

---

# Reviews in Computational Chemistry Volume 21

---

Edited by

**Kenny B. Lipkowitz, Raima Larter,  
and Thomas R. Cundari**

Editor Emeritus

**Donald B. Boyd**

 **WILEY-VCH**



---

**Reviews in  
Computational  
Chemistry  
Volume 21**

---



---

# Reviews in Computational Chemistry Volume 21

---

Edited by

**Kenny B. Lipkowitz, Raima Larter,  
and Thomas R. Cundari**

Editor Emeritus

**Donald B. Boyd**

 **WILEY-VCH**

Kenny B. Lipkowitz  
Department of Chemistry  
Ladd Hall 104  
North Dakota State University  
Fargo, North Dakota 58105-5516  
U.S.A.  
kenny.lipkowitz@ndsu.nodak.edu

Thomas R. Cundari  
Department of Chemistry  
University of North Texas  
Box 305070,  
Denton, Texas 76203-5070, U.S.A.  
tomc@unt.edu

Raima Larter  
Department of Chemistry  
Indiana University-Purdue University  
at Indianapolis  
402 North Blackford Street  
Indianapolis, Indiana 46202-3274, U.S.A.  
rlarter@nsf.gov

Donald B. Boyd  
Department of Chemistry  
Indiana University-Purdue University  
at Indianapolis  
402 North Blackford Street  
Indianapolis, Indiana 46202-3274, U.S.A.  
boyd@chem.iupui.edu

Copyright © 2005 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.  
Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400, fax 978-646-8600, or on the web at [www.copyright.com](http://www.copyright.com). Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008.

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor the author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print, however, may not be available in electronic format.

ISBN 0-471-68239-X  
ISSN 1069-3599

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

---

# Preface

---

For the past half-century, chemistry and the practitioners of the chemical sciences existed in a state of stasis; even though noteworthy advances were being made in this science, most of chemistry was divided into the traditional domains of analytical, inorganic, organic, and physical chemistry (except in industrial settings), so little cross-fertilization between disciplines took place. Biochemistry and medicinal chemistry departments existed, but to the average academic chemist, those departments were deemed as not being mainstream and were viewed as being remote or distant from “chemistry.” Chemical engineering likewise seemed far away to most chemists, and other subdisciplines in the chemical sciences such as nuclear chemistry were even more remote to our eyes. At the end of the last century, though, beginning especially in the early 1990s, things began shifting quickly. The blossoming of biological chemistry has put health-related chemical sciences front-and-center, inorganic chemistry is quickly being replaced by materials science, and the evolution of physical chemistry from a science that was once dominated by gas phase reaction dynamics to one that is beginning to focus on biological and materials systems in the condensed phase, all represent an abrupt change in our discipline.

Computational chemistry is itself multidisciplinary, transcending traditional boundaries that separate chemistry from biology, physics, mathematics, and computer science. It took a small leap of faith for the editors of this book series to accept that the future of chemistry might be different from its then present state, and that new opportunities in both science and technology could appear—as they soon did. Along with the need to be informed about these new opportunities came an understanding that people from different subdisciplines of the chemical sciences who wanted to use computing in their work would also need to learn, quickly, about how others solved computational problems in their areas of study. A quantum chemist would need to learn about protein folding, for example, because of the new funding opportunities and interest in that area, whereas an organic chemist would need to learn about drug design strategies to keep pace with the changes in that discipline. Ongoing studies focusing on both the basic science and the technological spinoffs of chemical biology and nanoscience along with the imminent seamless merger of chemistry with chemical engineering indicate to us that this cross-fertilization of ideas

and multidisciplinary approaches will continue at an even faster pace in the near future.

To expedite this learning process, we set out to develop a series of pedagogically driven reviews on various topics in computational chemistry that are of use to the scientific community. The focus of these reviews is not directed toward the theorist, but instead, to the novice computational chemist wanting to learn applied theory as well as to bench chemists wanting to use computation as an adjunct to their laboratory work. In addition, this book series appeals to the seasoned computational chemist who desires a quick introduction to methodology from other areas of computational science and to scientists in other disciplines who want to begin their studies at the molecular level and intend to use computational methods to assist them in their work. In this, the 21<sup>st</sup> volume of *Reviews in Computational Chemistry*, we present six tutorials on a diverse set of topics typical of the new breadth of the modern science of chemistry, covering material from solid state chemistry to the modeling of gene regulatory networks.

In Chapter 1, Roberto Dovesi, Bartolomeo Civalleri, Roberto Orlando, Carla Roetti, and Victor Saunders provide a tutorial on calculating structures, properties, and energies of solid state chemical systems using ab initio quantum methods. Concepts covering the direct lattice, reciprocal lattice, Bloch theorem and periodic boundaries are followed by the development of the one-electron Hamiltonian in a lucid section on invariance properties in a crystal. A discussion of band structure follows with examples of one-, two-, and three-dimensional systems. Cohesive energy, polymorphism, magnetic phases, and positional isomorphous phases are described. Modeling surfaces and interfaces is described next, which is then followed by a tutorial on modeling of defects in solids. What a novice modeler is allowed to do and pitfalls for that person to avoid are given throughout the chapter. As crystal engineering and nanoscience continues to grow, aspects of simulating solid state behavior will become more and more important and the material described in Chapter 1 will become more relevant to many bench chemists.

The theme of quantum mechanics is continued by Patrick Bultinck, Xavier Gironés, and Ramon Carbó-Dorca in Chapter 2. This tutorial focuses on molecular quantum similarity measures (MQSM). The authors describe why similarity (and dissimilarity) measures are important and useful, especially from a quantum mechanical perspective. The authors begin by providing a general overview of molecular similarity along with the associated vocabulary, and then they compare and contrast quantum similarity to more traditional molecular descriptors like feature counts, fragment descriptors, field-based descriptors, and topologic/topographic indices. The electron density as a molecular descriptor is introduced, and quantum molecular similarity measures are defined and described. Then manipulation and representations of the similarity matrix, details about electron densities for use in MQSM, and issues related to alignment of molecules for such measurements follow. The physical connotations of self-similarity measures and aspects of chirality are



described before a section detailing the mathematics associated with the topic is presented.

In Chapter 3, Jean-Loup Faulon, Donald Visco, Jr, and Diana Roe introduce us to the important topic of enumerating molecules. Enumeration means two things: first to list things separately, one by one, and second to determine the number of, or to count things; both aspects of enumeration in chemistry are described by these authors. The foremost application of enumeration is in structure determination. But for other applications, such as molecular design (especially with the advent of combinatorial technologies), enumeration takes a central role for, e.g., constructing virtual libraries, testing hypotheses, and optimizing experiments. The authors begin by describing how enumeration is accomplished. Graph theory and algorithms that rely on it are introduced. The counting of labeled and unlabeled graphs, with particular attention to Pólya's theorem, show the reader how to determine the number of isomers given a molecular formula. The number of isomers of acyclic compounds (including the possible number of stereoisomers), the number of benzenoids, and the number of molecular cages (fullerenes and nanotubes) are examples provided by the authors. Elementary definitions used to count, enumerate, and sample molecules are presented, replete with simple examples for the novice to study. In the second part of the chapter, the authors review the practical applications of molecular enumeration and give the reader pointers to relevant techniques and existing codes. In particular, the number of isomers for a specific molecular series is given, popular structure elucidation codes are reviewed, computer-aided structure-elucidation successes are surveyed, and the connections between structure enumeration and combinatorial library design are established.

A foundational method that computational chemistry groups are built on in pharmaceutical companies is quantitative structure-activity relationships (QSAR). QSAR, until recently, has not been part of academia, but that is beginning to change because of impetus from patent and technology transfer offices that are encouraging academics to consider fiduciary responsibilities; many academicians are now interested in topics previously of concern only to business, and QSAR can help them design new materials to meet their financial goals. The average bench chemist may be able to understand what a molecular descriptor is, and even how to generate various descriptors for QSAR and molecular design studies, but the subtleties of deciding how to select those variables from an enormous list of possibilities is usually beyond their current area of expertise. In Chapter 4, David J. Livingstone and David W. Salt clarify for us the many intricacies of variable selection. They begin with the concept of dimension reduction (reducing the dimensionality of a multivariate dataset while retaining most of the information contained therein), then they describe variable reduction, which seeks to reduce redundancy in a dataset by identifying variables that are correlated with one another, before they finally describe variable selection. In this chapter, the novice modeler will learn how to identify multicollinearity, decide which variables to eliminate, and learn about

supervised and unsupervised elimination methods. Variable selection techniques are introduced and compared with one another. Methods including ordinary least squares, ridge regression, principal components regression, partial least squares, and continuum regression techniques are described in a clear, cogent manner. Other methods needed for variable selection are also described, including best variable subset selection, forward inclusion, backward elimination, and stepwise regression. This chapter removes the “black-box” mentality associated with many commercially available software packages for QSAR and for molecular design by highlighting many of the pitfalls one should avoid when interpreting the results from those programs.

In Volume 19 of this book series, Gene Lamm introduced the theoretical underpinning and associated mathematics of the Poisson–Boltzmann equation (PBE). Here, in Chapter 5, Nathan A. Baker discusses biomolecular applications of these methods. An introduction to biomolecular electrostatics is first given, where Baker highlights the factors influencing those interactions during molecular simulations. A review of explicit versus implicit aqueous environments is presented, and a historical account of implicit solvent methods is given. Warnings about when not to use such simplifications are stressed, and situations where implicit solvent treatments are applicable are presented. How to treat polar and nonpolar interactions is given next; leading references to scaled particle theory methods, solvent accessible surface area methods, and the like are abundant in this chapter. After a brief overview of the PBE, Baker describes the common methods that solve it, with a major emphasis placed on discretization techniques, including finite difference, boundary element, and finite element methods. Multilevel solvers and parallel methods for solving the PBE are also presented, and a list of available software for computational electrostatics is given. The author then takes the reader through applications of PB methods to determine solvation-free energies, conformational-free energies, binding-free energies, titration calculations, and other applications.

Chapter 6 extends even further computational methods in the area of biology. Data sources and computational approaches for generating models of gene regulatory networks (GRNs) are described by Baltazar Aguda, Georghe Craciun, and Rengul Cetin-Atalay. A GRN refers to a set of molecules and interactions that affect the expression of genes located in the DNA of a cell. The authors begin with a formal representation of GRNs and then use the *Lac* Operon as an example. Hierarchies of GRN models are described, and a guide to databases and knowledge bases on the Internet is provided. Because bioinformatics is moving toward the creation of tools, languages, and software for the integration of heterogeneous biological data, something chemists will soon need to deal with in our collaboration with biologists, the authors provide ontologies, a set of controlled and unambiguous vocabulary for describing objects and concepts, for GRN modeling. How to extract models from databases of pathways is then described, and pathway and dynamic analysis tools for GRNs are introduced. The analysis of complex gene regulatory networks is in its infancy. As chemists interact

more with biologists and as chemists take a more global view of biological systems, the need to comprehend how parts of a cell are integrated and interact to determine the system's behavior becomes more important to us. This chapter provides a doorway to that collaborative world.

We are delighted to report that the Institute for Scientific information, Inc. (ISI) rates the *Reviews in Computational Chemistry* book series in the top ten in the category of "general" journals and periodicals. The reason for these accomplishments rests firmly on the shoulders of the authors whom we have contacted to provide the pedagogically driven reviews that have made this ongoing book series so popular. To those authors we are especially grateful.

We are also glad to note that our publisher has plans to make our most recent volumes available in an online form through Wiley InterScience. Please check the Web (<http://www.interscience.wiley.com/onlinebooks>) or contact [reference@wiley.com](mailto:reference@wiley.com) for the latest information. For readers who appreciate the permanence and convenience of bound books, these will, of course, continue.

We thank the authors of this and previous volumes for their excellent chapters.

Kenny B. Lipkowitz  
Fargo

Raima Larter  
Indianapolis

Thomas R. Cundari  
Denton

September 2004



---

# Contents

---

<b>1. Ab Initio Quantum Simulation in Solid State Chemistry</b>	<b>1</b>
<i>Roberto Dovesi, Bartolomeo Civalleri, Roberto Orlando, Carla Roetti, and Victor R. Saunders</i>	
Introduction	1
Translation Invariance Properties in a Crystal	6
The Direct Lattice	7
The Reciprocal Lattice	11
Bloch Theorem and Periodic Boundary Conditions	12
One-Electron Electrostatic Hamiltonian	16
Discussion of Band Structure Through a Few Simple Examples	21
A Monoatomic Linear Chain	21
A Two-Dimensional Periodic Example: Graphite	23
Three-Dimensional Periodic Examples	34
From the Band Structure to the Total Energy	37
Use of Symmetry in Reciprocal Space	40
Total Energy, Energy Differences, and Derivatives	43
Cohesive Energy	44
Polymorphism	52
Magnetic Phases	54
Positional Isomorphous Phases	56
Energy Derivatives	57
Modeling Surfaces and Interfaces	66
The Slab Model	66
Specifying the Surface Plane—Miller Indices	67
Choosing the Surface Termination	68
Surface Formation Energy and Stability	70
Surface Relaxation and Reconstruction	71
Vicinal Surfaces—Modeling Steps and Kinks	74
Adsorption on Surfaces	74
Interfaces	77
Modeling Defective Systems	80
Defects in Solids	80
How to Model a Defect	81

The Supercell Approach	83
Defect Formation Energy	85
Examples	85
Acknowledgments	104
Appendix 1: Available Periodic Programs	104
Appendix 2: Performance of the Periodic Program Crystal	106
Appendix 3: Acronyms	111
References	112
 2. <b>Molecular Quantum Similarity: Theory and Applications</b>	<b>127</b>
<i>Patrick Bultinck, Xavier Gironés, and Ramon Carbó-Dorca</i>	
Introduction	127
Basic Elements of Molecular Similarity	128
The Electron Density as Molecular Descriptor	132
Molecular Quantum Similarity	134
Extension to Other Operators	137
Stochastic Manipulations and Graphical Representations of the Similarity Matrix	140
Electron Densities for Molecular Quantum Similarity	143
The Alignment Issue in Molecular Quantum Similarity	154
Statement of the Problem	154
Quantum Similarity Maximization—MaxiSim and QSSA	157
Structural Alignment	161
Comparison of Alignment Techniques	163
Quantum Similarity Indices	164
Quantum Atoms-in-Molecules Similarity	167
The Hirshfeld Approach	168
AIM-Based Methods	169
Atom-Centered Basis Function Approach	170
Physical Connotations of (Self) Similarity Measures	171
Chirality and the Holographic Electron Density Theorem	177
Mathematical Aspects of Quantum Similarity	180
The Cramer Steroid Set—A Worked Out Example of MQS	191
Conclusions	196
Acknowledgments	196
References	197
 3. <b>Enumerating Molecules</b>	<b>209</b>
<i>Jean-Loup Faulon, Donald P. Visco, Jr., and Diana Roe</i>	
Enumerating Molecules: Why	209
Enumerating Molecules: How	211
From Graph Theory to Chemistry	211
Counting Structures: How Many Isomers Has Decane?	215

Enumerating Structures: Are There any Isomers of Decane Having Seven Methyl Groups?	233
Enumerating Labeled and Unlabeled Graphs	233
Enumerating Molecules	237
Sampling Structures: What is the Decane Isomer With the Highest Boiling Point?	255
Sampling Labeled and Unlabeled Graphs	256
Sampling Molecules	257
Enumerating Molecules: What are the Uses?	261
Chemical Information	261
Structure Elucidation	266
Combinatorial Library Design	270
Molecular Design with Inverse-QSAR	272
Conclusion and Future Directions	274
Acknowledgments	275
References	275
 <b>4. Variable Selection—Spoilt for Choice?</b>	 <b>287</b>
<i>David J. Livingstone and David W. Salt</i>	
Introduction	287
The Problem	290
Dimension Reduction	291
Variable Elimination	295
Why We Eliminate Variables	296
More Variables Than Objects	296
Multicollinearity	299
Identification of Multicollinearity	304
Which Variables Do We Eliminate?	307
Unsupervised Elimination	307
Supervised Elimination	309
Variable Selection	309
Ordinary Least Squares	311
Ridge Regression	311
Principal Component Regression, Partial Least Squares, and Continuum Regression	314
Best Variable Subset Selection	318
Forward Inclusion	323
Backward Elimination	323
Stepwise Regression	324
All Subset Regression	326
Other Stopping Rules	327
Case Study	329
Stepwise Regression	330
Best Subset Regression	332

Supervised and Unsupervised Variable Selection	334
Published Variable Selection Methods	339
Conclusions	341
Acknowledgments	342
Appendix	342
References	343
 <b>5. Biomolecular Applications of Poisson–Boltzmann Methods</b>	 <b>349</b>
<i>Nathan A. Baker</i>	
Introduction to Biomolecular Electrostatics	349
Simplifying the System: Implicit Solvent Methods	351
Polar Interactions	351
Nonpolar Interactions	352
Poisson–Boltzmann Theory: A Brief Overview	354
Solving the Poisson–Boltzmann Equation	357
Discretization Methods	357
Multilevel Solvers	359
Parallel Methods	360
Software for Computational Electrostatics	360
Applying Poisson–Boltzmann Methods	362
Solvation Free Energies	362
Conformational Free Energies	365
Binding Free Energies	367
Titration Calculations	369
Other Applications	370
Conclusions	371
Acknowledgments	372
References	372
 <b>6. Data Sources and Computational Approaches for</b>	
<b>Generating Models of Gene Regulatory Networks</b>	<b>381</b>
<i>Baltazar D. Aguda, Georghe Cracium, and Rengul Cetin-Atalay</i>	
Introduction	381
Formal Representation of GRNs	383
An Example of a GRN: The Lac Operon	384
Hierarchies of GRN Models: From Probabilistic Graphs to	
Deterministic Models	388
A Guide to Databases and Knowledgebases on the Internet	390
Pathway Databases and Platforms	393
Ontologies for GRN Modeling	395
Current Gene, Interaction, and Pathway Ontologies	395
Whole-Cell Modeling Platforms	396
Ontology for Modeling Multiscale and Incomplete	
Networks	397



---

An Ontology for Cellular Processes	397
The PATIKA Pathway Ontology	400
Extracting Models from Pathways Databases	401
Pathway and Dynamic Analysis Tools for GRNs	402
Global Network Properties	403
Recurring Network Motifs	403
Identifying Pathway Channels in Networks:	
Extreme Pathway Analysis	404
Network Stability Analysis	404
Predicting Dynamics and Bistability from Network	
Structure Alone	406
Concluding Remarks	407
References	408
<b>Author Index</b>	<b>413</b>
<b>Subject Index</b>	<b>431</b>



# Contributors

---

**Baltazar Aguda**, Bioinformatics Institute, 30 Biopolis Street, #07-01 Matrix, Singapore 138671. (Electronic mail: [baltazar@bii.a-star.edu.sg](mailto:baltazar@bii.a-star.edu.sg))

**Nathan A. Baker**, Department of Biochemistry and Molecular Biophysics, Center for Computational Biology, Washington University in St. Louis School of Medicine, 700 S. Euclid Ave., Campus Box 8036, St. Louis, MO 63110, U.S.A. (Electronic mail: [baker@biochem.wustl.edu](mailto:baker@biochem.wustl.edu))

**Patrick Bultinck**, Department of Inorganic and Physical Chemistry, Ghent University, Krijgslaan 281 (S-3), B-9000, Gent, Belgium. (Electronic mail: [patrick.bultinck@ugent.be](mailto:patrick.bultinck@ugent.be))

**Ramon Carbó-Dorca**, Institute of Computational Chemistry, University of Girona, Campus de Montilivi, 17005 Girona, Spain. (Electronic mail: [director@iqc.udg.es](mailto:director@iqc.udg.es))

**Rengul Cetin-Atalay**, MBG Department, Faculty of Science B Building, Bilkent University main campus, 06533 Bilkent, Ankara, Turkey. (Electronic Mail: [rengul@bilkent.edu.tr](mailto:rengul@bilkent.edu.tr))

**Bartolomeo Civalieri**, Dipartimento di Chimica IFM, Università di Torino, via Giuria 5, Torino I-10125, Italy. (Electronic mail: [bartolomeo.civalieri@unito.it](mailto:bartolomeo.civalieri@unito.it))

**Georghe Craciun**, Mathematical Biosciences Institute, The Ohio State University, 231 W. 18th Avenue, Columbus, Ohio 43210 U.S.A. (Electronic mail: [gcraciun@mbi.osu.edu](mailto:gcraciun@mbi.osu.edu))

**Roberto Dovesi**, Dipartimento di Chimica IFM, Università di Torino, via Giuria 5, Torino I-10125, Italy. (Electronic mail: [roberto.dovesi@unito.it](mailto:roberto.dovesi@unito.it))

**Jean-Loup Faulon**, Computational Biology and Evolutionary Computing Department, Sandia National Laboratories, P.O. Box 969, MS 9951, Livermore, CA 94551-0969 U.S.A. (Electronic mail: [jfaulon@sandia.gov](mailto:jfaulon@sandia.gov))

**Xavier Gironés**, Computational Chemistry, Medicinal Chemistry Department, AstraZeneca R&D, Pepparedsleden 1, S-43183 Mölndal, Sweden. (Electronic mail: [xavier.girones@astrazeneca.com](mailto:xavier.girones@astrazeneca.com))

**David J. Livingstone**, ChemQuest, Delamere House, 1 Royal Crescent, Sandown, Isle of Wight, U.K. PO36 8LZ. (Electronic mail: [davel@chemquest.uk.com](mailto:davel@chemquest.uk.com))

**Roberto Orlando**, Dipartimento di Scienze e Tecnologie Avanzate, Università del Piemonte Orientale, Corso Borsalino 54, I-15100 Alessandria, Italy. (Electronic mail: [roberto.orlando@unito.it](mailto:roberto.orlando@unito.it))

**Diana Roe**, Biosystems Research Department, Sandia National Laboratories, PO Box 969, Livermore, CA 94551, U.S.A. (Electronic mail: [dcroe@sandia.gov](mailto:dcroe@sandia.gov))

**David W. Salt**, Department of Mathematics, Buckingham Building, University of Portsmouth, Portsmouth, U.K. P01 3HE. (Electronic mail: [david.salt@port.ac.uk](mailto:david.salt@port.ac.uk))

**Victor R. Saunders**, Dipartimento di Chimica IFM, Università di Torino, via Giuria 5, Torino I-10125, Italy. (Electronic mail: [v.r.saunders@dl.ac.uk](mailto:v.r.saunders@dl.ac.uk))

**Donald P. Visco, Jr.**, Department of Chemical Engineering, Tennessee Technological University, PO Box 5013, Cookeville, TN 38505, U.S.A. (Electronic mail: [dvisco@tntech.edu](mailto:dvisco@tntech.edu))

---

# Contributors to Previous Volumes

---

## Volume 1 (1990)

**David Feller and Ernest R. Davidson**, Basis Sets for Ab Initio Molecular Orbital Calculations and Intermolecular Interactions.

