P. Kowalcyzk, R. A. Pollak, and D. A. Shirley, *Phys. Rev. Lett.* 29, 1088 (1972).
94. J. C. Hensel, H. Hasegawa, and M. Nakayama, *Phys. Rev.* 138, A225 (1965).
95. G. Dzesselhaus, A. F. Kip, and C. Kittel, *Phys. Rev.* 98, 368 (1955).
96. F. G. Allen, *J. Phys. Chem. Solids* 8, 119 (1959).
97. L. W. Wang and A. Zunger, in "Semiconductor Nanoclusters" (P. V. Kamat and D. Meisel, Eds.), Studies in Surface Science and Catalysis, Vol. 103, p. 161. Elsevier Science, Amsterdam, 1996.
98. D. R. Hamann, M. Schluter, and C. Chiang, *Phys. Rev. Lett.* 43, 1494 (1979); D. R. Hamann, *Phys. Rev. B* 40, 2980 (1989).
99. R. W. Goodby, M. Schulter, and L. J. Sham, *Phys. Rev. B* 37, 10159 (1988).
100. M. S. Hybertsen and S. Louie, *Phys. Rev. B* 34, 5390 (1986).

*101.* K. C. Pandey, *Phys. Rev. B* 14, 1557 (1976).

*102.* J. J. Boland, *Phys. Rev. Lett.* 65, 3325 (1990).

*103.* E. Kaxiras and J. D. Joannopoulos, *Phys. Rev. B* 37, 8842 (1988).

*104.* J. E. Northrup, *Phys. Rev. B* 44, 1419 (1991).

*105.* T. Sakurai and H. D. Hagstrum, *Phys. Rev. B* 12, 5349 (1975).

*106.* T. Sakurai and H. D. Hagstrum, *J. Vac. Sci. Technol.* 13, 807 (1976).

*107.* S. Maruno, H. Iwasaki, K. Horioka, S. T. Li, and S. Nakamura, *Phys. Rev. B* 27, 4110 (1983).

*108.* L. W. Wang and A. Zunger, *Phys. Rev. B* 51, 17398 (1995).

*109.* H. Fu and A. Zunger, *Phys. Rev. B* 55, 1642 (1997).

*110.* C. Y. Yeh, S. B. Zhang, and A. Zunger, *Appl. Phys. Lett.* 63, 3455 (1993).

*111.* A. Canning, L. W. Wang, A. Williamson, and A. Zunger, *J. Comp. Phys.* 160, 29 (2000).

*112.* L. W. Wang, J. N. Kim, and A. Zunger, *Phys. Rev. B* 59, 5678 (1999).

*113.* G. Allan, Y. M. Niquet, and C. Delerue, *Appl. Phys. Lett.* 77, 639 (2000).

*114.* P. Hohenberg and W. Kohn, *Phys. Rev.* 136, B864 (1964).

*115.* W. Kohn and L. J. Sham, *Phys. Rev.* 140, A1133 (1965).

*116.* N. Trouillier and J. L. Martins, *Phys. Rev. B* 43, 1993 (1991).

*117.* M. C. Payne, M. P. Teter, D. C. Allan, T. A. Arias, and J. D. Joannopoulos, *Rev. Mod. Phys.* 64, 1045 (1992).

*118.* B. Delley and E. F. Steigmeier, *Phys. Rev. B* 47, 1397 (1993).

*119.* J. R. Chelikowsky, N. Troubllier, and Y. Saad, *Phys. Rev. Lett.* 72, 1240 (1994).

*120.* S. Ogut, J. R. Chelikowsky, and S. G. Louie, *Phys. Rev. Lett.* 79, 1770 (1997).

*121.* S. Goedecker, *Rev. Mod. Phys.* 71, 1085 (1999).

*122.* J.-L. Fattebert and J. Bernholc, *Phys. Rev. B* 62, 1713 (2000).

*123.* L. W. Wang, *Appl. Phys. Lett.* 78, 1565 (2001).

*124.* L. W. Wang, *Phys. Rev. B* 65, 153410 (2002).

*125.* L. W. Wang, *Phys. Rev. Lett.* 88, 256402 (2002).

*126.* B. Delley and E. F. Steigmeier, *Appl. Phys. Lett.* 67, 2370 (1995).

*127.* M. Hirao, T. Uda, and Y. Murayama, *Mater. Res. Soc. Symp. Proc.* 283, 425 (1993).

*128.* A. J. Read, R. J. Needs, K. J. Nash, L. T. Canham, P. D. J. Calcott, and A. Qteish, *Phys. Rev. Lett.* 69, 1232 (1992).

*129.* F. Buda, J. Kohanoff, and M. Parrinello, *Phys. Rev. Lett.* 69, 2400 (1992).

*130.* E. K. U. Gross, F. J. Dobson, M. Petersilka, "Density Functional Theory," Springer, New York, 1996.

*131.* R. van Leeuwen, *Int. J. Mod. Phys. B* 15, 1969 (2001).

*132.* K. Burke, M. Petersilka, and E. K. U. Gross, in "Recent Advances in Density Functional Methods" (P. Fantucci, A. Bencini, Eds.), Vol. 3, World Scientific, Singapore, 2002.

*133.* K. Yabana and G. F. Bertsch, *Phys. Rev. B* 54, 4484 (1996).

*134.* M. E. Casida, in "Recent Advances in Density Functional Methods, Part I," (D. P. Chong, Ed.), p. 155, World Scientific, Singapore, 1995.

*135.* L. X. Benedict, A. Puzder, A. J. Williamson, J. C. Grossman, G. Galli, J. E. Klepeis, J. Y. Raty, and O. Pankratov, *Phys. Rev. B* 68, 85310 (2003).

*136.* I. V. Tokatly and O. Pankratov, *Phys. Rev. Lett.* 86, 2078 (2001); I. V. Tokatly, R. Stubner, and O. Pankratov, *Phys. Rev. B* 65, 113107 (2002).

*137.* P. J. Stephens, F. H. Devlin, C. F. Chabalowski, and M. J. Frisch, *J. Phys. Chem.* 98, 11623 (1994).

*138.* A. D. Decke, in "The Challenge of *d* and *f* Electrons: Theory and Computation" (D. R. Salahub and M. C. Zerner, Eds.), Chap. 12, pp. 165–179, American Chemical Society, Washington, DC, 1989.

*139.* C. Lee, W. Yang, and R. G. Parr, *Phys. Rev. B* 37, 785 (1988).

*140.* S. H. Vosko, L. Wilk, and M. Nusair, *Can. J. Phys.* 58, 1200 (1980).

*141.* J. Muscat, A. Wander, and N. M. Harrison, *Chem. Phys. Lett.* 342, 397 (2001).

*142.* C. S. Garoufalis, A. D. Zdetsis, and S. Grimme, *Phys. Rev. Lett.* 87, 276402 (2001).

*143.* A. J. Williamson, J. C. Grossman, R. Q. Hood, A. Puzder, and G. Galli, *Phys. Rev. Lett.* 89, 196803 (2002).

*144.* Y. Abe, Y. Osaka, and A. Morita, *J. Phys. Soc. Jpn.* 17, 1576 (1962).

*145.* L. J. Sham and T. M. Rice, *Phys. Rev.* 144, 708 (1966).

*146.* G. Strinati, *Phys. Rev. B* 29, 5718 (1984).

*147.* L. X. Benedict, *Phys. Rev. B* 66, 193105 (2002).

*148.* A. Franceschetti, L. W. Wang, H. Fu, and A. Zunger, *Phys. Rev. B* 58, R13367 (1998).

*149.* L. C. Andreani, F. Bassani, and A. Quattropani, *Lett. Nuovo Cimento* 10, 1473 (1988).

*150.* K. Cho, *Solid State Commun.* 33, 911 (1980).

*151.* V. A. Kiselev and A. G. Zhilich, *Fiz. Tverd. Tela (Leningrad)* 14, 1438 (1972) [*Sol. Phys. Solid State* 14, 1233 (1972)]; 15, 2024 (1972) [15, 1351 (1972)].

*152.* G. Bester, S. Nair, and A. Zunger, *Phys. Rev. B* 67, 161306 (2003).

*153.* L. W. Wang, M. Califano, A. Zunger, and A. Franceschetti, *Phys. Rev. Lett.* 91, 56404 (2003).

*154.* B. Farid, in "Electron Correlation in the Solid State," (N. H. March, Ed.), p. 103, World Scientific, Singapore, 1999.

*155.* G. Onida, L. Reining, and A. Rubio, *Rev. Mod. Phys.* 74, 601 (2002).

*156.* L. Hedin, *Phys. Rev.* 139, A796 (1965).

*157.* L. Hedin, S. Lundqvist, in "Solid State Physics, Advances in Research and Application" (F. Seitz, D. Turnbull and H. Ehrenreich, Eds.), Vol. 23, p. 1, Academic Press, New York, 1969.

*158.* P. Sanchez-Friera and R. W. Godby, *Phys. Rev. Lett.* 85, 5611 (2000).

*159.* P. Garcia-Gonzalez and R. W. Godby, *Phys. Rev. Lett.* 88, 56406 (2002).

*160.* F. Aryasetiawan, T. Miyak. and K. Terakura. *Phys. Rev. Lett.* 88, 166401 (2002).

*161.* C. Verdozzi, R. W. Godby, and S. Holloway, *Phys. Rev. Lett.* 74, 2327 (1995).

*162.* T. J. Pollehn, A. Schindlmayr, and R. W. Godby, *J. Phys. Condens. Matter* 10, 1273 (1998).

*163.* A. Schindlmayr, T. J. Pollehn, and R. W. Godby, *Phys. Rev. B* 58, 12684 (1998).

*164.* M. Rohlfing and S. G. Louie, *Phys. Rev. B* 62, 4927 (2000).

*165.* A. Franceschetti and A. Zunger, *Phys. Rev. B* 62, 2614 (2000).

*166.* C. Delerue, M. Lannoo, and G. Allan, *Phys. Rev. Lett.* 84, 2457 (2000).

*167.* M. Rohlfing and S. G. Louie, *Phys. Rev. Lett.* 80, 3320 (1998).

*168.* R. W. Godby and I. D. White, *Phys. Rev. Lett.* 80, 3161 (1998); S. Ogut, J. R. Chelikowsky, and S. G. Louie, *Phys. Rev. Lett.* 80, 3162 (1998); A. Franceschetti, L. W. Wang, and A. Zunger, *Phys. Rev. Lett.* 83, 1269 (1999); S. Ogut, J. R. Chelikowsky, and S. G. Louie, *Phys. Rev. Lett.* 83, 1270 (1999).

*169.* W. M. C. Foulkes, L. Mitas, R. J. Needs, and G. Rajagopal, *Rev. Mod. Phys.* 73, 33 (2001).

*170.* P. J. Reynolds, D. M. Ceperley, B. J. Alder, and W. A. Lester, Jr., *J. Chem. Phys.* 77, 5593 (1982).

*171.* C. J. Umrigar. K. G. Wilson, and J. W. Wilkins, *Phys. Rev. Lett.* 60, 1719 (1988).

*172.* D. Ceperley, G. V. Chester, and M. H. Kalos, *Phys. Rev. B* 16, 3081 (1971).

*173.* B. L. Hammond, W. A. Lester, Jr., and P. J. Reynolds, "Monte Carlo Methods in *Ab Initio* Quantum Chemistry." World Scientific, Singapore, 1994.

*174.* S. Fahy, X. W. Wang, and S. G. Louie, *Phys. Rev. B* 42, 3503 (1990).

*175.* A. J. Williamson, R. Q. Hood, R. J. Needs, and G. Rajagopal. *Phys. Rev. B* 57, 12140 (1998).

*176.* J. C. Grossman, M. Rohlfing, L. Mitas, S. G. Louie, and M. L. Cohen, 86, 472 (2001).

*177.* A. J. Williamson, R. Q. Hood, and J. C. Grossman. *Phys. Rev. Lett.* 87, 246406 (2001).

*178.* A. Puzder, A. J. Williamson, J. C. Grossman, and G. Galli, *Phys. Rev. Lett.* 88, 97401 (2002).

*179.* I. Vasiliev, S. Ogut, and J. R. Chelikowsky, *Phys. Rev. Lett.* 86, 1813 (2001).

*180.* T. van Buuren, L. N. Dinh, L. L. Chase, W. J. Siekhaus, and L. J. Terminello, *Phys. Rev. Lett.* 80, 3803 (1998).

*181.* T. van Buuren, T. Tiedje, J. R. Dahn, and B. M. Way, *Appl. Phys. Lett.* 63, 2911 (1993).

*182.* L. W. Wang and A. Zunger, *Phys. Rev. Lett.* 73, 1039 (1994).

*183.* R. Tsu, L. Ioriatti, J. F. Harvey, H. Shen, and R. A. Lux, *Mater. Res. Soc. Symp. Proc.* 283, 437 (1993).

*184.* D. R. Penn, *Phys. Rev.* 128, 2093 (1962).

*185.* H. Ehrenreigh and M. H. Cohen, *Phys. Rev.* 115, 786 (1959).

*186.* D. E. Aspnes and A. A. Studna, *Phys. Rev. B* 27, 985 (1983).

*187.* M. Lannoo, C. Delerue, and G. Allan, *Phys. Rev. Lett.* 74, 3415 (1995).

*188.* T. D. Krauss and L. E. Brus, *Phys. Rev. Lett.* 83, 4840 (1999).

*189.* L. W. Wang and A. Zunger, *Phys. Rev. B* 53, 9579 (1996).

*190.* S. Ogut, R. Burdick, Y. Saad, and J. R. Chelikowsky, *Phys. Rev. Lett.* 90, 127401 (2003).

*191.* Y. H. Xie, W. L. Wilson, H. M. Ross, J. A. Mucha, M. A. Fitzgerald, J. M. Macaulay, and T. D. Hariss, *J. Appl. Phys.* 71, 2403 (1992).

*192.* J. C. Vial, A. Bsiesy, F. Gaspard, R. Herino, M. Ligeon, F. Muller, and R. Romestain, *Phys. Rev. B* 45, 14171 (1992).

*193.* M. S. Hybertsen, *Phys. Rev. Lett.* 72, 1514 (1994).

*194.* C. Delerue, G. Allan, and M. Lannoo, *Phys. Rev. B* 64, 193402 (2001).

*195.* P. D. J. Calcott, K. J. Nash, L. T. Canham, M. J. Kane, D. Brumhead, *J. Phys. Condens. Matter* 5, L91 (1993).

*196.* E. Martin, C. Delerue, G. Allan, and M. Lannoo, *Phys. Rev. B* 50, 18258 (1994).

*197.* M. L. Brongersma, P. G. Kik, A. Polman, K. S. Min, and H. A. Atwater, *App. Phys. Lett.* 76, 351 (2000).

*198.* F. A. Reboredo, A. Franceschetti, and A. Zunger, *Appl. Phys. Lett.* 75, 2972 (1999).

*199.* J. P. Wilcoxon, G. A. Samara, and P. N. Provencio, *Phys. Rev. B* 60, 2704 (1999).

*200.* M. J. Caldas, *Phys. Stat. Sol. B* 217, 641 (2000).

*201.* G. Allan, C. Delerue, and M. Lannoo, *Phys. Rev. Lett.* 76, 2961 (1996).

*202.* A. Puzder, A. J. Williamson, J. C. Grossman, and G. Galli, *J. Chem. Phys.* 117, 6721 (2002).

*203.* M. Luppi and S. Ossicini, *J. Appl. Phys.* 94, 2130 (2003).

*204.* A. B. Filonov, S. Ossicini, F. Bassani, and F. Arnaud d'Avitaya, *Phys. Rev. B* 65, 195317 (2002).

*205.* I. Vasiliev, J. R. Chelikowsky, and R. M. Martin, *Phys. Rev. B* 65, 121302 (2002).

*206.* S. Schuppler, S. L. Friedman, M. A. Marcus, D. L. Adler, Y. H. Xie, F. M. Ross, T. D. Harris, W. L. Brown, Y. J. Chabal, L. E. Brus, and P. H. Citrin, *Phys. Rev. Lett.* 72, 2648 (1994).

*207.* L. N. Dinh, W. McLean. II. M. A. Schildbach, and M. Balooch, *Phys. Rev. B* 59, 15513 (1999).

# CHAPTER 13

# Nanoscale Device Modeling

## Massimo Macucci, Luca Bonci

*Dipartimento di Ingegneria dell'Informazione, Università degli studi di Pisa, Pisa, Italy*

## CONTENTS

## 1. INTRODUCTION

Modeling of nanoscale devices is currently acquiring growing importance and is reaching maturity, in the sense that quantitative predictive capabilities have been achieved. In particular, modeling has become instrumental in the design of new devices, in the interpretation of experimental results, and in the evaluation of new device proposals.

In particular, this latter issue, the evaluation of new device proposals, has become extremely relevant in the last few years, after realizing that very significant research efforts have been spent in the development of technologies whose intrinsic weaknesses could have been detected from the very beginning, if realistic models had been considered instead of very simple, idealized models, not including real-life nonidealities, such as fabrication tolerances, external interferences, and asymmetries.

It is at the same time apparent that a wide range of simulation tools is needed at different levels of approximation: realistic models are required for single devices, to verify their operability in a variety of conditions and in the presence of imperfections. Higher-level, more approximate models must instead be developed to investigate circuits including from a few to millions or billions of devices. Because of computational constraints, detailed models cannot treat more than a few devices at a time; therefore, a hierarchy of simulation tools must be created, ranging from atomistic approaches to logic level simulations.

In Section 2 a broad overview of the modeling techniques used for nanodevices will be presented. It is indeed impossible, within the space of this chapter, to cover all of these techniques in depth, therefore, we have chosen to discuss in detail several selected topics in the remaining sections.

In Section 3 we discuss a few techniques for the analysis of ballistic transport, starting from those initially developed for the investigation of disordered materials and of quantum interference devices. In particular, we present approaches for the calculation of the transmission and reflection matrices of a generic nanostructure, based on the recursive Green's function, the mode-matching, and the recursive scattering matrix techniques. We treat also the inclusion of the effect of a magnetic field, discussing advantages and shortcomings of a few methods that can be found in the literature.

Section 4 deals with the simulation of ballistic quantum wires defined by realistic potentials, with a discussion of some simplified semianalytic techniques allowing the calculation of the potential in a gated heterostructure. The subtle issue of the boundary conditions at the exposed surface will be examined in detail, presenting the different approaches that have been proposed in the literature.

Section 5 will focus on the simulation of quantum dots and generic nanostructures with confinement in all spatial directions. Techniques based on finite differences will be discussed for the solution of the Schrödinger equation. Then the problem of achieving self consistency with the solution of the Poisson equation will be addressed, introducing also the underrelaxation technique.

In Section 6 we move on to the simulation of "artificial molecules" (i.e., systems made up of multiple interacting quantum dots). These structures require particular attention because of the presence of quasidegenerate states and the relevance of the electrostatic interaction, which cannot be simply considered as a perturbation of the confinement energy. In such conditions, iterative methods often fail to converge, and if they do not include a properly constructed many-body wave function (as in the case of a density functional approximation), they may yield erroneous results. We discuss the method based on a Hubbard-like Hamiltonian and then a more general and complex approach based on an implementation of the configuration–interaction technique commonly used in molecular chemistry.

Finally, Section 7 deals with the simulation of single-electron circuits on the basis of the "orthodox" Coulomb blockade theory. In particular, we focus on the Monte Carlo approach, the master equation technique, and the approximate introduction of cotunneling into simulation codes.

## 2. OVERVIEW OF MODELING TECHNIQUES

In comparison to the simulation of traditional electronic circuits, such as those based on CMOS (complementary metal–oxide–semiconductor) transistors, in the case of nanoscale devices the hierarchy of tools must extend further toward the bottom end (i.e., more refined, physics-based tools are necessary). Atomistic *ab initio* simulation [1], starting from first principles, would be ideal, but it takes up huge computational resources as soon as structures are larger than 10–20 nm (corresponding to over 100,000 atoms); therefore, even with the most powerful supercomputers, it is applicable only to a very limited number of cases. A step further in simplification is represented by some tight-binding approaches [2], which are still very difficult to apply to an actual device, although at the nanoscale, and then by effective mass approximations, which can be single-band or multi-band, as in the case of the $k \cdot p$ method [3].

Any of the just-mentioned approaches requires in general the solution of a self-consistent Schrödinger-Poisson problem, with a level of complexity depending on the specific

characteristics of the device and on the relative importance of the electrostatic interaction with respect to the quantum confinement energy. Significant simplification in transport problems can be achieved using so called hard-wall models, based on the assumption that the device is delimited by infinite potential barriers. In such cases, it is often possible to use analytical solutions of the Schrödinger equation in each region, which significantly reduces the computational effort while preserving most of the relevant physics.

There are also techniques that, although introducing further simplifications in the description of the device structure, include a more precise treatment of specific physical aspects. An example is represented by the Luttinger liquid formulation of quantum wire problems [4, 5], in which a very abstract representation of the involved nanostructure is accompanied by a detailed treatment of electron correlations.

Another important issue, for which a completely satisfactory solution has not been found yet, is the inclusion of dephasing and dissipation phenomena. Basic models are ballistic (i.e. they include an elastic treatment), which is adequate as long as the device is very small and the temperature very low. Actual devices are, however, significantly affected by dephasing and dissipation phenomena [6], with the former being particularly important in the growing field of quantum computing. Dephasing and dissipation are most often the results of the interaction between the quantum system that we are analyzing and the environment, which is very difficult to model in a realistic way. Approaches such as those based on the Caldeira–Leggett model [7] for dissipation are usually not feasible from a computational point of view because they would require us to introduce interaction with too large a number of harmonic oscillators; other approaches, based on nonequilibrium Green's function [8–10] and on the Keldish formalism do indeed find some computational application, but they are limited to very simple structures. Other, more phenomenological, methods have also been proposed. Büttiker has introduced a technique based on the insertion of a voltage probe [11], which acts as a phase-breaking element; others have included partial dephasing by adding a random phase to the elements of the scattering matrix of the device or considering a complex confinement potential.

Moving up in the ladder of simulation levels, circuits must be treated on the basis of simplified models, as in the case of classical electron devices, with the difference that often compact models (such as for MOS transistors) cannot be derived or cannot be used in a SPICE-like architecture [12]. A typical example is represented by single-electron circuits, for which a SPICE description is not possible, unless capacitances much larger than those of the tunnel junctions are present at the nodes connecting devices (so that Coulomb Blockade effects are quenched at nodes interconnecting different devices), and new concepts had to be derived for circuit simulation, based on a Monte Carlo [13] or on a master equation [14] approach.

The same is true for innovative architectures, such as the Quantum Cellular Automaton (QCA) concept, which was proposed by Craig Lent and coworkers [15] for the implementation of logic circuits based on a bistable building block. In this and in analogous cases, simulation programs must be developed ad hoc, based on the specific characteristics of the circuit. As the number of basic cells in the circuit grows, increasingly radical semiclassical approximations must be introduced.

## 3. TRANSPORT IN BALLISTIC SYSTEMS

The methods treated in this section were initially developed in the 1980s for the analysis of quantum interference devices or of the effects of disorder in conductors. In particular, the recursive Green's function approach was first introduced to study weak localization effects. All of the methods that we will discuss in this section are used to essentially solve a scattering problem, computing the transmission and the reflection matrices through the nanostructure being considered. From the knowledge of such matrices it is possible to obtain the conductance of the device in quasiequilibrium conditions, on the basis of the Landauer-Büttiker formula [16, 17]:

$$G = 2\frac{e^2}{h} \sum_n |t_n|^2 \tag{1}$$

where $e$ is the electron charge. $h$ is Planck's constant and $t_{i,j}$ is the element of the transmission matrix between input mode $j$ and output mode $i$. In this case, the transmission matrix is represented over a basis of transverse eigenmodes of the input and output channels. The assumption that is usually made for numerical treatment is that semiinfinite ideal channels are attached to the input and output ports of the device, to implement properly transmitting boundary conditions.

It has been shown that further information can be obtained from the transmission matrix, such as the shot noise power spectral density. Following the pioneering work by Khlus [18] and Lesovik [19], Büttiker provided a detailed treatment of the relationship between the shot noise current power spectral density $S_I$ and the transmission matrix $T$ [20], the result of which is expressed by the relationship

$$S_I = 4\frac{e^2}{h}|eV|\mathrm{Tr}\big[t^\dagger t(I - t^\dagger t)\big] \qquad (2)$$

where $V$ is the voltage applied across the sample, and that can be also written as

$$S_I = 4\frac{e^2}{h}|eV|\sum_i T_i(1 - T_i) \qquad (3)$$

with the $T_i$s being the eigenvalues of the $tt^\dagger$ matrix.

The importance of having available efficient methods for the numerical calculation of the transmission matrix is therefore apparent.

## 3.1. Recursive Green's Function Technique

The foundations for the recursive Green's function technique were set by Thouless and Kirkpatrick [21] and by Lee and Fisher [22] for the investigation of disordered conductors. Guinea and Vergés [23] extended this study to the specific case of a one-dimensional chain with a random distribution of linear side branches and loops. They introduced a formalism based on the analytical derivation of the Green's functions for an infinite linear chain and for a finite segment of a linear chain, and on their numerical combination, with an approach derived from Dyson's equation. Such a technique was then extended to the investigation of two-dimensional structures by Sols et al. [24] and by Macucci et al. [25]. In the next subsection, we will discuss this further evolution of the method in detail.

For the application of the recursive Green's function method we consider the whole surface (in the two-dimensional case) or volume (in the three-dimensional case) of the device to be filled with a rectangular or parallelepiped tight-binding lattice, substantially equivalent to a lattice for a finite-difference scheme. In the methods originally derived for the analysis of disordered conductors, the Green's function for a semiinfinite lead at the output is derived analytically and then a layer of the discretization lattice is added at each step of the recursive procedure [26]. If, however, there are sections of the device characterized by a longitudinally constant transverse potential, they can be treated as single sections, for which the Green's functions can be derived analytically.

Therefore, in our approach the device is further subdivided into sections, each of which is assumed to be characterized by a longitudinally constant transverse potential, as shown in Fig. 1, to have Dirichlet boundary conditions at the ends (therefore to be "sealed"), and to



Figure 1. Subdivision of a structure into slices with constant transverse potential.

overlap the neighboring section by one lattice unit. If we consider two neighboring sections, the overall Hamiltonian can be written as

$$H = H_0 + V \tag{4}$$

where $H_0$ is the Hamiltonian for the two decoupled sections and $V$ is the perturbation corresponding to coupling them together (thereby opening up the facing ends of the two sections). This is quite a large perturbation, as it corresponds to a substantial variation of the system geometry. Thus, it cannot be treated with perturbation theory including only low-order terms. All terms must be taken into account. This is possible using Dyson's equation, which connects the Green's function of the unperturbed system $G_0$ (decoupled sections) with the Green's function for the perturbed system $G$ (joint sections):

$$G = G_0 VG \tag{5}$$

This is an implicit equation, as the Green's function $G$ appears on both sides of the equal sign. It is possible, with some algebra, to obtain an explicit expression for $G$. We will discuss it with a specific choice of the Green's function representation. Specifically, we choose a mixed representation: in real space in the longitudinal direction and in the space of transverse eigenmodes in the transverse direction. Thus, the matrix representing the Green's function $G_{ij}$ between two locations $i$ and $j$ along the device has elements $\langle i, s|G|t, j \rangle$, where $s$ is a transverse eigenmode at location $i$ and $t$ is a transverse eigenmode at location $j$. If each section has, as previously stated, a longitudinally constant transverse potential, there is no mode mixing within it; therefore, the matrix representing its Green's function will be diagonal (i.e. $\langle i, s|G|t, j \rangle = 0$ for $s \neq t$). Each element $\langle i, s|G|s, j \rangle$ can be treated as the Green's function for a one-dimensional discrete chain extending from $i$ to $j$, with an energy corresponding to the longitudinal energy component available for that specific mode. Because of the separability of the Schrödinger equation along the longitudinal and the transverse direction, it is possible to subdivide the total energy of the impinging particles into a transverse component, corresponding to the energy eigenvalue associated with the particular transverse mode we are considering and into a longitudinal component, equal to the difference between the total energy and the just-mentioned transverse component. For each section, the problem is thereby transformed into a collection of noninteracting one-dimensional problems.

Before discussing the one-dimensional Green's functions in detail, let us examine the procedure for the determination of the explicit expression of the Green's function for a chain $a - d$ resulting from the connection of two sections $a - b$ and $c - d$ (see Fig. 2), for each of which the unperturbed Green's functions $G_{aa}^0$, $G_{ab}^0$, $G_{cc}^0$, $G_{bb}^0$, and $G_{cd}^0$ have been obtained analytically. Let us start from a concrete representation of the Dyson equation for our specific case:

$$\langle a|G|d \rangle = \langle a|G^0|d \rangle + \langle a|G^0 VG|d \rangle \tag{6}$$

which, using the completeness relationships (i.e., $\sum_i |i\rangle\langle i| = I$, with $|i\rangle$ being the elements of a complete basis) can be rewritten

$$\langle a|G|d \rangle = \langle a|G^0|d \rangle + \sum_{m,n} \langle a|G^0|m\rangle\langle m|V|n\rangle\langle n|G|d \rangle \tag{7}$$



Figure 2. Green's functions in two neighboring sections.

Because the perturbation $V$ acts only between sites $b$ and $c$, there are only two nonzero terms in the sum over $m, n$:

$$\langle a|G|d\rangle = \langle a|G^0|d\rangle + \langle a|G^0|b\rangle\langle b|V|c\rangle\langle c|G|d\rangle + \langle a|G^0|c\rangle\langle c|V|b\rangle\langle b|G|d\rangle \tag{8}$$

which, considering that, before $V$ is applied, there is no connection between $a$ and $d$ and between $a$ and $c$, simplifies to

$$\langle a|G|d\rangle = \langle a|G^0|b\rangle\langle b|V|c\rangle\langle c|G|d\rangle \tag{9}$$

as $\langle a|G^0|d\rangle = \langle a|G^0|c\rangle = 0$. In shorthand notation we can write

$$G_{ad} = G^0_{ab}V_{bc}G_{cd} \tag{10}$$

which, however, is still an implicit equation, as it contains the unknown term $G_{cd}$. To obtain an explicit expression for $G_{ad}$, we need to do some more algebra and in particular to expand the expressions for $G_{cd}$ and $G_{bd}$:

$$\langle c|G|d\rangle = \langle c|G^0|d\rangle + \langle c|G^0 VG|d\rangle$$
$$= \langle c|G^0|d\rangle + \sum_{m,n}\langle c|G^0|m\rangle\langle m|V|n\rangle\langle n|G|d\rangle$$
$$= \langle c|G_0|d\rangle + \langle c|G|c\rangle\langle c|V|b\rangle\langle b|G|d\rangle \tag{11}$$

which can be written as

$$G_{cd} = G^0_{cd} + G^0_{cc}V_{cb}G_{bd} \tag{12}$$

As far $G_{bd}$ is concerned, we get

$$\langle b|G|d\rangle = \langle b|G^0|d\rangle + \langle b|G^0 VG|d\rangle$$
$$= \langle b|G^0|d\rangle + \sum_{m,n}\langle b|G^0|m\rangle\langle m|V|n\rangle\langle n|G|d\rangle$$
$$= \langle b|G^0|b\rangle\langle b|V|c\rangle\langle c|G|d\rangle \tag{13}$$

which, in simplified notation, reads

$$G_{bd} = G^0_{bb}V_{bc}G_{cd} \tag{14}$$

Substituting Eq. (14) into Eq. (12), we obtain

$$G_{cd} = G^0_{cd} + G^0_{cc}V_{cb}G^0_{bb}V_{bc}G_{cd} \tag{15}$$

which can be rewritten as

$$(I - G^0_{cc}V_{cb}G^0_{bb}V_{bc})G_{cd} = G^0_{cd} \tag{16}$$

If we multiply both sides by the inverse of $I - G^0_{cc}V_{cb}G^0_{bb}V_{bc}$, we get an explicit expression for $G_{cd}$

$$G_{cd} = (I - G^0_{cc}V_{cb}G^0_{bb}V_{bc})^{-1}G^0_{cd} \tag{17}$$

which can be substituted into Eq. (10), yielding an explicit expression for $G_{ad}$

$$G_{ad} = G^0_{ab}V_{bc}(I - G^0_{cc}V_{cb}G^0_{bb}V_{bc})^{-1}G^0_{cd} \tag{18}$$

To be able to attach one further section to the left of $a - d$ and to compute the overall Green's function, we also need to evaluate $G_{aa}$ an explicit expression of which can be derived with a procedure completely analogous to the one just presented. We start from the implicit expression for $G_{aa}$ obtained from the Dyson equation

$$G_{aa} = G^0_{aa} + G^0_{ab}V_{bc}G_{ca} \tag{19}$$

then obtain the expression for $G_{ca}$

$$G_{ca} = G_{cc}^{0} V_{cb} G_{ba}$$ (20)

and for $G_{ba}$

$$G_{ba} = G_{ba}^{0} + G_{bb}^{0} V_{bc} G_{ca}$$ (21)

Substituting Eq. (21) into Eq. (20), we get

$$G_{ca} = G_{cc}^{0} V_{cb} G_{ba}^{0} + G_{cc}^{0} V_{cb} G_{bb}^{0} V_{bc} G_{ca}$$ (22)

which can be rewritten as

$$(I - G_{cc}^{0} V_{cb} G_{bb}^{0} V_{bc}) G_{ca} = G_{cc}^{0} V_{cb} G_{ba}^{0}$$ (23)

so that, multiplying by the inverse of $I - G_{cc}^{0} V_{cb} G_{bb}^{0} V_{bc}$, we obtain an explicit expression for $G_{ca}$

$$G_{ca} = (I - G_{cc}^{0} V_{cb} G_{bb}^{0} V_{bc})^{-1} G_{cc}^{0} V_{cb} G_{ba}^{0}$$ (24)

which can be inserted into Eq. (19), yielding an explicit expression for $G_{aa}$

$$G_{aa} = G_{aa}^{0} + G_{ab}^{0} V_{bc} (I - G_{cc}^{0} V_{cb} G_{bb}^{0} V_{bc})^{-1} G_{cc}^{0} V_{cb} G_{ba}^{0}$$ (25)

Let us now discuss the derivation of the analytical expressions for the one-dimensional Green's functions for a semiinfinite chain and for a finite segment. In particular, we need the Green's function for the semiinfinite chain from its first site back to itself ($G_{11}^{s}$) and from its first site to a generic site $l$ ($G_{1l}^{s}$), as well as the Green's function for a finite chain with $N$ sites from its first site back to itself ($G_{11}^{f}$) and between its ends ($G_{1N}^{f}$).

We start from the Green's function between two generic sites ($l$ and $m$) in an infinite chain. For a discretized chain, it is given [27] by

$$G_{lm} = \frac{ie^{i|l-m|\theta}}{2V \sin \theta}$$ (26)

where $\theta = k\delta$, with $k$ being the wave vector for propagation, and $\delta$ is the discretization step.

To derive the Green's function for a semiinfinite chain, we can apply Dyson's equation backward, based on the structure shown in Fig. 3, in which an infinite chain is split between $a$ and $b$.

If we indicate as $G^{s}$ the Green's function for the semiinfinite chain and as $G$ the Green's function for the infinite chain, the Dyson equation reads

$$G = G^{s} + G^{s} V G$$ (27)

which, focusing on the pair of sites at $l$ and $m$, becomes

$$\langle l|G|m \rangle = \langle l|G^{s}|m \rangle + \langle G^{s} V G|m \rangle$$ (28)

Inserting the completeness relation and considering that the perturbation $V$ acts only between $a$ and $b$, we obtain

$$\langle l|G|m \rangle = \langle l|G^{s}|m \rangle + \langle l|G^{s}|a \rangle \langle a|V|b \rangle \langle b|G|m \rangle + \langle l|G^{s}|b \rangle \langle b|V|a \rangle \langle a|G|m \rangle$$ (29)



Figure 3. Infinite chain split into two semiinfinite sections.

and, as $\langle l|G^{\cdot}|a\rangle = 0$,

$$\langle l|G|m\rangle = \langle l|G^{\cdot}|m\rangle + \langle l|G^{\cdot}|b\rangle\langle b|V|a\rangle\langle a|G|m\rangle \tag{30}$$

which, with the already-used shorthand notation, becomes

$$G_{lm} = G^{\cdot}_{lm} + G^{\cdot}_{lb}V_{ba}G_{am} \tag{31}$$

Let us now assume that $m$ is coincident with $b$. In such a case

$$G_{lb} = G^{\cdot}_{lb} + G^{\cdot}_{lb}V_{ba}G_{ab} \tag{32}$$

Thus

$$G^{\cdot}_{lb} = G_{lb}(1 + VG_{ab})^{-1} \tag{33}$$

so that

$$
\begin{aligned}
G^{\cdot}_{lb} &= \frac{ie^{i(l-b)\theta}}{2V\sin\theta}\left(1 + V\frac{ie^{i(a-b)\theta}}{2V\sin\theta}\right)^{-1} \\
&= \frac{ie^{i(l-b)\theta}}{2V\sin\theta}\left(1 + \frac{ie^{i\theta}}{2\sin\theta}\right)^{-1} \\
&= \frac{1}{V}e^{i\theta}e^{i(l-b)\theta}
\end{aligned}
\tag{34}
$$

Therefore, the Green's function for a semiinfinite chain from the first location back to itself is given by

$$G^{\cdot}_{11} = \frac{1}{V}e^{i\theta} \tag{35}$$

whereas that between the first and the $l$th site will read

$$G^{\cdot}_{1l} = \frac{1}{V}e^{il\theta} \tag{36}$$

Let us now move to the Green's functions for a segment. We will follow a procedure analogous to that used to obtain the Green's functions for a semiinfinite chain from that for the infinite chain. We consider the structure shown in Fig. 4a: The semiinfinite chain extending to the right of site $b$ is interrupted between $l$ and $l + 1$. We can write the Green's function for the semiinfinite chain $G^{\cdot}$ by means of Dyson's equation:

$$G^{\cdot}_{lb} = \langle l|G''|b\rangle + \langle l|G''VG^{\cdot}|b\rangle \tag{37}$$

where $G_{0}$ is the Green's function for the structure disconnected between $l$ and $l + 1$. As usual, we apply the completeness relation and consider only the nonzero terms of the summation, obtaining

$$G^{\cdot}_{lb} = G^{0}_{lb} + G^{0}_{ll}VG^{\cdot}_{l+1,b} \tag{38}$$



Figure 4. (a) Semiinfinite chain with an interruption between sites $l$ and $l+1$

Let us now focus on the evaluation of the Green's function of a segment of length 1, which corresponds to setting $l = b$ and therefore allows an immediate derivation of an explicit expression:

$$G_{ll}^{\cdot} = G_{ll}^{0} + G_{ll}^{0} V G_{l+1,l}^{\cdot} \tag{39}$$

which yields

$$G_{ll}^{0} = G_{ll}^{\cdot}(1 + V G_{l+1,l}^{\cdot})^{-1} \tag{40}$$

By substituting Eqs. (35) and (36) into Eq. (40), we obtain

$$G_{ll}^{\cdot} = \frac{1}{2V\cos\theta} = \frac{\sin\theta}{V\sin 2\theta} \tag{41}$$

We can then evaluate the Green's function for a chain of two sites by combining those for two one-site chains. With reference to Fig. 4b, Eq. (18), the Green's function for the two combined one-site sections, will read

$$G_{aa} = \frac{1}{2V\cos\theta} V \left( 1 - \frac{1}{2V\cos\theta} V \frac{1}{2V\cos\theta} V \right)^{-1} \frac{1}{2V\cos\theta} = \frac{\sin\theta}{V\sin 3\theta} \tag{42}$$

It is also possible to show, using again Eq. (18), that if the Green's function between the ends $c$ and $d$ of a finite chain of length $N - 1$ is

$$G_{cd} = \frac{\sin\theta}{V\sin N\theta} \tag{43}$$

the Green's function for a chain obtained adding one site (and therefore of length $N$) is

$$G_{ad} = \frac{\sin\theta}{V\sin(N + 1)\theta} \tag{44}$$

which completes the proof of Eq. (44) by recursion. Actually, to combine the Green's function for a one-site chain with that for a $N - 1$-site chain, the Green's function for finite a chain from one end back to itself is also needed. It can be shown, with an analogous recursive procedure, that for a chain of length $M$

$$G_{11} = \frac{\sin M\theta}{V\sin(M + 1)\theta} \tag{45}$$

As far as the coupling potential $V$ is concerned, it is dependent from the discretization in the longitudinal direction and is given simply by

$$V = -\frac{\hbar^2}{2ma^2} \tag{46}$$

where $\hbar$ is the reduced Planck constant, $m$ is the effective mass of the electron, and $a$ is the discretization step. The matrix elements of the coupling potential are thus just the overlap integrals of the transverse modes in the two facing sections times $V$.

Such modes are computed by solving the one-dimensional (for a two-dimensional wire) or two-dimensional (for a three-dimensional wire) Schrödinger equation in the transverse direction with the confinement potential of each section.

If the effective mass varies along the transverse direction, the evaluation of the matrix elements of $V$ must be done considering $m$ not as a constant but as a function of position, and the transverse eigenfunctions must be evaluated with a somewhat more complex procedure, as detailed in Ref. [28].

Once the Green's function matrix for the whole structure has been obtained, the transmission and reflection matrices can be computed on the basis of the following relationships [24], which have been derived adapting to the discrete case the continuum relationships proposed by Stone and Szafer [29]. For the transmission matrix between locations $j$ and $l$, we get

$$t_{nm} = -i2V(\sin\theta_n \sin\theta_m)^{1/2}e^{i(\theta_m l - \theta_n j)}\langle n|G_{jl}|m\rangle \tag{47}$$

while the reflection matrix is given by

$$r_{nm} = -(\sin\theta_n/\sin\theta_m)^{1/2}e^{i2(\theta_n+\theta_m)j}(i2V\sin\theta_m\langle n|G_{ll}|m\rangle + \delta_{mn}) \tag{48}$$

In both these equations, $\theta_n$ and $\theta_m$ are the products of the longitudinal wave vector for the $n$th mode in section $j$ and for the $m$th mode in section $l$ times the longitudinal discretization step, and $\delta_{nm}$ is a Kronecker $\delta$, which differs from zero only if $n = m$. Clearly, only propagating modes must be considered in this calculation.

## 3.2. Recursive Scattering Matrix Technique

A different approach to the calculation of the transmission matrix, which allows inclusion of the effects of a magnetic field in a relatively straightforward way, is the one based on the recursive scattering matrix formalism. The scattering matrices of elementary sections are computed, and they are then combined starting from the right and adding one section at a time, moving backward. To apply this approach, the structure to be investigated must be subdivided into sections, each of which is characterized by a single discontinuity of the transverse confinement potential. To this purpose, we start from a partitioning of the structure into slices, within each of which the transverse potential can be assumed constant, as shown in Fig. 5 for a hard-wall structure (slice boundaries are represented with dashed lines). The sections for the calculation of the elementary scattering matrices straddle from the middle of one slice to the middle of the neighboring one (their boundaries are marked with dash-dot lines in Fig. 5), thereby including a single discontinuity. The scattering matrix for such sections can be computed applying the mode-matching technique (i.e. writing a general solution for the Schrödinger equation of each of the two facing sections, and then evaluating the unknown coefficients by enforcing the continuity of the wave function and of its normal derivative at the interface). In general, the scattering matrix $S$ is defined by

$$\begin{pmatrix} \mathbf{B} \\ \mathbf{C} \end{pmatrix} = \begin{pmatrix} \rho^+ & \tau^- \\ \tau^+ & \rho^- \end{pmatrix} \begin{pmatrix} \mathbf{A} \\ \mathbf{D} \end{pmatrix} = S \begin{pmatrix} \mathbf{A} \\ \mathbf{D} \end{pmatrix} \tag{49}$$

where $\mathbf{B}$ and $\mathbf{C}$ are the vectors of the amplitudes of the modes coming out of the scattering region from the left and from the right, respectively, and $\mathbf{A}$ and $\mathbf{D}$ are the vectors of the amplitudes of the modes going into the scattering region from the left and from the right, respectively, as shown in Fig. 6. In terms of the transmission and reflection matrices, we can interpret $\tau^+$ and $\rho^+$ as the "forward" transmission and reflection matrices (i.e., the matrices representing transmission and reflection of the forward-impinging modes). Correspondingly, $\tau^-$ and $\rho^-$ can be seen as the "backward" transmission and reflection matrices.



Figure 5. Subdivision of a device into slices with constant transverse potential.

Figure 6. Definition of the elements of the scattering matrix.

Let us consider the section shown in Fig. 7. On the left of the interface, the wave function can be written, assuming that a single mode is impinging from the left and that $x$ is the direction of electron propagation,

$$\psi_l(x, y) = \frac{1}{\sqrt{k_i}} \exp(ik_i x)\chi_i(y) + \sum_j \frac{1}{\sqrt{k_j}} r_{ij} \exp(-ik_j x)\chi_j(y) \tag{50}$$

where $\chi_i(y)$ are the transverse eigenfunctions in the left region. To the right of the interface the wave function will contain only right-going modes:

$$\psi_r(x, y) = \sum_m \frac{1}{\sqrt{k_m}} t_{im} \exp(-ik_m x)\phi_m(y) \tag{51}$$

where $\phi_m(y)$ are the transverse eigenfunctions in the right region.

Let us assume that the origin of the coordinates for the $x$ axis coincides with the interface (this may not be the best choice in actual scattering matrix calculations, but it is certainly convenient for the presentation of the matching procedure). With such a choice, enforcing the continuity of the wave function at the interface $[\psi_l(x, y) = \psi_r(x, y)|_{x=0}]$ we obtain

$$\frac{1}{\sqrt{k_i}}\chi_i(y) + \sum_j \frac{1}{\sqrt{k_j}} r_{ij}\chi_j(y) = \sum_m \frac{1}{\sqrt{k_m}} t_{im}\phi_m(y) \tag{52}$$



Figure 7. Elementary section for the calculation of the scattering matrix.

The matching condition for the normal derivative will read

$$\frac{k_i}{\sqrt{k_i}}\chi_i(y) - \sum_j \frac{k_j}{\sqrt{k_j}}r_{ji}\chi_j(y) = \sum_m \frac{k_m}{\sqrt{k_m}}t_{im}\phi_m(y) \tag{53}$$

In principle. summations in Eqs. (52) and (53) should run over an infinite number of elements, but in practical calculations they are limited to $N$ transverse modes on the left and to $M$ transverse modes on the right, with $N$ and $M$ chosen in such a way as to include a reasonable number of evanescent modes, in addition to all propagating modes. In general, $M$ and $N$ should be increased until no change is observed in the solution.

Both Eq. (52) and Eq. (53) depend on the transverse variable $y$ (or they would depend on the transverse variables $y$ and $z$, if we considered a three-dimensional structure) and contain a total of $N + M$ unknowns. We therefore need $N + M$ algebraic equations. which could be obtained, for example. by choosing a sufficient number of points along $y$ on which to enforce Eq. (52) and Eq. (53). There is, however, a simpler approach, which allows a simplification of the equations and automatically yields linearly independent equations: We "project" the mode-matching equations onto a basis of transverse modes; that is, we integrate over the transverse dimension the product of a set of transverse modes times the equations

$$\frac{1}{\sqrt{k_i}}\langle\eta_p(y)|\chi_i(y)\rangle + \sum_j \frac{1}{\sqrt{k_j}}r_{ij}\langle\eta_p(y)|\chi_j(y)\rangle = \sum_m \frac{1}{\sqrt{k_m}}t_{im}\langle\eta_p(y)|\phi_m(y)\rangle$$

$$\frac{k_i}{\sqrt{k_i}}\langle\eta_p(y)|\chi_i(y)\rangle - \sum_j \frac{k_j}{\sqrt{k_j}}r_{ij}\langle\eta_p(y)|\chi_j(y)\rangle = \sum_m \frac{k_m}{\sqrt{k_m}}t_{im}\langle\eta_p(y)|\phi_m(y)\rangle \tag{54}$$

In the case of "soft walls," that is, of a generic confinement potential, we can choose equivalently the transverse eigenfunctions of the left or of the right region for the projection. This is not the case for a hard-wall potential, such as that represented in Fig. 7. If we project both the wave function and the normal derivative matching equations onto the transverse modes of the narrower section, we will not enforce any condition for the wave function in the outer segments of the interface belonging to the wider region, where a zero value should be enforced. If we instead project both equations onto the wave functions of the wider region, we will enforce the wrong condition for the normal derivative in the outer segments of the wider region: the normal derivative will be set to zero, although no condition should be set on it.

It is therefore apparent that the correct approach consists in projecting the equations for the wave function matching onto the transverse modes of the wider region and the equations for the normal derivative matching onto the transverse modes of the narrower region:

$$\frac{1}{\sqrt{k_i}}\langle\phi_p(y)|\chi_i(y)\rangle + \sum_j \frac{1}{\sqrt{k_j}}r_{ij}\langle\phi_p(y)|\chi_j(y)\rangle = \sum_m \frac{1}{\sqrt{k_m}}t_{im}\langle\phi_p(y)|\phi_m(y)\rangle$$

$$\frac{k_i}{\sqrt{k_i}}\langle\chi_p(y)|\chi_i(y)\rangle - \sum_j \frac{k_j}{\sqrt{k_j}}r_{ij}\langle\chi_p(y)|\chi_j(y)\rangle = \sum_m \frac{k_m}{\sqrt{k_m}}t_{im}\langle\chi_p(y)|\phi_m(y)\rangle \tag{55}$$

which, considering the orthonormality of the elements belonging to a basis of eigenvectors, can be rewritten

$$\frac{1}{\sqrt{k_i}}\langle\phi_p(y)|\chi_i(y)\rangle + \sum_j \frac{1}{\sqrt{k_j}}r_{ij}\langle\phi_p(y)|\chi_j(y)\rangle = \sum_m \frac{1}{\sqrt{k_m}}t_{im}\delta_{pm}$$

$$\frac{k_i}{\sqrt{k_i}}\delta_{pi} - \sum_j \frac{k_j}{\sqrt{k_j}}r_{ij}\delta_{pj} = \sum_m \frac{k_m}{\sqrt{k_m}}t_{im}\langle\chi_p(y)|\phi_m(y)\rangle \tag{56}$$

where $\delta_{ij}$ is the Kronecker delta, which equals 1 if $i = j$ and 0 otherwise.

In this way, we obtain exactly $N + M$ equations. if $N$ are the modes considered in the narrower region and $M$ are those in the wider region. From these equations, we get the $r_{ii}$ and $t_{im}$ coefficients, which correspond to the elements of the $i$th row of the submatrices $\rho^-$ and $\tau^+$. To determine all rows, the calculation must be repeated $M$ times, considering each of the possible impinging modes.

We are then faced with the problem of combining the scattering matrices of two adjacent sections. With reference to Fig. 8, we can write the following two relationships from the definition of the scattering matrices for the two sections:

$$\begin{pmatrix} \mathbf{B} \\ \mathbf{C'} \end{pmatrix} = \begin{pmatrix} \rho_1^+ & \tau_1^- \\ \tau_1^+ & \rho_1^- \end{pmatrix} \begin{pmatrix} \mathbf{A} \\ \mathbf{D'} \end{pmatrix} \tag{57}$$

$$\begin{pmatrix} \mathbf{D'} \\ \mathbf{C} \end{pmatrix} = \begin{pmatrix} \rho_2^+ & \tau_2^- \\ \tau_2^+ & \rho_2^- \end{pmatrix} \begin{pmatrix} \mathbf{C'} \\ \mathbf{D} \end{pmatrix} \tag{58}$$

To obtain the scattering matrix between the input of the first section and the output of the second section, one must eliminate $C'$ and $D'$ from the above equations, as discussed, for example, in Ref. [30]. In extended form, Eqs. (57) and (58) will read

$$B = \rho_1^+ A + \tau_1^- D' \tag{59}$$

$$C' = \tau_1^+ A + \rho_1^- D' \tag{60}$$

$$D' = \rho_2^+ C' + \tau_2^- D \tag{61}$$

$$C = \tau_2^+ C' + \rho_2^- D \tag{62}$$

If we substitute Eq. (62) into Eq. (61), we obtain an expression for $C'$ as a function of $A$ and $D$:

$$C' = (I - \rho_1^- \rho_2^+)^{-1}(\tau_1^+ A + \rho_1^- \tau_2^- D) \tag{63}$$

We then substitute Eq. (62) into Eq. (60), obtaining

$$B = \rho_1^+ A + \tau_1^-(\rho_2^+ C' + \tau_2^- D) \tag{64}$$

which, on its substitution into Eq. (63), yields

$$B = [\rho_1^+ + \tau_1^- \rho_2^+ (I - \rho_1^- \rho_2^+)^{-1}\tau_1^+]A + [\tau_1^- \rho_2^+ (I - \rho_1^- \rho_2^+)^{-1}\rho_1^- \tau_2^- + \tau_1^- \tau_2^-]D \tag{65}$$

The coefficient of $D$ can be further simplified:

$$\tau_1^- \rho_2^+ (I - \rho_1^- \rho_2^+)^{-1}\rho_1^- \tau_2^- + \tau_1^- \tau_2^- = \tau_1^-[\rho_2^+ (I - \rho_1^- \rho_2^+)^{-1}\rho_1^- + I]\tau_2^-$$
$$= \tau_1^- \{[(\rho_2^+ \rho_1^-)^{-1} - I]^{-1} + I\}\tau_2^- \tag{66}$$

In the square bracket we have an expression of the type

$$(M^{-1} - I)^{-1} + I \tag{67}$$



**Figure 8.** Combination of the scattering matrices relative to the two cascaded sections.

which can be simplified as

$$
\begin{aligned}
(M^{-1} - I)^{-1} + I &= (IM^{-1} - MM^{-1})^{-1} + I \\
&= [(I - M)M^{-1}]^{-1} + I \\
&= M(I - M)^{-1} + I \\
&= M(I - M)^{-1} + (I - M)(I - M)^{-1} \\
&= (M + I - M)(I - M)^{-1}
\end{aligned}
\tag{68}
$$

Thus

$$
\tau_1^-\{[(\rho_2^+\rho_1^-)^{-1} - I]^{-1} + I\}\tau_2^- = \tau_1^-(I - \rho_2^+\rho_1^-)^{-1}\tau_2^-
\tag{69}
$$

so that

$$
B = [\rho_1^+ + \tau_1\,\rho_2^+(I - \rho_1^-\rho_2^+)^{-1}\tau_1^+]A + [\tau_1^-(I - \rho_2^+\rho_1^-)^{-1}\tau_2^-]D
\tag{70}
$$

Therefore, from Eq. (70) we get the expressions for two of the submatrices of the scattering matrix of the two joint sections, $\rho^+$ and $\tau^-$ :

$$
\begin{aligned}
\rho^+ &= \rho_1^+ + \tau_1^-\rho_2^2(I - \rho_1\rho_2^1)^{-1}\tau_1^+ \\
\tau^- &= \tau_1\,(I - \rho_2^+\rho_1)^{-1}\tau_2^-
\end{aligned}
\tag{71}
$$

The other two submatrices, $\tau^+$ and $\rho^-$ can be quickly obtained substituting Eq. (63) into Eq. (62):

$$
\begin{aligned}
C &= \tau_2^+(I - \rho_1^-\rho_2^+)^{-1}(\tau_1^+A + \rho_1^-\tau_2^-D) + \rho_2\,D \\
&= \tau_2^1(I - \rho_1^-\rho_2^+)\tau_1^+A + [\tau_2^1(I - \rho_1^-\rho_2^+)^{-1}\rho_1\,\tau_2^- + \rho_2^+]
\end{aligned}
\tag{72}
$$

so that

$$
\begin{aligned}
\tau' &= \tau_2^+(I - \rho_1^-\rho_2^+)\tau_1^+ \\
\rho &= \tau_2^1(I - \rho_1\rho_2^1)^{-1}\rho_1\,\tau_2 + \rho_2^1
\end{aligned}
\tag{73}
$$

We can easily show that the scattering matrix is unitary. Let us define the vectors $\mathbf{u}_{out}$ and $\mathbf{u}_{in}$ as follows:

$$
\begin{aligned}
\mathbf{u}_{out} &= \begin{pmatrix} \mathbf{B} \\ \mathbf{C} \end{pmatrix} \\
\mathbf{u}_{in} &= \begin{pmatrix} \mathbf{A} \\ \mathbf{D} \end{pmatrix}
\end{aligned}
\tag{74}
$$

With this definition

$$
\mathbf{u}_{out} = S\mathbf{u}_{in}
\tag{75}
$$

and, because for the conservation of the probability current $\sum_n |u_{out_n}|^2 = \sum_n |u_{in_n}|^2$, which can also be written

$$
\mathbf{u}_{out}^\dagger\mathbf{u}_{out} = \mathbf{u}_{in}^\dagger\mathbf{u}_{in}
\tag{76}
$$

we have

$$
\begin{aligned}
(S\mathbf{u}_{in})^*(S\mathbf{u}_{in}) &= \mathbf{u}_{in}^\dagger\mathbf{u}_{in} \\
\mathbf{u}_{in}^\dagger S^\dagger S\mathbf{u}_{in} &= \mathbf{u}_{in}^\dagger\mathbf{u}_{in} \\
S^\dagger S &= I
\end{aligned}
\tag{77}
$$

which completes the proof of the scattering matrix unitarity.

Although the method we have described is characterized by high numerical stability, it is necessary to represent numbers at least in double precision to keep a sufficient accuracy in calculations for extended or complex structures. One good check of the numerical accuracy is represented by the verification of the unitarity of the scattering matrix, which should be satisfied within at least $10^{-6}$.

## 3.3. Transport in the Presence of a Magnetic Field

Let us now consider the case of transport in the presence of a magnetic field orthogonal to the plane of the device. The Hamiltonian reads

$$\frac{(-i\hbar\vec{\nabla} + e\vec{A}(\vec{r}))^2}{2m^*}\psi(\vec{r}) + V(\vec{r})\psi(\vec{r}) = E\psi(\vec{r}) \tag{78}$$

where $\vec{A}$ is the vector potential, and $V(\vec{r})$ is the usual scalar potential.

Let us consider a two-dimensional device lying in the $x$-$y$ plane, with a magnetic field along the $z$ direction: there are several possible gauge choices for the representation of the vector potential, and we will focus in particular on two transverse gauges, one with a nonzero component only along the direction of electron propagation ($\vec{A} = [-By, 0, 0]^T$) and the other with a nonzero component in the transverse direction ($\vec{A} = [0, Bx, 0]^T$).

For the specific two-dimensional problem that we are considering, the Schrödinger equation reads

$$\frac{1}{2m^*}(-i\hbar\vec{\nabla} + q\vec{A})^2\psi(x, y) + V(x, y)\psi(x, y) = E\psi(x, y)$$

$$\frac{1}{2m^*}(-\hbar^2\nabla^2\psi - i\hbar e\vec{\nabla} \cdot (\vec{A}\psi) - i\hbar e\vec{A} \cdot \vec{\nabla}\psi + e^2|\vec{A}|^2\psi) + V\psi = E\psi \tag{79}$$

Because for both gauges we are considering $\vec{\nabla} \cdot \vec{A} = 0$, the previous equation becomes

$$-\frac{\hbar^2}{2m^*}\frac{\partial^2\psi}{\partial x^2} - \frac{\hbar^2}{2m^*}\frac{\partial^2\psi}{\partial y^2} - i\frac{\hbar e}{2m^*}A\frac{\partial\psi}{\partial y} + \frac{e^2 A^2}{2m^*}\psi + V\psi = E\psi$$

$$-\frac{\hbar^2}{2m^*}\frac{\partial^2}{\partial x^2}\psi + \frac{1}{2m^*}\left(-i\hbar\frac{\partial}{\partial y} + eA\right)^2\psi + V\psi = E\psi \tag{80}$$

Let us first discuss the solution of this Schrödinger equation for the choice of gauge with a nonzero component only in the transverse direction, $\vec{A} = [0, Bx, 0]^T$, which has been worked out in detail by Governale and Böse [31]. Because the term $\left(-i\hbar\frac{\partial}{\partial y} + eA\right)^2$ yields a cross term with the partial derivative of $\psi$ with respect to $y$ and $A$, we need to consider regions in which the dependence of $A$ on $x$ can be neglected if we want to separate the Schrödinger equation into a longitudinal equation (with only an $x$ dependence) and a transverse equation (with only a $y$ dependence). This result can be achieved considering slices in which not only the $x$ dependence of the transverse confinement potential is negligible but also the dependence of $A$ on $x$ can be ignored. To this purpose, we need to consider slices that are rather thin and numerous, even if there is a very slow variation along $x$ of $V$. We will then compute the scattering matrix of sections containing a discontinuity in the vector potential and, possibly, in the transverse potential. In other words, each slice will be assigned a constant value of the vector potential, based on its $x$ coordinate, and each section, straddling between the centers of nearby slices, will thus include a discontinuity of the vector potential. With the above hypotheses, the Hamiltonian can be split into a longitudinal and a transverse part (for each coordinate $x_i$):

$$H_{long} = -\frac{\hbar^2}{2m^*}\frac{\partial^2}{\partial x^2}$$

$$H_{transv} = \frac{1}{2m^*}\left(-i\hbar\frac{\partial}{\partial y} + eA(x_i)\right)^2 + V(x_i, y) \tag{81}$$

Since the eigenfunctions of the longitudinal Hamiltonian are simply plane waves, the generic eigenfunction of the complete Hamiltonian for each section will have the expression

$$\psi(x_i, y) = \sum_j \alpha_j \chi_j(x_i, y)e^{ik_{i,x_j}x} + \sum_i \beta_j \chi_j(x_i, y)e^{-ik_{i,x_j}x} \tag{82}$$

where $\chi_j(x_i, y)$ is the $j$th eigenfunction of the transverse Hamiltonian in $x_i$. It is possible to verify by substitution that the transverse eigenfunctions $\chi_j(x_i, y)$ equal the transverse eigenfunctions in the case of no magnetic field times a phase factor $\exp[(-i(e/\hbar)A(x_i)y)]$. This is quite a simplification, because no special approach is needed to compute the transverse eigenfunctions.

The requirement that must be satisfied to obtain reliable results is, as already stated, that the vector potential undergoes little variation within each slice and, more quantitatively, that the magnetic flux threaded by each slice is less than a flux quantum. The main drawback of this approach is represented by the very large number of slices that are needed to perform a calculation in the case of a large value of the magnetic field, or of an extended structure, even if the transverse potential can be assumed constant in long sections (and just a few slices would be needed from this point of view). This approach is quite convenient in the case of a magnetic field variable along the longitudinal direction or in the case of a transverse potential that undergoes rapid fluctuations along $x$, thereby forcing a discretization with a very large number of slices.

With this particular gauge, the group velocity of the modes is unchanged with respect to the case of no magnetic field and is given by $v_i = \hbar k_{x_i}/m^*$. Thus, the same normalization procedures used for the calculations without magnetic field do apply in this case.

A different approach consists in choosing the gauge with nonzero components only along the longitudinal direction ($\vec{A} = [-By, 0, 0]'$). In such a case, the Schrödinger equation is not separable in the sense that the transverse solution depends on the longitudinal wave vector (thus, the transverse eigenfunctions must be recomputed for each value of the energy of the impinging particle), although the wave function can still be written as a product of a function of $x$ and of a function of $y$. In this case, the group velocity of the modes does not have the same expression as in the absence of magnetic field and is given by [32]

$$v_n = \frac{\hbar}{m^*} \int \psi_n^*(y)\psi_n(y)\left(k_{x_n} - \frac{y}{l_B^2}\right)dy \tag{83}$$

where $l_B$ is the so-called magnetic length: $l_B = \sqrt{eB/\hbar}$. The value of $v_n$ is relevant for the correct calculation of the transmission and reflection coefficients. If, in the mode-matching equations for the determination of the scattering matrix of each section, modes without any normalization coefficient for the longitudinal component are used, such as $\exp(ik_{x_n})\phi_n(y)$, the correct reflection and transmission coefficients are obtained multiplying the results from mode-matching by the ratio $v_m/v_n$, where the $n$th mode is the impinging mode and the $m$th mode is transmitted or reflected.

In the literature, several approaches have been proposed for the calculation of the transverse wave functions with this choice of gauge. In particular, we mention the methods proposed by Palacios and Tejedor [33] and by Tamura and Ando [34], focusing on the latter.

Palacios and Tejedor propose a solution based on discretizing the Schrödinger equation once a trial solution of the form $\exp(-ik_{x_i}x)\phi_i(y)$ has been substituted into it. A system of linear equations is obtained that has a nonzero solution only if the coefficient matrix is singular. The determinant of the coefficient matrix can be expressed in the form of a polynomial in $k_x$; therefore, the allowed values of $k_x$ correspond to the roots of the polynomial. Unfortunately, finding the roots of a polynomial is a problem that is very sensitive to numerical precision, and serious numerical problems appear as soon as the number of transverse discretization points goes beyond a few tens. Some improvement can be achieved when transforming the problem of finding the roots of a polynomial of order $N$ into an eigenvalue problem of dimension $2N$, but the achievable precision is still unsatisfactory for large values of the magnetic field.

The technique proposed by Tamura and Ando is extremely stable from a numerical point of view. It consists of writing the transverse eigenfunctions as linear combinations of the transverse eigenfunctions in the absence of a magnetic field:

$$\phi_i^+(y) = \sum_{j=1}^{\Lambda} c_{ij}^{\pm} \chi_i^{0}(y)$$ (84)

equation where $\phi_i^+(\phi_i^-)$ is the transverse eigenfunction associated with the $i$th right- (left-) going mode and $\chi_j^0(y)$ is the $j$th transverse mode in the absence of a magnetic field.