**James J. P. Stewart**, Semiempirical Molecular Orbital Methods.

**Clifford E. Dykstra, Joseph D. Augspurger, Bernard Kirtman, and David J. Malik**, Properties of Molecules by Direct Calculation.

**Ernest L. Plummer**, The Application of Quantitative Design Strategies in Pesticide Design.

**Peter C. Jurs**, Chemometrics and Multivariate Analysis in Analytical Chemistry.

**Yvonne C. Martin, Mark G. Bures, and Peter Willett**, Searching Databases of Three-Dimensional Structures.

**Paul G. Mezey**, Molecular Surfaces.

**Terry P. Lybrand**, Computer Simulation of Biomolecular Systems Using Molecular Dynamics and Free Energy Perturbation Methods.

**Donald B. Boyd**, Aspects of Molecular Modeling.

**Donald B. Boyd**, Successes of Computer-Assisted Molecular Design.

**Ernest R. Davidson**, Perspectives on Ab Initio Calculations.

## Volume 2 (1991)

**Andrew R. Leach**, A Survey of Methods for Searching the Conformational Space of Small and Medium-Sized Molecules.

**John M. Troyer** and **Fred E. Cohen**, Simplified Models for Understanding and Predicting Protein Structure.

**J. Phillip Bowen** and **Norman L. Allinger**, Molecular Mechanics: The Art and Science of Parameterization.

**Uri Dinur** and **Arnold T. Hagler**, New Approaches to Empirical Force Fields.

**Steve Scheiner**, Calculating the Properties of Hydrogen Bonds by Ab Initio Methods.

**Donald E. Williams**, Net Atomic Charge and Multipole Models for the Ab Initio Molecular Electric Potential.

**Peter Politzer** and **Jane S. Murray**, Molecular Electrostatic Potentials and Chemical Reactivity.

**Michael C. Zerner**, Semiempirical Molecular Orbital Methods.

**Lowell H. Hall** and **Lemont B. Kier**, The Molecular Connectivity Chi Indexes and Kappa Shape Indexes in Structure-Property Modeling.

**I. B. Bersuker** and **A. S. Dimoglo**, The Electron-Topological Approach to the QSAR Problem.

**Donald B. Boyd**, The Computational Chemistry Literature.

## Volume 3 (1992)

**Tamar Schlick**, Optimization Methods in Computational Chemistry.

**Harold A. Scheraga**, Predicting Three-Dimensional Structures of Oligopeptides.

**Andrew E. Torda** and **Wilfred F. van Gunsteren**, Molecular Modeling Using NMR Data.

**David F. V. Lewis**, Computer-Assisted Methods in the Evaluation of Chemical Toxicity.

## **Volume 4   (1993)**

**Jerzy Cioslowski**, Ab Initio Calculations on Large Molecules: Methodology and Applications.

**Michael L. McKee** and **Michael Page**, Computing Reaction Pathways on Molecular Potential Energy Surfaces.

**Robert M. Whitnell** and **Kent R. Wilson**, Computational Molecular Dynamics of Chemical Reactions in Solution.

**Roger L. DeKock**, **Jeffrey D. Madura**, **Frank Rioux**, and **Joseph Casanova**, Computational Chemistry in the Undergraduate Curriculum.

## **Volume 5   (1994)**

**John D. Bolcer** and **Robert B. Hermann**, The Development of Computational Chemistry in the United States.

**Rodney J. Bartlett** and **John F. Stanton**, Applications of Post-Hartree-Fock Methods: A Tutorial.

**Steven M. Bachrach**, Population Analysis and Electron Densities from Quantum Mechanics.

**Jeffrey D. Madura**, **Malcolm E. Davis**, **Michael K. Gilson**, **Rebecca C. Wade**, **Brock A. Luty**, and **J. Andrew McCammon**, Biological Applications of Electrostatic Calculations and Brownian Dynamics Simulations.

**K. V. Damodaran** and **Kenneth M. Merz, Jr.**, Computer Simulation of Lipid Systems.

**Jeffrey M. Blaney** and **J. Scott Dixon**, Distance Geometry in Molecular Modeling.

**Lisa M. Balbes**, **S. Wayne Mascarella**, and **Donald B. Boyd**, A Perspective of Modern Methods in Computer-Aided Drug Design.

## **Volume 6   (1995)**

**Christopher J. Cramer** and **Donald G. Truhlar**, Continuum Solvation Models: Classical and Quantum Mechanical Implementations.

**Clark R. Landis, Daniel M. Root, and Thomas Cleveland,** Molecular Mechanics Force Fields for Modeling Inorganic and Organometallic Compounds.

**Vassilios Galiatsatos,** Computational Methods for Modeling Polymers: An Introduction.

**Rick A. Kendall, Robert J. Harrison, Rik J. Littlefield, and Martyn F. Guest,** High Performance Computing in Computational Chemistry: Methods and Machines.

**Donald B. Boyd,** Molecular Modeling Software in Use: Publication Trends.

**Eiji Ōsawa and Kenny B. Lipkowitz,** Appendix: Published Force Field Parameters.

## **Volume 7 (1996)**

**Geoffrey M. Downs and Peter Willett,** Similarity Searching in Databases of Chemical Structures.

**Andrew C. Good and Jonathan S. Mason,** Three-Dimensional Structure Database Searches.

**Jiali Gao,** Methods and Applications of Combined Quantum Mechanical and Molecular Mechanical Potentials.

**Libero J. Bartolotti and Ken Flurchick,** An Introduction to Density Functional Theory.

**Alain St-Amant,** Density Functional Methods in Biomolecular Modeling.

**Danya Yang and Arvi Rauk,** The A Priori Calculation of Vibrational Circular Dichroism Intensities.

**Donald B. Boyd,** Appendix: Compendium of Software for Molecular Modeling.

## **Volume 8 (1996)**

**Zdenek Slanina, Shyi-Long Lee, and Chin-hui Yu,** Computations in Treating Fullerenes and Carbon Aggregates.



**Gernot Frenking, Iris Antes, Marlis Böhme, Stefan Dapprich, Andreas W. Ehlers, Volker Jonas, Arndt Neuhaus, Michael Otto, Ralf Stegmann, Achim Veldkamp, and Sergei F. Vyboishchikov,** Pseudopotential Calculations of Transition Metal Compounds: Scope and Limitations.

**Thomas R. Cundari, Michael T. Benson, M. Leigh Lutz, and Shaun O. Sommerer,** Effective Core Potential Approaches to the Chemistry of the Heavier Elements.

**Jan Almlöf and Odd Gropen,** Relativistic Effects in Chemistry.

**Donald B. Chesnut,** The Ab Initio Computation of Nuclear Magnetic Resonance Chemical Shielding.

## **Volume 9    (1996)**

**James R. Damewood, Jr.,** Peptide Mimetic Design with the Aid of Computational Chemistry.

**T. P. Straatsma,** Free Energy by Molecular Simulation.

**Robert J. Woods,** The Application of Molecular Modeling Techniques to the Determination of Oligosaccharide Solution Conformations.

**Ingrid Pettersson and Tommy Liljefors,** Molecular Mechanics Calculated Conformational Energies of Organic Molecules: A Comparison of Force Fields.

**Gustavo A. Arteca,** Molecular Shape Descriptors.

## **Volume 10    (1997)**

**Richard Judson,** Genetic Algorithms and Their Use in Chemistry.

**Eric C. Martin, David C. Spellmeyer, Roger E. Critchlow, Jr., and Jeffrey M. Blaney,** Does Combinatorial Chemistry Obviate Computer-Aided Drug Design?

**Robert Q. Topper,** Visualizing Molecular Phase Space: Nonstatistical Effects in Reaction Dynamics.

**Raima Larter and Kenneth Showalter,** Computational Studies in Nonlinear Dynamics.

**Stephen J. Smith** and **Brian T. Sutcliffe**, The Development of Computational Chemistry in the United Kingdom.

## **Volume 11 (1997)**

**Mark A. Murcko**, Recent Advances in Ligand Design Methods.

**David E. Clark**, **Christopher W. Murray**, and **Jin Li**, Current Issues in De Novo Molecular Design.

**Tudor I. Oprea** and **Chris L. Waller**, Theoretical and Practical Aspects of Three-Dimensional Quantitative Structure–Activity Relationships.

**Giovanni Greco**, **Ettore Novellino**, and **Yvonne Connolly Martin**, Approaches to Three-Dimensional Quantitative Structure–Activity Relationships.

**Pierre-Alain Carrupt**, **Bernard Testa**, and **Patrick Gaillard**, Computational Approaches to Lipophilicity: Methods and Applications.

**Ganesan Ravishanker**, **Pascal Auffinger**, **David R. Langley**, **Bhuvabhootla Jayaram**, **Matthew A. Young**, and **David L. Beveridge**, Treatment of Counterions in Computer Simulations of DNA.

**Donald B. Boyd**, Appendix: Compendium of Software and Internet Tools for Computational Chemistry.

## **Volume 12 (1998)**

**Hagai Meirovitch**, Calculation of the Free Energy and the Entropy of Macromolecular Systems by Computer Simulation.

**Ramzi Kutteh** and **T. P. Straatsma**, Molecular Dynamics with General Holonomic Constraints and Application to Internal Coordinate Constraints.

**John C. Shelley** and **Daniel R. Bérard**, Computer Simulation of Water Physisorption at Metal–Water Interfaces.

**Donald W. Brenner**, **Olga A. Shenderova**, and **Denis A. Areshkin**, Quantum-Based Analytic Interatomic Forces and Materials Simulation.

**Henry A. Kurtz** and **Douglas S. Dudis**, Quantum Mechanical Methods for Predicting Nonlinear Optical Properties.

**Chung F. Wong**, **Tom Thacher**, and **Herschel Rabitz**, Sensitivity Analysis in Biomolecular Simulation.

**Paul Verwer and Frank J. J. Leusen**, Computer Simulation to Predict Possible Crystal Polymorphs.

**Jean-Louis Rivail and Bernard Maigret**, Computational Chemistry in France: A Historical Survey.

### **Volume 13    (1999)**

**Thomas Bally and Weston Thatcher Borden**, Calculations on Open-Shell Molecules: A Beginner's Guide.

**Neil R. Kestner and Jaime E. Combariza**, Basis Set Superposition Errors: Theory and Practice.

**James B. Anderson**, Quantum Monte Carlo: Atoms, Molecules, Clusters, Liquids, and Solids.

**Anders Wallqvist and Raymond D. Mountain**, Molecular Models of Water: Derivation and Description.

**James M. Briggs and Jan Antosiewicz**, Simulation of pH-dependent Properties of Proteins Using Mesoscopic Models.

**Harold E. Helson**, Structure Diagram Generation.

### **Volume 14    (2000)**

**Michelle Miller Francl and Lisa Emily Chirlian**, The Pluses and Minuses of Mapping Atomic Charges to Electrostatic Potentials.

**T. Daniel Crawford and Henry F. Schaefer III**, An Introduction to Coupled Cluster Theory for Computational Chemists.

**Bastiaan van de Graaf, Swie Lan Njo, and Konstantin S. Smirnov**, Introduction to Zeolite Modeling.

**Sarah L. Price**, Toward More Accurate Model Intermolecular Potentials for Organic Molecules.

**Christopher J. Mundy, Sundaram Balasubramanian, Ken Bagchi, Mark E. Tuckerman, Glenn J. Martyna, and Michael L. Klein**, Nonequilibrium Molecular Dynamics.

**Donald B. Boyd and Kenny B. Lipkowitz**, History of the Gordon Research Conferences on Computational Chemistry.

**Mehran Jalaie** and **Kenny B. Lipkowitz**, Appendix: Published Force Field Parameters for Molecular Mechanics, Molecular Dynamics, and Monte Carlo Simulations.

## **Volume 15 (2000)**

**F. Matthias Bickelhaupt** and **Evert Jan Baerends**, Kohn-Sham Density Functional Theory: Predicting and Understanding Chemistry.

**Michael A. Robb**, **Marco Garavelli**, **Massimo Olivucci**, and **Fernando Bernardi**, A Computational Strategy for Organic Photochemistry.

**Larry A. Curtiss**, **Paul C. Redfern**, and **David J. Frurip**, Theoretical Methods for Computing Enthalpies of Formation of Gaseous Compounds.

**Russell J. Boyd**, The Development of Computational Chemistry in Canada.

## **Volume 16 (2000)**

**Richard A. Lewis**, **Stephen D. Pickett**, and **David E. Clark**, Computer-Aided Molecular Diversity Analysis and Combinatorial Library Design.

**Keith L. Peterson**, Artificial Neural Networks and Their Use in Chemistry.

**Jörg-Rüdiger Hill**, **Clive M. Freeman**, and **Lalitha Subramanian**, Use of Force Fields in Materials Modeling.

**M. Rami Reddy**, **Mark D. Erion**, and **Atul Agarwal**, Free Energy Calculations: Use and Limitations in Predicting Ligand Binding Affinities.

## **Volume 17 (2001)**

**Ingo Muegge** and **Matthias Rarey**, Small Molecule Docking and Scoring.

**Lutz P. Ehrlich** and **Rebecca C. Wade**, Protein-Protein Docking.

**Christel M. Marian**, Spin-Orbit Coupling in Molecules.

**Lemont B. Kier**, **Chao-Kun Cheng**, and **Paul G. Seybold**, Cellular Automata Models of Aqueous Solution Systems.

**Kenny B. Lipkowitz** and **Donald B. Boyd**, Appendix: Books Published on the Topics of Computational Chemistry.

## **Volume 18    (2002)**

**Geoff M. Downs** and **John M. Barnard**, Clustering Methods and Their Uses in Computational Chemistry.

**Hans-Joachim Böhm** and **Martin Stahl**, The Use of Scoring Functions in Drug Discovery Applications.

**Steven W. Rick** and **Steven J. Stuart**, Potentials and Algorithms for Incorporating Polarizability in Computer Simulations.

**Dmitry V. Matyushov** and **Gregory A. Voth**, New Developments in the Theoretical Description of Charge-Transfer Reactions in Condensed Phases.

**George R. Famini** and **Leland Y. Wilson**, Linear Free Energy Relationships Using Quantum Mechanical Descriptors.

**Sigrid D. Peyerimhoff**, The Development of Computational Chemistry in Germany.

**Donald B. Boyd** and **Kenny B. Lipkowitz**, Appendix: Examination of the Employment Environment for Computational Chemistry.

## **Volume 19    (2003)**

**Robert Q. Topper**, **David L. Freeman**, **Denise Bergin** and **Keirnan R. LaMarche**, Computational Techniques and Strategies for Monte Carlo Thermodynamic Calculations, with Applications to Nanoclusters.

**David E. Smith** and **Anthony D. J. Haymet**, Computing Hydrophobicity.

**Lipeng Sun** and **William L. Hase**, Born-Oppenheimer Direct Dynamics Classical Trajectory Simulations.

**Gene Lamm**, The Poisson-Boltzmann Equation.

## **Volume 20    (2004)**

**Sason Shaik** and **Philippe C. Hibert**, Valence Bond Theory: Its History, Fundamentals and Applications. A Primer.

**Nikita Matsunaga** and **Shiro Koseki**, Modeling of Spin Forbidden Reactions.

**Stefan Grimme**, Calculation of the Electronic Spectra of Large Molecules.

**Raymond Kapral**, Simulating Chemical Waves and Patterns.

**Costel Sârbu** and **Horia Pop**, Fuzzy Soft-Computing Methods and Their Applications in Chemistry.

**Sean Ekins** and **Peter Swaan**, Development of Computational Models for Enzymes, Transporters, Channels and Receptors Relevant to ADME/Tox.

# Ab Initio Quantum Simulation in Solid State Chemistry

Roberto Dovesi,<sup>\*</sup> Bartolomeo Civalleri,<sup>\*</sup>  
Roberto Orlando,<sup>#</sup> Carla Roetti,<sup>\*</sup> and  
Victor R. Saunders<sup>\*</sup>

<sup>\*</sup>*Dipartimento di Chimica IFM, Università di Torino, Via Giuria 5, I-10125 Torino, Italy*

<sup>#</sup>*Dipartimento di Scienze e Tecnologie Avanzate, Università del Piemonte Orientale, Corso Borsalino 54, I-15100 Alessandria, Italy*

---

## INTRODUCTION

Molecular quantum chemistry and quantum mechanical simulation of solids have followed substantially independent paths and strategies for many years, with almost no reciprocal influence. In the implementation of computational schemes and formalisms, they started from different elementary models: either the hydrogen or helium atom like, for example, the parameterization of a correlation functional based on accurate He atom calculations by Colle and Salvetti,<sup>1</sup> or the electron gas, which is the reference system of the local density approximation<sup>2-7</sup> (LDA) to density functional theory (DFT). Moreover, if we compare the simplest *real* crystals, like lithium metal or sodium chloride, with the smallest molecule, H<sub>2</sub>, the much greater complexity of the solid system is

sufficient to explain the long delay of about 20 years, or more, in the development of ab initio simulation strategies in the two directions.

Molecular quantum chemistry evolved to maturity in many respects in the early 1970s, where the ab initio calculation of the molecular *total energy* became the key to understanding the chemical behavior of molecules, within a well-established, proper methodology, including the use of Gaussian functions as a basis set,<sup>8,9</sup> sophisticated approximations to the wave function,<sup>10–14</sup> and analytical gradients for geometry optimization.<sup>15,16</sup> Computer programs, such as IBMOL<sup>17</sup> and GAUSSIAN70,<sup>18</sup> were already available to the scientific community, at that time or even before, and many molecular properties could be predicted with excellent accuracy, although at the beginning severe limitations, involving the algorithm efficiency and scaling with the system size, restricted the applicability of ab initio quantum-chemical methods to small molecules with relatively poor basis sets.

The approach to solving problems in the solid state was completely different and coincided essentially with developments in solid state physics, which at the time focused on comprehending fundamental properties such as the band structure, the effective mass, the Fermi surface shape, and their relationship to the electrical behavior of materials or to the interpretation of excitation spectra. The popular textbooks by Bassani<sup>19</sup> and Moruzzi et al.<sup>20</sup> document well the state-of-the-art in solid state simulation during that period. Computer programs were mainly based on semi-empirical methods using “muffin-tin” potentials and the analytical simplicity of plane-wave (PW) basis sets. In late 1970s, ab initio pseudopotentials (PP), determined with reference to atomic calculations with the same Hamiltonian (see, for example, references 21–24), replaced previous empirical and semi-empirical PPs. Regarding the Hamiltonian, in the same years, the popular  $X\text{-}\alpha$ <sup>25</sup> method was replaced by the parameter-free LDA.<sup>3–5,7</sup> As a matter of fact, the combination PP-PW-LDA became, and remains, the most popular “recipe” for the calculation of the electronic structure of crystalline compounds, although other schemes were also largely adopted, such as Korringa–Kohn–Rostoker (KKR),<sup>26,27</sup> orthogonalized plane waves (OPW),<sup>28</sup> augmented plane waves (APW),<sup>29,30</sup> linearized augmented plane waves (LAPW),<sup>31</sup> spherical cellular schemes,<sup>32</sup> and diophantine integration schemes.<sup>33</sup>

Conversely, structural and elastic properties of ionic and semi-ionic solids were studied successfully in a completely different context with semi-classical methods,<sup>34</sup> based on force-field model potentials.

Preliminary attempts at introducing the quantum-chemical viewpoint into solid state modeling date to the late 1960s through the generalization of the Hartree–Fock (HF) equations for crystalline systems with a local basis set.<sup>35–41</sup> These were, however, in most cases, only formal equations or partial solutions to some of the many computational problems implicit in these equations. Only at the beginning of the 1970s were Ewema<sup>42,43</sup> and collaborators able to run a fully ab initio all-electron calculation for a crystalline compound



in a local basis, with reasonably good results for binding energies and lattice parameters of diamond,<sup>42</sup> boron nitride,<sup>43</sup> and a few other systems. Unfortunately, this research project was abandoned and the related experience was lost for a while. Reliable *ab initio* algorithms, capable of computing not only the band structure but also relatively accurate binding energies, equilibrium geometries, and elastic properties were implemented shortly before 1980. Most of them were based on LDA with PW combined with PP. CRYSTAL<sup>44</sup> was the only periodic *ab initio* all-electron program based on the HF Hamiltonian and the use of Gaussian functions at that time.

In all cases, the access to programs for solid state simulation was exclusive to the research groups developing them. For this reason and because of the differences in the computational programs and their implementation, a comparison of the methods by performance was difficult. CRYSTAL was the first periodic *ab initio* code to be distributed to the scientific community beginning in 1989.<sup>45</sup> Afterward, the evolution in the field was rapid, and now several *ab initio* codes are available to users (see Appendix 1 for a list of some of these codes, with short indication of their main features).