Tamura and Ando have shown that the coefficients $c_{ij}$, as well as the longitudinal eigenvalues $k_{x_i}$, can be obtained as the first $N$ elements of the $N$ eigenvalue problems

$$\left[\begin{pmatrix} 0 & I \\ A & B \end{pmatrix}\right]\left[\begin{pmatrix} \vec{c}_i^{\pm} \\ \vec{d}_i^{\pm} \end{pmatrix}\right] = \frac{k_{x_i}^{\pm} W}{\pi}\left[\begin{pmatrix} \vec{c}_i^{\pm} \\ \vec{d}_i^{\pm} \end{pmatrix}\right]$$ (85)

where $\vec{c}_i^{\pm} = [c_{i1}^{\pm}, \ldots, c_{iN}^{\pm}]^T$ and $\vec{d}_i^{\pm} = (k_{x_i}^{\pm} W/\pi)\vec{c}_i^{\pm}$. Furthermore, the elements of the $N \times N$ matrices $A$ and $B$ are given by

$$A_{jr} = \left[\left(\frac{k_F W}{\pi}\right)^2 - \frac{E_j}{E_1^*}\right]\delta_{jr} - \left(\frac{\hbar\omega_c}{2E_1^*}\right)^2\left\langle\chi_j^0\left|\left(\frac{\pi y}{W}\right)^2\right|\chi_r^0\right\rangle$$

$$B_{jr} = \frac{\hbar\omega_c}{E_1^*}\left\langle\chi_j^0\left|\left(\frac{\pi y}{W}\right)\right|\chi_r^0\right\rangle$$ (86)

where $\delta_{jr}$ is the Kronecker delta, $k_F$ is the Fermi wave vector of the impinging electron, $E_j$ is the transverse eigenvalue associated with the $j$th mode in the absence of a magnetic field, $E_1^* = (\hbar^2\pi^2)/(2m^*W^2)$, and $\omega_c = eB/m^*$.

With this approach the transverse eigenfunctions can be obtained efficiently even for large values of the magnetic field and relatively wide structures, as long as a large enough number of basis elements is considered (transverse modes for the case of $B = 0$). For example, in a structure 8 $\mu$m wide in the presence of a 3-T magnetic field, about 800 basis elements are needed.

## 4. REALISTIC SIMULATION OF QUANTUM WIRES

So far, we have considered the calculation of the transmission matrix for a given confinement potential without touching on the numerical techniques that are needed for the self-consistent determination of such a confinement potential. We will now discuss a few approaches for the computation of the potential in quasi-one-dimensional nanostructures, including rather approximate and simplified techniques as well as detailed self-consistent methods.

A very simple, although effective, technique for a first-order estimate of the confinement potential in a nanostructure defined by means of depletion gates was introduced by Davies, Larkin, and Sukhorukov [35]. They derived semianalytical expressions for the calculation of the bare confinement potential (i.e., the potential resulting from the electrostatic contribution of the gates, without the term from electron–electron interactions) in a heterostructure with polygonally shaped metal gates on its surface.

In particular, they present a very straightforward method for the evaluation of the electrostatic potential in the approximation of "pinned surface"; that is, if we assume that the Fermi-level at the surface is constant and independent of the bias voltages applied to the gates. This is a commonly accepted approximation that can be found quite often in the literature, although it is not exact (it is apparent that, to keep the potential at the surface constant and equal to zero as the gate voltages are changed, charge should be transferred from the gates to the surface to compensate for their electrostatic action). Under the hypothesis of pinned surface, it is possible to show that the potential in the heterostructure can simply be expressed as a linear superposition of the contributions of each gate, which greatly simplifies

the problem. To find the expression of the potential resulting from each single gate, we add the further condition that the electric field vanishes deep inside the heterostructure. This implies that, if we define as $\tilde{\phi}(\vec{q}, 0)$ the Fourier transform of the potential $\phi(\vec{r}, 0)$ at the surface (where $\vec{q}$ is the wave vector), the dependence on $z$ will introduce an exponential behavior:

$$\tilde{\phi}(\vec{q}, z) = \tilde{\phi}(\vec{q}, 0)\exp(-|qz|) \tag{87}$$

Because a multiplication in the transformed domain is equivalent to a convolution in the real-space domain, the spatial dependence of the potential is given by the convolution of the potential at the surface times the inverse Fourier transform of $\exp(-|qz|)$:

$$\phi(\vec{r}, z) = \int \frac{|z|}{2\pi(z^2 + |\vec{r} - \vec{r}\,'|^2)^{3/2}}\phi(\vec{r}\,', 0)d\vec{r}\,' \tag{88}$$

Davies et al. provide expressions of the contribution to the electrostatic potential at the level of the 2DEG (2-dimensional electron gas) for several basic geometric shapes (such as triangles, finite and infinite rectangles, etc.) and for generic polygons. Let us discuss specifically this last case, which is of more common interest.

Let us assume we have $N$ polygonal gates, each covering a surface $S_i$. Our boundary condition will be expressed by $\phi(\vec{r}, 0) = V_i$ if $\vec{r}$ belongs to one of the regions $S_i$; otherwise, $\phi(\vec{r}, 0) = 0$, as expected from the hypothesis of surface pinning. Following Ref. [35], we now consider the Green's function of the Laplace equation in the half space for $z > 0$, with the boundary condition that it will vanish for $z = 0$ and $r \to \infty$. Using the method of images [36], such a Green's function, satisfying the definition $\nabla^2 G(\vec{R}, \vec{R}\,') = -\delta(\vec{R} - \vec{R}\,')$, reads

$$G(\vec{R}, \vec{R}\,') = \frac{1}{4\pi\sqrt{|\vec{r} - \vec{r}\,'|^2 + (z - z')^2}} - \frac{1}{4\pi\sqrt{|\vec{r} - \vec{r}\,'|^2 + (z + z')^2}} \tag{89}$$

where $\vec{R} = (\vec{r}, z)$. We know the potential $\phi(\vec{r}, 0)$ at the surface of the heterostructure, the only region in our simplified model in which the source term for the Laplace equation (i.e., the charge $\rho$) does not vanish. Thus, we can extend $\phi$ to the whole semispace with $z > 0$ using Green's theorem [37, 38]

$$\phi(\vec{R}) = \iint_D d\vec{s}\,'\left[\phi(\vec{R}\,')\nabla'G(\vec{R}, \vec{R}\,') - G(\vec{R}, \vec{R}\,')\nabla'\phi(\vec{R}\,')\right] \tag{90}$$

where the surface $D$ is a hemisphere of very large radius closed at the bottom by a circle belonging to the plane $z = 0$, with the vector $d\vec{s}\,'$ aimed inside the surface. Because the electric field vanishes at infinity, the contribution from the integral over the hemisphere vanishes too, if we let its radius go to infinity and we are left simply with an integral over the plane $z = 0$:

$$\phi(\vec{R}) = \iint d\vec{r}\,'\phi(\vec{r}\,', 0)\left[\frac{\partial}{\partial z'}G(\vec{R}, \vec{R}\,')\right]_{z'=0} = \frac{1}{2\pi}\sum_i V_i I_i(\vec{R}) \tag{91}$$

where the integrals $I_i(\vec{R})$ are given by

$$I_i(\vec{R}) = \iint_{S_i} \frac{z d\vec{r}\,'}{(|\vec{r} - \vec{r}\,'|^2 + z^2)^{3/2}} \tag{92}$$

thereby corresponding to an integration over the surface of each gate [since everywhere else $\phi(\vec{r}\,', 0) = 0$]. The next step consists of using the divergence theorem to recast such surface integrals into contour integrals along the boundary of each gate $\partial S_i$. To this purpose, let us consider the divergence in two dimensions of the function $\vec{r}/r^2$:

$$\text{div}\frac{\vec{r}}{r^2} = 2\pi\delta(\vec{r}) \tag{93}$$

This implies that if we consider a generic function $M(r)$,

$$\text{div}\frac{\vec{r}M(r)}{r^2} = \frac{\frac{dM(r)}{dr}}{r} + 2\pi M(0)\delta(\vec{r}) \tag{94}$$

which, with the application of the divergence theorem, leads to

$$\iint_{S_i} \frac{\frac{dM(r)}{dr}}{r}d\vec{r} = \iint_{S_i} \text{div}\frac{\vec{r}M(r)}{r^2} - 2\pi M(0)\delta(\vec{r})$$

$$= -\int_{\partial S_i} \frac{M(r)\vec{r}\cdot d\vec{l}}{r^2} - 2\pi M(0)\theta(\vec{r}) \tag{95}$$

where $\theta(\vec{r}) = 1$ if the point $\vec{r} = 0$ belongs to $S_i$ and $\theta(\vec{r}) = 0$ otherwise, and the vector $d\vec{l}$ is perpendicular to the boundary $\partial S_i$ and directed inward.

If we now assume $M(r) = -z(r^2 + z^2)^{-1/2}$, we get

$$\frac{\frac{dM(r)}{dr}}{r} = z(r^2 + z^2)^{-3/2} \tag{96}$$

which, combined with Eqs. (92) and (96), yields

$$I_i(\vec{R}) = \int_{\partial S_i} \frac{z(\vec{r}' - \vec{r})\cdot d\vec{l}'}{|\vec{r} - \vec{r}'|^2\sqrt{|\vec{r}' - \vec{r}|^2 + z^2}} + 2\pi\gamma_i(\vec{r}) \tag{97}$$

where $\gamma_i(\vec{r}) = 1$ if $\vec{r} \in S_i$ and $\gamma_i(\vec{r}) = 0$ otherwise. Considering that the regions $S_i$ are polygons, a contour integral around their boundaries is the result of the sum of line integrals along the straight segments forming each edge. Thus Eq. (97) can be rewritten as

$$I_i(\vec{R}) = \sum_k \int_{E_{ik}} \frac{zD_{ik}dl}{(D_{ik}^2 + l^2)\sqrt{D_{ik}^2 + z^2 + l^2}} + 2\pi\gamma_i(\vec{r}) \tag{98}$$

where $E_{ik}$ is the $k$th edge of the $i$th gate, and $D_{ik}$ is the distance between such an edge and the point $\vec{r}$, as shown in Fig. 9. This integral can be performed analytically [35, 39] and we obtain

$$I_i(\vec{R}) = \sum_k \arctan\frac{zb_{ik}}{D_{ik}\sqrt{z^2 + D_{ik}^2 + b_{ik}^2}} + \arctan\frac{zc_{ik}}{D_{ik}\sqrt{z^2 + D_{ik}^2 + c_{ik}^2}} \tag{99}$$



Figure 9. Definition of the quantities used for the calculation of the field resulting from a polygonal electrode.

where $b_{ik}$ and $c_{ik}$ are the distances between the vertices delimiting the $k$-th edge and the point $A$ of the edge closest to $\vec{r}$, both measured outward from $A$. It is to be noted that $a_{ij}$, $b_{ik}$, and $D_{ik}$ can be negative, and in such cases the arctangents are to be intended according to the usual definition (i.e., with an angle limited to the domain $[-\pi/2, \pi/2]$) and not to that of the four-quadrant arctangent (atan2 in Fortran).

Equations (91) and (99) can be easily implemented into a computer code (an example is available on the Phantoms Hub [40]) to compute the bare confinement potential from an arbitrary configuration of polygonal gates.

Let us now briefly discuss how the contributions from the ionized donors and from the electrons can be introduced into this semianalytical approach. We shall follow the method proposed by Martorell et al. [41, 42], specifically for the case of an electrostatically defined quantum wire.

We consider a gate layout such as the one represented in Fig. 10: a split gate that defines an indefinitely long quantum wire. The $z$ axis is perpendicular to the surface of the heterostructure, as in the other examples, whereas $y$ is along the wire and $x$ is in the transverse direction. In such a case, with gates characterized by an infinite extension, the bare confinement potential can be computed with a semianalytical expression analogous to the ones we have previously discussed, given by Davies et al.:

$$\phi_g(x, z) = V_g\left\{1 - \frac{1}{\pi}\left[\arctan\left(\frac{w/2 - x}{z}\right) + \arctan\left(\frac{w/2 + x}{z}\right)\right]\right\} \tag{100}$$

where $V_g$ is the voltage applied to both sides of the split gate, and $w$ is the width of the gap. We notice that, as expected from the translational invariance, there is no dependence of $\phi$ on the coordinate $y$ along the wire. The other two contributions to the potential actually seen by the electrons in the wire derive from the other electrons and from the ionized donors. Another result of the translational invariance will be that the Poisson and Schrödinger equations (to be solved self consistently with the inclusion of the global electrostatic potential) will be separable. In particular, the Schrödinger equation can be separated into a longitudinal one-dimensional equation and a transverse two-dimensional equation. Because of translational invariance along $x$, the solution of the longitudinal equation is simply a plane wave. The solution of the transverse two-dimensional equation must instead be obtained numerically.

As far as the contribution of the donors is concerned, a very simple calculation is possible if we assume complete ionization (the case of incomplete ionization is treated, too, in Ref. [42]). Let us consider a donor layer with a uniform charge density $\rho_d$ located between $z = c$ and $z = c + d$. In such a case the electric field in each region can be evaluated by means of Gauss's theorem and the method of images. In the region for $0 < z < c$, there is no charge; therefore, the electric field is the result of the action of the charge contained in the slab between $c$ and $c + d$ and its image located between $-c$ and $-c - d$. Each slab



Figure 10. Indefinitely long wire defined in a heterostructure by means of a split gate.

(the actual one and its image) contributes with a term $-\rho_d d/(2\varepsilon)$. Thus, for $0 < z < c$ the donor contribution $\phi_d$ to the potential can be obtained simply from the integration of the electric field:

$$\phi_d(x, z) = -\frac{\rho d}{\varepsilon} z \qquad \text{for } 0 < z < c \qquad (101)$$

For $z$ intermediate between $c$ and $c + d$, we have the action of the image, which gives a contribution $E = -\rho_d d/(2\varepsilon)$ to the potential, of the charge in the region between $c$ and $z$ $[E = (z - c)\rho_d/(2\varepsilon)]$, and of the charge in the region between $z$ and $c + d$ $[E = -(c + d - z)\rho_d/(2\varepsilon)]$, where the electric field expressions have been obtained by applying Gauss's theorem to the parallelepipeds shown in Fig. 11. Parallelepiped $A$ has a side parallel to semiconductor surface at $z$ and one at $c$, and parallelepiped $B$ has a side at $z$ and the other at $c + d$. In the limit of a very large extension along $x$ and $y$, these are the only sides with a nonnegligible flux of the electric field, as the other sides are so far away that the electric field can be assumed to be null there. Thus, the flux per unit of surface through each of the two mentioned surfaces is one half of the charge per unit of surface contained in the parallelepiped ($\rho_d$ times the parallelepiped thickness) divided by twice the permittivity. Therefore, the total electric field in the region being considered reads

$$E = -\frac{\rho_d d}{2\varepsilon} + \frac{(z - c)\rho_d}{2\varepsilon} - \frac{c + d - z}{2\varepsilon}\rho_d = -\frac{\rho_d(d + c)}{\varepsilon} + \frac{z\rho_d}{\varepsilon} \qquad (102)$$

Finally the electric field for $z > c + d$ is zero, as the contribution from the donor charge and the images cancel each other.

We can now compute the potential by integrating the electric field. For $0 < z < c$ we have

$$\phi_d(z) = -\int_0^z -\frac{e\rho_d d}{\varepsilon} dz' = \frac{e\rho_d d}{\varepsilon} z \qquad (103)$$

for $c < z < c + d$ we get

$$\phi_d(z) = \phi_d(c) - \int_c^{c+d} -\frac{\rho_d(d + c)}{\varepsilon} + \frac{z'\rho_d}{\varepsilon} dz' = \frac{e\rho_d}{\varepsilon}\left[dz - \frac{1}{2}(z - c)^2\right] \qquad (104)$$

and for $z > c + d$ the potential is constant and equal to the value of the previous expression computed for $z = c + d$:

$$\phi_d(z) = \frac{e}{2\varepsilon}\rho_d d(2c + d) \qquad (105)$$

The calculation of the contribution from the donors is slightly more complex if partial ionization of the donors is included; such a case is treated in detail in Ref. [42].

Another commonly used boundary condition at the exposed surface of the semiconductor is the so-called "frozen-surface" hypothesis, in which, contrary to the "Fermi-level pinning" hypothesis, the surface charge is supposed to be frozen; that is, not to vary as a consequence of variations in the gate voltages and in the electric field at the surface, which is assumed to be zero (or in some cases constant), thereby enforcing a Neumann boundary condition in the free surface region and a Dirichlet boundary condition under the gates for the solution of the Poisson equation. It is to be noted that with a frozen surface boundary condition,



Figure 11. Surfaces for the application of Gauss's theorem to the calculation of the potential in the donor region.

the potential cannot be computed as a superposition of the contributions from the different gates, as in the case of Fermi-level pinning. Therefore, the treatment is significantly more complex, as outlined in Ref. [35].

Neither Fermi-level pinning nor the frozen surface hypothesis yield results that are in good quantitative agreement with the experimental data, because the real behavior of the semiconductor surface is intermediate between the two assumptions we have discussed so far.

A better, although not excellent in all cases, agreement with experimental data can be obtained with the surface treatment proposed by Iannaccone et al. [43]. This is a two-step process based on the knowledge of the experimental sheet charge density at equilibrium and on the assumption that the surface charge does not change when voltages are applied to the gates. The first step consists of assuming that all gates are grounded and that there is Fermi-level pinning at the exposed surface at a level within the bandgap that is computed by fitting the experimental sheet charge density. By solving the Poisson equation, the electric field at the surface is computed, and thus the charge density. We then assume that such a charge density remains the same when gate voltages are applied, and a Neumann boundary condition with the same value of the electric field is used. With this approach, a reasonable agreement for the pinch-off voltages of several quantum point contacts of different width can be achieved, but problems appear for split gates with a gap wider than 200 nm.

A more complete treatment can be obtained with the method proposed by Chen and Porod [44]. They assume the potential to be continuous across the surface, with a variation of the normal electric field depending on the surface charge density:

$$\varepsilon_s E_s = \varepsilon_{air} E_{air} + Q_{surf} \tag{106}$$

where $\varepsilon_s$ and $E_s$ are the permittivity and the electric field in the semiconductor, $\varepsilon_{air}$ and $E_{air}$ are the same quantities in the air, and $Q_{surf}$ is the surface charge density. The surface charge is computed in the equilibrium condition, as in the previously mentioned approach, whereas the solution of the Poisson equation is performed not only in the semiconductor but also in the air above it, with zero field boundary condition at infinity.

A further technique for the treatment of surface boundary conditions is presented by Fiori et al. in Ref. [45], considering a constant density of surface states per unit energy per unit area $D_s$ with an effective work function $\Phi^*$. Indicating the vacuum level with $E_0$, the surface states with energy below $E_0 - q\Phi^*$ act as acceptors, whereas those above act as donors. The occupancy of such states is determined self-consistently, based on the solution of the Poisson equation. This approach has provided quite a good agreement with experimental data on the pinch-off voltages for quantum point contacts of several dimensions.

## 5. SIMULATION OF QUANTUM DOTS

Many computational techniques for the investigation of quantum dots have been developed during the last two decades, starting from the early work by Kumar, Laux and Stern [46], the models for circular dots by Maksym et al. [47], and then moving to the sophisticated approaches by Stopa [48], Scholze [49], Mataigne [50], and Hawrylak [51]. Different authors have focused on limiting the approximations in specific portions of the calculation. Some consider a very simplified confinement potential and apply refined methods to the treatment of the electron–electron interaction within the dot, whereas others spend a larger effort on the computation of a realistic confinement potential consistent with the layout of the electrodes and use less advanced methods for the solution of the many-body problem. In general, there is no approach that is preferable over the others. The choice between the many techniques that have been developed for quantum dot simulation depends on the specific characteristics of the problem, on the quantities that one wants to compute, and on the level of detail with which the parameters of the device are known.

The ideal solution would involve usage of the least approximate method (i.e., the exact diagonalization of the Hamiltonian), which, however, is possible only for a very small number of electrons (below 10) with the currently available computational facilities. The second least

approximate approach is configuration–interaction, which will be discussed in detail in the following section, dedicated to the numerical simulation of multiple quantum dot systems, coupled with the solution of the Poisson equation. With configuration–interaction, the many-electron wave function is properly symmetrized, being represented as a linear combination of Slater determinants (whereas, for example, Hartree–Fock uses a single, optimized Slater determinant), but if the electrostatic interaction is more complex (as in almost any situation of interest), knowledge of the complete Green's function of the Poisson equation is needed.

In the other methods, self-consistency between the solution of the Schrödinger and of the Poisson equations is achieved with an iterative procedure, solving repeatedly for the wave functions and for the electrostatic potential. In these approaches, the electrostatic interaction between electrons is usually introduced with some type of mean-field approximation (Hartree, density-functional theory). Convergence is achieved swiftly if the electrostatic interaction leads only to a small perturbation of the quantum confinement energy, although it may become hard or even impossible to attain when it accounts for a significant part of the particle energies.

Let us first discuss approaches in which a so-called "bare confinement potential" is postulated (sometimes it is parabolic, such as in Ref. [47], and other times it is the result of a distribution of positive charge on a finite surface [52]). This corresponds to the potential that would be felt by a single electron occupying the quantum dot. As we have more than one electron in the dot, each electron will be subject also to the interaction with the other electrons, which will be represented as a "mean" field contribution, averaged over their probability distribution. The original mean-field theory was the result of Hartree and involved a different Schrödinger equation for each orbital, as each electron should see the electrostatic contribution of all other electrons except for its own. The solution of such a system of equations is somewhat time consuming; therefore, many authors resort to the so called Common Hartree Hamiltonian (CHH), a Hamiltonian that is common to all orbitals, because the self-interaction is not removed. Oaknin et al. propose [53] a correction to the usage of the CHH by subtracting the self-interaction contributions from the eigenvalues obtained at the end of the self-consistent procedure.

Other commonly used approaches are based on the local density functional approximation (LDA); that is, on the assumption that the electron density in the quantum dot varies slowly enough to allow the application of density-functional theory [54, 55]. Within these approaches, the Hamiltonian is the same for all orbitals and, in addition to the Coulomb interaction term, contains exchange and correlation terms:

$$-\frac{\hbar^2}{2m^*}\nabla^2\psi_i(\mathbf{r}_i) + \left[V_c(\mathbf{r}_i) + V_{ex}(\mathbf{r}_i) + V_{corr}(\mathbf{r}_i)\right]\psi_i(\mathbf{r}_i) = \epsilon_i\psi_i(\mathbf{r}_i) \tag{107}$$

where $\psi_i(\mathbf{r})$ is the wave function for the $i$th orbital, $\epsilon_i$ is the eigenvalue associated with the $i$th orbital, and $V_c(\mathbf{r}_i)$ is the Coulomb interaction term, which in the simple case of an isolated quantum dot will read

$$V_c(\mathbf{r}_i) = \sum_{j=1}^{N}\frac{1}{4\pi\varepsilon_0\varepsilon_r}\int_V \frac{e^2|\psi_j(\mathbf{r}_j)|^2}{|\mathbf{r}_i - \mathbf{r}_j|}\,d\mathbf{r}_j \tag{108}$$

where $N$ is the total number of electrons in the dot; $\varepsilon_0$ and $\varepsilon_r$ are the absolute and relative, respectively, dielectric permittivities; and $V$ is the volume of the quantum dot.

If the dot is not isolated, the term $V_c(\mathbf{r}_i)$ must be obtained, solving the Poisson equation in the structure. The only exception is represented by the particular case of a quantum dot defined in a medium of constant permittivity and located close to a conducting plane, in which case the term $V_c(\mathbf{r}_i)$ can be obtained summing the contributions from the images to the expression given above [52].

As far as the exchange and correlation terms are concerned, their expressions depend on the dimensionality of the problem. In three dimensions they are given, in atomic units

$(\hbar = 1, e = 1, m^* = 1)$ and in the spin unpolarized case, by the Kohn-Sham theory [54] and Ceperley's parametrization [55] as

$$V_{\rm ex} = -\frac{1}{\pi}[3\pi^2\rho(\mathbf{r})]^{1/3} \tag{109}$$

and

$$V_{\rm corr} = \frac{\gamma}{(1 + \beta_1\sqrt{r_s} + \beta_2 r_s)^2}(1 + \gamma/6\beta_1\sqrt{r_s} + 4/3\beta_2 r_s) \qquad \text{for } r_s \geq 1$$

$$V_{\rm corr} = A\ln r_s + (B - A/3) + 2/3Cr_s\ln r_s + 1/3(2D - C)r_s \qquad \text{for } r_s < 1 \tag{110}$$

where $r_s = a/a_0$, with $a = 1/\sqrt{\pi\rho}$ ($\rho$ being the local electron density) and $a_0$ the Bohr radius in the considered material; $\gamma = -0.1423$; $\beta_1 = 1.0529$; $\beta_2 = 0.3334$; $A = 0.0311$; $B = -0.048$; $C = 0.0020$; and $D = -0.0116$.

In two dimensions (very often quantum dots are studied as two-dimensional structures, because their vertical dimension is significantly less than those in the horizontal plane), following Ref. [56] and taking the functional derivatives [52, 54], they read

$$V_{\rm ex} = -\frac{4\sqrt{2}}{\pi}\frac{1}{r_s} \tag{111}$$

and

$$V_{\rm corr} = -C_0\frac{1 + d_1 w + d_2 w^2 + d_3 w^3 + d_4 w^4}{(1 + C_1 w + C_2 w^2 + C_3 w^3)^2} \tag{112}$$

where $C_0 = -0.3578$, $C_1 = 1.13$, $C_2 = 0.9052$, $C_3 = 0.4165$, $w = \sqrt{r_s}$, $d_1 = 2.26$, $d_2 = 2.635$, $d_3 = 2.007$, and $d_4 = 0.70597$.

Equation (107) can be discretized and solved numerically in a number of ways. We will discuss specifically the approach based on finite differences in one and two dimensions (the three-dimensional case is just a trivial extension of the two-dimensional case). With finite differences in two dimensions and assuming a constant discretization step $\Delta x = x_i - x_{i-1}$, we get

$$-\frac{\hbar^2}{2m^*}\frac{\psi(x_{i-1}) - 2\psi(x_i) + \psi(x_{i+1})}{\Delta x^2} + [V_c(x_i) + V_{\rm ex}(x_i) + V_{\rm corr}(x_i)]\psi(x_i) = \epsilon\psi(x_i) \tag{113}$$

which corresponds to as many algebraic equations as the number of discretization points considered in the solution domain. Let us consider a domain extending from $x = 0$ to $x = \bar{x}$, as shown in Fig. 12. Let us also assume that it is divided into $N$ mesh intervals, each of which will have a length $\Delta x = \bar{x}/(N - 1)$. If Dirichlet boundary conditions (i.e., vanishing wave function), are enforced at $x = 0$ and $x = \bar{x}$, the first equation, for $i = 1$ will read

$$-\frac{\hbar^2}{2m^*}\frac{-2\psi(x_1) + \psi(x_2)}{\Delta x^2} + [V_c(x_1) + V_{\rm ex}(x_1) + V_{\rm corr}(x_1)]\psi(x_1) = \epsilon\psi(x_1) \tag{114}$$

where the term $\psi(x_0)$ has disappeared, as it vanishes as a result of the boundary condition at $x = 0$. The second equation will read

$$-\frac{\hbar^2}{2m^*}\frac{\psi(x_1) - 2\psi(x_2) + \psi(x_3)}{\Delta x^2} + [V_c(x_2) + V_{\rm ex}(x_2) + V_{\rm corr}(x_2)]\psi(x_2) = \epsilon\psi(x_2) \tag{115}$$



Figure 12. Discretization in one dimension for the solution of the Schrödinger equation

and the same format will be valid for all other equations, except for the last, which will read

$$-\frac{\hbar^2}{2m^*}\frac{\psi(x_{N-2})-2\psi(x_{N-1})}{\Delta x^2}+[V_c(x_{N-1})+V_{ex}(x_{N-1})+V_{corr}(x_{N-1})]\psi(x_{N-1})=\epsilon\psi(x_{N-1})$$

(116)

where the term $\psi(xN)$ does not appear, because of the Dirichlet boundary condition at $x = N$.

This system of linear equations can be written in matrix form:

$$\begin{pmatrix} 2\eta+V_1 & -\eta & 0 & 0 & \cdots & 0 & 0 & 0 \\ -\eta & 2\eta+V_2 & -\eta & 0 & \cdots & 0 & 0 & 0 \\ 0 & -\eta & 2\eta+V_3 & -\eta & \cdots & 0 & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & -\eta & 2\eta+V_{N-3} & -\eta & 0 \\ 0 & 0 & 0 & \cdots & 0 & -\eta & 2\eta+V_{N-2} & -\eta \\ 0 & 0 & 0 & \cdots & 0 & 0 & -\eta & 2\eta+V_{N-1} \end{pmatrix}\begin{pmatrix} \psi_1 \\ \psi_2 \\ \psi_3 \\ \cdots \\ \cdots \\ \psi_{N-3} \\ \psi_{N-2} \\ \psi_{N-1} \end{pmatrix}=\epsilon\begin{pmatrix} \psi_1 \\ \psi_2 \\ \psi_3 \\ \cdots \\ \cdots \\ \psi_{N-3} \\ \psi_{N-2} \\ \psi_{N-1} \end{pmatrix}$$

(117)

where $\eta = \hbar^2/(2m^*\Delta x^2)$ and $V_i = V_c(x_i) + V_{ex}(x_i) + V_{corr}(x_i)$. Equation (117) represents an algebraic eigenvalue problem in which the matrix for which the eigenvalues are to be computed is tridiagonal; that is, it has only the three middle diagonals that are nonzero. This simplifies the task of computing the eigenvalues remarkably, because one of the many available routines for the efficient solution of tridiagonal eigenvalue problems can be used, such as TQLI in the EISPACK package. It is also possible to enforce different boundary conditions at $x = 0$ or $x = \bar{x}$; in particular, we can enforce a Neumann boundary condition (i.e., a condition involving the derivative being zero at the boundary). If the derivative is zero at the boundary, it means that, as shown in Fig. 13a, the wave function must have symmetric values around it. In particular, if we are considering the left boundary, we should enforce that the wave function in $x_{-1}$ has the same value as in $x_1$ (i.e., $\psi_{-1} = \psi_1$). Because we are using a three-point formula for the derivatives, it is not easy to enforce such a condition directly. We can rather enforce the condition $\psi_0 = \psi_1$, which is equivalent to placing the Neumann boundary condition in the middle between $x_0$ and $x_1$. In such a case, the first equation will read

$$-\frac{\hbar^2}{2m^*}\frac{\psi(x_0) - 2\psi(x_1) + \psi(x_2)}{\Delta x^2} + [V_c(x_1) + V_{ex}(x_1) + V_{corr}(x_1)]\psi(x_1) = \epsilon\psi(x_1)$$

(118)



Figure 13. Neumann boundary condition in the origin with uniform discretization mesh (a) and with nonuniform discretization mesh (b).

which, being $\psi(x_0) = \psi(x_1)$, becomes

$$-\frac{\hbar^2}{2m^*}\frac{-\psi(x_1) + \psi(x_2)}{\Delta x^2} + \left[V_c(x_1) + V_{ex}(x_1) + V_{corr}(x_1)\right]\psi(x_1) = \epsilon\psi(x_1) \qquad (119)$$

The other equations remain unchanged, as they are not affected by the boundary condition. It is also possible to enforce the Neumann boundary condition exactly in $x_0$, using a nonuniform mesh, and in particular a mesh that is twice the size between $x_0$ and $x_1$, as shown in Fig. 13b. This issue will be discussed in greater detail in the following, when treating the solution of the Schrödinger equation in polar coordinates. The discretized equations undergo some changes as a consequence of the nonuniform mesh. The second derivative in $x_1$ will be discretized as

$$\left.\frac{\partial^2\psi}{\partial x^2}\right|_{x=x_1} \simeq \frac{\frac{\psi(x_2)-\psi(x_1)}{\Delta x_1} - \frac{\psi(x_1)-\psi(x_0)}{\Delta x_0}}{\frac{\Delta x_0}{2} + \frac{\Delta x_1}{2}}$$

$$= \frac{\psi(x_2) - \psi(x_1)}{\Delta x_1\left(\frac{\Delta x_0}{2} + \frac{\Delta x_1}{2}\right)} - \frac{\psi(x_1) - \psi(x_0)}{\Delta x_0\left(\frac{\Delta x_0}{2} + \frac{\Delta x_1}{2}\right)}$$

$$= \frac{-2\psi(x_1)}{\Delta x_1\Delta x_0} + \frac{\psi(x_2)}{\Delta x_1\left(\frac{\Delta x_0 + \Delta x_1}{2}\right)} + \frac{\psi(x_0)}{\Delta x_0\left(\frac{\Delta x_0 + \Delta x_1}{2}\right)} \qquad (120)$$

where $\Delta x_0 = x_1 - x_0$ and $\Delta x_1 = x_2 - x_1$. This procedure can be straightforwardly extended to include a generic nonuniform mesh.

Let us now consider the two-dimensional case. The $\nabla^2$ operator will be discretized as

$$\nabla^2\psi(x, y)\Big|_{x=x_i, y=y_j} \simeq \frac{\psi_{i-1,j} - 2\psi_{i,j} + \psi_{i+1,j}}{\Delta x^2} + \frac{\psi_{i,j-1} - 2\psi_{i,j} + \psi_{i,j+1}}{\Delta y^2} \qquad (121)$$

which, in the simplified case of $\Delta x = \Delta y$ yields the discretized Schrödinger equation

$$-\frac{\hbar^2}{2m^*}\frac{\psi_{i-1,j} + \psi_{i+1,j} - 4\psi_{i,j} + \psi_{i,j-1} + \psi_{i,j+1}}{\Delta x^2} + V_{i,j}\psi_{i,j} = \epsilon\psi_{i,j} \qquad (122)$$

If we consider a rectangular $N \times M$ domain and order the mesh points in the so-called lexicographic order (i.e., moving from left to right along each row, starting with the top row and ending with the bottom row), we obtain, writing the discretized equations corresponding to each mesh point, an eigenvalue problem with a five-diagonal matrix in which the three middle diagonals are nonzero, as well as the two that are $N$ elements away from the middle on both sides (which correspond to the $\psi_{i,j-1}$ and $\psi_{i,j+1}$ terms). Considering that the dimension of the resulting eigenvalue problem is $N \times M$ by $N \times M$, its numerical solution may result quite challenging, in particular if the determination of all eigenvalues is attempted. Because we are in general interested only in the lowest eigenvalues and eigenfunctions, it is convenient to use iterative solvers that compute only a given number of the lowest eigenvalues instead of attempting the full solution. There are a number of such solvers available, based on the Lanczos method, on the Ritz procedure, or on conjugate gradient approaches.

## 5.1. The Solution of the Schrödinger Equation in Polar Coordinates

It is possible to simplify the solution of the Schrödinger equation in the presence of symmetries; in particular, circular symmetry allows a reduction of the dimensionality of the problems via separation of the variables.

As an example, let us examine the simple case of a two-dimensional circular region delimited by hard walls. The first step is to write the Laplacian in polar coordinates:

$$\nabla^2_{\rho,\phi} = \frac{1}{\rho}\frac{\partial}{\partial\rho}\left[\rho\frac{\partial\psi(\rho,\phi)}{\partial\rho}\right] + \frac{1}{\rho^2}\frac{\partial^2\psi(\rho,\phi)}{\partial\phi^2} \qquad (123)$$

The polar representation of the Schrödinger equation is thus given by

$$-\frac{\hbar^2}{2m}\left\{\frac{1}{\rho}\frac{\partial}{\partial\rho}\left[\rho\frac{\partial\psi(\rho,\phi)}{\partial\rho}\right]+\frac{1}{\rho^2}\frac{\partial^2\psi(\rho,\phi)}{\partial\phi^2}\right\}+V(\rho)\psi(\rho,\phi)=E\psi(\rho,\phi) \qquad (124)$$

where we have included a potential $V(\rho)$ that only depends on the radial coordinate (there is actually no advantage in using the polar coordinate if the potential also depends on the angle $\phi$). To solve the previous equation, we can insert a trial solution written as $\psi(\rho,\phi)=P(\rho)X(\phi)$ (i.e., a solution factorized into the product of a function of $\rho$ and a function of $\phi$). In this way we obtain:

$$-\frac{\hbar^2}{2m}\left\{\frac{1}{\rho}\frac{\partial}{\partial\rho}\left[\rho\frac{\partial}{\partial\rho}P(\rho)X(\phi)\right]+\frac{1}{\rho^2}\frac{\partial^2 P(\rho)X(\phi)}{\partial\phi^2}\right\}+V(\rho)P(\rho)X(\phi)=EP(\rho)X(\phi)$$

$$-\frac{\hbar^2}{2m}\left\{\frac{1}{\rho}\frac{\partial}{\partial\rho}\left[\rho\frac{\partial P(\rho)}{\partial\rho}\right]X(\phi)+\frac{1}{\rho^2}P(\rho)\frac{\partial^2 X(\phi)}{\partial\phi^2}\right\}+V(\rho)P(\rho)X(\phi)=EP(\rho)X(\phi) \qquad (125)$$

$$-\frac{\hbar^2}{2m}\left\{\rho\frac{\partial P}{\partial\rho}X+\rho^2 X\frac{\partial^2 P}{\partial\rho^2}+P\frac{\partial^2 X}{\partial\phi^2}\right\}+V\rho^2 PX=E\rho^2 PX$$

By dividing all terms by $-\hbar^2/(2m)PX$ we obtain

$$\rho\frac{1}{P}\frac{\partial P}{\partial\rho}+\rho^2\frac{1}{P}\frac{\partial^2 P}{\partial\rho^2}+\frac{1}{X}\frac{\partial^2 X}{\partial\phi^2}-\frac{2m}{\hbar^2}\rho^2 V=-\frac{2m}{\hbar^2}E\rho^2 \qquad (126)$$

The previous equation can be divided in two equations, each one depending only on one variable. Before doing that, let us define $K\equiv 2m/(\hbar^2)E$ and add and subtract the constant $\nu^2$ in Eq. (126). We obtain an equation for the angular dependence

$$\frac{1}{X}\frac{\partial^2 X}{\partial\phi^2}+\nu^2=0 \qquad (127)$$

and an equation for the radial dependence

$$\frac{\rho^2}{P}\frac{\partial^2 P}{\partial\rho^2}+\frac{\rho}{P}\frac{\partial P}{\partial\rho}-\frac{2m}{\hbar^2}\rho^2 V-\nu^2=-\rho^2 K^2$$

$$\frac{\partial^2 P}{\partial\rho^2}+\frac{1}{\rho}\frac{\partial P}{\partial\rho}-\frac{2m}{\hbar^2}VP+\left(K^2-\frac{\nu^2}{\rho^2}\right)P=0 \qquad (128)$$

We start the analysis of the solutions from the angular equation [Eq. (127)]. The generic solution can be written as $X(\phi)=Ae^{i\nu\phi}+Be^{-i\nu\phi}$, but because of the space periodicity, $X(\phi)=X(\phi+2\pi N)$ and the previous expression can be valid only for integer values of $\nu$. Moreover, the assumption that the potential does not depend on the angular variable implies that the modulus of the angular solution does not depend on $\phi$, and thus we can write

$$X(\phi)=e^{i\nu\phi}\qquad \nu\in I \qquad (129)$$

Once the angular solution has been decided, a numerical solution to the radial equation Eq. (128) can be obtained by enforcing the proper boundary conditions. In the case we are considering, the wave function must vanish at $\rho=R$ (i.e., on the disk edge). The other boundary that has to be considered is the point $\rho=0$, at which the condition to be satisfied depends on the value of $\nu$. If $\nu$ is even, after a $\pi$ rotation, the wave function assumes the same values, and then on the radial origin we need to impose a Neumann boundary condition, i.e. the vanishing of the first derivative of the wave function. In the other case, a $\pi$ rotation produces a change of the sign of the angular function, and this is reflected in

a Dirichlet condition for the radial function in $\rho = 0$. These results are summarized in the following:

$$\nu \text{ even} \rightarrow \left. \frac{\partial P}{\partial \rho} \right|_{\rho=0} = 0$$

(130)

$$\nu \text{ odd} \rightarrow P(0) = 0$$

Notice that the previous considerations are valid also in the case of soft-wall radial confinement. In such a case it is sufficient to impose the Dirichlet conditions at a distance at which the value of all the considered wave functions equals zero. As an example, let us consider the hard-wall case with $V(\rho) = 0$ within the circular domain $\rho < R$ ($V = \infty$ for $\rho \geq R$). The radial equation becomes

$$\rho^2 \frac{\partial^2 P}{\partial \rho^2} + \rho \frac{\partial P}{\partial \rho} + (K^2 \rho^2 - \nu^2) P(\rho) = 0$$

(131)

This equation is formally identical to the Bessel equation, the solutions of which are called Bessel functions, with the correspondence $P(\rho) = J_\nu(K\rho)$. Let us point out that the solutions of the previous expression do not depend on the sign of $\nu$. A sketch of the first three Bessel functions is shown in Fig. 14. Let us point out that the previously discussed boundary condition at $\rho = 0$ is satisfied: The first derivative vanishes for $\nu = 0$ and $\nu = 2$, whereas the function itself vanishes for $\nu = 1$. To satisfy the boundary condition at $\rho = R$, we need to impose $J_\nu(KR) = 0$. and thus, if we call $Z_{\nu_i}$ the $i$th zero of the Bessel function of order $\nu$, we must have

$$KR = Z_\nu$$

(132)

This last expression leads to the following eigenvalues:

$$E_{\nu_i} = \frac{\hbar^2}{2m} \left( \frac{Z_{\nu_i}}{R} \right)^2$$

(133)

and to the eigenfunctions:

$$\psi_{\nu_i}(\rho, \phi) = A e^{i\nu\phi} J_\nu(K_{\nu_i}\rho)$$

(134)

where $K_{\nu_i} = Z_{\nu_i}/R$, and $A$ is a proper normalization constant. Notice that for $\nu = 0$ the system is nondegenerate, whereas for different values of $\nu$ there is a double degeneracy as a result of the two possible values of $\nu$ with opposite sign.

Let us now discuss the numerical solution of the Schrödinger equation for a realistic system. Quantum dots can be obtained by embedding tiny regions of a semiconductor material within another semiconductor that has a higher conduction band. A simplified geometrical representation is that of a lens-shaped dot with cylindrical symmetry and with the cross section as shown in Fig. 15. Because of the system symmetry, we can work in cylindrical



Figure 14. The first three Bessel functions.

**Figure 15.** Cross section of a lens-shaped dot.

coordinates. By writing the solution as the product of an angular function and a function depending on the radius $\rho$ and the height $z$ only, $\psi(\rho, z, \phi) = U(\rho, z)X(\phi)$, and by using the previous results, it is possible to separate the Schrödinger equation into the following two equations:

$$\frac{1}{X}\frac{\partial^2 X}{\partial \phi^2} + \nu^2 = 0$$

$$\frac{\partial^2 U}{\partial \rho^2} + \frac{1}{\rho}\frac{\partial U}{\partial \rho} + \frac{\partial^2 U}{\partial z^2} - \frac{2m^*}{\hbar^2}VU - \frac{\nu^2}{\rho^2}U = \frac{2m^*}{\hbar^2}EU$$

(135)

where $m^*$ indicates the effective mass of the electrons that, for simplicity, we assume to be the same in the two materials (this approximation will be released in the following). To numerically solve the previous equations, a proper discretization of the three-dimensional space is needed. The goal is to reduce the Schrödinger equation to a standard eigensystem in which the matrix to be diagonalized is as simple as possible. This can be achieved by accurately choosing the discretization mesh, but let us start with a natural, although not optimal, discretization. The $\{\rho, z\}$ space is chosen as shown in Fig. 16, with Dirichlet boundary conditions on all the borders but $\rho = 0$. This last condition must be carefully chosen. For the previously mentioned symmetry reasons, we need a Neumann condition for even values of $\nu$ and a Dirichlet condition for odd $\nu$. There is a problem in enforcing one or the other boundary condition on the same mesh, which can be easily explained with the help of Fig. 17. If we choose a discretization mesh $h$ in such a way that the first discretization point coincides with $\rho = 0$, it is easy to enforce a Dirichlet boundary condition by requiring $U(0, z) = 0$. However, a Neumann condition can be enforced only by setting $U(0, z) = U(1, z) = 0$, which is equivalent to enforcing the boundary condition in $0.5h$ and not in $0$. The Neumann boundary condition is thus displaced by half a mesh, a problem that could be solved by using



**Figure 16.** The discretization space.

Figure 17. Boundary condition problem.

a nonuniform mesh, but that can also be neglected if $h$ is small enough. With the stated approximations, the discretization of the equation for $U(\rho, z)$ reads

$$\frac{U_{i+1,j} + U_{i-1,j} - 2U_{i,j}}{h^2} + \frac{1}{ih}\left[\frac{U_{i+1,j} - U_{i-1,j}}{2h}\right]$$

$$+ \frac{U_{i,j+1} + U_{i,j-1} - 2U_{i,j}}{h^2} - \frac{2m^*}{\hbar^2}V_{i,j}U_{i,j} - \frac{\nu^2}{(ih)^2}U_{i,j} = -\frac{2m^*}{\hbar^2}EU_{i,j} \qquad (136)$$

where $U_{i,j} \equiv U(ih, jk)$ and $h, k$ are the discretization steps along the two directions. This can be represented as a pentadiagonal matrix that, unfortunately, is not symmetric, making its diagonalization more difficult.

It is, however, possible to symmetrize the problem, with a simple transformation. Let us consider the second equation in Eq. (136) and multiply both sides by $\rho$:

$$\frac{\partial}{\partial\rho}\left(\rho\frac{\partial U}{\partial\rho}\right) + \rho\frac{\partial^2}{\partial z^2}U = \rho\frac{2m^*}{\hbar^2}(V - E)U - \frac{\nu^2}{\rho}U \qquad (137)$$

which can be rewritten

$$-\frac{\hbar^2}{2m^*}\frac{\partial}{\partial\rho}\left(\rho\frac{\partial U}{\partial\rho}\right) - \frac{\hbar^2}{2m^*}\rho\frac{\partial^2}{\partial z^2}U + \rho V - \frac{\hbar^2}{2m^*}\frac{\nu^2}{\rho}U = -\rho E \qquad (138)$$

Let us first focus on the discretization of the operator

$$\frac{\partial}{\partial\rho}\left(\rho\frac{\partial U}{\partial\rho}\right) \qquad (139)$$

which leads to the asymmetry in the previous treatment. If we center the derivatives over half-mesh points, the discretization of such a term will read

$$\frac{\rho_{i+1/2}\frac{U_{i+1}-U_i}{h} - \rho_{i-1/2}\frac{U_i-U_{i-1}}{h}}{h} = \rho_{i+1/2}\frac{U_{i+1}-U_i}{h^2} - \rho_{i-1/2}\frac{U_i-U_{i-1}}{h^2} \qquad (140)$$

where $h$ is the discretization step along $\rho$. This is a symmetric operator, as are those deriving from the discretization of the other terms in Eq. (138). However, Eq. (138) does not represent an eigenvalue problem because $\rho$, which is a vector in the discretized version, appears in the right-hand side; that is, we have a problem that can be written as

$$AU = \rho EU \qquad (141)$$

where $A$ is the discretization coefficient matrix. Let us left-multiply both sides by $\rho^{-1/2}$:

$$\rho^{-1/2}AU = \rho^{1/2}EU$$

$$(\rho^{-1/2}A\rho^{-1/2})(\rho^{1/2}U) = E(\rho^{1/2}U) \qquad (142)$$

$$(\rho^{1/2}A\rho^{-1/2})W = EW$$

which is a regular eigenvalue problem: once the eigenvectors $W$ have been computed, the actual wave function $U$ can be obtained by left-multiplying it by $\rho^{-1/2}$.

Let us now look for a more general approach, in which a spatially variable effective mass can be treated as well as a nonuniform mesh.

If we have regions with different effective masses, one way of writing the Schrödinger equation that preserves the probability flux is the following:

$$-\frac{\hbar^2}{2}\nabla\left(\frac{1}{m^*}\nabla\psi\right) + V\psi = E\psi \qquad (143)$$

We should, however, point out that it is not at all rigorous to define an effective mass close to the boundary between two different materials, as the effective mass is a bulk property.

Considering that the effective mass can depend on the position, we obtain a new expression for the action of the Laplacian over the wave function; namely,

$$\nabla\left(\frac{1}{m^*}\nabla\psi\right) = \frac{1}{\rho}\frac{\partial}{\partial\rho}\left(\rho\frac{1}{m^*}\frac{\partial\psi}{\partial\rho}\right) + \frac{1}{\rho}\frac{\partial}{\partial\phi}\left(\frac{1}{m^*}\frac{\partial\psi}{\partial\phi}\frac{1}{\rho}\right) + \frac{\partial}{\partial z}\left(\frac{1}{m^*}\frac{\partial\psi}{\partial z}\right) \qquad (144)$$

Because our system has cylindrical symmetry, $m^*$ does not depend on $\phi$. Therefore, we can write

$$\nabla\left(\frac{1}{m^*}\nabla\psi\right) = \frac{1}{\rho}\frac{\partial}{\partial\rho}\left(\rho\frac{1}{m^*}\frac{\partial\psi}{\partial\rho}\right) + \frac{1}{\rho^2}\frac{1}{m^*}\frac{\partial^2\psi}{\partial\phi^2} + \frac{\partial}{\partial z}\left(\frac{1}{m^*}\frac{\partial\psi}{\partial z}\right) \qquad (145)$$

By inserting the previous expression in the complete Schrödinger equation, and by factoring the wave function $[\psi(\rho, z, \phi) = U(\rho, z)X(\phi)]$, we can separate the problem into the two following equations:

$$\frac{1}{X}\frac{\partial^2 X}{\partial\phi^2} + m^*\Theta = 0 \qquad (146)$$

where $m^*\Theta = \nu^2$ with $\nu$ integer, and

$$\frac{1}{\rho}\frac{\partial}{\partial\rho}\left[\rho\frac{1}{m^*}\frac{\partial U}{\partial\rho}\right] + \frac{\partial}{\partial z}\left(\frac{1}{m^*}\frac{\partial U}{\partial z}\right) - \frac{2}{\hbar^2}VU - \frac{\nu^2}{\rho^2 m^*}U = \frac{2}{\hbar^2}EU \qquad (147)$$

The angular equation has the usual exponential solutions [Eq. (129)] whereas the solution of the equation for $U(\rho, z)$ can be numerically approached if we refer to the discretization scheme of Fig. 18.

The resulting set of equations can be written in a concise notation as

$$\alpha_{i,j}U_{i+1,j} + \beta_{i,j}U_{i-1,j} + \gamma_{i,j}U_{i,j+1} + \delta_{i,j}U_{i,j-1} + \eta_{i,j}U_{i,j} = \frac{2}{\hbar^2}EU_{i,j} \qquad (148)$$

where the definition of the coefficients depends on the boundary conditions. At the left vertiica boundary (i.e., at $\rho = 0$), we need to enforce either Dirichlet or Neumann boundary conditions, depending on the value of $\nu$, as discussed before. To avoid the mentioned problem in enforcing the Neumann condition on the first point, we shall use a uniform mesh eveywhere except for $\rho = 0$. The first mesh interval starting from the origin will be identical

**Figure 18.** Discretization scheme.

to the others in the case of Neumann boundary conditions, and only half a mesh in the case of a Dirichlet boundary condition. Moreover, we shall shift all the mesh points of half a mesh toward the radial origin. This ensures, as shown in Fig. 19, that the boundary condition is always enforced in the same position. With this choice we have

$$\alpha_{i,j} = \frac{2}{h_{i-1} + h_i} \frac{1}{\rho_i} \frac{\rho_{i+1/2}}{m^*_{i+1/2,j} h_i} \qquad \text{for } 1 \le i < N, \quad 0 \text{ otherwise} \tag{149}$$

$$\beta_{i,j} = \frac{2}{h_{i-1} + h_i} \frac{1}{\rho_i} \frac{\rho_{i-1/2}}{m^*_{i-1/2,j} h_{i-1}} \qquad \text{for } 1 < i \le N, \quad 0 \text{ otherwise} \tag{150}$$

$$\gamma_{i,j} = \frac{2}{k_{j-1} + k_j} \frac{1}{k_j} \frac{1}{m^*_{i,j+1/2}} \qquad \text{for } 1 \le j < M, \quad 0 \text{ otherwise} \tag{151}$$

$$\delta_{i,j} = \frac{2}{k_{j-1} + k_j} \frac{1}{k_{j-1}} \frac{1}{m^*_{i,j-1/2}} \qquad \text{for } 1 < j \le M, \quad 0 \text{ otherwise} \tag{152}$$



**Figure 19.** The mesh for Neumann and Dirichlet b.c.

$$\eta_{i,i} = -\frac{2}{h_{i-1}+h_i}\frac{1}{\rho_i}\frac{\rho_{i-1/2}}{m^*_{i+1/2,i}/h_i} - \frac{2}{h_{i-1}+h_i}\frac{1}{\rho_i}\frac{\rho_{i-1/2}}{m^*_{i-1/2,i}/h_{i-1}}$$

$$-\frac{2}{k_{j-1}+k_j}\frac{1}{k_j}\frac{1}{m^*_{i,j+1/2}} - \frac{2}{k_{j-1}+k_j}\frac{1}{k_{j-1}}\frac{1}{m^*_{i,j-1/2}}$$

$$-\frac{2}{\hbar^2}V_{i,i} - \frac{\nu^2}{\rho_i^2} \qquad \text{for } 1 < i \le N \text{ and for } i = 1 \text{ and Dirichlet b.c.} \tag{153}$$

$$\eta_{i,i} = -\frac{2}{h_{i-1}+h_i}\frac{1}{\rho_i}\frac{\rho_{i+1/2}}{m^*_{i+1/2,j}/h_i} - \frac{2}{k_{j-1}+k_j}\frac{1}{k_j}\frac{1}{m^*_{i,j+1/2}}$$

$$-\frac{2}{k_{j-1}+k_j}\frac{1}{k_{j-1}}\frac{1}{m^*_{i,j-1/2}} - \frac{2}{\hbar^2}V_{i,j} - \frac{\nu^2}{\rho_i^2} \qquad \text{for } i = 1 \text{ and Neumann b.c.} \tag{154}$$

## 5.2. Self-Consistent Solution

As previously stated, self-consistency between the solution of the Schrödinger equation and that of the electrostatic problem must be achieved with an iterative procedure. As the dimensions of a quantum dot increase, the electrostatic interaction becomes more and more important compared to the quantum confinement, because the Coulomb energy scales with the inverse of the distance between charges (and therefore with the inverse of the dot size), whereas the confinement energy scales approximately with the square of the dot size. It is thus apparent that achieving convergence will be increasingly difficult for bigger dots. Some improvement can be obtained with the application of underrelaxation techniques, which consist of considering at each iteration a weighed average of the potential obtained from the current solution of the Poisson equation and of that at the previous iteration. In mathematical terms, the potential $V^p(x_i)$ used in the Schrödinger equation at the $p$th iteration is obtained as

$$V^p(x_i) = V^{poiss}(x_i)\alpha + V^{p-1}(x_i)(1-\alpha) \tag{155}$$

where $V^{poiss}(x_i)$ is the potential from the solution of the Poisson equation (or from the evaluation of the Coulomb term, in the simplest cases), $V^{p-1}(x_i)$ is the potential at the previous iteration. and $\alpha$ is the underrelaxation parameter, which varies between 0 and 1 (no underrelaxation). If small values of $\alpha$ are used, convergence will be very slow and particular care must be exercised when choosing the convergence criterion to stop the iterative procedure. Indeed, the most common stopping criterion consists of checking when the modulus of the difference between the eigenvalues or the mean square error between the eigenfunctions at two consecutive iterations is less than a given threshold. If a small value of $\alpha$ is used, such parameters may be very small even if convergence is not achieved at all.

Calculations on quantum dots may have the objective of determining the chemical potential in the quantum dot once the occupancy of the dot is given or, vice versa, of determining the number of electrons in the dot once the Fermi-level in the reservoirs connected to the quantum dot is known. Let us examine the two cases in detail, with the assumption, for the time being, that the temperature is 0 K (finite temperature situations will be discussed in the following), which is often a reasonable approximation, as experiments are usually performed in the tens of millikelvin range.

If the number of electrons in the dot is given, our aim is to compute the chemical potential in the dot, which can be achieved either using the definition itself of chemical potential or applying Slater's transition rule. By definition, the chemical potential $\mu$ of a system is given by the variation of its free energy $F$, as the number of particles in the system is varied by 1:

$$\mu(N) = F(N) - F(N-1) \tag{156}$$

The free energy corresponds, at 0 K, to the internal energy $E(N)$ of the system, which has the expression, if the LDA approximation has been used for the calculation of the eigenenergies [54],

$$E(N) = \sum_{i=1}^{N} \epsilon_i - \frac{1}{2} \iint \frac{e^2}{4\pi\varepsilon_0\varepsilon_r} \frac{n(\mathbf{r})n(\rho)}{|\mathbf{r} - \rho|} d\mathbf{r}\, d\rho$$

$$+ \int n(\mathbf{r}) \Big\{ E_{ex}[n(\mathbf{r})] + E_{corr}[n(\mathbf{r})] - V_{ex}[n(\mathbf{r})] - V_{corr}[n(\mathbf{r})] \Big\} d\mathbf{r} \qquad (157)$$

where $N$ is the total number of electrons, $\epsilon_i$ are the energy eigenvalues for each electron, $n(\mathbf{r})$ is the total electron density, and $E_{ex}[n(\mathbf{r})]$, $E_{corr}[n(\mathbf{r})]$, $V_{ex}[n(\mathbf{r})]$, $V_{corr}[n(\mathbf{r})]$ are the exchange and correlation energies and potentials, respectively. The problem with this approach is that the calculation of the chemical potential involves a subtraction between two terms, $E(N)$ and $E(N - 1)$, which can be rather close to each other, and therefore the result may be affected by significant numerical error.

An alternative approach consists of applying Slater's transition rule [55], which states that a good approximation of the chemical potential for $N$ electrons is given by the eigenvalue corresponding to a fictitious half electron added to a system with $N - 1$ electrons:

$$\mu(N) = \epsilon(N - 0.5) \qquad (158)$$

This approach, in addition to being much less heavy from the computational point of view, leads in general to much better precision.

The determination of the number of electrons in the dot, if the chemical potential of the external reservoirs is given, requires a procedure for the minimization of the free energy of the system, which, in this case, includes also the external voltage sources. Stopa [48] has worked out this problem in detail while investigating, with a three-dimensional approach, conductance through a quantum dot defined electrostatically in a GaAs/AlGaAs heterostructure. The expression for the free energy provided by Stopa includes the sum of the energy eigenvalues of the electrons, two terms for the removal of the double counting of the interaction energies, and a term representing the work done by the external voltage sources:

$$F(n_j, N, V_i) = \sum_j n_j\epsilon_j - \frac{1}{2}\int d\mathbf{r}\, \rho_{el}(\mathbf{r})\phi_{sc}(\mathbf{r}) + \frac{1}{2}\int d\mathbf{r}\, \rho_{ion}(\mathbf{r})\phi_{sc}(\mathbf{r}) - \frac{1}{2}\sum_i Q_i V_i \qquad (159)$$

where the sum over $j$ is over orbitals, $n_j$ represents the occupancy of the $j$th orbital, $\rho_{el}(\mathbf{r})$ is the local electron density in $\mathbf{r}$, $\phi_{sc}(\mathbf{r})$ is the self-consistent electric potential, $\rho_{ion}(\mathbf{r})$ is the charge density resulting from the donor and impurity ions, and $Q_i$ and $V_i$ are the charge and voltage on the $i$th electrode, respectively. Because electrons will tunnel into the quantum dot until its chemical potential equals that $E_f$ of the external reservoirs, the equilibrium condition will be reached when $\mu(N, V_i) = E_f$; that is,

$$F(N, V_i) - F(N - 1, V_i) = E_f \qquad (160)$$

Sometimes, in quantum dot calculations at a finite temperature, a Fermi–Dirac distribution function is assumed to compute the energy levels of the quantum dot. This is in principle not correct, although it is often not too far from the exact result. Indeed, the Fermi–Dirac distribution function is obtained in statistical physics as the limit for a system containing an infinite number of fermions. The number of electrons in a quantum dot is in general relatively small, and therefore the actual occupancy of the levels would be given by the Gibbs distribution, whose calculation, however, is quite challenging, because it requires knowledge of the partition function of the system, and therefore of the energy associated with all possible electron configurations. In the literature there are some proposals to compute the Gibbs distribution function by means of Monte Carlo techniques, which allow an estimation of the configuration energies with an acceptable computational burden.

# 6. SIMULATION OF MULTIPLE INTERACTING QUANTUM DOTS

Multiple quantum dot systems pose new simulation challenges compared to single quantum dots, in particular because of increased difficulties in achieving a self-consistent solution and of the larger size of the problem.

A simplified model that can already provide interesting results consists in the Hubbard formalism, which is usually known as "occupation number Hamiltonian," as it neglects the detailed electronic structure of each single quantum dot. Each dot, and its interaction with the others, is described by means of phenomenological parameters: the tunneling energy, the confinement energy in the dots, and the electrostatic interaction. The success of such a description is in capturing the overall behavior and in providing a qualitative understanding of the underlying physics, but quantitative predictions lie beyond the possibilities of such methods. Hubbard-like formulations have been used, for example, by Tougaw et al. [15] to study the cell-to-cell response function for a Quantum Cellular Automaton (QCA) cell, or by Das Sarma et al. [57] and Klimeck et al. [58] to analyze the appearance of collective Coulomb Blockade phenomena in arrays of quantum dots, or more general transport phenomena in coupled quantum dots.

As we have seen in the previous section, rather refined techniques based on iterative self-consistent procedures can be developed for the analysis of single quantum dots. These approaches, however, fail in reaching convergence when the electrostatic interaction in the system is comparable to the confinement energy and when strong degeneracies resulting from symmetries are present, as in the case of QCA cells.

To overcome these difficulties, noniterative methods have been developed, based on techniques typical of molecular chemistry, that allow the detailed simulations of realistically described multiple quantum dot systems. In particular, we shall focus on the configuration–interaction technique.

## 6.1. Occupation-Number Hamiltonian

The occupation-number Hamiltonian represents a relatively straightforward approach requiring limited computational resources, but it has the drawback of relying on the introduction of the tunneling energy $t$ and of other quantities as phenomenological parameters.

The system is described by means of an occupation number Hamiltonian, $|n_{1,\uparrow}, n_{1,\downarrow};$ $n_{2,\uparrow}, n_{2,\downarrow}; \dots n_{N,\uparrow}, n_{N,\downarrow}\rangle$, where $n_{i,\uparrow\downarrow}$ indicates the occupation number of the $i$th dot by electrons with up or down spin. On this basis the Hamiltonian of the system can be written in terms of the creation and annihilation operators $b_{j,\sigma}$ and $b_{j,\sigma}^\dagger$ that create or annihilate, respectively, an electron of spin $\sigma$ in the $i$th dot:

$$H_0 = \sum_{i,\sigma} E_{0,i} n_{i,\sigma} + \sum_{i>j,\sigma} t(b_{i,\sigma}^\dagger b_{j,\sigma} + b_{j,\sigma}^\dagger b_{i,\sigma}) + \sum_i E_{Qi} n_{i,\uparrow} n_{i,\downarrow} + \sum_{i>j,\sigma,\sigma'} V_Q \frac{n_{i,\sigma} n_{j,\sigma'}}{|\vec{R}_i - \vec{R}_j|} \quad (161)$$

where $E_{0,i}$ is the ground-state energy of the single, isolated $i$th dot; $n_{j,\sigma} \equiv b_{j,\sigma} b_{j,\sigma}^\dagger$ is the number operator for electrons in the $i$th dot with spin $\sigma$; $t$ is the tunneling energy between adjacent dots; $V_Q$ is equal to $e^2/(4\pi\epsilon)$, with $e$ being the electron charge and $\epsilon$ the dielectric permittivity of the medium; $E_{Qi}$ is the on-site charging energy for the $i$th dot; and $\vec{R}_i$ is the position of the center of the $i$th dot.

The nearest-neighbor assumption for tunneling can be lifted, but this would imply the definition of tunneling energies that depends on the distance, making the phenomenological description of $t$ more difficult. In this simplified case, the tunneling energy $t$ can be related to the actual potential confining the dots by means of simple approximations. It is usually estimated to be equal to one half of the level splitting caused by coupling between neighboring dots.

If we rewrite this Hamiltonian in the representation corresponding to the occupation number basis, we obtain a sparse matrix that can be diagonalized by means of standard procedures. The lowest eigenvalue obtained from the diagonalization corresponds to the

ground state $|\psi_0\rangle$ of the system, and the associated electronic configuration can be evaluated, computing the average number of electrons in each dot as

$$\rho_i = \langle \psi_0 | n_{i\uparrow} + n_{i\downarrow} | \psi_0 \rangle \tag{162}$$

## 6.2. Configuration–Interaction Method

Other, more accurate, "one-shot" methods are based on the representation of the many-electron wave function on a basis of Slater determinants. Within this representation, the computation of the many-electron system ground state can be written as an algebraic eigenvalue problem. Techniques of this kind are typical of molecular chemistry [59], among which we describe here the configuration–interaction method and its specific application to the study of multiple-dot systems [60].

Let us consider the following $N$-electron Hamiltonian, corresponding to a system of quantum dots located at a distance $z$ from a conducting plane:

$$\widehat{H}_1 = \sum_{i=1}^{N} \left[ -\frac{\hbar^2}{2m^*} \nabla_i^2 + V(\mathbf{r}_i) \right] + \sum_{i,j} \frac{1}{4\pi\varepsilon} \frac{e^2}{|\mathbf{r}_i - \mathbf{r}_j|} - \frac{1}{4\pi\varepsilon} \frac{e^2}{\sqrt{|\mathbf{r}_i - \mathbf{r}_j|^2 + (2z)^2}} - \frac{1}{4\pi\varepsilon} \frac{e^2}{2z} \tag{163}$$

where $N$ is the number of electrons; $V(\mathbf{r}_i)$ is the confinement potential seen by the $i$th electron; $\varepsilon$ is the dielectric permittivity; $\mathbf{r}_i$ and $\mathbf{r}_j$ are the coordinates of the $i$th and of the $j$th electron, respectively; $z$ is the distance between the plane containing the quantum dots (i.e., the 2DEG plane) and the surface of the heterostructure (for which we assume Fermi-level pinning); and $e$ is the electron charge. The last two terms correspond to the contribution to the total energy from the images resulting from the presence of the conducting plane.

Dealing with confined systems, we can consider a numerable complete basis $\{\varphi_i(q)\}$, where $q$ denotes all the electron coordinates (both spatial and spin), to the elements which we shall refer as spin-orbitals. Using such a basis, all possible independent Slater determinants can be built:

$$\Phi_k = \frac{1}{\sqrt{N!}} \begin{vmatrix} \varphi_{n_{1k}}(q_1) & \varphi_{n_{2k}}(q_1) & \cdots & \varphi_{n_{Nk}}(q_1) \\ \varphi_{n_{1k}}(q_2) & \varphi_{n_{2k}}(q_2) & \cdots & \varphi_{n_{Nk}}(q_2) \\ \cdots & \cdots & \cdots & \cdots \\ \varphi_{n_{1k}}(q_N) & \varphi_{n_{2k}}(q_N) & \cdots & \varphi_{n_{Nk}}(q_N) \end{vmatrix} \tag{164}$$

where $k$ labels the Slater determinants and the integer $n_{ik}$ specifies which spin-orbital appears in the $i$th column of the $k$th Slater determinant. The set $\{\Phi_k\}$ is a complete orthonormal basis for the $N$-electron eigenfunctions $\Psi_i$ of the Hamiltonian (163) [61]:

$$\Psi_i = \sum_{k=1}^{\infty} c_{ik} \Phi_k \tag{165}$$

The eigenfunctions of $\widehat{H}$ can be obtained by solving the secular equation

$$\mathscr{H} \mathbf{c}_i = E_i \mathbf{c}_i \tag{166}$$

where the infinite-dimensional "Hamiltonian matrix" is

$$\mathscr{H}_{kk'} = \langle \Phi_k | \widehat{H} | \Phi_{k'} \rangle \tag{167}$$

$E_i$ is the $i$th eigenvalue of $\mathscr{H}$, and the vector $\mathbf{c}_i$ contains the coefficients $c_{ik}$ of the expansion. To treat the problem numerically, this approach cannot be implemented "exactly": We must consider a finite set of $M$ spin-orbitals $\{\varphi_i(q)\}$, with $i = 1 \ldots M$. From $N$-electrons and $M$ spin-orbitals (with $M \geq N$), $\mathscr{N}_{SD}$ different Slater determinants will be built, where

$$\mathscr{N}_{SD} = \binom{M}{N} \tag{168}$$

The secular equation [Eq. (166)] thus represents a Hermitian $\mathcal{N}_{SD} \times \mathcal{N}_{SD}$ eigenvalue problem. Computation of the matrix elements of the Hamiltonian [Eq. (163)] between two Slater determinants is not a trivial task [61]. For the diagonal elements one finds:

$$\langle \Phi_k | \hat{H} | \Phi_k \rangle = \sum_i \langle \varphi_{n_{ik}} | h | \varphi_{n_{ik}} \rangle + \frac{1}{2} \sum_{ij} (\langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{ik}} \varphi_{n_{jk}} \rangle - \langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{jk}} \varphi_{n_{ik}} \rangle) \qquad (169)$$

where $h$ represents the "one-body" terms of Eq. (163), $g$ represents the "two-body" terms, and in general,

$$\langle \varphi_i \varphi_j | g | \varphi_l \varphi_m \rangle = \int dq_1 \, dq_2 \, \varphi_i^*(q_1) \varphi_j^*(q_2) g(\vec{r}_1, \vec{r}_2) \varphi_l(q_1) \varphi_m(q_2) \qquad (170)$$

"Selection rules" (Slater's rules [59, 61]) exist and simplify the computation of the off-diagonal matrix elements between two different Slater determinants $\Phi_k$, $\Phi_{k'}$. There are only two possible cases in which $\langle \Phi_k | \hat{H} | \Phi_{k'} \rangle$ does not vanish (i.e., when $\Phi_k$, $\Phi_{k'}$ either differ by one spin-orbital or by two):

1. one spin-orbital difference ($\varphi_{n_{ik}} \neq \varphi_{n_{ik'}}$)

$$\langle \Phi_k | \hat{H} | \Phi_{k'} \rangle = \langle \varphi_{n_{ik}} | h | \varphi_{n_{ik'}} \rangle + \sum_{j \neq i} (\langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{ik'}} \varphi_{n_{jk}} \rangle - \langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{jk}} \varphi_{n_{ik'}} \rangle) \qquad (171)$$

2. two spin-orbital difference ($\varphi_{n_{ik}} \neq \varphi_{n_{ik'}}$ and $\varphi_{n_{jk}} \neq \varphi_{n_{jk'}}$)

$$\langle \Phi_k | \hat{H} | \Phi_{k'} \rangle = \langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{ik'}} \varphi_{n_{jk'}} \rangle - \langle \varphi_{n_{ik}} \varphi_{n_{jk}} | g | \varphi_{n_{jk'}} \varphi_{n_{ik'}} \rangle \qquad (172)$$

The expressions in Eqs. (171) and (172) refer to the case in which the spin-orbitals that are common to both Slater determinants occur in the same columns. If this is not the case, it is possible to perform a permutation of the columns of one determinant so that the above condition is satisfied; the permutation has the effect of changing the sign of the matrix element if it is of an odd order.

Finally, it is worth pointing out that Eqs. (169), (171) and (172) are valid only if the orthonormality condition on the spin-orbitals is satisfied.

We notice that, in virtue of Eqs. (169), (171) and (172), the number of nonzero matrix elements of the total Hamiltonian between two generic Slater determinants is less than $(\mathcal{N}_{SD})^2$ [see Eq. (168)] and is given by the following expression:

$$\binom{M}{N} \times \left[ \binom{M-N}{1} \binom{N}{1} + \binom{M-N}{2} \binom{N}{2} + 1 \right]$$

$$= \binom{M}{N} \times \left[ (M-N)N + \frac{(M-N)(M-N-1)N(N-1)}{4} + 1 \right] \qquad (173)$$

Once the $N$-electron wave function $\Psi_i$ has been obtained, the corresponding electron density is simply given by

$$\rho_i(\vec{r}) = \sum_{s_1} N \int |\Psi_i(\vec{r}, s_1, q_2, \ldots, q_N)|^2 \, dq_2 \ldots dq_N \qquad (174)$$

For a system with no more than one electron per dot, a limited number of spin orbitals, built from the single-particle eigenfunctions of each dot, is needed for a proper representation of the wave function $\Psi_i$, because, in the absence of a strong electrostatic interaction (such as that of two electrons in the same dot) the components of the wave functions do not undergo too large a deformation. The situation changes if there are more than one electron in each dot. In such a case, a larger basis is needed, because the wave function components are significantly distorted and the number of determinants to be considered grows very quickly. For example, to study a QCA cell with four dots and single dot occupancy of one electron or less, 24 determinants are sufficient to achieve good precision, whereas

if we allow the presence of two electrons in each dot, a number of determinants of the order of 600 must be included. There are several selection techniques that help in selecting the determinants which give a more significant contribution to the correct solution, but the combinatorial growth of the problem as the number of electrons is increased remains the main limitation of this otherwise powerful technique.

Another limitation is represented by the fact that, if the electrostatic interaction cannot be simply represented with terms such as those appearing in Eq. (163) but has instead to be treated with the solution of the Poisson equation, the calculation of the matrix elements relative to the two-body terms $g$ becomes much more involved. Such terms correspond to the Green function of the Poisson equation, which must be evaluated numerically for each pair of grid points, an extremely complex numerical test.

# 7. SIMULATION OF SINGLE-ELECTRON CIRCUITS

Single-electron circuits are characterized by the dominant role played by Coulomb blockade effects, which, to become observable, require that the charging energy be large compared to the thermal energy $kT$ and that strong electron localization be present on the nodes, as a consequence of tunneling resistances being much larger than the inverse of the conductance quantum $2e^2/h$.

These conditions are required for the applicability of the so called "orthodox theory" of the Coulomb blockade pioneered by Kulik and Shekhter [62], which has played a key role in the development of single electronics by Likharev and others.

When the orthodox theory is applicable, electrons behave as classical particles localized on the nodes and moving from one node to another, if this is energetically allowed, with instantaneous transitions. This completely classical (except for the tunneling effect) picture is at the basis of the application of the Monte Carlo method to the numerical simulation of single-electron devices [13, 14, 63–65].

A typical single-electron circuit can be described as a network of dots (nodes) that are connected, among them or to external leads, by standard capacitors and by tunneling capacitors (i.e. capacitors across which tunneling is possible with a nonvanishing probability). Tunneling capacitors are described as standard capacitors with a tunneling resistance in parallel. An example of a single-electron circuit is shown in Fig. 20, where the tunneling junctions are represented as double boxes.

The behavior of such a circuit cannot be studied with SPICE-like simulators (i.e., simulators that derive currents and voltages in the circuit from the current-voltage characteristics of the single devices). SPICE simulations of circuits containing single-electron devices yield



Figure 20. Single-electron device: a single-electron transistor.

reliable results only if the capacitances on the nodes connecting different devices are much larger than those present in the devices themselves, so that no Coulomb Blockade effect takes place outside the device (whose internal $I$–$V$ characteristics have been determined analytically or with some other technique based on the orthodox theory).

The mentioned conditions for the applicability of SPICE are in general not satisfied in a circuit like the one in Fig. 20, and a different approach to circuit simulation is thus needed, one that is capable of properly treating Coulomb blockade effects.

Two main approaches have been developed to treat this problem based on the Monte Carlo technique or on the master equation formalism.

Let us first examine the approach based on the Monte Carlo method, which is more intuitive. The relevant quantity for the determination of the probabilities of electron transitions through the tunnel junctions, to be used in the Monte Carlo calculation, is the free energy: Transitions that lower the free energy of the system are favored, whereas those that raise it are not (they are not allowed if the temperature is zero).

As already discussed in the previous sections, the free energy of a system corresponds to the total energy of the system (in this case the electrostatic energy stored in the capacitors) minus the work done by the voltage sources:

$$E = \frac{1}{2}(\mathbf{q} \quad \mathbf{q'})\,\mathbf{C}^{-1}\begin{pmatrix}\mathbf{q}\\\mathbf{q'}\end{pmatrix} - \mathbf{v}\cdot\mathbf{q'} \tag{175}$$

where $\mathbf{v}$ is the vector of the voltages on the leads and $\mathbf{q}$, $\mathbf{q'}$ are the vectors of the charge on the nodes and on the leads, respectively. The term $\mathbf{C}^{-1}$ is the inverse of the capacitance matrix, which describes the capacitive couplings among all circuit nodes. The capacitance matrix is defined in the following way: Once the circuit nodes are numbered, the diagonal term $i$ has the value of the sum of all the capacitances that are connected to the $i$th node. Each off-diagonal term $C_{ij}$ has the value of the capacitance connecting node $i$ to node $j$, multiplied by $-1$.

Once the free energy associated with the initial configuration of the circuit has been determined, we can consider all the configurations that can be reached from the initial configuration as a consequence of single electron transitions and, by using Eq. (175), the energy difference between the new configuration and the previous one is evaluated. Thus we obtain a set of energy differences $\Delta E_i$, each corresponding to a single electron transition through a different junction. From these free-energy variations we can compute the corresponding probability rates $\Gamma_i$, using the orthodox Coulomb blockade theory [66]:

$$\Gamma_i = \frac{1}{e^2 R_T}\frac{\Delta E_i}{1 - e^{-\frac{\Delta E_i}{k_B T}}} \tag{176}$$

where $R_T$ is the tunneling resistance of the junction relative to the considered transition.

We thus have all the information needed for treating the system evolution via the Monte Carlo approach, but the details of the simulation will differ, depending on the simulation goal. Let us describe two different examples: in the first example, we are interested in the stationary behavior, whereas in the second one we want to study the temporal evolution.

In both cases the first step is to enumerate all the possible configuration changes associated with a single-electron transition. For each new configuration we compute the transition rate $\Gamma_i$ and we indicate the sum of all the transition rates with $\Gamma_{tot}$. With this choice, the probability density function for the time interval $\Delta t$ between two consecutive transitions can be written as

$$f(\Delta t) = \Gamma_{tot}\,e^{-\Gamma_{tot}\Delta t} \tag{177}$$

Thus the $\Delta t$ for each simulation step can be obtained statistically by extracting a random number distributed following Eq. (177). The generation of random numbers with a given distribution is the core of any Monte Carlo algorithm and can be approached in different ways. The methods most commonly used are the Transformation method and the Rejection method [67, 68]. The Transformation method is usually the faster between the two, but it

can be used only for a restricted set of distribution functions; namely, when the inverse function of the indefinite integral of the distribution function is known analytically or at least, is easily numerically computable. This is the case in our situation, and thus we shall give a brief description of the method.

Let us start from the transformation law for probability distributions. Consider a random variable $x$ distributed with a probability distribution $p(x)$ and a change of variable described by a certain function $y(x)$. The distribution of the new variable $y$ will be connected to the old distribution by the Jacobian of the transformation; namely, in the one-dimensional case, $p(y) = p(x)|dx/dy|$. Let us now suppose that we want to obtain a random variable $y$ distributed according to a given function $p(y) = f(y)$ and that the initial variable $x$ is chosen as a uniform deviate $[p(x) = 1$ if $0 < x < 1; p(x) = 0$ otherwise]. Then the relation between the new and the old variable can be obtained by solving the following differential equation: $f(y) = dx/dy$ that leads to

$$y(x) = F^{-1}(x) \tag{178}$$

where $F(y)$ is the indefinite integral of $f(y)$.

The inverse of the indefinite integral of the probability function described by Eq. (177) can be calculated analytically, and thus the time between two consecutive transitions is obtained starting from a uniformly distributed random number $r$ in the $[0, 1]$ interval and applying the transformation

$$\Delta t = -\frac{\ln(r)}{\Gamma_{tot}} \tag{179}$$

Once the time step has been obtained, we must decide which transition actually occurred. This is done by normalizing the rates in a way that their sum $\Gamma_{tot}$ is equal to one. The rates of the different transitions can thus be represented with intervals in the $[0, 1]$ range, each one with a length proportional to the associated probability. Another uniformly distributed random number in the $[0, 1]$ interval will be used to choose one of the intervals and the associated transition.

This procedure is repeated a sufficient number of times to obtain the quantity of interest with the desired statistical precision. Typical quantities of interest are the charge on the nodes, which is calculated via a simple average over the different realizations, or the current through a junction or a lead, which will be obtained taking into account also the transition times for each event, as obtained by the previously described Monte Carlo procedure.

This procedure is quite efficient, especially in the case in which the transition rates are small, because the appropriate time step is dictated by the transition probability itself. However, it is clear that the described approach cannot be used in the nonstationary case, as, for example, when the externally applied voltages change on a timescale that is comparable with the transition time.

In such a case, the evolution of the external voltages sets the timescale, and thus a different approach has to be considered that is closer to the standard solution of equations of motion. A proper time step $\delta t$ must be chosen that is smaller than the characteristic time between two transitions (a proper choice would be to consider an adaptive time step) because changes in the external conditions (the external voltages) can strongly affect the transition rates. Starting from the initial conditions, the rates of all possible transitions are computed as shown before, and then by comparing the total transition probability $P = \Gamma_{tot}\delta t$ with a random number, it is possible to decide whether a transition has occurred. If it has, we have to choose again among the various transitions, as described earlier, and then to increase the time and update the electronic configuration; moreover, the external voltages are to be changed according to their temporal dependence. If no transition has occurred, only the external voltages are updated. In any case, new transition rates must be computed, because of the changed conditions. Each single-electron transition can be observed, but each simulation represents a different realization of a random process. Averaged values must be obtained by repeating the same simulation several times, starting from the same initial condition (but using a different seed to start the random number generator); that is, performing ensemble averages.

As it should be apparent from the previous discussion, single-electron circuits are completely classical stochastic systems whose evolution is driven by a quantum process: electron tunneling. The just-described Monte Carlo method takes directly into account this randomness, generating random transitions and thus obtaining the desired quantities by means of time or ensemble averaging.

Another approach to the simulation of single-electron circuits is based on the solution of the master equation describing the stochastic system [69]. This equation defines the time evolution of the occupation probability for each state, in our case the occupation probability of each node in the circuit, and can be written as

$$\frac{\partial P_i(t)}{\partial t} = \sum_{j \neq i} \left[ \Gamma_{ij} P_j(t) - \Gamma_{ji} P_i(t) \right] \tag{180}$$

where $P_i(t)$ indicates the occupation probability of the $i$th node as a function of time and $\Gamma_{ij}$ the rate for transition from node $i$ to node $j$.

The main assumption behind this equation is that we deal with a Markov process (i.e., transition processes possess no memory), which is exactly the same assumption underlying the implementation of the Monte Carlo method. Equation (180) can be solved numerically, as was done for example in Ref. [14]. The main problem to be faced when trying to solve a master equation is the inclusion of all the relevant states, which can be a very large number if the system is not trivial. Typical circuits actually possess an infinite number of states, because the number of unbalanced electrons in a node is not bounded. It is clear that states with a very large number of excess electrons are quite unlikely, but it is almost impossible to exclude states a priori. Moreover, even if we consider only charge configurations consisting of nodes filled with a single charge, or hole, or empty, the number of states scales as $2^{3N/2}$. The best way to proceed is to use adaptive techniques, in which the set of states, starting from a reduced set, is changed at every step by excluding the states that have too small a probability to be populated and by adding new states until convergence in the results is reached. To proceed in this way, it is, for example, possible to choose a probability threshold $P_{th}$ and to disregard all the states that have a occupation probability smaller than such a value. This naturally depends on the actual configuration, and thus at every iteration new states may be added to the initial set, which could quickly increase the number of states, although usually convergence is easily reached.

If we define the vector of state probability $p \equiv [P_1(t), P_2(t) \ldots P_N(t)]$, Eq. (180) can be written as $\dot{p} = \Gamma p$, where $\Gamma$ is the matrix of transition rates. The formal solution can be written as $p(t) = \exp[\Gamma t] p(0)$, and thus the numerical calculation, once the basis set has been built, reduces to the evaluation of the exponential of a matrix. This is not a simple task: One of the best algorithms to obtain this result relies on the use of the Padé approximant [70], but reliable results are produced only when the exponent is small. In our case this implies choosing a very small time step, and thus an increase of the computational time. A better solution is to use Krylov subspace techniques [71], as such methods allow us to isolate the dominant eigenvalues, and thus to reduce the computational complexity. The main idea behind the Krylov technique to compute $f(A)v$ for analytical function $f$ of the matrix $A$ is to approximate this expression by projecting onto a small subspace defined by the repeated application of the matrix to the initial vector: $K_m = \{v, Av, \ldots, A^{m-1}v\}$. In this way, the problem is reduced to the calculation of $f(H_m)$, where $H_m$ is the tridiagonal matrix obtained applying the Lanczos method with $v/\|v\|^2$ as initial vector. Once the tridiagonal matrix has been computed, the Krylov approximation reads

$$f(A)v \approx \|v\|^2 V_m f(H_m) e_1 \tag{181}$$

where $V_m$ is the orthonormal basis of the subspace obtained with the Lanczos method and $e_1$ is the first column of the $m \times m$ identity matrix. The advantage of such method, compared to other polynomial approximations of the exponential of matrices, is that no information on the matrix spectrum is needed.

## 7.1. Many-Electron Processes

A further step in the simulation of circuits made up of quantum dots is the inclusion of cotunneling effects [14, 72, 73]. Cotunneling is a quantum coherent process whereby two or more electrons simultaneously change their state. This process is clearly less common than the single tunneling event and can be usually neglected in the simulation of single-electron devices, but it can become important or even dominant in particular cases in which the single electron tunneling processes are not energetically viable (i.e., in the Coulomb blockade regime), whereas simultaneous many-electron tunneling transitions can lead to a lower energy state.

The rate of an $N$th order cotunneling process can be written as [74, 75]:

$$\Gamma^{(N)} = \frac{2\pi}{\hbar} \left( \prod_{i=1}^{N} \alpha_i \right) \int_0^{\infty} S^2(\omega_1, \ldots, \omega_{2N}) \times \delta\left( \Delta E_N + \sum_{i=1}^{2N} \omega_i \right) \prod_{i=1}^{2N} [1 - f(\omega_i)]\, d\omega, \qquad (182)$$

where $\Delta E_N$ is the change in electrostatic energy in the entire cotunneling process and $f$ is the Fermi function. $S$ is defined as

$$S(\omega_1, \ldots, \omega_{2N}) = \sum_{\text{permutation}\{(k_1, \ldots, k_N)\}} \prod_{i=1}^{N-1} \frac{1}{\epsilon_k} \qquad (183)$$

where $\epsilon_k = \Delta E_k + \sum_{l=1}^{k}(\omega_{2l-1} + \omega_{2l})$ are the increments in the total energy of the system during cotunneling to an intermediate state, $(\omega_{2l-1} + \omega_{2l})$ is the energy of the electron-hole excitation created, and permutations are over all the possible sequences of intermediate states.

A rigorous simulation of this effect is quite a difficult task, but approximate solutions can be used if we restrict ourselves to the simplest case, namely, the two-electron case. An approximate two-electron tunneling expression has been proposed by Fonseca et al. [14], who wrote the two-electron cotunneling probability as

$$\Gamma = \frac{\hbar}{12\pi e^4} \frac{1}{R_T^{(1)}} \frac{1}{R_T^{(2)}} \left[ \frac{1}{\Delta E^{(1)} - \Delta E/2} + \frac{1}{\Delta E^{(2)} - \Delta E/2} \right]^2 \frac{\Delta E}{\exp(\frac{\Delta E}{k_B T}) - 1} [(\Delta E^2 + (2\pi k_B T)^2)]$$

$$(184)$$

where $\Delta E^{(1,2)}$ represents the energy change resulting from the transitions and $R_T^{(1,2)}$ indicates the respective tunneling resistances. This is not an expression suitable for numerical simulation: divergences occur when one of the intermediate transitions matches one half of the total energy difference. This problem can be solved by means of accurate numerical procedures [14] or, if we are interested in the order of magnitude of the cotunneling, by resorting to an approximate expression:

$$\Gamma = \frac{16}{3} \frac{\hbar}{\pi} \frac{C^2}{e^8 R_T^2} \frac{\Delta E}{\exp(\frac{\Delta E}{k_B T}) - 1} [(\Delta E^2 + (2\pi k_B T)^2)] \qquad (185)$$

where $R_T$ and $C$ represent, respectively, the average value of the tunneling resistance and that of the capacitance of the junctions involved in the cotunneling event. This expression has been shown [14] to give the correct order of magnitude of the effect.

## REFERENCES

1. G. Klimeck, F. Oyafuso, R. C. Bowen, T. B. Boykin, T. A. Cwik, E. Huang, and E. S. Vinyard, *Superlattices Microstruct.* 31, 171 (2002)
2. Y. M. Niquet, C. Delerue, G. Allan, and M. Lannoo, *Phys. Rev. B* 65, 165334 (2002).
3. M. Holm, M. E. Pistol, and C. Pryor, *J. Appl. Phys.* 92, 932 (2002).
4. C. Chamon, M. Oshikawa, and I. Affleck, *Phys. Rev. Lett.* 91, 206403 (2003).
5. J. U. Kim, I. V. Krive, and J. M. Kinaret, *Phys. Rev. Lett.* 90, 176401 (2003).
6. T. Dittrich, P. Hänggi, G.-L. Ingold, B. Kramer, G. Schon, and W. Zwerger, "Quantum Transport and Dissipation." Wiley, Weinheim, 1998.

7. A. O. Caldeira and A. J. Leggett, *Ann. Phys. (N.Y.)* 149, 374 (1983).

8. Y. Yahin, T. C. Au-Yeung, W. Z. Shangguan, and C. H. Kam, *J. Phys. Condensed Matter* 14, 703 (2002).

9. R. Venugopal, M. Paulsson, S. Goasguen, S. Datta, and M. S. Lundstrom, 93, 5613 (2003).

10. I. Knezevic and D. K. Ferry, *Physica E* 19, 71 (2003).

11. C. Texier and M. Büttiker, *Phys. Rev. B* 62, 7454 (2000).

12. T. Quarles, A. R. Newton, D. O. Pederson, and A. Sangiovanni-Vincentelli, "SPICE 3 Version 3F5 User's Manual," University of California, Berkeley, 1994.

13. C. Wasshuber, "Computational Single-Electronics," Springer, Berlin, 2001.

14. L. R. C. Fonseca, A. N. Korotkov, K. K. Likharev, and A. A. Odinstov, *J. Appl. Phys.* 78, 3238 (1995).

15. P. D. Tougaw, C. S. Lent, and W. Porod, *J. Appl. Phys.* 74, 3558 (1993).

16. R. Landauer, *IBM J. Res. Dev.* 1, 223 (1957).

17. M. Büttiker, Y. Imry, R. Landauer, and S. Pinhas, *Phys. Rev. B* 31, 6207 (1985).

18. V. A. Khlus, *Sov. Phys. JETP* 66, 1243 (1987).

19. G. B. Lesovik, *JETP Lett.* 49, 592 (1989).

20. M. Büttiker, *Phys. Rev. Lett.* 65, 2901 (1990).

21. D. J. Thouless and S. Kirkpatrick, *J. Phys. C.* 14, 235 (1981).

22. P. A. Lee and D. S. Fisher, *Phys. Rev. Lett.* 47, 882 (1981).

23. F. Guinea and J. A. Vergés, *Phys. Rev. B* 35, 979 (1987).

24. F. Sols, M. Macucci, U. Ravaioli, and Karl Hess, *J. Appl. Phys.* 66, 3892 (1989).

25. M. Macucci, A. Galick, and U. Ravaioli, *Phys. Rev. B* 52, 5210 (1995).

26. T. Kawamura and J. P. Leburton, *Phys. Rev. B* 48, 8857 (2003).

27. E. N. Economou, "Green's Functions in Computational Physics," Springer, Berlin, 1983.

28. L. Bonci, M. Macucci, D. Guan, and U. Ravaioli, *J. Comp. Electronics*, 2, 127 (2003).

29. A. D. Stone and A. Szafer, *IBM J. Res. Develop.* 32, 384 (1988).

30. S. Datta, "Electronic Transport in Mesoscopic Systems," p. 125. Cambridge University Press, Cambridge, 1995.

31. M. Governale and D. Böse, *Appl. Phys. Lett.* 77, 3215 (2000).

32. Y. Takagaki and K. Ploog, *Phys. Rev B* 53, 3885 (1996).

33. J. J. Palacios and C. Tejedor, *Phys. Rev. B* 48, 5386 (1993).

34. H. Tamura and T. Ando, *Phys. Rev. B* 44, 1792 (1991).

35. J. H. Davies, I. A. Larkin, and E. V. Sukhorukov, *J. Appl. Phys.* 77, 4504 (1995).

36. J. D. Jackson, "Classical Electrodynamics," p. 20. Wiley, New York, 1962.

37. I. A. Larkin and E. V. Sukhorukov, *Phys. Rev. B* 49, 5498 (1994).

38. P. M. Morse and H. Feshbach, "Methods of Theoretical Physics," McGraw-Hill, New York, 1953.

39. I. S. Gradshtein and I. M. Ryhzik, "Tables of Integrals, Series and Products," Academic Press, New York, 1980.

40. Phantoms Computational hub: at http://www.vonbiber.iet.unipi.it

41. J. Martorell and D. W. L. Sprung, *Phys. Rev. B* 49, 13750 (1994).

42. J. Martorell, H. Wu, and D. W. L. Sprung, *Phys. Rev. B* 50, 17298 (1994).

43. G. Iannaccone, M. Macucci, E. Amirante, Y. Jin, H. Launois, and C. Vieu, *Superlattices Microstruct.* 27, 369 (2000).

44. M. Chen and W. Porod, *J. Appl. Phys.* 75, 2545 (1994).

45. G. Fiori, G. Iannaccone, and M. Macucci, *J. Comp. Electronics* 1, 39 (2002).

46. A. Kumar, S. E. Laux, and F. Stern, *Phys. Rev. B* 42, 5166 (1990).

47. P. A. Maksym and T. Chakraborty, *Phys. Rev. Lett.* 65, 108 (1990).

48. M. Stopa, *Phys. Rev. B* 48, 18340 (1993).

49. A. Scholze, A. Schenk, and W. Fichtner, *IEEE Trans. Electron Devices* 47, 1811 (2000).

50. P. Matagne and J. P. Leburton, *Phys. Rev. B* 65, 155311 (2002).

51. M. Korkusinski and P. Hawrylak, *Phys. Rev. B* 63, 195311 (2001).

52. M. Macucci, K. Hess, and G. J. Iafrate, *Phys. Rev. B* 48, 17354 (1993).

53. J. H. Oaknin, J. J. Palacios, L. Brey, and C. Tejedor, *Phys. Rev. B* 49, 5718 (1994).

54. W. Kohn and L. J. Sham, *Phys. Rev.* 140, A1133 (1965).

55. J. P. Perdew and Alex Zunger, *Phys. Rev. B* 23, 5048 (1981).

56. B. Tanatar and D. M. Ceperley, *Phys. Rev. B* 39, 5005 (1989).

57. C. A. Stafford and S. Das Sarma, *Phys. Rev. Lett.* 72, 3590 (1994); R. Kotlyar and S. Das Sarma, *Phys. Rev. B* 56, 13235 (1997); R. Kotlyar and S. Das Sarma, *Phys. Rev. B* 55, R10205, (1997).

58. G. Klimeck, G. Chen, and S. Datta, *Phys. Rev. B* 50, 2316 (1994).

59. R. MacWeeny, "Methods of Molecular Quantum Mechanics." Academic Press, London, 1989.

60. M. Governale, M. Macucci, G. Iannaccone, C. Ungarelli, and J. Martorell, *J. Appl. Phys.* 85, 2962 (1999).

61. J. C. Slater, "Quantum Theory of Matter." McGraw-Hill, New York, 1968.

62. I. O. Kulik and R. I. Shekhter, *Sov. J. Low Temp. Phys.* 3, 532 (1977).

63. K. K. Likharev, N. S. Bakhvalov, G. S. Kazacha, and S. I. Serdyukova, *IEEE Trans. Magn.* 25, 1436 (1989).

64. M. Macucci, M. Gattobigio, and G. Iannaccone, *J. Appl. Phys.* 90, 6428 (2001).

65. L. Bonci, M. Gattobigio, G. Iannaccone, and M. Macucci, *J. Appl. Phys.* 92, 3169 (2002).

66. D. V. Averin and K. K. Likharev, in "Mesoscopic Phenomena in Solids" (B. L. Altshuler, P. A. Lee, and R. A. Webb, Eds.). Elsevier, Amsterdam, 1991.

67. W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, "Numerical Recipes in Fortran," 2nd edn. Cambridge University Press, New York, 1992.

68. D. E. Knuth, "Seminumerical Algorithm. The Art of Computer Programming." 2nd edn. Addison-Wsley, Reading, MA, 1981.

69. C. W. Gardiner, "Handbook of Stochastic Methods. Springer Series in Synergetics." 2nd edn. Springer, Erlin, 1985.

70. P. R. Graves-Morris, "Padé Approximation and Its Applications." Lectures Notes in Mathematics, Vol 765. Springer, Berlin, 1979.

71. Y. Saad, *SIAM J. Numerical Anal.* 29, 209 (1992).

72. H. Pothier, Ph.D. Thesis University of Paris 6, 1991.

73. H. D. Jensen and J. M. Martinis, *Phys. Rev. B* 46, 13407 (1992).

74. D. V. Averin and A. A. Odinstov, *Phys. Lett. A* 140, 251 (1989).

75. L. I. Glazman and K. A. Matveev, *Sov. Phys.-JETP* 71, 1031 (1990).

# CHAPTER 14

# Wigner Function–Based Device Modeling

## Hans Kosina, Mihail Nedjalkov

*Institute for Microelectronics, TU Vienna, Vienna, Austria*

## CONTENTS

## 1. INTRODUCTION

Modeling of electronic transport in mesoscopic systems requires a theory that describes open, quantum-statistical systems driven far from thermodynamic equilibrium. Several formulations of quantum transport have been employed practically, such as those based on the density matrix, nonequilibrium Green's functions, and the Wigner function.

A quantum-mechanical phase-space distribution was introduced by Eugene Wigner n 1932 [1]. The purpose was the formulation of a quantum correction for the thermodynamic equilibrium of a many-body system by means of a quasiprobability function. In more recent times, the definition of the Wigner function has been generalized as a Fourier transform of a many-body Green's function [2].

The Wigner function is a real-valued but not necessarily positive definite quasidistribution and represents a quantum generalization of Boltzmann's N-particle distribution. The Wigner function formalism is attractive as it allows the expression of quantum dynamics n a phase-space formulation, directly comparable with the classical analogue. A phase-space approach may appear more intuitive compared with the more abstract density matrix and Green's function approaches. The method of quasidistributions has proved especially useful in providing reductions to classical physics and kinetic regimes under suitable conditions.

To discuss the physical interpretation of a quasidistribution, let us consider the simple case of a one-particle distribution. Starting with the classical case, the distribution $f_{cl}(\mathbf{p}, \mathbf{r}, )$ is proportional to the probability density of finding a particle of momentum $\mathbf{p}$ and position $\mathbf{r}$ in the phase-space volume $d^3p\,d^3r$. This is a purely classical interpretation, directly conflicting with the uncertainty principle. The quantum mechanical quasidistribution $f_w(\mathbf{p}, \mathbf{r}, t)$, however, is not positive definite and has to be interpreted as a joint density of $\mathbf{p}$ and $\mathbf{r}$ [3]. Only the marginal distributions are positive definite, that is, integrating $f_w(\mathbf{p}, \mathbf{r}, t)$ over momentum space gives the probability density in $\mathbf{r}$-space, and vice versa.

An excellent review of quantum-mechanical phase-space distributions in scattering theory has been given by Carruthers and Zachariason [4]. This work deals with potential scattering, the two-body problem, and the N-body problem. A coupled hierarchy for reduced distribution functions and its truncation to the Boltzmann-Vlasov equation is presented. Tatarskii [5] concentrates on quantum-mechanical systems in a pure state and investigates the representation of quantum mechanics by phase-space distributions. He points out that not every function that solves the Wigner equation describes a pure state. Therefore, initial conditions for the Wigner equation have to be subjected to a supplementary restriction. Today, phase-space quantization is considered to be a third autonomous and logically complete formulation of quantum mechanics beyond the conventional ones based on operators in Hilbert space or path integrals [5, 6]. This formulation is free of operators and wave functions. Observables and matrix elements are computed through phase-space integrals of c-number functions weighted by a Wigner function.

Important quantum mechanical properties of electronic transport in semiconductor structures are often those associated not with the degeneracy of the Fermi system but rather with quantum interference effects [7]. A wide variety of electronic quantum transport problems of interest are essentially one-particle in nature. In such cases, a full many-body description of the problem is not necessary, and a description of electronic transport that makes use of the one-particle approximation can be used from the very outset. However, even when the electron–electron interaction effects are of interest, certain approximations do exist, allowing their description on a one-particle level [7]. Therefore, we shall consider in the following only electronic systems with one-particle degrees of freedom.

## 1.1. History and State of the Art Review

Reports on finite-difference solutions of the one-particle Wigner equation for device applications are due to Ravaioli [8], Kluksdahl [9], and coworkers, and date back to the mid 1980s. Frensley [10–12] was the first who introduced boundary conditions on the Wigner function to model open quantum systems. Later, self-consistency was added to the Wigner equation solvers [13, 14]. Main and Haddad included a reduced Boltzmann scattering operator in transient Wigner function–based simulations [15]. Research on finite-difference solution methods for the Wigner equation culminated in 1990 when the review articles of Frensley [16] and Buot and Jensen [17] appeared.

The 1990s have seen further extensions and applications of the finite-difference Wigner function method. High-frequency operation of resonant tunneling diodes has been studied by Jensen and Buot [18, 19], and the transient response by Gullapalli [20] and Biegel [21], and later by [22]. A new finite-difference discretization scheme has been proposed in [23].

In 2002, implementations of Monte Carlo methods for solving the Wigner device equation were reported [24, 25]. Although with the finite-difference method, scattering was restricted to the relaxation time approximation and the momentum space to one dimension, the Monte Carlo method allows scattering processes to be included on a more detailed level, assuming a three-dimensional momentum-space [26, 27]. Issues such as choosing proper up-winding schemes, restrictions on matrix size and momentum space resolution are largely relaxed or do not exist when using the Monte Carlo method. Construction of new Monte Carlo algorithms is complicated by the fact that the kernel of the integral equation to solve is not positive semidefinite. As a consequence, the commonly applied Markov chain Monte Carlo method shows a variance exponentially increasing with time, prohibiting its application to realistic structures or larger evolution times [25, 28, 29]. Because of this so-called negative sign problem, the concept of Wigner paths alone [30, 31] is not sufficient to construct a stable Monte Carlo algorithm. Instead, additional measures have to be introduced that prevent a runaway of the particle weights and hence of the variance [26, 32]. Note that in [26], the statistical weights are termed affinities.

Large basic research efforts on the Monte Carlo modeling of electron–phonon interaction based on the Wigner function formalism have been reported in [28, 31, 33–35].

The effect of a spatially varying effective mass in Wigner device simulations has been demonstrated in [36] and [37]. A nonparabolic version of the Wigner equation has been derived by Bufler [38]. Multiband models have been reported in [39–41].

A Wigner equation including a magnetic field has been solved in [42]. The gauge-invariant formulation of the Wigner equation has been given by Levinson [43], and a discussion can be found in various works [4, 44–47]. Two-time and frequency-dependent Wigner functions are considered in [2, 47–49].

Finally, we note that the Wigner function formalism is often used to derive reduced transport models, such as the quantum hydrodynamic model [50, 51–53], or to find quantum corrections to classical models, such as the ensemble Monte Carlo method [54] or the spherical harmonics expansion method [55, 56].

## 2. THE WIGNER FUNCTION FORMALISM

In the Schrödinger picture, a physical system is quantum-mechanically described by a state vector $|\Psi(t)\rangle$ as function of time $t$. Often, the precise quantum-mechanical state of a system is not known, but rather some statistical information about the probabilities for the system being in one of a set of states. Suppose that there is a set of ortho-normal states $\{|\Psi_1\rangle, |\Psi_2\rangle, \ldots\}$, and that the probabilities that the system is in one of these states are $\{p_1, p_2, \ldots\}$. Then, the expectation value of operator $\hat{A}$ associated with the observable $A$ is given by

$$\langle A \rangle = \sum_i p_i \langle \Psi_i | \hat{A} | \Psi_i \rangle \tag{1}$$

which is a quantum and statistical average. Introducing the density operator $\hat{\rho}$ as

$$\hat{\rho} = \sum_i p_i |\Psi_i\rangle\langle\Psi_i| \tag{2}$$

the expectation value becomes

$$\langle A \rangle = \mathrm{Tr}(\hat{\rho}\hat{A}) = \mathrm{Tr}(\hat{A}\hat{\rho}) \tag{3}$$

Formulations (1) and (3) require the operator $\hat{A}$ to be self-adjoint. Equation (3) can be easily verified by expressing the trace of some operator $\hat{X}$ in the basis $\{|\Psi_i\rangle\}$.

$$\mathrm{Tr}\langle\hat{X}\rangle = \sum_i \langle\Psi_i|\hat{X}|\Psi_i\rangle \tag{4}$$

The fact that the probabilities sum up to unity, $\sum_i p_i = 1$, is expressed by the fact that the trace of the density operator is also unity, $\mathrm{Tr}(\hat{\rho}) = 1$. If the system is in a pure state $|\Psi_i\rangle$ it

holds $p_i = 1$ and $p_j = 0 \; \forall \; j \neq i$, and the density operator is idem-potent, $\hat{\rho}^2 = \hat{\rho}$. Otherwise, the system is in a mixed state, and $\hat{\rho}$ does not obey the idem-potency condition. From the Schrödinger equation for the state vector and the definition of $\hat{\rho}$, we immediately obtain the Liouville-von Neumann equation for the evolution of the density operator.

$$ i\hbar \frac{\partial \hat{\rho}}{\partial t} = [\hat{H}, \hat{\rho}] \tag{5} $$

Introducing the one-particle approximation [7] implies that the electron system is modeled as consisting of many, noninteracting electrons. In the next step, one chooses the coordinate representation, where the set of basis vectors is given by the electron position eigenstates $|r\rangle$. The eigenstates of the system are then represented by the wavefunctions $\Psi_i(\mathbf{r}, t) = \langle \mathbf{r} | \Psi_i(t) \rangle$, and the density operator by the density matrix $\rho(\mathbf{r}_1, \mathbf{r}_2, t)$.

$$ \rho(\mathbf{r}_1, \mathbf{r}_2, t) = \langle \mathbf{r}_1 | \hat{\rho}(t) | \mathbf{r}_2 \rangle = \sum_i p_i \Psi_i(\mathbf{r}_1, t) \Psi_i^*(\mathbf{r}_2, t) \tag{6} $$

The Liouville-von Neumann equation in coordinate representation is found as

$$ \frac{\partial \rho(\mathbf{r}_1, \mathbf{r}_2, t)}{\partial t} = (H_{\mathbf{r}_1} - H_{\mathbf{r}_2}) \rho(\mathbf{r}_1, \mathbf{r}_2, t) \tag{7} $$

## 2.1. The Wigner Function

The Wigner function is obtained from the density matrix by means of the Wigner-Weyl transformation. This transformation consists of a change of independent coordinates to diagonal and cross-diagonal coordinates

$$ \mathbf{r} = \frac{1}{2}(\mathbf{r}_1 + \mathbf{r}_2), \qquad \mathbf{s} = \mathbf{r}_1 - \mathbf{r}_2 \tag{8} $$

followed by a Fourier transformation with respect to $\mathbf{s}$ [16]. The variables $\mathbf{r}_1$ and $\mathbf{r}_2$ may be expressed in terms of the new ones.

$$ \mathbf{r}_1 = \mathbf{r} + \frac{\mathbf{s}}{2}, \qquad \mathbf{r}_2 = \mathbf{r} - \frac{\mathbf{s}}{2} \tag{9} $$

Then, the elementary definition of the Wigner distribution is given by the following transformation of the density matrix.

$$ f_w(\mathbf{k}, \mathbf{r}, t) = \int \rho \left( \mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{10} $$

The Wigner function (10) is real-valued, but not positive semidefinite. In terms of the wave functions, the definition (10) becomes

$$ f_w(\mathbf{k}, \mathbf{r}, t) = \sum_i p_i \int \Psi_i \left( \mathbf{r} + \frac{\mathbf{s}}{2}, t \right) \Psi_i^* \left( \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{11} $$

The normalization of the Wigner function results from the normalization of the wave functions.

$$ \frac{1}{(2\pi)^3} \int d\mathbf{r} \int d\mathbf{k} \, f_w(\mathbf{k}, \mathbf{r}, t) = 1 \tag{12} $$

Here, the $\mathbf{k}$-integration can be performed first, giving $\int e^{-i\mathbf{k}\cdot\mathbf{s}} d\mathbf{k} = (2\pi)^3 \delta(\mathbf{s})$. The normalization (12) ensures that the quantity $N f_w$, where $N$ is the number of electrons in the system, will approach the classical distribution function $f_{cl}$ in the classical limit [35].

Sometimes it is convenient to use the inverse Fourier transform of (10).

$$ \rho \left( \mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) e^{i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{k} \tag{13} $$

Changing variables gives a transformation that inverts the Wigner-Weyl transformation.

$$\rho(\mathbf{r}_1, \mathbf{r}_2, t) = \frac{1}{(2\pi)^3} \int f_w\left(\mathbf{k}, \frac{\mathbf{r}_1 + \mathbf{r}_2}{2}, t\right) e^{i\mathbf{k}\cdot(\mathbf{r}_1 - \mathbf{r}_2)} d\mathbf{k} \tag{14}$$

An important feature of the phase-space approach is the possibility of expressing quantum-mechanical expectation values in the same way as it is done in classical statistical mechanics, employing integration over the phase-space. The expectation values of operators of the form $A(\hat{\mathbf{r}})$ and $B(\hat{\mathbf{k}})$, where $\hat{\mathbf{k}} = \hat{\mathbf{p}}/\hbar$, are given as follows.

$$\langle A(\hat{\mathbf{r}}) \rangle = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) A(\mathbf{r}) \, d\mathbf{k} \, d\mathbf{r} = \sum_i p_i \int A(\mathbf{r}) \, |\Psi_i(\mathbf{r}, t)|^2 \, d\mathbf{r} \tag{15}$$

$$\langle B(\hat{\mathbf{k}}) \rangle = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) B(\mathbf{k}) \, d\mathbf{k} \, d\mathbf{r} = \sum_i p_i \int B(\mathbf{k}) \, |\Phi_i(\mathbf{k}, t)|^2 \, d\mathbf{k} \tag{16}$$

If the classical observable $C(\mathbf{k}, \mathbf{r})$ is a function of both momentum and position, the definition of a corresponding Hermitian operator $\widehat{C}$ is not unique. In this case, the Weyl quantization can be applied. Thereby, the function $C$ is expressed through its Fourier transform $c$.

$$C(\mathbf{k}, \mathbf{r}) = \int c(\mathbf{a}, \mathbf{b}) e^{i(\mathbf{k}\cdot\mathbf{a} + \mathbf{r}\cdot\mathbf{b})} \, d\mathbf{a} \, d\mathbf{b} \tag{17}$$

The operator $\widehat{C}$ is defined by the following rule of correspondence.

$$\widehat{C} = \int c(\mathbf{a}, \mathbf{b}) e^{i(\hat{\mathbf{k}}\cdot\mathbf{a} + \hat{\mathbf{r}}\cdot\mathbf{b})} \, d\mathbf{a} \, d\mathbf{b} \tag{18}$$

Then, the expectation value of $\widehat{C}$ is given by the phase-space integral.

$$\mathrm{Tr}\left(\widehat{C}\hat{\rho}\right) = \int C(\mathbf{k}, \mathbf{r}) f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{k} \, d\mathbf{r} \tag{19}$$

To proceed with (18), one may employ the Baker-Campbell-Hausdorff formula,

$$e^{\hat{A} + \hat{B}} = e^{\hat{A}} e^{\hat{B}} e^{-\frac{[\hat{A}, \hat{B}]}{2}} \tag{20}$$

which is generally valid when $[\hat{A}, [\hat{A}, \hat{B}]] = [\hat{B}, [\hat{A}, \hat{B}]] = 0$, or in particular when $[\hat{A}, \hat{B}]$ is a c-number.

## 2.2. Marginal Distributions

The Wigner function (10) can assume negative values. Only the marginal distributions of $f_w(\mathbf{k}, \mathbf{r}, t)$ are positive semidefinite and have the meaning of probability distributions in real space and momentum space, respectively.

$$n(\mathbf{r}) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{k} = \sum_i p_i \, |\Psi_i(\mathbf{r}, t)|^2 \tag{21}$$

$$p(\mathbf{k}) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{r} = \sum_i p_i \, |\Phi_i(\mathbf{k}, t)|^2 \tag{22}$$

Here, $\Phi_i(\mathbf{k},\,)$ denotes the momentum representation of the state vector $|\Psi_i\rangle$. The integration in (22) can easily be carried out after changing variables, using (8).

$$\int d\mathbf{r} \int d\mathbf{s} \, \Psi_i\left(\mathbf{r} + \frac{\mathbf{s}}{2}, t\right) \Psi_i^*\left(\mathbf{r} - \frac{\mathbf{s}}{2}, t\right) e^{-i\mathbf{k}\cdot\mathbf{s}}$$

$$= \int d\mathbf{r}_1 \int d\mathbf{r}_2 \Psi_i(\mathbf{r}_1, t) \Psi_i^*(\mathbf{r}_2, t) e^{-i\mathbf{k}\cdot(\mathbf{r}_1 - \mathbf{r}_2)} = (2\pi)^3 |\Phi_i(\mathbf{k}, t)|^2 \tag{23}$$

The marginal distributions (21) and (22) can also be expressed as the diagonal elements of the density matrix.

$$\frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{k} = \langle \mathbf{r}|\hat{\rho}|\mathbf{r}\rangle = \rho(\mathbf{r}, \mathbf{r}) \tag{24}$$

$$\frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{r} = \langle \mathbf{k}|\hat{\rho}|\mathbf{k}\rangle = \sigma(\mathbf{k}, \mathbf{k}) \tag{25}$$

Here, $|\mathbf{k}\rangle$ denotes the electron momentum eigenstate with eigenvalue $\hbar\mathbf{k}$ and $\sigma$ the density matrix in momentum representation. Note that the latter can be used for a dual definition of the Wigner function [28, 57].

$$f_w(\mathbf{k}, \mathbf{r}, t) = \int \sigma\left(\mathbf{k} + \frac{\mathbf{l}}{2}, \mathbf{k} - \frac{\mathbf{l}}{2}, t\right) e^{i\mathbf{r}\cdot\mathbf{l}} \, d\mathbf{l} \tag{26}$$

This definition follows, for example, from (11), when the $\Psi_i$ are replaced by

$$\Psi_i(\mathbf{r}, t) = (2\pi)^{-3/2} \int \Phi_i(\mathbf{k}', t) e^{i\mathbf{k}'\cdot\mathbf{r}} \, d\mathbf{k}' \tag{27}$$

Other marginal distributions than the elementary ones, (21) and (22), have to be constructed with care. Only Hermitian operators give real marginal distributions. For the current density, this operator would be $(\hat{\mathbf{k}}\hat{\rho} + \hat{\rho}\hat{\mathbf{k}})/2$. Expressing $\hat{\rho}$ in terms of the wave functions, we get the elementary current definition from wave mechanics.

$$\mathbf{j}(\mathbf{r}) = \frac{\hbar}{2m^*} \langle \mathbf{r}|\hat{\mathbf{k}}\hat{\rho} + \hat{\rho}\hat{\mathbf{k}}|\mathbf{r}\rangle$$

$$= \frac{\hbar}{2m^*} \sum_i p_i\left(\langle \mathbf{r}|\hat{\mathbf{k}}|\Psi_i\rangle\langle \Psi_i|\mathbf{r}\rangle + \langle \mathbf{r}|\Psi_i\rangle\langle \Psi_i|\hat{\mathbf{k}}|\mathbf{r}\rangle\right)$$

$$= \frac{\hbar}{2im^*} \sum_i p_i\left[\Psi_i^*(\mathbf{r})\nabla\Psi_i(\mathbf{r}) - \Psi_i(\mathbf{r})\nabla\Psi_i^*(\mathbf{r})\right] \tag{28}$$

Choosing the momentum representation of $\hat{\rho}$, we get the current density expressed in terms of the Wigner function.

$$\mathbf{j}(\mathbf{r}) = \frac{\hbar}{2m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 \left(\langle \mathbf{r}|\hat{\mathbf{k}}|\mathbf{k}_1\rangle\langle \mathbf{k}_1|\hat{\rho}|\mathbf{k}_2\rangle\langle \mathbf{k}_2|\mathbf{r}\rangle + \langle \mathbf{r}|\mathbf{k}_1\rangle\langle \mathbf{k}_1|\hat{\rho}|\mathbf{k}_2\rangle\langle \mathbf{k}_2|\hat{\mathbf{k}}|\mathbf{r}\rangle\right)$$

$$= \frac{\hbar}{2m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 \, \sigma(\mathbf{k}_1, \mathbf{k}_2)(\mathbf{k}_1 + \mathbf{k}_2) e^{i(\mathbf{k}_1 - \mathbf{k}_2)\cdot\mathbf{r}} = \frac{1}{(2\pi)^3} \int \frac{\hbar}{m^*}\mathbf{k} \, f_w(\mathbf{k}, \mathbf{r}, t) \, d\mathbf{k} \tag{29}$$

Here, the Wigner function has been introduced using (26). The current density is given by the first-order moment of the Wigner function, in full analogy with the classical phase space definition.

For the definition of the energy density we discuss several options. Starting from the trace operation for the statistical average, one would consider the symmetrized operator $(\hat{\mathbf{k}}^2\hat{\rho} + \hat{\rho}\hat{\mathbf{k}}^2)/2$ and derive the marginal distribution.

$$w_1(\mathbf{r}) = \frac{\hbar^2}{4m^*}(\langle \mathbf{r}|\hat{\mathbf{k}}^2\hat{\rho}|\mathbf{r}\rangle + \langle \mathbf{r}|\hat{\rho}\hat{\mathbf{k}}^2|\mathbf{r}\rangle)$$

$$= -\frac{\hbar^2}{4m^*} \sum_i p_i\left[\Psi_i^*(\mathbf{r})\nabla^2\Psi_i(\mathbf{r}) + \Psi_i(\mathbf{r})\nabla^2\Psi_i^*(\mathbf{r})\right]$$

$$= \sum_i p_i[E_i - V(\mathbf{r})]|\Psi_i(\mathbf{r})|^2 \tag{30}$$

The last expression in (30) is obtained with the help of the stationary Schrödinger equation. Apparently, $w_1$ describes the kinetic energy density, as the potential energy term $V(\mathbf{r})n(\mathbf{r})$ is subtracted from the total energy term. This energy density can become negative in tunneling

regions, where for one or more states $E_i < V(\mathbf{r})$ holds. In a derivation similar to (29), one finds the Wigner representation of $w_1$.

$$w_1(\mathbf{r}) = \frac{\hbar^2}{4m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2\, \sigma(\mathbf{k}_1, \mathbf{k}_2)\,(\mathbf{k}_1^2 + \mathbf{k}_2^2)\, e^{i(\mathbf{k}_1 - \mathbf{k}_2)\mathbf{r}}$$

$$= \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} \left( |\mathbf{k}|^2 - \frac{1}{4}\nabla_r^2 \right) f_w(\mathbf{k}, \mathbf{r}, t)\, d\mathbf{k} \tag{31}$$

To ensure positiveness of the energy density, in [58] the Hermitian operator $\hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}}$ is considered. Its marginal distribution can be shown to be positive semidefinite.

$$w_2(\mathbf{r}) = \frac{\hbar^2}{2m^*} \langle \mathbf{r}|\hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}}|\mathbf{r}\rangle = \frac{\hbar^2}{2m^*} \sum_i p_i |\langle \mathbf{r}|\hat{\mathbf{k}}|\Psi_i\rangle|^2 \tag{32}$$

$$= \frac{\hbar^2}{2m^*} \sum_i p_i |\nabla \Psi_i(\mathbf{r})|^2 \ge 0 \tag{33}$$

The Wigner representation of $w_2$ is obtained as

$$w_2(\mathbf{r}) = \frac{\hbar^2}{4m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2\, \sigma(\mathbf{k}_1, \mathbf{k}_2)\,(\mathbf{k}_1^2 - \mathbf{k}_2^2)\, e^{i(\mathbf{k}_1 - \mathbf{k}_2)\mathbf{r}}$$

$$= \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} \left( |\mathbf{k}|^2 + \frac{1}{4}\nabla_r^2 \right) f_w(\mathbf{k}, \mathbf{r}, t)\, d\mathbf{k} \tag{34}$$

Conditions for obtaining non-negative marginal distributions are theoretically discussed in [59]. The Weyl correspondence (18) gives the definition of the energy density as the second-order moment of the Wigner function.

$$w_3(\mathbf{r}) = \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} |\mathbf{k}|^2 f_w(\mathbf{k}, \mathbf{r}, t)\, d\mathbf{k} \tag{35}$$

It can be seen that (35) is just the arithmetic mean of (31) and (34), $w_3 = (w_1 + w_2)/2$. Therefore, (35) represents the marginal distribution of the symmetrized operator $(\hat{\mathbf{k}}^2\hat{\rho} + 2\hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}} + \hat{\rho}\hat{\mathbf{k}}^2)/4$.

All three definitions of the energy density give the same statistical average $\langle \hat{\epsilon} \rangle = \mathrm{Tr}[\epsilon(\hat{\mathbf{k}})\hat{\rho}]$. The differences among the definitions are in the $\nabla^2$ term, which vanishes after the $\mathbf{r}$-integration. However, only the density $w_1$ seems to have a clear physical interpretation as the kinetic energy density.

## 2.3. The Wigner Equation

In this section, we consider a system consisting of one electron interacting with a potential distribution $V_{tot}(\mathbf{r})$. This potential is assumed to be a superposition of some potential $V(\mathbf{r})$ and a uniform electric field: $V_{tot}(\mathbf{r}) = V(\mathbf{r}) - \hbar\mathbf{F}\cdot\mathbf{r}$, with $\hbar\mathbf{F} = -e\mathbf{E}$. Although the existence of a field term is not physically motivated at this point, it is introduced here to demonstrate its treatment in the Wigner function formalism. The potential $V(\mathbf{r})$ comprises the electrostatic potential and the band-edge profile of the semiconductor. A uniform effective mass $m^*$ is assumed. In the usual coordinate representation, the Hamiltonian of the system is then given by

$$H = H_0 + V(\mathbf{r}) - \hbar\mathbf{F}\cdot\mathbf{r} \tag{36}$$

with

$$H_0 = -\frac{\hbar^2}{2m^*}\nabla_r^2 \tag{37}$$

The electron phonon interaction neglected here will be discussed in detail in Section 3. The evolution equation for the Wigner function is found by taking the time derivative of the defining Eq. (10) and substituting the Liouville-von Neumann Eq. (7) on the right-hand side.

$$\frac{\partial}{\partial t} f_w(\mathbf{k}, \mathbf{r}, t) = \frac{1}{i\hbar} \int (H_{r_1} - H_{r_2}) \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t\right) e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{38}$$

In the following, the three parts of the Hamiltonian (36) will be separately transformed. Unlike in Section 2.2, were calculations where done in momentum representation, we choose below the configuration representation to carry out the transformations [33].

The free-electron Hamiltonian is given by $H_0$. To calculate the Wigner transform of $H_0$, we have to transform the gradients first. Differentiating the density matrix with respect to the new variables $\mathbf{r}$ and $\mathbf{s}$

$$\nabla_r \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t\right) = \nabla_{r_1} \rho + \nabla_{r_2} \rho, \quad \nabla_s \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t\right) = \frac{1}{2}\nabla_{r_1} \rho - \frac{1}{2}\nabla_{r_2} \rho \tag{39}$$

gives the relations

$$\nabla_{r_1} + \nabla_{r_2} = \nabla_r, \qquad \nabla_{r_1} - \nabla_{r_2} = 2\nabla_s, \qquad \nabla_{r_1}^2 - \nabla_{r_2}^2 = 2\nabla_r \cdot \nabla_s \tag{40}$$

Now the free-electron term transforms to a diffusion term. For the sake of brevity, we write $\rho_{r,s} = \rho(\mathbf{r} + \mathbf{s}/2, \mathbf{r} - \mathbf{s}/2, t)$ in the following.

$$\frac{1}{i\hbar} \int -\frac{\hbar^2}{2m^*}\left(\nabla_{r_1}^2 - \nabla_{r_2}^2\right)\rho_{r,s} e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} = -\frac{\hbar}{im^*}\nabla_r \cdot \int (\nabla_s \rho_{r,s}) e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{41}$$

$$= -\frac{\hbar\mathbf{k}}{m^*} \cdot \nabla_r \int \rho_{r,s} e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{42}$$

$$= -\frac{\hbar\mathbf{k}}{m^*} \cdot \nabla_r f_w(\mathbf{k}, \mathbf{r}, t) \tag{43}$$

Next, we transform the potential term $V(\mathbf{r})$.

$$\frac{1}{i\hbar} \int \left[V\left(\mathbf{r} + \frac{\mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\mathbf{s}}{2}\right)\right] \rho_{r,s} e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) \, d\mathbf{k}' \tag{44}$$

This transformation is readily found by replacing $\rho_{r,s}$ on the left-hand side by the inverse Fourier transformation (13). The remaining integral over $\mathbf{s}$ is denoted by $V_w$ and referred to as the Wigner potential.

$$V_w(\mathbf{q}, \mathbf{r}) = \frac{1}{(2\pi)^3 i\hbar} \int \left[V\left(\mathbf{r} + \frac{\mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\mathbf{s}}{2}\right)\right] e^{-i\mathbf{q}\cdot\mathbf{s}} \, d\mathbf{s} \tag{45}$$

Using the simple relation $-(\mathbf{F}\cdot\mathbf{r}_1 - \mathbf{F}\cdot\mathbf{r}_2) = -\mathbf{F}\cdot\mathbf{s}$, the constant-field term transforms as

$$\frac{1}{i\hbar} \int (-\hbar\mathbf{F}\cdot\mathbf{s})\rho_{r,s} e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} = -\frac{1}{\hbar}\mathbf{F} \cdot \nabla_k f_w(\mathbf{k}, \mathbf{r}, t) \tag{46}$$

Collecting the above results gives the Wigner equation for the system Hamiltonian (36).

$$\left(\frac{\partial}{\partial t} + \frac{\hbar\mathbf{k}}{m^*} \cdot \nabla_r + \mathbf{F} \cdot \nabla_k\right) f_w(\mathbf{k}, \mathbf{r}, t) = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) \, d\mathbf{k}' \tag{47}$$

The terms are arranged so to form the classical Liouville operator on the left-hand side. The interaction of the electron with the potential distribution $V(\mathbf{r})$ is described by the potential operator on the right-hand side. As can be seen, the Wigner function in $\mathbf{k}$ and $\mathbf{r}$ depends in a nonlocal manner on the Wigner function in all other momentum points $\mathbf{k}'$ and through $V_w$ also on the potential at all other locations $\mathbf{r} \pm \mathbf{s}/2$.

## 3.. ELECTRON–PHONON INTERACTION

The Wigner equation has frequently been solved using the finite-difference method [16, 60], assuming the phenomenological relaxation time approximation for dissipative transport. Recently developed Monte Carlo methods allowed phonon scattering to be included semi-classically in quantum device simulations [24, 27]. Use of a Boltzmann scattering operator acting on the Wigner distribution was originally suggested by Frensley [16]. In this section, the Wigner equation with a Boltzmann scattering operator is rigorously derived, using a many-phonon single-electron Wigner function formalism as the starting point.

### 3..1. The System Hamiltonian

The Hamiltonian (36) is now extended to describe a system consisting of one electron interacting with a many-phonon system and a given potential distribution.

$$H = H_0 + V(\mathbf{r}) - \mathbf{F} \cdot \mathbf{r} + H_p + H_{ep} \tag{48}$$

The additional components of this Hamiltonian are given by [34]

$$H_p = \sum_q \hbar \omega_q \, b_q^\dagger \, b_q \tag{49}$$

$$H_{ep} = i\hbar \sum_q \mathcal{F}(\mathbf{q}) \, (b_q \, e^{i\mathbf{q}\cdot\mathbf{r}} - b_q^\dagger \, e^{-i\mathbf{q}\cdot\mathbf{r}}) \tag{50}$$

Here, $H_p$ is the Hamiltonian of the free phonon-system, $H_{ep}$ the electron–phonon interaction Hamiltonian, $b_q$ and $b_q^\dagger$ denote the annihilation and creation operators for a phonon with momentum $\hbar\mathbf{q}$ and energy $\hbar\omega_q$, and $\hbar\mathcal{F}(\mathbf{q})$ is the interaction matrix element.

We introduce a set of basis vectors $|\mathbf{r}, \{n\}\rangle$ in the occupation number representation. A set of occupation numbers is defined as $\{n\} = n_{q_1}, n_{q_2}, \ldots n_{q}, \ldots$, where $n_q$ is the number of phonons with momentum $\mathbf{q}$. The Wigner-Weyl transformation of the density matrix $\rho(\mathbf{r}_1, \{n\}, \mathbf{r}_2, \{m\})$ gives the generalized Wigner function $f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$ [28, 33].

$$f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) = \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \left| \hat{\rho}(t) \right| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{51}$$

Note that only the electron coordinates are transformed, such that $f_g$ is a Wigner function on the electron phase-space, but still is the density matrix for the phonon system.