Nowadays, simulation of infinite systems relies on an ensemble of strategies and methods differing in many respects. By simply looking at the list of solid state programs reported in Appendix 1, one has an idea of the large variety of approaches available. Illustrating the features of the various codes, or their merits and limits, is not the aim of this chapter. Instead, we provide here only a brief summary of the main “ingredients” in the “recipe” of a code for solid state simulation that includes:

1. The model. Many different models can be proposed for the simulation of a single physical or chemical phenomenon. For example, a point defect in a crystalline system can be simulated either with a finite cluster with a defect at the center of the cluster and by assuming that the cluster is big enough and border effects are small, or with a periodic supercell approach, with the defects repeated periodically in such a way that the defect–defect interaction is small, if the supercell is big enough.
2. The Hamiltonian. Although most of the periodic calculations are performed with reference to DFT, the debate is still open about the most appropriate functional to use for different systems and properties, ranging from LDA (that is still popular in solid state physics) to various generalized gradient approximation (GGA) formulations and hybrid schemes like B3LYP. In a few cases, HF is still preferred.
3. The basis set. Codes based on plane waves, local functions, and mixed (local functions in atomic spheres, plane waves in the interstices) or numeric basis sets are available.
4. The overall computational scheme, in all its features, such as direct or reciprocal space representation, all-electron versus pseudo-potential

formulation, and analytical versus numerical calculation of matrix elements and relevant integrals.

A reader is probably interested in finding answers to the following questions: What additional basic information is needed for proper use of periodic codes by a scientist with a molecular quantum chemistry background? Are there features peculiar to the solid state, with no analogy to the gas phase? In this chapter, we shall provide answers to these questions as well as provide a tutorial for the nonspecialist wanting to learn about solid state calculations.

The solid under study with a periodic program is infinite and translation invariant; it is a perfect crystal. Despite that no real crystal is a perfect crystal, this model is suitable in most cases, and indeed, experimental evidence of crystal periodicity exists in x-ray, neutron, and electron diffraction patterns, which are hardly affected by the presence of the surface, unless the experiment is done in special conditions. Translation invariance has a series of interesting properties with important consequences on simplification of the problem and the implementation of efficient algorithms.

Even in those cases where the model of a perfect crystal appears as inappropriate does one try to simulate partially nonperiodic systems with some nearly equivalent, formally periodic structure, whenever possible, as happens in the descriptions of local defects with the supercell approximation (see the section on defects) or in the treatment of substitutionally disordered systems.<sup>46,47</sup> Therefore, the use of a periodic program by a scientist requires basic knowledge of crystallography, such as the definitions of lattice, direct and reciprocal space, unit cell, Brillouin zone, and the main concepts of the solid state language. These ideas are described briefly in the next section. Other more specific points will be mentioned with almost no discussion, because comprehension of their details is beyond the scope of this chapter. For example, the evaluation of electrostatic interactions in a solid<sup>48-51</sup> is more complicated than would appear at first sight and it represents one of the more crucial aspects of the computational problem. The formulation of a convenient method to compute the electrostatic potential generated by a three-dimensional array of charge distributions<sup>41,52-55</sup> required more than 50 years' work; coverage of this topic is thus ill-advised.

Apart from the methodological aspects, solid state systems possess many interesting properties that are immaterial for single molecules. In single molecules, point symmetry usually decreases as the size of the molecule increases. Molecules with more than, say, 20 atoms often lack symmetry. Crystalline systems, contrarily, usually maintain high point symmetry, even in the case of large unit cell systems, like zeolites and garnets.

Another important difference between nonamorphous solids and single molecules is anisotropy (different space directions are not equivalent). No anisotropic effect is observed in the gas phase with no applied field because of the averaging process caused by the random orientation of molecules. In contrast,

crystals are macroscopic objects that can be oriented with respect to a reference frame and their properties generally depend on orientation. All properties related to crystal anisotropy are then described by tensors of various rank. For example, the relationship between stress and strain (second-order tensors each) cannot be expressed with a single constant, but it can be expressed as a fourth-order tensor, whereas piezoelectricity is described by a third-order tensor. In most cases, the physical and technological interest in materials science is focused on the possibility of increasing or reducing anisotropy in materials.

Throughout this chapter, we will illustrate some of the possibilities offered by *ab initio* simulation in the area of solid state chemistry, physics, materials science, surface science, and catalysis. The examples are mainly focused on simple properties like energy and its derivatives, band structure, and charge density, to give the reader who is not acquainted with solid state simulation an introductory overview. For consistency of the data and their representations, all examples have been generated with the CRYSTAL code,<sup>56</sup> implemented by the present authors and collaborators. All cases reported here refer to the static limit, and temperature effects are not discussed. Temperature effects can be taken into account by calculating thermodynamic functions from the vibration spectrum following a methodology common to most molecular codes (see, for example, reference 57 for a recent review). Alternatively, in solid state physics, the Car–Parrinello<sup>58</sup> methodology is popular, because it is an efficient way of finding equilibrium electron and nuclear coordinates at once.

Many important and interesting systems and properties could not be considered in this presentation either for conciseness, as they would require some preliminary long explanation, or because they are not yet available in CRYSTAL, although they have been implemented in other codes. The calculation of NMR tensors,<sup>59</sup> Raman intensity tensors,<sup>60</sup> and electro-optic tensors<sup>61</sup> are only a few examples from this long list of omissions.

Two relevant topics have been ignored completely in this short chapter: the treatment of electron correlation with more sophisticated methods than DFT (that remains unsatisfactory from many points of view) and the related subject of excited states. Wave function-based methods for the calculation of electron correlation, like the perturbative Møller–Plesset (MP) expansion or the coupled cluster approximation, have registered an impressive advancement in the molecular context. The computational cost increases with the molecular size (as the fifth power in the most favorable cases), especially for molecules with low symmetry. That increase was the main disadvantage of these electron correlation methods, and it limited their application to tiny molecules. This scaling problem has been improved dramatically by modern reformulation of the theory by localized molecular orbitals, and now a much more favorable scaling is possible with the appropriate approximations. Linear scaling with such low prefactors has been achieved<sup>62,63</sup> with MP schemes that the

feasibility of this kind of calculation has been extended to molecules of medium and large size.

In principle, this electron correlation strategy is transferable from single molecules to solids, after the crystalline orbitals have been transformed to an equivalent set of well-localized functions (Wannier functions). Procedures for orbital localization have been proposed and implemented only recently,<sup>64,65</sup> and the first MP2 calculations are becoming possible in the case of simple crystalline compounds.<sup>66</sup>

Alternative strategies have also been proposed for estimating correlation energies, including quantum Monte Carlo methods (see reference 67 and references therein), MP2 schemes, either canonical<sup>68,69</sup> or based on the Laplace transform algorithm,<sup>70</sup> and the molecular-like incremental method applied by Stoll.<sup>71–75</sup> However, none of these methods seems to have arrived at a sufficiently advanced stage of development to be of general use to the scientist at the moment.

Regarding excited states, time-dependent density functional theory<sup>76–78</sup> (TDDFT) is considered a relatively accurate method (see, for example, reference 79) for the study of the low-lying excited states, with results by far superior to the simple virtual-occupied DFT energy difference. The most recent formulations of the GW formalism, originally proposed about 20 years ago,<sup>80</sup> seem to provide good band gaps, optical spectra, and electron-hole excitations<sup>81–84</sup> (the GW acronym arises from the name of the two matrices the method is based on: the Green function matrix, G, and the screened Coulomb matrix, W).

---

## TRANSLATION INVARIANCE PROPERTIES IN A CRYSTAL

Because a crystal can be regarded as a huge molecule consisting of about as many as Avogadro's number of atoms or ions, calculation of the crystalline electronic structure and properties may appear as an unattainable problem. Fortunately, however, crystals exhibit a very important symmetry property: They are translation invariant by definition. In fact, a perfect crystal consists of a three-dimensional array of atoms, ions, or molecules, a few of which form a spatial pattern that is repeated identically throughout the crystal. Clever exploitation of this symmetry property makes the computational problem solvable, and the theory, on which the solution is based, is known as band theory. The application of band theory to the study of periodic systems requires the knowledge of a specific language and some understanding of the properties of translation symmetry. In this section, a few basic concepts of crystallography and band theory will be introduced with reference to some elementary definitions as well as the discussion of a few simple examples, which are aimed to show how band structure and properties originate and to provide a little



Figure 1 Ammonia molecule.

insight into the methods applied in the calculation of the electronic structure of matter in the condensed phase.

What is the main difference between studying the electronic structure of a molecule in the gas phase and in condensed phase? In the gas phase, because of the low density and large kinetic energies, molecules interact only during collisions, which may promote them to an excited state. However, either before or after such a brief collision, a molecule is essentially not influenced by the other molecules. Thus, as far as we are not interested in molecular dynamics and thermodynamic properties, the electronic structure in the ground or in any excited state can be studied for only one isolated molecule.

For example, if we are interested in studying ammonia in the gas phase, we can consider only one ammonia molecule, like in Figure 1. The positions of the nitrogen and hydrogen atoms, defined in the Cartesian coordinates or through a set of internal coordinates, are the only information necessary to compute the molecular wave function by *ab initio* methods.

On the contrary, at very high pressure or low temperature, ammonia molecules interact with each other and pack together to form a crystalline phase, known as phase I,<sup>85</sup> where the number of molecules involved is indeed large. Even the definition of the composition and geometry of a crystal is not as simple as for molecules. However, the arrangement of the molecules in a crystal must satisfy the condition of maximizing intermolecular attractive interactions, which imposes some severe constraints on their mutual orientations. In fact, observing an ammonia crystal carefully (Figure 2), it is possible to identify a set of four molecules that, when translated along each side of the cube by a multiple of the entire length, overlap with another set of four identical molecules exactly, because of translation invariance. Crystallography provides a mathematical description of this kind of object along with the tools for managing such complex systems.

## The Direct Lattice

The crystallographer's view of a crystal<sup>86</sup> starts from the definition of a lattice: A *lattice* is a collection of points repeated at intervals of length  $a_1$ ,  $a_2$  and  $a_3$  along three non-coplanar directions, indefinitely. The three constants  $a_1$ ,  $a_2$ , and  $a_3$  are called *lattice parameters*, and the vectors  $\mathbf{a}_1$ ,  $\mathbf{a}_2$ , and  $\mathbf{a}_3$ , oriented in the same three non-coplanar directions with the lattice parameters

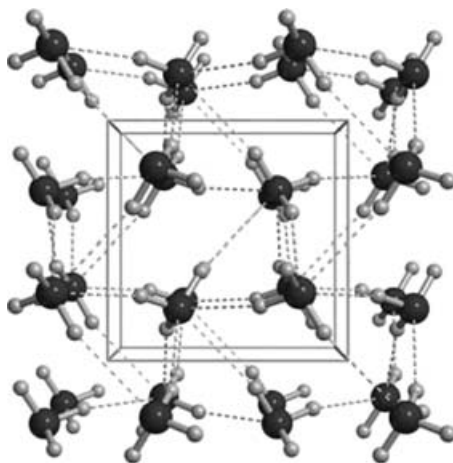


Figure 2 Crystal packing for phase I of solid ammonia.

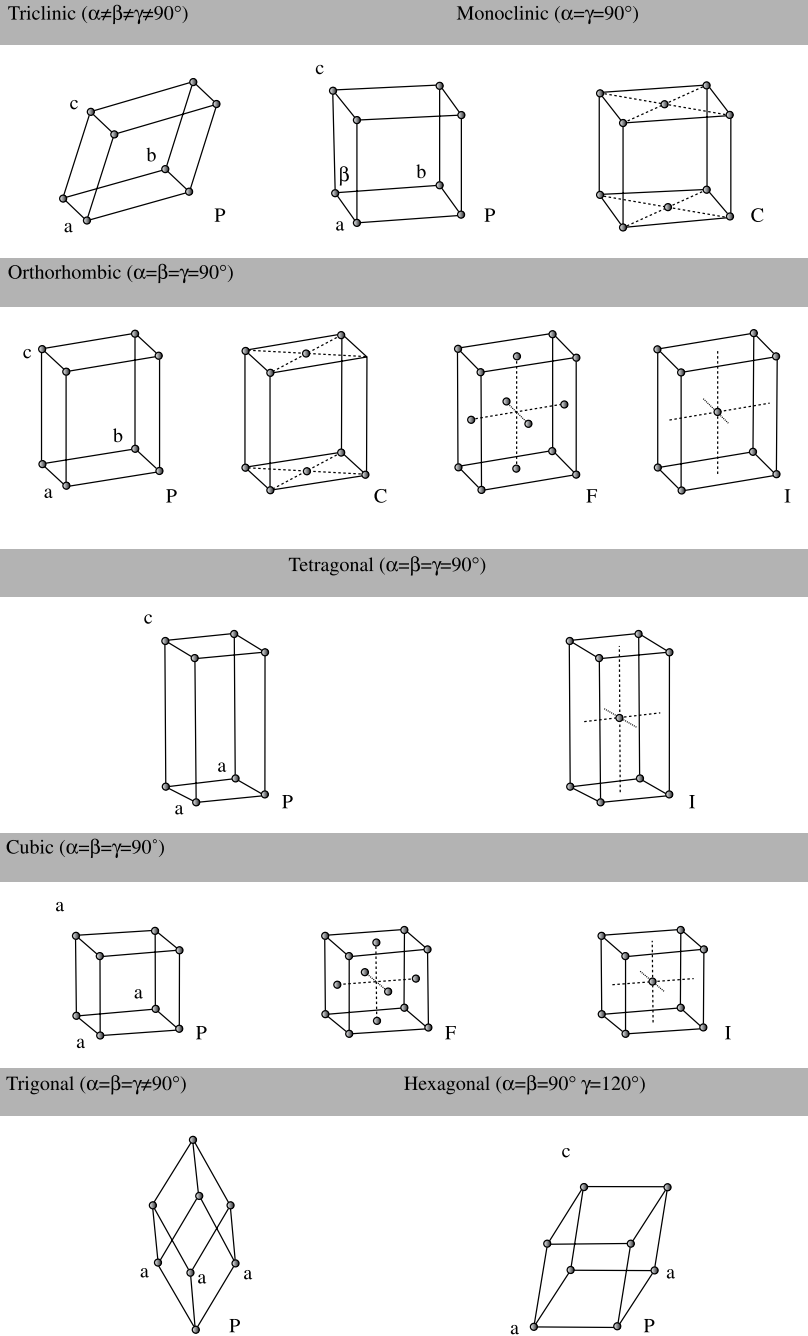
as norms, are the *basis vectors*. Lattice parameters and angles between the lattice vectors are collectively called *cell parameters*.

A vector  $\mathbf{g}$  joining any two lattice points is a *lattice vector*. Every lattice vector can always be expressed by the basis vectors and three integer coefficients  $n_1$ ,  $n_2$ , and  $n_3$ :

$$\mathbf{g} = n_1\mathbf{a}_1 + n_2\mathbf{a}_2 + n_3\mathbf{a}_3 \quad [1]$$

Basis vectors  $\mathbf{a}_1$ ,  $\mathbf{a}_2$ , and  $\mathbf{a}_3$  define a parallelepiped called the *unit cell*, which is *primitive*, because it contains one lattice point. All cells that are obtained by translation of this unit cell, the origin cell, through the application of all vectors  $\mathbf{g}$  in Eq. [1], fill the space completely. Then, the entire lattice can be subdivided into cells and every vector  $\mathbf{g}$  can be used to label a cell with respect to the origin cell, or 0-cell. Actually, the definition of a unit cell is arbitrary, and many (an infinite number) different possible choices exist, because all cells containing the same number of lattice points are equivalent. The actual shape of a unit cell depends on the lattice type.

Primitive three-dimensional lattices have been classified into seven crystalline systems: *triclinic*, *monoclinic*, *orthorhombic*, *tetragonal*, *cubic*, *trigonal*, and *hexagonal*. They are different in the relative lengths of the basis vectors as well as in the angles they form. An additional seven nonprimitive lattices, belonging to the same crystalline systems, are added to the seven primitive lattices, which thus completes the set of all conceivable lattices in ordinary space. These 14 different types of lattices are known as Bravais lattices (Figure 3).



**Figure 3** Bravais lattices. Symbols *P*, *F*, and *I* denote primitive, face-centered, and body-centered lattices, respectively.  $\alpha$ ,  $\beta$ , and  $\gamma$  are the angles between the **b** and **c**, **a** and **c**, and **a** and **b** basis vectors.

Filling the unit cell of a lattice with matter in a well-defined geometrical arrangement and applying the translation pattern permits the creation of an ideal crystal. Crystals usually exhibit point symmetry in addition to the set of translations. Point and translation symmetries combine to form a *space group*. It has been demonstrated that in ordinary space, only 230 of these different possible combinations exist and every space group refers to only one particular Bravais lattice. Space groups are fully characterized in the International Tables of Crystallography.<sup>87</sup> Every group is identified by a symbol (Hermann–Mauguin) that specifies it completely (the Schönflies notation, which is more frequently used for molecules, is also available, but it is less adequate to describe translational properties). Some points in the cell are invariant to one or more symmetry operations of the space group. In those cases, the number of symmetry equivalent points in a cell, or *multiplicity*, is smaller than the total number of symmetry operations. Such a point is called a *special position*, whereas every other point is referred to as a *general position*. The minimal set of atoms, either in special or general positions, which generates the complete unit cell after application of all space group symmetry operations, is referred to as the *asymmetric unit*.

In summary, specifying the geometry of a crystal requires the following information:

- Space group
- Cell parameters
- Type and position of the atoms in the asymmetric unit

The position  $\mathbf{r}$  of an atom in the unit cell is usually not expressed in terms of Cartesian coordinates, but in terms of *fractional coordinates*  $x_1, x_2, x_3$  such that

$$\mathbf{r} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + x_3\mathbf{a}_3 \quad [2]$$

$x_1, x_2$ , and  $x_3$  are pure numbers, as coefficients of the lattice basis vectors. Atoms with fractional coordinates all in the range between 0 and 1 belong to the 0-cell.

The parallelepiped in Figure 2 is the unit cell of the ammonia crystal phase I. Thus, the ammonia crystal can be regarded as the combination of a pattern of four ammonia molecules (16 atoms) in the unit cell with all possible translations in a cubic primitive lattice. Considerations about crystalline symmetry lead to the conclusion that ammonia in phase I crystallizes according to space group  $P2_13$ . Letter  $P$  in the symbol stands for primitive lattice, and the other symbols denote the main symmetry operations. The last element in the symbol, 3, indicates the presence of a three-fold axis not aligned with the principal rotation axis (if it was, it would follow letter  $P$ ), which further indicates that the lattice is cubic. A cubic unit cell is completely specified by just one



lattice parameter, with the basis vectors all at right angles having the same norm. Because of the symmetry, only two atoms are in the asymmetric unit: one N and one H atom, so that only six fractional coordinates need to be specified. N is in a special position (along a three-fold axis) with multiplicity 4, whereas H is in a general position and has multiplicity 12, i.e., the number of point symmetry operations in the group.

When modeling a polymer or a surface, translation invariance is restricted to only one or two independent directions, instead of three. Space groups cannot characterize the symmetry of one-dimensional and two-dimensional periodic systems, and we need to refer to special subgroups of the space groups, the 75 *rod groups* and the 80 *layer groups*, respectively, which include the symmetry of all possible arrangements of three-dimensional objects (molecules or sets of atoms) in one-dimensional and two-dimensional lattices. Most of the considerations about the space groups are still valid for the rod and layer groups, with the exception of the classification of lattices, which is intimately related to the type of periodicity.

From here on, we will refer to ordinary space as the *direct space*, in order to contrast it to the *reciprocal space*, which is introduced in the next paragraph.

## The Reciprocal Lattice

Every direct lattice admits a geometric construction, the *reciprocal lattice*, by the prescription that the reciprocal lattice basis vectors ( $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ ) obey the following important orthogonality rules relative to the direct lattice basis vectors ( $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ ):

$$\mathbf{a}_i \cdot \mathbf{b}_j = 2\pi \delta_{ij} \quad [3]$$

which implies that every reciprocal lattice basis vector (normalization to  $2\pi$ ) is orthogonal to the plane of the corresponding direct lattice basis vectors with unequal indices ( $\delta_{ij}$  is 1 if  $i$  equals  $j$  and 0 when  $i$  is different from  $j$ ).

Like in direct space, any reciprocal lattice vector can be expressed as a linear combination of the basis vectors with integer coefficients such as

$$\mathbf{K} = K_1\mathbf{b}_1 + K_2\mathbf{b}_2 + K_3\mathbf{b}_3 \quad [4]$$

Among all possible equivalent choices of a unit cell in the reciprocal lattice, one is particularly useful. It can be obtained by connecting one reciprocal lattice point to all its nearest neighbors and letting orthogonal planes pass through their midpoints. The volume within these planes is known as the *first Brillouin zone*. It includes all points that are closer to that reciprocal lattice point than to any other lattice point.

## Bloch Theorem and the Periodic Boundary Conditions

A real crystal is a finite macroscopic object made of a finite, although extremely large, number of atoms. However, the ratio of the number of atoms at the surface to the total number of atoms in the crystal,  $N$ , is very small, and proportional to  $N^{-1/3}$ . When  $N$  is large and the surface is neutral, the perturbation caused by the presence of the boundary is limited to only a few surface layers and, therefore, has no influence on the bulk properties. For this reason, a macroscopic crystal mostly exhibits properties and features of the bulk material, and unless attention is deliberately focused onto the crystal boundary, surface effects can be thoroughly neglected. If this is the case, the crystallographic model of an infinite and translation-invariant crystal fits in the aim of studying bulk properties.