The evolution of the generalized Wigner function is found by taking the time derivative of (51) and using the Liouville-von Neumann equation for the evolution of the density matrix.

$$\frac{\partial}{\partial t} f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) = \frac{1}{i\hbar} \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \left| [\hat{H}, \hat{\rho}(t)] \right| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k}\cdot\mathbf{s}} \, d\mathbf{s} \tag{52}$$

To continue, one may express the density matrix in the state vectors of the system.

$$\rho(\mathbf{r}_1, \{n\}, \mathbf{r}_2, \{m\}, t) = \sum_i p_i \, \Psi_i(\mathbf{r}_1, \{n\}, t) \, \Psi_i^*(\mathbf{r}_2, \{m\}, t) \tag{53}$$

The creation and annihilation operators, and the occupation number operator $b_q^\dagger b_q$ satisfy the following well-known eigenvalue equations.

$$b_q^\dagger \Psi(\mathbf{r}, \{n\}, t) = \sqrt{n_q + 1} \, \Psi(\mathbf{r}, \{n_{q_1}, n_{q_2}, \ldots n_q + 1, \ldots\}, t)$$

$$b_q \Psi(\mathbf{r}, \{n\}, t) = \sqrt{n_q} \, \Psi(\mathbf{r}, \{n_{q_1}, n_{q_2}, \ldots n_q - 1, \ldots\}, t) \tag{54}$$

$$b_q^\dagger b_q \Psi(\mathbf{r}, \{n\}, t) = n_q \, \Psi(\mathbf{r}, \{n_{q_1}, n_{q_2}, \ldots n_q, \ldots\}, t)$$

With the help of these equations and the representation (53), the transformation of the free-phonon Hamiltonian is readily found.

$$\frac{1}{i\hbar} \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \left| [\widehat{H}_{\mathrm{p}}, \hat{\rho}(t)] \right| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k}\cdot\mathbf{s}} \, ds$$

$$= \frac{1}{i\hbar} (\epsilon(\{n\}) - \epsilon(\{m\})) f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$$

The energy of the phonon state $|\{n\}\rangle$ is denoted by $\epsilon(\{n\})$.

$$\epsilon(\{n\}) = \sum_{\mathbf{q}} n_{\mathbf{q}} \hbar \omega_{\mathbf{q}} \tag{55}$$

The electron-phonon interaction Hamiltonian is transformed following the same lines [33]. Combining the two terms of the Hamiltonian (50) and the two terms of the commutator in (52) results in four terms related to the electron–phonon interaction. In the equation for the generalized Wigner function shown below, these four terms appear under the sum.

$$\left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r + \frac{1}{\hbar} \mathbf{F} \cdot \nabla_{\mathbf{k}}\right) f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$$

$$= \int V_{\mathrm{w}}(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_{\mathrm{g}}(\mathbf{k}', \mathbf{r}, \{n\}, \{m\}, t) \, d\mathbf{k}' + \frac{1}{i\hbar}(\epsilon(\{n\}) - \epsilon(\{m\})) f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$$

$$+ \sum_{\mathbf{q}} \mathcal{F}(\mathbf{q}) e^{i\mathbf{q}\cdot\mathbf{r}} \sqrt{n_{\mathbf{q}} + 1} f_{\mathrm{g}}\left(\mathbf{k} - \frac{\mathbf{q}}{2}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \ldots n_{\mathbf{q}} + 1, \ldots\}, \{m\}, t\right)$$

$$- e^{-i\mathbf{q}\cdot\mathbf{r}} \sqrt{n_{\mathbf{q}}} f_{\mathrm{g}}\left(\mathbf{k} + \frac{\mathbf{q}}{2}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \ldots n_{\mathbf{q}} - 1, \ldots\}, \{m\}, t\right)$$

$$- e^{i\mathbf{q}\cdot\mathbf{r}} \sqrt{m_{\mathbf{q}}} f_{\mathrm{g}}\left(\mathbf{k} + \frac{\mathbf{q}}{2}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \ldots m_{\mathbf{q}} - 1, \ldots\}, t\right)$$

$$+ e^{-i\mathbf{q}\cdot\mathbf{r}} \sqrt{m_{\mathbf{q}} + 1} f_{\mathrm{g}}\left(\mathbf{k} - \frac{\mathbf{q}}{2}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \ldots m_{\mathbf{q}} + 1, \ldots\}, t\right) \tag{56}$$

Each term under the sum represents a phonon interaction event that changes only one set of phonon variables, increasing or decreasing the occupation number of the single-phonon state $|\mathbf{q}\rangle$ by one and changing the electron momentum by $\pm\mathbf{q}/2$.

## 3.2. A Hierarchy of Transport Equations

The equation for the generalized Wigner function (56) is too complex for the purpose of mesoscopic device simulation. Several approximations need to be introduced in order to arrive at a more feasible quantum transport equation. In the following, these approximations are discussed.

### 3.2.1. Weak Scattering Limit

The generalized Wigner equation couples one element of the phonon density matrix, $f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$, with four neighboring elements,

$$f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \ldots, n_{\mathbf{q}} \pm 1, \ldots\}, \{m\}, t) \tag{57}$$

$$f_{\mathrm{g}}(\mathbf{k}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \ldots, m_{\mathbf{q}} \pm 1, \ldots\}, t) \tag{58}$$

The equations for the four nearest neighbor elements couple to second nearest neighbors of the element $\{n\}, \{m\}$, and so forth. In the weak scattering limit, all couplings between elements of the first and the second off-diagonals are neglected. Only the main diagonal terms and the first off-diagonal terms remain, as shown in Fig. 1. Higher order electron–phonon interactions are neglected in this way.

**Figure 1.** Terms of of the phonon density matrix retained in the weak scattering limit.

### 3.2.2. The Reduced Wigner Function

The reduced Wigner function, $f_w(\mathbf{k}, \mathbf{r}, t)$, is defined as the trace of the generalized Wigner function over all phonon states [28, 61].

$$f_w(\mathbf{k}, \mathbf{r}, t) = \sum_{\{n\}} f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{n\}, t) \tag{59}$$

Further approximations are needed to evaluate this trace and hence to derive a closed equation for the reduced Wigner function [62]. One approximation is to replace any occupation number $n_q$ involved in a transition by the equilibrium phonon number, $N_q$, and to assume that the phonon system stays in equilibrium during the evolution of the electron state. With these assumptions, the trace operation can be performed, and a closed equation set for the reduced Wigner function can be obtained. The set consists of an equation for the reduced Wigner function coupled to two auxiliary equations.

$$\left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r - \Theta_w\right) f_w(\mathbf{k}, \mathbf{r}, t)$$

$$= 2\operatorname{Re} \sum_q \mathcal{F}^2(q) \left\{ e^{iq\cdot r} f_1\left(\mathbf{k} - \frac{q}{2}, \mathbf{r}, t\right) - e^{-iq\cdot r} f_2\left(\mathbf{k} + \frac{q}{2}, \mathbf{r}, t\right) \right\} \tag{60}$$

In this equation, we denote the Wigner potential operator by $\Theta_w$ and set the classical force to $\mathbf{F} = 0$.

$$\Theta_w[f_w](\mathbf{k}, \mathbf{r}, t) = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) \, d\mathbf{k}' \tag{61}$$

The auxiliary equations arise from the first off-diagonal terms of the equation for the generalized Wigner function. In the following equation, the lower sign gives $f_1$ and the upper $f_2$:

$$\left(\frac{\partial}{\partial t} + \frac{\hbar}{m^*}\left(\mathbf{k} \pm \frac{q}{2}\right) \cdot \nabla_r \mp i\omega_q - \Theta_w\right) f_{1,2}\left(\mathbf{k} \pm \frac{q}{2}, \mathbf{r}, t\right)$$

$$= \pm e^{\pm iq\cdot r}\left\{ \left(N_q + \frac{1}{2} \mp \frac{1}{2}\right) f_w(\mathbf{k}, \mathbf{r}, t) - \left(N_q + \frac{1}{2} \pm \frac{1}{2}\right) f_w(\mathbf{k} \pm q, \mathbf{r}, t) \right\} \tag{62}$$

Although the equation for the reduced Wigner function is real-valued, the two auxiliary equations are complex-valued. Note that $f_w$ depends either on some initial momentum $\mathbf{k}$ or

the momentum after a completed electron–phonon interaction, $k \pm q$. On the other hand, $f_1$ and $f_2$ depend on intermediate states $k \pm q/2$, where only half of the phonon momentum has been transferred.

### 3.2.3. Mean Field Approximation

To simplify the equation system, one may assume a mean field over the length scale of an electron–phonon interaction. This mean field can be set to the local force field $\hbar F(r) = -\nabla V(r)$. Note that this field is kept constant during an electron–phonon interaction event, even though the electron moves on an $r$-space trajectory. For a uniform electric field, the potential operator becomes local, $\Theta_w[f_w] = -F \cdot \nabla_k f_w$, and the two auxiliary equations (62) can explicitly be solved. The solutions $f_{1,2}$ are expressed as path integrals over the reduced Wigner function. In this way, a single equation for the reduced Wigner function is derived from (60).

$$\left(\frac{\partial}{\partial t} + \frac{\hbar k}{m^*} \cdot \nabla_r - \Theta_w\right) f_w(k, r, t) = \int_0^t d\tau \int dk' \left[S(k, k', \tau) f_w(k' - F\tau, R(k, k', \tau), t - \tau)\right.$$
$$\left. - S(k', k, \tau) f_w(k - F\tau, R(k, k', \tau), t - \tau)\right] \qquad (63)$$

The scattering kernel is of the form

$$S(k', k, \tau) = \frac{2V}{(2\pi)^3} \bar{F}^2(q) \sum_{\nu=\pm 1} \left(N_q + \frac{1}{2} - \frac{\nu}{2}\right) \cos \int_0^\tau \frac{1}{\hbar}\left(\epsilon_{(k-F\tau)} - \epsilon_{(k-F\tau)} + \nu\hbar\omega_q\right) d\tau' \qquad (64)$$

and the $r$-space trajectory defined as

$$R(k, k', \tau) = r - \frac{\hbar(k + k')}{2m^*}\tau + \frac{\hbar F}{2m^*}\tau^2 \qquad (65)$$

To interpret the above equations, we assume some phase space point $k$, $r$ and some time $t$ to be given. A transition from $k$ to $k'$ as described by (64) starts in the past, at time $t - \tau$, where the retarded momentum $k - F\tau$ has to be considered [see (63)]. At the beginning of the electron–phonon interaction, half of the phonon momentum is transferred, which determines the initial momentum $k - F\tau \pm q/2$ of a phase space trajectory. With $k' = k \pm q$, the initial momentum becomes

$$k - F\tau \pm \frac{q}{2} = \frac{k + k'}{2} - F\tau \qquad (66)$$

During the interaction duration $\tau$, the particle drifts over a phase space trajectory and arrives at $r$ and $k \pm q/2$ at time $t$. At this time, the electron–phonon interaction is completed by the transfer of another $\pm q/2$, which produces the final momentum $k \pm q$. Also included are virtual phonon emission and absorption processes, where the initial momentum transfer $\pm q/2$ at $t - \tau$ is compensated by $\mp q/2$ at $t$. This model thus includes effects due to a finite collision duration, such as collisional broadening and the intra-collisional field effect. A discussion of the integral form of (63) can be found in [63].

### 3.2.4. Levinson Equation

For a uniform electric field and an initial condition independent of $r$, (63) simplifies to the Levinson equation [43].

$$\left(\frac{\partial}{\partial t} + F \cdot \nabla_k\right) f_w(k, t) = \int_0^t d\tau \int dk' \left[S(k, k', \tau) f_w(k' - F\tau, t - \tau)\right.$$
$$\left. - S(k', k, \tau) f_w(k - F\tau, t - \tau)\right] \qquad (67)$$

$S$ is given by (64). This equation is equivalent to the Barker-Ferry equation [64] with an infinite electron lifetime. Recently, Monte Carlo methods for the solution of the Levinson equation have been developed, which allow the numerical study of collisional broadening, retardation effects, and the intracollisional field effect [65, 66].

### 3.2.5. Classical Limit

The classical limit of the scattering operator in (63) is obtained by an asymptotic analysis. For this purpose, the equation is written in a dimensionless form. The primary scaling factors are $k_0$ for the wave-vector $\mathbf{k}$ and $t_0$ for the time $t$. Additional scaling factors to be introduced are $s_0$ for the scattering rate $S$, $\epsilon_0$ for the energy $\epsilon$, $F_0$ for the force $\mathbf{F}$, $r_0$ for the real-space vector $\mathbf{r}$, and $\omega_0$ for the interaction matrix element $\tilde{\jmath}$.

The key issue is now to choose an appropriate scale $k_0$. Scaling the phonon energy to unity gives $\hbar k_0^2 = m^* \omega_q$. The kinetic equation is now considered on a timescale that is much larger than the timescale of the lattice vibrations. Therefore, one sets $t_0 = (\varepsilon \omega_q)^{-1}$, where $\varepsilon \ll 1$ denotes a dimensionless parameter. The remaining scaling factors are found as

$$ s_0 = \frac{1}{t_0^2 k_0^3}, \qquad \epsilon_0 = \frac{\hbar^2 k_0^2}{m^*} \tag{68} $$

$$ F_0 = \frac{k_0}{t_0}, \qquad r_0 = \frac{\hbar k_0 t_0}{m^*} \tag{69} $$

The frequency scale of the electron–phonon interaction can be chosen as $\omega_j = \tilde{\jmath}(q_{th})$, where $q_{th}$ is the wave number of a thermal electron. The scaled Levinson equation has the same form as the unscaled equation (67). The scaled scattering rate varies on a time scale of order $\varepsilon^{-1}$. To keep the time integral of order $O(1)$, the amplitude of the scattering rate should be of order $\varepsilon^-$ as well, which is obtained by setting [67]

$$ \varepsilon = \frac{2V\omega_j^2}{(2\pi)^3} \sqrt{\frac{m_*^3}{\hbar^3 \omega_q}} \tag{70} $$

This gives a scaled scattering rate of the form

$$ S(\mathbf{k}',\mathbf{k},\cdot) = \frac{\tilde{\jmath}^2(q)}{\varepsilon} \sum_{\nu=\pm 1} \left(N_q + \frac{1}{2} - \frac{\nu}{2}\right) \cos \int_0^\tau \frac{1}{\varepsilon} \left(\frac{(\mathbf{k}-\mathbf{F}t')^2}{2} - \frac{(\mathbf{k}'-\mathbf{F}t')^2}{2} + \nu\right)dt' \tag{71} $$

The classical limit is valid in the regime where the quantity defined by (70) is small, and thus for timescales $t_0 = (\varepsilon \omega_q)^{-1}$ much larger than the inverse phonon frequency. The scattering operator in (67) converges for $\varepsilon \to 0$ to the Fermi golden rule operator in the weak sense. From the asymptotic analysis also a first-order correction to the Fermi golden rule is found [67]. Using parameters for GaAs at room temperature, one computes $\varepsilon = 0.011$, which suggests that assuming the asymptotic regime is appropriate.

A heuristic argument for the convergence to the golden rule is as follows. Changing variables in the scattering operator in (67) gives

$$ Q[f_w](\mathbf{k},t) = \int_0^{t/\varepsilon} d\tau \int d\mathbf{k}' \big[\varepsilon S(\mathbf{k},\mathbf{k}',\varepsilon\tau)f_w(\mathbf{k}' - \varepsilon\mathbf{F}\tau, t - \varepsilon\tau) $$
$$ - \varepsilon S(\mathbf{k}',\mathbf{k},\varepsilon\tau) f_w(\mathbf{k} - \varepsilon\mathbf{F}\tau, t - \varepsilon\tau)\big] $$

and

$$ \varepsilon S(\mathbf{k}',\cdot,\varepsilon\tau) = \tilde{\jmath}^2(q) \sum_{\nu=\pm 1} \left(N_q + \frac{1}{2} - \frac{\nu}{2}\right) \cos\left[(\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu)\tau - \varepsilon\mathbf{F}\cdot(\mathbf{k}-\mathbf{k}')\frac{\tau^2}{2}\right] $$

Expanding $f_w$ and $\varepsilon S$ into a Taylor series in $\varepsilon$ and keeping only terms of zeroth order leads to the integral,

$$ \int_0^\infty \cos[(\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu)\tau]d\tau = \pi\delta[\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu] \tag{72} $$

which evaluates to the energy-conserving $\delta$-function of the golden rule. Undoing the scaling gives the well-known form of the scattering rate.

$$S(\mathbf{k}', \mathbf{k}) = \frac{V}{(2\pi)^3} \sum_{r=\pm1} \frac{2\pi}{\hbar} M^2(\mathbf{q}) \left( N_q + \frac{1}{2} - \frac{\nu}{2} \right) \delta\left[\epsilon(\mathbf{k}') - \epsilon(\mathbf{k}) + \nu\hbar\omega_q\right] \qquad (73)$$

The interaction matrix element is denoted here by $M = \hbar\mathcal{F}$. Introducing the total scattering rate $\lambda(\mathbf{k}) = \int S(\mathbf{k}', \mathbf{k})\,d\mathbf{k}'$, the Boltzmann scattering operator takes on the following form.

$$Q[f_w](\mathbf{k}, \mathbf{r}, t) = \int S(\mathbf{k}, \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t)\,d\mathbf{k}' - \lambda(\mathbf{k}, \mathbf{r}) f(\mathbf{k}, \mathbf{r}, t) \qquad (74)$$

Finally, we consider the classical limit of the potential operator. Scaling the $\mathbf{r}$-dependent equation (63) gives the scaled form

$$\Theta_w[f](\mathbf{k}, \mathbf{r}, t) = \frac{1}{i(2\pi)^3\varepsilon} \int d\mathbf{k}' \int d\mathbf{s}\, f(\mathbf{k}', \mathbf{r}, t) \left[ V\left(\mathbf{r} + \frac{\varepsilon\mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\varepsilon\mathbf{s}}{2}\right) \right] e^{-i(\mathbf{k}-\mathbf{k}')\cdot\mathbf{s}}$$

This expression converges for $\varepsilon \to 0$ to the classical drift term of the Boltzmann equation.

$$\Theta_{cl}[f](\mathbf{k}, \mathbf{r}, t) = \nabla_r V(\mathbf{r}) \cdot \nabla_k f(\mathbf{k}, \mathbf{r}, t) \qquad (75)$$

### 3.2.6. Wigner Equation with Boltzmann Scattering Operator

To obtain a model more suitable for device simulation, the nonlocal potential operator is maintained, whereas in the scattering operator the classical limit is introduced. The result is a Wigner equation with a Boltzmann scattering operator. It is convenient to introduce formally a classical force field $\mathbf{F}(\mathbf{r})$ in this equation to make the form of the Liouville operator equal to that of the Boltzmann equation. This is accomplished by redefining the potential operator.

$$\tilde{V}_w(\mathbf{k}, \mathbf{r}) = \frac{1}{(2\pi)^3 i\hbar} \int \left( V\left(\mathbf{r} + \frac{\mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\mathbf{s}}{2}\right) + \hbar\mathbf{F} \cdot \mathbf{s} \right) e^{-i\mathbf{k}\cdot\mathbf{s}}\,d\mathbf{s} \qquad (76)$$

Substituting $\Theta_w[f] = \tilde{\Theta}_w[f] - \mathbf{F} \cdot \nabla_k f$ into (63) gives the following equation,

$$\left( \frac{\partial}{\partial t} + \frac{\hbar\mathbf{k}}{m^*} \cdot \nabla_r + \mathbf{F}(\mathbf{r}) \cdot \nabla_k \right) f_w = Q[f_w] + \tilde{\Theta}_w[f_w] \qquad (77)$$

From a formal point of view, the classical force field $\mathbf{F}$ can be chosen arbitrarily, as the corresponding terms in (77) cancel each other. Typical choices are the mean electric field in a device region, the local electric field, or, of course, $\mathbf{F} = 0$. Alternatively, an equation of the form (77) can also be obtained by using an approximation. The potential is decomposed as $V = V_{cl} + V_{qm}$, where $V_{cl}$ is a smooth potential such as the electrostatic potential, that can be treated in the classical limit (75), and $V_{qm}$ represents a rapidly varying component that has to be treated quantum mechanically.

### 3.3. Integral Form of the Wigner Equation

From the integro-differential form of the Wigner equation, a path-integral formulation can be derived. The equation to be transformed reads

$$\left( \frac{\partial}{\partial t} + \mathbf{v}(\mathbf{k}) \cdot \nabla_r + \mathbf{F}(\mathbf{r}) \cdot \nabla_k \right) f_w(\mathbf{k}, \mathbf{r}, t)$$

$$= \int [S(\mathbf{k}, \mathbf{k}') + \tilde{V}_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r})\delta(\mathbf{k} - \mathbf{k}')] f_w(\mathbf{k}', \mathbf{r}, t)\,d\mathbf{k}'$$

$$- [\lambda(\mathbf{k}, \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r})] f_w(\mathbf{k}, \mathbf{r}, t) \qquad (78)$$

At this point, we introduced a fictious scattering mechanism $\alpha\delta(\mathbf{k} - \mathbf{k}')$, referred to as self-scattering [68]. Because of the $\delta$-function, this mechanism does not change the state of the

electron and hence does not affect the solution of the equation. For the sake of brevity, we define an incgral kernel $\Gamma$ and the symbols $\mu$ and $U'$.

$$\mu(\mathbf{k}, \mathbf{r}) = \lambda(\mathbf{k}, \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r}) \tag{79}$$

$$\Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r}) = \frac{S(\mathbf{k}, \mathbf{k}') + \tilde{V}_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r})\,\delta(\mathbf{k} - \mathbf{k}')}{\mu(\mathbf{k}', \mathbf{r})} \tag{80}$$

$$U(\mathbf{k}, \mathbf{r}, t) = \int \Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r})\,\mu(\mathbf{k}', \mathbf{r})\,f_w(\mathbf{k}', \mathbf{r}, t)\,d\mathbf{k}' \tag{81}$$

The Liouvile operator in (78) is treated by the method of characteristics. One introduces paith variabes $\mathbf{K}(t)$ and $\mathbf{R}(t)$ and takes the total time derivative of $f_w$.

$$\frac{d}{dt} f_w(\mathbf{K}(t), \mathbf{R}(t), t) = \left( \frac{\partial}{\partial t} + \frac{d\mathbf{K}(t)}{dt} \cdot \nabla_k + \frac{d\mathbf{R}(t)}{dt} \cdot \nabla_r \right) f_w \tag{82}$$

The right-hand side equals the Liouville operator if the path variables satisfy the following equations of motion.

$$\frac{d}{dt}\mathbf{K}(t) = \mathbf{F}(\mathbf{R}(t)) \qquad \frac{d}{dt}\mathbf{R}(t) = \mathbf{v}(\mathbf{K}(t)) \tag{83}$$

Now we assume some phase-space point $\mathbf{k}, \mathbf{r}$ and some time $t$ to be given. A phase-space trajectory with the initial condition $\mathbf{K}(t' = t) = \mathbf{k}$ and $\mathbf{R}(t' = t) = \mathbf{r}$ is obtained by formal integration

$$\mathbf{K}(t') = \mathbf{k} + \int_t^{t'} \mathbf{F}(\mathbf{R}(y))\,dy \qquad \mathbf{R}(t') = \mathbf{r} + \int_t^{t'} \mathbf{v}(\mathbf{K}(y))\,dy \tag{84}$$

Note that $\mathbf{k}, \mathbf{r}, t$ are treated as constants in the following derivation, only $t'$ is a variable. Introducin, the functions

$$\tilde{f}_w(t') = f_w(\mathbf{K}(t'), \mathbf{R}(t'), t'), \quad \tilde{\mu}(t') = \mu(\mathbf{K}(t'), \mathbf{R}(t')), \quad \tilde{U}(t') = U(\mathbf{K}(t'), \mathbf{R}(t'), t') \tag{85}$$

alllows (78 to be rewritten as an ordinary differential equation of first order.

$$\frac{d}{dt'}\tilde{f}_w(t') + \tilde{\mu}(t')\tilde{f}_w(t') = \tilde{U}(t') \tag{86}$$

Iff multiplid by an integrating factor $\exp[\int_0^{t'} \tilde{\mu}(y)dy]$, the equation takes on a form that can be easily itegrated in time.

$$\frac{d}{dt'}\exp\left[\int_0^{t'} \tilde{\mu}(y)dy\right]\tilde{f}_w(t') = \exp\left[\int_0^{t'} \tilde{\mu}(y)dy\right]\tilde{U}(t') \tag{87}$$

The choic of the upper and lower bounds of time integration depends on whether the problem nder consideration is time-dependent or stationary.

The orcnary differential equation (87), which is the result of treating the Liouville operator by the nethod of characteristics, has the same structure as the corresponding differential equation or the Boltzmann equation. Therefore, we can refer to the work on the Boltzmann equation egarding the details of the time integration of (87) [69, 70].

### 3.3.1. The Time-Dependent Equation

The upper bound of the time integration should be $t' = t$ to obtain $\tilde{f}_w(t) = f_w(\mathbf{k}, \mathbf{r}, t)$, the value of the unknown at the given phase space point. At $t' = 0$, an initial distribution $f_i(\mathbf{k}, \mathbf{r})$ is assumed to be given. In analogy with the Boltzmann equation [70], the integral form of the Wigner equation is obtained.

$$f_w(\mathbf{k}, \mathbf{r}, t) = \int_0^t dt' \int d\mathbf{k}' \ \exp\left\{ -\int_{t'}^t \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\}$$

$$\times \ \Gamma[\mathbf{K}(t'), \mathbf{k}', \mathbf{R}(t')] \mu[\mathbf{k}', \mathbf{R}(t')] f_w[\mathbf{k}', \mathbf{R}(t'), t']$$

$$+ \exp\left\{ -\int_0^t \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} f_i(\mathbf{K}(0), \mathbf{R}(0)) \tag{88}$$

This equation states that the Wigner function at time $t$ depends on the Wigner function at some previous time $t'$. Using (88) in an iterative procedure, with each iteration the time variable would move to smaller values. Therefore, another equation is desirable that describes the evolution of the system in forward time direction. Such an equation is given by the adjoint equation of (88).

$$g_w(\mathbf{k}', \mathbf{r}', t') = \int_{t'}^{\infty} d\tau \int d\mathbf{k} \ g_w[\mathbf{K}(\tau), \mathbf{R}(\tau), \tau] \exp\left\{ -\int_t^\tau \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\}$$

$$\times \ \Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r}') \mu(\mathbf{k}', \mathbf{r}') + g_0(\mathbf{k}', \mathbf{r}', t') \tag{89}$$

The derivation of the adjoint equation (89) is discussed in detail in [69, 70].

### 3.3.2. The Stationary Equation

In a stationary system, the potential and all material parameters are independent of time. A phase-space trajectory is invariant under time translations. This property can be conveniently used to adjust the time reference of each trajectory [71, 72]. In the stationary case, we assume the phase-space point $\mathbf{k}, \mathbf{r}$ to be given at $t' = 0$. So the initial condition for the phase-space trajectory is $\mathbf{K}(0) = \mathbf{k}$ and $\mathbf{R}(0) = \mathbf{r}$. For the upper bound of time integration of (87), we choose now $t' = 0$ to obtain $\tilde{f}_w(0) = f_w(\mathbf{k}, \mathbf{r})$. The lower time bound has to be chosen such that the functions $\mathbf{K}(t)$ and $\mathbf{R}(t)$ take on values at which the Wigner function is known. In the steady-state, this function is known only at the domain boundary. An appropriate lower time bound is therefore the time when the trajectory enters the simulation domain. This time is denoted by $t_b^-$ and depends on the point $\mathbf{k}, \mathbf{r}$ under consideration. The case that the real space trajectory $\mathbf{R}(t)$ never intersects the domain boundary can occur for a classically bound state. Then the trajectory forms a closed loop and the appropriate choice is $t_b^- = -\infty$. Integration of (87) in the time bounds discussed above results in the integral form of the stationary Wigner equation (cf. [71]).

$$f(\mathbf{k}, \mathbf{r}) = f_0(\mathbf{k}, \mathbf{r}) + \int_{t_b^-(\mathbf{k}, \mathbf{r})}^{0} dt' \int d\mathbf{k}' \ \exp\left\{ -\int_{t'}^0 \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\}$$

$$\times \ \Gamma[\mathbf{K}(t'), \mathbf{k}', \mathbf{R}(t')] \mu(\mathbf{k}', \mathbf{r}') f_w[\mathbf{k}', \mathbf{R}(t')] \tag{90}$$

$$f_0(\mathbf{k}, \mathbf{r}) = f_b \{\mathbf{K}[t_b^-(\mathbf{k}, \mathbf{r})], \mathbf{R}[t_b^-(\mathbf{k}, \mathbf{r})]\} \exp\left\{ -\int_{t_b^-(\mathbf{k}, \mathbf{r})}^{0} \lambda[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \tag{91}$$

Here, $f_b$ denotes the boundary distribution. The integral form (90) represents a backward equation. The corresponding forward equation is given by the adjoint equation.

$$g_w(\mathbf{k}, \mathbf{r}) = g_0(\mathbf{k}, \mathbf{r}) + \int_0^{t_b^-(\mathbf{k}', \mathbf{r})} dt \int d\mathbf{k}' g_w[\mathbf{K}(\tau), \mathbf{R}(\tau)]$$

$$\times \ \exp\left\{ -\int_t^\tau \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \Gamma(\mathbf{k}', \mathbf{k}, \mathbf{r}) \mu(\mathbf{k}, \mathbf{r}) \Theta_D(\mathbf{r}) \tag{92}$$

$\Theta_D$ denotes the indicator function of the simulation domain $D$. The initial conditions for the phase space trajectory are $\mathbf{K}(t) = \mathbf{k}'$ and $\mathbf{R}(t) = \mathbf{r}$.

# 4. THE MONTE CARLO METHOD

Monte Carlo s a numerical method that can be applied to solve integral equations. Applying this method to the various integral formulations of the Wigner equation gives rise to a variety of Monte Carlo algorithms, as discussed in the following.

## 4.1. The General Scheme

This section ntroduces the general scheme of the Monte Carlo method and outlines its application to the solution of integrals and integral equations. To calculate some unknown value $m$ by the Monte Carlo method, one has to find a random variable $\xi$ whose expectation value equals $E\{\xi\} = m$. The variance of $\xi$ is designated $\sigma^2$, with $\sigma$ being the standard deviation.

Now consider $N$ independent random variables $\xi_1, \xi_2, \ldots, \xi_N$ with distributions identical to that of $\xi$. Consequently, their expectation values and their variance are equal.

$$E\{\xi_i\} = m, \qquad \text{Var}\{\xi_i\} = \sigma^2, \qquad i = 1, 2, \ldots, N \tag{93}$$

Expectation value and variance of the sum of all these random variables are given by

$$E\{\xi_1 + \xi_2 + \cdots + \xi_N\} = E\{\xi_1\} + E\{\xi_2\} + \cdots + E\{\xi_N\} = Nm \tag{94}$$

$$\text{Var}\{\xi_1 + \xi_2 + \cdots + \xi_N\} = \text{Var}\{\xi_1\} + \text{Var}\{\xi_2\} + \cdots + \text{Var}\{\xi_N\} = N\sigma^2 \tag{95}$$

Using the properties $E\{c\xi\} = cE\{\xi\}$ and $\text{Var}\{c\xi\} = c^2\text{Var}\{\xi\}$, one obtains from (94) and (95)

$$E\left\{\frac{1}{N}(\xi_1 + \xi_2 + \cdots + \xi_N)\right\} = m \tag{96}$$

$$\text{Var}\left\{\frac{1}{N}(\xi_1 + \xi_2 + \cdots + \xi_N)\right\} = \frac{\sigma^2}{N} \tag{97}$$

Therefore, the random variable

$$\bar{\xi} = \frac{1}{N}\sum_{i=1}^{N}\xi_i \tag{98}$$

has the same expectation value as $\xi$ and an $N$ times reduced variance. A Monte Carlo simulation of the unknown $m$ consists of drawing one random number $\xi$. Indeed, this is equivalent to drawing $i$ values of the random variable $\xi$, and evaluating the sample mean (98).

The Monte Carlo method gives an estimate of both the result and the error. According to the central limit theorem, the sum $\rho_N = \xi_1 + \xi_2 + \cdots + \xi_N$ of a large number of identical random variables is approximately normal. For this reason, the following three-sigma rule holds only approximately

$$P\{|\rho_N - Nm| < 3\sqrt{N\sigma^2}\} \approx 0.997 \tag{99}$$

In this equation, the expectation value and the variance of $\rho_N$ are given by (94) and (95), respectively Dividing the inequality by N and using $\bar{\xi} = \rho_N/N$ we arrive at an equivalent inequality a d the probability will not change:

$$P\left\{|\bar{\xi} - m| < 3\frac{\sigma}{\sqrt{N}}\right\} \approx 0.997 \tag{100}$$

This formul indicates that the sample mean $\bar{\xi}$ will be approximately equal to $m$. The error of this approximation will most probably not exceed the value $3\sigma/\sqrt{N}$. This error evidently approaches zero as $N$ increases [73].

### 4.1.1. Monte Carlo Integration

We apply the Monte Carlo method to the evaluation of an integral.

$$m = \int_a^b \phi(x)\,dx \tag{101}$$

For this purpose, the integrand has to be decomposed into a product $\phi = p\psi$, where $p$ is a density function, which means that $p$ is non-negative and satisfies $\int_a^b p(x)\,dx = 1$. Integral (101) becomes

$$m = \int_a^b p(x)\psi(x)\,dx \tag{102}$$

and denotes the expectation value $m = E\{\Psi\}$ of some random variable $\Psi = \psi(X)$. Now the general scheme described in the previous section can be applied. First, a sample $x_1, \ldots, x_N$ is generated from the density $p$. Then the sample $\psi_1, \ldots, \psi_N$ is obtained by evaluating the function $\psi$: $\psi_i = \psi(x_i)$. The sample mean

$$m \simeq \bar{\psi} = \frac{1}{N}\sum_{i=1}^N \psi_i \tag{103}$$

approximates the expectation value. To employ the error estimation (100), the variance of $\Psi$ can be approximately evaluated by the sample variance

$$\sigma^2 \simeq \bar{\sigma}^2 = \frac{1}{N-1}\sum_{i=1}^N (\psi_i - \bar{\psi})^2 \tag{104}$$

Because the factorization of the integrand is not unique, different random variables can be introduced depending on the choice of the density $p$. All of them have the same expectation value but different variance.

### 4.1.2. Integral Equations

The kinetic equations considered in this work can be formulated as integral equations of the form

$$f(x) = \int K(x, x')f(x')\,dx' + f_0(x) \tag{105}$$

where the kernel $K$ and the source term $f_0$ are given functions. Equations of this form are known as Fredholm integral equations of the second kind. In the particular cases of the Boltzmann equation and the Wigner equation, the unknown function $f$ represents the phase-space distribution function. The multidimensional variable $x$ stands for $(k, r, t)$ in the transient case and for $(k, r)$ in the steady state.

Substituting (105) recursively into itself gives the Neumann series, which if convergent, is a formal solution to the integral equation [74].

$$f = f^{(0)} + f^{(1)} + f^{(2)} + \ldots \tag{106}$$

The iteration terms are defined recursively beginning with $f^{(0)}(x) = f_0(x)$.

$$f^{(n+1)}(x) = \int K(x, x')f^{(n)}(x')\,dx', \qquad n = 0, 1, 2, \ldots \tag{107}$$

The series (106) yields the function value in some given point $x$. However, in many cases one is interested in mean values of $f$ rather than in a point-wise evaluation. Such a mean value represents a linear functional and can be expressed as an inner product.

$$(f, A) = \int f(x)A(x)\,dx \tag{108}$$

It is to note that (105) is a backward equation. The corresponding forward equation is given by the adjoint equation.

$$g(x') = \int_{J} K^{*}(x', x)g(x)\,dx + A(x') \tag{109}$$

where the kernel is defined by $K^{*}(x', x) = K(x, x')$. Multiplying (105) by $g(x)$ and (109) by $f(x')$, and integrating over $x$ and $x'$, respectively, results in the equality

$$(f, A) = (g, f_0) \tag{110}$$

By means of (110), one can calculate a statistical mean value not only from $f$, but also from $g$, the solution of the adjoint equation. The given function $A$ has to be used as the source term of the adjoint equation. The link with the numerical Monte Carlo method is established by evaluating the terms of the Neumann series by Monte Carlo integration, as pointed out in the previous section.

Note that usage of (110) precludes a point-wise evaluation of the distribution function using a forward algorithm, because $A(x) = \delta(x)$ cannot be treated by the Monte Carlo method. The probability for a continuous random variable $x'$ to assume a given value $x$ is zero. Only the probability of finding $x'$ within a small but finite volume around $x$ is non-zero.

## 4.2. Particle Models

Each term of the Neumann series of the adjoint equation describes a sequence of alternating free flight and scattering events. A transition consisting of a free flight with initial state $\mathbf{k}_i$ at time $t_i$ and a scattering process to the final state $\mathbf{k}_f$ at time $t_f$ is described by the following expression. For the sake of brevity, the $\mathbf{r}$-dependence of $\Gamma$ and $\mu$ is omitted in the following.

$$P(\mathbf{k}_f, t_f, \mathbf{k}_i, t_i) = \Gamma[\mathbf{k}_f, \mathbf{K}_i(t_f)]\mu[\mathbf{K}_i(t_i)]\exp\left\{ -\int_{t_i}^{t_f} \mu(\mathbf{K}_i(\tau))\,d\tau \right\} \tag{111}$$

In a Monte Carlo simulation, $t_f$, the time of the next scattering event, is generated from an exponential distribution, given by the terms $\mu \exp()$ in (111). Then, a transition from the trajectory end point $\mathbf{K}_i(t_f)$ to the final state $\mathbf{k}_f$ is realized using the kernel $\Gamma$. In contrast to the classical case, where $P$ would represent a transition probability, such an interpretation is not possible in the case of the Wigner equation because $P$ is not positive semidefinite. The problem originates from the Wigner potential, which assumes positive and negative values. However, because of its antisymmetry with respect to $\mathbf{q}$, the Wigner potential can be reformulated in terms of one positive function $V_w^+$ [27].

$$V_w^{+}(\mathbf{q}, \mathbf{r}) = \max(0, V_w(\mathbf{q}, \mathbf{r})) \tag{112}$$

$$V_w(\mathbf{q}, \mathbf{r}) = V_w^{+}(\mathbf{q}, \mathbf{r}) - V_w^{+}(-\mathbf{q}, \mathbf{r}) \tag{113}$$

Then, the kernel $\Gamma$ is rewritten as a sum of the following conditional probability distributions.

$$\Gamma(\mathbf{k}, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}' - \mathbf{k}) + \frac{\gamma}{\mu} [w(\mathbf{k}, \mathbf{k}') - w^{*}(\mathbf{k}, \mathbf{k}')] \tag{114}$$

$$s(\mathbf{k}, \mathbf{k}') = \frac{S(\mathbf{k}', \mathbf{k})}{\lambda(\mathbf{k}')}, \qquad w(\mathbf{k}, \mathbf{k}') = \frac{V_w^{+}(\mathbf{k} - \mathbf{k}')}{\gamma}, \qquad w^{*}(\mathbf{k}, \mathbf{k}') = w(\mathbf{k}', \mathbf{k}) \tag{115}$$

The normalization factor associated with the Wigner potential is defined as

$$\gamma(\mathbf{r}) = \int V_w^{+}(\mathbf{q}, \mathbf{r})\,d\mathbf{q} \tag{116}$$

In the following, different variants of generating the final state $\mathbf{k}_f$ from the kernel $\Gamma$ will be discussed.

## 4.2.1. The Markov Chain Method

In analogy to the simple integral (102), we have now to decompose the kernel $P$ into a transition probability $p$ and the remaining function $P/p$. More details on the Markov chain method can be found in [75, 76]. With respect to (111), one could use the absolute value of $\Gamma$ as a transition probability. Practically, it is more convenient to use the absolute values of the components of $\Gamma$, giving the following transition probability.

$$p(\mathbf{k}_\mathrm{f}, \mathbf{k}') = \frac{\lambda}{\nu}\, s(\mathbf{k}_\mathrm{f}, \mathbf{k}') + \frac{\alpha}{\nu}\, \delta(\mathbf{k}_\mathrm{f} - \mathbf{k}') + \frac{\gamma}{\nu}\, w(\mathbf{k}_\mathrm{f}, \mathbf{k}') + \frac{\gamma}{\nu}\, w^*(\mathbf{k}_\mathrm{f}, \mathbf{k}') \qquad (117)$$

The normalization factor is $\nu = \lambda + \alpha + 2\gamma$. In the first method considered here, the free-light time is generated from the exponential distribution appearing in (111).

$$p_\mathrm{t}(t_\mathrm{f}, t_\mathrm{i}, \mathbf{k}_\mathrm{i}) = \mu[\mathbf{K}_\mathrm{i}(t_\mathrm{i})] \exp\left\{-\int_{t_\mathrm{i}}^{t_\mathrm{f}} \mu[\mathbf{K}_\mathrm{i}(\tau)]\,d\tau\right\} \qquad (118)$$

For the sake of brevity, the state at the end of the free flight is labeled $\mathbf{k}' = \mathbf{K}_\mathrm{i}(t_\mathrm{i})$ in he following. To generate the final state $\mathbf{k}_\mathrm{f}$, one of the four terms in (117) is selected with he associated probabilities $\lambda/\nu$, $\alpha/\nu$, $\gamma/\nu$, and $\gamma/\nu$, respectively. Apparently, these probabiliies sum up to one. If classical scattering is selected, $\mathbf{k}_\mathrm{i}$ is generated from $s$. If self-scatterin is selected, the state does not change and $\mathbf{k}_\mathrm{f} = \mathbf{k}'$ holds. If the third or fourth term are selected, the particle state is changed by scattering from the Wigner potential and $\mathbf{k}_\mathrm{i}$ is selected from $w$ or $w^*$, respectively. The particle weight has to be multiplied by the ratio

$$\frac{\Gamma}{p} = \pm\left(1 + \frac{2\gamma}{\lambda + \alpha}\right) \qquad (119)$$

where the minus sign applies if $\mathbf{k}_\mathrm{i}$ has been generated from $w^*$. For instance, for a quantim mechanical system, where the classical scattering rate $\lambda$ is less than the Wigner scatterng rate $\gamma$, the self-scattering rate $\alpha$ can be chosen in such a way that $\lambda + \alpha = \gamma$. Then, the nultiplier (119) evaluates to $\pm 3$. An ensemble of particles would evolve as shown schematically in Fig. 2.

In the second method, we again use the transition rate (117), but now the free flight tine is generated with rate $\nu$ rather than with $\mu$. In this case, (111) can be rewritten as

$$P(\mathbf{k}_\mathrm{f}, t_\mathrm{f}, \mathbf{k}_\mathrm{i}, t_\mathrm{i}) = \frac{\Gamma(\mathbf{k}_\mathrm{f}, \mathbf{k}')\,\mu(\mathbf{k}')}{p(\mathbf{k}_\mathrm{f}, \mathbf{k}')\,\nu(\mathbf{k}')}\, p(\mathbf{k}_\mathrm{f}, \mathbf{k}')$$

$$\times \nu(\mathbf{k}') \exp\left\{-\int_{t_\mathrm{i}}^{t_\mathrm{f}} \nu[\mathbf{K}_\mathrm{i}(\tau)]\,d\tau\right\} \exp\left\{2\int_{t_\mathrm{i}}^{t_\mathrm{f}} \gamma[\mathbf{R}_\mathrm{f}(\tau)]\,d\tau\right\} \qquad (120)$$

The exponential distribution distribution is used to generate $t_\mathrm{f}$ and the distribution $p$ to generate $\mathbf{k}_\mathrm{i}$. The remaining terms form the factor by which the particle weight changes



Figure 2. With the Markov chain method, the number of numerical particles is conserved. The magnitude of the particle weight increases with each event, and the sign of the weight changes randomly according to a gven probability distribution.

during one free flight. Because of $(\Gamma\mu)/(pr) = \pm 1$, the multiplier for the $i$th free flight evaluates to

$$m_i = \pm\exp\left\{2\int_{t_i}^{t_{i+1}} \gamma[\mathbf{R}_i(\tau)]\,d\tau\right\} \tag{121}$$

Note that the absolute values of both multipliers, (119) and (121), are always greater than one. With each transition of the Markov chain, the particle weight is multiplied by such factor. Thereore, the absolute value of the particle weight will inevitably grow with the number of transitions on the trajectory. To solve the problem of growing particle weights, one can split particles. In this way, an increase in particle weight is transformed to an increase in particle number.

### 4.2.2. Pair Generation Methods

The basic idea of splitting is refined so to avoid fractional weights. Different interpretations of the kernel are presented that conserve the magnitude of the particle weight. Choosing the initial weight to be $+1$, all generated particles will have weight $+1$ or $-1$. This is achieved by interpreting the potential operator in (77) as a generation term of positive and negative particles. We consider the kernel (114).

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_f, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}_i - \mathbf{k}') + \frac{\gamma}{\mu}\left[w(\mathbf{k}_i, \mathbf{k}') - w^*(\mathbf{k}_i, \mathbf{k}')\right] \tag{122}$$

If the Wigner scattering rate $\gamma$ is larger than the classical scattering rate $\lambda$, the self-scattering rate $\alpha$ has to be chosen large enough to satisfy the inequality $\gamma/\mu \leq 1$. Typical choices are $\mu = \text{Max}(\lambda, \gamma)$ or $\mu = \lambda + \gamma$. These expressions also hold for the less interesting case $\gamma < \lambda$, where quantum interference effects are less important than classical scattering effects.

As in the classical Monte Carlo method, the distribution of the free-flight duration is given by the exponential distribution (118). At the end of a free flight, the complementary probabilities $p_c = \lambda/\mu$ and $1 - p_s = \alpha/\mu$ are considered. With probability $p_s$, classical scattering is selected The final state is generated from $s$. The complementary event is self-scattering. In addition, with probability $p_w = \gamma/\mu$ a pair of particle states is generated from the distributions $w$ and $w^*$. The multiplier of the weight is $+1$ for a state generated from one of first three terms and $-1$ for a state generated from $w^*$. Therefore, the magnitude of the initial particle weight is conserved, as shown in Fig. 3.

**Method G1:** In the following, we discuss the case $\gamma > \lambda$, where quantum effects are dominant. We begin with the smallest possible value for $\mu$: $\mu = \text{Max}(\lambda, \gamma) = \gamma$. Because $p_w = \gamma/\mu = 1$, a particle pair is generated after each free flight as shown in Fig. 4. At the same instances, classical or self-scattering events occur. In Fig. 4 and the following figures, only the trajectory of a sample particle is shown and not the whole cascade of trajectories of the generated particles.

**Method G2:** Choosing the self-scattering rate to be $\alpha = \gamma$, the kernel can be regrouped as

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_i, \mathbf{k}') + \left(1 - \frac{\lambda}{\mu}\right)\left[\delta(\mathbf{k}_f - \mathbf{k}') + w(\mathbf{k}_f, \mathbf{k}') - w^*(\mathbf{k}_f, \mathbf{k}')\right] \tag{123}$$



Figure 3. With the pair generation method, the magnitude of the particle weight is conserved, but one initial particle generates a cascade of numerical particles. At all times, mass is exactly conserved.

**Figure 4.** Trajectory of a sample particle resulting from method G1.

With probability $p_s = \lambda/\mu$, classical scattering is selected. Otherwise, a self-scattering event and a pair generation event occur. In this algorithm, classical scattering and pair generation cannot occur at the same time, as shown in Fig. 5. Compared to method G1, the average free flight time is now reduced, because $\mu$ has been increased from $\gamma$ to $\lambda + \gamma$.

### 4.2.3. Single-Particle Generation Methods

The idea of this method is to further reduce the free-flight time. We rewrite the kernel as

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_f, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}_f - \mathbf{k}') + \frac{2\gamma}{\mu} \left[ \frac{1}{2} w(\mathbf{k}_i, \mathbf{k}') - \frac{1}{2} w^*(\mathbf{k}_f, \mathbf{k}') \right] \tag{124}$$

In this case, the self-scattering rate $\alpha$ has to be chosen large enough to satisfy the inequality $2\gamma/\mu \leq 1$. Typical choices are $\mu = \text{Max}(\lambda, 2\gamma)$ and $\mu = \lambda + 2\gamma$. As in method G1, classical scattering is selected with probability $p_s = \lambda/\mu$, whereas the complementary event is self-scattering. In addition, with probability $p_w = 2\gamma/\mu$, particle generation is selected. If selected, with equal probability either the distribution $w$ or $w^*$ is chosen to generate the final state $\mathbf{k}_i$. If $w^*$ has been chosen, the weight is multiplied by $-1$.

**Method G3:** Assuming $\gamma > \lambda/2$ and $\mu = 2\gamma$ gives $p_w = 1$. Therefore, after each free flight, either a positive or negative particle is generated, as depicted in Fig. 6. At the same instances, classical or self-scattering events occur.

Note that in method G3 ($\mu = 2\gamma$), the free-flight time is reduced by a factor of two compared to method G1 ($\mu = \gamma$), which means that now the kernel is applied twice as frequently. In method G3, single particles are generated at a rate of $2\gamma$, whereas in method G1 particle pairs are generated at half of this rate.

**Method G4:** In this method, we set $\alpha = 2\gamma$ and obtain $\mu = \lambda + 2\gamma$. In analogy with method G2, classical scattering and particle generation are now complementary events. Figure 7 indicates that these two types of events occur at different times. From all methods discussed above, this method uses the shortest free flight time.

From a numerical point of view, method G1 and method G2 have the advantage that they exactly conserve charge as they generate particles pairwise with opposite sign. Method G3 and method G4 generate only one particle each time. Because the sign of the weight is selected randomly, charge is conserved only on average. Simulation experiments, however, have shown that the quality of the pseudo random number generator is good enough to generate almost equally many positive and negative particles even during long simulator.



**Figure 5.** Trajectory of a sample particle resulting from method G2.

**Figure 6.** Trajectory of a sample particle resulting from method G3.

times, such that the small difference of net generated particles has no visible effect on the solution.

### 4.2.4. Other Methods

In method G1 to method G4, the weight of the generated particles is $\pm 1$, because the generation rate used equals $2\gamma$. If a generation rate larger than $2\gamma$ or a fixed time-step less than $1/2\gamma$ were used, the magnitude of the generated weight would be less than one. This approach has been followed in [24], where the resulting fractional weights are termed affinities. On the other hand, a generation rate less than $2\gamma$ would result in an under-sampling of the physical process. Then, the magnitude of the generated weights would be generally greater than one.

### 4.3. The Negative Sign Problem

In the following, we analyze the growth rates of particle weights and particle numbers associated with the different Monte Carlo algorithms. In the first Markov chain method discussed in Section 4.2.1, the weight increases at each scattering event by the multiplier (119). The growth rate of the weight can be estimated for the case of constant coefficients $\gamma$ and $\mu$. Because free-flight times are generated with rate $\mu$, the mean free-flight time will be $1/\mu$. During a given time interval $t$, on-average $n = \mu t$ scattering events will occur. The total weight is then estimated asymptotically for $t \gg 1/\mu$.

$$|W(t)| = \left(1 + \frac{2\gamma}{\mu}\right)^n = \left(1 + \frac{2\gamma t}{n}\right)^n \simeq \exp(2\gamma t) \tag{125}$$

This expression shows that the growth rate is determined by the Wigner scattering rate $\gamma$ independently of the classical and the self-scattering rates.

With the second Markov chain method, one readily obtains that the total weight after $n$ free flights grows as a function of the path integral over $\gamma[\mathbf{R}(\tau)]$.

$$|W(t_n)| = \prod_{i=0}^{n-1} |m_i| = \exp\left\{2 \int_0^{t_n} \gamma(\mathbf{R}(\tau)) \, d\tau\right\} \tag{126}$$

In this equation, the $m_i$ are given by (121). This result generalizes (125) for a position-dependent $\gamma$. The growth rate $2\gamma$ is equal to the $L_1$ norm of the Wigner potential.



**Figure 7.** Trajectory of a sample particle resulting from method G4.

In the pair generation methods, the potential operator

$$\Theta_w[f_w](\mathbf{k}) = \int V^+(\mathbf{q})[f_w(\mathbf{k} - \mathbf{q}) - f_w(\mathbf{k} + \mathbf{q})]\,d\mathbf{q} \tag{127}$$

has been interpreted as a generation term. It describes the creation of two new states, $\mathbf{k} - \mathbf{q}$ and $\mathbf{k} + \mathbf{q}$. The generation rate is equal to $\gamma$. When generating the second state, the sign of the statistical weight is changed. It should be noted that the Wigner equation strictly conserves mass, as can be seen by taking the zero-order moment of (77).

$$\frac{\partial n}{\partial t} + \text{div}\,\mathbf{J} = 0 \tag{128}$$

Looking at the number of particles regardless of their statistical weights, that is, counting each particle as positive, would correspond to using the following potential operator.

$$\Theta_w^*[f_w](\mathbf{k}) = \int V_w^+(\mathbf{q})[f_w(\mathbf{k} - \mathbf{q}) + f_w(\mathbf{k} + \mathbf{q})]\,d\mathbf{q} \tag{129}$$

Using (129), a continuity equation for numerical particles is obtained.

$$\frac{\partial n^*}{\partial t} + \text{div}\,\mathbf{J}^* = 2\gamma(\mathbf{r})n^* \tag{130}$$

Assuming a constant $\gamma$, the generation rate in this equation will give rise to an exponential increase in the number of numerical particles $N^*$.

$$N^*(t) = N^*(0)\exp(2\gamma t) \tag{131}$$

This discussion shows that the appearance of an exponential growth rate is independent of the details of the particular Monte Carlo algorithm, and must be considered to be a fundamental consequence of the non-positive kernel.

## 4.4. Particle Annihilation

The discussed particle models are instable, because either the particle weight or the particle number grows exponentially in time. Using the Markov chain method, it has been demonstrated that tunneling can be treated numerically by means of a particle model [25]. However, because of the exponentially increasing particle weight at the very short timescale $(2\gamma)^{-1}$, application of this algorithm turned out to be restricted to single-barrier tunneling and small barrier heights only. This method can be useful for devices where quantum effects are weak, and the potential operator is a small correction to the otherwise classical transport equation.

A stable Monte Carlo algorithm can be obtained by combining one of the particle generation methods with a method to control the particle number. One can assume that two particles of opposite weight and a sufficiently small distance in phase space annihilate each other. The reason is that the motions of both particles are governed by the same equation. Therefore, when they come close to each other at some time instant, the two particles have approximately the same initial condition and thus a common probabilistic future. In an ensemble Monte Carlo method, a particle removal step should be performed at given time steps. During the time step, the ensemble is allowed to grow to a certain limit, then particles are removed and the initial size of the ensemble is restored. In this work, the problem has been solved for the stationary transport problem. In the algorithm, the trajectory of only one sample particle is followed, whereas other numerical particles are temporarily stored on a phase space grid. Due to the opposite sign, particle weights annihilate to a large extent in the cells of the grid. The total residual weight in each cell has to be minimized, as it represents a measure for the numerical error of the method [32].

# 5. SIMULATION RESULTS

Virtually all published results of Wigner function–based device modeling focus on resonant tunneling diodes [77, 78]. In this section, three different devices are discussed. Their parameter values are collected in Table 1, where RTD1 [36] and RTD2 [24] are devices from literature. The semiclassical scattering model includes polar optical, acoustic deformation potential, and ionized-impurity scattering. Parameter values for GaAs have been assumed.

## 5.1. Comparison with Other Numerical Methods

RTD1 has been used as a benchmark device to compare different numerical approaches to quantum transport. In this device, the potential is assumed to vary linearly only in the double-barrier region and to be constant in the two contact regions. Results of the Monte Carlo method outined in this work have been compared to nonequilibrium Green's function–based results [79].The latter have been obtained by NEMO-1D, a one-dimensional nanoelectronic modeling tool [80]. NEMO-1D has served as a quantitatively predictive design and analysis tool for resonant tunneling diodes [81–83].

RTD1 shows a rather large coherent off-resonant valley current. Therefore, phonon scattering has only little effect on the current–voltage characteristics of this device. Both simulators predict only a slight increase in valley current due to inelastic scattering (Fig. 8). The resonance voltages predicted by the two solvers agree very well.

A compaison between finite-difference results and Monte Carlo results is shown in Fig. 9. An important parameter is the cutoff length $L_c$ used in the numerical Wigner transformation. Assuming only one spatial coordinate, $L_c$ is introduced as follows.

$$V_w(k_x, x) = \frac{1}{2\pi i \hbar} \int_{-\infty}^{\infty} \left[ V\left(x + \frac{s}{2}\right) - V\left(x - \frac{s}{2}\right) \right] e^{-i k_x s} \, ds \tag{132}$$

$$\approx -\frac{1}{\pi i \hbar} \int_{0}^{L_c/2} \left[ V\left(x + \frac{s}{2}\right) - V\left(x - \frac{s}{2}\right) \right] \sin(k_x s) \, ds \tag{133}$$

The cutoff length has to be selected carefully when solving the Wigner equation numerically. The comparison of current–voltage characteristics shown in Fig. 9 demonstrates that only a sufficiently large value for $L_c$ gives a realistic result. A too small value results in an overestimation of the valley current.

## 5.2. The Effect of Scattering

In RTD2, the potential changes linearly in a region of 40 nm length, starting 10 nm before the emitter barrier and extending 19 nm after the collector barrier, as shown in Fig. 10. The Wigner potential is discretized using $N_k = 640$ equidistant $k_x$ points and $\Delta x = 0.5$ nm spacing in $x$-direction. Assuming a cutoff length of $L_c = 80$ nm, one would require at least $N_k = L_c/\Delta x = 160$. This minium value is often used in finite-difference simulations for the Wigner equation, but in the Monte Carlo simulation we use the considerably larger value stated above in order to get a better resolution of the energy domain. The annihilation mesh is three-dimensional. In $x$-direction, the grid covers the region where the Wigner potential is nonzero. Because of the cylindrical symmetry of the Wigner function, only two momentum coordinate have to be considered. The mesh extends to an energy of 6 eV in both axial and radial $k$-diection.

Table 1. Parameter values of the simulated resonant tunneling diodes.[a]

| Device name | Barrier height (eV) | Barrier width (nm) | Well width (nm) | Device length (nm) | Contact doping (cm⁻³) |
|---|---|---|---|---|---|
| RTD1 | 0.27 | 2.83 ($5a_0$) | 4.52 ($8a_0$) | 100.6 ($178a_0$) | $2 \times 10^{18}$ |
| RTD2 | 0.3 | 3.0 | 5.0 | 200.0 | $10^{17}$ |
| RTD3 | 0.47 | 3.0 | 4.0 | 270.0 | $10^{18}$ |

[a] The lattice constant of GaAs is $a_0 = 0.565$ nm.

**Figure 8.** Current-voltage characteristics of RTD1 at 300 K obtained from Wigner Monte Carlo and NEMO-1D. Transport is coherent (coh.) or dissipative (scatt.).

The Wigner generation rate (127) is of the order $10^{15} s^{-1}$ for RTD2 (Fig. 11). The relation of this rate to the typically much smaller semiclassical scattering rate is a quantitative measure of the fact that quantum interference effects are dominant. The zero-field contact regions have been chosen sufficiently large, such that the Wigner potential drops to zero within these regions.

Figure 12 shows the electron concentration in RTD2 at voltages below the resonance voltage. Classical behavior is observed before and after the double barrier, whereas in the



**Figure 9.** Effect of the cutoff length on the current-voltage characteristics in Wigner simulations. The finite-difference (FD) result is taken from [36].

**Figure 10.** Conduction band edge of RTD2 for different voltages. A linear voltage drop is assumed.

quantum well the behavior of the solution is nonclassical. In front of the barrier an accumulation layer forms, with its maximum concentration increasing with the band bending. In the quantum well, the concentration increases as the resonance is approached. After the barrier a depletion layer forms, which grows with applied voltage. In this region, the concentration at 0.15 V varies exponentially in response to the linear potential (see Fig. 10), which is again a classical property.

For voltages above the resonance voltage, the concentration in the well drops, whereas the depletion layer continues to grow (Fig. 13). The mean kinetic energy of the electrons is depicted in Fig. 14. The energy density has been calculated from the second-order moment of the Wigner function (35) and divided by the electron density to get the mean energy per electron. In the zero-field regions, an energy close to the equilibrium energy is obtained, which demonstrates that the energy conservation property of the Wigner potential operator is also satisfied by the numerical Monte Carlo procedure. One has to keep in mind that the Wigner potential can produce a rather large momentum transfer. For the chosen value



**Figure 11.** Pair generation rate $\gamma(x)$ in RTD2 caused by the Wigner potential for two different voltages.

**Figure 12.** Electron concentration in RTD2 for voltages less than the resonance voltage.

for $\Delta x$, the related energy transfer can reach values as large as 5 eV, which shows that a large degree of cancellation occurs in the estimator for the mean energy. Electrons injected from the second barrier into the collector space charge region show initially a high kinetic energy.

Phonon scattering strongly affects the current–voltage characteristic of RTD2 (Fig. 15). As compared to the coherent case, phonon scattering leads to an increase in the valley current and a resonance voltage shift. The large difference in the valley current can be explained by the electron concentration in off-resonance condition (Fig. 16). With phonon scattering included, a significantly higher concentration forms in the emitter notch, and injection in the double barrier is increased. This indicates that a quasi bound state forms in the emitter notch. The population of this state increases when scattering is switched on. On the other hand, in resonance condition where the applied voltage is lower, such a bound state does not form and very similar electron concentrations are observed for the coherent and noncoherent case (Fig. 17).



**Figure 13.** Electron concentration in RTD2 for voltages greater than the resonance voltage.

**Figure 14.** Mean kinetic energy in RTD2 for two different voltages.

## 5.3. Inclusion of Extended Contact Regions

As discussed in Section 3.2.5, the Wigner equation simplifies to the Boltzmann equation when the potential variation is sufficiently smooth. The proposed quantum Monte Carlo method turns into the semiclassical Monte Carlo method for vanishing Wigner potential. Therefore, one can simulate a quantum region embedded in an extended classical region with the interface between the regions correctly treated in an implicit way. By means of the Wigner generation rate $\gamma$. the simulation domain can be decomposed into quantum regions ($\gamma' > 0$) and classical regions ($\gamma \simeq 0$). In Fig. 18, these regions within RTD2 are marked. The electron concentration and the mean energy are smooth in the extended contact regions and not affected by the strong onset of the Wigner generation rate, as shown in Fig. 11.

In the simulation of RTD3, the Wigner potential $V_w^i(k_x, x)$ is discretized using $N_k = 1200$ equidistant $k_x$ points and $\Delta x = 0.5$ nm spacing in the $x$-direction. A cutoff length of



**Figure 15.** Influence of phonon scattering on the current–voltage characteristics of the RTD2.

**Figure 16.** Electron concentration in RTD2 in off-resonance condition.



**Figure 17.** Electron concentration in RTD2 in resonance condition.



**Figure 18.** Electron concentration and mean electron energy in RTD2 at $T = 300$ K and 0.1 V applied voltage.

**Figure 19.** Electron concentration profiles in RTD3.

$L_c = 60$ nm is assumed. The annihilation mesh consists of 480 points in the longitudinal and 120 points in the perpendicular momentum direction, and the real space coordinate is discretized using $\Delta x = 0.5$ nm. The electrostatic potential has been computed using the self-consistent Schrödinger-Poisson solver NANOTCAD-1D [84]. Figure 19 shows the electron concentration profile in the device. At the resonance voltage of 1.2 V, the concentration in the quantum well is considerably higher than in the off-resonance condition at 1.6 V. The concentration in the depletion region left of the barrier depends on the injected current and is thus correlated with the concentration in the well.

## 6. CONCLUSION

The examples presented in Section 5 demonstrate that a numerical solver for the Wigner equation can provide quantitatively correct results. One requirement is that the cutoff length is chosen sufficiently large. The completeness relation of the discrete Fourier transform reflecting Heisenberg's uncertainty principle, $\Delta k_A = \pi/L_c$, shows that a small $L_c$ will result in a coarse grid in momentum space, and resonance peaks in the transmission coefficient might not be resolved properly. In the past, the Wigner equation has been solved most frequently by finite-difference methods. Due to the nonlocality of the potential operator, all points in momentum space are coupled, resulting in a very poor sparsity pattern of the matrix. Therefore, increasing the number of grid points in $k$-space, related to the cutoff length by $N_k = L_c/\Delta x$, is limited by prohibitive memory and computation time requirements. This might be one reason why quantitatively correct solutions were difficult to obtain in the past. We believe that the frequently reported accuracy problems with finite-difference Wigner function–based device simulations result from a too coarse $k$-space discretization. As this problem occurs already for one-dimensional geometries, higher dimensional simulations using the finite-difference method are probably out of reach. It is interesting to note that Frensley, who pioneered the finite-difference method for the Wigner equation [16], later abandoned this method and developed the quantum-transmitting boundary method to describe coherent transport in open systems [85].

The Monte Carlo method allows the number of $k$-points to be increased. In this work, the Wigner potential has been discretized using $N_k$ of the order $10^3$. However, high-performance resonant tunneling diodes with very high peak-to-valley current ratio pose still a problem for the Monte Carlo method. In such a device, the density can vary over several orders of magnitude, which often cannot be resolved by the Monte Carlo method. This problem is also well-known from the classical Monte Carlo method. As a solution, one could apply statistical enhancement techniques in such cases. At present, an equidistant $k$-grid is used for the discretization of the Wigner potential. Because the transmission coefficient of double-barrier structures may show very narrow resonance peaks, using an equidistant $k$-grid may not be the optimal choice. However, because of the discrete Fourier transform of the potential

involved in the computation of the Wigner potential, the use of a nonequidistant $k$-grid appears to be problematic.

In a Wigner function–based simulation of one-dimensional heterostructures, fundamental simulation parameters such as the cutoff length are closely linked to physical device paramters such as the spacing from the contacts. This property stems from the choice of plane-wace basis sets in a quantum mechanical regime of broken translational symmetry. Although an.-lytically appealing, this basis set can cause numerical difficulties. Other approaches such s the nonequilibrium Green's function formalism may have the advantage that other basis ses can be used more straightforwardly.

These considerations indicate that from a numerical point of view, the Wigner functin formalism might not be the optimal choice for resonant tunneling simulation. Howeve, because the formalism describes quantum effects and scattering effects with equal accurac, it appears well suited especially when a quasi-ballistic transport condition without energe-ically sharp resonances is present. One strength of the Wigner function approach is tle treatment of contact regions. Nonequilibrium transport can be simulated in the whole devie formed by a central quantum region embedded in extended classical regions. The presented Wigner Monte Carlo method can bridge the gap between classical device simulation ard pure quantum ballistic simulations.

Development of Monte Carlo methods for the solution of the Wigner equation is stil in the beginning. Research efforts are needed especially with respect to the negative sin problem. The particle generation–annihilation algorithm developed by the authors is just oe solution to that problem. Improved variants of this algorithm or even new solution strategis are yet to be devised. Extension of the Monte Carlo methods to higher dimensional devie geometries is straightforward.