The potential energy of such a crystal must be a periodic function with the same periodicity as the lattice, so that for a translation by any direct lattice vector  $\mathbf{g}$ , the potential energy does not change

$$V(\mathbf{r} - \mathbf{g}) = V(\mathbf{r}) \quad [5]$$

Because of symmetry requirements, the Schrödinger equation

$$\hat{H}(\mathbf{r})\Psi(\mathbf{r}) = E\Psi(\mathbf{r}) \quad [6]$$

must also be translation invariant, which is equivalent to the requirement that, after a translation of the entire crystal by  $\mathbf{g}$ , the solutions of equation

$$\hat{H}(\mathbf{r} - \mathbf{g})\Psi(\mathbf{r} - \mathbf{g}) = E\Psi(\mathbf{r} - \mathbf{g}) \quad [7]$$

coincide with those of Eq. [6]. It has been demonstrated<sup>88</sup> that eigenfunctions with the correct symmetry relative to a potential of the form of Eq. [5] must obey the Bloch theorem, stating that

$$\Phi(\mathbf{r} + \mathbf{g}; \mathbf{k}) = e^{i\mathbf{k}\cdot\mathbf{g}} \Phi(\mathbf{r}; \mathbf{k}) \quad [8]$$

and providing a relation between the values of an eigenfunction at equivalent points in the lattice, which indicates that its periodicity is generally different from that of the lattice. As  $\Phi$  verifies the Bloch theorem, it is called the **Bloch function** and is a function of the position in space  $\mathbf{r}$  and the **wave vector**  $\mathbf{k}$ . Parameter  $\mathbf{k}$  labels the different solutions to Eq. [6].

The equivalence of Eq. [6] and Eq. [7] can be verified easily, by supposing that the Bloch function  $\Phi(\mathbf{r}; \mathbf{k})$  be an eigenfunction of the Hamiltonian in Eq. [6]. In this case, Eq. [7] can be rewritten as

$$\hat{H}(\mathbf{r} - \mathbf{g})\Phi(\mathbf{r} - \mathbf{g}; \mathbf{k}) = E(\mathbf{k})\Phi(\mathbf{r} - \mathbf{g}; \mathbf{k}) \quad [9]$$

However, by using the Bloch theorem and considering that  $\hat{H}(\mathbf{r} - \mathbf{g})$  is equivalent to  $\hat{H}(\mathbf{r})$  as an obvious consequence of the form of the crystalline potential (Eq. [5]), we obtain Eq. [6] again

$$\hat{H}(\mathbf{r})e^{-i\mathbf{k}\cdot\mathbf{g}}\Phi(\mathbf{r}; \mathbf{k}) = E(\mathbf{k})e^{-i\mathbf{k}\cdot\mathbf{g}}\Phi(\mathbf{r}; \mathbf{k}) \quad [10]$$

where  $e^{-i\mathbf{k}\cdot\mathbf{g}}$  is simply a constant factor with unitary module (the eigenfunctions of an operator can always be multiplied by any arbitrary constant factor).

What is the form of Bloch functions? Equation [8] implies that a Bloch function can be written as the product of a plane wave and a periodic function  $u(\mathbf{r}; \mathbf{k})$  with the same periodicity of the lattice:

$$\Phi(\mathbf{r}; \mathbf{k}) = e^{i\mathbf{k}\cdot\mathbf{r}} u(\mathbf{r}; \mathbf{k}) \quad [11]$$

In fact, Eq. [8] is immediately verified in this case:

$$\Phi(\mathbf{r} + \mathbf{g}; \mathbf{k}) = e^{i\mathbf{k}\cdot(\mathbf{r}+\mathbf{g})} u(\mathbf{r} + \mathbf{g}; \mathbf{k}) = e^{i\mathbf{k}\cdot\mathbf{g}} e^{i\mathbf{k}\cdot\mathbf{r}} u(\mathbf{r}; \mathbf{k}) = e^{i\mathbf{k}\cdot\mathbf{g}} \Phi(\mathbf{r}; \mathbf{k}) \quad [12]$$

Bloch functions span an infinite crystal and do not decay to zero at infinity. To circumvent the problem of normalizing a wave function with infinite extent, it is easier to consider a finite crystal consisting of  $N = N_1 \times N_2 \times N_3$  cells and then let  $N$  grow to infinity. To preserve periodicity, *periodic boundary conditions* are imposed, which can be stated in the following form: If  $N_j$  cells exist along the  $j$ -th direction ( $j = 1, 2, 3$ ) in the macroscopic crystal, it must happen that for any integer  $m$  and every  $j$

$$\Phi(\mathbf{r} + m N_j \mathbf{a}_j; \mathbf{k}) = \Phi(\mathbf{r}; \mathbf{k}) \quad [13]$$

as if the crystal was a three-dimensional infinite array of identical and contiguous finite crystals with the shape of a parallelepiped, each consisting of  $N$  primitive cells. However, according to the Bloch theorem (Eq. [8])

$$\Phi(\mathbf{r} + m N_j \mathbf{a}_j; \mathbf{k}) = e^{im N_j \mathbf{k} \cdot \mathbf{a}_j} \Phi(\mathbf{r}; \mathbf{k}) \quad [14]$$

and, on comparing these two expressions, it is evident that the phase factor must be equal to one

$$e^{im N_j \mathbf{k} \cdot \mathbf{a}_j} = e^{im N_j \mathbf{k}_j \cdot \mathbf{a}_j} = 1 \quad [15]$$

But, if the component  $\mathbf{k}_j$  of the wave vector is defined as

$$\mathbf{k}_j = \frac{n_j}{N_j} \mathbf{b}_j \quad [16]$$

with  $n_j$  being an integer, by Eq. [3],  $\mathbf{k}$  can be interpreted as a point in the reciprocal lattice. Therefore,  $N$  of  $\mathbf{k}$  points exist in every reciprocal lattice cell, each of which can be written by the reciprocal lattice basis vectors as

$$\mathbf{k} = \left( \frac{n_1}{N_1} \mathbf{b}_1 + \frac{n_2}{N_2} \mathbf{b}_2 + \frac{n_3}{N_3} \mathbf{b}_3 \right) \quad [17]$$

If  $n_j$  is such that  $\phi \leq n_j < N_j$  for every  $j$   $\mathbf{k}$  belongs to the origin cell of the reciprocal lattice.

When  $N_1, N_2, N_3$  are allowed to approach infinity, the number of  $\mathbf{k}$  points in every reciprocal lattice cell also tends to infinity until they completely fill the space, so that  $\mathbf{k}$  can be considered as a continuous variable.

Bloch functions also have interesting translational properties in the reciprocal space, which can be investigated by considering a new point  $\mathbf{h} = \mathbf{k} + \mathbf{K}$ , obtained by a translation of the wave vector  $\mathbf{k}$  by any reciprocal lattice vector  $\mathbf{K}$  (Eq. [4]), and then applying Eq. [8] to the corresponding Bloch function  $\Phi(\mathbf{r}; \mathbf{h})$ . By comparison with Eq. [8] and with Eq. [3], it is evident that  $\Phi(\mathbf{r}; \mathbf{h})$  exhibits the same translational properties as  $\Phi(\mathbf{r}; \mathbf{k})$

$$\Phi(\mathbf{r} + \mathbf{g}; \mathbf{h}) = e^{i(\mathbf{k} + \mathbf{K}) \cdot \mathbf{g}} \Phi(\mathbf{r}; \mathbf{k} + \mathbf{K}) = e^{i\mathbf{k} \cdot \mathbf{g}} \Phi(\mathbf{r}; \mathbf{h}) \quad [18]$$

so that both  $\Phi(\mathbf{r}; \mathbf{h})$  and  $\Phi(\mathbf{r}; \mathbf{k})$  can be referred to the same  $\mathbf{k}$  and are acceptable eigenfunctions for that  $\mathbf{k}$  in Eq. [6]. This behavior of Bloch functions in reciprocal space has the important consequence that the analysis can be restricted to the first Brillouin zone.

Another very important property of Bloch functions is related to the evaluation of the following integral extended to the entire space, which involves a function  $f(\mathbf{r})$  with the same periodicity of the lattice:

$$f(\mathbf{k}, \mathbf{k}') = \int [\Phi(\mathbf{r}; \mathbf{k}')]^* f(\mathbf{r}) \Phi(\mathbf{r}; \mathbf{k}) d\mathbf{r} \quad [19]$$

with  $\mathbf{k}$  and  $\mathbf{k}'$  being points in the first Brillouin zone. In accordance with the property of Bloch functions in the reciprocal space just shown, the periodic component of a Bloch function referred to as  $\mathbf{k}$  (Eq. [11]) can be expanded into a linear combination of those plane waves for which the wave vector is obtained by addition of all reciprocal lattice vectors to the corresponding  $\mathbf{k}$ :

$$\begin{aligned} u(\mathbf{r}; \mathbf{k}) &= \sum_{\mathbf{K}} c_{\mathbf{K}} e^{i(\mathbf{k} + \mathbf{K}) \cdot \mathbf{r}} \\ u(\mathbf{r}; \mathbf{k}') &= \sum_{\mathbf{K}'} c_{\mathbf{K}'} e^{i(\mathbf{k}' + \mathbf{K}') \cdot \mathbf{r}} \end{aligned} \quad [20]$$

where  $c_{\mathbf{K}}$  and  $c_{\mathbf{K}'}$  are the expansion coefficients. Similarly, also  $f(\mathbf{r})$ , having the same periodicity as  $u(\mathbf{r}; \mathbf{k})$  and  $u(\mathbf{r}; \mathbf{k}')$ , can be expanded in terms of plane waves, with  $\mathbf{k} = 0$

$$f(\mathbf{r}) = \sum_{\mathbf{K}''} d_{\mathbf{K}''} e^{i\mathbf{K}'' \cdot \mathbf{r}} \quad [21]$$

Now the integral can be calculated through the evaluation of the following three infinite sums of integrals involving plane waves only:

$$f(\mathbf{k}, \mathbf{k}') = \sum_{\mathbf{K}} c_{\mathbf{K}} \sum_{\mathbf{K}'} c_{\mathbf{K}'} \sum_{\mathbf{K}''} d_{\mathbf{K}''} \int e^{i(\mathbf{k} + \mathbf{K} + \mathbf{K}'') \cdot \mathbf{r}} e^{i(\mathbf{k}' + \mathbf{K}') \cdot \mathbf{r}} d\mathbf{r} \quad [22]$$

and for the orthogonality of plane waves, these terms are zero unless  $\mathbf{k} + \mathbf{K} + \mathbf{K}'' = \mathbf{k}' + \mathbf{K}'$ . By Eq. [4] and Eq. [17], it is clear that this condition is fulfilled only if  $\mathbf{k} = \mathbf{k}'$ .

All integrals that need to be calculated in Eq. [6] are of this kind, as the potential energy term is a periodic function of the lattice, like  $f(\mathbf{r})$ , and the kinetic energy term involves second derivatives of  $u(\mathbf{r}; \mathbf{k})$  with respect to  $\mathbf{r}$ , which have the same periodicity as  $u(\mathbf{r}; \mathbf{k})$ .

We could also arrive at the same conclusion in a different way by observing that Bloch functions are the eigenfunctions of translation operators and of all operators commuting with translation operators, like the Hamiltonian for a periodic system. Then, Bloch functions are bases for the irreducible representations for the group of the lattice translations, each one corresponding to one wave vector  $\mathbf{k}$ , and it is known from group theory that basis functions belonging to different irreducible representations are mutually orthogonal.

Therefore, great advantage exists in representing the Hamiltonian of a periodic system, where the potential energy operator has the form of Eq. [5], in Bloch functions. In fact, in this basis, the Hamiltonian matrix is block-diagonal (Figure 4), with each block referring to one particular point  $\mathbf{k}$  in the reciprocal space.

Suppose we have a finite basis set of  $n_f$  Bloch functions. The Hamiltonian matrix represented in this basis, then, consists of diagonal blocks of



**Figure 4** Transformation of the infinite Hamiltonian matrix when expressed in the basis of Bloch functions.

$n_f \times n_f$  elements, with each block referring to an individual  $\mathbf{k}$  point and being completely independent of all the others, such that the elements of a block do not interact with those of others blocks and can, therefore, be treated separately. Unfortunately, an infinite number of such factorized finite-sized blocks exists. In other words, Bloch functions as a basis set allows us to transform a problem of infinite size into an infinite number of problems of finite size. Nevertheless, what may appear as a poor advantage actually represents a great improvement, owing to the usually smooth change of the eigenvalues and the eigenvectors with  $\mathbf{k}$ . Therefore, it is generally possible to sample matrix  $\mathbf{H}$  at a finite number of points and solve the Schrödinger equation for a periodic system at different points in the first Brillouin zone:

$$\hat{H} \Psi_n(\mathbf{r}; \mathbf{k}) = E_n(\mathbf{k}) \Psi_n(\mathbf{r}; \mathbf{k}) \quad [23]$$

If sampling is convenient, the number of  $\mathbf{k}$  points to be considered is usually relatively small and solving the Schrödinger equation in the reciprocal space is a feasible method.

### One-Electron Electrostatic Hamiltonian

If  $\hat{H}$  is the one-electron electrostatic Hamiltonian, based on the Born-Oppenheimer approximation, the solutions to Eq. [23] are called crystalline orbitals (CO). They are linear combinations of one-electron Bloch functions (Eq. [8])

$$\Psi_n(\mathbf{r}; \mathbf{k}) = \sum_j c_{jn}(\mathbf{k}) \Phi_j(\mathbf{r}; \mathbf{k}) \quad [24]$$

with coefficients  $c_{jn}$  to be determined. In the basis of Bloch functions, Eq. [23] can be written in the form of a matrix equation:

$$\mathbf{H}(\mathbf{k}) \mathbf{C}(\mathbf{k}) = \mathbf{S}(\mathbf{k}) \mathbf{C}(\mathbf{k}) \mathbf{E}(\mathbf{k}) \quad [25]$$

where the size of all matrices is equal to the number of Bloch functions in the basis and  $\mathbf{S}(\mathbf{k})$  is the overlap matrix, which accounts for nonorthogonal basis sets.  $\mathbf{C}(\mathbf{k})$ , the matrix of coefficients, is constrained by the following orthonormalization condition ( $\mathbf{I}$  is the identity matrix):

$$\mathbf{C}(\mathbf{k}) \mathbf{S}(\mathbf{k}) \mathbf{C}^\dagger(\mathbf{k}) = \mathbf{I} \quad [26]$$

Atomic orbitals (AO) and plane waves are common choices to represent Bloch functions. Both choices would be equivalent, in principle, if an infinite basis set was considered, but they are not equivalent in the practical case of a finite

basis set. The use of AOs is better linked to the chemical experience of molecular codes and is particularly suitable to the description of crystals with chemical bonds. On the contrary, the description of free or nearly free electrons in conductors is hard to achieve in local functions and the addition of new functions to the AO basis set rapidly leads to saturation, which causes numerical instability because of quasi-linear dependence problems. Plane waves are more suitable to the case of metals and, in general, to the description of delocalized electrons. Another advantage of using plane waves is that the mathematics involved in the use of plane waves is usually much easier.

Any method of solution of Eq. [25] is specific of the kind of basis set used. In the remaining part of this chapter, we will always refer to the use of one-electron local basis sets within the linear combination of atomic orbitals (LCAO) method. Accordingly,  $n_f$  AOs in the 0-cell are chosen and replicated in the other cells of the crystal to form the periodic component  $u(\mathbf{r}; \mathbf{k})$  of  $n_f$  Bloch functions. In particular, by denoting the  $\mu$ -th AO, with the origin at  $\mathbf{r}_\mu$  in the 0-cell, as  $\chi_\mu(\mathbf{r} - \mathbf{r}_\mu)$  and the corresponding AO in a different cell, the  $\mathbf{g}$ -cell, as  $\chi_\mu(\mathbf{r} - \mathbf{r}_\mu - \mathbf{g})$  or, equivalently,  $\chi_\mu^{\mathbf{g}}(\mathbf{r} - \mathbf{r}_\mu)$ , the expression used for  $u_\mu(\mathbf{r}; \mathbf{k})$  consists of a linear combination of the equivalent AOs in all  $N$  cells of the crystal:

$$u_\mu(\mathbf{r}; \mathbf{k}) = \frac{1}{\sqrt{N}} e^{-i\mathbf{k} \cdot \mathbf{r}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} \chi_\mu^{\mathbf{g}}(\mathbf{r} - \mathbf{r}_\mu) \quad [27]$$

The translation invariance of  $u_\mu(\mathbf{r}; \mathbf{k})$  is obvious because the sum is extended to all cells in the crystal. In fact, if a translation by lattice vector  $\mathbf{l}$  is applied

$$\begin{aligned} u_\mu(\mathbf{r} - \mathbf{l}; \mathbf{k}) &= \frac{1}{\sqrt{N}} e^{-i\mathbf{k} \cdot \mathbf{r}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot (\mathbf{g} + \mathbf{l})} \chi_\mu^{\mathbf{g} + \mathbf{l}}(\mathbf{r} - \mathbf{r}_\mu) = \frac{1}{\sqrt{N}} e^{-i\mathbf{k} \cdot \mathbf{r}} \sum_{\mathbf{m}} e^{i\mathbf{k} \cdot \mathbf{m}} \chi_\mu^{\mathbf{m}}(\mathbf{r} - \mathbf{r}_\mu) \\ &= u_\mu(\mathbf{r}; \mathbf{k}) \end{aligned} \quad [28]$$

$u_\mu(\mathbf{r}; \mathbf{k})$  is verified to be periodic throughout the direct lattice (the equivalence of the sum over lattice vectors  $\mathbf{m} = \mathbf{g} + \mathbf{l}$  and the sum over  $\mathbf{g}$  originates from translation invariance and the periodic boundary conditions).

The corresponding Bloch function is immediately obtained from Eq. [11] after substitution of  $u_\mu(\mathbf{r}; \mathbf{k})$ :

$$\Phi_\mu(\mathbf{r}; \mathbf{k}) = \frac{1}{\sqrt{N}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} \chi_\mu^{\mathbf{g}}(\mathbf{r} - \mathbf{r}_\mu) \quad [29]$$

Apart from a few starting elementary examples, all results that will be presented in this chapter have been obtained with the approximations presently available in the CRYSTAL code.<sup>56</sup> The method<sup>89</sup> of solving the Schrödinger

equation in CRYSTAL is similar in many respects to that used in molecular codes based on the usage of Gaussian basis sets. For example, finding Hartree–Fock eigenvalues and eigenvectors for a molecule requires the following steps:

1. Forming the basis of the AOs from contractions of Gaussian functions (linear combinations of a set of functions with constant coefficients) times angular functions followed by evaluating the overlap matrix  $S$  in this basis set.
2. Evaluating Fock matrix elements ( $F_{\mu\nu}$ ) for all pairs of AOs in the local basis set, consisting of a sum of the following contributions: kinetic energy (T) terms, electron–nuclear (Z) interactions, and electron–electron Coulomb (C) and exchange (X) interactions

$$F_{\mu\nu} = T_{\mu\nu} + Z_{\mu\nu} + C_{\mu\nu} + X_{\mu\nu} \quad [30]$$

3. Solving Roothaan equations,  $FC = SCE$ , for  $E$  and  $C$  with the normalization condition

$$C^\dagger SC = I \quad [31]$$

4. Forming the density matrix from the eigenvectors of the occupied states with generic element

$$P_{\mu\nu} = \sum_n^{\text{occ.}} C_{\mu n}^* C_{\nu n} \quad [32]$$

5. Calculating the total energy according to the formula

$$E_t = N + \frac{1}{2} \sum_{\mu,\nu} P_{\mu\nu} (T_{\mu\nu} + Z_{\mu\nu} + F_{\mu\nu}) \quad [33]$$

which contains the internuclear repulsion energy  $N$  and a double sum over the AOs of one-electron terms.

Because the calculation of the electron–electron contributions to  $F$  in step 2 involves knowledge of the density matrix, the Roothaan equations are solved iteratively by repeating steps 2–4 to self-consistency.

It is now possible to compare the molecular scheme with the main steps of the CRYSTAL program:

1. Forming the basis of Bloch functions as linear combinations of the local basis of the AOs (Eq. [29]), in turn expressed as contractions of Gaussian



functions times angular functions, followed by evaluating the overlap matrix in the local basis set.

2. Evaluating Fock ( $F_{\mu\nu}^g$ ) matrix elements in direct space in the local basis set; the average value of the Fock operator with respect to the AOs  $\chi_\mu(\mathbf{r} - \mathbf{r}_\mu)$  in the 0-cell and  $\chi_\nu^g(\mathbf{r} - \mathbf{r}_\nu)$  in the  $\mathbf{g}$ -cell is calculated as a sum of the following contributions:

$$F_{\mu\nu}^g = \langle \chi_\mu | \hat{F} | \chi_\nu^g \rangle = T_{\mu\nu}^g + Z_{\mu\nu}^g + C_{\mu\nu}^g + X_{\mu\nu}^g \quad [34]$$

Every matrix element is properly identified by the three indices  $\mu$ ,  $\nu$ , and  $\mathbf{g}$ , which specify the two AOs and the direct lattice vector  $\mathbf{g}$  labeling the cell where the  $\nu$ -th AO is centered because, in principle, the origin of  $\chi_\nu^g(\mathbf{r} - \mathbf{r}_\nu)$  can be anywhere in the crystal. The possibility of always referring  $\chi_\mu$  to the 0-cell is because of translation invariance of the integrals in the local basis, for example:

$$\langle \chi_\mu^{g'} | \hat{F} | \chi_\nu^g \rangle = \langle \chi_\mu^0 | \hat{F} | \chi_\nu^{g-g'} \rangle = \langle \chi_\mu^0 | \hat{F} | \chi_\nu^m \rangle \quad [35]$$

with  $\mathbf{m} = \mathbf{g} - \mathbf{g}'$  being a direct lattice vector.