# REFERENCES

1.  E. Wigner, *Phys. Rev.* 40, 749 (1932).
2.  G. Mahan, "Many-Particle Physics." Plenum Press, New York, 1983.
3.  V. Tatarskii, *Soviet Physics Uspekhi* 26, 311 (1983).
4.  P. Carruthers and F. Zachariasen, *Rev. Mod. Phys.* 55, 245 (1983).
5.  T. Cutright and C. Zachos, *Mod. Phys. Lett. A* 16, 2381 (2001).
6.  C. Zachos, *Int. J. Mod. Phys. A* 17, 297 (2002).
7.  J. Rammer, *Rev. Mod. Phys.* 63, 781 (1991).
8.  U. Ravaioli, M. Osman, W. Pötz, N. Kluksdahl, and D. Ferry, *Physica B* 134, 36 (1985).
9.  N. Kluksdahl, W. Pötz, U. Ravaioli, and D. Ferry, *Superlattices Microstruct.* 3, 41 (1987).
10. W. Frensley, *Phys. Rev. Lett.* 57, 2853 (1986).
11. W. Frensley, in "International Electron Devices Meeting." p. 571. Institute of Electrical and Electronis Engineers. Los Angeles, CA, 1986.
12. W. Frensley, *Phys. Rev. B* 36, 1570 (1987).
13. N. Kluksdahl, A. Kriman, D. Ferry, and C. Ringhofer, *Phys. Rev. B* 39, 7720 (1989).
14. W. Frensley, *Solid-State Electronics* 32, 1235 (1989).
15. R. Mains and G. Haddad, *J. Appl. Phys.* 64, 5041 (1988).
16. W. Frensley, *Rev. Mod. Phys.* 62, 745 (1990).
17. F. Buot and K. Jensen, *Phys. Rev. B* 42, 9429 (1990).
18. K. Jensen and F. Buot, *Phys. Rev. Lett.* 66, 1078 (1991).
19. F. Buot and K. Jensen, *COMPEL* 10, 241 (1991).
20. K. Gullapalli, D. Miller, and D. Neikirk, *Phys. Rev. B* 49, 2622 (1994).
21. B. Biegel and J. Plummer. *IEEE Trans. Electron Devices* 44, 733 (1997).
22. D. Woolard, P. Zhao, and H. Cui, *Physica B* 314, 108 (2002).
23. R. Mains and G. Haddad, *J. Comput. Phys.* 112, 149 (1994).
24. L. Shifren and D. Ferry, *Physica B* 314, 72 (2002).
25. M. Nedjalkov, R. Kosik, H. Kosina, and S. Selberherr, in "Simulation of Semiconductor Processes ard Devices." p. 187. Business Center for Academic Societies Japan, Kobe, Japan, 2002.
26. L. Shifren, C. Ringhofer, and D. Ferry, *Phys. Lett. A* 306, 332 (2003).
27. H. Kosina, M. Nedjalkov, and S. Selberherr, in "Nanotech" (M. Laudon and B. Romanowicz, Eds.), p 191. Computational Publications, San Francisco, CA, 2003.
28. F. Rossi, C. Jacoboni, and M. Nedjalkov, *Semicond. Sci. Technol.* 9, 934 (1994).
29. M. Nedjalkov, I. Dimov, F. Rossi, and C. Jacoboni, *J. Math. Comp. Modelling* 23, 159 (1996).
30. M. Pascoli, P. Bordone, R. Brunetti, and C. Jacoboni, *Phys. Rev. B* 58, 3503 (1998).
31. P. Bordone, A. Bertoni, R. Brunetti, and C. Jacoboni, *Math. Comp. Simulation* 62, 307 (2003).
32. H. Kosina, M. Nedjalkov, and S. Selberherr, *J. Comput. Electronics* 2, 147 (2003).

33. A. Bertoni, P. Bordone, R. Brunetti, and C. Jacoboni, *J. Phys. Condens. Matter* 11, 5999 (1999).

34. P. Bordone, M. Pascoli, R. Brunetti, A. Bertoni, and C. Jacoboni, *Phys. Rev. B* 59, 3060 (1999).

35. C. Jacoboni, R. Brunetti, P. Bordone, and A. Bertoni, *Int. J. High Speed Electronics Systems* 11, 387 (2001).

36. H. Tsuchiya, M. Ogawa, and T. Miyoshi, *IEEE Trans Electron Devices* 38, 1246 (1991).

37. J.-J. Shih, H.-C. Huang, and G. Wu, *Phys. Rev. B* 50, 2399 (1994).

38. E. Bufler and J. Schlösser, *J. Phys. Condens. Matter* 6, 7445 (1994).

39. D. Miller and D. Neikirk, *Appl. Phys. Lett.* 58, 2803 (1991).

40. L. Demaio, L. Barletti, A. Bertoni, P. Bordone, and C. Jacoboni, *Physica B* 314, 104 (2002).

41. M. B. Inlu, B. Rosen, H.-L. Cui, and P. Zhao, *Phys. Lett. A* 327, 230 (2004).

42. G. Wu and K.-P. Wu, *J. Appl. Phys.* 71, 1259 (1992).

43. I. Levinson, *Soviet Phys. JETP* 30, 362 (1970).

44. P. Holland and K. Kypriandis, *Phys. Rev. A* 33, 4380 (1986).

45. S. Sonego, *Phys. Rev. A* 44, 5369 (1991).

46. D. Ferry and S. Goodnick, "Transport in Nanostructures," Cambridge University Press, Cambridge, UK, 2001.

47. M. Levanda and V. Fleurov, *Ann. Phys.* 292, 199 (2001).

48. W. Hänsch, "The Drift-Diffusion Equation and Its Applications in MOSFET Modeling, Computational Microelectronics," Springer Verlag, Vienna, 1991.

49. R. Brunetti, A. Bertoni, P. Bordone, and C. Jacoboni, in "International Workshop on Computational Electronics" (J. R. Barker and J. R. Watling, Eds.), p. 131. University of Glasgow, Glasgow, Scotland, 2000.

50. J. Zhou and D. Ferry, *IEEE Trans. Electron Devices* 39, 473 (1992).

51. C. Gardner, *SIAM J. Appl. Math.* 54, 409 (1994).

52. C. L. Gardner and C. Ringhofer, *Phys. Rev. E* 53, 157 (1996).

53. P. Degond and C. Ringhofer, *J. Statistical Phys.* 112, 587 (2003).

54. H. Tsuchiya and U. Ravaioli, *J. Appl. Phys.* 89, 4023 (2001).

55. Z. Hai, N. Goldsman, and C. Lin, in "International Electron Devices Meeting," p. 62. Institute of Electrical and Eectronics Engineers, San Francisco, CA, 2000.

56. N. Godsman, C.-K. Lin, Z. Han, and C.-K. Huang, *Superlattices Microstruct.* 27, 159 (2000).

57. A. Faınjiang, S. Jin, and G. Papanicolaou, *SIAM J. Appl. Math.* 63, 1328 (2002).

58. R. Koik, Ph.D. Thesis, Vienna University of Technology, 2004.

59. J. J. Wlodarz, *Phys. Lett. A* 264, 18 (1999).

60. B. Biegel and J. Plummer, *Phys. Rev. B* 54, 8070 (1996).

61. R. Brunetti, C. Jacoboni, and F. Rossi, *Phys. Rev. B* 39, 10781 (1989).

62. M. Neljalkov, R. Kosik, H. Kosina, and S. Selberherr, *Microelectron. Eng.* 63, 199 (2002).

63. M. Nedjalkov, H. Kosina, R. Kosik, and S. Selberherr, *J. Computat. Electronics* 1, 27 (2002).

64. J. Barker and D. Ferry, *Phys. Rev. Lett.* 42, 1779 (1979).

65. M. Neljalkov, H. Kosina, S. Selberherr, and I. Dimov, *VLSI Design* 13, 405 (2001).

66. T. Gurov, M. Nedjalkov, P. Whitlock, H. Kosina, and S. Selberherr, *Physica B* 314, 301 (2002).

67. C. Ringhofer, M. Nedjalkov, H. Kosina, and S. Selberherr, *SIAM J. Appl. Math.* (2004) (in press).

68. C. Jacoboni and L. Reggiani, *Rev. Mod. Phys.* 55, 645 (1983).

69. H. Kosina, M. Nedjalkov, and S. Selberherr, *IEEE Trans. Electron Devices* 47, 1898 (2000).

70. H. Kosina and M. Nedjalkov, *Int. J. High Speed Electronics Systems* 13, 727 (2003).

71. H. Kosina, M. Nedjalkov, and S. Selberherr, *J. Appl. Phys.* 93, 3553 (2003).

72. M. Nedjalkov, H. Kosina, and S. Selberherr, *J. Appl. Phys.* 93, 3564 (2003).

73. I. Sobol, "The Monte Carlo Method," Mir Publishers, Moscow, 1984.

74. F. Byron and R. Fuller, "Mathematics of Classical and Quantum Physics," Dover, New York, 1992.

75. S. Ermakow, "Die Monte-Carlo-Methode und verwandte Fragen." R. Oldenburg Verlag, Munich, 1975.

76. J. Hammersley and D. Handscomb, "Monte Carlo Methods." John Wiley, New York, 1964.

77. J. Sun, G. Haddad, P. Mazumder, and J. Schulman, *Proc. IEEE* 86, 641 (1998).

78. H. Mizuta and T. Tanoue, "The Physics and Applications of Resonant Tunneling Diodes." Cambridge University Press, Cambridge, UK, 1995.

79. H. Kosina, G. Klimeck, M. Nedjalkov, and S. Selberherr, in "Simulation of Semiconductor Processes and Devices," p. 171. Institute of Electrical and Electronics Engineers, Boston, MA, 2003.

80. R. Lake, G. Klimeck, R. Bowen, and D. Jovanovic, *J. Appl. Phys.* 81, 7845 (1997).

81. R. Bowen, G. Klimeck, R. Lake, W. Frensley, and T. Moise, *J. Appl. Phys.* 81, 3207 (1997).

82. G. Klmeck, R. Lake, D. Blanks, C. Fernando, R. C. Bowen, T. Moise, and Y. Kao, *Phys. Status Solidi (b)* 204, 408 (997).

83. G. Kmeck, R. Lake, and D. Blanks, *Phys. Rev. B* 58, 7279 (1998).

84. G. Iannaccone, M. Macucci, P. Coli, G. Curatola, G. Fiori, M. Gattobigio, and M. Pala, in "IEEE Conference on Nnotechnology," p. 117. IEEE, Maui, Hawaii, 2001.

85. W. Frensley, *Superlattices Microstruct.* 11, 347 (1992).

# CHAPTER 15

# Logic Design of Nanodevices

## Svetlana N. Yanushkevich

Department of Electrical and Computer Engineering,
University of Calgary, Calgary, Alberta, Canada

## CONTENTS

# 1. INTRODUCTION

Design of computer devices is experiencing a transition between microelectronics, which s reaching its physical limits at the submicron scale, and the upcoming era of nanotechnolog'. This is because at the nanoscale level, comparably with atomic wave length, microscale phy:- ical processes do not work and, therefore, give place to quantum effects. This, eventuall', yields to a global modeling of computing processes by appropriate physical (quantum) phe- nomena in the foreseeable future.

So far, technology offers electronic devices which exhibit local quantum effects, that s nanoscale components. For example, devices such as resonant-tunneling diodes [1] and single-electron devices [2] have been known for a while. However, our understanding cf information flow and of the computer as a "processor-memory-Datapath-Input/Output" sy:- tem remains a traditional one, and it is global logic designs that are a networked system cf devices with local quantum effects. Since today's quantum devices are weak and sensitive, they are not suited to conventional logic gate architectures, which require robust devices.

This survey explains connections between contemporary approaches to computer hardware design based on traditional Boolean algebra and the new requirements posed by nanc- technology. It abstracts from the traditional gate-level or programmable logic implemer- tation of switching functions, and characterizes the structural requirements of existing and predictable nanoscale devices.

The key features of the proposed view of logic design of nanodevices are the following:

1. A central role is reserved for topological models (embedding in hypercubes, assembling topology).
2. The technique of advanced logic design is deployed and flavored with elements of mas- sively parallel, distributed and tolerant computing, taking into account new possibilities of processing in spatial dimensions.

This approach to the logic design of nanodevices is introduced through:

1. Spatial data structures and the corresponding topological models that satisfy the criteria of massive parallel processing, homogeneity, and fault tolerance: parallel arrays, cellular arrays and neural-like networks.
2. Fault tolerance computing in spatial structures.
3. Methods for analysis of data structures and topologies in nanodimensions.
4. Information theoretical measures in spatial structures.

The rest of the sections are organized as follows. Section 2 overviews the directions and methodology of logic design in nanodimensions. Section 3 covers data structures for logic

design emphasizing on three-dimensional (3D) topological structures. Section 4 presents the technique of a multilevel circuit design based on hypercube structures. Section 5 introduces the technique of analysis in spatial dimensions based on the concept of change. It shows that logic difference is a useful model in some cases for understanding the relations between different data structures and representations of switching functions. Section 6 considers the technique of information-theoretical measures in spatial structures. The reason that the Shannon theory of information should be one of the most important measurement characteristics in nanospace is that it reflects the physical nature and restrictions that nanostructures on a molecular level pose to information carriers. Section 7 discusses the problem of computation using nonreliable elements which nanoscale devices are. due to quantum phenomena and other features of ultrasmall structures.

## 2. LOGIC DESIGN IN NANODIMENSIONS

This section aims at a global perspective on how contemporary logic design techniques can be incorporated into nanosystem design. This can be accomplished by taking advantage of the interdisciplinary approach that involves:

1. Approprate spatial topologies.
2. Information or entropy measures.
3. A distributed parallel processing paradigm.
4. Fault-tolerant computing.

### 2.1. Selected Methods of Advanced Logic Design

Traditional logic design models and techniques may not satisfy the requirements and properties of nanoscale computing devices:

1. While traditional, gate-level, randomly networked circuits are made up of AND/OR/NOT or other gates, working on the principle of voltage state logic, the nanocircuits are supposed to be locally connected arrays of elements (e.g., molecular ones [3-8]). This means that corresponding data structures such as directed acyclic graphs of the netlists or symbolic structures may not work for the purposes of optimization and manipulation of logic functions implemented on nanodevices.
2. While information flow in today's semiconductor devices is associated with surges of electron, and voltage state logic can still be acceptable in some types of nanodevices, in nanodevices this is likely to be associated with states, count of electrons, and so on [9-11].
3. Because the size of the ultrasmall devices is compared to wave length. the nature of the signals and processes is supposed to be stochastic. These devices are very sensitive to many physical factors (thermal fluctuation, wave coherence, random tunneling, etc. [9]) Thus, the need for fault-tolerance computation increases as device to device fluctuations become larger at the furthest limits of integration. While fault-tolerance is achieved in today's microelectronics by duplication of blocks at the block level, this it is not an acceptable practice for silicon devices on the transistor level. Because of power dissipation and clocking and area constraints, it is likely to be achieved in nanostructures though introducing redundant hardware [12, 13].

The fundamentals of switching theory, or Boolean algebra, cannot be changed while technology and even carriers of logic data are being changed. However, an appropriate choice of data structure is the way to adjust implementation of the logic function to the existing technology. That is why this survey will give particular attention to the data structures suitable for logic data processing in nanodimensions.

Two distinct approaches to logic design in nanodimensions can be observed so far:

1. Using the advanced logic design techniques and methods from other disciplines, in particular, fault-tolerant computing.
2. Development of a new theory and technique for logic design in spatial dimensions.

The first direction adapts architectural approach such as, in particular. array-based logic [14], parallel and distributed architectures, methods from fault-tolerant computing and adaptive structures such as neural networks.

The second direction can be justified, in particular, by nanotechnologies that implement devices on the reversibility principle. It is rather relevant to design of adiabatic circuits, which is not associated only with nanotechnology. An example is Q3M (quantum 3D mesh) model of a parallel quantum computer in which each cell has a finite number of quantum bits (qubits) and interacts with neighbors, exchanging particles by means of Hamiltonian derivable from Schrödinger equation [15].

## 2.2. Spatial Nanostructures

Three-dimensional design has been explored at the macrolevel for a long time, for example, for design of distributed systems [16], and "connection machines" [17]. It was inspired by nature; for example, the brain, with its "distributed computing and memory." was the prototype in the case of the "connection machine." The 3D computing architecture concept has been employed by the creators of the supercomputers Cray T3D, NEC's Ncube, and cosmic cube [18]. The components of these supercomputers, as systems, are designed based on classical paradigms of 2D architecture that becomes 3D because of the 3D topology (of interconnects), or 3D data structures and corresponding algorithms, or 3D communication flow.

The topology of today's silicon-integrated circuits at a system level varies: one-dimensional (1D) arrays (e.g., pipelines, linear systolic processors [19]), 2D arrays (e.g., matrix processors, systolic arrays [19]), and 3D arrays (e.g., of hypercube architecture [16]).

At the physical level, very large scale integration circuits, for instance, are 3D devices because they have a layered structure (i.e., interconnection between layers while each layer has a 2D layout). On the way to the top of VLSI hierarchy (the most complex VLSI systems). [20] linear and 2D arrays eventually evolved to multi-unit architectures such as 3D arrays [21]. Their properties can be summarized as follows:

1. As stated in the theory of parallel and distributed computing, processing units are packed together and communicate best only with their nearest neighbors.
2. In the optimal organization, each processing unit will have a diameter comparable to the maximum distance a signal can travel within the time required for some reasonable minimal amount of processing of a signal, for example to determine how it should be routed or whether it should be processed further.
3. 3D architectures need to contain a fair number of cells before the advantages of the multi-cell organization become significant compared with competing topologies.

In current multiprocessor designs, 3D structures are not favored, since they suffer from gate and wire delay. Internally to each processing unit, 2D architecture is preferable, since at that smaller scale, communication times will be short compared with the cost of computation.

Nanostructures are associated with a molecular/atomic physical platform [3, 5]. This has a truly 3D structure instead of the 3D layout of silicon integrated circuits composed of 2D layers with interconnections forming the third dimension. At the nanoscale level, the advantages of the 3D architecture will begin to become apparent for the following reasons:

1. The distance light can travel in 1 ns in a vacuum is only around 30 cm and two times less in a solid (for example. 1 ns is the cycle of a computer with clock speed of 1 GHz), which means that components of such a computer must be packed in a single chip of several centimeter size; thus, a reasonable number of 3D array elements of nanometer size must be integrated on a single tiny chip [22–24].
2. They are desirable for their ideal nature for large scale computing.
3. There are many 3D algorithms and designs for existing microscale components that are arranged in 3D space, which computer designers already have experience with [25, 26].
4. There are limits to information density as well that imply a direct limit on the size of, in particular, a one-cycle-latency random access memory [22].

Thus, the speed of light limit (that is information transfer speed limit) and information density pose the following implications for computer architecture:

- Traditional architecture will be highly inefficient because most of the processor and memory will not be accessed at any given time due to the size limit,
- Interconnection topology must be scalable, most of the existing multiprocessor topologies (hypercubes, binary trees) do not scale, since communication times start to dominate as the machine sizes are scaled up.

The solutions to that are

1. Use of a parallel, multiprocessing architecture, where each processor is associated with some local memory, in contrast to the traditional von Neumann architecture.
2. The number of processing nodes reachable in $n$ hops from any node cannot grow faster than order $n^3$ and still embed the network in 3D space with a constant time per hop [27].

This leads us to the conclusion that there is only one class of network topologies that is asymptotically optimal as the machine size is scaled up: namely, a 3D structure, where each processing element is connected to a constant number of physical neighbors [28]. In fact the processing elements of this 3D mesh can simulate both the processors and wires of the alternative architecture, such as, in particular, randomly connected network of logic gates. The processing elements must be spread through the structure at a constant density.

Therefore, as processor speeds increase, the speed-of-light limit will cause communication distances to shrink, and the idea of mesh-connected processing elements and memory is, perhaps, the most reasonable and feasible solution.

### 2.2.1. Data Structures and Circuit Topology

Data structure plays a crucial role in logic design of nanoICs. We adopt certain methods of advanced logic design including techniques for representation and manipulation of different data structures (algebraic, matrix, decision trees, etc. [29–31]. These methods are selected based on the following criteria:

1. Graph-based models suitable for embedding and manipulation.
2. Technique for massive parallel computing.
3. Testability and observability.
4. Fault tolerance and reliable computing.

There are a number of particular characteristics of representing logic functions in nanospace:

1. The logic functions have to be represented by a spatial data structure in which information about the function satisfies the requirements of the implementation technology.
2. An information flow has to comply with the implementation topology.
3. This data structure and information flow must be effectively embedded into the 3D topological structure.

Therefore, designing architectures for computing logic functions supposes a resolution to the problem of finding the appropriate data structure and topology while taking into account computational and implementation aspects.

In this survey, we focus on hypercube topologies as the most suitable data structures for representation of a logic function, for the following reasons:

1. They are 3D, and thus meet the requirements of distributed spatial topology.
2. They correspond to functional (Shannon expansion) and dataflow organization (information relation of variables and function values) requirements, since data structures such as decision trees can be embedded into.
3. They meet the requirement of certain nanotechnologies with local quantum effects [32–34] and charge-state logic [10].

In Fig. 1 two hypercubes are represented: the left is referred to as a classical hypercube, the right is referred to as a hypercube-like topology called $V$-hypercube. Both hypercubes represent the same switching function $f = \bar{x}_1 x_2 \vee x_1 \bar{x}_2 \vee x_1 x_2 x_3$ but in different ways.

011   111

011    111

010    110

010    110

001    101

001    101

000    100

000    100

**Input f**

(a)            (b)

0110    1110

0100    1100

0111    1111

0101

010    110

010   011   111

001

000

100

0010

0000    1000

0001   1001   1011

000   100   101

1010

*Input*

(c)            (d)

**Figure 1.** Representation of a switching function $f = \bar{x}_1 x_2 \vee x_1 \bar{x}_2 \vee x_1 x_2 x_3$ by the classical hypercube (a) and $\Lambda$-hypercube (b) in 3D and 4D (c, d).

## 2.2.2. Assembling

The assembly philosophy of nanodevice design in spatial dimensions differs significantly from the usual ideas of building complex computer systems:

*Assembly* means the construction of more complex systems from the components provided. in particular, with features identical to the components which began the process.

*Self-assembly* is the process of construction of a unity from components acting under forces/motives internal or local to the components themselves, and arising through their interaction with the environment. Self-assembling structures create their own representations of the information they receive. A self-assembling system is able to process noisy, distorted, incomplete or imprecise data.

*Self-organizing.* The organism, for example, is a self-organizing system. In the organism. self-organizing is implemented through the local physical and chemical interactions of the individual elements themselves.

*Adaptive self-assembling* is the ability of a structure to learn how to perform assembling (appropriate architecture) aiming to solve certain tasks by being presented with examples.

## 2.2.3. Computer-Aided Design of Nanodevices

The goal of CAD of switching circuits is to automatically transform a description of circuits in the algorithmic or behavioral domains to one in the physical domain (i.e., down to a layout mask for chip production). In today's electronics this process is divided into

*System level* (major blocks of information processing)

*Behavioral level* (information flows)

*Logic level* (the behavior of the circuit is described by switching functions)

*Layout level* (mapping of logic network to physical layout topology).

Today. design tends to one-pass synthesis from behavioral description down to layout. and the most popular data structure for switching functions is the decision diagram. A CAD system also implements verification (i.e., verifying if the circuits). as a results of the complex design process. are logically equivalent to the initial description in the form of logic

equations, networks Usually. formal verification techniques deal with various data structures and descriptions.

It is expected tha in nanotechnology, behavioral. logic. and sometimes, layout levels of design will be evenually merged [35]. The efficiency of the design algorithms applied in these levels dependslargely on the chosen data structure. An efficient representation of logic functions is of fundamental importance for CAD of electronics and nanodevices. For example, in deep-submicon technology, which precedes nanotechnology, some levels of design are merged, and deision diagrams are the integrated data structure in this unified design process, called one-rass synthesis [36].

In the design of ranodevices, which are of 3D nature. each of the facets of CAD should employ spatial topoogical structures. Thus, these CAD tools must do manipulation, transformation, measurenents in spatial dimensions and design of topological structures. The corresponding CADalgorithms must satisfy a number of requirements: scalability (i.e., algorithms for constructor and manipulation of nanoscale topologies must be scalable for the size of the circuit tht can be processed), suitability for nanotechnology (i.e., implementation at a minimal cost);and parallelism and recurrency of calculations because nanostructures are distributed archtectures.

## 2.3. Distributed and Parallel Processing Paradigm

Massive and parallel computation of logic functions can be accomplished in nanodimensional structures viaword-level representation and also through borrowing some approaches developed in the theory of parallel computing on the "macroscale."

It should be noted that parallel and distributed computation on arrays (on the macrolevel) has been well-studed for example, systolic arrays of cells [19] and programmable-logic devices (PLD, FPCA. Most of them are 2D, however, 3D ones have been proposed as the most cutting-ecge models, for instance, hypercube-configured networks of computers, and hypercube supercomputers. Three-dimensional architectures have been used at the macrolevel, the level cf computer systems and networks, for a long time, especially in the area of network conmunication (communication hypercubes) [37, 38], and parallel and distributed algorithmsimplemented in supercomputers [25, 26].

Hypercube-like topologies inherit high parallelism of computing due to their

1. Regular and homogeneous structures.
2. Local connectvit/.
3. Ability to assemble and extend the structure.

The hypercubes depicted in Fig. 1 are called 3D hypercubes. They can carry limited information because the number of nodes and links is limited. The unique property of hypercubes is the possibilities for extension. This extension is expressed by notation of dimension (i.e., by the 3D nature of the physical world, we design many-dimensional hypercubes). For example, the four-dimensional (4D) classical hypercube is designed by multiple copies of the 3D hypercube illustrated n Fig. 1(c). Hypercube-like topology inherits this property from the hypercube: the 4D $N$-hypercube is designed by connections of root nodes (Fig. 1[d]).

On this basis, assenbly of complex multidimensional structures can be accomplished. The natural parallelism of the $N$-hypercube can be increased by the appropriate data structure of logic functions, socalled word-level representation. In this case each node of the $N$-hypercube computes a set of logic functions. Hence, the hypercube topology is very flexible in this extension.

The term cellular array refers to networks composed of some regular interconnection of logic cells. These arrays may be either 1D. 2D. or theoretically of any higher dimension than two. Practical considerations usually constrain them to 1D and 2D cases. Cellular automata are a reasonable model for the study of self-assembling and self-reproducing phenomena [39]. Further deployment of this concept led to cellular neural networks which can be possibly implemented on nanelectronic basis [40].

Systolic array is another name for parallel-pipelined computing structures [19]. In these structures, data inpu and output are organized in a sequential or partially parallel way,

and processing is accomplished by parallel computing on the array of the unified processing elements. The topology is usually linear (e.g.. for matrix-vector multiplication) or 2D (for matrix-matrix multiplication). Locality of data input/output and pipelining computing makes this organization of data processing attractive to implementation on nanoarrays [41].

### 2.3.1. Relevance to the Hierarchical FPGA

A conventional FPGA consists of an array of logic blocks that can be connected by routing blocks. Logic blocks are grouped into clusters which are recursively grouped together. With some simplification, the FPGA can be represented by the multirooted $k$-ary decision tree, while the routing, or switch blocks consist of wire segments and switches which can be configured to connect wire segments and logic blocks into a network. Array-based programmable logic for implementation of switching functions can be scaled to nanoelectronics (see, for example [14]).

EXAMPLE 1. *Hierarchical, or binary-treelike FPGA is based on single-input two-output switches. A $2 \times 2$ cluster of logic blocks can be connected using switches, and four copies of the cluster are organized into a "macro" cluster. The structure of this FPGA is represented by a binary decision tree of depth 4 in which the root and levels correspond to switching blocks and the 16 terminal nodes correspond to logic blocks (Fig. 2, where ■ denotes a logic block and o denotes a switch block). This tree is embedded in a hypercube-like structure of two dimensions.*

Hierarchical FPGA is based on a cluster of logic blocks connected by single-input four-output switch blocks [20]. Possible configurations can be described by a complete binary tree or a multirooted $k$-ary tree. This tree can be embedded into the hypercube-like structures [31].

EXAMPLE 2. *In Fig. 2 two topologies of FPGA are illustrated where four copies of the cluster are organized into a "macro cluster" of different configurations. The first topology (Fig. 2a) is described by a complete binary tree. The spatial structure includes two 3D hypercubes. The second topology (Fig. 2b) is described by a complete quaternary tree and also a hypercube-like structure in three dimensions.*

### 2.4. Fault-Tolerant Computing

Traditional logic circuits are very sensitive to failure: if even a single gate or single wire malfunctions, then the computation may be completely wrong. If one is worried about the physical possibility of such failures, then it is desirable to design circuits that are more resilient (i.e., design reliable circuits from unreliable elements [42, 43]).



(a)



(b)

Figure 2. Topological representations of hierarchical FPGAs by 2D structure, tree. and hypercube-like structure (examples 1 and 2).

Any computer with nanoscale components will contain a significant number of defects, as well as massive numbers of wires and switches for communication purposes. It therefore makes sense to consider architectural issues and defect tolerance early in the development of a new paradigm [12, 44].

An incorrect result is defined as a fault. In the presence of faults, a fault-tolerant system reconfigures itself to exclude the faulty elements from the system. A system so reconfigured may or may not change its topology. Ideally, a fault-tolerant system retains the same topology after faults arise.

Two main aspects are critical in the design of nanodevices:

*The high defect rates* at nanoscale; this means that defects of fabrication cause the distortion of logic correctness.

*The stochastic nature of computing* in nanodevices (electrons, molecules); this means that the calculated switching function is valid with some probability.

This problem is tackled from two directions:

*The first direction:* to compensate the technological defects, redundancy, through additional hardware resources or resources at time of computing, is introduced)

*The second direction:* the designed architecture must be fault-tolerant.

The first direction, in addition to double or triple hardware, may also take advantage of additional, or repetitive computations. The latter is relevant to stochastic computing [45, 46]. The methods of stochastic computing provide another approach to overcoming the problem of design of reliable computers from unreliable elements, which nanodevices are. For example, a signal is represented by the probability that a logic level is 1 or 0 at a clock pulse. In this way, random noise is being deliberately introduced into the data. A quantity is represented by a clocked sequence of logic levels generated by a random process. Operations are performed via the completely random data.

The second direction, fault-tolerant computing architecture, is designed in order to detect and correct errors. If certain nanodevices or parts of the network are destroyed, it will continue to function properly. The input, internal and output data can be noisy [47, 48], distorted or incomplete, and also the physical degradation of the system itself. When damage becomes extensive, it must only affect the system's performance, as opposed to a complete failure. Self-assembling nanosystems must be capable of this type of fault tolerance because they store information in a distributed (redundant) manner, in contrast to traditional storage of data in a specific memory location in which data can be lost in case of the hardware fault. An example of reconfigurable architecture is a massively parallel computer "Teramac," built at Hewlett-Packard Laboratories [44].

## 2.5. Information-Theoretical Measures

In the previous section, we already considered the benefit of 3D packing instead of 2D layout because three-dimensional structures are the only way to reduce propagation delays (as information propagation speed is limited to the speed of light). There are other fundamentally important characteristics of the computations in nanodimensions [49]:

1. The density of information to be packed into a cube, or space.
2. The bandwidth of information that can be transferred within the system per unit space.
3. The number of bit operations per second in the system.

These characteristics are limited by physical constants such as speed of light (already discussed in the sections above) and Boltzmann's constant and can be evaluated in terms of physical entropy and information.

Generally, the amount of information in a system is the logarithm of the number of this system's distinguished states. In a binary system, the information is measured in bits, and the density is the number of bits per volume.

EXAMPLE 3. *An empirically estimated information density in a cubic angstrom,* $\overset{\circ}{A}$ $(1 \ \overset{\circ}{A} = 10^{-8}$ m; $1 \ \overset{\circ}{A}^3$ *is roughly a hydrogen-atom-sized volume) for matter at normal pressure and temperature is about* 1 *bit* [49].

Information density is dependent on the material's atomic size, pressure and temperature.

EXAMPLE 4.   *Copper is capable of the storage of 6 bits/atom, and 0.5 b/Å³ at room temperature; this number increases to 1.5 b/Å³ at its boiling temperature (1365 K) [49]. This affects physical size and, thus, the propagation delay across memories and processors.*

The limit of scaling of operating frequencies is relevant to the maximum level of computational temperatures, which is higher than thermal temperature, and is the total energy per bit of information, that is, overall physical clock speed. For example, at room temperature (300 K), the thermal energies of individual degrees of freedom are only about 26 meV, which corresponds to a maximum frequency of 12.5 THz, according to the Margolus-Levitin theorem [23]. Note that if a molecular energy barrier which is of order 1 eV, corresponding to a computational temperature of 11,600 K, can be achieved, then a frequency of about 500 THz might be reached [50].

In a binary system, physical, or thermodynamic, entropy becomes Shannon entropy [51].

In terms of logic design of nanodevices, where symbolic and graph-based models of calculation are employed, Shannon information theory is useful because

1. Signal streams, or information, transferred within the system at logic level can be measured in terms of information, or entropy.
2. In the coherence of information measures at a physical level (thermodynamic), Shannon theory can be useful in evaluation of the characteristics of nanodevices, such as the density of information to be packed into a 3D nanospace, the number of bit operations per second.
3. Processing of information by nanodevices is described by probabilistic and statistical methods. This is where the information theoretic approach is united with fault-tolerant logic design.
4. The information-theoretical methods have been reasonable heuristic approach to minimization of decision diagrams. This is especially important for technology-independent and some technology-dependent design where information flows, specified by decision diagrams embedded into 3D space, are adequate to physical information flow, specified by 3D topologies.

## 3. DATA STRUCTURES FOR LOGIC DESIGN

Regardless of the choice of technology, unless computing principles are radically different (such as in the case of reversible computing [11]), the selected methods of classical logic design, with certain revision, offer appropriate forms of data representation and methods of data manipulation, which comply with the particular properties of 3D implementation. Selected methods of logic design include:

1. Algebraic-, matrix- and hypercube-based methods of computing a switching function.
2. Decision trees and decision diagram design, as a basis for 3D embedding techniques [31].
3. Event-driven analysis of dynamic changes in logic [52], which is relevant to testability and fault-tolerant computing.

Data structure is the term used to define an abstract data type. Data structure for a switching function must show the dependance of the function on its variables. In other words, it is a mathematical model of a switching function. Data structures for switching functions are the following:

1. Algebraic (sum-of-products, Reed-Muller, arithmetic, and word-level expressions).
2. Tabular representation, or binary arrays, such as truth tables, or look-up-tables (LUTs).
3. Graph-based representation, including cube, flowgraph, decision tree, decision diagram, and switching gate-level description.

The gate-level description is a classical graph-based representation of a switching function as a random network of interconnected gates called a circuit. A collection of wires that always carry the same electrical signal is called a switching net. The tabulation of gate inputs and outputs and the nets to which they are connected is called the netlist. The fan-in of a gate

is its number of inputs. The fan-out of a gate is the number of inputs to which the gate's output is connected.

The functional descriptions of multi-input multi-output switching functions, employed intensively in today's logic design, are decision trees and diagrams. Fundamentals of decision diagram technique can be found in [53–57]. Useful details are considered in [58–60]. In [61], a linear decision diagrams are developed. In [32–34], decision diagrams were used in single-electron circuit design.

## 3.1. Decision Tee

A decision tree is a rooted acyclic graph in which every vertex but the root has an indegree of 1. The decisiontree is characterized by

1. The *size*, the number of vertices.
2. The *depth*, the number of levels;
3. The *width*, the maximum number of vertices for a level.
4. The *area*, *depth* × *width*.

EXAMPLE 5. In Fig. 3, the complete ternary ($p = 3$) 3-level ($n = 3$) tree is given. The root corresponds to the evel (depth) 0. Its three children are associated with level 1 ($3^1 = 3$). Level 2 includes $3^2 = 9$ nodes. Finally, there are $3^3 = 27$ terminal nodes.

## 3.2. Binary Deision Diagrams

A decision tree can be reduced to a decision diagram. This is possible by applying certain reduction rules, for instance, if the decision tree contains any vertex $v$ whose successors lead to the same node and if it contains any distinct vertices $v$ and $v'$ such that the subgraphs rooted in $v$ and $v'$ are isomorphic. A binary decision diagram (BDD) is a directed acyclic graph that represents a switching function $f$ in the following way:

1. It has exactly one root and internal vertices, or nodes, are labeled by a Boolean variable $x_i$ and have exactly two outgoing edges, a 0-edge and 1-edge.
2. It has two leves, or terminal nodes, labeled by 0 and 1, which are the function value.
3. Each assignment to the input variables $x_i$ defines a uniquely determined path from the root of the graph to one of the terminal nodes which is the function value for this input assignment.

An ordered BDD (OBDD) is a rooted directed acyclic graph that represents a switching function. A linear variable order is placed on the input variables. The variables' occurrences on each path of this diagram have to be consistent with this order. An OBDD is called reduced if it does not contain any vertex $v$ such that the 0-edge and 1-edge of $v$ lead to the same node, and i does not contain any distinct vertices $v$ and $v'$ such that the subgraphs rooted in $v$ and $v'$ are isomorphic.

A decision diagram is characterized, similarly to a decision tree, by the size, depth, width, area, and the efficiency of reduction. Size of decision diagram strongly depends on variable order [62, 63].

EXAMPLE 6. Figure 4 illustrates the reduction of a decision tree to an OBDD of the switching function $f = x_1 x_2 \lor x_3$.



| | Size=15 |
|---|---|
| | Width=8 |
| | Depth=3 |
| | Indegree=1 |
| | Outdegree=2 |
| | Area=8 × 3=24 |
| | Terminal nodes=8 |
| | Intermedial nodes=7 |

Figure 3. The complete binary tree (example 5).

Figure 4. The complete (a) and reduced (b) decision tree for the switching function $f = x_1 x_2 \vee x_3$; the shared ROBDD for the switching function $f_1 = x_1 x_2 \vee x_3$, $f_2 = x_1 \vee x_2 \vee x_3$ (c) (examples 6 and 7).

A multi-output switching function is represented by a multirooted decision diagram, which is called a shared decision diagram.

EXAMPLE 7. *A two-output switching function* $f_1 = x_1 x_2 \vee x_3$, *and* $f_2 = x_1 \vee x_2 \vee x_3$ *is represented by a shared OBDD given in Figure* 4(c).

In a multiplexer tree (network), each internal tree node is represented as a 2-to-1 multiplexer controlled by the node variable and each terminal node is implemented as a constant logical value (wired at 0 or wired at 1); the interconnection scheme is that of the decision tree (diagram). The evaluation of a function then proceeds from the terminal nodes (the constant values) to the root multiplexer, the function variables, used as control variables, select a unique path from the root to one terminal node, and the value assigned to that terminal node propagates along the path to the output of the root multiplexer.

In [31], decision trees and diagrams for representation of switching functions are related to the spatial structures, reflecting the requirement of distributed three-dimensional topologies assumed in nanodimensions.

It should be noted that generalization of switching algebra, multiple-valued logic that is a theoretical framework for design of multivalued logic circuits, which are considered as a resolution to information density problem (as multivalued signal possess larger information capacity). The corresponding data structures for multi-valued functions representation, including decision diagrams and embeddings in hypercube, are considered in [31, 52].

## 3.3. Other Spatial Structures

Spatial structures have been used in distributed and parallel network design at the macro-level for a long time. However, most of the known topologies for massive parallel computing have not been considered relevant to spatial logic design.

### 3.3.1. Requirement for Representation in Spatial Dimensions

The following characteristics are critical to a spatial computing network topology:

1. Good embedding capabilities.
2. Minimal degree.
3. Ability to extend the size of structure with minimal changes to the existing configuration.
4. Ability to increase reliability and fault tolerance with minimal changes to the existing configuration.
5. Flexibility of design methods.
6. Flexibility of technology.

On the basis of this criteria. a number of topologies can be considered relevant to the problems of spatial logic design, for example, hypercube topology [64], cube-connected cycles known as $CCC$-topology [65], and pyramid topology [25].

EXAMPLE 8. *Examples of criteria for choosing a topology:*

1. *If the distance (diameter) is small, then computing elements are likely to able to communicate more quickly.*

2. *It is desirable that all pairs of computing elements communicate with equal ease or at least that traffic patterns between all computing elements be reasonably balanced.*

3. *If the network can be efficiently embedded into 2D or 3D space such that all the wires are relatively short, then information can propagate quickly between computing elements.*

4. *If there many possible paths between each pair of computing elements, a partially defective network may continue to function.*

The singular binary $n$-hypercube is a special case of the family of $k$-ary $n$-hypercubes, which are hypercubes with $n$ dimensions and $k$ nodes in each dimension. The total number of nodes in such a hypercube is $N = k^n$.

### 3.3.2. The CCC-Hypercube

Figure 5(b) is created from a hypercube by replacing each node with a cycle of $s$ nodes. It hence increases the total number of nodes from $2^n$ to $s \cdot 2^n$ and preserves all features of the hypercube [65]. The CCC-hypercube is closely relevant to the butterfly network. It should be noted that "butterfly" flowgraphs are the nature of most transforms of switching functions in matrix form [5].

### 3.3.3. Pyramid Topology

Figure 5(c) is suitable for many computations based on the principle of hierarchical control, for example, of nonbinary decision trees and decision diagrams. An arbitrarily large pyramid can be embedded into the hypercube with a minimal load factor, and it is flexible for extension. Pyramid topology is relevant also to fractal-based computation, which is effective for symmetric functions and is used in digital signal processing, image processing, and pattern recognition [25].

### 3.3.4. Hypercube Topology

Figure 5(a) has received considerable attention in classical logic design due mainly to its ability to interpret logic formulas and logic computation (small diameter, regularity, high connectivity, symmetry) [64]. Hypercube-based structures are at the forefront of massive parallel computation because of the unique characteristics of hypercubes (fault tolerance, ability to efficiently permit the embedding of various topologies, such as lattices and trees) [37, 38].

A hypercube is an extension of a graph. The dimensions are specified by the set $\{0, 1, \ldots, n - 1\}$. An $n$-dimensional binary hypercube is a network with $N = 2^n$ nodes and diameter $n$. There are $d \times 2^{d-1}$ edges in a hypercube of $d$ dimensions. Hypercube $Q_n$ can be defined recursively as the graph product.

### 3.3.5. Gray Code

Gray code is used for encoding the indexes of the nodes in graphs. There are several reasons to encode the indexes. The most important of them is to simplify analysis, synthesis and embedding of topological structures. Gray code is referred to as unite-distance codes [66].



Figure 5. Spatial configurations: hypercube (a), CCC-hypercube (b), and pyramid (c).

Let $b_n...b_1b_0$ be a binary representation of an integer positive number $B$ and $g_n...g_1g_0$ be its Gray code. Figure 6 illustrates the above transformation for $n = 3$.

EXAMPLE 9. *Binary to Gray and Gray to binary transformation is illustrated by Fig. 6 for $n = 3$.*

In the hypercube, two nodes are connected by a link (edge) if their Hamming distance is equal to 1, that is, if they have labels that differ by exactly one bit. The number of bits by which labels $g_i$ and $g_j$ differ is denoted by $h(g_i, g_j)$; this is the Hamming distance, or sum, between the nodes. The Hamming sum is defined as the bitwise operation $(g_{d-1}...g_0) \oplus (g'_{d-1}...g'_0) = (g_{d-1} \oplus g'_{d-1}),...,(g_1 \oplus g'_1), (g_0 \oplus g'_0)$, where $\oplus$ is an exclusive or operation.

The hypercubes can be flexibly extended to any size, or number of dimensions, that corresponds to the number of variables $n$ of a switching function. This is accomplished through assembling.

## 3.3.6. Assembling

Assembling is the basic topological operation that we apply to the synthesis of hypercube and hypercube-like data structure. Assembling is the first phase of the development of self-assembling, that is, the process of construction of a unity from components acting under forces internal or local to the components themselves.

The assembly procedure is applied once the following data is specified:

1. The structural topological components.
2. Formal interpretation of the structural topological components in terms of the problem.
3. The rules of assembly.

The assembling is a key philosophy of building complex systems. For example, assembling a circuit after configuration.

## 3.4. Embedding of a Guest Graph into a Host Graph

An embedding, $\langle \varphi, \alpha \rangle$, of a graph $G$ into a graph $H$ is a one-to-one mapping $\varphi: V(G) \rightarrow V(H)$, along with a mapping $\alpha$ that maps an edge $(u, v) \in E(G)$ to a path between $\varphi(u)$ and $\varphi(v)$ in $H$ [?]. The embedding of a guest graph $G$ into a host graph $H$ is a one-to-one mapping of the vertices in $G$ to the vertices in $H$. An embedding is characterized by a set

| Binary code | Gray code | | Binary code | Gray code |
|---|---|---|---|---|
| 000 | 000 | | 000 | 000 |
| 001 | 001 | | 001 | 001 |
| 010 | 011 | | 011 | 010 |
| 011 | 010 | | 010 | 011 |
| 100 | 110 | | 110 | 100 |
| 101 | 111 | | 111 | 101 |
| 110 | 101 | | 101 | 110 |
| 111 | 100 | | 100 | 111 |

Suppose that $B = b_n...b_1b_0$ is given, then the corresponding binary Gray code representation $g_i = b_i \oplus b_{i+1}$. Given Gray code $G = g_n...g_1g_0$, the corresponding binary representation is derived by

$$b_i = g_0 \oplus g_1 \oplus \cdots g_n = \bigoplus_{i=0}^{n} g_i$$



Figure 6. Flowgraph and formal equation for binary to Gray code (a) and inverse transformations (b) (example 9)

of parameters:

1. The *expansion* is the ratio $|V(H)|/|V(G)|$.
2. The *dilation* cost of an embedding of $G$ into $H$ is the maximum distance in $H$ between any two neighboring vertices in $G$.
3. The *congestion* of the embedding is the maximum of the congestions (parallel edges) of all edges of $H$.

EXAMPLE 10. *Details of embedding graph $G$ into a host graph $H$ are given in Fig. 7.*

## 3.5. Hypercubes in Logic Design

Hypercubes are used in classical logic design to interpret logic formulas and manipulation of them. A hypercube used to represent a switching function is called a singular hypercube [64]. In the singular hypercube, each node is labeled by a product of the canonical SOP of a switching function. The topology and the labels of the neighboring nodes are specified through the assembly rules.

Assembling a singular hypercube of switching functions is accomplished by

1. Generating the products as enumerated points (nodes) in the plane.
2. Encoding the nodes by Gray code.
3. Generating links using Hamming distance.
4. Assembling the nodes and links.
5. Joining a topology of a hypercube in $n$ dimensions.

Let $x_j^i$ be a literal of a Boolean variable $x_j$ such that $x_j^0 = \bar{x}_j$, and $x_j^1 = x_j$, and $x_1^{i_1} x_2^{i_2} \ldots x_n^{i_n}$ is a product of literals. Topologically, it is a set of points on the plane enumerated by $i = 0, 1, \ldots, n$. To map this set into the hypercube, the numbers must be encoded by Gray code and represented by the corresponding graphs based on Hamming distance. The example below demonstrates the assembly procedure.

EXAMPLE 11. *Figure 8(a) demonstrates the assembly of hypercubes.*

The 0-dimensional hypercube ($n = 0$) represents constant 0 or 1. The 1D hypercube includes the line segment connecting vertices 0 and 1, and this segment is called the face and denoted by x. A 2D hypercube has four faces, 0x, 1x, x0, and x1. The total 2D hypercube can be denoted by xx.

EXAMPLE 12. *Six faces of the hypercube, xx0, xx1, 0xx, 1xx, x1x, and x0x (Fig. 9) represent 1-term products for a switching function of three variables.*

Rings, meshes, pyramids, shuffle-exchange networks, and complete binary trees can be embedded into hypercubes. The latter is of particular interest for logic design in nanodimensions.

## 3.6. $N$-Hypercube Definition

In this section the extension of the traditional hypercube is considered. This extension, introduced in [31], is called the $N$-hypercube.

The key to the extension of a hypercube to an $n$-dimensional $N$-hypercube is embedding a complete binary tree of an $n$-variable switching function into a hypercube.

The vertex mapping:
$1 \rightarrow 1, 2 \rightarrow 2,$
$3 \rightarrow 3, 4 \rightarrow 4$
The edge to path mapping:
$(1,2) \rightarrow 1, 2; (2,4) \rightarrow 2,4;$
$(3,4) \rightarrow 3, 4; (1,3) \rightarrow 1,3;$
Expansion = 2
Dilation = 1
Congestion = 1

*Graph G*          *Graph H*

Figure 7. Embedding graph $G$ into a host graph $H$ (example 10).

1D: Product $x_1^{c_1}$, 2 points $(n = 1)$

2D: Product $x_1^{c_1}x_2^{c_2}$, 4 points $(n = 2)$

3D: Product $x_1^{c_1}x_2^{c_2}x_3^{c_3}$, 8 points $(n = 3)$

4D: Product $x_1^{c_1}x_2^{c_2}x_3^{c_3}x_4^{c_4}$, 16 points $(n = 4)$

5D: Product $x_1^{c_1}x_2^{c_2}x_3^{c_3}x_4^{c_4}x_5^{c_5}$, 32 points $(n = 5)$

Figure 8. Assembling a hypercube for representation of product terms (example 11).

### 3.6.1. Embedding a Binary Decision Tree into an $N$-Hypercube

A binary decision tree embedded into an $N$-Hypercube is accomplished as follows (Fig. 10):

1. The nodes of the $n$-dimensional singular hypercube are replaced with the $2^n$ terminal nodes of the decision tree.
2. *Intermediate* nodes are embedded into the edges and faces of a singular hypercube.
3. The *root* node is embedded into the center of the singular hypercube.
4. The edges of the binary decision tree correspond to the edges of the new $N$-hypercube.

For example, the terminal node of the complete binary decision tree with $q$ levels can be embedded into a hypercube with $2^q$ vertices and $q \times 2^{q-1}$ edges. This is because the complete binary decision tree with $q$ levels has $2^q$ leaves. This is exactly the number of nodes in the hypercube structure, where each node is connected to $q - 1$ neighbors and assigned the $q$-bit binary code that satisfies the Hamming encoding rule and, thus, has $q \times 2^{q-1}$ edges.

The number of vertices of the binary tree embedded into the middle of each edge of the $N$-hypercube is equal to $2^{q-1}$, whereas the possible number of such embeddings (the number of all wires) is $q \times 2^{q-1}$. The number of inner nodes of the binary tree embedded into the middle of each edge of the $N$-hypercube is equal to $2^{q-2}$ while the possible number of such embeddings (the number of all edges) is $q \times 2^{q-2}$.

The total number of nodes and the total number of edges (connections) between nodes in the $n$-dimensional $N$-hypercube is specified as below

$$N_d = \sum_{i=0}^{n} 2^{n-i} C_i^n \quad \text{and} \quad N_c = \sum_{i=0}^{n} 2i \cdot 2^{n-i} C_i^n$$

The total number of internal nodes is equal to the number of all nodes except leaves in the complete binary decision tree that represents a switching function of $n$ variables, and the total number of edges is equal to the number of edges in the complete binary decision tree.



$xx0: \bar{x}_3(\bar{x}_1\bar{x}_2 \vee \bar{x}_1x_2 \vee x_1\bar{x}_2 \vee x_1x_2)$

$xx1: x_3(\bar{x}_1\bar{x}_2 \vee \bar{x}_1x_2 \vee x_1\bar{x}_2 \vee x_1x_2)$

$0xx: \bar{x}_1(\bar{x}_2\bar{x}_3 \vee \bar{x}_2x_3 \vee x_2\bar{x}_3 \vee x_2x_3)$

$1xx: x_1(\bar{x}_2\bar{x}_3 \vee \bar{x}_2x_3 \vee x_2\bar{x}_3 \vee x_2x_3)$

$x0x: \bar{x}_2(\bar{x}_1\bar{x}_3 \vee \bar{x}_1x_3 \vee x_1\bar{x}_3 \vee x_1x_3)$

$x1x: x_2(\bar{x}_1\bar{x}_3 \vee \bar{x}_1x_3 \vee x_1\bar{x}_3 \vee x_1x_3)$

Figure 9. Faces of the hypercube interpretation in the sum-of-products of a switching function of three variables (example 12).

**Figure 10.** Correspondence of the attributes of a binary tree and an X-hypercube.

EXAMPLE 13.  *Given a switching function of one variable, $f = \bar{x}$, the corresponding binary decision tree has $q = 1$ level (Fig. 11). The function is equal to 1 while $x = 0$ and equal to 0 while $x = 1$. These values assign two leaves of the binary decision diagram. The hypercube of the function is a two-node graph, and the corresponding X-hypercube has a root and two leaves. Given $q = 2$, the four leaves of the complete binary decision tree of the two-variable function $f = \bar{x}_1 x_2$ can be embedded into a X-hypercube with the root, two intermediate nodes and four leaves. The corresponding singular hypercube consists of four nodes.*

The embedding is as follows (Fig. 11[b]):

1.  Embed four leaves of the binary tree into a four-node singular hypercube; the codes 00, 01, 10, 11 of nodes are assigned in order the Hamming distance between the neighbor nodes is equal to one.

2.  Embed 2 inner nodes of the binary tree into edges connecting the existing nodes of the singular hypercube; note that two of the edge-embedded nodes must be considered at a time; the first hypercube in Fig. 11(c) corresponds to the order $x_2$, $x_1$, and the second one describes the order $x_1$, $x_2$; the axes are associated with the polarity of variables (complemented, uncomplemented) and explain the meaning of the bold edges.

3.  Embed the root of the tree into the center of the facet of the hypercube and connect it to the edge-embedded nodes.

EXAMPLE 14.  *Given a complete binary decision tree of a switching function of three variables, $q = 3$, the eight leaf nodes of the tree are embedded into the 3D singular hypercube, which is transformed to a 3D X-hypercube (Fig. 12).*

Additional embedded nodes and links assignments correspond to embedding decision trees in a hypercube and thus, convert a hypercube from the passive representation of a function to a connection-based structure (i.e., a structure in which calculations can be accomplished). In other words, information connectivity is introduced into the hypercube.

Because the intermediate nodes of the decision tree perform Shannon expansion, they are associated with demultiplexors. Terminal nodes carry information about the results of



**Figure 11.** Embedding a binary decision tree into a 2-D X-hypercube (example 13).

Figure 12. Embedding the complete binary tree into the 3D .\-hypercube (example 14).

computing. The nodes (their functions and coordinates), and edges not only carry information about the function implemented by the .\-hypercube but allow to do all the manipulations and calculation the decision tree allows.

### 3.6.2. Degree of Freedom and Rotation

Each intermediate node in the .\-hypercube is associated with a so-called degree of freedom. This degree of freedom can be used for variable order manipulation, as the order of variables is a parameter to adjust in decision trees and diagrams.

The term "degree of freedom" is associated with the order of variables in decomposition, and hence, to the order of variables in decision trees and diagrams. The polarity of variables influences the above characteristics too.

Only intermediate and root nodes in an .\'-hypercube can be characterized by a degree of freedom. An intermediate node in the 1D .\'-hypercube has two degrees of freedom. The 2D .\-hypercube is assembled from two 1D .\'-hypercubes, and includes three intermediate nodes. The .\-hypercube in 2-D has $2 \times 2 \times 2 = 8$ degrees of freedom. There are four decision trees with different orders of variables.

EXAMPLE 15. *Consider an .\-hypercube in 3D. This .\-hypercube is assembled of two two-dimensional .\-hypercubes and includes seven intermediate nodes and has $8 \times 8 \times 2 = 128$ degrees of freedom. The degree of freedom of an intermediate node at i-th dimension. $i = 2, 3, \ldots, n$, is equal to $DF_i = 2^{n-i} + 1$. In general, the degree of freedom of the n-dimensional .\-hypercube is defined as $\sum_i DF_i = \sum_i (2^{n-i} + 1)$.*

### 3.7. .\'-Hypercube Design for $n > 3$ Dimensions

To design a 4D .\'-hypercube, two .\-hypercubes must be joined by links. The number of bits in the coordinate description of both .\'-hypercubes must be increased by one bit.

EXAMPLE 16. *Figure 13 shows the possibilities for assembling a given .\'-hypercube to 4D .\-hypercube, 4D to 5D, etc. This connection property follows from the properties of intermediate nodes.*



Figure 13. Connections between .\-hypercubes in n-dimensional space (example 16).

Figure 14. \-hypercube design by embedding a complete binary tree in the \-hypercube (example 17).

Summarizing the above characteristics, the \-hypercube can be specified as an undirected hypercube with the following properties:

1. \-hypercube corresponds to an $n$-level $2^n$-leaves complete binary tree.
2. $k = 2^n$ terminal nodes labelled from 0 to $2^n - 1$ so that there is an edge between any two vertices if and only if the binary representations of their labels differ by one and only one bit.
3. Each edge is assigned with an intermediate node which corresponds to binary representation of both edge-ending constant nodes with the do not care value for the only different bit.
4. There is an edge between two intermediate nodes if and only if the binary representations of their labels differ by one and only one bit.

EXAMPLE 17. Figure 14 illustrates the 3D \-hypercube assembling (assembling step by step from (a) to (d)).

# 4. OTHER SPATIAL MODELS

A multi-input multi-output switching function can be represented by

1. A gate-level model, or network of gates, described by a direct acyclic graph (DAG).
2. A functional level model, that is a two-level sum-of-products expression, decision tree or decision diagram.

A multilevel circuit is defined as a DAG which describes interconnections of single-output combinational logic gates under the assumption that the interconnection provides a unidirectional (no feedback) flow of signals from primary inputs to primary outputs. Multilevel network implementations can be restricted to a gate type (NAND, NOR) with a fixed fan-in.

## 4.1. DAG-Based Representation of Switching Circuits

The well-known techniques that are used in today's multilevel circuit design utilize DAG construction and optimization [29]. In design at the physical level, Boolean mapping is required, which means covering of the DAG by DAGs of elementary gates from the library of gates. A DAG is the way contemporary multilevel circuits are represented and it is a gate-level model that is compatible with a library of traditional gates.

## 4.2. Hybrid Representation of Switching Circuits

Decision diagrams correspond to multiplexer (MUX)-based implementation that is not widely used in today's logic design. The situation is quite opposite in the design of new

devices: the appropriate logic for nanowire wrap-gate single-electron devices is based on MUX-gate or T-gate implementation [33, 34]. This technique suffers from exponential complexity, which is reduced by transforming decision trees to decision diagrams. In nanodimensions, however, the regularity of the tree structures may hold more benefits than the compactness of decision diagrams.

A certain trade-off is to combine $N$-hypercube with DAG. This can be accomplished by embedding a DAG into a hypercube, or assembling primitive $N$-hypercubes into a DAG topology, that is, a random network of logic gates. This would require representation of each elementary gate by an $N$-hypercube [31].

## 4.3. Library of $N$-Hypercubes

In this section, the focus is generation of elementary $N$-hypercubes, and evaluation of elementary $N$-hypercube structures. We introduce the library of $N$-hypercubes that implement elementary switching and multivalued functions. This representation is based on the two-level form of a switching function of a gate, corresponding to a two-level tree, and characterization, analysis, and study of recombination (while the order of variables is changed).

### 4.3.1. Structure of the Library

In Boolean mapping, a library is understood as a set of logical elements. Design of a logic network over a given library of gates is accomplished by covering a DAG by DAGs of elementary gates from the library of gates. The library contains the set of logic gates that are available in the desired design style. Each element is characterized by its function, inputs, outputs, and some parameters such as area, delay, and capacity load.

### 4.3.2. Metrics of $N$-Hypercube

A primitive $N$-hypercube is a one-dimensional $N$-hypercube of (Fig. 15). This is a unit to be assembled into $N$-hypercube models of gates. The simplest topological characteristics of the $N$-hypercube gate include:

1. Number of terminal nodes. The number of terminal nodes for an $n$-input logic gate is $2^n$ (for example, the number of terminal nodes is equal to four and eight for the two-input and three-input logic gates, respectively). The number of active terminal nodes is equal to the number of 1's in the truth table of the implemented switching function, which is the current state of the $N$-hypercube structure. The active terminal nodes are used to evaluate power dissipation characteristics.

2. Number of intermediate nodes. The intermediate nodes correspond to the number of nodes in the levels of the decision tree of the switching function, except for terminal nodes. For instance, the number of intermediate nodes is equal to 2 plus one root node for a two-input gates, and 6 plus one root node for a three-input gate. This number can be used for evaluation of circuit complexity.

3. Connectivity is the number of links between root, terminal and intermediate nodes. This characteristic describes the complexity in the number of links. Given the coordinates



- A primitive $N$-hypercube is a one-dimensional $N$-hypercube (a root node and two terminal nodes).

- The primitive $N$-hypercube corresponds to the decision tree of a single variable (node).

- The primitive is a unit metric unit in the $N$-hypercube structure.

Figure 15. The primitive one-dimensional $N$-hypercube and the corresponding node of a decision tree.

of a link, it is possible to measure the sizes of links and compare them. Based on the notation of connectivity, the topological distributions of links are calculated, as are the distances from an arbitrary point to a set of points in space.

4. Diameter is the maximum number of links between any two nodes in the $N$-hypercube. It is the global characteristic of a circuit in space.

5. Dimension. The dimension can be calculated given the number of levels of the embedded decision tree, that is the number of variables in the corresponding function.

In Table 1, the topological characteristics of the $N$-hypercube model of two- and three-input $N$-hypercube gates are given.

### 4.3.3. Manipulation of $N$-Hypercube

Two attributes of a $N$-hypercube can be changed by reconfiguration: polarity of variables and the order of variables in a decision tree and decision diagram. The reconfiguration of the $N$-hypercube is defined as rotation.

Any standard gate of the conventional combinational logic (i.e., an $n$-input AND, OR, NOR, NAND, NOT, and EXOR gate), can be represented by the $N$-hypercube model in two steps:

1. A decision tree of the gate is derived.
2. The tree is embedded in an $N$-hypercube of $n$-dimensions.

An $N$-hypercube of a two-input gate includes the root node, two intermediate nodes, four terminal nodes, and the edges (connections).

EXAMPLE 18. *Figure 16 illustrates designing an $N$-hypercube model for the two-input NAND gate. Given the two-input NAND function $f = \overline{x_1 x_2}$ (a), the Shannon decision tree is derived (b). Next, this tree is embedded in an $N$-hypercube (c). The number of active terminal nodes in this 2D $N$-hypercube is 3, and the total number of terminal nodes is 4, while the number of intermediate nodes is 3, the connectivity is 6 and the diameter is 4.*

The $N$-hypercube in 2D, 3D, 4D, and 5D is denoted by a cube with two, three, four, and five inputs, respectively (Fig. 17).

### 4.4. Hybrid Design Paradigm: Embedding a DAG in $N$-Hypercube

An arbitrary switching function can be represented by an $N$-hypercube derived from the decision tree of a function. However, for a traditional approach based on covering the DAG of a circuit by the library of gates, netlist must be also considered. In terms of implementation on nanodevices (nanowire networks), the problem of interconnections arises.

We focus on two approaches for deriving 3D data structures for switching functions: a logic network (direct acyclic graph) presented, i.e., netlist is treated as a tree; and embedding the tree in an $N$-hypercube.

Deriving a 3D structure given a circuit netlist is implemented in two steps:

1. The circuit netlist is represented by a DAG and levelized to obtain an incomplete decision tree.
2. The tree is embedded in an $N$-hypercube.

Table 1. Metrics of 2-input and 3-input $N$-hypercube gates.

| Metrics | 2-input | 3-input |
|---|---|---|
| Number of terminal nodes | 4 | 8 |
| Number of intermediate nodes plus the root node | 3 | 7 |
| Total number of nodes | 7 | 15 |
| Connectivity | 6 | 14 |
| Dimension | 2 | 3 |
| Diameter | 4 | 6 |

**Figure 16.** Fragment of a library of $X$-hypercube gates: a two-input NAND gate (a), its decision tree (b), and the corresponding $X$-hypercube (c) (example 18).

Thus, there are two levels of embedding:

1. Embedding decision trees in elementary $X$-hypercubes of gates.
2. Embedding a tree, i.e., DAG, in a "macro" $X$-hypercube.

The technique of transformation of circuit models in spatial dimensions is based on algebraic simplifications of switching functions, topological simplifications, and logic-topological transformations.

Algebraic simplifications can modify the topology, and topological simplifications can change the switching function described by this topology. Therefore, algebraic and topological transformations must be carefully combined. In this section, the rules for manipulation of an $X$-hypercube are introduced. These rules allow reduction of dimensions and simplify circuit representation in spatial dimensions, and are also useful in verification of $X$-hypercubes.

*Dimension reduction.* If the input $i$ of an $X$-hypercube to implement an $n$-input OR function is equal to 0, reduce this input and replace this hypercube with the $(n - 1)$-dimensional $X$-hypercube. By analogy, the dimensions of the $X$-hypercube are decreased by 1 for an $n$-input AND function if the input is equal to 1.

*$X$-hypercube deleting.* If the input of an $X$-hypercube implementing an AND function is equal to 0, replace the hypercube with the constant 0. By analogy, replace an OR $X$-hypercube with constant 1, if one of the outputs is equal to 1.

*$X$-hypercube merging.* Two AND (OR) $X$-hypercubes of dimensions $i$ and $j$ correspondingly connected in a series can be merged into one AND (OR) $X$-hypercube of $i + j$ dimensions.

*Deleting of duplicated $X$-hypercubes.* If there are two $X$-hypercubes whose inputs and outputs are the same, remove one and create a fan-out.



**Figure 17.** Denotation of a multidimensional $X$-hypercube: (a) 2D, (b) 3D, and (c) 5D.

*Reduction of redundant connections.* Consider the manipulation of two $V$-hypercubes connected in a series, A and B, with respect to the input $x_j$:

$$f = \underbrace{\underbrace{(\overline{x_j x_2 x_3})}_{\text{Hypercube } A}\, x_j x_4 x_5}_{\text{Hypercube } B} = \underbrace{\underbrace{(\overline{x_2 x_3})}_{\text{Hypercube } A}\, x_j x_4 x_5}_{\text{Hypercube } B}$$

It follows from the above that the input $x_j$ can be deleted from $N$-hypercube A.

On the basis of properties of switching functions, the remaining rules for simplification and manipulation of $N$-hypercube topology structures can be derived.