3. Representing  $\mathbf{S}$  and  $\mathbf{F}$  matrices in the Bloch function basis set at every  $\mathbf{k}$  point of the sampling set. In this basis, the expression of the matrix elements contains a double sum over the direct lattice vectors. For example, a generic element of the Fock matrix represented in the reciprocal space is given by

$$F_{\mu\nu}(\mathbf{k}) = \langle \Phi_\mu(\mathbf{k}) | \hat{F} | \Phi_\nu(\mathbf{k}) \rangle = \frac{1}{N} \sum_{\mathbf{g}'} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot (\mathbf{g} - \mathbf{g}')} \langle \chi_\mu^{g'} | \hat{F} | \chi_\nu^g \rangle$$

Nevertheless, by taking Eq. [35] into account, the double sum reduces to  $N$  times a single sum and the expression can be simplified to

$$F_{\mu\nu}(\mathbf{k}) = \sum_{\mathbf{m}} e^{i\mathbf{k} \cdot \mathbf{m}} \langle \chi_\mu^0 | \hat{F} | \chi_\nu^m \rangle \quad [36]$$

This last equation can be interpreted as a Fourier transform of the Fock matrix from direct to reciprocal space.

4. Solving the Schrödinger equation with the orthonormality condition [26] at every  $\mathbf{k}$

$$\mathbf{F}(\mathbf{k}) \mathbf{C}(\mathbf{k}) = \mathbf{S}(\mathbf{k}) \mathbf{C}(\mathbf{k}) \mathbf{E}(\mathbf{k}) \quad [37]$$

5. Determining the Fermi energy,  $E_F$ , which is the highest energy value of an occupied state in the system inside the first Brillouin zone.

6. Forming the density matrix  $\mathbf{P}$  and Fourier anti-transforming it to direct space

$$p_{\mu\nu}^g = \frac{1}{V_{BZ}} \sum_n \int_{BZ} e^{i\mathbf{k}\cdot\mathbf{g}} C_{\mu n}^*(\mathbf{k}) C_{\nu n}(\mathbf{k}) \theta(E_F - E_n(\mathbf{k})) d\mathbf{k} \quad [38]$$

Here the sum over  $\mathbf{k}$  points has become an integral over the first Brillouin zone (with volume  $V_{BZ}$ ), because it has already been shown that  $\mathbf{k}$  can be considered as a continuous variable. By limiting the integration to states with energy below  $E_F$ , a Heaviside step function  $\theta$  permits us to exclude the eigenvectors relative to empty states from the sum. The reason why this cannot be achieved by simply truncating the sum over the eigenvectors, like in the molecular case, will be clear when the main features of band structure are illustrated later on in this chapter.

7. Calculating the total energy per cell as

$$E_t = N + \frac{1}{2} \sum_{\mu,\nu} \sum_g p_{\mu\nu}^g (T_{\mu\nu}^g + Z_{\mu\nu}^g + F_{\mu\nu}^g) \quad [39]$$

The total energy of an infinite crystal is obviously infinite and has no physical meaning, but the total energy per cell, which includes the interaction of the nuclei and electrons in the 0-cell with all nuclei and electrons in the crystal, is finite. In this expression, a new sum over the infinite direct lattice vectors appears.

Again, steps 2–6 are iterated to self-consistency. Basically, two aspects are specific for the application of this method to solids: the calculation of matrices in direct space, which involve multiple sums over all the infinite direct lattice vectors, and the integration in reciprocal space. This latter aspect will be discussed with reference to a few specific examples in the next sections.

As an example of the problems involved in the evaluation of matrix elements for a periodic system, we consider the explicit form of the Coulomb electron–electron repulsion term in Eq. [34]

$$C_{\mu\nu}^g = \sum_{\lambda} \sum_1 \sum_{\rho} \sum_m P_{\lambda\rho}^{m-1} \left\langle \chi_{\mu} \chi_{\lambda}^1 \left| \frac{1}{r_{12}} \right| \chi_{\nu}^g \chi_{\rho}^m \right\rangle \quad [40]$$

where the three-index notation can be used for the density matrix elements, as a consequence of translation invariance, and  $\chi_{\mu}$  and  $\chi_{\nu}^g$  in the two-electron integrals refer to electron 1, whereas  $\chi_{\lambda}^1$  and  $\chi_{\rho}^m$  refer to electron 2, with  $r_{12}$  being the mutual distance between electrons. The presence of the two infinite sums over the direct lattice vectors in Eq. [40] complicates the calculation of these terms dramatically. Thus, it is really the system size and the amount of

long-range interactions involved that make high accuracy and numerical stability important and more demanding targets than in molecular cases. Finding a solution to this problem has required a deep analysis of the convergence properties<sup>48–50,54,55</sup> of the series to be evaluated. This analysis has resulted both in the formulation of convenient truncation criteria and in the application of Ewald's<sup>52</sup> method to the calculation of long-range interactions in the slowly convergent Coulomb series. Fortunately, a local basis set is suitable to the definition of series truncation schemes, because an estimate of the importance of interparticle interactions is relatively easy in direct space.

## DISCUSSION OF BAND STRUCTURE THROUGH A FEW SIMPLE EXAMPLES

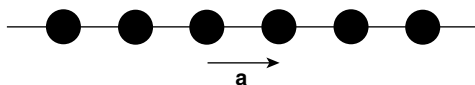
### A Monoatomic Linear Chain

The simplest case of a periodic model that one can imagine is a polymer with an atom per cell and one electron per atom, as depicted in Figure 5.

All atoms are equivalent by translation  $a$ . Although this is only a simplified ideal situation, it is used as a first example in many introductory textbooks on physical chemistry and solid state physics, because it is a problem simple enough to be treated analytically, especially if an easy approximation such as Hückel's model Hamiltonian is applied. Several important simplifications in Hückel's model make the calculation very easy, while preserving the main topological characteristics of the system. In this simple model, only one  $p_z$  AO is considered for each atom. The different orbitals will be identified by the  $g$  lattice vector of the cell in which they are centered and denoted as  $p_z^g$ . Hückel's approximation prescribes simple rules for the determination of the overlap and the Hamiltonian matrices with two parameters,  $\alpha$  and  $\beta$ :

$$\langle p_z^g | p_z^l \rangle = \delta_{gl} \quad \langle p_z^g | \hat{H} | p_z^l \rangle = \begin{cases} \alpha & g = l \\ \beta & |l - g| = a \\ 0 & |l - g| > a \end{cases} \quad [41]$$

Before evaluating  $H(k)$  and  $S(k)$  following Hückel's prescriptions, a basis set of Bloch functions must be defined in the local basis of the  $p_z$  AOs in the polymer



**Figure 5** Schematic representation of a Bloch function for a monoatomic linear chain. A black circle represents a  $p_z$  AO.

according to Eq. [29]. However, this case is particularly simple, because there is only one atom in each cell and we are considering only one AO per atom (Figure 5). Therefore, only one Bloch function can be generated

$$\Phi(\mathbf{r}; \mathbf{k}) = \frac{1}{\sqrt{N}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} p_z^{\mathbf{g}}(\mathbf{r}) \quad [42]$$

and all matrices are one-dimensional, i.e., all terms in Eq. [25] are actually single-variable functions of  $\mathbf{k}$ .

Because the basis set is minimal (only one Bloch function per  $\mathbf{k}$  point), this Bloch function is itself an eigenfunction of the Hamiltonian. With this basis set,  $\mathbf{S}$  takes the form:

$$\mathbf{S}(\mathbf{k}) = \langle \Phi(\mathbf{k}) | \Phi(\mathbf{k}) \rangle = \frac{1}{N} \sum_{\mathbf{g}, \mathbf{l}} e^{i\mathbf{k} \cdot (\mathbf{l} - \mathbf{g})} \langle p_z^{\mathbf{g}} | p_z^{\mathbf{l}} \rangle \quad [43]$$

but, because of translation invariance, the overlap depends only on the distance between the AOs, so that any integral can be translated into the reference cell and any linear combination of lattice vectors is, again, a lattice vector, say  $\mathbf{m} = \mathbf{l} - \mathbf{g}$ . For these reasons, the double sum in the expression of  $\mathbf{S}(\mathbf{k})$  reduces to  $N$  times a single sum (like in Eq. [36]), and on the basis of Hückel's rules,  $\mathbf{S}(\mathbf{k})$  results to be constant in  $\mathbf{k}$

$$\mathbf{S}(\mathbf{k}) = \frac{1}{N} \sum_{\mathbf{g}, \mathbf{l}} e^{i\mathbf{k} \cdot (\mathbf{l} - \mathbf{g})} \langle p_z^0 | p_z^{\mathbf{l} - \mathbf{g}} \rangle = \sum_{\mathbf{m}} e^{i\mathbf{k} \cdot \mathbf{m}} \langle p_z^0 | p_z^{\mathbf{m}} \rangle = \sum_{\mathbf{m}} e^{i\mathbf{k} \cdot \mathbf{m}} \delta_{\mathbf{m}0} = 1 \quad [44]$$

Function  $\mathbf{H}(\mathbf{k})$  is obtained in a similar way

$$\begin{aligned} \mathbf{H}(\mathbf{k}) &= \langle \Phi(\mathbf{k}) | \hat{H} | \Phi(\mathbf{k}) \rangle = \frac{1}{N} \sum_{\mathbf{g}, \mathbf{l}} e^{i\mathbf{k} \cdot (\mathbf{l} - \mathbf{g})} \langle p_z^{\mathbf{g}} | \hat{H} | p_z^{\mathbf{l}} \rangle = \\ &= e^{i\mathbf{k} \cdot 0} \alpha + e^{i\mathbf{k} \cdot \mathbf{a}} \beta + e^{-i\mathbf{k} \cdot \mathbf{a}} \beta = \alpha + 2\beta \cos(ka) \end{aligned} \quad [45]$$

In this simple case, both  $\mathbf{C}(\mathbf{k})$  and  $\mathbf{S}(\mathbf{k})$  are constant and equal to 1; thus, the eigenvalue  $\mathbf{E}(\mathbf{k})$  coincides with  $\mathbf{H}(\mathbf{k})$ . It is apparent that the eigenvalue spectrum of a periodic system like this polymer does not consist of discrete energy levels as occurring in the case of an atom or a molecule (it would consist of a single level for one atom of this polymer). Instead, the eigenvalue spectrum includes all possible energy values within a definite range (between  $\alpha + 2\beta$  and  $\alpha - 2\beta$ ), forming a band, as represented in Figure 6 (parameters  $\alpha$  and  $\beta$  are both negative).

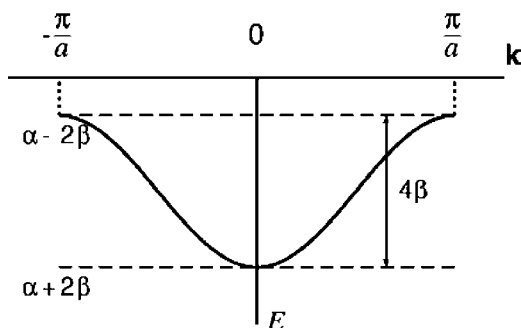


Figure 6 Hückel  $\pi$ -band of a monoatomic linear chain in the first Brillouin zone.

## A Two-Dimensional Periodic Example: Graphite

As a second example, we consider graphite. Although graphite in a pencil consists of a large number of layers of the type represented in Figure 7, it is well known that we can use a pencil for writing or drawing because the inter-layer interactions are weak, much weaker than the intralayer interactions, and difficult to render computationally if we are not using sophisticated methods. For this reason, a single layer of graphite is a good model, which will be used in this example.

A layer of graphite can be considered as infinite and periodic in two dimensions. Graphite has a planar hexagonal structure. It belongs to a layer group that is derivable from the  $P6/mmm$  space group, containing 24 symmetry operations.

The unit cell has the shape of a rhomb with angles of  $60^\circ$  and  $120^\circ$  (Figure 7). The point symmetry elements (six-fold axis normal to the plane, six two-fold rotation axes in the plane, six mirrors normal to the plane, and one in the plane, inversion) pass through the unit cell origin, at the center of the hexagons. There are two symmetry-related carbon atoms in the unit cell, labeled as A and B in Figures 7 and 9, with fractional coordinates  $(1/3, 2/3)$  and

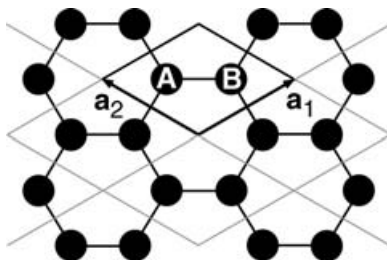


Figure 7 Graphite layer.

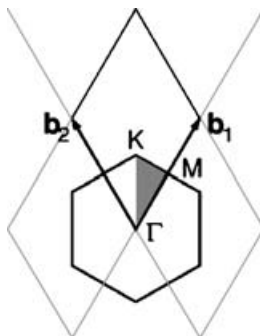
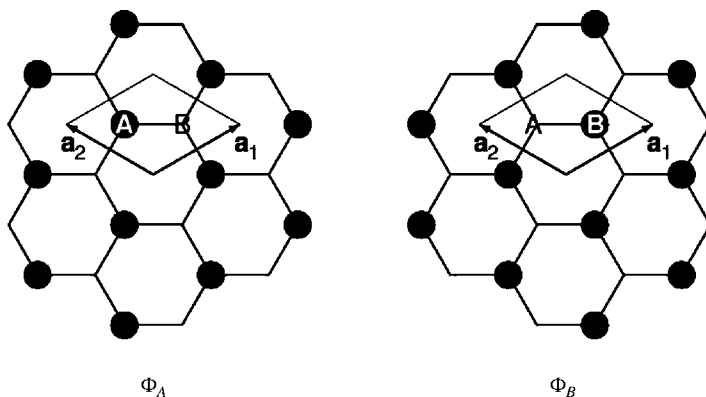


Figure 8 Graphite reciprocal lattice.

$(2/3, 1/3)$ , respectively. Each atom of type A is surrounded by three atoms of type B, and vice versa.

The reciprocal lattice is again hexagonal, as an obvious consequence of the basis vectors  $\mathbf{b}_1$  and  $\mathbf{b}_2$  being orthogonal to  $\mathbf{a}_1$  and  $\mathbf{a}_2$ , according to Eq. [3]. The first Brillouin zone has the shape of a hexagon as in Figure 8, with its center at the lattice origin, where  $\mathbf{k} = 0$ . The shaded triangle in the picture represents the asymmetric unit in the Brillouin zone. Special symbols have been assigned to the points at the vertices of the triangle, which correspond to special positions in the reciprocal lattice. Denoting  $\mathbf{k}$  points by their components along  $\mathbf{b}_1$  and  $\mathbf{b}_2$  as  $(b_1, b_2)$ , M identifies the point at the top of the triangle with components  $(\frac{1}{2}, 0)$ , K denotes point  $(\frac{1}{3}, \frac{1}{3})$ , and  $\Gamma$  is the lattice origin  $(0, 0)$ . Actually,  $\Gamma$  identifies the origin in any reciprocal lattice.

Also in this case,  $\pi$  electron bands can be studied with Hückel's approximation by representing all matrices in a basis set of two Bloch functions from

Figure 9 Representation of Bloch functions for graphite. A circle represents a  $p_z$  AO.

$p_{zA}$  and  $p_{zB}$  AOs according to Eq. [29]

$$\Phi_A(\mathbf{r}; \mathbf{k}) = \frac{1}{\sqrt{N}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} p_{zA}^{\mathbf{g}}(\mathbf{r}) \quad \Phi_B(\mathbf{r}; \mathbf{k}) = \frac{1}{\sqrt{N}} \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} p_{zB}^{\mathbf{g}}(\mathbf{r}) \quad [46]$$

At each  $\mathbf{k}$  point,  $\mathbf{H}$ ,  $\mathbf{S}$ ,  $\mathbf{C}$ , and  $\mathbf{E}$  in Eq. [25] are  $2 \times 2$  matrices in the basis of  $\Phi_A$  and  $\Phi_B$ . Again,  $\mathbf{H}$  and  $\mathbf{S}$  must be computed to find the unknown matrices  $\mathbf{C}$  and  $\mathbf{E}$ . For symmetry reasons, only the asymmetric part of the Brillouin zone needs to be explored instead of the entire zone, and because the full representation of a  $\pi$ -type band structure for graphite would be three-dimensional, we can start by exploring a representative monodimensional path (as is usually done for three-dimensional structures). The most obvious choice is the triangle perimeter, and in particular, the special positions  $\Gamma$ ,  $M$ , and  $K$  are expected to be topologically interesting for their symmetry properties.

During the calculation of  $\mathbf{S}$  and  $\mathbf{H}$ , it must be taken into account that, because of the symmetry equivalence of the carbon atoms,  $S_{AA}(\mathbf{k}) = S_{BB}(\mathbf{k})$  and  $H_{AA}(\mathbf{k}) = H_{BB}(\mathbf{k})$ , whereas  $H_{AB}(\mathbf{k}) = H_{BA}^{\dagger}(\mathbf{k})$  for hermiticity. Moreover, Hückel's rules imply that  $\mathbf{S}(\mathbf{k})$  is the identity matrix of order 2, completely independent of  $\mathbf{k}$ . In fact, diagonal elements can be computed in exactly the same way as for the linear chain, even though the geometry is different in this case, and the off-diagonal elements are zero because of Hückel's orthogonality assumption. Consequently, Eq. [25] becomes

$$\mathbf{H}(\mathbf{k})\mathbf{C}(\mathbf{k}) = \mathbf{C}(\mathbf{k})\mathbf{E}(\mathbf{k}) \quad [47]$$

We start by computing matrix  $\mathbf{H}$  at  $\Gamma$ . This is a peculiar point because each of the two Bloch functions in Eq. [46] reduces to a simple sum of all AOs of that type (A or B) in the lattice (all factor phases are 1 when  $\mathbf{k} = 0$ ).

$H_{AA}$  can be calculated in a similar way as for the linear chain, because a carbon atom does not have any nearest neighbor of the same type, as can be easily seen in Figure 9. Applying Hückel's rules, we find:

$$\begin{aligned} H_{AA}(0) &= \langle \Phi_A(0) | \hat{H} | \Phi_A(0) \rangle = \frac{1}{N} \sum_{\mathbf{g}, \mathbf{l}} \langle p_{zA}^{\mathbf{g}} | \hat{H} | p_{zA}^{\mathbf{l}} \rangle = \sum_{\mathbf{m}} \langle p_{zA}^0 | \hat{H} | p_{zA}^{\mathbf{m}} \rangle \\ &= \langle p_{zA}^0 | \hat{H} | p_{zA}^0 \rangle = \alpha \end{aligned} \quad [48]$$

This result is actually unrelated to the choice of  $\Gamma$ , but it comes from the fact that all three nearest neighbours of atom A are type-B atoms. Therefore, the same value of  $H_{AA}$  is to be expected at any  $\mathbf{k}$  point. For this reason, the calculation of this element will not be repeated in the following cases.

The value of  $H_{AB}(0)$  is, again, determined by the number of neighboring atoms of carbon A. In fact, only the  $p_{zB}$  AOs in the  $0$ ,  $-\mathbf{a}_1$ , and  $\mathbf{a}_2$  cells can

contribute a nonzero interaction with  $p_{zA}$  in the 0-cell:

$$H_{AB}(0) = \langle \Phi_A(0) | \hat{H} | \Phi_B(0) \rangle = \sum_m \langle p_{zA}^0 | \hat{H} | p_{zB}^m \rangle = 3\beta \quad [49]$$

Equation [25] at  $\Gamma$  can now be solved by imposing the condition that the following determinant annihilates:

$$|\mathbf{H}(0) - E^\Gamma \mathbf{I}| = \begin{vmatrix} \alpha - E^\Gamma & 3\beta \\ 3\beta & \alpha - E^\Gamma \end{vmatrix} = 0 \quad [50]$$

which implies that two possible energy values exist at  $\Gamma$

$$E_\pm^\Gamma = \alpha \pm 3\beta \quad [51]$$

The corresponding eigenfunctions are obtained by replacing each of the two values of  $E^\Gamma$  in Eq. [47] in turn ( $E_+$  and  $E_-$  to obtain  $\Psi_+$  and  $\Psi_-$ , respectively) and imposing the normalization condition of

$$c_A^2 + c_B^2 = 1 \quad [52]$$

The solution of this two-equation system with two unknown coefficients is a fully constrained algebraic problem, which leads to a symmetric and an anti-symmetric linear combination of Bloch functions, as a consequence of the equivalence of carbon atoms of types A and B. These solutions correspond to  $\pi$ -bonding and  $\pi$ -antibonding COs having the following form:

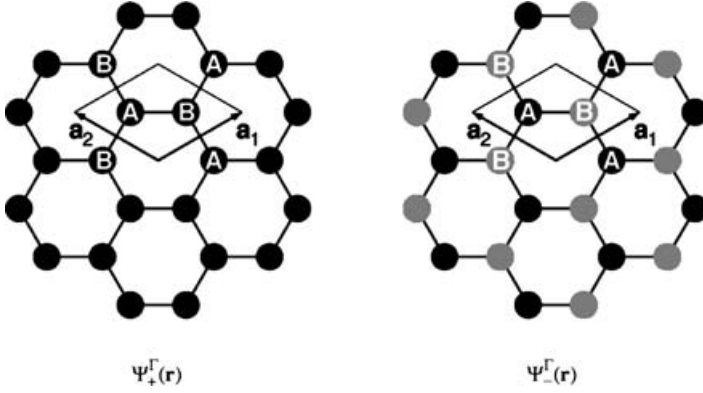
$$\Psi_+^\Gamma(\mathbf{r}) = \frac{1}{\sqrt{2}} [\Phi_A(\mathbf{r}; \mathbf{0}) + \Phi_B(\mathbf{r}; \mathbf{0})] \quad \Psi_-^\Gamma(\mathbf{r}) = \frac{1}{\sqrt{2}} [\Phi_A(\mathbf{r}; \mathbf{0}) - \Phi_B(\mathbf{r}; \mathbf{0})] \quad [53]$$

which are represented schematically in Figure 10. The AOs in  $\Psi_+^\Gamma$  are all in phase and define a totally symmetric combination, which is known to correspond to the most stable state. Contrarily, the antibonding CO represents the most unstable state, where the number of nodal planes is maximum.