EXAMPLE 19. *In Fig. 18(a), the initial circuit is cascaded into two subcircuits with outputs $f_1$ and $f_2$. Because of the gates can be considered as nodes, the subcircuits are equivalent to the two logic networks. The embedding of these networks results in two incomplete $N$-hypercubes over the library of 3D gates (Fig. 18[b]).*

Hence, instead of representation of this circuit by a complete decision tree, that is, a tree with $2^6$ terminal nodes, 6 levels, 62 intermediate nodes and a root, or by 6-dimensional $N$-hypercube, the circuit is described by two logic networks with the following characteristics:

1. The first network includes five terminal nodes and four intermediate nodes.
2. The second network includes four terminal nodes and three intermediate nodes

that correspond to two incomplete 3D $N$-hypercubes.



Figure 18. Cascading of a circuit (a) to represent the outputs $f_1$ and $f_2$ by incomplete $N$-hypercubes (b) (example 19).

## 4.5. Numerical Evaluation of 3D Structures

In this section, the results of numerical evaluation of $N$-hypercubes' combinational circuits from the ISCAS85 standard database are given. In experiments, circuits with from six gates to more than 3500 were used.

### 4.5.1. Experiment on Evaluating the $N$-hypercube

In this experiment, parameters of the $N$-hypercube derived by the decision diagrams were evaluated. We used the parameters acquired from the shared reduced ordered binary decision diagrams (BDDs) built from the circuit netlist.

The results of evaluation are summarized in Table 2 in which "Circuit" is the benchmark type, "I/O" is the number of inputs and outputs, "#G" is the number of gates in the circuit, "#Dim" is the maximum number of dimensions (which is the numbers of variables and can be less for each separate function represented by a BDD); next distribution of the dimensions in 3D that is the size of the solid in "X," "Y," and "Z" coordinates is given, and the "Nodes" includes the number of all intermediate nodes plus active terminal nodes in the shared BDD representing the multi-output functions; this number is too large for some circuits and is not shown in the table.

Note that the number of active terminal nodes is upperbounded by $2^n$ for an $n$-variable function.

### 4.5.2. Experiment on Evaluating the Hybrid Approach

In this experiment, the hybrid $N$-hypercube based approach was evaluated for a selected output in each test circuit. Each gate in the network was replaced by an $N$-hypercube model. An $N$-hypercube of each output has been derived from the tree obtained for each output by scanning from output to inputs, and next, by levelization.

The results of evaluation are summarized in Table 3, where "Circuits" is the benchmark type, "#O" is the selected output number (out of all the outputs), and "#G" is the number of gates in the subnetwork implementing the selected output. We have selected the output whose implementation involves the maximum number of gates, and considered the subnetwork that involves the inputs and gates to implement this function. The next three columns include distribution of the dimensions in "X," "Y," and "Z" coordinates. Note that the total number of dimensions is equal to the maximal number of level in the selected subnetwork. The last two columns contain the total number of terminal nodes "#T" (which is the total number of terminal nodes of the 2- and 3D $N$-hypercubes of the gates), and the total number of intermediate nodes "#N" in the $N$-hypercube.

It follows from the evaluation that

1. The hybrid approach does not suffer from exponential complexity of the number of terminal nodes, and demonstrates less values of the terminal nodes than is relevant to the number of gates in the benchmark circuit,

2. The space sizes $X$, $Y$, and $Z$ are higher than in the case of shared BDD models (Table 2).

Table 2. Fragment of an experiment on evaluation of $N$-hypercube parameters derived from a decision diagram.

| Circuit | I/O | #G | #Dim | X | Y | Z | Node |
|---|---|---|---|---|---|---|---|
| 27-channel interrupt controller | 36 7 | 160 | 36 | 88 | 12 | 12 | 1460 |
| 32-bit SEC circuit | 41/32 | 202 | 41 | 14 | 14 | 13 | 45922 |
| 8-bit ALU | 60 26 | 383 | 60 | 20 | 26 | 20 | 1011676 |
| 32-bit SEC circuit | 41 32 | 546 | 41 | 14 | 14 | 13 | 45922 |
| 16-bit SEC/DEC circuit | 33 25 | 889 | 33 | 11 | 11 | 11 | 42427 |
| 12-bit ALU and controller | 233/140 | 1193 | 233 | 78 | 78 | 77 | --- |
| 8-bit ALU | 50 22 | 1669 | 50 | 17 | 17 | 16 | 422803 |
| 9-bit ALU | 178 123 | 2307 | 178 | 60 | 59 | 59 | — |
| 16x16 multiplier | 32 32 | 2416 | 32 | 11 | 11 | 10 | — |
| 32-bit adder comparator | 207 108 | 3512 | 207 | 69 | 69 | 69 | --- |

**Table 3.** Fragment of an experimental study of the DAG based 3D models.

| Circuit | #O | #G | X | Y | Z | #I | #Node |
|---|---|---|---|---|---|---|---|
| 27-channel interrupt controller | 5 | 126 | 66 | 64 | 66 | 2022 | 1896 |
| 32-bit SEC circuit | 1 | 102 | 28 | 24 | 20 | 468 | 366 |
| 8-bit ALU | 24 | 130 | 70 | 72 | 70 | 612 | 482 |
| 32-bit SEC circuit | 1 | 322 | 58 | 52 | 54 | 1346 | 1024 |
| 16-bit SEC/DEC circuit | 25 | 522 | 100 | 104 | 92 | 2526 | 2004 |
| 12-bit ALU and controller | 139 | 828 | 82 | 80 | 78 | 3594 | 2766 |
| 8-bit ALU | 21 | 1458 | 132 | 132 | 140 | 9462 | 8004 |
| 9-bit ALU | 122 | 937 | 138 | 132 | 126 | 3750 | 2813 |
| 16×16 multiplier | 32 | 2327 | 248 | 248 | 244 | 9246 | 6916 |
| 32-bit adder/comparator | 107 | 474 | 114 | 112 | 106 | 1916 | 1442 |

# 5. NOTATION OF CHANGE IN COMPUTATIONAL NANOSTRUCTURES

## 5.1. Change at the Physical and Logical Level

Change is the fundamental concept of system at any level of observation and description:

1. Change at the physical level represents the behavior of nanostructure in time. The unit of change is derived from the physical nature of a nanostructure (molecular activity, atom dynamics, etc.).

2. Change at the logical level represents the behavior of the nanostructure during computation. A unit of logical change is derived from a computational model (elementary logic operations, data structure representation and manipulation, etc.).

Physical change is modeled by logical change.

Elementary transition in physical signal results in logical changes from 0 to 1 or vice versa. If physical change in a nanostructure can be uniquely transferred to logical change, a set of logical changes can describe the deterministic computational machine. Unfortunately, in today's nanostructures these transitions are not unique, that is, physical change can be transferred into logical change with a probability $p$

$$\langle \text{Physical change} \rangle \rightarrow \langle p(\text{Logical change}) \rangle$$

This physical nature of nanostructures results in various non-deterministic paradigms of computation, in particular, stochastic computing. Thermodynamic information corresponds to a physical unit of change, and Shannon information is rather associated with logical change:

1. $\langle \text{Unit of physical change} \rangle \Rightarrow \langle \text{Unit of thermodynamic information} \rangle$
2. $\langle \text{Unit of logical change} \rangle \Rightarrow \langle \text{Unit of Shannon information} \rangle$

To accommodate these dependencies, the design of nanostructures must comprise

1. Measurement of logical change in computational nanostructure.
2. Formal description of logical change.
3. Data structure representation in terms of logical change.

Measures on nanostructures that use the concept of change provide tools for their analysis at the following levels:

'Topology of computational nanostructure⟩ ⟺ ⟨Graphical data structure of computation⟩, i.e., from a given 2D or 3D topology of nanostructure (meshes, hypercube-like, pyramid-like, etc.), graphical data structure for computing can be derived (decision trees and diagrams, Boolean networks) using logical measures based on the concept of change.

Logical measures in computational nanostructure⟩ ⟺ ⟨Data structure representation⟩, that is, from logical measures using the concept of change, various data representations for computing can be derived (sum-of-products, Reed-Muller, Taylor-like expansions, word-level descriptions, hypercubes, and 2D and 3D decision diagrams).

⟨Probabilistic change⟩ ⇔ ⟨Stochastic computing⟩ states that physical change in nanostructure corresponds to logical change with some probability and any derived data structure from probabilistic changes will be of a stochastic nature.

⟨Many-valued change⟩ ⇔ ⟨Many-valued computational nanostructures⟩ states that physical changes in some nanotechnologies can be described at a logical level over a finite set of values, for example, 0, 1, and 2 for ternary logic, or $0, 2, \ldots, 7$ for octal logic. The computation addresses multivalued logic systems [52].

In this section, we give the formal definition of change, and introduce a technique of detection of change using various data structures.

## 5.2. Detection of Change

### 5.2.1. Detection of a Change in a Binary System

A signal in a binary system is represented by two logical levels, 0 and 1. Let us formulate the task as detection of the change in this signal. The simplest solution is to deploy an EXOR operation, modulo 2 sum of the signal $s_{i-t}$ before an "event" and the signal $s_t$ after the "event" (e.g., a faulty signal); that is, $s_{t-i} \oplus s_i$.

Example 20.   *For the signal depicted in Fig. 19, four possible combinations of the logical values or signals 0 and 1 are analyzed.*

If not change itself but direction of change is the matter, then two logical values 0 and 1 can characterize the behavior of the logic signal $s_i \in \{0, 1\}$ in terms of change, where 0 means any change of a signal, and 1 indicates that one of two possible changes has occurred $0 \to 1$ or $1 \to 0$.

### 5.2.2. Detection of Change in a Switching Function

Let the $i$th input of a switching function have been changed from the value $x_i$ to the opposite value, $\bar{x}_i$. This causes the circuit output to be changed from the initial value. Note that values $f(x_i)$ and $f(\bar{x}_i)$ are not necessarily different. The simplest way to recognize whether or not they are different is to try to find a difference between $f(x_i)$ and $f(\bar{x}_i)$.

### 5.2.3. Model of Single Change: Boolean Difference

Akers introduced the concept of Boolean difference [67]. Fundamentals of Boolean differential calculus are developed in [68, 69]. In [70], Boolean differences are used to find tests for switching circuits. Theoretical and applied aspects have been studied in [52].

The Boolean difference of a switching function $f$ of $n$ variables with respect to a variable $x_i$ is defined by the equation

$$\frac{\partial f}{\partial x_i} = \underbrace{f(x_1, \ldots, x_i, \ldots, x_n)}_{\text{Initial function}} \oplus \underbrace{f(x_1, \ldots, \bar{x}_i, \ldots, x_n)}_{\text{Function with } x_i \text{ complemented}} \tag{1}$$

It follows from the definition of Boolean difference that

$$\frac{\partial f}{\partial x_i} = \underbrace{f(x_1, \ldots, 0, \ldots, x_n)}_{x_i \text{ is replaced with } 0} \oplus \underbrace{f(x_1, \ldots, 1, \ldots, x_n)}_{x_i \text{ is replaced with } 1}$$

$$= f_{x_i=0} \oplus f_{x_i=1} \tag{2}$$



Figure 19.  The change of a binary signal and its detection (example 20).

Therefore, the simplest (but optimal) algorithm to calculate the Boolean difference of a switching function with respect to a variable $x_i$ includes two steps:

1. Replace $x_i$ in the switching function with 0 to get a cofactor $f_{x_i=0}$; similarly, replacement of $x_i$ with 1 yields $f_{x_i=1}$.
2. Find the modulo 2 sum of the two cofactors.

EXAMPLE 21.   *There are four combinations of possible changes of the output function $f = x_1 \vee x_2$ with respect to input $x_1$ ($x_2$). The Boolean differences of a switching function $f$ with respect to $x_1$ and $x_2$ are calculated by Eq. (2) in Fig. 20.*

The Boolean difference (Eq. [1]) possesses the following properties:

1. The Boolean difference is a switching function calculated by the Exclusive OR operation of the primary function and the function derived by complementing variable $x_i$; otherwise, it can also be calculated as EXOR of cofactors $f_{x_i=0}$ and $f_{x_i=1}$.
2. The Boolean difference is a switching function of $n - 1$ variables $x_1, x_2, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n$; that is, it does not depend on variable $x_i$.
3. The value of the Boolean difference reflects the fact of local change of the switching function $f$ with respect to changing the $i$th variable $x_i$: the Boolean difference is equal to 0 when such change occurs, and it is equal to 1 otherwise.

The Boolean difference (Eq. [1]) has a number of limitations, in particular: it cannot recognize the direction of change and cannot recognize the change in a function while changing a group of variables. This is the reason for extending the class of differential operators.

### 5.2.4. Model for Simultaneous Change

Consider the model of change with respect to simultaneously changed values of input signals. This model is called Boolean difference with respect to vector of variables. For a switching function $f$ Boolean difference of $n$ variables $x_1, \ldots, x_n$ with respect to the vector of $k$ variables $x_{i_1}, \ldots, x_{i_k}, i_1, \ldots, i_n \in \{1, \ldots, n\}$, is defined as follows

$$\frac{\partial f}{\partial(x_{i_1}, x_{i_2}, \ldots, x_{i_k})} = \overbrace{f(x_1, \ldots, x_{i_1}, x_{i_2}, \ldots, x_{i_k}, \ldots, x_n)}^{\text{Initial function}}$$

$$\oplus \underbrace{f(x_1, \ldots, \bar{x}_{i_1}, \bar{x}_{i_2}, \ldots, \bar{x}_{i_k}, \ldots, x_n)}_{\text{Function while } \bar{i}_{i_1}, \bar{x}_{i_2}, \ldots, \bar{i}_{i_k}} \tag{3}$$



Boolean difference with respect to $x_1$

$$\frac{\partial f}{\partial x_1} = \frac{\partial(x_1 \vee x_2)}{\partial x_1}$$

$$= (x_1 \vee x_2) \oplus (\bar{x}_1 \vee x_2)$$

$$= (0 \vee x_2) \oplus (1 \vee x_2) = \bar{x}_2$$

Boolean difference with respect to $x_2$

$$\frac{\partial f}{\partial x_2} = \frac{\partial(x_1 \vee x_2)}{\partial x_2}$$

$$= (x_1 \vee x_2) \oplus (x_1 \vee \bar{x}_2)$$

$$= (\bar{x}_1 \vee 0) \oplus (x_1 \vee 1) = \bar{x}_1$$

Figure 20.  Computing Boolean differences for a two-input OR gate (example 21).

Given $k = 2$, it follows from Eq. (3) that

$$\frac{\partial f}{\partial(x_i, x_j)} = f(x_1, \ldots, x_i, x_j, \ldots, x_n) \oplus f(x_1, \ldots, \bar{x}_i, \bar{x}_j, \ldots, x_n) \tag{4}$$

EXAMPLE 22.   *Calculate the Boolean difference of the switching function $f = x_1 x_2 \vee x_3$ with respect to a vector of variables using Eq. (4):*

$$\frac{\partial f}{\partial(x_1, x_2)} = f(x_1, x_2, x_3) \oplus f(\bar{x}_1 \bar{x}_2 x_3) = (x_1 x_2 \vee x_3) \oplus (\bar{x}_1 \bar{x}_2 \vee x_3)$$

$$= x_1 x_2 \oplus x_3 \oplus x_1 x_2 x_3 \oplus \bar{x}_1 \bar{x}_2 \oplus x_3 \oplus \bar{x}_1 \bar{x}_2 x_3 = (x_1 x_2 \oplus \bar{x}_1 \bar{x}_2) \bar{x}_3$$

$$= x_1 x_2 \bar{x}_3 \vee \bar{x}_1 \bar{x}_2 \bar{x}_3$$

## 5.2.5. Model of Multiple Change: k-Ordered Boolean Differences

Multiple, or $k$-ordered, Boolean difference is defined as

$$\frac{\partial^k f}{\partial x_{i_1} \partial x_{i_2} \ldots \partial x_{i_k}} = \frac{\partial}{\partial x_{i_1}} \left( \frac{\partial}{\partial x_{i_2}} \left( \ldots \frac{\partial f}{\partial x_{i_k}} \right) \ldots \right) \tag{5}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{Either way}}$$

It follows from Eq. (5) that

1. High-order differences can be obtained from single-order differences.
2. The order of calculation of the Boolean differences does not influence the result.

Let $k = 2$, then the second-order Boolean difference with respect to variables $x_i$ and $x_j$ will be

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left( \frac{\partial}{\partial x_j} \right) = \frac{\partial}{\partial x_j} \left( \frac{\partial}{\partial x_i} \right) \tag{6}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{Either way}}$$

## 5.2.6. Relationship of a Boolean Differences

There is a relationship between the second order Boolean difference (Eq. [6]) and Boolean difference with respect to a vector of two variables (Eq. [4]):

$$\left\{ \begin{array}{l} \dfrac{\partial f}{\partial(x_i, x_j)} = \dfrac{\partial f}{\partial x_i} \oplus \dfrac{\partial f}{\partial x_j} \oplus \dfrac{\partial^2 f}{\partial x_i \partial x_j} \\[3mm] \dfrac{\partial^2 f}{\partial x_i \partial x_j} = \dfrac{\partial f}{\partial x_i} \oplus \dfrac{\partial f}{\partial x_j} \oplus \dfrac{\partial f}{\partial(x_i, x_j)} \end{array} \right. \tag{7}$$

This relationship for two variables can be generalized for $k \leq n$ variables, that is, between multiple or $k$-ordered Boolean difference (Eq. [5]) and Boolean difference with respect to vector of $k$ variables (Eq. [3]).

EXAMPLE 23.   *Calculation of the two-ordered Boolean difference of the switching function $f = x_1 x_2 \vee x_3$ with respect to variables $x_1$, $x_2$ and the vector of variables $(x, x_2)$, is shown in Figure 21. To calculate Boolean difference $\partial f / \partial(x_1, x_2)$, Eq. (7) was used.*

## 5.2.7. Symmetric Properties of Boolean Difference

By inspection of Eq. (1), one can observe the symmetry in the computation:

$$\frac{\partial f_{x_i=0}}{\partial x_i} = \frac{\partial f_{x_i=1}}{\partial x_i}$$

The signal graph of the computation has a symmetrical structure well-known as "butterfly" (in signal processing). The graph input is the truth vector $\mathbf{F}$ of the given switching function $f$, and the result is the truth vector of the Boolean difference.

$$\text{Boolean difference } \frac{\partial f}{\partial(x_1, x_2)}$$
$$\text{with respect to vector of}$$
$$\text{variables } (x_1, x_2)$$

$$\frac{\partial f}{\partial x_1} \oplus \frac{\partial f}{\partial x_2} \oplus \frac{\partial^2 f}{\partial x_1 \partial x_2}$$

$$= (x_2 \bar{x}_3 \oplus x_1 \bar{x}_3 \oplus \bar{x}_3)$$

$$= \overline{(x_1 \oplus x_2)} \bar{x}_3$$

$$= x_1 x_2 \bar{x}_3 \vee \bar{x}_1 \bar{x}_2 \bar{x}_3$$

Figure 21. Interpretation of a Boolean difference with respect to a vector of variables by $N$-hypercube (example 23).

## 5.3. Computing of Change

In this section, the technique for computation of Boolean differences using a decision tree and an $N$-hypercube is introduced. There are two approaches:

*The first approach* is based on interpretation of the decision tree and $N$-hypercube whose nodes implement Shannon expansion. This attractive technique allows us to get values of Boolean differences without extra manipulation of the data structure (tree or hypercube);

*The second approach* is oriented to the Davio tree and a corresponding $N$-hypercube structure.

Both approaches include two phases: computing of Boolean differences and analysis of behavior of the switching function in terms of change.

## 5.4. Boolean Difference and $N$-Hypercube

The problem of computation of a decision tree or $N$-hypercube is formulated as the analysis of the behavior of data structure in terms of change. The example below introduces this technique.

EXAMPLE 24. *The $N$-hypercube in Fig. 22 represents the switching function* $x_1 x_2 \vee x_3$. *To analyze the behavior of this function, let us detect the changes as follows.*



Figure 22. Interpretation of Boolean differences by $N$-hypercube: Boolean difference with respect to $x_1$ (a), $x_2$ (b), and $x_3$ (c) (example 24).

1. *Boolean difference with respect to variable* $x_1$ *is* $\partial f/\partial x_1 = x_2\bar{x}_3$. *The logic equation* $x_2x_3 = 1$ *yields the solution* $x_2x_3 = 10$. *This specifies the conditions for detecting the changes at* $x_1$: *when* $x_2x_3 = 10$, *a change at* $x_1$ *cause a change at* $f$. *This can be seen on the decision tree and on the* $N$*-hypercube (Fig. 22[a]).*

2. *Boolean difference with respect to variable* $x_2$ *is* $\partial f/\partial x_1 = x_1\bar{x}_3$. *The logic equation* $x_2\bar{x}_3 = 1$ *specifies the condition of observation as a change at* $f$ *while changing* $x_1$: $x_2x_3 = 10$ *(Fig. 22[b]).*

3. *Boolean difference with respect to variable* $x_3$ *is* $\partial f/\partial x_1 = \overline{x_1x_2}$. *The logic equation* $\overline{x_1x_2} = 1$ *determines the condition:* $x_1x_2x_3 = \{00, 01, 10\}$ *(Fig. 22[c]).*

It was shown in [31] that the Davio decision tree can be embedded in $N$-hypercube which implements positive Davio expansion in the nodes.

To compute Boolean differences, let us rewrite positive Davio expansion in the form

$$f = f|_{x_i=0} \oplus x_i(f|_{x_i=0} \oplus f|_{x_i=1})$$

$$= \underbrace{f|_{x_i=0}}_{\text{Left branch}} \oplus \underbrace{x_i \frac{\partial f}{\partial x_i}}_{\text{Right branch}}$$

It follows from this form that

1. Branches of the Davio decision tree carry information about Boolean differences.
2. Terminal nodes are the values of Boolean differences for corresponding variable assignments.
3. Computing of Reed-Muller coefficients can be implemented on the Davio decision tree as a data structure.
4. Representation of a switching function in terms of change is a unique representation; it means that the corresponding decision diagram is canonical.
5. The values of terminal nodes correspond to coefficients of logic Taylor expansion.

The Davio tree can be embedded into an $N$-hypercube, and the previously mentioned properties are valid for that data structure as well. In addition, the $N$-hypercube enables computing of the Reed-Muller coefficients/Boolean differences, assuming that the processing is organized using parallel-pipelined, or systolic, processing.

EXAMPLE 25. *Figure 23 shows a Davio decision tree and corresponding* $N$*-hypercube for an arbitrary switching function of two and three variables.*

EXAMPLE 26. *Let* $f = x_1 \vee x_2$. *The values of Boolean differences given assignments* $x_1x_2 = \{00, 01, 10, 11\}$ *are:* $f(00) = 0$, $\partial f(01)/\partial x_1 = x_2 = 1$. $\partial f(10)/\partial x_2 = x_1 = 1$, *and* $\partial^2 f(11)/\partial x_1\partial x_2 = 1$. *They correspond to the terminal nodes of the Davio tree and* $N$*-hypercube (Fig. 24).*



Figure 23. Computing Boolean differences by Davio decision tree and $N$-hypercube for a switching function of three variables (example 25).

**Figure 24.** Computing Boolean differences of the switching function $f = x_1 \vee x_2$ (example 26).

One can conclude from example 26 that data structure in the form of a Davio decision tree carries information about Reed-Muller representation of switching functions and representation of switching functions in terms of change. The edges and values in terminal nodes of a Davio decision tree and $N$-hypercube carry information about the behavior of a switching function. Manipulation of a decision tree can be interpreted in terms of change: reduction of the decision tree to a decision diagram leads to minimization of Reed-Muller expression and can be used as a behavioral model of this function in terms of change.

## 6. FAULT-TOLERANT COMPUTATION

In the deterministic models of gates and circuits that were considered in previous chapters, the basic statements are

1. The input and output signals are deterministic.
2. The implemented logic function is performed correctly.

In nanocircuits, defects and faults arise from instability and noise-proneness on nanometer scales. The nature of noise signals in nanocircuits varies from thermal fluctuation to wave interface. Hence, different models are needed for investigation of the effects of noise, and development of methods for protection. For example, a model can be developed based on the assumption that desired signals in circuits are very noisy. In this model, the signals are modeled by the average of stochastic pulses generated by special devices.

This can be achieved in nanotechnology using probabilistic models.

There are various approaches to the development of probabilistic models, in particular:

1. *Stochastic models* for noise-making signals, in particular, Markov chain models and stochastic pulse stream models [45, 48, 71, 72].
2. *Neural networks* that use resources for optimization, and feedforward networks for computing logic functions over threshold elements [73].
3. *Computational techniques that are inspired by biology*. Some common examples of logic function calculation based on biologically inspired techniques include evolutionary algorithms [74–77], fuzzy logic [78], and artificial immune systems (immunological computation) [79]. The similarity between all known applications of algorithms based on biological paradigms is that they utilize the pattern-matching and learning mechanisms of the immune system to perform desired system functions. Biological immune system models are parallel and distributed structures that can be viewed as a multiagent system (separate functions are carried out by individual agents). The immune system model is a model of adaptive processes at the local level, resulting in useful behavior at the global level.

Below, some of state-of-the-art models based on the assumption of random factors are listed:

1. *Models for faults detection in wires.* For example, stuck-at-0 or stuck-at-1 is a fault type that causes a wire to be stuck at zero or one respectively [70].
2. *Stochastic* (probabilistic) models of behavior of gates and circuits [71, 72]. In these models, to estimate signal probabilities, one has to calculate the switching activity of the internal nodes of the circuit.
3. *Error correction codes* correct errors in order to ensure data fidelity. The random error correction codes refers to its ability to correct random bit errors within a code word [66].
4. *Model of switching activity,* or transition density model is based on the concept of change [80]. The model is represented in terms of Boolean differences.

## 6.1. Von Neumann's Model of Reliable Computation with Unreliable Components

The study of reliable computation by unreliable devices originates with von Neumann [43, 81]. He developed the multiplexing technique known as von Neumann's model of computing. In this model, each wire in a circuit is replaced by a bundle of wires on which a majority vote is conducted to establish its value. The classic von Neumann's model is the focus of many recent investigations [44, 82]. In this model, it is assumed that NAND is an unreliable gate, the following technique is implemented:

1. Replace each input of the NAND gate as well as its output by a bundle of $N$ lines.
2. Duplicate the NAND $N$ times.
3. Perform a random permutation of the input signals: each signal from the first input bundle is randomly paired with a signal from the second input bundle to form the input pair of one of the duplicated NANDs.

The key to this approach is modifying the NAND gate (an arbitrary network, in general) by replacing each interconnect with a parallel bundle of interconnects and a strategy of random interconnections that prevents the propagation of errors. In other words, parallelization by bundles and random interconnections can be viewed as a method for increasing the reliability of the NAND element.

## 6.2. Probabilistic Behavior of Nanodevices

In the probabilistic models, it is assumed that

1. The input and output signals are performed within some probability because of noisy signals.
2. The implemented logic function is performed within some probability because of the nature of nanodevices.

Noise in digital circuits is defined as any deviation of a signal from its stable value and can stem from sources as varied as physical and chemical processes in devices, measurement limitations, and stochastic simulation procedures. Noise can affect timing, causing a delay in failure, increase power consumption, and cause function failure because of signal deviation. It is important to understand and predict the effects of noise. As noise can have a variety of sources; different noise models that are effective in different situations are desirable, in particular,

1. Mutual inductance noise: when signal switching causes transient current to flow through the loop formed by the signal wire and current return path, a changing magnetic field is created and mutual inductance noise occurs.
2. Thermal noise: electrical power distribution and signal transmission through interconnections are always accompanied by thermal noise due to self-heating caused by the current flow. Thermal noise affects both interconnect design and reliability.

It is reasonable to distinguish the effects of noise in transmission and storage of information and processing of information. Figure 25 illustrates the random factors that influence the performance of a nanodevice.

Figure 25. Random factors that influence the performance of a nanodevice.

Deterministic models operate with noise-free signals, gates, networks, and fault-free hyper-cube structures.

Probabilistic models assume that input signals are applied to gates with some level of probability and correct output signals are calculated with some level of probability. When noise is allowed, the switching function is replaced with a random function and the configuration is a set of random variables.

It is essential that input and output signals of nanogates are described by additive or multiplicative expression of both noise and the desired signal. Often the desired signal is very noisy. This means that a special technique must be utilized to extract the desired signal from the noisy signal. These methods are well known and widely used in communication [47, 48]. However, these are costly methods, and their technical implementation is complicated. It is rather impractical to apply these methods in nanocircuit design, and other models are needed to solve the problem of fault tolerant computation in nanocircuits, even at the level of a single nanogate.

### 6.2.1. Neural Networks

Neural Network can be defined as a computational paradigm alternative to the conventional von Neumann model. The computational potential and limits of conventional computing models are well understood in terms of classical models such as the Turing machine. Many important results have been achieved in investigation of the computational power of neural networks by comparison with conventional computational tools such as finite automata, Turing machines, and logic circuits. Examples of deployment of this approach toward nano-electronic circuits are cellular neural networks [40] and neuromorphic networks [83].

### 6.2.2. Threshold Networks

Logic circuit design based on threshold gates can be considered as an alternative to traditional logic gate design procedure. The implementation of a massively interconnected network of threshold gates is possible [73]. It was shown, that multiple-addition, multiplication, division, and sorting can be implemented by polynomial-size threshold circuits of small depth. Formally, a threshold gate is described by a threshold decision (linearly separable) function. This principle is a general one in and of itself, so simple logic gates, such as AND and OR gates, are merely special cases of the threshold gate. An example of implementation of threshold logic circuit using nanoelectronics is adders based on resonant tunneling devices [84].

### 6.2.3. Stochastic Feedforward Neural Networks

An example of these is a stochastic feedforward neural network built on noisy spiking neurons, which have been developed to model biological neurons [85]. The important properties of these networks are as follows:

1. An arbitrary switching function can be implemented by a sufficient large network of noisy spiking neurons with an arbitrarily high probability of correctness.
2. An arbitrary deterministic finite automation can be simulated by a network of noisy spiking neurons with an arbitrarily high probability of correctness.

## 6.3. Fault-Tolerant Computing

The basic terminology of fault-tolerant computing includes:

*Robustness to errors* is the ability of a computer system to operate correctly in the presence of errors.

*Fault tolerance* is the ability of a computer system to recover from transient errors during computing.

*Defect tolerance* is the ability of a computer system to operate correctly in the presence of permanent hardware errors that emerged in the manufacturing process.

### 6.3.1. Techniques

Fault-tolerant techniques are basically built on two approaches:

1. Redundancy (*R*-fold modular and reconfiguration) achieves tolerance to faults by employing *R* copies of a unit [82, 86].
2. Stochastic computing achieves tolerance to faults by employing statistical models in which deterministic logic signals are replaced by random variables [71, 72, 87].

EXAMPLE 27. *Let Boolean variables $x_1$ and $x_2$ correspond to stochastic pulse signals with averages $E(x_1)$ and $E(x_2)$. Suppose these pulse streams are independent. It is possible to find logic operations that correspond to the sum $E(x_1) + E(x_2)$ and product $E(x_1) \times E(x_2)$ of these averages.*

So far, the first, architectural, approach has been mainly investigated in research on the design of fault-tolerant nanodevices [12, 13, 44, 88]. A few attempts have been made in the second direction, stochastic computing on nanocircuits [31, 89].

### 6.3.2. Hierarchical Levels

Hierarchical levels of fault-tolerant computing consist of

1. The basic primitives of a system. In the simplest case, the primitives are similar to a library of cells in conventional design. However, the complexity of primitives depends on the technology.
2. A finite set of primitives makes up macroprimitives, which are the smallest processors possible within the associated memory. This is similar to the microprocessor in conventional systems.
3. A finite set of macroprimitives makes up a system similar to the organization of multiprocessor systems.
4. The system makes up a distributed set of systems. This is the highest level of organization of conventional computer systems.

EXAMPLE 28. *A 4D hypercube can be recognized as a distributed (two connected 3D hypercubes) and multiprocessor system (each node corresponds to processor).*

Noise is but one aspect of the effect of errors on the practical implementation of computing circuits and systems. Permanent defects affecting computing resources during the manufacture of the system and within their subsequent lifetime are an engineering problem. Reconfigurable and self-repairing architectures are used to solve this problem.

In particular, error-adaptive architectural approaches such as. reconfigurable architectures, and neural-like network with training can be employed.

The next level of this hierarchy is the level of self-replicated. self-repairing. and self-assembling systems. At this level, a system, for example, can replicate itself, giving rise to a population of identical systems.

From the above follow different fault-tolerant computing models:

1. *Probabilistic fault-tolerant computing models of nanogates* which rely on the observation of the mechanism of nanodevices. based on the probabilistic behavior of nanostructures (electrons, molecules) [3].

2. *Probabilistic behavior of circuits and systems.* At these levels the probability of getting failed components becomes higher. This approach is based on the idea of incorporating into the circuit and system a "guard" against failures [12, 44, 88].

EXAMPLE 29. *In the presence of faults, a fault-tolerant circuit or system reconfigures itself to exclude the faulty elements. Normally, it is expected for a circuit and system, upon reconfiguration, to encompass all the healthy elements whenever possible. A system so reconfigured may or may not change its topology. Ideally, a fault-tolerant design retains the same system topology after faults arise.*

An example of an adaptively configured architecture, called neuromorphic networks with single-electron latching switches as nanoscale synapses [83]. The network require training to be useful for solving problems of pattern recognition and signal processing.

Another approach, a formal probabilistic framework for a reconfigurable architecture without training is based on Markov random field [89]. In this model (not grounded in physical device structure as yet), a logic circuit is mapped into a Markov random field that is a graph indicating the neighborhood relations of the circuit nodes. This graph is used in probability maximization process, aimed at characterization of circuit configurations for the best thermal noise reliability (expressed in terms of logic signal errors).

## 6.4. Stochastic Models of Switching Gates

In this section, models based on reliable gates with stochastic input streams are considered [31, 45, 46]. Figure 26 illustrates this model. If the input stochastic streams are independent (technically this means that independent generators of random pulses are used with some additional tools for decorrelation of signals) with $E(x_1)$ and $E(x_2)$, the output is described by the equation $E(f) = E(x_1) \times E(x_2)$. Then follows the transformation of the values to the range [0, 1].

EXAMPLE 30. *Given the deterministic signal $x \in \{0, 1\}$, generate this signal with probability $p(x)$. The simplest model is $p(x) = x \cdot r$ where $r \in \{0, 1\}$ is a random variable with probability $p(r)$.*

EXAMPLE 31. *Implementation of the model for generating a signal with a given probability is shown in Fig. 27 (synchronization is not shown). This illustrates the possibility of generating a signal with a given probability. For example, if $p(r) = 1$, the output is a signal $x$ with probability $p(x) = 1$. On the basis of this model it is possible to study the simplest features of probabilistic computation.*

### 6.4.1. The Model of a Gate for Input Random Pulse Streams

Let us analyze the output of a gate that implements an elementary switching function for input independent random pulse streams as 0s and 1s. A binary stochastic pulse stream is defined as a sequence of binary digits, or bits. The information in a pulse stream is contained in the primary statistics of the bit stream, or the probability of any given bit in the stream being a logic 1. Hence, the output of a gate will generally be in the form of a nonstationary Bernoulli sequence. Such a sequence can be considered in probabilistic terms as a deterministic signal with superimposed noise. Suppose that statistical characteristics of these streams are known (i.e., can be measured). In other words, these streams carry a signal by statistical characteristics (a single event carries very little information, it is not enough for decision making).



Figure 26. Stochastic pulse model of computing.

**Signal**



*Deterministic signal x is transmitted to
output with probability p(r).*

**Random pulse stream**

**Figure 27.** Model for generating the signal binary $x$ with probability $p(r)$.

The stochastic pulse stream model states that

1. Input signals are modeled by stochastic pulse streams with known characteristics.
2. The output signals are calculated as an average of statistical characteristics.

The generated probability of a sequence of logic levels corresponds to the relative frequency of 1 logic levels in a sufficiently long sequence. A probability cannot be measured exactly but only estimated as the relative frequency of 1 logic level in a sufficiently long sample.

The stochastic computer introduces its own errors in the form of random variance. If we observe a sequence of $N$ logic levels and $k$ of them are 1, then the estimated generating probability is $\hat{p} = k/N$. The sampling distribution of the value of $k$ is binomial, and hence the standard deviation of the estimated probability $\hat{p}$ from the true probability $p$ is $\sigma(\hat{p}) = [p(1-p)/N]^{1/2}$.

Hence the accuracy in estimation of a generated probability increases as the square root of the length of the sequence examined, that is, the square root of the length, or time, of computation.

There are several features that distinguish classical computation and stochastic computation:

1. A signal is represented by the probability that a logic level be 1 or 0 at a clock pulse.
2. Random noise is being deliberately introduced into the data; usually, noise distribution is normal.
3. A quantity is represented by a clocked sequence of logic levels generated by a random process: the successive levels are statistically independent, and the probability of the logic level being ON is a measure of that quantity.
4. Arithmetic operations are performed via the completely random data, and the probability that a logic level will be ON or OFF is determined. Its actual value is a chance event which cannot be predicted, and repetition of a computation will give rise to a different sequence of logic levels.

In a conventional computer, logic levels represent data change deterministically from value to value as the computation proceeds. If the computation is repeated, the same sequence of logic levels will occur. Note that unlike binary radix arithmetic, stochastic arithmetic is robust in the presence of noise/single bit fault, and accuracy may be controlled using the dimension of time.

If the input distribution is unconstrained, Bernoulli sequences can be used for formal modeling. This means that

1. The probability of a given bit being a 1 is independent of the values of any previous bits.
2. Elements' processing functions are evaluated only with respect to their outputs' primary statistics. The outputs are not, in general, Bernoulli sequences.
3. In the case of processing elements with multiple inputs, the inputs are uncorrelated with each other.

EXAMPLE 32. In Table 4 probabilistic parameters of elementary switching functions for stochastic computing are given, where autocorrelation function is defined as $K_f(\tau) = E[(f(t) - E(f))(f(t-\tau) - E(f))]$.

It follows from Table 4 that, for example, for stochastic computing of a NOT function, we assume that input signal $x$ is a stochastic pulse stream characterized by the probability

Table 4. Stochastic computing of elementary switching functions.

| | | | |
|---|---|---|---|
| AND | $E(f) = \begin{cases} E(x_1)E(x_2) + K_{1,2} & x_1 \text{ and } x_2 \text{ are dependent} \\ E(x_1)E(x_2) & \text{otherwise} \end{cases}$ | | |
| OR | $E(f) = \begin{cases} E(x_1) + E(x_2) - E(x_1)E(x_2) - K_{1,2} & x_1 \text{ and } x_2 \text{ are dependent} \\ E(x_1) + E(x_2) - E(x_1)E(x_2) & \text{otherwise} \end{cases}$ | | |
| NAND | $E(f) = \begin{cases} 1 - E(x_1)E(x_2) - K_{1,2} & x_1 \text{ and } x_2 \text{ are dependent} \\ 1 - E(x_1)E(x_2) & \text{otherwise} \end{cases}$ | | |
| NOR | $E(f) = \begin{cases} 1 - E(x_1) - E(x_2) + E(x_1)E(x_2) + K_{1,2} & x_1 \text{ and } x_2 \text{ are dependent} \\ 1 - E(x_1) - E(x_2) + E(x_1)E(x_2) & \text{otherwise} \end{cases}$ | | |
| EXOR | $E(f) = \begin{cases} E(x_1) + E(x_2) - 2E(x_1)E(x_2) - 2K_{1,2} & x_1 \text{ and } x_2 \text{ are dependent} \\ E(x_1) + E(x_2) - 2E(x_1)E(x_2) & \text{otherwise} \end{cases}$ | | |

$p(x)$ and the autocorrelation function $K_x(\tau_x)$. The mean of the output signal is $E(f) = p(f)$, and hence $p(\bar{x}) = 1 - p(x) = 1 - E(x)$. By analogy, if the input pulse streams are independent

AND gate $f = x_1 x_2$ is modeled by $E(f) = p_1 p_2$

OR gate $f = x_1 \vee x_2$ is modeled by $E(f) = p_1 + p_2 - p_1 p_2$

NOR gate $f = \overline{x_1 \vee x_2}$ is modeled by $E(f) = 1 - p_1 - p_2 + p_1 p_2$

EXOR gate $f = x_1 \oplus x_2$ is modeled by $E(f) = p_1 + p_2 - 2p_1 p_2$

NOT-EXOR gate $f = \overline{x_1 \oplus x_2}$ is modeled by $E(f) = 1 - p_1 - p_2 + 2p_1 p_2$

where $p_1 = E(x_1)$ and $p_2 = E(x_2)$.

EXAMPLE 33. *Let the inputs of OR gate $x_1 \in \{0, 1\}$ and $x_2 \in \{0, 1\}$ be mutually independent with probabilities $p_1 = p(x_1)$ and $p_2 = p(x_2)$ correspondingly (Fig. 28). The output probability can be evaluated as the probability of at least one event $x_1$ and $x_2$, that is,*

$$p = 1 - (1 - p_1)(1 - p_2) = p_1 + p_2 - p_1 p_2.$$

*Supposing $p_1 = 0.8$, $p_2 = 0.9$, correct output is produced with a probability of $p = 0.8 + 0.9 - 0.8 \cdot 0.9 = 0.98$. If $p_1 = p_2 = 1$, the inputs become deterministic and $f = x_1 + x_2 - x_1 x_2$, that is, $f = x_1 \vee x_2$.*

Note that in the above example, in general, the mean $E(f)$ and variance $D(y)$ of the output $y$ are equal to

$$E(f) = p_1 + p_2 - p_1 p_2 = p$$

$$D(f) = p \cdot (1 - p)$$

Stochastic computing can be interpreted by decision trees and $N$-hypercubes.



(a)                                    (b)

Figure 28. Deterministic (a) and probabilistic (b) models of computing.

## 6.5. Fault-Tolerant Hypercube-Like Computing Structures

Fault-tolerant properties of a hypercube-like structures are well studied and used in computing systems. Here, we briefly review the basic principles of fault-tolerant computing via hypercube-like structures.

Below, the basic definitions of fault tolerance computing in hypercube and hypercube-like structures are given:

1. A hypercube computing structure is called a faulty hypercube if it contains any faulty node (computing device) or communication link. For hypercube-like structures of large dimensions, the number of processing elements is very large and hence the probability of occurrence of faults increases.

2. A network is robust if its performance does not decrease significantly when its topology changes.

3. Since efficient cooperation between nonfaulty computing devices is desirable, one measure for robustness is the network connectivity, which is defined as the number of node or link failures that can be allowed without disrupting the system.

4. *Fault models* of a hypercube computing system are defined from subcube and node reliability.

5. *Multiple fault models* of a hypercube computing system are calculated based on the probability that many faults in the node or subcube exist.

6. The *reliability* of a hypercube-based computing system is defined as the probability that the system has survived the interval $[0, t]$ given that it was operational at time $t = 0$, where $t$ is the time. Usually, reliability is used in models of computing systems in which repair cannot take place.

7. The *terminal reliability* of a computer system is defined as the reliability of computer devices in nodes of a hypercube computer system. Terminal reliability can be also defined as task-based reliability, which is the probability that some minimal set of connected nodes are available in the hypercube structure.

8. The *fault-tolerance computing* of a hypercube-based computing system is provided by reconfiguration and application of error correcting codes.

Fault-tolerance technique for hypercube and hypercube-like structures is based on the principle of Reconfiguration and Error correcting.

These techniques are well studied and widely used in hypercube-like system design [37, 38].

EXAMPLE 34. *Consider the model of nanodevice based on transmitting a binary signal through a channel. Suppose that the probability that the receiver will get one (zero) when zero (one) is sent is $0 < p < 1/2$. If this channel is used only once, the probability of a correct transition is $1 - p$. To improve the reliability of transmission, the sender transmits the sequence of either three zeros or three ones. Then, the probability of a correct transmission is $1 - (p^3 + 3p^2(1 - p))$. For example, $1 - (p^3 + 3p^2(1 - p)) = 0.896$ for $p = 0.2$ compared with $1 - p = 0.8$.*

Several algorithms have been developed for reconfiguring a hypercube with faults. These algorithms aim to achieve different characteristics after reconfiguration, in particular, acceptable performance and connectivity. The crucial idea is to identify maximum subcubes in a faulty hypercube, retaining as many healthy nodes as possible to keep performance degradation to a minimum.

There are several assumptions and additional data that must be known to apply the above characteristics: in particular, the reliability function of node computing devices (usually assumed to be homogeneous), and the characteristic of statistical dependence of failures of computing devices in nodes (usually assumed to be independent).

The hypercube-like network has been proved to be very robust and to divide it into two components requires at least $n$ faults.

In the probability fault model, the reliability of each node at time $t$ is a random variable. The probability that a hypercube-like network is operational is represented by the reliability of the computing devices in the hypercube-like network. The reliability of computing hypercube-like structure can be formulated as the union of probabilistic events that all the possible hypercubes of lower dimensions are operational.

# 7. INFORMATION THEORETICAL MEASURES FOR COMPUTATIONAL NANOSTRUCTURES

Information theory has been mentioned in the "Introduction" as a way to characterize and evaluate the parameters of a nanosystem, such as information transfer per unit space, and the number of operations per second.

Entropy is interpreted as the amount of disorder in the system. Indeed, in thermodynamics, entropy is defined as the thermodynamic probability of the internal particles of a system while holding the external properties constant. A hundred years later, in 1948, Shannon suggested a measure to represent the information by a numerical value, nowadays known as Shannon entropy [51] with respect to this transformation. Since then, the term "uncertainty" is interchangable with the term "entropy." The Shannon information theory has been developed for many applications in circuit design. The bit strings of information are understood as messages to be communicated from a messenger to a receiver.

## 7.1. Overview

Physicists have emphasized thermodynamics entropy. The relationship of these different measures has been considered in many papers, for instance, in [90]. In a binary system, physical entropy becomes Shannon entropy.

Information theoretical measures in Shannon notation have been used in [91–93] in decision trees, and in diagram design. Entropy-based strategies for minimization of logic functions have been studied in [75, 94]. These results are related to the earlier work [95–99] on conversion of decision tables (truth tables of logic functions) into decision trees.

Existing techniques for power estimation at gate and circuit levels rely on probabilistic information of the input stream. The average switching activity per node (gate) is the main parameter that needs to be correctly determined. These and related problems are the focus of many researchers. For example, in [100, 101], it is demonstrated that the average switching activity in a circuit can be calculated using either entropy or information energy averages.

Most of the algorithms for minimization of state assignments in finite state machines target reduced average switching per transition, that is, average Hamming distance between states [102].

There have already been some approaches to evolutionary circuit design [76]. The main idea is that an evolutionary strategy would inevitably explore a much richer set of possibilities in the design space than are within the scope of traditional methods. In [74, 75, 77, 94] an evolutionary strategy and information theoretical measures were used in circuit design.

A deep and comprehensive analysis of computing systems' information engine has been done in [103]. The relationship between function complexity and entropy is conjectured in [104].

This section focuses on information measures in nanosystems. Applying the notation to a physical system (hardware), information, in a certain sense, is a measurable quantity, which is independent of the physical medium by which it is conveyed. The most appropriate measure of information is mathematically similar to the measure of entropy, but there is good reason for reversing the sign and stating that information is the negative of entropy in nature as well as in mathematical formulation. The technique of information theory is applied to problems of the extraction of information from systems containing an element of randomness.

## 7.2. Information Theoretical Measures on Data Structures

The measures can be made on different data structures that carry information about a logic function: formal representation (sum-of-products, Reed-Muller, arithmetic, and word-level forms), logic network, flowgraphs, and decision trees and diagrams, including spatial representations. Information-theory measures are sensitive to data structure.

In this section, we focus on the entropy of spatial measurement in $N$-hypercube space. The information content of a switching function is an inherent attribute of a function and is technology independent. Information content defines the complexity of function implementation

and can be used to estimate a lower bound on some physical (topological) parameters with respect to various implementations. Therefore, it reflects the fundamental characteristic of function behavior. Entropy of spatial measurement in $N$-hypercube space can be viewed as a contribution to information content, over all nodes of the embedded decision diagram.

The above is the basis for approaches to estimating the complexity involved when data are transmitted from various points in a circuit. The estimated attributes here are information flow, information amount, and entropy measures on the hypercube. Finally, we describe other information-theoretical definitions and outline their application to the problem of synthesis of 3D structures.

## 7.3. Measures in Logic Design

The most basic information-theoretical measure is entropy. Many useful additional characteristics are derived from the entropy, namely, the conditional entropy, mutual information, joint information, and relative information. Figure 29 illustrates the basic principles of input and output information measures in a logic circuit, where the shared arrows mean that the value of $x_i(f)$ carries the information; the shared arrow therefore indicates the direction of the information stream. Obviously, we can compare the results of the input and output measures and calculate the loss of information.

## 7.4. Information-Theoretical Standpoint

A computing system can be seen as a process of communication between computer components. The classical concept of information advocated by Shannon is insufficient to capture a number of features of the design and processing of a computing system. The information-theoretical standpoint on computing is based on the following notations:

1. *Source of information*, a stochastic process where an event occurs at time point $i$ with probability $p_i$. In other words, the source of information is defined in terms of the probability distribution for signals from this source. Often the problem is formulated in terms of sender and receiver of information and used by analogy with communication problems [90, 97, 101].

2. *Information engine*, the machine that deals with information [103].

3. *Quantity of information*, a value of a function that occurs with the probability $p$ carries a quantity of information equal to $(-\log_2 p)$ [104, 105].

4. *Entropy*, $H(f)$, the measure of the information content of $X$. The greater the uncertainty in the source output, the higher is its information content. A source with zero uncertainty would have zero information content and, therefore, its entropy would be likewise equal to zero [51].



Figure 29. Information measures at the input and output of a logic circuit, and computing input/output relationships of information.

The information and entropy. in their turn. can be calculated with respect to the given sources:

*Information* carried by the value of a variable or function.
*Conditional entropy* of function $f$ values given function $g$.
*Relative information* of the value of a function given the value of a variable.
*Mutual information* between the variable and function,
*Joint entropy* over a distribution of jointly specified functions $f$ and $g$.

## 7.4.1. Quantity of Information

Let us assume that all combinations of values of variables occur with equal probability. A value of a logic function that occurs with the probability $p$ carries a quantity of information equal to

$$\langle \text{Quantity of information} \rangle = -\log_2 p \text{ bit}$$

where $p$ is the probability of that value occuring.
The information carried by the value of $a$ of $x_i$ is equal to

$$I(x_i)|_{x_i=a} = -\log_2 p \text{ bit}$$

where $p$ is the quotient between the number of tuples whose $i$th components equal $a$ and the total number of tuples. Similarly, the information carried by a value $b$ of $f$ is

$$I(f)|_{f=b} = -\log_2 q \text{ bit}$$

where $q$ is the quotient between the number of tuples in the domain of $f$ and the number of tuples for which $f$ takes the value $b$.

EXAMPLE 35. *The information carried by the values of variable $x_i$ and function $f$ for a switching function given by a truth table is calculated in Figure 30.*

## 7.4.2. Conditional Entropy and Relative Information

Conditional entropy is a measure of a random variable $f$ given a random variable $x$. To compute the conditional entropy, the conditional probability of $f$ must be calculated. The



*Probabilities of the values of variable $x_i$:*

$$p(x_i = 0) = 3/5, \quad p(x_i = 1) = 2/5$$

*The information carried by the variable $x_i$:*

$$\begin{cases} I(x_i)|_{x_i=0} = -\log_2 3/5 = 0.737 \text{ bit} \\ I(x_i)|_{x_i=1} = -\log_2 2/5 = 1.322 \text{ bit} \end{cases}$$

*Probabilities of the values of function $f$:*

$$p(f = 0) = 4/5, \quad p(f = 1) = 1/5$$

*while*

$$\begin{cases} p(f = 0)|_{x_i=0} = 3/5 \quad p(f = 0)|_{x_i=1} = 1/5 \\ p(f = 1)|_{x_i=0} = 0 \quad p(f = 1)|_{x_i=1} = 1/5 \end{cases}$$

*Information carried by the function $f$:*

$$\begin{cases} I(f)|_{f=0} = \log_2 4/5 = 0.322 \text{ bit} \\ I(f)|_{f=1} = \log_2 1/5 = 2.322 \text{ bit} \end{cases}$$

| Input | Output |
|-------|--------|
| $x_i$ | $f$ |
| 0 | 0 |
| 1 | 0 |
| 0 | 0 |
| 1 | 1 |
| 0 | 0 |

**Figure 30.** The information carried by values of variable $x_i$ and switching function $f$ (example 35).

conditional probability of value $b$ of logic function $f$, given input value $a$ of $x_i$ is

$$p(f = b|x_i = a) = \frac{p_{|f=b \atop |x_i=a}}{p_{|x_i=a}}$$

Similarly, the conditional probability of value $a$ of $x_i$, given value $b$ of the function $f$, is

$$p(x_i = a|f = b) = \frac{p_{|f=b \atop |x_i=a}}{p_{|f=b}}$$

Conditional entropy $H(f|g)$ of function $f$ values given logic function $g$ is

$$H(f|g) = H(f, g) - H(g) \tag{8}$$

In circuit analysis and decision tree design, the so-called chain rule is useful (Fig. 31):

$$H(f_1, \ldots, f_n|g) = \sum_{i=1}^{n} H(f_i|f_1, \ldots, f_{i-1}, g) \tag{9}$$

The relative information of the value $b$ of a logic function $f$ given the value $a_i$ of the input variable $x_i$ is

$$I(f = b|x_i = a) = -\log_2 p(f = b|x_i = a)$$

The relative information of the value $a_i$ of the input variable $x_i$ given the value $b$ of the logic function $f$ is

$$I(x_i = a|f = b) = -\log_2 p(x_i = a|f = b)$$

Once the probability is equal to 0, we suppose that the relative information is equal to 0.

EXAMPLE 36. *Figure 32 illustrates the calculation of the conditional and relative information given the truth table of a switching function.*

### 7.4.3. Entropy of a Variable and a Function

Let the input variable $x_i$ be the outcome of a probabilistic experiment and the random logic function $f$ represent the output of some step of computation. Each experimental outcome results in different conditional probability distributions on the random $f$. Shannon's entropy of the variable $x_i$ is defined as

$$H(x_i) = \sum_{l=0}^{m-1} p_{|x_i=a_l} \log_2 p_{|x_i=a_l} \tag{10}$$

where $m$ is the number of distinct values assumed by $x_i$. Shannon's entropy of the logic function $f$ is

$$H(f) = \sum_{k=0}^{n-1} p_{|f=b_k} \log_2 p_{|f=b_k} \tag{11}$$

where $n$ is the number of distinct values assumed by $f$.

This definition of the measure of information implies that the greater the uncertainty in the source output, the smaller is its information content. In a similar fashion, a source with



The conditional entropy $H(f|g)$ is non-negative. The value is zero if and only if a function $f$ such that $f = u(g)$ exists with a probability of one. Additivity of entropy or the chain rule for entropies is defined as

$$H(f, g) = H(f) + H(g|f)$$

Figure 31. The additivity of entropy.

*Conditional probabilities:*

$$p_{|f=0}= 3/5 \quad p_{f=0}= 1/5 \quad p_{|f=1}=0 \quad p_{|f=1}= 1/5$$
$$|x=0 \qquad |x=1 \qquad |x=0 \qquad |x=1$$

*Then*

$$p(f=0|x_i=0) = p_{|f=0} \cdot p_{|x=0}= 3/5 \cdot 3/5 = 1$$
$$|x=0$$

$$p(f=0|x_i=1) = p_{|f=0} \cdot p_{|x=1}= 1/5 : 2/5 = 1/2$$
$$|x=1$$

$$p(f=1|x_i=0) = p_{|f=1} \cdot p_{|x=0}= 0$$
$$|x=0$$

$$p(f=1|x_i=1) = p_{|f=1} \cdot p_{|x=1}= 1/5 : 2/5 = 1/2$$
$$|x=1$$

*Conditional entropy $H(f|x_i)$:*

$$H(f|x_i) = -p(f=0|x_i=0) \log p(f=0|x_i=0)$$
$$-p(f=0|x_i=1) \log p(f=0|x_i=1)$$
$$-p(f=1|x_i=0) \log p(f=1|x_i=0)$$
$$-p(f=1|x_i=1) \log p(f=1|x_i=1)$$
$$= -1 \log 1 - 1/2 \log 1/2 - 0 \log - 1/2 \log 1/2$$

*Relative information $I(f=b|x_i=a)$:*

$$I(f=0|x_i=0) = -\log_2 1 = 0$$
$$I(f=0|x_i=1) = -\log_2 1/2 = 1$$
$$I(f=1|x_i=0) = 0$$
$$I(f=1|x_i=1) = -\log_2 1/2 = 1$$

| Input | Output |
|-------|--------|
| xi | f |
| 0 | 0 |
| 1 | 0 |
| 0 | 0 |
| 1 | 1 |
| 0 | 0 |

**Figure 32.** Computing conditional entropy and relative information (example 36).

zero uncertainty would have zero information content and, therefore, its entropy would be likewise equal to zero.

EXAMPLE 37. *Figure 33 illustrates the calculation of entropy for the variable and function. The entropy of the variable $x_i$ and switching function $f$ are 0.971 bits and 0.722 bits.*

Therefore,

1. For any variable $x_i$ it holds that $0 \le H(x_i) \le 1$; similarly, for any function $f$, $0 \le H(x_i) \le 1$.



$H(x_i)$       $H(f)$

*Shannon's entropy*

$$H(x_i) = -3/5 \cdot \log_2 3/5 - 2/5 \cdot \log_2 2/5$$
$$= 0.971 \ bit$$
$$H(f) = -4/5 \cdot \log_2 4/5 - 4/5 \cdot \log_2 4/5$$
$$= 0.722 \ bit$$

| Input | Output |
|-------|--------|
| xi | f |
| 0 | 0 |
| 1 | 0 |
| 0 | 0 |
| 1 | 1 |
| 0 | 0 |

*The mutual information*

$$I(f:x_i) = \sum_{k=1}^{4} \sum_{l=1}^{5} p_{|f=b} \times I_{|f=b}$$
$$|x=a \quad |x=a$$

$$= 0.6 \cdot 0.322 - 0.2 \cdot 0.678 + 0 + 0.2 \cdot 1.322$$
$$= 0.322 \ bit$$

**Figure 33.** Shannon's entropy and mutual information (examples 37 and 38).

2. The entropy of any variable in a completely specified function is 1.
3. The entropy of a constant is 0.

## 7.4.4. Mutual Information

Mutual information is used to measure the dependence of the function $f$ on the values of the variable $x_i$ and vice versa, that is, how statically distinguishable distributions of $f$ and $x_i$ are. If the distributions are different, then the amount of information $f$ carries about $x_i$ is large. If $f$ is independent of $x_i$, then $f$ carries zero information about $x_i$. Figure 34 illustrates the mutual information between two variables $f$ and $g$.

The mutual information between the value $b$ of the function and the value $a$ of the input variable $x_i$ is:

$$I(f; x_i) = I(f; x_i)_{|f=b} - I(f = b|x_i = a)$$

$$= -\log_2 p_{|f=b} + \log_2 \frac{p_{|f=b}^{|x_i=a}}{p_{.x_i=a}}$$

By analogy, the mutual information between the input variable $x_i$ and the function $f$ is

$$I(f; x_i) = \sum_k \sum_l p_{|f=b_k}^{|x_i=a_l} \times I(f; x_i)_{|f=b_k}^{|x_i=a_l}$$

$$= \sum_k \sum_l p_{|f=b_k}^{|x_i=a_l} \times \log_2 \frac{p_{|f=b_k}^{|x_i=a_l}}{p_{|x_i=a_l}}$$

Useful relationships are

$$I(g; f) = I(f; g) = H(f) - H(f|g)$$

$$= H(g) - H(g|f)$$

$$= H(f) + H(g) - H(f, g);$$

$$I(g; f_1, \ldots, f_n|z) = \sum_{i=1}^{n} I(g; f_i|f_1, \ldots, f_{i-1}, z)$$

where $I(g; f|z)$ is the conditional mutual information between $g$ and $f$ given $z$. If $g$ and $f$ are independent, then $I(g; f) \geq 0$. The mutual information is a measure of the correlation between $g$ and $f$. For example, if $g$ and $f$ are equal with high probability, then $I(g; f)$ is large. If $f_1$ and $f_2$ carry information about $g$ and are independent given $g$ then $I(z(f_1, f_2); g) \leq I(f_1; g) + I(f_2; g)$ for any switching function $z$.

EXAMPLE 38.   *Figure 34 illustrates the calculation of the mutual information. The variable $x_i$ carries 0.322 bits of information about the switching function $f$.*



*The mutual information is defined as the difference between the entropy and the conditional entropy*

$$I(f; g) = H(f) - H(f|g)$$

$$= H(f) + H(g) - H(f, g).$$

*i.e., the difference of the uncertainty of $f$ and the remaining uncertainty of $f$ after knowing $g$. This quantity is the information of $f$ obtained by knowing $g$. The mutual information is a degree of the dependency between $f$ and $g$ and always takes positive values.*

*The additivity for mutual information of three random variables:*

$$I(f; g, z) = I(f; g) + I(f; z|g)$$

Figure 34. The mutual information.

## 7.4.5. Joint Entropy

Let the distribution of jointly specified functions $f$ and $g$'s values be known. Then, the joint entropy $H(f, g)$ given this distribution is defined as follows:

$$H(f, g) = - \sum_{a=0}^{m-1} \sum_{b=0}^{m-1} p_{\substack{f=a \\ g=b}} \cdot \log p_{\substack{f=a \\ g=b}}, \tag{12}$$

where $p_{\substack{f=a \\ g=b}}$ denotes the probability that $f$ takes value $a$ and $g$ takes value $b$, simultaneously.

## 7.5. Information Measures of Elementary Switching Functions

There are two approaches to information measures of elementary functions of two variables:

1. The values of input variables are considered as random patterns; for a two-input elementary function there are four random patterns $x_1 x_2 \in \{00, 01, 10, 11\}$.
2. The values of input variables are considered as noncorrelated random signals; for a two-input elementary function, there are random signals $x_1 \in \{0, 1\}$ and $x_2 \in \{0, 1\}$.

## 7.5.1. Information Measures Based on Pattern

Consider a two-input AND function with four random combinations of input signals: 00 with probability $p_{00}$, 01 with probability $p_{01}$, 10 with probability $p_{10}$, and 11 with probability $p_{11}$ (Fig. 35[a]).

Using Shannon's formula (Eq. [10]), we can calculate the entropy of the input signals, denoted by $H_{in}$ as follows

$$H_{in} = - p_{00} \times \log_2 p_{00} - p_{01} \times \log_2 p_{01}$$

$$- p_{10} \times \log_2 p_{10} - p_3 \times \log_2 p_{11} \text{ bit/pattern}$$

Maximum entropy of the input signals can be calculated by inserting into the above equation $p_i = 0.25$, $i = 0, 1, 2, 3$ (Fig. 36).

The output of the AND function is equal to 0 with probability 0.25, and equal to 1 with probability 0.75. The entropy of the output signal, $H_{out}$, is calculated by Eq. (11)

$$H_{out} = -0.25 \times \log_2 0.25 - 0.75 \times \log_2 0.75 = 0.81 \text{ bit/pattern}$$

The example below demonstrates a technique of computing information measures with which input signals are not correlated.

## 7.5.2. Information Measures Based on Noncorrelated Signals

Let the input signal be equal to 1 with probability $p$, and 0 with probability $1 - p$ (Fig. 35[b]). The entropy of the input signals is

$$H_{in} = -(1 - p)^2 \times \log_2 (1 - p)^2 - 2(1 - p) \times \log_2 (1 - p)p - p^2 \times \log_2 p^2$$

$$= -2(1 - p) \times \log_2 (1 - p) - 2p \times \log_2 p \text{ bit}$$



(a)                    (b)

Figure 35. Measurement of probabilities: random patterns (a) and noncorrelated signals (b).

*Method 1:*

The entropy of the input pattern (probability $p_i = 0.25$)

$$H_{in} = -4 \times 0.25 \times \log_2 0.25 = 2 \ bit/pattern$$

Entropy of the output signal

$$H_{out} = -0.25 \times \log_2 0.25 - 0.75 \times \log_2 0.75 = 0.81 \ bit/pattern$$

Loss of information

$$H_{loss} = H_{out} - H_{in} = 2.0 - 0.81 = 1.189 \ bit$$

*Method 2:*

The entropy of the input signal (probability $p = 0.707$)

$$H_{in} = -2(1-p) \times \log_2(1-p) - 2p \times \log_2 p$$

$$= -2(1-0.707) \times \log_2(1-0.707) - 2 \times 0.707 \times \log_2 0.707$$

$$= 1.745 \ bit$$

Output entropy

$$H_{out} = -0.707^2 \times \log_2 0.707^2$$

$$-(1-0.707^2) \times \log_2(1-0.707^2) = 0.804 \ bit$$

Loss of information

$$H_{loss} = H_{out} - H_{in}$$

$$= 1.745 - 0.804 = 0.941 \ bit$$

$j = x_1 x_2$

|        | $x_1$ | $x_2$ | $f$ |
|--------|-------|-------|-----|
| Pattern 1 | 0 | 0 | 0 |
| Pattern 2 | 0 | 1 | 0 |
| Pattern 3 | 1 | 0 | 0 |
| Pattern 4 | 1 | 1 | 1 |

**Figure 36.** Information measures of AND functions of two variables.

The output of the AND function is equal to 1 with probability $p^2$, and equal to 0 with probability $1 - p^2$. Hence, the entropy of the output signal is

$$H_{out} = -p^2 \times \log_2 p^2 - (1-p)^2 \times \log_2 (1-p)^2 \ \text{bit}$$

The maximum value of the output entropy is equal to 1, when $p = 0.707$. Hence, the input entropy of the AND function is 0.745 bit (Fig. 36). We observe that in the case of noncorrelated signals, information losses are less.