At point M, taking Eq. [3] into account, Bloch functions can be written as

$$\Phi_A(\mathbf{r}; \mathbf{M}) = \frac{1}{\sqrt{N}} \sum_{n_1, n_2} e^{i\mathbf{b}_1 \cdot (n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2)} p_{zA}^{n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2}(\mathbf{r}) = \frac{1}{\sqrt{N}} \sum_{n_1, n_2} e^{in_1 \pi} p_{zA}^{n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2}(\mathbf{r}) \quad [54]$$





**Figure 10** Representations of bonding  $\Psi_+^{\Gamma}(\mathbf{r})$  and antibonding  $\Psi_-^{\Gamma}(\mathbf{r})$   $\pi$ -crystalline orbitals at  $\Gamma$  in graphite. Black and gray circles represent the positive and negative signs of each  $p_z$  AO in the linear combination, respectively.

Because  $n_1$  and  $n_2$  are integers, the phase factor can be evaluated straightforwardly, as

$$\Phi_A(\mathbf{r}; \mathbf{M}) = \frac{1}{\sqrt{N}} \sum_{n_1, n_2} (-1)^{n_1} p_{zA}^{n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2}(\mathbf{r}) \quad [55]$$

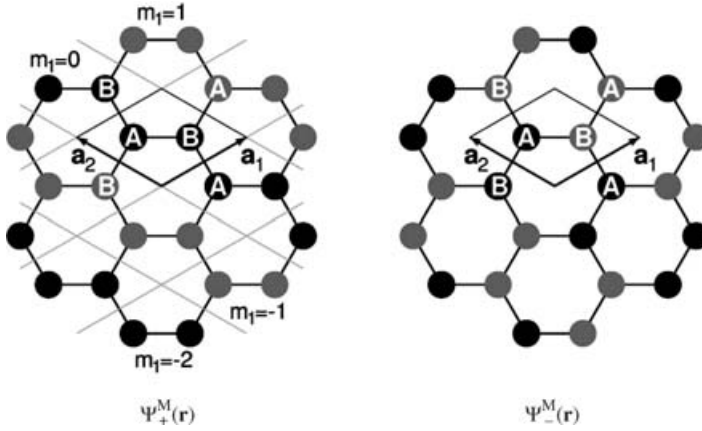
and similarly for  $\Phi_B$ . The form of  $\Phi_A$  implies an inversion of the sign of the AOs at  $\mathbf{M}$  when  $n_1$  is odd, so that the AOs change their signs every time  $n_1$  is incremented by 1. This change corresponds to the alternation of the rows of the AOs with positive and negative signs along the direction of  $\mathbf{a}_1$ , like a wave, as shown in the left picture in Figure 11, where  $\Phi_A$  and  $\Phi_B$  are combined in phase.

In this basis set and considering that two of the nearest neighbors of atom A (see Figure 11) are associated with a phase factor with the same sign, whereas the third nearest neighbor has opposite sign:

$$\begin{aligned} H_{AB}(\mathbf{M}) &= \langle \Phi_A(\mathbf{M}) | \hat{H} | \Phi_B(\mathbf{M}) \rangle = \sum_{m_1, m_2} (-1)^{m_1} \langle p_{zA}^0 | \hat{H} | p_{zB}^{m_1 \mathbf{a}_1 + m_2 \mathbf{a}_2} \rangle \\ &= (2 - 1)\beta = \beta \end{aligned} \quad [56]$$

In this case, the solution of Eq. [47] implies the fulfilment of the following condition:

$$\begin{vmatrix} \alpha - E^M & \beta \\ \beta & \alpha - E^M \end{vmatrix} = 0 \quad [57]$$



**Figure 11** Representation of bonding and antibonding  $\pi$ -crystalline orbitals at M in graphite.

leading to the following eigenvalues at point M:

$$E_{\pm}^M = \alpha \pm \beta \quad [58]$$

The corresponding bonding and antibonding eigenfunctions are, again, a symmetric and an antisymmetric combination of  $\Phi_A$  and  $\Phi_B$ , depicted schematically in Figure 11. The in-phase interactions among AOs of chains along the direction of  $\mathbf{a}_2$  confer stability to the bonding CO, although not as large as the completely in-phase combination in  $\Gamma$ . Conversely, the antibonding CO is less destabilized at M than at  $\Gamma$ , which owes to one in-phase interaction between  $p_{zA}$  and  $p_{zB}$  in the former.

The third highly symmetric point is K. In this case

$$\begin{aligned} \Phi_A(\mathbf{r}; \mathbf{K}) &= \frac{1}{\sqrt{N}} \sum_{n_1, n_2} e^{i(\mathbf{b}_1 + \mathbf{b}_2) \cdot (n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2)} p_{zA}^{n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2}(\mathbf{r}) \\ &= \frac{1}{\sqrt{N}} \sum_{n_1} e^{i \frac{2\pi}{3}(n_1 + n_2)} p_{zA}^{n_1 \mathbf{a}_1 + n_2 \mathbf{a}_2}(\mathbf{r}) \end{aligned} \quad [59]$$

The phase factor in this expression can only have three possible values, depending on the sum of the two integers  $n_1$  and  $n_2$ . In fact, for any integer  $m$ ,

$$e^{i \frac{2\pi}{3}(n_1 + n_2)} = \begin{cases} 1 & n_1 + n_2 = 3m \\ -\frac{1}{2} + i \frac{\sqrt{3}}{2} & \forall \quad n_1 + n_2 = 3m + 1 \\ -\frac{1}{2} - i \frac{\sqrt{3}}{2} & n_1 + n_2 = 3m + 2 \end{cases} \quad [60]$$

At point K, the Bloch function basis also contains an imaginary part, in addition to the real part, which is not the case at  $\Gamma$  or M. The evaluation of the  $H_{AB}$  interaction element of the Hamiltonian matrix leads to the following result:

$$\begin{aligned} H_{AB}(K) &= \langle \Phi_A(K) | \hat{H} | \Phi_B(K) \rangle = \sum_{m_1, m_2} e^{i\frac{2\pi}{3}(m_1+m_2)} \langle p_{zA}^0 | \hat{H} | p_{zB}^{m_1 a_1 + m_2 a_2} \rangle = \\ &= \left[ -\frac{1}{2} - \frac{1}{2} + 1 + i \left( \frac{\sqrt{3}}{2} - \frac{\sqrt{3}}{2} \right) \right] \beta = 0 \end{aligned} \quad [61]$$

and the Hamiltonian is already in the diagonal form. The equation to be solved is

$$\begin{vmatrix} \alpha - E^K & 0 \\ 0 & \alpha - E^K \end{vmatrix} = 0 \quad [62]$$

Thus, clearly a degeneracy of the bonding and the antibonding states at K occurs

$$E_{\pm}^K = \alpha \quad [63]$$

The real and imaginary parts of the corresponding COs are represented in Figure 12.

So far, the energy of bonding and antibonding  $\pi$ -COs are degenerate at K and split at  $\Gamma$  and M, with the splitting being larger at  $\Gamma$ .

What happens at other points, for instance, at  $\mathbf{k} = (\frac{1}{4}, 0)$ , midway the  $\Gamma$ -M path? In this case, it is easily seen that

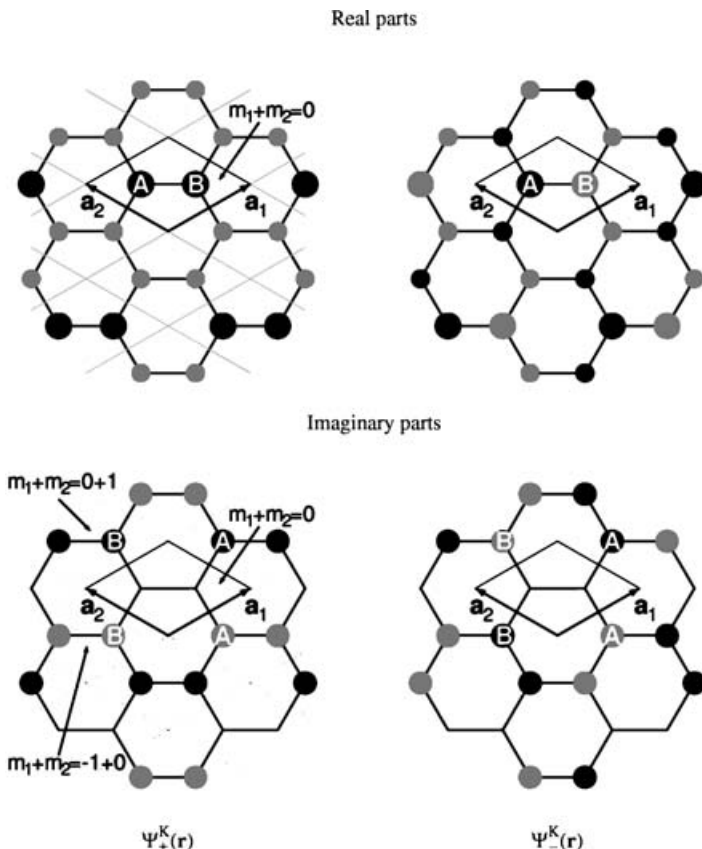
$$\Phi_A\left(\mathbf{r}; \frac{1}{4}\mathbf{b}_1\right) = \frac{1}{\sqrt{N}} \sum_{n_1, n_2} i^{n_1} p_{zA}^{n_1 a_1 + n_2 a_2}(\mathbf{r}) \quad [64]$$

and

$$H_{AB}\left(\frac{1}{4}\mathbf{b}_1\right) = \sum_{m_1, m_2} i^{m_1} \langle p_{zA}^0 | \hat{H} | p_{zB}^{m_1 a_1 + m_2 a_2} \rangle = (2 - i)\beta \quad [65]$$

The eigenvalues can then be found by solving the following equation:

$$\begin{vmatrix} \alpha - E^{(\frac{1}{4}, 0)} & (2 - i)\beta \\ (2 + i)\beta & \alpha - E^{(\frac{1}{4}, 0)} \end{vmatrix} = 0 \quad [66]$$



**Figure 12** Representation of bonding and antibonding  $\pi$ -crystalline orbitals at K in graphite.

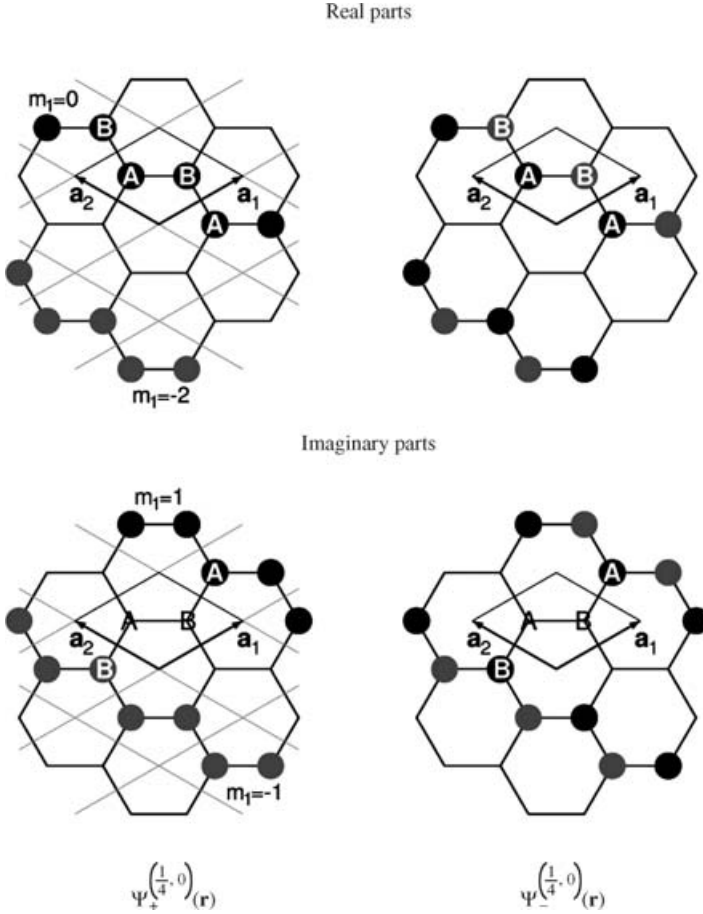
The two possible energy values at this point

$$E(\frac{1}{2}0) = \alpha \pm \sqrt{5}\beta \quad [67]$$

are intermediate between the corresponding bonding and antibonding  $E^M$  and  $E^\Gamma$ , in fact  $|\beta| < \sqrt{5}|\beta| < 3|\beta|$ . The topology of the wave function at this point (see Figure 13) resembles that in M, with the main differences being that in this case, an imaginary component exists, as well, and that the period of the wave doubles, when  $\mathbf{k}$  is midway to the  $\Gamma$ -M path.

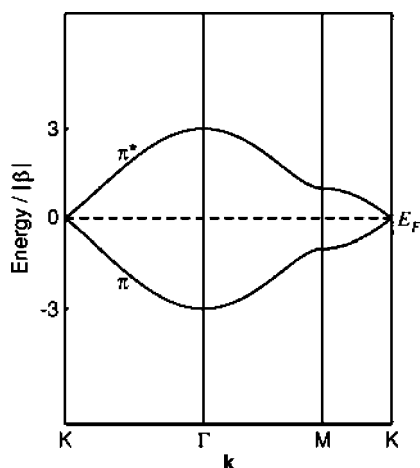
The values of  $E_\pm$  found at special  $\mathbf{k}$  points in the reciprocal lattice are actually special forms of a general expression for  $E_\pm$ , which is a simple function of  $\mathbf{k} = k_1\mathbf{b}_1 + k_2\mathbf{b}_2$  also in this case:

$$E_\pm^{\mathbf{k}} = \alpha \pm \beta \sqrt{[1 + \cos(2k_1\pi) + \cos(2k_2\pi)]^2 + [\sin(2k_2\pi) - \sin(2k_1\pi)]^2} \quad [68]$$



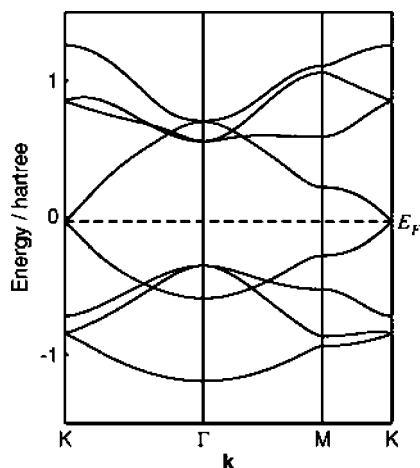
**Figure 13** Representation of bonding and antibonding  $\pi$ -crystalline orbitals at  $(\frac{1}{4}, 0)$  in graphite.

The graphical representation of Eq. [68] along the contour of the triangle  $\Gamma\text{KM}$  (Figure 14) clearly shows that the energy of the  $\pi$  bonding and antibonding COs is a continuous and smooth function of  $\mathbf{k}$ , with  $\text{K}$  being the only point with degeneracy and  $\Gamma$  being the point with maximum energy splitting, where  $E_+$  is the minimum and  $E_-$  is the maximum energy. For these characteristics, graphite is reported to be a zero-gap semiconductor, because in fact bands do not cross but are tangent in  $\text{K}$ . Equation [68] allows us to compute  $E_{\pm}$  at any other point in the triangle, which confirms that all these values vary continuously within the range between  $E^{\text{K}}$  and  $E^{\Gamma}$ . Thus, the chosen path is really representative of the band structure, as expected. The shape of these bands is qualitatively preserved, even when we abandon the crude parametric Hückel's approximation and study the graphite band structure ab initio. The



**Figure 14** Hückel  $\pi$ -bands of graphite along the K- $\Gamma$ -M-K path in reciprocal space. Energy is measured in units of  $\beta$  (absolute value) and referred to as  $\alpha$ .  $E_F$  stands for the Fermi energy.

band structure of graphite in Figure 15 has been obtained with the HF approximation and an extended, all-electron basis set. Direct comparison with bands in Figure 14 is possible because of the one-electron character of the HF approximation, which permits the assignment of each band to a well-defined one-electronic state. The corresponding  $\pi$  and  $\pi^*$  bands compare fairly well,



**Figure 15** Hartree-Fock valence and conduction band structure of graphite along the K- $\Gamma$ -M-K path.

although it must be noted that the energy scales used in Figures 14 and 15 are arbitrarily different. The reason why this correspondence happens is that the shape of bands is largely determined by topological features; here, where symmetry is high, a poor approximation of the interactions affects only the energy range and the bandwidth. For example, an extended basis set removes the artifact of complete symmetry between  $\pi$  and  $\pi^*$  bands in Figure 14. However, qualitative aspects, such as band maxima and minima or points of degeneracy, are essentially determined by the symmetry properties of the system. Another remarkable point is that the *bandwidth*, which is a measure of *dispersion* in  $k$ -space, clearly depends on the magnitude and the range of the interactions within the crystal (corresponding to the value of  $\beta$  and the number of interactions among the nearest neighbors in Hückel's theory).

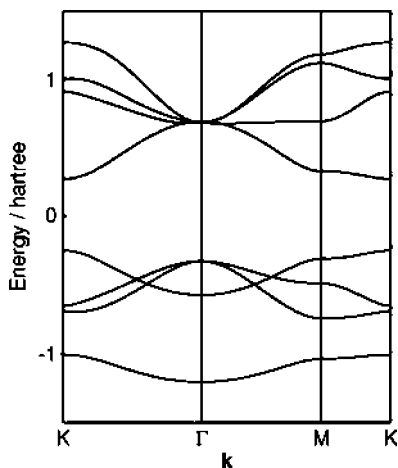
The band structure reported in Figure 15 has been obtained by considering all electrons in graphite. Core bands are not represented in the plot because they are separated by a large energy gap from valence bands. However, they are completely flat, which implies that core electrons are effectively screened by valence electrons and do not couple to the rest of the crystal. The four lowest bands in the figure are valence bands, and above them are the first virtual (or *conduction*) bands. The highest energy associated with a populated state of the crystal in its fundamental state is called the *Fermi energy*, the analogue of the highest occupied molecular orbital (HOMO) energy level in molecular cases. The separation between the top of the valence band and the bottom of the conduction band is known as the *conduction-valence gap* (it is zero in this case). Obviously, transitions from the valence to the conduction band are not restricted to this energy, but they can also occur, for example, from the valence band bottom to the conduction band top, so that electron excitation spectra also exhibit bands of possible electronic transitions.

A comparison of graphite with hexagonal boron nitride illustrates more clearly how significant topology is in determining the band structure of a compound and in affecting its properties. Planar BN is isostructural and isoelectronic to graphite, with the obvious exception that some symmetry operations are lost because atoms A and B in the unit cell are no longer equivalent. It belongs to group  $P\bar{6}m2$ , consisting of 12 symmetry operations. The great similarity of the band structure of BN (Figure 16) to that of graphite (Figure 15) is remarkable: The only important difference is the appearance of a gap between the valence and conduction bands, because degeneracy is lost at K.

This behavior of  $\pi$  and  $\pi^*$  bands can easily be interpreted in terms of simple Hückel's model. In fact, in this case, two different  $\alpha$  parameters would appear in Eq. [62],  $\alpha_B$  and  $\alpha_N$ , and the two solutions would necessarily be distinct:

$$E_+ = \alpha_N \quad E_- = \alpha_B \quad [69]$$

generating a band gap and making BN a semiconductor with more usual characteristics.



**Figure 16** Hartree-Fock upper valence and lower conduction band structure of BN along the K- $\Gamma$ -M-K path.

### Three-Dimensional Periodic Examples

Three-dimensional periodicity makes real crystals more complex structures than polymers and layers, owing to the larger number of the interactions and geometric arrangements involved. Nevertheless, the basic ideas underlying the way band structures originate are essentially the same as those previously mentioned. As it was pointed out, crystals exhibit different properties depending on the different characteristics of their band structure, and in particular, one major classification is based on the extension of the gap between the valence and the conduction bands. Three examples that are representative of different behaviors from this point of view have been included in Figure 17.

Magnesium oxide is a cubic, almost fully ionic oxide with a large band gap, and for this reason, it is classified as an insulator because, even if the real extension of the gap is not as large as is predicted by the HF approximation, the amount of energy that is required by electrons to undergo a transition to virtual states is far beyond the thermal energy at room temperature. Also silicon, which is a covalent crystal, exhibits a band gap. However, this conduction-valence gap is smaller than in the case of MgO, so that silicon is a semiconductor, where virtual states are more easily accessible to electrons. Conversely, no gap is detectable in beryllium, because valence and conduction bands intercross each other. Fermi level passes across them, and a large amount of available empty states are accessible to valence electrons, so that the conduction phenomenon can be easily induced by applying some potential difference through the crystal. For this peculiarity, beryllium is a conductor.



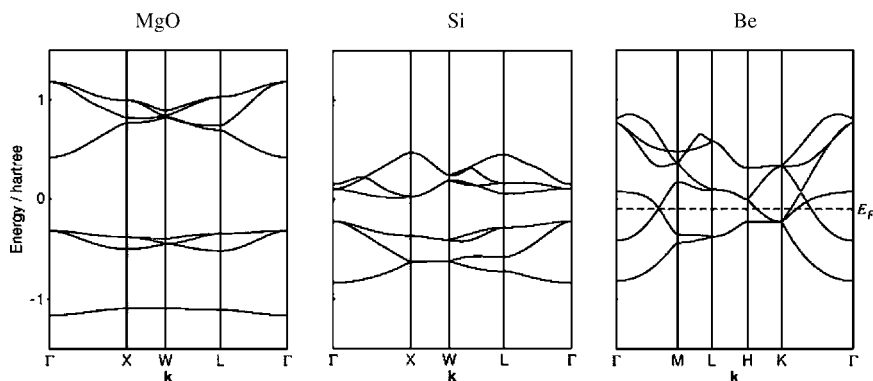


Figure 17 Hartree-Fock upper valence and lower conduction band structure of magnesium oxide, silicon, and beryllium.

The role of crystal symmetry properties in determining the shape of the bands has been emphasized, but the few examples reported have also shown that the existence of a gap and the energy range of bands depend on the mutual interactions of all particles, electrons, and nuclei, in the lattice. Therefore, the correctness of a calculation is largely dependent on the kind of approximation used in the evaluation of such interactions. In fact, different approximations of the Hamiltonian can produce a variety of results and, in particular, band structures that are not only quantitatively but also qualitatively different in some cases. In Figure 18, the HF band structure of silicon is compared with that obtained with DFT methods, both in the LDA, in the form of Slater-Vosko-Wilk-Nusair<sup>6,90</sup> functional, and with the Becke 3 (B3) parameter-Lee-Yang-Parr (LYP) approximation,<sup>91</sup> which incorporates a part of the exact exchange

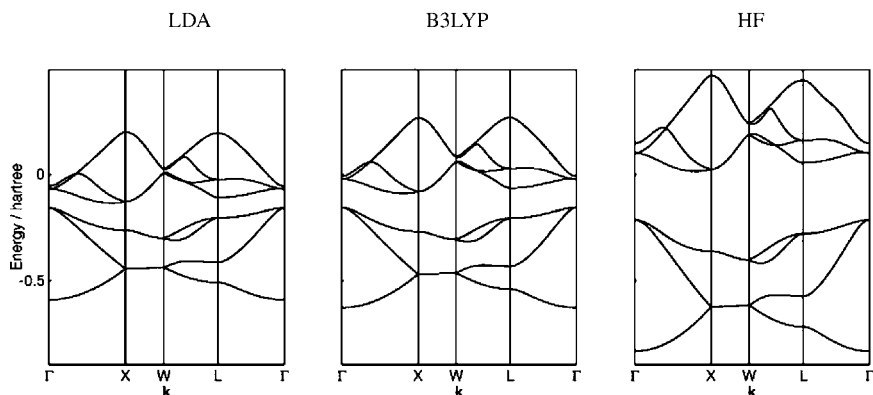


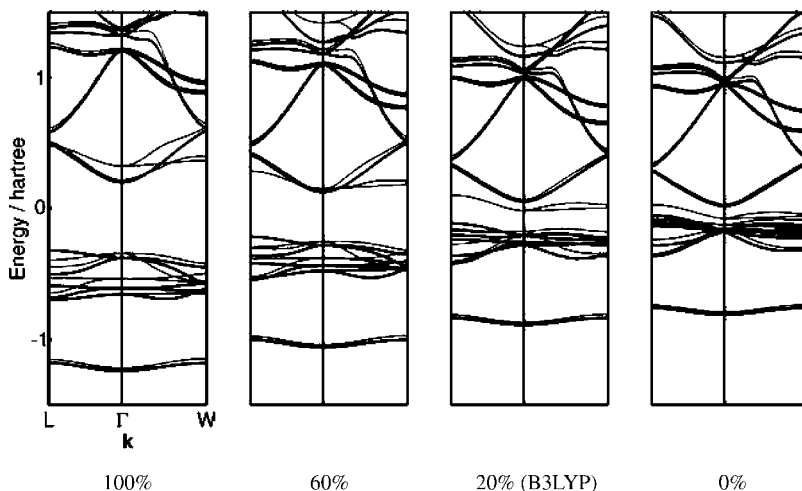
Figure 18 LDA, B3LYP and HF upper valence and lower conduction band structure of silicon.

into the exchange-correlation functional based on the generalized gradient approximation. In spite of the similarity of DFT and HF bands in regard to their shape, the difference between these band structures concerns the basic physical properties of the system. Silicon is predicted to be nearly a conductor by LDA, whereas stimulating conductivity in the HF silicon would be harder. The band gap predicted by LDA for silicon is extremely narrow (0.59 eV), but it is much larger (6.25 eV) according to the HF approximation. The experimental <sup>92</sup> value is 1.17 eV. Also, bandwidths scale differently in HF and DFT. B3LYP bandwidths are close to LDA, but the inclusion of the exact exchange and gradient corrections increases the gap (1.81 eV).

Please note at this point that the HF and Kohn–Sham (KS) eigenvalues (and the resulting band structure) can be related to the excitation spectrum and the conducting properties of a system only in a loose and qualitative way. In HF theory, Koopmans theorem<sup>93</sup> identifies ionization energy  $I_m$  with the  $m$ -th eigenvalue (with negative sign) of the Hamiltonian, under the assumption that relaxation effects can be neglected, which is a rough approximation. In KS theory, only the eigenvalue corresponding to the highest occupied pseudo-orbital has a rigorous physical meaning, in the limit of an “exact” exchange-correlation functional: It is  $-I_m$ .<sup>94,95</sup> Virtual “exact” KS pseudo-orbitals have been shown<sup>94,95</sup> to represent good approximations of the excitation energies in finite systems in some cases, but unfortunately none of the exchange-correlation functionals currently available can reproduce the real potentials equally well.

More accurate techniques exist for the calculation of excitation energies, which apply the HF and KS solutions just as the starting point in the calculation. They are usually indicated as time-dependent DFT<sup>76–78</sup> and density functional perturbation theory.<sup>96,97</sup> As was already mentioned in the Introduction, this matter falls beyond the scope of the present chapter.

The overestimation of bandwidths and gaps obtained with HF compared with the underestimation with LDA is well known. Correct evaluation of the exchange interactions appears as important in band structure calculations, when well-localized electrons are involved. An interesting example is nickel oxide, a compound exhibiting magnetic properties, because of the  $d$ -type unpaired electrons in the configuration of the transition metal ions. The different number of spin-up ( $\alpha$ ) and spin-down ( $\beta$ ) electrons per unit cell causes a polarization in the band structure of NiO. The band structures reported in Figure 19 account for the effect of incorporating different percentages of exact exchange into the Becke–Lee–Yang–Parr exchange-correlation potential.<sup>98,99</sup> B3LYP parameterization corresponds to the inclusion of 20% exact exchange. The edge cases correspond to a pure GGA approximation (on the right end) and to HF corrected by the inclusion of the LYP correlation potential. The width of the gap is clearly proportional to the amount of exact exchange in the exchange-correlation potential, and this term is particularly important in this case to get a correct characterization of NiO as an insulator.



**Figure 19** Upper valence and lower conduction band structure of nickel oxide corresponding to different percentages of exact exchange in Becke–Lee–Yang–Parr exchange–correlation functional. Black and gray lines correspond to spin-up and spin-down states, respectively.

These behaviors and the different performance of the different approximations in this respect are well known. Nevertheless, research over the last 20 years has shown that, despite these large errors in the determination of gaps and bandwidths, these methods perform well in predicting a large variety of observables within an error bar that is in most cases acceptable and helping to draw conclusions about interesting physical and chemical properties of matter in the solid state.

### From the Band Structure to the Total Energy

Solving the HF or KS equations in the present (CRYSTAL) scheme requires numerical integration over the first Brillouin zone because, in general, we do not possess an analytic expression for the eigenvalues and eigenvectors, as is the case of Hückel’s approximation. The question then becomes: How many points need to be sampled, that is, in how many points must Eq. [25] be solved to get sufficiently accurate values of the observables of interest?

The total energy is important and useful to us for answering this question. As discussed, the total energy of an infinite crystal, like in our model, is infinite. Therefore, the total energy per cell is definitely a preferable choice, for it is a finite well-defined property, because of translation invariance. For the sake of clarity, we remind the reader that, although the total energy per cell is defined within the direct lattice context (Eq. [39]), its calculation depends on knowing the density matrix, which in our scheme is obtained from Eq. [38].

Pack and Monkhorst<sup>100</sup> have suggested that a commensurate grid of points is a suitable option for this purpose. In their method, the grid size depends on a parameter, the shrinking factor  $s$ , that specifies how many equidistant  $\mathbf{k}$  points must be taken along each direction of  $\mathbf{b}_1$ ,  $\mathbf{b}_2$ , and  $\mathbf{b}_3$  inside one reciprocal lattice unit cell so that the total number of points in the grid,  $n_s$ , is equal to  $s^n$ , with  $n$  denoting the order of periodicity ( $n = 3$  for three-dimensional crystals).

Magnesium oxide, silicon, and beryllium are three simple and convenient cases for analyzing the accuracy of ab initio periodic calculations in connection with the density of points in the grid. Analyzing their band structures, we can take these systems as representative examples of insulators (MgO), semiconductors (Si), and metals (Be) that are expected to show differences in convergence of various properties as a function of the number of  $\mathbf{k}$  points in the Brillouin zone.

The values of the total energy per cell reported in Table 1 as a function of the shrinking factor show that the size of the Pack–Monkhorst grid is related to the extent of the conduction-valence band gap (16.02, 6.25, 0 eV for the three cases considered, respectively).

By recalling that 1 mhartree ( $\sim 2.6$  kJ/mol) is less than the amount of energy involved in weak hydrogen-bonding, it is clear that a coarse grid in reciprocal space is sufficient to compute the total energy (and wave function) of an insulator, such as MgO, with high accuracy.

When the gap is smaller, like in silicon, a finer grid is needed to obtain comparable accuracy in the estimate of total energy. For example, uncertainty in the determination of the total energy in the order of  $10^{-5}$  hartree is obtained

**Table 1** Total Energy (hartree) per Cell of Magnesium Oxide, Silicon, and Beryllium as a Function of the Shrinking Factor  $s$

$s$	MgO		Si		Be	
	$n_s$	$E$	$n_s$	$E$	$n_s$	$E$
4	8	−274.67867182	8	−577.87453872	12	−29.28470884
5	10	−274.67995386	10	−577.87691374	15	−29.27944675
6	16	−274.67996773	16	−577.87760126	28	−29.26513569
7	20	−274.67996874	20	−577.87777915	32	−29.27586234
8	29	−274.67996884	29	−577.87782962	50	−29.27312660
9	35	−274.67996884	35	−577.87784429	60	−29.27158322
10			47	−577.87784847	84	−29.27088006
12			72	−577.87785017	133	−29.27244875
16			145	−577.87785036	270	−29.27240692
20			256	−577.87785036	484	−29.27264824
24					793	−29.27278307
28					1200	−29.27276667
32					1734	−29.27279225

Bold digits reflect the level of inaccuracy in the calculated energy values.

with  $s = 8$  for silicon and  $s = 5$  for magnesium oxide, whereas  $s$  is 16 versus 8 for accuracy below  $10^{-7}$  hartree.

When no gap exists, as in the case of beryllium, convergence of the total energy with respect to the grid size is much slower. Data in the last column of the table show that a gain of one order of magnitude in the total energy accuracy implies approximately a factor 3 in  $s$ . An additional difficulty in dealing with conducting systems is related to the determination of the Fermi energy, which needs to be known accurately enough for the integration over all populated states, as in the reconstruction of the density matrix (see Eq. [38]). In these cases, band energy is interpolated at points of a finer grid (Gilat's grid<sup>101,102</sup>) and the new approximated values are used to integrate the number of electrons per cell. In spite of the grid size, this calculation is still easily manageable from the computational point of view because Be is a light atom and the unit cell is small and symmetric. However, a more complicated case could become excessively time consuming. A technique of smearing<sup>103</sup> the Fermi surface can be helpful in making convergence faster in these cases, where the sharp cutoff in occupancy at  $E_F$  would otherwise cause unphysical oscillations in the charge density.

No data are reported for  $s < 4$  in Table 1 for a reason. This reason is connected with the tight correlation of the sets of  $\mathbf{g}$  direct lattice vectors and  $\mathbf{k}$  reciprocal space points selected in the calculation, when using a local basis set. Iterative Fourier transforms of matrices from direct to reciprocal space, like in Eq. [36], and vice versa (Eq. [38]), are the price to be paid for the already mentioned advantage of determining the extent of the interparticle interactions to be evaluated in direct space on the basis of simple criteria of distance. Consequently, the sets of the selected  $\mathbf{g}$  vectors and  $\mathbf{k}$  points must be well balanced. The energy values reported in Table 1 were all obtained for a particular set of  $\mathbf{g}$  vectors, corresponding to the selection of those AOs in the lattice with an overlap of at least  $10^{-6}$  with the AOs in the 0-cell. This process determines the  $\mathbf{g}$  vectors for which  $F^{\mathbf{g}}$ ,  $S^{\mathbf{g}}$ , and the  $P^{\mathbf{g}}$  matrices (Eqs. [34] and [38]) need to be calculated, and if the number of  $\mathbf{k}$  points included in the calculation is too small compared with the number of the direct lattice vectors, the determination of the matrix elements is poor and numerical instabilities occur.

The number of  $\mathbf{k}$  points required to reach a given accuracy for total energy decreases when the unit cell is larger than the ones considered so far. In fact, the adjective "reciprocal" before "space" qualifies the relation of inverse proportionality between direct and reciprocal space (Eq. [3]), so that the bigger a unit cell in real space, the smaller the volume of the corresponding cell in reciprocal space. In those cases where the volume of the first Brillouin zone is small, it is sufficient to solve Schrödinger's equation only at a few  $\mathbf{k}$  points. To illustrate this point, we consider what happens when we repeat the calculation for magnesium oxide with nonprimitive cells. In particular, we refer to unit cells with volumes of 4, 16, and 64 times the primitive cell

**Table 2** Magnesium Oxide Total Energy per Pair of Mg and O Ions as Determined with Unit Cells Containing  $n_{\text{MgO}}$  Ion Pairs Corresponding to Grids with Different Values of the Shrinking Factor  $s$ 

$n_s$	$n_{\text{MgO}} = 1$	$n_{\text{MgO}} = 4$	$n_{\text{MgO}} = 16$	$n_{\text{MgO}} = 64$
1	—	—	—	-274.67867185
2	—	—	-274.67996667	-274.67996887
4	-274.67867182	-274.67996888	-274.67996886	-274.67996887
8	-274.67996884	-274.67996885	-274.67996886	-274.67996887

volume. Although only one pair of ions exists in the MgO primitive cell, whereas 64 of these pairs (128 ions) are in the biggest unit cell considered in Table 2, we are modeling the same crystal in all cases and the total energy per MgO pair is invariant to the size of the unit cell, provided the method used is numerically stable and accurate. It is apparent that similar levels of accuracy in the total energy value can be obtained with increasingly poorer grids as the direct lattice unit cell volume increases. In particular, when  $n_{\text{MgO}} = 64$ , using one  $\mathbf{k}$  point (the  $\Gamma$  point) is acceptable.

Artificially increasing the unit cell volume to profit from the reduced number of necessary diagonalizations involved in solving Eq. [25] is not a good strategy because of the many more interactions and, consequently, more one- and two-electron integrals to be computed. Nevertheless, this point is indeed relevant in those cases where the size of the primitive cell of a crystal is considerable and the legitimacy for considering a rare grid is a real advantage. For example, in the case of faujasite, a zeolite mineral containing 144 atoms in the primitive cell, 3  $\mathbf{k}$  points in the first Brillouin zone are sufficient to get a total energy value affected by an error of  $3 \cdot 10^{-9}$  hartree (the error is just one order of magnitude higher when considering the  $\Gamma$  point alone).

## Use of Symmetry in Reciprocal Space

The relative cost of ab initio calculations depends on many variables, such as the Hamiltonian, basis set, accuracy requirement, size, and density of the system (see Appendix 2). The Fock or KS matrix diagonalization step during the solution of Eq. [25] can become the calculation bottleneck with a large basis set, when, for example, more than 1000 basis functions are used. Such a number of functions may correspond to about 100 atoms per cell, when a local basis set is used, but this is the usual size of plane wave calculations, even with a small unit cell. As many crystalline systems are highly symmetric, taking advantage of symmetry is, therefore, important for reducing computational time.

Symmetry properties can be used both in the direct and in the reciprocal space, for example, to form matrices in direct space, such as  $\mathbf{F}^{\mathbf{g}}$  and  $\mathbf{P}^{\mathbf{g}}$ , or to diagonalize  $\mathbf{F}(\mathbf{k})$  more efficiently. The application of symmetry to direct space

matrices is discussed, for example, in Dovesi<sup>104</sup> and will not be reconsidered here. Instead, we briefly illustrate the use of symmetry properties in reciprocal space, because it is less widely known.

The application of a point symmetry operator of the space group to a given point  $\mathbf{k}$  in reciprocal space has two possible consequences:

1.  $\mathbf{k}$  is moved to another equivalent point  $\mathbf{k}'$ .
2.  $\mathbf{k}$  is not moved.

By accounting for the totally symmetric character of the Hamiltonian, it can be demonstrated that in case 1, the eigenvalues of  $F(\mathbf{k})$  and  $F(\mathbf{k}')$  coincide, and the eigenvectors of  $F(\mathbf{k}')$  are directly obtained from the corresponding eigenvectors of  $F(\mathbf{k})$  by the action of that symmetry operator. Moreover,  $F(\mathbf{k})$  and  $F(-\mathbf{k})$  are always related by symmetry, even if inversion is not present in the space group, owing to the so-called “time reversal symmetry.”<sup>105</sup> In fact, by taking the complex conjugate of Eq. [25] and considering the structure of the COs (Eq. [24]), it is seen that the eigenvalues are the same at  $\mathbf{k}$  and  $-\mathbf{k}$ , because they are real, and the eigenvectors are the complex conjugate of one another. On the basis of these considerations, it is, therefore, recommended that Eq. [25] be solved only at points belonging to the asymmetric part of the first Brillouin zone (the minimal set of symmetry unrelated  $\mathbf{k}$  points), as was done in previous examples. Indeed, the use of this kind of symmetry is so easy that it is probably implemented in all periodic codes. In Table 1, we can appreciate how far calculations benefit from symmetry type 1 by comparing  $n_s$ , the actual total number of  $\mathbf{k}$  points where  $F(\mathbf{k})$  was diagonalized, with  $s^3$ , the total number of  $\mathbf{k}$  points in the entire first Brillouin zone: The ratio of  $n_s$  to  $s$  for cubic systems like MgO and Si, when  $s = 8$ , is close to 1/20th (1/18.2).

As a second step, symmetry type 2 can be applied to the set of the  $n_s$   $\mathbf{k}$  points, which allows one to further reduce the Fock matrix into a block-diagonal form. By transforming the basis set into an equivalent set of symmetry adapted basis functions, every block of the transformed matrix in Figure 4, which corresponds to one particular point  $\mathbf{k}_j$ , reveals, in turn, a block-diagonal structure, for example, of the kind depicted in Figure 20.

Each block in the matrix on the right corresponds to a different irreducible representation (IR) of the so-called “little co-group,”<sup>106</sup> and the number



Figure 20 Block-factorization of the Fock matrix corresponding to a specific  $\mathbf{k}$  point,  $\mathbf{k}_j$ .

of these blocks varies with the symmetry invariance properties of  $\mathbf{k}_f$ . As a matter of fact, the application of this kind of symmetry properties to the solution of Eq. [25] is not trivial,<sup>107,108</sup> but the gain in computational efficiency is dramatic in the case of highly symmetric systems with large unit cells, when diagonalization dominates the calculation.

As an example, we consider pyrope, a garnet that crystallizes according to the cubic space group  $Ia\bar{3}d$ , with four  $\text{Mg}_3\text{Al}_2\text{Si}_3\text{O}_{12}$  formula units per cell consisting of 80 atoms. Shrinking factor  $s = 2$  has already been indicated as convenient for such a large unit cell, and symmetry type 1 reduces  $n_s$ , the number of  $\mathbf{k}$  points to be accounted for, from 8 to just 3:  $\Gamma$ ,  $H(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$ , and  $N(\frac{1}{2}, 0, 0)$ . The eigenvectors at these points are all real. If we use a local basis set consisting of  $n_f = 1440$  AOs, which corresponds to the choice of a split-valence basis set plus polarization functions for every atom in the unit cell, the size of the matrices in Eq. [25] would be  $1440 \times 1440$ . But because of symmetry type 2, they can be decomposed into lower size independent blocks as reported in Table 3.