## 7.6. Measures in Decision Trees

In this section, we address the design of decision trees with nodes of three types: Shannon ($S$), positive Davio ($pD$), and negative Davio ($nD$) based on the information theoretical approach. An approach revolves around choosing the "best" variable and the "best" expansion type with respect to this variable for any node of the decision tree in terms of information measures. This means that in any step of the decision making strategy, we have an opportunity to choose both the variable and the type of expansion based on the criterion of minimum entropy. The entropy-based optimization strategy can be described as the generating of the optimal paths in a decision tree, with respect to the minimum entropy criterion [96–98, 106].

Calculation of entropy and information on decision trees is best known as the induction of decision trees (ID3) algorithm for optimization.

EXAMPLE 39. *Figure 37 illustrates the calculation of entropy on a decision tree.*

In the process of decision tree design two information measures are used:

$$\langle \text{Conditional entropy} \rangle = H(f|\text{Tree})$$

$$\langle \text{Mutual information} \rangle = I(f; \text{Tree})$$

The initial state of this process is characterized by the maximum value for the conditional entropy:

$$H(f|\text{Tree}) = H(f)$$

**Figure 37.** Measure of entropy on a decision tree (example 39).

Nodes are recursively attached to the decision tree by using the top-down strategy. In this strategy, the entropy $H(f|\text{Tree})$ of the function is reduced, and the information $I(f;\text{Tree})$ increases, because the variables convey the information about the function. Each intermediate state can be described in terms of entropy by the equation

$$I(f;\text{Tree}) = H(f) - H(f|\text{Tree}) \tag{13}$$

We maximize the information $I(f;\text{Tree})$ that corresponds to the minimization of entropy $H(f|\text{Tree})$, in each step of decision tree design. The final state of the decision tree is characterized by $H(f|\text{Tree}) = 0$ and $I(f;\text{Tree}) = H(f)$, that is, Tree represents the switching function $f$.

The decision tree design process is a recursive decomposition of a switching function. A step of this recursive decomposition corresponds to the expansion of switching function $f$ with respect to the variable $x$. Assume that the variable $x$ in $f$ conveys information that is, in some sense, the rate of influence of the input variable on the output value for $f$.

### 7.6.1. Information Notation of S, pD, and nD Expansion

The designed decision tree based on the $S$ expansion is mapped into a sum-of-products expression as follows: a leaf with the logic value 0 is mapped into $f = 0$, and with the logic value 1 into $f = 1$; a nonterminal node is mapped into $f = \bar{x} \cdot f_{|x=0} \vee x \cdot f_{|x=1}$. The information measure of $S$ expansion for a switching function $f$ with respect to the variable $x$ is represented by the equation

$$H^S(f|x) = p_{x=0} \cdot H(f_{x=0}) + p_{x=1} \cdot H(f_{|x=1}) \tag{14}$$

The information measure of $S$ expansion is equal to the conditional entropy $H(f|x)$ (Fig. 38):

$$H^S(f|x) = H(f|x) \tag{15}$$

The information measure of $pD$ expansion of a switching function $f$ with respect to the variable $x$ is represented by

$$H^{pD}(f|x) = p_{x=0} \cdot H(f_{|x=0}) + p_{|x=1} \cdot H(f_{x=0} \oplus f_{x=1}) \tag{16}$$

The information measure of the $nD$ expansion of a switching function $f$ with respect to the variable $x$ is

$$H^{nD}(f|x) = p_{x=1} \cdot H(f_{x=1}) + p_{x=0} \cdot H(f_{|x=0} \oplus f_{x=1}) \tag{17}$$

Shannon expansion $f = \bar{x} \cdot f|_{x=0} \oplus x \cdot f|_{x=1}$
Information-theoretical notation

$$H^S(f|x) = p|_{x=0}H(f|_{x=0}) + p|_{x=1}H(f|_{x=1})$$

Left leaf · Right leaf

Positive Davio expansion $f = f|_{x=0} \oplus x \cdot (f|_{x=0} \oplus f|_{x=1})$
Information-theoretical notation

$$H^{pD}(f|x) = p|_{x=0}H(f|_{x=0}) + p|_{x=1}H(f|_{x=0} \oplus f|_{x=1})$$

Left leaf · Right leaf

Negative Davio expansion $f = f|_{x=1} \oplus \bar{x} \cdot (f|_{x=0} \oplus f|_{x=1})$
Information-theoretical notation

$$H^{nD}(f|x) = p|_{x=1}H(f|_{x=1}) + p|_{x=0}H(f|_{x=0} \oplus f|_{x=1})$$

Left leaf · Right leaf

Figure 38. Shannon and Davio expansions and their information measures for a switching function.

## 7.6.2. Optimization of Variable Ordering in a Decision Tree

The entropy-based optimization of decision tree design can be described as the optimal (with respect to the information criterion) node selection process. A path in the decision tree starts from a node and finishes in a terminal node. Each path corresponds to a term in the final expression for $f$.

EXAMPLE 40. *Design of the Shannon tree based on a sum-of-products expression given the hidden weighted bit function. The Shannon tree is shown in Fig. 39.*

## 7.7. Information Measures in the $N$-Hypercube

It has been shown that information-theoretical measures for logic networks can be evaluated by decision trees. In this section we focus on the details of information measures in $N$-hypercube based on information measures in decision trees.

A useful property of an $N$-hypercube is that compared with decision trees and diagrams is that it is possible to obtain information measure without recalculation after changing the order of variables. The example below illustrates this property.

EXAMPLE 41. *Figure 40 illustrates the calculation of entropy on an $N$-hypercube. Starting with the root, where entropy is maximal, we approach variables in sequence. Approaching $x_1$ reveals information about this variable. Approaching terminal nodes means that the entropy becomes 0.*



Entropy of the function

$H(f) = -1/2 \cdot \log 1/2 - 1/2 \cdot \log 1/2$

$= 1$ bit/pattern

Conditional entropy with respect to the variable $x_1$

$H(f|x_1) = -3/8 \cdot \log 3/8 - 1/8 \cdot \log 1/8$
$\qquad -1/8 \cdot \log 1/8 - 3/8 \cdot \log 3/8$

$= 0.81$ bit/pattern

Sum-of-products expression

$f = \bar{x}_1 \cdot x_1 \vee x_1 \cdot x_2$

| $x_1$ | $x_2$ | $x_3$ | $f$ |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |

Figure 39. Shannon decision tree design (example 40).

Figure 40. Measure of entropy on an .\-hypercube (example 41).

Here we suppose that input patterns are generated with equal probabilities. An alternative approach is based on the calculation of input and output entropy assuming that input patterns are generated with different probabilities.

## 8. CONCLUSIONS

We witness a "transition" from state-of-the-art design methods in two-dimensional space to three-dimensional space of parallel and distributed ultrasmall devices, or nanodevices. Several physical constraints of nanotechnologies considered in the introduction motivate the study of 3D structures and algorithms for Boolean function manipulation and calculation. The data structure used to represent a switching function is an important factor in mapping the functional description at physical structure topology. This is because at a nanoscale the traditional steps of logic design may be merged toward explicit mapping of the data structure and algorithms to each form of data structure, corresponding to a certain level of abstraction, and useful at certain phases of logic design.

Technological restrictions imply that nanoelements compute elementary logic functions with some probability. This fact addresseses the problem of the reliability of computing using nonreliable elements. There are many methods in state-of-the-art of logic design that deal with computation with unreliable elements and noisy signals, in particular, methods based on: probabilistic description of signals, stochastic pulse stream organization, fuzzy logic, probabilistic logic, residue number system, paradigms inspired by biological systems, stochastic logic neural networks, Von Neumann multiplexing, $R$-fold modular, and error-correcting codes.

Several particular results of the study on logic design in nanodimensions can be summarized as follows:

1. Graph-based data structures are the "bridge" between logic design and the 3D topology of nanomaterials. Among them, hypercube topology is a useful model of computing in spatial dimensions. This topology can be used at all levels of abstraction from the gate level (nodes implement the simplest logic functions) to the macrolevel (nodes implement complex devices).

2. Treelike and hypercube-like topology is common in parallel and distributed architectures; this fact reflects that the principle of an optimal computing scheme is being preserved while scaling it down to a molecular/atomic structure.

3. Fault-tolerance computing is the central problem of nanosystem design because of the probabilistic nature of nanodevices.

4. In 3D structures on the molecular/atomic level, information carriers become compatible with partially distributed sources/receivers/transmitters, and so measures that are inherent to the nature of information processing on a nanoscale level are required.

5. The information content of a logic function is a natural attribute of the function and it is technology-independent. The information content defines the complexity of function implementation, and thus can be used to estimate a lower bound on some physical (topological) parameters with respect to various implementations. Thus, it captures the fundamental characteristic of logic function behavior. Entropy, as spatial measurement in $N$-hypercube space, can be viewed as a contribution to the information content, with respect to all nodes of the embedded decision diagram. Based on the information-theoretic approach, an arbitrary decision tree or decision diagram can be designed. In each step of the decision-making process, the variable and type of expansion is chosen based on the information estimations.

## ACKNOWLEDGMENTS

## REFERENCES

1. E. Özbay, D. M. Bloom, D. H. Chow, and J. N. Schulman, *IEEE Electron Device Lett.* 14, 400 (1993).
2. K. K. Likharev, *Proc. IEEE* 87, 606 (1999).
3. C. P. Collier, E. W. Wong, M. Belohradsk, F. M. Raymo, J. F. Stoddart, P. J. Kuekes, R. S. Williams, and J. R. Heath, *Science* 285, 391 (1999).
4. V. Derycke, R. Martel, J. Appenzeller, and P. Avouris, *Nano Lett.* 1, 453 (2001).
5. J. C. Ellenbogen and J. C. Love, *Proc. IEEE* 88, 386 (2000).
6. M. A. Reed and J. M. Tour, *Sci. Am.* 86, 282 (2000).
7. J. M. Tour, W. L. V. Zandt, C. P. Husband, S. M. Husband, L. S. Wilson, P. D. Franzon, and D. P. Nackashi, *IEEE Trans. Nanotechnol.* 1, 100 (2002).
8. J. M. Tour, *Acc. Chem. Res.* 33, 791 (2000).
9. S. M. Goodnick and J. Bird, *IEEE Trans. Nanotechnol.* 2, 368 (2003).
10. A. Korotkov and K. Likharev, *J. Appl. Phys.* 84, 6114 (1998).
11. H. Lo, T. Spiller, and S. Popescu, "Introduction to Quantum Computation and Information." World Scientific, Hackensack, NJ, 1998.
12. J. Han and P. Jonker, *IEEE Trans. Nanotechnol.* 1, 201 (2002).
13. A. S. Sadek, K. Nikolić, and M. Forshaw, *Nanotechnology*, 15, 192 (2004).
14. A. DeHon, *IEEE Trans. Nanotechnol.* 2, 23 (2003).
15. M. P. Frank. Ph.D. Thesis. Massachusetts Institute of Technology, 1999.
16. J. P. Hayes and T. Mudge, Q. F. Stout, S. Colley, and J. A. Palmer, *IEEE Micro.* 6, 6 (1986).
17. W. D. Hillis, "The Connection Machine." MIT Press, Cambridge, MA. 1985.
18. C. L. Seitz, *Commun. ACM* 28, 22 (1985).
19. H. T. Kung and C. E. Leiserson, "in Introduction to VLSI Systems" (C. Mead and L. Conway, Eds.), p. 260. Addison-Wesley, Reading, MA. 1980.
20. Y. T. Lai and P. T. Wang, *IEEE Trans. VLSI Sys.* 5, 186 (1997).
21. T. Endoh, H. Sakuraba, K. Shinmei, and F. Masuoka, *Proc. 22nd Int. Conf. Microelectron.* 442, 447 (2000).
22. M. P. Frank, *Comput. Sci. Eng.* 4, 16 (2002).
23. N. Margolus and L. Levitin, *Physica D* 120, 188 (1998).
24. W. Smith, Technical report. NECI, 1995. http://www.neci.nj.nec.com/-homepages/wds/fundphys.ps
25. R. I. Greenberg, *IEEE Trans. on Comput.* 43, 1358 (1994).
26. C. H. Leiserson, *IEEE Trans. Comput.* 34, 892 (1985).
27. P. M. B. Vitanyi, "Proceedings of the 2nd IEEE Workshop on Physics and Computation, PhysComp'94." Dallas (Texas), 1994, p. 24.
28. G. Bilardi and F. P. Preparata. *Theory Comput. Syst.* 30, 523 (1997).
29. S. Hassoun and T. Sasao, (Eds.) "Logic Synthesis and Verification" (R. Brayton, Ed.). Kluwer Academic Publishers. Boston, 2002.
30. T. Sasao, "Switching Theory for Logic Synthesis." Kluwer Academic, 1999.
31. S. N. Yanushkevich, V. P. Shmerko, and S. E. Lyshevski, "Logic Design of NanoICs." CRC Press. Boca Raton, FL. 2005.
32. N. Asahi, M. Akazawa, and Y. Amemiya, *IEEE Trans. Electron Devices* 44, 1109 (1997).
33. S. Kasai and H. Hasegawa, *IEEE Electron Device Lett.* 23 (2002).
34. T. Yamada, Y. Kinoshita, S. Kasai, H. Hasegawa, and Y. Amemiya, *J. Appl. Phys.* 40, 4485 (2001).
35. K. Goser, C. Pacha, A. Kanstein, and M. L. Rossmann, *Proc. IEEE.* 85, 558 (1997).
36. R. Drechsler and W. Guenter "Towards One-Pass Synthesis." Kluwer Academic Publishers. Boston, 2002.
37. H. L. Chen and N. F. Tzeng, *IEEE Trans. Parallel Distributed Sys.* 8, 1171 (1997).
38. H. L. Chen and N. F. Tzeng, *IEEE Trans. Comput.* 46, 871 (1997).

39. G. Vichniac, *Physica D* 10, 96 (1984).

40. C. Gerousis, S. M. Goodnick, and W. Porod, *Int. J. Circuits Theory Appl.* 28, 523 (2000).

41. M. G. Ancona, *Superlattices Microstruct.* 20, 461 (1996).

42. E. F. Moore and C. E. Shannon, *J. Franklin Inst.* 262, 191 (1956); 281 (1956).

43. J. Von Neumann, "in Automata Studies" (C. E. Shannon and J. McCarthy, Eds.), Princeton University Press, Princeton, NJ, 1955.

44. J. R. Heath, P. J. Kuekes, G. S. Snider, and R. S. Williams, *Science* 280, 1716 (1998).

45. B. R. Gaines, "in Advances in Information Systems Science" (J. T. Tou, Ed.), Vol. 2, Chap. 2, p. 37, Plenum, New York, 1969.

46. V. V. Yakovlev and R. F. Fedorov, "Stochastic Computing," Mashinostroenie Publishers, Moscow, 1974.

47. W. S. Evans and L. J. Schulman, *IEEE Trans. Inform. Theory*, 45, 2367 (1999).

48. S. Winograd and J. D. Cowan, "Reliable Computation in the Presence of Noise," MIT Press, Cambridge, MA, 1963.

49. M. P. Frank, Nanocomputers: theoretical models, In "Encyclopedia of Nanoscience and Nanotechnology" (Hari Singh Nalwa, Ed.), pp. 249–300. American Scientific Publishers, Stevenson Ranch, CA, 2004.

50. M. P. Frank and T. F. Knight, Jr., *Nanotechnology* 9, 162 (1998).

51. C. Shannon *Bell Syst. Tech. J.* 27, 379 623 (1948).

52. S. N. Yanushkevich, "Logic Differential Caluclus in Multi-Valued Logic Design." Technical University of Szczecin Academic Publishers, Poland, 1998.

53. S. B. Akers, *IEEE Trans. Comput.* C 27, 509 (1978).

54. R. E. Bryant, *IEEE Trans. Comput.* C 35, 677 (1986).

55. C. Meinel and T. Theobald, "Algorithms and Data Structures in VLSI Design." Springer, 1998.

56. S. Minato, "Binary Decision Diagrams and Applications for VLSI Design." Kluwer Academic Publishers, Dordrecht, 1996.

57. R. S. Stankovic and J. T. Astola, "Spectral Interpretation of Decision Diagrams." Springer, 2003.

58. B. Becker, *Integration* 26, 20 (1998).

59. R. E. Bryant and Y.-A. Cen, "Proceedings of 32rd ACM/IEEE Design Automation Conference." 1995, p. 535.

60. N. Takahashi, N. Ishiura, and S. Yajima, "Proceedings IEEE/ACM International Conference on Computer-Aided Design." 1991, p. 550.

61. P. Dziurzanski, V. Malyugin Shmerko, and S. Yanushkevich, *Automation and Remote Control* 63, 960 (2002).

62. S. J. Friedman and K. J. Supowit, *IEEE Trans. Comput.* 39, 710 (1990).

63. R. Rudell, "Proceedings the International Conference on Computer-Aided Design." Santa Clara, CA, 1993, p. 42.

64. J. P. Roth, *Trans. Am. Math. Soc.* 88, 301 (1958).

65. F. R. Preparata and J. Vuillemin. *Commun. ACM* 24, 300 (1981).

66. T. R. N. Rao, "Error Coding for Arithmetic Processors." Academic Press, New York, 1974.

67. S. B. Akers, *Soc. Ind. Appl. Math.* 7, 487 (1959).

68. D. Bochmann and Ch. Posthoff, "Binäre Dynamishe Systeme." Akademieverlag, Berlin, 1981.

69. M. J. Davio, P. Deschamps, and A. Thayse, "Discrete and Switching Functions." McGraw-Hill, New York, 1978.

70. F. F. Sellers, M. Y. Hsiao, and L. W. Bearson, *IEEE Trans. Comput.* 1, 676 (1968).

71. K. P. Parker and E. J. McCluskey, *IEEE Trans. Comput.* 24, 573 (1975).

72. K. P. Parker and E. J. McCluskey, *IEEE Trans. Comput.* 24, 668 (1981).

73. V. Beiu, J. M. Quintana, and M. J. Avedillo, *IEEE Trans. Neural Networks* 14, 1217 (2003).

74. A. H. Aguirre and C. A. Coello, in "Artificial Intelligence in Logic Design" (S. N. Yanushkevich, Ed.), pp. 285–311. Kluwer, Dordrecht, 2004.

75. V. A. Cheushev, S. N. Yanushkevich, V. P. Shmerko, C. Moraga, and J. Kolodziejczyk, "Proceedings IEEE 31st International Symposium on Multiple-Valued Logic." 2001, p. 201.

76. H. Iba, M. Iwata, and T. Higuchi, "in Lecture Notes in Computer Science." Springer, Heidelberg, 1997, Vol. 1259, p. 327.

77. T. Luba, C. Moraga, S. N. Yanushkevich, M. Opoka, and V. P. Shmerko, "Proceeding IEEE 30th International Symposium on Multiple-Valued Logic." 2000, pp. 253–258.

78. F. N. Marinos, *IEEE Trans. Comput.* 18, 343 (1969).

79. L. N. de Castro and J. I. Timmis, "Artificial Immune Systems: A New Computational Intelligence Approach." Springer, London, 2002.

80. F. N. Najm, *IEEE Trans. VLSI* 2, 446 (1994).

81. J. Von Neumann "The Theory of Self-Reproducing Automata." University of Illinois Press, Urbana, 1966.

82. S. Mitra, N. R. Saxena, and E. J. McCluskey, *IEEE Trans. Reliability* 49, 285 (2000).

83. G. Türel and K. Likharev, *Int. J. Circuit Theory Appl.* 31, 37 (2003).

84. C. Pacha, U. Auer, C. Burwick, P. Glosekoter, K. Goser, W. Prost, A. Brennemann, and F.-J. Tegude, *IEEE Trans. VLSI Syst.* 8, 558 (2000).

85. F. Rodriguez and H. C. Tuckwell, *Biosystems* 48, 187 (1998).

86. S. Amarel and J. A. Brzozowski, in "Redundancy Techniques for Computing Systems" (R. H. Wilcox and W. C. Mann, Eds.). Spartan Books, Washington DC. 1962.

87. A. Majumdar and S. B. K. Vrudhula, *IEEE Trans. VLSI* 1, 365 (1993).

88. T. Peper, J. Lee, F. Abo, T. Isokawa, S. Adachi, N. Matsui, and S. Mashiko, *IEEE Trans. Nanotechnol.* 3, 187 (2004).

89. R. I. Bahar, J. Chen. and J. Mundy. in "Nano Quantum and Molecular Computing: Implications to High Level Design and Validation" (S. Shukla and R. I. Bahar. Eds.). Kluwer. 2004.

90. N. Gershenfeld, *IBM Sys. J.* 35, 577 (1996).

91. V. Cheushev, V. Shmerko, D. Simovici, and S. Yanushkevich. "Proceedings IEEE 28th International Symposium on Multiple-Valued Logic." Japan, 1998. p. 357.

92. R. M. Goodman and P. Smyth. *IEEE Trans. Inf. Theory* 34, 979 (1988).

93. C. R. P. Hartmann, P. K. Varshney. K. G. Mehrotra, and C. L. Gerberich. *IEEE Trans. Inf. Theory* 28, 565 (1982).

94. V. A. Cheushev, S. N. Yanushkevich, C. Moraga, and V. P. Shmerko, in "Research Report 741." Forschungs-bericht, University of Dortmund, Germany. 2000.

95. V. Agraval. *IEEE Trans. Comput.* 30, 582 (1981).

96. S. Ganapathy and V. Rajaraman, *Commun. ACM* 16, 532 (1973).

97. R. M. Goodman and P. Smyth. *IEEE Trans. Inf. Theory* 34, 979 (1988).

98. C. R. P. Hartmann, P. K. Varshney. K. G. Mehrotra, and C. L. Gerberich, *IEEE Trans. Inf. Theory* 28, 565 (1982).

99. P. Varshney, C. Hartmann, and J. De Faria. *IEEE Trans. Comput.* 31, 164 (1982).

100. D. Marculescu, R. Marculesku, and M. Pedram, *IEEE Trans. Comput. Aided Design Integr. Circuits Syst.* 15, 599 (1996).

101. S. Ramprasad, N. R. Shanbhag, and I. N. Hajj, *IEEE Trans. VLSI Syst.* 7, 359 (1999).

102. A. Tyagi, *IEEE Trans. VLSI Syst.* 5, 456 (1997).

103. H. Watanabe, *IEICE Trans. Fund.* E76-A, 265 (1993).

104. R. W. Cook and M. J. Flynn, *IEEE Trans. Comput.* 22, 823 (1973).

105. N. Zhong and S. Ohsuga, *Trans. Inf. Process. Soc. Jpn.* 38, 687 (1997).

106. V. Shmerko, D. V. Popel, R. S. Stankovic, V. A. Cheushev, and S. Yanushkevich. "Proceedings IEEE International Conference on Telecommunications." Yugoslavia, 1999, p. 444.

# CHAPTER 16

# Nanoelectromechanical Systems and Modeling

## Changhong Ke, Horacio D. Espinosa

Department of Mechanical Engineering, Northwestern University, Illinois, USA

## CONTENTS

## 1. INTRODUCTION

Nanoelectromechanical systems (NEMS) are made of electromechanical devices that have critical dimensions from hundreds to a few nanometers. By exploring nanoscale effects, NEMS present interesting and unique characteristics, which deviate greatly from their predecessor microelectromechanical systems (MEMS). For instance, NEMS-based devices can have fundamental frequencies in microwave range ($\sim$100 GHz) [1]; mechanical quality factors in the tens of thousands, meaning low-energy dissipation; active mass in the femtogram range; force sensitivity at the attonewton level; mass sensitivity up to attogram [2] and sub-attogram [3] levels; heat capacities far below a "yoctocalorie" [4]; power consumption in the order of 10 attowatts [5]; extreme high integration level, approaching $10^{12}$ elements per square centimeter [1]. All these distinguished properties of NEMS devices pave the way to applications such as force sensors, chemical sensors, biological sensors, and ultrahigh-frequency resonators.

The interesting properties of the NEMS devices typically arise from the behavior of the active parts, which, in most cases, are in the forms of cantilevers or doubly clamped beams with dimensions at nanometer scale. The materials for those active components include

Figure 1. SEM Image of an undercut Si beam. with length of 7.7 $\mu$m. width of 0.33 $\mu$m and height of 0.8 $\mu$m. Reprinted with permission from [6]. A. N. Cleland and M. L. Roukes, *Appl. Phys. Lett.* 69, 2653 (1996). © 1996. American Institute of Physics.

silicon and silicon carbide, carbon nanotubes, and gold and platinum, to name a few. Silicon is the basic material for integrated circuit (IC) technology during the past few decades, and MEMS, and is widely used to build NEMS. Figure 1 is a scanning electron microscopy (SEM) image of a double-clamped resonator fabricated from a bulk, single-crystal silicon substrate [6]. However, ultrasmall silicon-based NEMS fail to achieve desired high-quality factors because of the dominance of surface effects, such as surface oxidation and reconstruction, and thermoelastic damping. Limitations in strength and flexibility also compromise the performance of silicon-based NEMS actuators. Instead, carbon nanotubes (CNTs) can well represent the ideas of NEMS, given their nearly one-dimensional structures with high-aspect ratio, perfect terminated surfaces, and excellent electrical and mechanical properties. Because of significant advances in growth, manipulation, and knowledge of electrical and mechanical properties, carbon nanotubes have become the most promising building blocks for the next generation of NEMS. Several carbon nanotube-based functional NEMS devices have been reported so far [1, 7–12]. Similar to carbon nanotubes, nanowires (NWs) are another type of one-dimensional novel nanostructure for building NEMS because of their size and controllable electrical properties.

This chapter provides a comprehensive review of NEMS devices to date and summarizes the modeling currently being pursued to gain insight into their performance. This chapter is organized as follows: in the first part, we review the carbon nanotubes and carbon nanotube-based NEMS. We also discuss nanowire-based NEMS. In the second part, we present the modeling of NEMS, including multiscale modeling and continuum modeling.

## 2. NANOELECTROMECHANICAL SYSTEMS

### 2.1. Carbon Nanotubes

Carbon nanotubes exist as a macromolecule of carbon, analogous to a sheet of graphite rolled into a cylinder. They were discovered by Sumio Iijima in 1991 and are a subset of the family of fullerene structures [13]. The properties of the nanotubes depend on the atomic arrangement (how sheets of graphite are rolled to form a cylinder), their diameter, and their length. They are light, stiff, flexible, thermally stable, and chemically inert. They have the ability to be either metallic or semiconducting depending on the "twist" of the tube, which is called "chirality" or "helicity." Nanotubes may exist as either single-walled or multiwalled structures. Multiwalled carbon nanotubes (MWNTs) (Fig. 2(B)) are simply composed of multiple concentric single-walled carbon nanotubes (SWNTs) (Fig. 2(A)) [14]. The spacing between the neighboring graphite layers in MWNTs is ~0.34 nm. These layers interact with each other via van der Waals forces.

The methods to synthesize carbon nanotubes include electric arc-discharge [15, 16], laser ablation [17], and catalytic chemical vapor deposition (CVD) methods [18]. During synthesis, nanotubes are usually mixed with residues, including various types of carbon particles.

Figure 2. High-resolution transmission electron microscopy image of typical single-walled carbon nanotubes (SWNT) (A) and multi-walled carbon nanotubes (MWNT) (B). Reprinted with permission from [14], P. Ajayan. *Chem. Rev.* 99, 1787 (1999). © 1999, American Chemical Society.

For most applications and tests, a purification process is required. In one of the most common approaches, nanotubes are ultrasonically dispersed in a liquid (e.g., isopropanol) and the suspension is centrifuged to remove large particles. Other methods, including dielectrophoretic separation, are being developed to provide improved yield.

The mechanical and electrical properties of carbon nanotubes have been under intensive study during the past decade. Qian et al. [19] contributed a comprehensive review article, "Mechanics of Carbon Nanotube," from the perspective of both experimentation and modeling. The electronics of carbon nanotubes is extensively reviewed by McEuen et al. [20]. Besides, the study of the coupled electromechanical properties, which are essential to NEMS, is rapidly progressing. Some interesting results have been reported regarding the fact that the electrical properties of carbon nanotubes are sensitive to the structure variation and can be changed dramatically because of the change of the atomic bonds induced by mechanical deformations. It is known that carbon nanotubes can even change from metallic to semiconducting when subjected to mechanical deformation [21–23].

## 2.2. Fabrication Methods

The fabrication processes of NEMS devices can be categorized according to two approaches. *Top-down* approaches, that evolved from manufacturing of MEMS structures, use submicron lithographic techniques, such as electron-beam lithography, to fabricate structures from bulk materials, either thin films or bulk substrates. *Bottom-up* approaches fabricate the nanoscale devices by sequentially assembling of atoms and molecules as building blocks. Top-down fabrication is size limited by facts such as the resolution of the electron-beam lithography, etching-induced roughness, and synthesis constraints in epitaxially grown substrates. Significant interest has been shown in the integration of nanoscale materials such as carbon nanotubes and nanowires, fabricated by bottom-up approaches, to build nanodevices. Most of the nanodevices reported so far in the literature are obtained by "hybrid" approaches, that is, combination of bottom-up (self assembly) and top-down (lithographic) approaches [24].

One of the key and most challenging issues of building carbon nanotube–based or nanowire-based NEMS is the positioning of nanotubes or nanowires at the desired locations with high accuracy and high throughput. Reported methods of manipulation and positioning of nanotubes are briefly summarized in the following section.

### 2.2.1. Random Dispersion Followed by E-Beam Lithography

After purification, a small aliquot of a nanotube suspension is deposited onto a substrate. The result is nanotubes randomly dispersed on the substrate. Nanotubes on the substrate are imaged inside a scanning electron microscope (SEM) and then this image is digitized and imported to a mask-drawing software, where the mask for the subsequent electron-beam lithography is designed. In the mask layout, pads are designed to superimpose over the carbon nanotubes. Wet etching is employed to remove the material under the carbon nanotubes to form freestanding nanotube structures. This process requires an alignment capability in the lithographic step with an accuracy of 0.1 $\mu$m or better. This method was firstly employed to make nanotube structures for mechanical testing [25, 26]. The reported NEMS devices based on this method include nanotube–based rotational actuators [9] and nanowire-based resonators [24].

## 2.2.2. Nanomanipulation

Manipulation of individual carbon nanotubes using piezo-driven manipulators inside electron microscope chambers is one of the most commonly used methods to build NEMS [8] and structures for mechanical testing [27–32]. In general, the manipulation and positioning of nanotubes is accomplished in the following manner: (1) a source of nanotubes is positioned close to the manipulator inside the microscope; (2) the manipulator probe is moved close to the nanotubes under visual surveillance of the microscope monitor until a protruding nanotube is attracted to the manipulator due to either van der Waals forces or electrostatic forces; (3) the free end of the attracted nanotube is "spot welded" by the electron-beam-induced deposition (EBID) of hydrocarbon [8, 31] or metals, like platinum [32] from adequate precursor gases.

Figure 3 shows a three-dimensional nanomanipulator (Klocke Nanotechnik Co.) having the capability of moving in X, Y, and Z directions with nanometer displacement resolution. The manipulation process of an individual carbon nanotube is illustrated in Fig. 4(A)–4(C).

## 2.2.3. External Field Alignment

DC/AC electric fields have been successfully used in the manipulation of nanowires [33], nanotubes [34, 35], and bioparticles [36–39]. Microfabricated electrodes are typically used to create an electric field in the gap between them. A droplet containing carbon nanotubes in suspension is dispensed into the gap with a micropipette. The applied electric field aligns the nanotubes, due to the dielectrophoretic effect, which results in the bridging of the electrodes by a single nanotube. The voltage drop that arises when the circuit is closed (DC component) ensures the manipulation of only one nanotube. Besides, AC dielectrophoresis has been employed to successfully separate metallic from semiconducting single-walled carbon nanotubes in suspension [40]. NEMS devices fabricated using this method include nanotube–based nanorelays [41].

Huang et al. [42] demonstrated another method for aligning nanowires. A laminar flow was employed to achieve preferential orientation of nanowires on chemically patterned surfaces. This method was successfully used in the alignment of silicon nanowires. Magnetic fields have also been used to align carbon nanotubes [43].

## 2.2.4. Direct Growth

Instead of manipulating and aligning carbon nanotubes after their manufacturing, researchers have also examined methods for controlled direct growth. Huang et al. [44] used the microcontact printing technique to directly grow aligned nanotubes vertically. Dai et al.



Figure 3. Klocke Nanotechnik nanomanipulator possessing nanometer resolution in the x, y, and z axes.

Figure 4. SEM images of the manipulation of carbon nanotubes using the three-dimensional Klocke Nanotechnik nanomanipulator. (A) Manipulator probe is approaching a protruding nanotube. The sample is dried nanotube solution on top of a TEM copper grid. (B) Manipulator probe makes contact with the free end of the nanotube and the nanotube is welded to the probe by EBID of platinum. (C) A single nanotube mounted to the manipulator probe.

[45–48] reported several patterned growth approaches developed in his group. The idea is to pattern the catalyst in an arrayed fashion and control the growth of carbon nanotubes from specific catalytic sites. The authors successfully carried out patterned growth of both MWNTs and SWNTs and exploited methods including self-assembly and external electric-field control. Figure 5 shows a SEM image of suspended single-walled nanotubes grown by electric-field-directed CVD method [47]. The carbon nanotube–based tunable oscillators, reported in Ref. [12], were fabricated using this method.

## 2.2.5. Self-Assembly

Self-assembly is a method of constructing nanostructures by forming stable bonds between the organic or nonorganic molecules and substrates. Recently, Rao et al. [49] reported an approach in large-scale assembly of carbon nanotubes with high throughput. Dip Pen Nanolithography (DPN), a technique invented by Mirkin's group [50], was employed to functionalize the specific surface regions either with polar chemical groups such as amino ($-NH_2/$ $-NH_3^+$) or carboxyl ($-COOH/-COO^-$), or with nonpolar groups such as methyl ($-CH_3$). When the substrate with functionalized surfaces was introduced into a liquid suspension of carbon nanotubes, the nanotubes were attracted toward the polar regions and self-assembled



Figure 5. Electric-field-directed freestanding single-walled nanotubes. Reprinted with permission from [47], Y. Zhang et al., *Appl. Phys. Lett.* 79, 3155 (2001). © 2001, American Institute of Physics.

to form predesigned structures, usually within 10 s, with a yield higher than 90%. The reported method is scalable to large arrays of nanotube devices by using high-throughput patterning methods such as photolithography, stamping, or massively parallel DPN.

## 2.3. Inducing and Detecting Motion

For nanostructures, both inducing and detecting motion are challenging. Some of the methods routinely used in MEMS face challenges when the size shrinks from microscale to nanoscale. For example, optical methods, such as simple-beam deflection schemes or more sophisticated optical and fiber-optical interferometry—both commonly used in scanning probe microscopy to detect the deflection of the cantilevers—generally fall beyond the diffraction limit, which means these methods cannot be applied to objects with cross-section much smaller than the wavelength of light [51].

### 2.3.1. Inducing Motion

Similar to MEMS, electrostatic actuation of nanostructures by an applied electrical field is commonly used for the actuation of NEMS (e.g., nanotweezers [7, 8]). The Lorenz force has been used to move small conducting beams [6, 24, 52], with alternating currents passing through them in a strong transverse magnetic field. The induced electromotive force, or voltage, can be detected as a measure of the motion. This method requires a fully conducting path and works well with a beam clamped at both ends [53]. Other actuation methods include piezoelectric actuation, thermal actuation using bilayers of materials with different thermal expansion, thermal in-plane actuation due to a specially designed topography [54], and scanning tunneling microscope (STM) [55].

### 2.3.2. Detecting Motion

The most straightforward method is by direct observation of the motion under electron microscopes [7, 8, 56, 57]. This visualization method, typically with resolution in the nanometer scale, projects the motion in the direction to be perpendicular to the electron beam. Limitations in depth of focus requires that the nano-object motion be primarily in a plane, which normally is coaxial with the electron beam. Electron tunneling is a very sensitive method that can detect subnanometer motion by the exponential dependence of the electron tunneling current on the separation between tunneling electrodes. Therefore, this technique is widely used in NEMS motion detection [5, 12]. Magnetomotive detection is a method based on the presence of an electrostatic field, either uniform or spatially inhomogeneous, through which a conductor is moved. The time-varying flux generates an induced electromotive force in the loop, which is proportional to the motion [24, 52, 58–60]. The displacement detection sensitivity of this technique is less than 1 Å [61]. It is known that carbon nanotubes can act as transistors; as such they can be used to sense their own motion [12, 62]. Capacitance sensors have been widely used in MEMS. They can also be used in NEMS motion sensing with a resolution of a few nanometers [54], and the resolution can be potentially increased to Angstrom range provided that the capacitance measurement can be improved by one order of magnitude.

## 2.4. Functional Nanoelectromechanical Systems Devices

In this section, we review the carbon nanotube- or nanowire-based NEMS devices reported in the literature with a special emphasis on fabrication methods, working principles, and applications.

### 2.4.1. Carbon Nanotube-Based Nanoelectromechanical Systems Devices

#### 2.4.1.1. Nonvolatile Random Access Memory    A carbon nanotube–based nonvolatile random access memory (NRAM) reported by Rueckes et al. [1] is illustrated in Fig. 6(A). The device is a suspended SWNT crossbar array for both I/O and switchable, bistable device elements with well-defined OFF and ON states. This crossbar consists of a set of parallel SWNTs or nanowires (lower) on a substrate composed of a conducting layer (e.g., highly

**823**

(B)



Figure 6. Schematics of freestanding nanotube device architecture with multiplex addressing (A). Three-dimensional view of a suspended crossbar array showing four junctions with two elements in the ON (contact) state and two elements in the OFF (separated) state. (B) Top view of an $n \times m$ device array. Reprinted with permission from [1], T. Rueckes et al., *Science* 289, 94 (2000). © 2000, American Association for the Advancement of Science.

doped silicon [dark gray]) that terminates in a thin dielectric layer (e.g., $SiO_2$ [light gray]) and a set of perpendicular SWNTs (upper) that are suspended on a periodic array of inorganic or organic supports. Each nanotube is contacted by a metal electrode. Each cross-point in this structure corresponds to a device element with a SWNT suspended above a perpendicular nanoscale wire. Qualitatively, bistability can be envisioned as arising from the interplay of the elastic energy and the van der Waals energy when the upper nanotube is freestanding or the suspended SWNT is deflected and in contact with the lower nanotube. Because the nanotube junction resistance depends exponentially on the separation gap, the separated upper-to-lower nanotube junction resistance will be orders of magnitude higher than that of the contact junction. Therefore, two states—OFF and ON—are well defined. For a device element, these two states can be read easily by measuring the resistance of the junction and, moreover, can be switched between OFF and ON states by applying voltage pulses to nanotubes at corresponding electrodes to produce attractive or repulsive electrostatic forces. A key aspect of this device is that the separation between top and bottom conductors must in the order of 10 nm. In such case, the van der Waals energy overcomes the elastic energy when the junction is actuated (ON state) and remains on this state even if the electrical field is turned off (nonvolatile feature).

The concept of the bitable device was demonstrated by current-voltage (I–V) behaviors of suspended, crossed nanotube devices made from SWNT ropes ($\sim$50 nm in diameter), with junction gap $\sim$150 nm by mechanical manipulation under an optical microscope. The resistances of the OFF state of the devices were found consistently 10-fold larger than the ON state.

In the integrated system, electrical contacts are made only at one end of each of the lower and upper sets of nanoscale wires in the crossbar array, and thus, many device elements can be addressed from a limited number of contacts (see Fig. 6(B)). This approach suggests a highly integrated, fast, and macroscopically addressable NRAM structure that could overcome the fundamental limitations of semiconductor random access memory in

size, speed, and cost. Integration levels as high as $1 \times 10^{12}$ elements per square centimeter and switching time down to ~5 ps (200-GHz operation frequency) using 5-nm device elements and 5-nm supports are envisioned while maintaining the addressability of many devices through the long (~10-$\mu$m) SWNT wires. However, such small dimensions, in particular, the junction gap size, impose significant challenges in the nanofabrication of parallel device arrays.

**2.4.1.2. Nanotweezers** There are two types of carbon nanotube–based nanotweezers reported by Kim and Lieber in 1999 [7] and Akita et al. in 2001 [8], respectively. Both nanotweezers employ MWNTs as tweezers' arms that are actuated by electrostatic forces. The applications of these nanotweezers include the manipulation of nanostructures and two-tip STM or atomic force microscope (AFM) probes [7].

The fabrication process of the carbon nanotube–based nanotweezers reported by Kim and Lieber [7] is illustrated in Fig. 7(A). Freestanding electrically independent electrodes were deposited onto tapered glass micropipettes with end diameters of 100 nm (Fig. 7(B)). MWNT or SWNT bundles with diameters 20–50 nm were attached to the two gold electrodes, under the direct view of an optical microscope operated in dark-field mode, using an adhesive [63, 64]. A SEM image of fabricated nanotube tweezers is shown in Fig. 7(C).

The electromechanical response of nanotube nanotweezers was investigated by applying bias voltages to the electrodes while simultaneously imaging the nanotube displacements under an optical microscope in dark-field mode. As the bias voltage increased from 0 to 8.3 V (see Fig. 8), the free ends of the tweezers' arms bent closer to each other from their relaxed position (at 0 V). The tweezers' arms relaxed to the original position when the applied voltage was removed, and this process could be repeated more than 10 times, producing the same displacement each time within the optical microscope resolution limit. These results demonstrated that the mechanical response was elastic and thus that neither the nanotubes nor the nanotube-electrode junctions deform inelastically. At 8.3 V, the distance between the tweezers' ends decreased by about 50% of the initial value, and as the voltage was increased further to 8.5 V, the tweezers' arms suddenly closed (Fig. 8(E)).

The nanotube nanotweezers have been demonstrated successfully to manipulate nanostructures, such as fluorescently labeled polystyrene spheres and $\beta$-SiC nanocluster (see Fig. 9) and GaAs nanowires [7].



Figure 7. Overview of the fabrication process of carbon nanotube nanotweezers. (A) Schematic illustrating the deposition of two independent metal electrodes and the subsequent attachment of carbon nanotubes to these electrodes. (B) SEM image of the end of a tapered glass structure after the two deposition steps. Scale bar, 1 $\mu$m. The higher-resolution inset shows clearly that the electrodes are separated. Scale bar, 200 nm. (C) SEM image of nanotweezers after mounting two MWNT bundles on each electrode. Scale bar, 2 $\mu$m. Reprinted with permission from [7], P. Kim and C. M. Lieber, *Science* 126, 2148 (1999). © 1999, American Association for the Advancement of Science.

The carbon nanotube–based nanotweezers, reported by Akita et al. [8], are shown in Fig. 10. Commercially available Si AFM cantilevers were employed as the device body. A Ti/Pt film was coated on the tip of the cantilever and connected to three Al interconnects that were patterned on the cantilever by a conventional lithographic technique as shown in Fig. 10(A). The Ti/Pt film was separated into two by focused ion beam etching. These two parts were independently connected to Al interconnects as shown in Fig. 10(B). DC voltage was applied between the separated Ti/Pt tips, through the Al interconnects, to operate the tweezers after attaching two arms of nanotubes on them.

The attachment of the nanotubes was carried out in a specially designed field-emission-type SEM with three independent movable stages. The Si cantilever and the nanotube cartridge were mounted on two different stages. A third stage, where a tungsten needle was installed, was used for the fine adjustment of the position of the nanotubes after being mounted to the Si tip. When the metal-coated Si tip was manipulated to be in contact with a target nanotube, an amorphous carbon film was deposited on this contact portion by the electron-beam dissociation of contaminants, mainly hydrocarbons, in the SEM chamber. The target nanotube was finally pulled away from the cartridge. Another nanotube was also attached in the same manner. The position of the two nanotube arms was adjusted to be parallel by using the stage with the tungsten needle and fixed by deposition of the carbon films at the base of the arms. The nanotube arms were also coated with an ultrathin carbon film (a few to several nanometers) to achieve insulation from the outside. This film coating prevents large current flows when the two nanotube arms close or pick up a conductive particle. Figure 11(a) shows a SEM image of a typical pair of nanotube nanotweezers prepared this way. Two arms of the nanotubes were fixed at the most appropriate position on the Si tips. Their length was 2.5 μm, and the separation between their tips was 780 nm.

Figure 10. SEM images of a Si cantilever as a base for nanotube nanotweezers. (a) A Ti/Pt film was coated on the tip and connected to three Al lines patterned on the cantilever. (b) The Ti/Pt film was separated into two by a focused ion beam. and these two were connected to the one and two Al lines. respectively. Reprinted with permission from [8]. S. Akita et al., Appl. Phys. Lett. 79. 1691 (2001). © 2001. American Institute of Physics.

The operation of the nanotube nanotweezers was examined by in situ SEM. Various voltages were applied between the two arms to get them to close because of the electrostatic attraction force. Figures 11(b)–11(d) show the motion of the nanotube arms as a function of the applied voltage $V$. It is clearly seen that the arms bent and the separation between the tips decreased with increasing applied voltage. The separation became 500 nm at $V = 4$ V and zero at $V > 4.5$ V. It is noted that the motion in Figs. 11(a)–11(d) could be repeated many times without any permanent deformation, showing that carbon nanotubes are ideal materials for building NEMS.

**2.4.1.3. Rotational Motors** A carbon nanotube–based rotational motor, reported by Fennimore et al. in 2003 [9], is conceptually illustrated in Fig. 12(a). The rotational element (R), a solid rectangular metal plate serving as a rotor, is attached transversely to a suspended support shaft. The support shaft ends are embedded in electrically conducting anchors (A1, A2) that rest on the oxidized surface of a silicon chip. The rotor plate assembly is surrounded



Figure 11. SEM images of the motion of nanotube arms in a pair of nanotweezers as a function of the applied voltage. Reprinted with permission from [8]. S. Akita et al., Appl. Phys. Lett. 79. 1691 (2001). © 2001. American Institute of Physics.

**Figure 12.** Integrated synthetic nanoelectromechanical systems (NEMS) actuator. (a) Conceptual drawing of nanoactuator. (b) SEM image of nanoactuator just prior to HF etching. Scale bar, 300 nm. Reprinted with permission from [9], A. M. Fennimore et al., *Nature* 424, 408 (2003). © 2003, Nature Publishing Group.

by three fixed stator electrodes: two "in-plane" stators (S1, S2) are horizontally opposed and rest on the silicon oxide surface, and the third "gate" stator (S3) is buried beneath the surface. Four independent (DC and/or appropriately phased AC) voltage signals, one to the rotor plate and three to the stators, are applied to control the position, speed, and direction of rotation of the rotor plate. The key component in the assembly is a MWNT, which serves simultaneously as the rotor plate support shaft and the electrical feed-through to the rotor plate; most importantly, it is also the source of rotational freedom.

The nanoactuator was constructed using lithographic methods. MWNTs in suspension were deposited on a doped silicon substrate covered with 1 $\mu$m of SiO$_2$. The nanotubes were located using an AFM or a SEM. The remaining actuator components (in-plane rotor plate, in-plane stators, anchors, and electrical leads) were then patterned using electron-beam lithography. An HF (hydro fluoric acid) etch was used to remove roughly 500 nm of the SiO$_2$ surface to provide clearance for the rotor plate. The conducting Si substrate here serves as the gate stator. Figure 12(b) shows an actuator device prior to etching. Typical rotor plate dimensions were 250–500 nm on a side.

The performance of the nanoactuator was examined *in situ* inside the SEM chamber. Visible rotation could be obtained by applying DC voltages up to 50 V between the rotor plate and the gate stator. When the applied voltage was removed, the rotor plate would rapidly return to its original horizontal position. To exploit the intrinsic low-friction-bearing behavior afforded by the perfectly nested shells of MWNTs, the MWNT supporting shaft was modified *in situ* by successive application of very large stator voltages. The process resulted in fatigue and eventually shear failure of the outer nanotube shells. In the "free" state, the rotor plate was still held in position axially by the intact nanotube core shells but could be azimuthally positioned, using an appropriate combination of stator signals, to any arbitrary angle between 0° and 360°. Figure 13 shows a series of still SEM images, recorded from an actuated device in the free state, being "walked" through one complete rotor plate revolution using quasi-static DC stator voltages. The stator voltages were adjusted sequentially to attract the rotor plate edge to successive stators. By reversing the stator voltage sequence, the rotor plate rotation could be reversed in an equally controlled fashion. Finite frequency operation of the actuator was also performed using a variety of suitably phased AC and DC voltage signals to the three stators and rotor plate. The rotor plate was successfully flipped between the extreme horizontal (90° and 270°) positions. The experiments show that the MWNT clearly serves as a reliable, presumably wear-free, NEMS element providing rotational freedom. No apparent wear or degradation in performance was observed after many thousands of cycles of rotations.

The potential applications of the MWNT-based actuators include ultra-high-density optical sweeping and switching elements, paddles for inducing and/or detecting fluid motion in microfluidics systems, gated catalysts in wet chemistry reactions, biomechanical elements in biological systems, or general (potentially chemically functionalized) sensor elements.

**2.4.1.4. Nanorelays**  Carbon nanotube–based nanorelays were first reported by Kinaret et al. in 2003 [10] and later experimentally demonstrated by Lee et al. in 2004 [41]. The nanorelay is a three-terminal device including a conducting carbon nanotube placed on a

**Figure 13.** Series of SEM images showing the actuator rotor plate at different angular displacements. The schematic diagrams located beneath each SEM image illustrate a cross-sectional view of the position of the nanotube/rotor-plate assembly. Scale bar, 300 nm. Reprinted with permission from [9], A. M. Fennimore et al., *Nature* 424, 408 (2003). © 2003, Nature Publishing Group.

terrace in a silicon substrate and connected to a fixed-source electrode (S), as shown in Fig. 14(A). A gate electrode (G) is positioned underneath the nanotube so that a charge can be induced in the nanotube by applying a gate voltage. The resulting capacitance force $q$ between the nanotube and the gate bends the tube and brings the tube end into contact with a drain electrode (D) on the lower terrace, thereby closing an electric circuit. Theoretical modeling of the device shows that there is a sharp transition from a nonconducting (OFF) to a conducting (ON) state when the gate voltage is varied at a fixed source-drain voltage. The sharp switching curve allows for amplification of weak signals superimposed on the gate voltage [10].

One fabricated nanorelay device is shown in Fig. 14(B). A multiwalled nanotube was positioned on top of the source, gate, and drain electrodes with PMMA (polymethyl-methacrylate) as sacrificial layer using AC-electrophoresis techniques [34]. Then, a top electrode was placed over the nanotube at the source to ensure good contact. The underlying PMMA layer was then carefully removed to produce a nanotube suspended over the gate and drain electrodes. The separation between gate and drain was approximately 250 nm, and the source drain distance was 1.5 $\mu$m.

The electromechanical properties of nanotube relays were investigated by measuring the current-gate voltage ($I-V_{sg}$) characteristics, while applying a source-drain voltage of 0.5 V. Figure 15 shows the $I-V_{sg}$ characteristics of one of the nanotube relays with an initial height



**Figure 14.** Schematic diagrams of a CNT nanorelay device (A). Reprinted with permission from [10], J. Kinaret et al., *Appl. Phys. Lett.* 82, 1287 (2002). © 2002, American Institute of Physics. SEM image of a fabricated nanorelay device (B). Reprinted with permission from [41], S. Lee et al., *Nano Lett.* 4, 2027 (2004). © 2004, American Chemical Society.

**Figure 15.** I-V$_{sg}$ characteristics of a nanotube relay initially suspended approximately 80 nm above the gate and drain electrodes. Reprinted with permission from [41], S. Lee et al., *Nano Lett.* 4, 2027 (2004). © 2004, American Chemical Society.

difference between the nanotube and drain electrode of approximately 80 nm. The drain current started to increase nonlinearly when the gate voltage reached 3 V (at this gate voltage, the current is on the order of 10 nA). The nonlinear current increase was a signature of electron tunneling as the distance between the nanotube and the drain electrode was decreased. Beyond V$_{sg}$ = 20 V, there was a change in the rate of current increase. With the current increase rate becoming more linear, strong fluctuations could be detected. The deflection of the nanotube was found to be reversible. The current decreased with the reduction of gate voltage, showing some hysteresis, until it reached zero for a gate voltage below 3 V. The current measured during the increasing V$_{sg}$ part of the second scan closely followed that of the first scan, especially in the region below V$_{sg}$ = 12 V.

The dynamics of nanorelays was recently investigated by Jonsson et al. [65]. The results show that the intrinsic mechanical frequencies of nanorelays are in the gigahertz regime, and the resonance frequency can be tuned by the biased voltage.

The potential applications of nanorelays include memory elements, pulse generators, signal amplifiers, and logic devices.

### 2.4.1.5. Feedback-Controlled Nanocantilevers
A feedback-controlled carbon nanotube–based NEMS devices reported by Ke and Espinosa 2004 [11], schematically shown in Fig. 16, is made of a multiwalled carbon nanotube placed as a cantilever over a microfabricated step. A bottom electrode, a resistor, and a power supply are parts of the device circuit.



**Figure 16.** Schematic of nanotube-based device with tunneling contacts. Reprinted with permission from [11], C.-H. Ke and H. D. Espinosa, *Appl. Phys. Lett.* 85, 681 (2004). © 2004, American Institute of Physics.

When the applied voltage $U < V_{PI}$ (pull-in voltage), the electrostatic force is balanced by the elastic force from the deflection of the nanotube cantilever. The nanotube cantilever remains in the "upper" equilibrium position. When the applied voltage exceeds a pull-in voltage, the electrostatic force becomes larger than the elastic force and the nanotube accelerates toward the bottom electrode. When the tip of the nanotube is very close to the electrode (i.e., gap $\Delta \approx 0.7$ nm) as shown in Fig. 16, a substantial tunneling current passes between the tip of the nanotube and the bottom electrode. Because of the existence of the resistor R in the circuit, the voltage applied to the nanotube drops, weakening the electric field. Because of the kinetic energy of the nanotube, it continues to deflect downward, and the tunneling current increases, weakening the electric field further. In this case, the elastic force is larger than the electrostatic force, and the nanotube decelerates and eventually changes the direction of motion. This decreases the tunneling current and the electrical field recovers. If there is damping in the system, the kinetic energy of the nanotube is dissipated and the nanotube stays at the position where the electrostatic force is equal to the elastic force, and a stable tunneling current is established in the device. This is the "lower" equilibrium position for the nanotube cantilever. At this point, if the applied voltage U decreases, the cantilever starts retracting. When U decreases to a certain value, called pull-out voltage $V_{PO}$, the cantilever is released from its lower equilibrium position and returns back to its upper equilibrium position. At the same time, the current in the device diminishes substantially. Basically, the pull-in and pull-out processes follow a hysteretic loop for the applied voltage and the current in the device. The upper and lower equilibrium positions correspond to ON and OFF states of a switch, respectively. Also the existence of the tunneling current and feedback resistor make the "lower" equilibrium states very robust, which is key to some applications of interest. The representative characteristic curve of the device is shown in Fig. 17: (a) shows the relation between the gap $\Delta$ and the applied voltage U; (b) shows the relation between the current $i$ in the circuit and the applied voltage U.

The current jump behavior at the pull-in has been observed for a nanotube cantilever freestanding above an electrode actuated by electrostatic forces [32], as shown in Fig. 18, and the I–V behavior after the pull-in has been demonstrated based on the good agreement between experimental measurements and theoretical prediction. The parameters used in the theoretical prediction includes the length of the nanotube $L = 3.8$ μm; the diameter of the nanotube $R_{ext} = 20$ nm; and the initial gap between the nanotube cantilever and the electrode $H = 200$ nm, $R = 0.98$ GΩ, and contact resistance $R_0 = 50$ Ω [11].

The potential applications of the device include ultrasonic wave detection for monitoring the health of materials and structures, gap sensing, NEMS switches, memory elements, and logic devices.



Figure 17. Representative characteristic of pull-in and pull-out processes for the feedback-controlled nanocantilever device. (a) Relationship between the gap $\Delta$ and the applied voltage $U$. (b) Relationship between the current $i$ in the circuit and the applied voltage U. Reprinted with permission from [11]. C.-H. Ke and H. D. Espinosa, *Appl. Phys. Lett.* 85, 681 (2004). © 2004, American Institute of Physics.

**Figure 18.** Comparison between theoretical prediction and I-V measurement of an electrostatically actuated free-standing nanotube cantilever with an electronic circuit incorporating a resistor.

In comparison to nanorelays [10, 41], the device reported in Ref. [11] is a two-terminal device, providing more flexibility in terms of device realization and control than the nano-relay. In comparison to the NRAM described in Ref. [1], the feedback-controlled device employs an electrical circuit incorporated with a resistor to adjust the electrostatic field to achieve the second stable equilibrium position. This feature reduces the constraints in fabricating devices with nanometer gap control between the freestanding CNTs or NWs and the substrate, providing more reliability and tolerance to variability in fabrication parameters. However, the drawback of the device in memory applications is that the memory becomes volatile. The working principle and the potential applications for these two devices are somewhat complementary.

**2.4.1.6. Tunable Oscillators** The fabrication and testing of a tunable carbon nanotube oscillator was reported by Sazonova et al. [12]. It consists of a doubly clamped nanotube, as shown in Fig. 19. They demonstrated that the resonance frequency of the oscillators can be widely tuned and that the devices can be used to transduce very small forces.



**Figure 19.** SEM image of a suspended device (top) and a schematic of device geometry (bottom). Scale bar, 300 nm. The sides of the trench, typically 1.2–1.5 $\mu$m wide and 500 nm deep, are marked with dashed lines. A suspended nanotube can be seen bridging the trench. Reprinted with permission from [12], V. Sazonova et al., *Nature* 431, 284 (2004). © 2004, Nature Publishing Group.

Single- or few-walled nanotubes with diameters in the range of 1–4 nm, grown by CVD were suspended over a trench (typically 1.2–1.5 μm wide, 500 nm deep) between two metal (Au/Cr) electrodes. A small section of the tube resided on the oxide on both sides of the trench; the adhesion of the nanotube to the oxide provided clamping at the end points. The nanotube motion was induced and detected using the electrostatic interaction with the gate electrode underneath the tube. In this device, the gate voltage has both a static (DC) component and a small time-varying (AC) component. The DC voltage at the gate produces a static force on the nanotube that can be used to control its tension. The AC voltage produces a periodic electric force, which sets the nanotube into motion. As the driving frequency approaches the resonance frequency of the tube, the displacement becomes large.

The transistor properties of semiconducting [66] and small-bandgap semiconducting [67, 68] carbon nanotubes were employed to detect the vibrational motion. Figure 20(a) shows the measured current through the nanotube as a function of driving frequency at room temperature. A distinctive feature in the current on top of a slowly changing background can be seen. This feature is due to the resonant motion of the nanotube, which modulates the capacitance, while the background is due to the modulating gate voltage.

The DC voltage on the gate can be used to tune the tension in the nanotube and therefore the oscillation frequency. Figure 20(b) and 20(c) show the measured response as a function of the driving frequency and the static gate voltage. The resonant frequency shifts upward as the magnitude of the DC gate voltage is increased. Several distinct resonances are observed, corresponding to different vibrational modes of the nanotube. Figure 20(d) shows the theoretical predictions for the dependence of the vibration frequency on gate voltage for a representative device. The predictions are based on finite element analysis, with the nanotube modeled as a long beam suspended over a trench. With the increase of the gap voltage, the deflection of the nanotube becomes larger and the stretching dominates the bending. Therefore, the stiffness of the nanotube beam increases and so does the resonance frequency. The theoretical predictions (Fig. 20(d)) show good qualitative agreement with experiments (Fig. 20(b) and 20(c)). The device showed a high-force sensitivity (below 5aN), which made it a small-force transducer.



Figure 20. Measurements of resonant response. (a) Detected current as a function of driving frequency. (b)–(c) Detected current as a function of gate voltage $V_g$ and frequency for devices 1 and 2. (d) Theoretical predictions for the dependence of vibration frequency on gate voltage for a representative device. Reprinted with permission from [12], V. Sazonova et al., *Nature* 431, 284 (2004). © 2004, Nature Publishing Group.

## 2.4.2. Nanowire-Based Nanoelectromechanical Systems Devices

Nanowires, like carbon nanotubes, are high-aspect-ratio, one-dimensional nanostructures. The materials of nanowires include silicon [52, 69–72], gold [73, 74], silver [75–77], platinum [24], germanium [71, 78–81], zinc oxide [82, 83], and so on. Besides their size, the advantages offered by nanowires when employed in NEMS are their electronic properties, which can be controlled in a predictable manner during synthesis. This has not been achieved yet for carbon nanotubes. In contrast to carbon nanotubes, nanowires do not exhibit the same degree of flexibility, which may be a factor concerning device fabrication and reliability. In the following section, two nanowire-based NEMS device are briefly reviewed.

### 2.4.2.1. Resonators

Figure 21 shows a suspended platinum nanowire resonator (a), reported by Husain et al. 2003 [24], and the circuit used for magnetomotive drive and detection of its motion (b).

Synthetized platinum nanowires were deposited on a Si substrate capped by a 300-nm-thick layer of thermally grown silicon dioxide and prepatterned with Au alignment marks. The location of the deposited wires was mapped, by means of optical microscopy, using their strong light-scattering properties [76, 84]. Metallic leads (5 nm Cr, 50 nm Au) to individual wires were subsequently patterned by electron-beam lithography, evaporation, and lift-off. Finally, the $SiO_2$ was removed by wet etching (HF) to form suspended nanowire structures. The suspended Pt nanowire shown in Fig. 21 has a diameter of 43 nm and a length of 1.3 $\mu$m. A magnetomotive detection scheme (see Fig. 21, right), in which an AC current drives a beam in a transverse magnetic field, was used to drive and read out the resonators. Figure 22 shows the measured motion-induced impedance of the nanowire device, $|Z_m(f)|$, versus frequency. The measured quality factor Q was approximately 8500 and decreased slightly with the increase in magnetic field. It was noted that the characteristic curve shown in Fig. 22 corresponds to a linear response of the beam. Badzey et al. [52] reported a doubly clamped nanomechanical Si beam working in the nonlinear response region. The nonlinear response of the beam displays notable hysteresis and bistability in the amplitude-frequency space when the frequency sweeps upward and downward. This particular behavior shows that the device can be used as mechanical memory elements.

### 2.4.2.2. Nanoelectromechanical Programmable Ready-Only Memory

A nanowire-based nanoelectromechanical programmable read-only memory (NEMPROM), reported by Ziegler et al. 2004 [80], is shown in Fig. 23(a). The germanium nanowire was synthesized directly onto a macroscopic gold wire (diameter = 0.25 mm). The combination of transmission electron microscope (TEM) and STM was used to control and visualize the nanowire under investigation. Figure 23(b)–23(g) illustrates how the device can work as NEMPROM. In equilibrium, the attractive van der Waals (vdW) force and electrostatic interactions between the nanowire and the gold electrode are countered by the elastic force from the deflection of the nanowire. Figure 23(b) shows the position of the nanowire with relative low applied voltage. With the increase in voltage, the nanowire moves closer to the electrode (Fig. 23(c)). When the applied voltage exceeds a certain value, a jump-to-contact happens, i.e., the nanowire makes physical contact with the electrode (Fig. 23(d)). The nanowire remains in contact with the electrode even when the electrostatic field is removed because the vdW force is larger than the elastic force (Fig. 23(e)). This is the ON state of the NEMPROM. The NEMPROM



(a)                                        (b)

**Figure 21.** (a) SEM image of the suspended nanowire device, 1.3 $\mu$m long and 43 nm in diameter. (b) Measurement circuit used for magnetomotive drive and detection. Reprinted with permission from [24], A. Husain et al., *Appl. Phys. Lett.* 83, 1240 (2003). © 2003, American Institute of Physics.

Figure 22. Measured mechanical impedance of a Pt nanowire device as a function of frequency, at a series of magnetic fields from 1 to 8 T. The left inset shows the characteristic $B^2$ dependence typical of magnetomotive detection. The right inset shows the quality factor $Q$ as a function of magnetic field. Reprinted with permission from [24]. A. Husain et al., *Appl. Phys. Lett.* 83, 1240 (2003). © 2003, American Institute of Physics.

device can be switched OFF by mechanical motion or by heating the device above the stability limit to overcome the vdW attractive forces. Figure 23(f) and 23(g) show the separation of the nanowire and the electrode after imposing a slight mechanical motion, resulting in a jump-off-contact event. This is the OFF state of NEMPROM. The working principle of NEMPROM is similar to that of NRAM [1] since both of them employ van der Waals energy to achieve the bistability behavior, although the usage of germanium may provide better control of size and electrical behaviors of the device than that of carbon nanotube.

## 2.5. Future Challenges

NEMS offer unprecedented and intriguing properties in the fields of sensing and electronic computing. Although significant advancement has been achieved, there are many challenges that will need to be overcome before NEMS can replace and revolutionize current



Figure 23. (a) TEM image of a Ge nanowire device. (b)-(d) TEM sequence showing the jump-to-contact of a Ge nanowire as the voltage is increased. (e) TEM image demonstrating the stability of device after removal of the electrostatic potential. (f) and (g) TEM sequence demonstrating the resetting behavior of the device. Reprinted with permission from [80]. K. J. Ziegler et al., *Appl. Phys. Lett.* 84, 4074 (2004). © 2004, American Institute of Physics.

technologies. Among the issues that need further research and development are

1. *Extremely high integration level*: For applications such as RAM and data storage, the density of the active components is definitely a key parameter. Direct growth and directed self-assembly are the two most promising methods to make NEMS devices with levels of integration orders of magnitudes higher than that of current microelectronics.

   A process for nanofabrication of the NEMS device developed by Ke and Espinosa [11], based on the directed self-assembly, is schematically shown in Fig 24.

   a. A 1-$\mu$m-thick $Si_3N_4$ dielectric film is deposited on a Si wafer by low pressure chemical vapor deposition (LPCVD). Then, a 50-nm-thick gold film (with 5 nm Cr film as adhesion layer) is deposited by e-beam evaporation and patterned by lithography to form the bottom electrodes. A 1-$\mu$m-thick $SiO_2$ layer is deposited by PECVD (plasma enhanced chemical vapor deposition).

   b. The fountain-pen nanolithography technique [85] is then employed to functionalize specific areas, with widths down to 40 nm, either with polar chemical groups (such as the amino groups ($-NH_2/-NH_3^+$) of cysteamine) or carboxyl ($-COOH/-COO^-$) or with nonpolar groups (such as methyl ($-CH_3$) from molecules like 1-octadecanethiol, ODT).

   c. The substrate is dipped into a solution containing prefunctionalized (with polar chemical groups) CNTs or NWs to adhere and self-assembly to the functionalized sites.

   d. The chip is patterned with e-beam lithography and e-beam evaporation of 100 nm gold film (lift-off, with 5 nm Cr film as adhesion layer) to form the top electrodes.

   e. Removal of the $SiO_2$ layer using wet etching (HF) to free one end of the CNT cantilever completes the process.

   The final product, a two-dimensional array of NEMS devices, with multiplexing capabilities is schematically shown in Fig. 25. The top and bottom electrodes are interconnected to the pads, correspondingly. By applying voltage between the corresponding pads, the individual NEMS devices can be independently actuated.

2. *Better understanding the quality factor*: One of the keys to realize the potential applications of NEMS is to achieve ultra-high-quality factors. However, it has been consistently observed that the quality factor of resonators decreases significantly with size scaling [5]. Defects in the bulk materials and interfaces, fabrication-induced surface damages, adsorbates on the surface, and thermoelastic damping are a few commonly listed factors that can dampen the motion of resonators. Unfortunately, the dominant energy dissipation mechanism in nanoscale mechanical resonators is still unclear.

3. *Reproducible and routine nanomanufacturing*: Fabrication reproducibility is key in applications such as mass sensors. Because the NEMS can respond to mass at the level of single atom or molecules, it places an extremely stringent requirement on the cleanness and precision of nanofabrication techniques. Likewise, devices that rely on



Figure 24. Schematic of the fabrication steps involving nano fountain probe (NFP) functionalization.

Figure 25. Schematics of two-dimensional array of the NEMS device with multiplexing.

van der Waals energy require dimensional control (e.g., gap dimension) in the order of a few nanometers.

4. *Quantum limit for mechanical devices*: The ultimate limit for NEMS is its operation at, or even beyond, the quantum limit [5]. In the quantum regime the individual mechanical quanta are of the same order of magnitude, or greater than the thermal energy. Quantum theory should be used to understand and optimize force and displacement measurements. Recently, position resolution with a factor of 4.3 above the quantum limit has been achieved for a single-electron transistor with high-quality factor at millikelvin temperature [86]. The pursuit of NEMS devices operating at the quantum limit will potentially open new fields in science at the molecular level.

## 3. MODELING OF NANOELECTROMECHANICAL SYSTEMS DEVICES

The design of NEMS depends on a thorough understanding of the mechanics of the devices themselves and the interactions between the devices and the external forces/fields. With the critical dimension shrinking from micron to nanometer scale, new physics emerges so that the theory typically applied to MEMS does not immediately translated to NEMS. For example, van der Waals forces from atomic interactions play an important role in NEMS, while they can be generally neglected in MEMS. The behavior of materials at nanometer scale begins to be atomistic rather than continuous, giving rise to anomalous and often nonlinear effects, for example,

- The roles of surfaces and defects become more dominant.
- The devices become more compliant than continuum models predict.
- Molecular interactions and quantum effects become key issues to the point that thermal fluctuation could make a major difference in the operation of NEMS.

For instance, the nanoresonators reported by the research groups of Roukes and Craighead are operated in the gigahertz range and usually have sizes within $200 \times 20 \times 10$ nm$^3$ [87].

Devices of this size and smaller are so minuscule that material defects and surface effects have a large impact on their performance.

In principle, atomic-scale simulations should well predict the behavior of NEMS devices. However, atomic simulations of the entire NEMS involve prohibitively expensive computational resources or exceed the current computational power. Alternatively, multiscale modeling, which simulates the key region of a device with an atomistic model and other regions with a continuum model, can well serve the purpose under the circumstance of limited computational resources. Besides, it has been demonstrated that the behavior of some nanostructures, like carbon nanotubes, can be approximated by continuum mechanics models, based on the same potentials governing molecular dynamics (MD) simulation [88], if the surface nonideality of the nanostructures is neglected. Thus, continuum mechanics models are still adequate to the design of NEMS, in particular, in the initial stages.

## 3.1. Multiscale Modeling

Multistage modeling can be pursued sequentially or concurrently. In the sequential method, information from each model at a given scale is passed to the next modeling level. In this fashion, "informed" or physically motivated models are developed at larger scales. In concurrent multiscale modeling, the system is split into primarily two domains: the atomistic domain and the continuum domain. In the atomistic domain, MD and quantum mechanics (QM) are typically employed, while in the continuum domain, the finite element method is often used. MD deals with the interaction of many thousands of atoms or more according to an interaction law. The "constitutive" behavior of each atom is governed by QM. QM involves the electronic structure, which in turn determines the interatomic force law—the "constitutive" behavior of each atom. However, in practice, the interatomic force laws have been determined empirically based on both QM and experiments. To model the response of NEMS devices, MD and continuum mechanics are generally adequate; hence, here we restrict our discussion to the basic ideas behind these two models. In some cases, QM modeling is required so the reader should consult the information on QM.

As an example of sequential multiscale modeling, we discuss how the mechanical properties of bulk tantalum were calculated using a multiscale modeling strategy. Moriarty et al. [89] started with fundamental atomic properties and used rigorous quantum-mechanical principles calculations to develop accurate interatomic force laws that were then applied to atomistic simulations involving many thousands of atoms. From these simulations, they derived the properties of individual dislocations in a perfect crystal and then, with a new microscale simulation technique, namely, dislocation dynamics, examined the behavior of large collections of interacting dislocations at the microscale in a grain-sized crystal. They modeled the grain interactions in detail with finite-element simulation, and from those simulations, they finally constructed appropriate models of properties such as yield strength in a macroscopic volume of tantalum. At each length scale, the models were experimentally tested and validated with available data. The concept of information passing between models, from quantum modeling to atomic to continuum scale, is quite general and can be applied in a variety of problems including NEMS.

### 3.1.1. Implementation of Concurrent Multiscale Modeling

MD computes the classical trajectories of atoms by integrating Newton's law, $F = ma$, for the system. In the MD domain, the interaction force follows an empirical potential. Consider a set of $n_M$ molecules with the initial coordinates $X_I$, $I = 1$ to $n_M$. Let the displacements be denoted by $d_I(t)$. The potential energy is then given by $W_M(d)$. For a given potential function $W_M(d)$, an equilibrium state is given by

$$dW_M(d) = 0 \tag{1}$$

From the continuum viewpoint, the governing equations arise from conservation of mass, momentum and energy. Using a so-called total Lagrangian description [90], the linear momentum equations are

$$\frac{\partial P_{ji}}{\partial X_j} + \rho_0 b_i = \rho_0 \ddot{u}_i \tag{2}$$

where $\rho_0$ is the initial density, $P$ is the nominal stress tensor, $b$ is the body force per unit mass, $u$ the displacement, and the superposed dots denote material time derivatives.

There are two approaches to building a multiscale model: domain decomposition with overlapping domains, often referred to as the "handshake" model [87], and the edge-to-edge decomposition method [91].

The main features of the overlapping domain decomposition include: (a) A Lagrange multiplier method and an augmented Lagrange method to impose consider the constraints on the motion; (b) The Lagrange multiplier field in the overlapping domain vanishes at the edges of the continuum domain so that the interaction forces between the continuum and molecular mechanics model are smooth if an atom exits the overlapping domain.

In edge-to-edge decomposition coupling method, there are three types of particles. Besides the nodes of the continuum domain and the atoms of the molecular domain, virtual atoms are defined to model the bond angle-bending for bonds between the continuum and the molecular domains. The virtual atoms are connected with the molecular domain by virtual bonds. An example showing the domains in difference of all these methods, when applied to the modeling of a graphite sheet, is shown in Fig. 26.

## 3.1.2. Examples of Concurrent Multiscale Modeling

Because of the computational power and efficiency, multiscale modeling has been used widely in the modeling and simulation of nanostructures and NEMS. Here three examples are highlighted: carbon nanotube fracture [91], carbon nanotube–based switch performance [88], and nanogears kinematics [87].

In the model developed by Belytschko et al. [91] for studying carbon nanotubes fracture, two shells of a nanotube interacting by van der Waals forces were considered. The molecular model was used only in a small subdomain surrounding a defect, while the finite element model was employed outside of the molecular model (Fig. 27). A modified Morse potential was used. The load was only applied to the outside shell. For the entire domain to be modeled by molecular mechanics, 46,200 atoms were required, which is very expensive computationally. The numerical results were compared with reported experimental values for the failure stress. The model with a certain number of defects agreed much better with the experimental measurements than did a perfect nanotube model.

As aforementioned, carbon nanotube–based electrostatic switches have the potential to operate in the gigahertz range and achieve much higher integration levels than currently possible. As such, modeling attempts to gain insight into the performance of the device pursued. Aluru's group developed various numerical models, including continuum models, and continuum/MD coupled models to analyze device behavior [92]. This work shows that continuum modeling, considering nonlinear beam theory, is in good agreement with molecular



**Figure 26.** Comparison of bent graphite sheet by means of (a) molecular mechanics, (b) overlapping coupling, (c) edge-to-edge coupling. Reprinted with permission from [91], T. Belytschko et al., *Int. J. Multiscale Comp. Engr.* 1, 115 (2003). © 2003, Begell House, Inc.

Figure 27. Carbon nanotube model for fracture study by overlapping coupling method. Reprinted with permission from [91]. T. Belytschko et al., *Int. J. Multiscale Comp. Engr.* 1, 115 (2003). © 2003, Begell House, Inc.

mechanics modeling. MD simulation showed that, with an increase of the gap between the nanotube and the ground, the nanotube locally buckles as it approaches the ground. The local buckling phenomenon was not captured by the continuum beam theory employed in the analysis. Hence a combined continuum/MD technique was used to overcome this limitation. Figure 28 compares the deformed shapes obtained from the combined continuum/MD model and the fully molecular mechanics model. The buckling that occurred at the two ends of the nanotube was captured by the MD subdomain. In the center of the nanotube, the nonlinear beam theory was able to predict well and greatly reduced the computational cost. The combined continuum/MD approach and the fully MD approach provided good agreement in terms of the static pull-in voltage.



Figure 28. Deformation plot of a fixed-fixed carbon nanotube-based switch. The top figure is the result from full molecular dynamics (MD) simulation and the lower is the multiscale result in which the central region is modeled by one-dimensional nonlinear beam theory. Reprinted with permission from [92]. M. Desquenes et al., *J. Eng. Mater. Tech.* 126, 230 (2004). © 2004, American Society of Mechanical Engineers.

Microgears have been one of the most successful MEMS devices so far. Such devices are presently made at the 100 $\mu$m scale and rotate at speeds of 150,000 rpm. Next generation of devices (nanogears), based on nanofabrication, are expected to be below the 1 $\mu$m level. The effects of wear, lubrication, and friction are expected to have significant consequences on the performance of the nanogears, where areas of contact are an important part of the systems. However, the process of nanogear teeth grinding against each other cannot be simulated accurately with FE because of the bond breaking and formation at the point of contact can only be treated empirically in FE. Alternatively, multiscale modeling provides a good tool to predict the mechanics-related issues for these devices. Figure 29 shows the multiscale decomposition for the modeling of nanogears. An inner region, including the shaft, is discretized by finite elements. The handshaking between the FE and MD region is accomplished by a self-consistent overlap region. In regions at the gear-gear contact point in the nonlubricated case, a tightbinding (TB) description is used as part of a QM simulation [87].

## 3.2. Continuum Mechanics Modeling

Many NEMS devices can be modeled either as biased cantilever beams or fixed-fixed beams freestanding over a ground substrate, as shown in Fig. 30. The beams can be carbon nanotubes, nanowires, or small nanofabricated parts. The electromechanical characterization of NEMS involves the calculation of the elastic energy ($E_{elas}$), from the deformation of active components, the electrostatic energy ($E_{elec}$), and van der Waals energy ($E_{vdw}$) from atomic interactions. In the following section, we summarize the continuum theory for each of these



Figure 29. Illustration of dynamic simulation zone and domain decomposition for coupling of length scales: from continuum (FE) to atomistic (MD) to electronic structure (TB). Reprinted with permission from [87], R. E. Rudd and J. Q. Broughton, R. E. Rudd and J. Q. Broughton, J. Model. Simul. Microsys. 1, 29 (1999). © 1999, Applied Computational Research Society.

**Figure 30.** Schematic of NEMS devices. (A) Cantilever beam configuration. (B) Doubly clamped beam configuration. Reprinted with permission from [103], C.-H. Ke and H. D. Espinosa, *J. Appl. Mech.* 72, 721 (2005). © 2005, American Society of Mechanical Engineers.

energy domains and the governing equations of equilibrium for both small deformation and finite deformation. We follow the work reported in Refs. [32, 88, 92, 103, 104, 106].

## 3.2.1. Continuum Theory

### 3.2.1.1. Van der Waals Interactions
The van der Waals (vdW) energy originates from the interaction between atoms. The Lennard-Jones potential is a suitable model to describe van der Waals interaction [93]. In the Lennard-Jones potential, there are two terms: one is repulsive and the other is attractive. The Lennard-Jones potential between two atoms $i$ and $j$ is given by

$$\phi_{ij} = \frac{C_{12}}{r_{ij}^{12}} - \frac{C_6}{r_{ij}^6} \tag{3}$$

where $r_{ij}$ is the distance between atoms $i$ and $j$ and $C_6$ and $C_{12}$ are attractive and repulsive constants, respectively. For the carbon-carbon interaction, $C_6 = 15.2$ eVÅ$^6$ and $C_{12} = 24.1$ keVÅ$^{12}$ and the equilibrium spacing $r_0 = 3.414$ Å [94]. From Eq. (3), we can see that the repulsive components of the potential decay extremely fast and play an important role only when the distance is close to or smaller than $r_0$. The total van der Waals energy can be computed by a pair-wise summation over all the atoms. The computational cost (number of operations) is proportional to the square of the number $n$ of atoms in the system. For a NEMS device with millions of atoms, this technique is prohibitively expensive. Instead, a continuum model was established to compute the van der Waals energy by the double-volume integral of the Lennard-Jones potential [95], that is,

$$E_{vdW} = \int_{v_1} \int_{v_2} n_1 n_2 \left( \frac{C_{12}}{r^{12}(v_1, v_2)} - \frac{C_6}{r^6(v_1, v_2)} \right) dv_1 dv_1 \tag{4}$$

where, $v_1$ and $v_2$ represent the two domains of integration, and $n_1$ and $n_1$ are the densities of atoms for the domains $v_1$ and $v_2$, respectively. The distance between any point on $v_1$ and $v_2$ is $r(v_1, v_2)$.

Let us consider SWNT freestanding above a ground plane consisting of layers of graphite sheets, with interlayer distance $d = 3.35$ Å, as illustrated in Fig. 31(A). The energy per unit length of the nanotube is given by

$$\frac{E_{vdW}}{L} = 2\pi\sigma^2 R \sum_{n=1}^{N} \int_{-\pi}^{\pi} \left( \frac{C_{12}}{10[(n-1)d + r_{init} + R + R\sin\theta]^{10}} - \frac{C_6}{4[(n-1)d + r_{init} + R + R\sin\theta]^4} \right) d\theta \tag{5}$$

where $L$ is the length of the nanotube, $R$ is the radius of the nanotube, $r_{init}$ is the distance between the bottom of the nanotube and the top graphene sheet, $N$ is the number of graphene sheets and $\sigma \cong 38$ nm$^{-2}$ is the graphene surface density. When $r_{init}$ is much larger than the equilibrium spacing $r_0$, the repulsive component can be ignored and Eq. (5) can be simplified [88] as

$$\frac{E_{vdW}}{L} = C_6 \sigma^2 \pi^2 R \sum_{r=r_{init}}^{(N-1)d + r_{init}} \frac{(R+r)[3R^2 + 2(r+R)^2]}{2[(r+R)^2 - R^2]^{7/2}} \tag{6}$$

**Figure 31.** Van der Waals integration of a SWNT (A) and MWNT (B) over a graphite ground plane. Reprinted with permission from [88], M. Desquenes et al., *Nanotechnology* 13, 120 (2002). © 2002, Institute of Physics.

The accuracy of Eq. (6) in approximating the continuum van der Waals energy of a SWNT placed over a graphite plane is verified by the comparison with the direct pair-wise summation of the Lennard-Jones potential given by Eq. (3) for a (16, 0) tube, which is shown in Fig. 32 [88].

For a MWNT, as illustrated in Fig. 31(B), the energy per unit length can be obtained by summing up the interaction between all separate shells and layers:

$$\frac{E_{vdW}}{L} = \sum_{R=R_{int}}^{R_{ext}} \sum_{r=r_{int}}^{(N-1)d+r_{int}} \frac{C_6 \sigma^{-2} \pi^2 R(R + r)[3R^2 + 2(r + R)^2]}{2[(r + R)^2 - R^2]^{7/2}}$$

(7)

where $R_{int}$ and $R_{ext}$ are the inner and outer radii of the nanotube, respectively.

The van der Waals force per unit length can be obtained as

$$q_{vdW} = \frac{d\left(\frac{E_{vdW}}{L}\right)}{dr}$$

(8)



**Figure 32.** Comparison of the continuum van der Waals (vdW) energy given by Eq. (6) with the discrete pair-wise summation given by Eq. (3). Reprinted with permission from [88], M. Desquenes et al., *Nanotechnology* 13, 120 (2002). © 2002, Institute of Physics.

Thus, inserting Eq. (7) into Eq. (8) and taking the derivative with respect to $r$, one obtains [88]

$$q_{dw} = \sum_{R=R_{int}}^{R_{ext}} \sum_{r=r_{int}}^{R_{ext}+r_{int}} -\frac{C_0 \sigma^2 \pi^2 R \sqrt{r(r+2R)}(8r^4 + 32r^3R + 72r^2R^2 + 80rR^3 + 35R^4)}{2[2r^3(r+2R)^5]^5} \quad (9)$$

### 3.2.1.2. Electrostatic Force

When a biased conductive nanotube is placed above a conductive substrate, there are induced electrostatic charges both on the tube and on the substrate. The electrostatic force acting on the tube can be calculated using a capacitance model [96].

Let us look at the electrostatic force for a conductive nanotube with finite length and round cross section above an infinite ground plane. Although nanotubes have hollow structures, carbon nanotubes with capped ends are more electrochemically stable than those with open ends [97]. Thus, nanotubes with finite length, as well as nanowires, can be geometrically approximated by conductive nanocylinders. For small-scale nanocylinders, the density of states on the surface is finite. The screening length, the distance that the "surface charge" actually penetrates into the cylinder interior, is found to be a nanometer-scale quantity [98]. For nanocylinders with transverse dimension (i.e., diameter approaching the screening length), such as SWNT, the finite size and density of states (quantum effects) have to be considered thoroughly when calculating the surface/volume charge distribution [99, 100]. For nanocylinders with transverse dimension much larger than the screening length, such as MWNT or nanowires with large outer diameter (e.g., 20 nm), this quantum effect can be considered negligible. Thus, the charge distribution can be approximated by the charge distribution on a metallic, perfectly conductive cylinder with the same geometry to which classical electrostatic analysis can be applied.

For infinitely long metallic cylinders, the capacitance per unit length [96] is given by

$$C_d(r) = \frac{\pi\varepsilon}{a\cosh(1 + \frac{r}{R})} \quad (10)$$

where $r$ is the distance between the lower fiber of the nanocylinder and the substrate, $R$ is the radius of the nanocylinder, and $\varepsilon$ is the permittivity of the medium. For vacuum, $\varepsilon_0 = 8.854 \times 10^{-12}$ $C^2N^{-1}m^{-2}$. Equation (10) can be applied for infinitely long MWNTs with large diameters ($R = R_{ext}$).

For the charge distribution on infinite long SWNT, Bulashevich and Rotkin [100], proposed a quantum correction, rendering the capacitance per unit length as

$$C = \frac{C_d}{1 + \frac{C_d}{C_Q}} \approx C_d\left(1 - \frac{C_d}{C_Q}\right) \quad (11)$$

where $C_Q = e^2 v_m$, and $v_M$ is the constant density of the states near the electroneutral level measured from the Fermi level.

For nanocylinders with finite length, there are two types of boundary surfaces—the cylindrical side surface and the planar end surface. Essentially classical distribution of charge density with a significant charge concentration at the cylinder end has been observed [99, 101, 102]. Here we discuss a model to calculate the electrostatic charge distribution on metallic cylindrical cantilevers based on a boundary element method (BEM), considering both the concentrated charge at the free end and the finite rotation due to the deflection of the cantilever [103].

Figure 33 shows the charge distribution along the length $L$ of a freestanding nanotube, subjected to a bias voltage of 1 V. The contour plot shows the charge density (side view), while the curve shows the charge per unit length along the nanotube. The calculation was performed using the CFD-ACE+ software (a commercial code from CFD Research Corporation based on finite and BEMs). There is significant charge concentration on the free ends and uniform charge distribution in the central of the cantilever, which is found to follow Eq. (10). The charge distribution along a deflected cantilever nanotube is shown in Fig. 34.

**Figure 33.** Charge distribution for a biased nanotube. The device parameters are $R_{ext}$ = 9 nm, $H$ = 100 nm, and $L$ = 1 $\mu$m. Reprinted with permission from [104]. C.-H. Ke et al., *J. Appl. Mech.* 72, 721 (2005). © 2005, American Society of Mechanical Engineers.

The parameters are $R_{ext}$ = 20 nm, $H$ = 500 nm, $L$ = 3 $\mu$m, and the gap between the free end and the substrate $r(L)$ = 236 nm. From Fig. 34, it is seen that, besides the concentrated charge on the free end, the clamped end imposes a significant effect to the charge distribution in the region close to it [98]. However, this effect can be considered negligible because its contribution to the deflection of the nanotube is quite limited. The charge distribution in regions other than the two ends closely follows Eq. (10). A formula for the charge distribution, including end charge effects and the deflection of the cantilever, is derived from a parametric analysis [103], as follows:

$$C(r(x)) = C_d(r(x))\{1 + 0.85[(H + R)^2 R]^{1/3}\delta(x - x_{tip})\} = C_d(r(x))\{1 + f_c\} \qquad (12)$$

where the first term in the bracket accounts for the uniform charge along the side surface of the tube and the second term, $f_c$, accounts for the concentrated charge at the end of the tube (for doubly clamped tube, $f_c$ = 0). $H$ is the distance between the cantilever and the substrate when the cantilever is in horizontal position, $R$ is the radius of the tube (for MWNT $R = R_{ext}$), $x = x_{tip} = L$ for small deflection (when considering the finite kinematics.



**Figure 34.** (Top) Two-dimensional side view of the charge distribution in a deflected nanotube cantilever. (Bottom) Charge distribution per unit length along a deflected nanotube cantilever. The solid line is plotted from Eq. (10); the dotted line is the simulation result performed with CFD-ACE+. Reprinted with permission from [103], C.-H. Ke and H. D. Espinosa, *J. Appl. Mech.* 72, 721 (2005). © 2005, American Society of Mechanical Engineers.

i.e., large displacement, $x = x_{tip} \neq L$), $\delta(x)$ is the Dirac function, and $r(x) = H - w(x)$, with $w$ being the tube deflection.

Thus, the electrostatic force per unit length of the nanotube is given by differentiation of the energy [104], as follows:

$$q_{elec} = \frac{1}{2}V^2\frac{dC}{dr} = \frac{1}{2}V^2\left(\frac{dC_d}{dr}\right)\{1 + f_i\} = \frac{-\pi\varepsilon_0 V^2}{\sqrt{r(r+2R)}a\cosh^2(1+\frac{r}{R})}(1+f_i) \qquad (13)$$

where $V$ is the bias voltage.

### 3.2.1.3. Elasticity

Continuum-beam theory has been widely used to model the mechanics of nanotubes [25, 32, 88, 92, 104–106]. The applicability and accuracy of the continuum theory have been evaluated by comparison with MD simulations [88]. Figure 35 shows the comparison of the deflection of a 20-nm-long, doubly clamped DWNT with a diameter of 1.96 nm, calculated by MD simulation and by the beam equation, respectively. The solid black curve—the deflection predicated by the beam equation—follows closely the shape predicted by MD calculations.

Because nanotubes have high flexibility with strain at tensile failure of the order of 30% [107], nonlinear effects such as finite kinematics, accounting for large displacement need to be considered in the modeling. This is particularly important for doubly clamped nanotube beams because the stretching from the finite kinematics stiffens the beam, resulting in a significant increase of the pull-in voltage, a key parameter in NEMS devices.

### 3.2.1.4. Governing Equations

The electromechanical characteristic of nanotube cantilevers or doubly clamped nanotube beams can be determined by coupling the van der Waals, electrostatic, and elastic forces. The governing equation under the small deformation assumption (considering only bending) [88] is given by

$$EI\frac{d^4 r}{dx^4} = q_{elec} + q_{vdW} \qquad (14)$$

where $r$ is the gap between the nanotube and the ground plane, $x$ is the position along the tube, $E$ is the Young's modulus (for carbon nanotube $E = 1 - 1.2$ TPa), $I$ is the moment of inertia (for nanotubes, $I = \frac{\pi}{4}(R_{ext}^4 - R_{int}^4)$, $R_{ext}$ and $R_{int}$ are the outer and inner radii of the nanotubes, respectively), and $q_{elec}$ and $q_{vdW}$ are given by Eqs. (13) and (9), respectively.

For cantilevers exhibiting large displacements, as shown in Fig. 36, the curvature of the deflection should be considered and the governing equation [104] changes into

$$EI\frac{d^2}{dx^2}\left(\frac{\frac{d^2 r}{dx^2}}{\left(1 + \left(\frac{dr}{dx}\right)^2\right)^{\frac{3}{2}}}\right) = (q_{vdW} + q_{elec})\sqrt{1 + \left(\frac{dr}{dx}\right)^2} \qquad (15)$$

For doubly clamped structures exhibiting finite kinematics, as shown in Fig. 37, stretching becomes significant as a consequence of the ropelike behavior of a doubly clamped nanotube. The corresponding governing equation [92, 104, 106] is expressed as

$$EI\frac{d^4 w}{dx^4} - \frac{EA}{2L}\int_0^L\left(\frac{dw}{dx}\right)^2 dx\frac{d^2 w}{dx^2} = q_{elec} + q_{vdW} \qquad (16)$$

where the term $\frac{EA}{2L}\int_0^L(\frac{dw}{dx})^2 dx$ is the tension along the axis of the tube due to stretching.



**Figure 35.** Comparison between MD and beam theory of the deflection of a 20-nm-long fixed-fixed DWNT (diameter 1.96 nm). The solid black curve is the deflection predicated by beam theory. Reprinted with permission from [88], M. Desquenes et al., *Nanotechnology* 13, 120 (2002). © 2002, Institute of Physics.

**Figure 36.** Schematic of finite kinematics configuration of a cantilever nanotube device subjected to electrostatic and van der Waals forces. Reprinted with permission from [104], C.-H. Ke et al., *J. Appl. Mech.* 72, 726 (2005). © 2005, American Society of Mechanical Engineers.

The aforementioned governing equations can be numerically solved by either direct integration or finite difference method. The effect of various factors, such as concentrated charge, finite kinematics, and stretching, on the prediction of pull-in voltages of devices can then be identified.

In the following, the effects of concentrated charge and finite kinematics on the prediction of the pull-in voltage for a cantilevered nanotube with $R_{ext} = 10$ nm, $R_{int} = 0$, $E = 1$ TPa, $H = 100$ nm, $l = 500$ nm are considered [104]. The displacement of the tip as a function of the applied voltage is shown in Fig. 38. As expected, the role of the finite kinematics becomes negligible. The pull-in voltages, corresponding to the vertical lines, differ by less than 1%. Both numerical solutions reported in Fig. 38 consider the charge concentration at the tip of the cantilevered nanotube. Figure 39 shows the error in the pull-in voltage in case the charge concentration is ignored. It is inferred that the error from neglecting charge concentration can be appreciable.

The effect of finite kinematics and stretching on the prediction of pull-in voltage for a doubly clamped nanotube is examined by investigating a device with the following characteristics $R_{ext} = 10$ nm, $R_{int} = 0$, $E = 1$ TPa, $H = 100$ nm, $L = 3000$ nm [104]. The central deflection of the nanotube as a function of the applied voltage is shown in Fig. 40 for both with and without stretching. The two vertical lines correspond to reaching unstable behavior (i.e., pull-in voltages). The role of tension stiffening due to the ropelike behavior is quite pronounced in this case.



**Figure 37.** Schematic of finite kinematics configuration of a doubly clamped nanotube device subjected to electrostatic and van der Waals forces. Reprinted with permission from [104], C.-H. Ke et al., *J. Appl. Mech.* 72, 726 (2005). © 2005, American Society of Mechanical Engineers.

**Figure 38.** The effect of finite kinematics on the characteristic of the cantilever nanotube based device (tip displacement vs. voltage). The solid lines illustrate the result accounting for finite kinematics, while the dashed line shows the result when finite kinematics is neglected. Both analyses account for charge concentration at the end of the cantilever nanotube. Reprinted with permission from [104], C.-H. Ke et al., *J. Appl. Mech.* 72, 726 (2005). © 2005, American Society of Mechanical Engineers.

## 3.2.2. Analytical Solutions

In this section, we discuss the analytical solutions of the electromechanical characteristic of the NEMS devices consisting of both cantilever and double-clamped nanotubes. In particular, the pull-in voltage calculations based on the energy method are reported [32, 106].

For nanotube cantilevers (singly clamped), the deflection of the cantilever nanotube can be approximated by the following quadratic function [32]:

$$w(s) \approx \frac{x^2}{L^2} c \qquad (17)$$



**Figure 39.** The effect of charge concentration on the characteristic of the cantilever nanotube based device (tip displacement vs. voltage). The solid line illustrates the deflection curve accounting for charge concentration. The dashed line shows the deflection curve in the absence of charge concentration. Both curves are based on the small deflection model. Reprinted with permission from [104], C.-H. Ke et al., *J. Appl. Mech.* 72, 726 (2005). © 2005, American Society of Mechanical Engineers.

**Figure 40.** Electromechanical characteristic (central displacement–voltage curve) for doubly clamped nanotube device. The dashed line is for small deformation model (pure bending), and the solid line is for finite kinematics model (bending plus stretching). Reprinted with permission from [106]. N. Pugno et al., *J. Appl. Mech.* 72, 445 (2005). © 2005, American Society of Mechanical Engineers.

where $L$ is the length of the nanotube, $c$ is a constant that represents the displacement of the end of the cantilever, and $x$ is the coordinate along the nanotube.

The total energy of the system $E_{total}$ is expressed as

$$E_{total}(c) = E_{elas}(c) + E_{elec}(c) + E_{vdW}(c) \tag{18}$$

where the elastic energy $E_{elas}(c)$, the electrostatic energy $E_{elec}(c)$, and van der Waals energy $E_{vdW}(c)$ can be obtained by integration as

$$E_{elas}(c) = \frac{EI}{2} \int_0^L \left(\frac{d^2w}{dx^2}\right)^2 dx \tag{19a}$$

$$E_{elec, vdW}(c) \approx \int_0^L \frac{dE_{elec, vdW}(r(w(x)))}{dx} dx \tag{19b}$$

The equilibrium condition is reached when the total energy reaches a minimum value, that is,

$$\frac{dE_{total}}{dc} = 0 \tag{20a}$$

Similarly, the instability of the devices (i.e., pull-in) happens when the second-order derivative of total energy equals zero, namely.

$$\frac{d^2E_{total}}{dc^2} = 0 \tag{20b}$$

The van der Waals interaction plays an important role only for a small gap between the nanotubes and substrate (i.e., a few nanometers). Thus it can be neglected in the analysis of NEMS with large gaps. We consider $E_{vdW} \approx 0$ in this analysis.

By assuming that the nanotube's (external) radius $R_{ext}$ is much smaller than the distance $r$ between nanotube and ground plane (i.e., $R_{ext}/r \ll 1$), the pull-in voltage [32], considering the nonlinear finite kinematics and the concentrated charges at the free end, is given by

$$V_{PI} \approx k_v \sqrt{\frac{1 + K_v^{LK}}{1 + K_v^{TP}} \frac{H}{L^2}} \ln\left(\frac{2H}{R_{ext}}\right) \sqrt{\frac{EI}{\varepsilon_0}} \tag{21a}$$

$$k_v \approx 0.85, \quad K_v^{LK} \approx \frac{8H^2}{9L^2}, \quad K_v^{TP} \approx \frac{2.55[R_{ext}(H + R_{ext})^2]^{1/3}}{L} \tag{21b}$$

where subscripts $S$ refer to singly clamped boundary conditions for cantilevers, superscript FK refers to finite kinematics, and TIP refers to the charge concentration.

For doubly clamped nanotubes, the deflection $w(x)$ is assumed such that it satisfies the boundary conditions $w(x = 0, L) = w'(x = 0, L) = 0$ [106], namely,

$$w(z) \approx 16\left[\left(\frac{x}{L}\right)^2 - 2\left(\frac{x}{L}\right)^3 + \left(\frac{x}{L}\right)^4\right]c \tag{22}$$

where $w(x = L/2) = c$ is here an unknown constant that represents the displacement of the central point. The pull-in voltage [106] can be expressed as

$$V_D^{PI} = k_D\sqrt{1 + k_D^{FK}}\frac{H+R}{L^2}\ln\left(\frac{2(H+R)}{R}\right)\sqrt{\frac{EI}{\varepsilon_0}} \tag{23a}$$

$$k_D = \sqrt{\frac{1024}{5\pi S'(c_{PI})}\left(\frac{c_{PI}}{H+R}\right)}, \quad k_D^{FK} = \frac{128}{3003}\left(\frac{c_{PI}}{\rho}\right)^2 \tag{23b}$$

$$\rho^2 = \frac{I}{A} = \frac{R_{ext}^2 + R_{int}^2}{4} \quad S(c) = \sum_{i=1}^{\infty}\left(\frac{1}{\left(\ln\left(\frac{2(H+R)}{R}\right)\right)^i}\sum_{j=i}^{\infty}a_{ij}\left(\frac{c}{(H+R)}\right)^j\right) \tag{23c}$$

Subscripts $D$ refer to double clamped boundary conditions, $c_{PI}$ is the central deflection of the nanotube at the pull-in, and the $\{a_{ij}\}$ in Eq. (23) are known constants [106].

The accuracy of the analytical solutions is verified by the comparison with both numerical integration of the governing equations [104, 106] and experimental measurements (see Section 3.2.3) [32]. The comparison between pull-in voltages evaluated numerically and theoretically for doubly $(D)$ and singly $(S)$ clamped nanotube devices is listed in Table 1 [104]. Columns 6 and 7 in Table 1 compare analytical and numerical pull-in voltage predictions under the assumption of small deformations. Columns 8 and 9 in Table 1 compare analytical and numerical pull-in voltage predictions under the assumption of finite kinematics. The agreement is good (with a maximum discrepancy of 5%).

### 3.2.3. Comparison Between Analytical Predictions and Experiments

In this section, a comparison between analytical predictions and experimental data, for both small deformation and finite kinematics regimes, is presented.

**3.2.3.1. Small Deformation Regime** The nanotweezers experimental data reported by Akita et al. 2001 [8], plotted in Fig. 41, is used to assess the model accuracy under small deformation. In this case, the nanotweezers are equivalent to a nanotube cantilever with length of 2.5 $\mu$m freestanding above an electrode with a gap of 390 nm. Symmetry is

**Table 1.** Comparison between pull-in voltages evaluated numerically (num.) and theoretically (theo.) for doubly $(D)$ and singly $(S)$ clamped nanotube devices, respectively; $E = 1$ TPa, $R_{int} = 0$. For cantilever nanotube device, the symbol (w) denotes that the effect of charge concentration has been included.

| Case | BC | $H$ (nm) | $L$ (nm) | $R = R_{ext}$ (nm) | $V_{PI}$ [V] (theo. linear) | $V_{PI}$ [V] (num. linear) | $V_{PI}$ [V] (theo. nonlinear) | $V_{PI}$ [V] (num. nonlinear) |
|---|---|---|---|---|---|---|---|---|
| 1 | D | 100 | 4000 | 10 | 3.20 | 3.18 | 9.06 | 9.54 |
| 2 | D | 100 | 3000 | 10 | 5.69 | 5.66 | 16.14 | 16.95 |
| 3 | D | 100 | 2000 | 10 | 12.81 | 12.73 | 36.31 | 38.14 |
| 4 | D | 150 | 3000 | 10 | 9.45 | 9.43 | 38.93 | 40.92 |
| 5 | D | 200 | 3000 | 10 | 13.53 | 13.52 | 73.50 | 77.09 |
| 6 | D | 100 | 3000 | 20 | 19.21 | 18.74 | 31.57 | 32.16 |
| 7 | D | 100 | 3000 | 30 | 38.57 | 37.72 | 51.96 | 50.63 |
| 8 | S | 100 | 500 | 10 | 27.28 (w) | 27.05 (w) | 27.52 (w) | 27.41 (w) |
| 9 | S | 100 | 500 | 10 | 27.28 (w) | 27.05 (w) | 30.87 | 31.66 |

Figure 41. Comparison between experimental data and theoretical prediction in the small deformation regime. Reprinted with permission from [32], C.-H. Ke et al., *J. Mech. Phys. Solids* 53, 1314 (2005). © 2005, Elsevier, Ltd.

exploited. In the same figure, a comparison between the analytically predicted nanotube cantilever deflection and the experimentally measured data are shown [32]. The analytical model includes the van der Waals force and charge concentration at the free end of the nanotube cantilever. Model parameters include Young's modulus, $E = 1$ TPa, external radius $R = R_{ext} = 5.8$ nm, and $R_{int} = 0$. The pull-in voltage from the analytical model is 2.34 V, while the experimentally measured pull-in voltage was 2.33 V. It is clear that the analytical prediction and experimental data for the deflection of the nanotube cantilever, as a function of applied voltage, are in very good agreement.

### 3.2.3.2. Finite Kinematics Regime
Experimental data corresponding to the deflection of carbon nanotube cantilevers in the finite kinematics regime were recently obtained by *in situ* SEM measurements [32].

The configuration of the *in situ* measurement is shown in Fig. 42. The electrode was made of silicon wafer coated with 50 nm Au film by e-beam evaporation. This Si chip was attached



Figure 42. Schematic of the experimental configuration employed in the electrostatic actuation of MWNTs. Reprinted with permission from [32], C.-H. Ke et al., *J. Mech. Phys. Solids* 53, 1314 (2005). © 2005, Elsevier, Ltd.

**Figure 43.** Scanning electron microscopy (SEM) images of the deformed carbon nanotube at various bias voltages. Reprinted with permission from [32], C.-H. Ke et al., *J. Mech. Phys. Solids* 53, 1314 (2005). © 2005, Elsevier, Ltd.

onto the side of a Teflon block and mounted to the SEM sample holder at an angle of 93° with respect to the holder plane. The nanotube cantilever fabricated by the method shown in Fig. 4 was placed horizontally and parallel to the electrode surface as schematically shown in Fig. 42. The distance between the top surface and the electron-beam gun was 5 mm, while the distance between the nanotube and the electron-beam gun was measured to be 6.8 mm. By focusing on the electrode surface and adjusting the working distance to be 6.8 mm, a feature on the electrode, which was on the same horizontal plane with the nanotube, was located. Such a feature is schematically marked as a line in Fig. 42. The horizontal distance between the nanotube and the line was controlled by the nanomanipulator and set to 3 $\mu$m. In the circuit, a resistor $R_0 = 1.7$ M$\Omega$ was employed to limit the current. Because the ratio between the length of the nanotube and the gap between the nanotube and electrode is 2.3, the deflection of the nanotube can be considered to be in the finite kinematics regime.

Figure 43(A)–43(E) shows the SEM images of the deflection of the carbon nanotube as it is subject to increasing applied voltages. The feature on the electrode, which is in the same



**Figure 44.** Comparison between experimental data and theoretical prediction in the finite kinematics regime. Reprinted with permission from [32], C.-H. Ke et al., *J. Mech. Phys. Solids* 53, 1314 (2005). © 2005, Elsevier, Ltd.

horizontal plane containing the cantilevered nanotube, is schematically marked as a solid black line in Fig. 43(A)–43(E). These images clearly reveal changes in nanotube deflection and local curvature as a function of applied voltage. A very noticeable effect, although difficult to quantify accurately, is the change in local curvature. The pull-in voltage, $V_{pi}$, was measured to be 48 V. Through digital image processing, the tip deflection as a function of voltage was measured.

The experimentally measured nanotube cantilever deflections, in the finite kinematics regime, are plotted in Fig. 44 [32]. The figure also shows a comparison between analytical prediction and experimental data. The analytical model includes finite kinematics, the van der Waals force, and charge concentration at the free end of the nanotube cantilever. For these predictions, the following parameters were employed: length of the nanotube, $L = 6.8$ $\mu$m; initial gap between nanotube and electrode, $H = 3$ $\mu$m; $R = R_{ext} = 23.5$ nm; $R_{int} = 0$, $E = 1$ TPa. The pull-in voltage given by the analytical analysis is 47.8 V, while the pull-in voltage experimentally measured was 48 V.

## ACKNOWLEDGMENTS

## REFERENCES

1. T. Rueckes, K. Kim, E. Joslevich, G. Y. Tseng, C. Cheung, and C. M. Lieber, *Science* 289, 94 (2000).
2. B. Ilic, H. G. Craighead, and S. Krylov, *J. Appl. Phys.* 95, 3694 (2004).
3. Z. J. Davis, G. Abadal, O. Kuhn, O. Hansen, F. Grey, and A. Boisen, *J. Vac. Sci. Technol. B* 18, 612 (2000).
4. M. L. Roukes, *Physica B* 1, 263 (1999).
5. M. L. Roukes, Technical Digest of the 2000 Solid-State Sensor and Actuator Workshop, 2000.
6. A. N. Cleland and M. L. Roukes, *Appl. Phys. Lett.* 69, 2653 (1996).
7. P. Kim and C. M. Lieber, *Science* 126, 2148 (1999).
8. S. Akita, Y. Nakayama, S. Mizooka, Y. Takano, T. Okawa, Y. Miyatake, S. Yamanaka, M. Tsuji, and T. Nosaka, *Appl. Phys. Lett.* 79, 1691 (2001).
9. A. M. Fennimore, T. D. Yuzvinsky, W. Q. Han, M. S. Fuhrer, J. Cummings, and A. Zettl, *Nature* 424, 408 (2003).
10. J. Kinaret, T. Nord, and S. Viefers, *Appl. Phys. Lett.* 82, 1287 (2003).
11. C.-H. Ke and H. D. Espinosa, *Appl. Phys. Lett.* 85, 681 (2004).
12. V. Sazonova, Y. Yaish, H. Üstünel, D. Roundy, T. Arias, and P. McEuen, *Nature* 431, 284 (2004).
13. S. Iijima, *Nature* 354, 56 (1991).
14. P. Ajayan, *Chem. Rev.* 99, 1787 (1999).
15. C. Journet, W. K. Maser, P. Bernier, A. Loiseau, M. L. delaChapelle, S. Lefrant, P. Deniard, R. Lee, and J. E. Fischer, *Nature* 388, 756 (1997).
16. T. W. Ebbesen and P. M. Ajayan, *Nature* 358, 220 (1992).
17. A. Thess, R. Lee, P. Nikolaev, H. J. Dai, P. Petit, J. Robert, C. H. Xu, Y. H. Lee, S. G. Kim, A. G. Rinzler, D. T. Colbert, G. E. Scuseria, D. Tomanek, J. E. Fischer, and R. E. Smalley, *Science* 273, 483 (1996).
18. W. Z. Li, S. S. Xie, L. X. Qian, B. H. Chang, B. S. Zou, W. Y. Zhou, R. A. Zhao, and G. Wang, *Science* 274, 1701 (1996).
19. D. Qian, G. J. Wagner, W. K. Liu, M. F. Yu, and R. S. Ruoff, *Appl. Mech. Rev.* 55, 495 (2002).
20. P. L. McEuen, M. S. Fuhrer, and H. Park, *IEEE Trans. Nanotechnol.* 1, 78 (2002).
21. T. W. Tombler, C. W. Zhou, L. Alexseyev, J. Kong, H. J. Dai, L. Lei, C. S. Jayanthi, M. J. Tang, and S. Y. Wu, *Nature* 405, 769 (2000).
22. B. Liu, H. T. Johnson, and Y. Huang, *J. Mech. Phys. Solids* 52, 1 (2004).
23. T. Kuzumaki and Y. Mitsuda, *Appl. Phys. Lett.* 85, 1250 (2004).
24. A. Husain, J. Hone, H. W. Ch. Postma, X. M. H. Huang, T. Drake, M. Barbic, A. Scherer, and M. L. Roukes, *Appl. Phys. Lett.* 83, 1240 (2003).
25. E. W. Wong, P. E. Sheehan, and C. M. Lieber, *Science* 277, 1971 (1997).
26. P. A. Williams, S. J. Papadakis, A. M. Patel, M. R. Falvo, S. Washburn, and R. Superfine, *Appl. Phys. Lett.* 82, 805 (2003).
27. R. M. Taylor II and R. Superfine, in "Advanced Interfaces to Scanning Probe Microscopes" (H. S. Nalwa, Ed.), Vol. 2, Academic, New York, 1999.

28. M. R. Falvo, G. J. Clary, R. M. Taylor, V. P. Chi, F. P. Brooks, S. Washburn, and R. Superfine, *Nature* 389, 582 (1997).

29. M. R. Falvo, R. M. Taylor, A. Helser, V. P. Chi, F. P. Brooks, S. Washburn, and R. Superfine, *Nature* 397, 236 (1999).

30. M. F. Yu, M. J. Dyer, G. D. Skidmore, H. W. Rohrs, X. K. Lu, K. D. Ausman, J. R. von Her, and R. S. Ruoff, *Nanotechnology* 10, 244 (1999).

31. M. F. Yu, O. Lourie, M. J. Dyer, K. Moloni, T. F. Kelly, and R. S. Ruoff, *Science* 287, 637 (2000).

32. C.-H. Ke, N. Pugno, B. Peng, and H. D. Espinosa, *J. Mech. Phys. Solids* 53, 1314 (2005).

33. P. A. Smith, C. D. Nordquist, T. N. Jackson, T. S. Mayer, B. R. Martin, J. Mbindyo, and T. E. Mallouk, *Appl. Phys. Lett.* 77, 1399 (2000).

34. X. Q. Chen, T. Saito, H. Yamada, and K. Matsushige, *Appl. Phys. Lett.* 78, 3714 (2001).

35. J. Chung and J. Lee, *Sens. Actuators A* 104, 229 (2003).

36. M. P. Hughes and H. Morgan, *J. Phys. D* 31, 2205 (1998).

37. M. P. Hughes, in "Handbook of Nanoscience, Engineering and Technology" (D. Brenner, S. Lyshevski, G. Iafrate, and W. A. Goddard III, Eds.). CRC Press, Boca Raton, FL, 2002.

38. A. Ramos, H. Morgan, N. G. Green, and A. Castellanos, *J. Phys. D* 31, 2338 (1998).

39. K. Yamamoto, S. Akita, and Y. Nakayama, *J. Phys. D* 31, L34 (1998).

40. R. Krupke, F. Hennrich, H. Löhneysen, and M. M. Kappes, *Science* 301, 344 (2003).

41. S. Lee, D. Lee, R. Morjan, S. Jhang, M. Sveningsson, O. Nerushev, Y. Park, and E. Campbell, *Nano Lett.* 4, 2027 (2004).

42. Y. Huang, X. F. Duan, Q. Q. Wei, and C. M. Lieber, *Science* 291, 630 (2001).

43. M. Fujiwara, E. Oki, M. Hamada, Y. Tanimoto, I. Mukouda, and Y. Shimomura, *J. Phys. Chem. A* 105, 4383 (2001).

44. S. M. Huang, L. M. Dai, and A. W. H. Mau, *J. Phys. Chem. B* 103, 4223 (1999).

45. J. Kong, H. T. Soh, A. M. Cassell, C. F. Quate, and H. Dai, *Nature* 395, 878 (1998).

46. H. J. Dai, *Physics World* 13, 43 (2000).

47. Y. Zhang, A. Chang, J. Cao, Q. Wang, W. Kim, Y. Li, N. Morris, E. Yenilmez, J. Kong, and H. Dai, *Appl. Phys. Lett.* 79, 3155 (2001).

48. H. J. Dai, *Acc. Chem. Res.* 35, 1035 (2002).

49. S. Rao, L. Huang, W. Setyawan, and S. Hong, *Nature* 425, 36 (2003).

50. R. D. Piner, J. Zhu, F. Xu, S. Hong, and C. A. Mirkin, *Science* 283, 661 (1999).

51. M. L. Roukes, *Physics World* 14(2) (2001).

52. R. L. Badzey, G. Zolfagharkhani, A. Gaidarzhy, and P. Mohanty, *Appl. Phys. Lett.* 85, 3587 (2004).

53. H. G. Craighead, *Science* 290, 1532 (2000).

54. Y. Zhu, N. Moldovan, and H. D. Espinosa, *Appl. Phys. Lett.* 86, 013506 (2005).

55. M. Zalalutdinov, B. Ilic, D. Czaplewski, A. Zehnder, H. G. Craighead, and J. M. Parpia, *Appl. Phys. Lett.* 77, 3287 (2000).

56. M. J. Treacy, T. W. Ebbesen, and J. M. Gibson, *Nature* 381, 678 (1996).

57. P. Poncharal, Z. L. Wang, D. Ugarte, and W. A. de Heer, *Science* 283, 1513 (1999).

58. D. S. Greywall, B. Yurje, P. A. Busch, A. N. Pargellis, and R. L. Willett, *Phys. Rev. Lett.* 72, 2992 (1994).

59. A. N. Cleland and M. L. Roukes, *Sens. Actuators A* 72, 256 (1999).

60. K. L. Ekinci, Y. T. Yang, X. M. H. Huang, and M. L. Roukes, *Appl. Phys. Lett.* 81, 2253 (2002).

61. P. Mohanty, D. A. Harrington, K. L. Ekinci, Y. T. Tang, M. J. Murphy, and M. L. Roukes, *Phys. Rev. B* 66, 085416 (2002).

62. M. S. Dresselhaus, G. Dresselhaus, and P. Avouris, "Carbon Nanotubes." Springer, Berlin, 2001.

63. H. Dai, J. H. Hafner, A. G. Rinzler, D. T. Colbert, and R. E. Smalley, *Nature* 384, 147 (1996).

64. S. S. Wong, E. Joselevich, A. T. Woolley, C. L. Cheung, and C. M. Lieber, *Nature* 394, 52 (1998).

65. L. M. Jonsson, S. Axelsson, T. Nord, S. Viefers, and J. M. Kinaret, *Nanotechnology* 15, 1497 (2004).

66. S. J. Tan, A. R. Verschueren, and C. Dekker, *Nature* 393, 49 (1998).

67. C. W. Zhou, J. Kong, and H. J. Dai, *Phys. Rev. Lett.* 84, 5604 (2000).

68. E. D. Minot, Y. Yaish, V. Sazonova, and P. L. McEuen, *Nature* 428, 536 (2004).

69. J. Hu, O. Min, P. Yang, and C. M. Lieber, *Nature* 399, 48 (1999).

70. J. Hu, T. W. Odom, and C. M. Lieber, *Acc. Chem. Res.* 32, 435 (1999).

71. A. M. Morales and C. M. Lieber, *Science* 279, 208 (1998).

72. D. P. Yu, C. S. Lee, I. Bello, X. S. Sun, Y. H. Tang, G. W. Zhou, Z. G. Bai, Z. Zhang, and S. Q. Feng, *Solid State Commun.* 105, 403 (1998).

73. C. Ji and P. C. Searson, *Appl. Phys. Lett.* 81, 4437 (2002).

74. T. C. Wong, C. P. Li, R. Q. Zhang, and S. T. Lee, *Appl. Phys. Lett.* 84, 407 (2004).

75. S. Bhattacharyya, S. K. Saha, and D. Chakravorty, *Appl. Phys. Lett.* 77, 3770 (2000).

76. M. Barbic, J. J. Mock, D. R. Smith, and S. Schultz, *J. Appl. Phys.* 91, 9341 (2002).

77. G. Malandrino, S. T. Finocchiaro, and I. L. Fragalà, *J. Mater. Chem.* 14, 2726 (2004).

78. J. R. Heath and F. K. Legoues, *Chem. Phys. Lett.* 208, 263 (1993).

79. A. B. Greytak, L. J. Lauhon, M. S. Gudiksen, and C. M. Lieber, *Appl. Phys. Lett.* 84, 4176 (2004).

80. K. J. Ziegler, D. M. Lyons, J. D. Holmes, D. Erts, B. Polyakov, H. Olin, K. Svensson, and E. Olsson, *Appl. Phys. Lett.* 84, 4074 (2004).

81. Y. Wu and P. Yang, *Chem. Mater.* 12, 605 (2000).

82. D. Banerjee, J. Y. Lao, D. Z. Wang, J. Y. Huang, Z. F. Ren, D. Steeves, B. Kimball, and M. Sennett, *Appl. Phys. Lett.* 83, 2061 (2003).
83. Y. Dai, Y. Zhang, Y. Q. Bai, and Z. L. Wang, *Chem. Phys. Lett.* 375, 96 (2003).
84. J. J. Mock, S. J. Oldenburg, D. R. Smith, D. A. Schultz, and S. Schultz, *Nano Lett.* 2, 465 (2002).
85. K.-H. Kim, N. Moldovan, and H. D. Espinosa, *Small* 1, 632 (2005).
86. M. D. LaHye, O. Buu, B. Camarota, and K. C. Schwab, *Science* 304, 74 (2004).
87. R. E. Rudd and J. Q. Broughton, *J. Model. Simul. Microsys.* 1, 29 (1999).
88. M. Desquenes, S. V. Rotkin, and N. R. Alaru, *Nanotechnology* 13, 120 (2002).
89. J. A. Moriarty, J. F. Belak, R. E. Rudd, P. Soderlind, F. H. Streitz, and L. H. Yang, *J. Phys.* 14, 2825 (2002).
90. T. Belytschko, W. K. Liu, and B. Moran, "Nonlinear Finite Elements for Continua and Structures." Wiley, Chichester, UK, 2000.
91. T. Belytschko and S. P. Xiao, *Int. J. Multiscale Comp. Engr.* 1, 115 (2003).
92. M. Desquenes, Z. Tang, and N. R. Aluru, *J. Eng. Mater. Tech.* 126, 230 (2004).
93. J. E. Lennard-Jones, *Proc. R. Soc. A* 129, 598 (1930).
94. L. A. Girifalco, M. Hodak, and R. S. Lee, *Phys. Rev. B* 62, 13104 (2000).
95. L. A. Girifalco, *J. Phys. Chem.* 96, 858 (1992).
96. W. Hayt and J. Buck, "Engineering Electromagnetics." 6th Edn. McGraw-Hill, New York, 2001.
97. L. Lou, P. Nordlander, and R. E. Smalley, *Phys. Rev. B* 52, 1429 (1995).
98. M. Krćmar, W. M. Saslow, and A. Zangwill, *J. Appl. Phys.* 93, 3495 (2003).
99. S. V. Rotkin, V. Shrivastava, K. A. Bulashevich, and N. R. Aluru, *Int. J. Nanosci.* 1, 337 (2002).
100. K. A. Bulashevich and S. V. Rotkin, *JETP Lett.* 75, 205 (2002).
101. W. R. Smythe, *J. Appl. Phys.* 27, 917 (1956).
102. P. Keblkinski, S. K. Nayak, P. Zapol, and P. M. Ajayan, *Phys. Rev. Lett.* 89, 255503 (2002).
103. C.-H. Ke and H. D. Espinosa, *J. Appl. Mech.* 72, 721 (2005).
104. C.-H. Ke, H. D. Espinosa, and N. Pugno, *J. Appl. Mech.* 72, 726 (2005).
105. J. P. Salvetat, A. J. Kulik, J. M. Bonard, G. A. D. Briggs, T. Stockli, K. Metenier, S. Bonnamy, F. Beguin, N. A. Burnham, and L. Forro, *Adv. Mater.* 11, 161 (1999).
106. N. Pugno, C.-H. Ke, and H. D. Espinosa, *J. Appl. Mech.* 72, 445 (2005).
107. B. I. Yakobson, M. P. Campbell, C. J. Brabec, and J. Bernholc, *Comp. Mater. Sci.* 8, 341 (1997).

# Index

## A

*Ab initio* calculations, 150, 180, 568, 651, 656
  atomic-level, 163
  isotropic, 193
*Ab initio* identification, 176
*Ab initio* methods, 439
*Ab initio* simulation, 688
*Ab initio* techniques, 446
Abrupt heterojunctions, 48
Absorbing boundary conditions, 81
Accelerated dynamics
  methods, 163–164
  simulations, 200
Acceleration theorem, 49
Acoustic phonons, 67, 69, 610
  energies, 69
  long wavelength, 68
Acoustic waves, 67
Activation barrier diffusion series, 177
Actuator rotor plate, 828
Adatom diffusion mechanism, 164
Adatom exchange mechanism, 164
Adders, 252
  4-bit adder, 259
  2-bit addition, 262
  operation of, 264
Addition energies, 577, 601
Advanced logic design, 767
  selected methods of, 767–768
Advanced micro devices (AMD), 447
Ag bulk, 393
Ag(001) system, 389
  chain of four atoms, 394
  convergence study in, 391
  geometrical setup, 389
  magnetic finite chains in, 394
  nonlocal conductivities, 395
  surface layer of, 389
Airy equation, 59
Airy functions, 486
  derivatives, 486
ALAMODE, 148, 188
AlGaAs barriers, 547
  region, 64
Algebraic equations, 555, 698

Alternating direction implicit (ADI) methods, 81
Alternative nonvolatile memory devices, 532
Amplified spontaneous emission (ASE), 614
Analog computation, 284
  example of, 284
Analog-digital converter (ADC), 357
Analogous circuits, 288
  energy function and local minima in, 288
  state transition, 290
Anderson impurity, 416
  current through, 419
  Green's function for, 416
AND function, 809–810
AND gate, 801
Angle-resolved electron-energy-loss
  spectroscopy (AR-EELS), 651
Angular momentum quantum numbers, 561, 590
Anharmonic interaction potential, 610
Anharmonicity error, 167
Anisotropic elastic continuum, 68
Annealing operation method, 289
  minimum energy state, 290
  steps, 289
Annihilation operators, 721, 739
Antisymmetrization operator, 569
Arbitrary energy wells, 493
  eigenvalues of, 493
Armchair tubes, 216
  smoothening effect, 216
ARPACK library, 556
ARPACK package, 560
Arrhenius expression, 178
Arrhenius-style plot, 166
Associative processing or associative memory, 295
  architecture of, 296
AsV diffusion, 196
  by ring mechanism, 196
Atomic basis states, 23
Atomic basis vectors, 14
Atomic Bloch states, 23
Atomic bonding environment, 656
Atomic force microscopy (AFM), 219
  measurements, 65
Atomic hopping, 151

# M

Macro cluster, 772
Macromodel, 327, 339
    parameters at various temperatures, 327
Macroscopic diffusion constants
    from atomic hopping, 151
Madelung constant, 181
Madelung-corrected Coulomb energy, 181
Magic numbers, 568
Magnetic length, 702
Magnetic vector potential, 548
Magnetoconductance
    of antidot embedded in quantum waveguide, 124
Magnetomotive detection, 834
Magnetomotive drive, 833
Main field propagation, 614
Majority gate, 269
    circuit, 271
    device, 272
    with a single-electron box, 270
    with a Tucker's inverter, 269
Majority logic, 266
    operation, 273
    sketch of, 266
    unit function of, 267
Malfunctioning transistors, 158
Manhattan distance, 310
Many-body
    effects, 38
    system, 10, 38
    theory, 38
Marginal distributions, 735-736
    non-negative, 737
Markov approximation, 609
Markov chain method, 750, 753-754
Markov chain Monte Carlo method, 733
Markov process, 327, 727
Master equation method, 328, 331
    comparison, 334
Matrix eigenvalue, 589, 594
Matrix elements, 66-70
    electron-phonon, 68
    for acoustic and nonpolar-optical phonon
        scattering, 67
    optical, 69
    screened, 75
    zero-order, 70
Matrix equation, 22
Matsubara frequency, 423
Matsubara-poles, 374
Max-Cut problem, 281
Maxwell-Bloch equations, 606, 620
    time-dependent effective, 620
Maxwell-Boltzmann distribution, 474
Maxwell-Boltzmann statistics, 45, 46
Maxwellian distribution, 470
    function, 77

Maxwell's equations, 78-81, 84, 92
    in SI units, 79
    quasi-static solutions, 81, 613
Maxwell's wave equation (MWEs), 606
Mean field approximation, 73, 742
Mean field theory, 383
Mean free path vector, 366
Mesh
    chaining, 87
    charge assignment to, 86
    compute forces on, 86
    force, 88-90
    lines, 85
    potentials, 88
    size, 85
    versus distance between two electrons, 89
Mesoscopic systems, 457
Mesotetraphenyl porphyrins (MTPP), 230
Metal 5,15-di-(4-thiophenyl)-porphyrin
    (MDTP), 228
    relative energy (eV) of, 230
    structure of, 228
Metal-oxide-semiconductor (MOS) capacitors, 470, 479
Metal oxide semiconductor field effect
    transistor (MOSFET), 49, 89, 94, 137-140, 157-158, 320, 356
    device, 139
    fifty-nanometer, 198-199
    nanoscale devices, 189
    silicon-based, 137
    simulated atomistic configurations of, 199
    structures, 139, 190, 193-194, 197
Metal semiconductor field effect transistor
    (MESFET), 84
    3D structure, 84
    example of particle distribution, 84
Methylene bridges, 462
Metropolis Monte Carlo scheme, 666
Microcanonical probability, 383
Microelectromechanical systems (MEMS), 817, 836
    devices, 840
Microscale process simulation, 145
Microscopic scattering probability, 365
Migration barrier, 178
    free-vacancy, 180
Migration energy, 188
MINDO (modified intermediate neglect of
    differential overlap) methods, 440
MINILASE, 91
Minimum-energy barrier diffusion path, 182
Minimum-energy configuration, 184
Minimum energy path (MEP), 169, 187
MNDO theory, 440
Mobile electrons, 547
Mobile extrinsic complexes, 148
Mobility factor, 153
Mode-following methods, 172
Model biological neurons, 797
Model comparison, 509

# R

# Handbook of
# THEORETICAL and COMPUTATIONAL
# NANOTECHNOLOGY

## Edited by Michael Rieth and Wolfram Schommers

The future applications of nanotechnology in high-tech industries require deep understanding of the theoretical and computational aspects of all kinds of materials and devices on a nanometer scale. *Handbook of Theoretical and Computational Nanotechnology* is the first single reference source ever published in the field that offers such an unified approach, covering all of the major topics dealing with theory, modeling, design, and simulations of nanostructured materials and nanodevices, quantum computing, computational chemistry, physics, and biology, nanomechanics, nanomachines, nanoelectronics, nanoprocesses, nanomagnetism, nanooptics, nanomedicines, nanobiotechnology, etc. This 10-volume handbook provides the first ideal introduction and an up-to-date survey of the fascinating new developments and interdisciplinary activities in the whole field presented by scientists working in different subject areas of science, engineering, and medicine. This handbook is the most profound publication on this topic-the first treatment of computational nanotechnology. This outstanding handbook, presented by the world's leading scientists, is the most significant academic title ever published in this research field. This handbook has been divided into 10 thematic volumes by documenting computational treatment of nanomaterials and nanodevices.

**Volume 1:** Basic Concepts, Nanomachines, and Medical Nanodevices
**Volume 2:** Atomistic Simulations—Algorithms and Methods
**Volume 3:** Quantum and Molecular Computing, Quantum Simulations
**Volume 4:** Nanomechanics and Multiscale Modeling
**Volume 5:** Transport Phenomena and Nanoscale Processes
**Volume 6:** Bioinformatics, Nanomedicine, and Drug Design
**Volume 7:** Magnetic Nanostructures and Nanooptics
**Volume 8:** Functional Nanomaterials, Nanoparticles, and Polymer Design
**Volume 9:** Nanocomposites, Nano-Assemblies, and Nanosurfaces
**Volume 10:** Nanodevice Modeling and Nanoelectronics

## KEY FEATURES

O The World's first handbook ever published in the field of theoretical and computational nanotechnology.

O The first comprehensive reference dedicated to all disciplines of science, engineering, and medicine.

O Most up-to-date reference source drawing on the past two decades of pioneering research.

O About 140 Review chapters written by world leading scientists familiar with the current trends of nanotechnology.

O Over 8,000 pages written by 265 authors from 30 countries, truly international.

O 26,000 references, 4124 figures, 374 tables, and thousands of mathematical equations and formula.

O Clearly written, self-contained, timely, authoritative, and most comprehensive contributions.

O Extensive cross-refereeing in each chapter provides reader with a broader range of knowledge.

O Multidisciplinary reference source for scientists, engineers, biologists, medical experts and related professionals.

## READERSHIP

This handbook is an invaluable reference source for scientists, engineers, and biologists working in the field of theoretical and computational nanotechnology. The handbook is intended for a broad audience working in the fields of quantum chemistry, physics, biology, materials science, electrical and electronics engineering, mechanical engineering, optical science, ceramic and chemical engineering, device engineering, aerospace engineering, computer science and technology, information technology, bioinformatics, biotechnology, medical sciences, medicine, surface science, and polymer science and technology.