$\Gamma$  is the most symmetric  $\mathbf{k}$  point in the reciprocal space, and all irreducible representations of the  $O_h$  point group are contained in the matrices. This situation is the most favorable from the point of view of computational efficiency. In fact, the size of the largest block is  $n_f^{T_{2g}} = 90$ . Moreover, knowing just one of the eigenvectors for every IR with  $d_{IR} > 1$  is enough to obtain the other  $(d_{IR} - 1)$  eigenvectors simply by application of the corresponding transfer operators. Thus, only one diagonalization is necessary for each IR. This property is particularly effective when  $d_{IR}$  is large, and in this respect, the appearance of IR with  $d_{IR}$  as large as 6 (see IR  $H$  at point  $H$  in Table 3) is an important peculiar feature of symmetry in reciprocal space, which does not occur in direct space. The application of symmetry type 2 has a smaller

**Table 3** Example of Decomposition of Fock, Overlap, and Eigenvector Matrices of Pyrope (with a Local Basis Set of  $n_f = 1440$  AOs) into Independent Irreducible Blocks at  $\mathbf{k}$  Points  $\Gamma$ ,  $H$ , and  $N$  in the First Brillouin Zone

$\Gamma$			$H$			$N$		
<i>IR</i>	$d_{IR}$	$n_f^{IR}$	<i>IR</i>	$d_{IR}$	$n_f^{IR}$	<i>IR</i>	$d_{IR}$	$n_f^{IR}$
$A_{1g}$	1	34	$E$	2	60	$E_1$	2	358
$A_{2g}$	1	32	$F$	2	114	$E_2$	2	362
$E_g$	2	60	$H$	6	182			
$T_{1g}$	3	88						
$T_{2g}$	3	90						
$A_{1u}$	1	32						
$A_{2u}$	1	30						
$E_u$	2	62						
$T_{1u}$	3	88						
$T_{2u}$	3	88						

$n_f^{IR}$  is the size of the block belonging to one row of the irreducible representation  $IR$  with dimension  $d_{IR}$ .



impact at N, where the decomposition results in two different blocks only. Nevertheless, even in this case, the size of the bigger matrix is  $n_f^{E_2} = 362$ , i.e., about one fourth of the original matrix.

---

## TOTAL ENERGY, ENERGY DIFFERENCES, AND DERIVATIVES

In the previous section, we have examined the dependence of the total energy of a system on the number of **k** points and have illustrated a convenient use of symmetry for reducing the computational effort. Many other elements influence the accuracy of the method, which for clarity and conciseness, can be grouped into three categories:

1. The Hamiltonian
2. The adopted basis set
3. The computational scheme (the implementation of the basic equations in a specific code).

We will further analyze the problem of accuracy in the calculation of the total energy and its derivatives through a few examples, where we must take into account that, from the physical or chemical point of view, we are never interested in the total energy of a system as such, but rather in energy differences, that might be as small as a few kcal/mol. It is with respect to this scale of energy that the overall accuracy of a calculation must be verified. In fact, systematic improvement of algorithms in molecular codes during the past 40 years now allows an accuracy of 1 kcal/mol for thermochemical data,<sup>109</sup> which is still far from being attained in solid state chemistry, although attention to the quantitative aspects of the calculation is increasing rapidly.

Why is the total energy such an important observable? Reactivity of a surface, formation of a defect, and structural modification of a material to enhance some of its properties can all be discussed with reference to the total energy easily and rigorously. Moreover, knowledge of the total energy derivatives is also informative concerning the equilibrium properties, lattice dynamics, and the response of materials to perturbations. Table 4 lists some of the many energy differences that are relevant in solid state chemistry. Most of them imply a comparison of systems with different periodicity, therefore, it is important that zero- (atoms or molecules), one- (polymeric chains), two- (slabs), and three-dimensional (bulk) systems all be handled in a consistent way, as concerning not only the method used but also the computational conditions, so that the comparability of the total energy be guaranteed to a high degree. As the systems involved in the comparison must be at equilibrium, finding the equilibrium geometry is preliminary to any energy difference calculation and selection of an efficient technique for geometry optimization is another requirement in this kind of calculation. The CRYSTAL code satisfies

**Table 4** List of Some Relevant Energy Differences in Solid State Chemistry

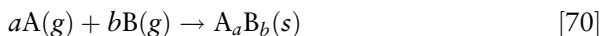
Computed energy	System 1	System 2	Example
Cohesive	Bulk	Atoms	Ionic, covalent crystals
Interaction	Bulk	Molecules	Molecular crystals
Relative stability	Bulk	Bulk	Polymorphism
Super-exchange	AFM bulk	FM bulk	NiO magnetic phases
Solid state reaction	Bulk	Bulk	$\text{MgO} + \text{Al}_2\text{O}_3 \rightarrow \text{MgAl}_2\text{O}_4$
Surface formation	Bulk	Slab	MgO(100)
Surface stability	Slab	Slab	MgO(100) vs MgO(110)
Adsorption	Slab + molecule	Slab, molecule	CO on MgO(100)
Adsorption	Bulk + molecule	Bulk, molecule	NH <sub>3</sub> in acidic zeolites
Interface	Slab	Slab	MgO monolayer on Ag(100)
Substitution defect	Bulk + defect	Bulk, atoms	C in Si

both these requirements. Four examples of energy calculations involving bulk systems with CRYSTAL are discussed in this section. Other cases will be presented in the sections devoted to surfaces and defects.

## Cohesive Energy

### *Ionic, Covalent, and Metallic Crystals*

The *cohesive energy* is the energy necessary to dissociate a solid into separated entities, generally the atoms. Following this definition, the cohesive energy ( $\Delta E$ ) of a crystalline compound with formula unit  $\text{A}_a\text{B}_b$ , for example, is associated with the reaction of formation of  $\text{A}_a\text{B}_b$  in the solid state from the noninteracting atoms A and B in the gas phase



and is computed as

$$\Delta E = aE^{\text{atom}}(\text{A}) + bE^{\text{atom}}(\text{B}) - E^{\text{bulk}}(\text{A}_a\text{B}_b) \quad [71]$$

Hence,  $\Delta E$  is positive for any thermodynamically stable crystal. This definition of cohesive energy is not unique. As alternative definitions, reference can be made to the ions in the case of ionic crystals or to the molecules for molecular crystals. Another expression, lattice energy, is also in use, either as a synonym of cohesive energy or to denote the energy difference relative to the ions or molecules, thereby distinguishing it from the cohesive energy referred to the atoms. Thus, some care must always be taken when analyzing data in the literature in regard to nomenclature. We will use the following symbols:  $\Delta E^{\text{atoms}}$ ,  $\Delta E^{\text{ions}}$ , and  $\Delta E^{\text{mol}}$ .

As occurs in the calculation of molecular binding energies, the expression for  $\Delta E$  implies basis sets that are complete. However, partial basis set

incompleteness is commonly accepted in molecular calculations, provided all terms involved in the expression of the binding energy are computed with the same basis sets. This procedure becomes critical in most cases in the solid state. In fact, convenient basis sets for atoms are generally overcomplete in crystals, where close packing makes atomic function tails unnecessary in the description of the crystalline wave function. Moreover, the use of diffuse AOs with a periodic potential is normally to be avoided, for it may introduce artificially unwanted numerical instability in the calculation of the eigenvalues and eigenvectors of the Hamiltonian. Therefore, in the evaluation of the cohesive energy, partially different basis sets, as concerns the description of the valence region, need to be used for the atoms or ions and the bulk, in most cases.

We will illustrate how these problems can be handled with a simple example: the cohesive energy calculation of bulk NaCl with the HF approximation. Following Pople's scheme for the construction of a basis set, we adopt the convention that every set of ( $p_x$ ,  $p_y$ ,  $p_z$ ) orbitals share the same gaussian function with one s-type AO and these four AOs collectively form an sp shell. An *all-electron* basis set optimized for the bulk is available (see Prencipe et al.<sup>110</sup>), which consists of one core s-type AO plus three and four sp shells with the origin at Na and Cl, respectively. The exponents of the outer Gaussian functions of both Na and Cl are reported in Table 5 (Case 1). The system is fully ionic, in accordance with the usual representation of NaCl as a salt consisting of  $\text{Na}^+$  and  $\text{Cl}^-$  ions, which explains why the exponent of the outermost Gaussian of the cation, i.e.,  $\alpha_b(\text{sp})$  in the first row of Table 5, is about three times larger than at the anion (the valence orbitals of  $\text{Na}^+$  are basically empty). When the same basis is used to compute the atomic and ionic energies, an extremely large value of  $\Delta E^{\text{atoms}}$  results (Case 1 in Table 6). This is clear evidence of the poor performance of that basis set in the calculation of the atomic energies, which are severely underestimated (Case 1 in Table 7), particularly for  $E(\text{Na})$ . Indeed, readjustment of  $\alpha_b(\text{sp})$  for the atoms and ions in the gas phase is necessary to improve  $\Delta E^{\text{atoms}}$  (Case 2). Obviously, the best

**Table 5** Exponents of the Most Diffuse Shells (in  $\text{Bohr}^{-2}$ ) of Na and Cl as Optimized for Bulk NaCl and for the Isolated Atoms and Ions at the HF Level

Case	System	Na				Cl			
		$\alpha_a(\text{sp})$	$\alpha_b(\text{sp})$	$\alpha_c(\text{sp})$	$\alpha_a(\text{d})$	$\alpha_a(\text{sp})$	$\alpha_b(\text{sp})$	$\alpha_c(\text{sp})$	$\alpha_a(\text{d})$
1	bulk	0.578	0.323			0.320	0.125		
2	atom	0.497	0.042			0.315	0.119		
	ion	0.542	0.229			0.294	0.090		
3	atom	0.509	0.089	0.030		0.297	0.248	0.116	
	ion	0.539	0.204	0.111		0.329	0.151	0.059	
4	bulk	0.578	0.323	0.125		0.320	0.125		
5	bulk	0.578	0.323	0.125	0.400	0.320	0.125		0.400

**Table 6** HF Computed Cohesive Energy for NaCl (in kJ/mol) with Respect to the Isolated Atoms and Ions

Case	$\Delta E^{\text{atoms}}$	$\Delta E^{\text{ions}}$
1	1499.8 (+132.8)	761.0 (−3.8)
2	521.0 (−19.1)	744.3 (−5.9)
3	512.0 (−20.5)	738.3 (−6.7)
4 <sup>a</sup>	513.3 (−20.3)	739.6 (−6.5)
5 <sup>a</sup>	515.8 (−19.9)	742.1 (−6.2)
Exp.	644.2	791.3

Percentage underestimation/overestimation of the experimental cohesive energies in parentheses. Cases 1–5 correspond to different basis sets in Table 5.

<sup>a</sup>Cohesive energy computed with respect to atomic and ionic energy of Case 3.

**Table 7** HF Total Energy (in Hartree) of the Isolated Atoms and Ions Obtained with the Basis Set of Cases 1–3 Reported in Table 5

	Case 1	Case 2	Case 3
Na (g)	−161.475779	−161.848514	−161.850828
Na <sup>+</sup> (g)	−161.669712	−161.670010	−161.670022
Cl (g)	−459.449667	−459.449740	−459.450841
Cl <sup>−</sup> (g)	−459.537121	−459.543201	−459.545457

exponents for the isolated atoms are smaller than in the bulk (Case 2 in Table 5) and the corresponding change in  $\Delta E^{\text{atoms}}$  (Case 2 in Table 6) is dramatic, at the same time being far from negligible (about 2%) in  $\Delta E^{\text{ions}}$ . The addition of one extra sp shell to the basis sets of the isolated atoms and ions (Case 3), followed by re-optimization of the outer valence shell exponents, still decreases both  $\Delta E^{\text{atoms}}$  and  $\Delta E^{\text{ions}}$  by about 1%. Actually, Case 3 corresponds to the hypothesis that extra functions added to the bulk basis set would have no effect on the bulk total energy, which is almost true, as results from a comparison with Case 4.

Also the bulk basis set can still be improved by the addition of polarization functions, which do not contribute to the energy of atoms and ions in the gas phase because of AO orthogonality, but they can be important in the expansion of the bulk wave function. However, in this particular case, the contribution from polarization functions to the bulk total energy (Case 5) is minimal, because of high symmetry and the nearly spherical shape of the closed-shell ions.

Table 6 demonstrates that separate optimizations of the basis sets for the bulk and the isolated atoms or ions is mandatory for obtaining the cohesive energies of ionic compounds. As a general rule, *variationally equivalent* basis sets are to be used for the bulk and the atoms or ions, rather than *equal* basis sets.

Apparently, further improvement of the basis sets would hardly affect the values of  $\Delta E$  in Table 6, which are much smaller than the corresponding

**Table 8** Cohesive Energies (in kJ/mol) of LiF, NaCl, KBr, MgO, Si, and Be Computed with Respect to the Atoms with Different Hamiltonians

	HF	LDA	PW91	B3LYP	CCSD(T)	Exp.
LiF	656.6 (−23.7)	958.1 (+11.3)	850.7 (−1.2)	825.5 (−4.1)	845.4 (−1.8)	861.1
NaCl	512.0 (−20.5)	685.3 (+6.4)	614.7 (−4.6)	589.9 (−8.4)	627.5 (−2.6)	644.2
KBr	477.1 (−21.2)	611.9 (+1.1)	554.8 (−8.4)	533.9 (−11.8)	579.7 (−4.2)	605.4
MgO	699.5 (−29.3)	1118.7 (+13.0)	966.4 (−2.4)	908.0 (−8.3)	969.3 (−2.1)	989.8
Si	643.9 (−29.7)	1101.5 (+20.2)	943.3 (+2.9)	850.9 (−7.2)	852.8 (−7.0)	916.6
Be	180.1 (−43.6)	372.7 (+16.8)	319.2 ( 0.0)	259.0 (−18.9)	—	319.2

Percentage difference between calculated and experimental data in parentheses.

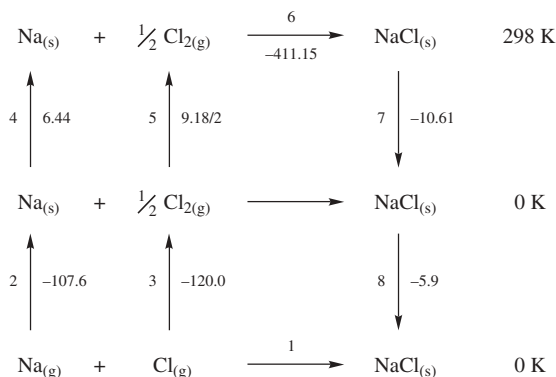
experimental values. Indeed, HF underestimates  $\Delta E^{\text{atoms}}$  of NaCl by about 20% and  $\Delta E^{\text{ions}}$  by about 6%. The origin of such a poor approximation is to be found in intra-ionic and interionic correlation effects, both disregarded at this level of theory, although the better agreement for  $\Delta E^{\text{ions}}$  is caused by cancellation of intra-ionic correlation effects when subtracting the energy of the ions from that of the bulk.

The underestimation of the cohesive energy is a general feature of HF, as can be seen in Table 8, where  $\Delta E^{\text{atoms}}$  is reported for a series of different simple crystalline compounds, which includes three alkali halides with increasing ion size, an ionic oxide with stronger electrostatic interactions caused by divalent ions (MgO), a covalent system (Si), and a metal (Be).

The error in the HF cohesive energy varies between −20% and −45% for this series. In accordance with the interpretation of electron correlation as the source of such an error, DFT calculations of  $\Delta E$  allow one to recover part of the contributions that are disregarded with HF, at about the same computational cost. LDA tends even to overestimate the cohesive energy, whereas GGA and B3LYP results are closer to the experimental measurements. At any rate, the performance of none of the Hamiltonians used is fully satisfactory, and the correct answer is always somewhere in between the two extremes represented by HF and LDA, but LDA results are generally improved when gradient corrections are included.

Properly correlated wave functions are obviously expected to perform better than HF and DFT methods. A systematic post-HF investigation of correlation effects in crystalline compounds has been obtained with the “incremental scheme” of Stoll<sup>71,72</sup> at both CCSD and CCSD(T) levels.<sup>73–75</sup> The CCSD(T) cohesive energies reported in Table 8, in fact, agree on average better with the experimental data than the one-electron results.

Another aspect of the comparison between the calculated and experimental cohesive energy is important to recall here. It is related to the difference between the definition of cohesive energy and the crystal formation energy that is reported in thermodynamic tables, the main point probably being that quantum mechanical calculations refer to the static limit ( $T = 0$  K and frozen nuclei), whereas experiments refer to some finite temperature. In fact, the



**Figure 21** Born–Haber cycle used to correct the experimental formation energy of NaCl from 298 K to the static limit.

comparison is never straightforward, and the original experimental datum is linked to the calculated cohesive energy values through a chain of thermodynamic transformations. As an example, the Born–Haber cycle for the formation of NaCl is reported in Figure 21. The calculated cohesive energies reported in Tables 4 and 5 refer to step 1, i.e., to the formation of the crystal at 0 K from the pure atoms in the gas phase, whereas the tabulated experimental datum refers to the standard formation reaction of NaCl at room temperature, step 6 ( $\Delta H_{298}^0 = -411.15$  kJ/mol) of the cycle. Hence, the experimental value reported in Table 8 corresponds, instead, to the path consisting of steps 2–8 of the Born–Haber cycle. These steps are as follows: (2) condensation of Na atoms from the gas phase (the inverse of the sublimation enthalpy); (3) formation of  $\text{Cl}_2$  molecules (the inverse of the enthalpy of dissociation); (4) heating solid Na atoms to room temperature; (5) heating  $\text{Cl}_2$  molecules to room temperature; (6) formation of crystalline NaCl from the elements in their standard states at room temperature; (7) cooling NaCl to 0 K. These enthalpy differences are tabulated<sup>111</sup> (ideal behavior of  $\text{Cl}_2$  is assumed in step 5). Step 8 refers to the zero point energy  $\varepsilon_0$ , which is usually not available experimentally. The Debye model (see, for instance, p. 100 of Hill<sup>112</sup>) relates  $\varepsilon_0$  to Debye’s temperature  $\Theta_D$ , through the following equation:

$$\varepsilon_0 = n9/8k_B\Theta_D \quad [72]$$

where  $k_B$  is the Boltzmann’s constant and  $n$  is the number of atoms in the unit cell. Debye temperatures for alkali halides can be found in Ashcroft and Mermin<sup>105</sup>.  $\Theta_D$  is 321 K for NaCl. In conclusion, the experimental *static* cohesive energy of NaCl reported in Tables 6 and 8 originates from the following sum:

$$\begin{aligned}
 -\Delta E_{\text{static}}^{\text{atoms}} &= -107.6 - 120.0 + 6.44 + \frac{9.18}{2} - 411.15 - 10.61 - 5.9 \\
 &= -644.2 \text{ kJ/mol}
 \end{aligned} \quad [73]$$

The experimental lattice energy at the static limit can also be evaluated with respect to the isolated ions, and it is easily obtained by including the ionization energy of the alkali metal (Na) and the electron affinity of the halogen (Cl). According to *The Handbook of Chemistry and Physics*,<sup>111</sup> the ionization energy of sodium is  $-495.8$  kJ/mol and the electron affinity of chlorine is  $348.7$  kJ/mol. Then, the experimental cohesive energy from ions is  $791.3$  kJ/mol at the static limit.

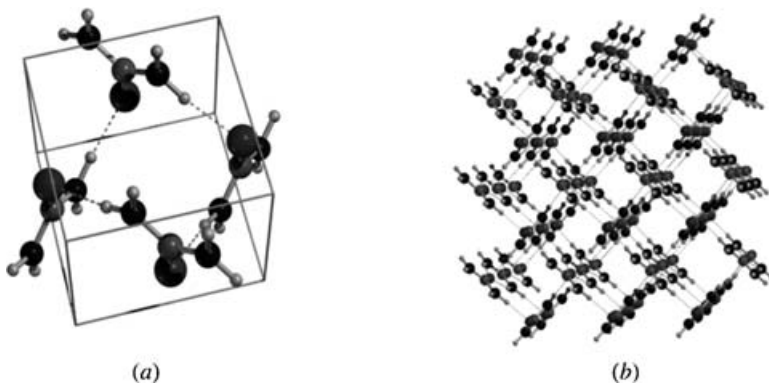
### Molecular Crystals

The cohesive energy of molecular crystals is usually computed with reference to the molecules in the gas phase, like in a sublimation process, so that calculated values of  $\Delta E^{\text{mol}}$  can be compared with experimental sublimation energies.  $\Delta E$  can be decomposed into two terms:

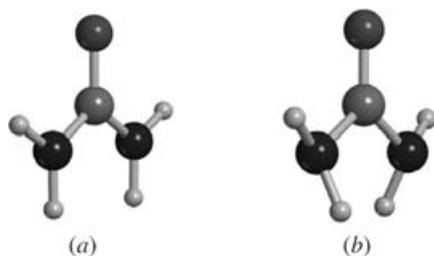
$$\Delta E^{\text{mol}} = -\Delta E_{\text{condensation}} - \Delta E_{\text{conformation}} \quad [74]$$

with the first term referring to the condensation of molecules from the gas phase, but with the same conformation as in the crystal, and the latter to the difference between the energies of the isolated molecules in the bulk and in the gas phase conformations.  $\Delta E_{\text{conformation}}$  is negligible in the case of rigid molecules, but it can be significant for floppy molecules.

Let us consider urea as an example. Crystalline urea is tetragonal, with two molecules in the unit cell (Figure 22a). The molecules are linked to each other through hydrogen bonds to form infinite planar tapes (Figure 22b), which are mutually orthogonal, the cohesion among them being provided by hydrogen bonds. Every oxygen atom is involved in four nearly equivalent hydrogen bonds, two within the tape and the other two linking neighboring tapes. Molecules in nearest neighboring tapes are oriented along opposite directions; this provides additional stabilization through dipole-dipole



**Figure 22** (a) Crystalline urea unit cell. (b) Arrangement of urea molecules in the crystalline structure.



**Figure 23** Urea molecular structures: (a)  $C_{2v}$  conformation as found in the crystalline structure; (b)  $C_2$  *anti* conformation as found in the gas phase.

interactions, and it annihilates the total dipole of the unit cell, as is always the case in molecular crystals

A molecule in bulk has  $C_{2v}$  symmetry, whereas the most stable structure in the gas phase corresponds to a  $C_2$ -symmetric *anti* conformation (see Figure 23). The  $C_{2v}$  geometry is a second-order saddle point on the potential energy surface.<sup>113</sup> Therefore,  $\Delta E_{\text{conformation}}$  between conformers must be taken into account when computing the cohesive energy.

The cohesive energy of urea with respect to the molecules in the gas phase is calculated as

$$\Delta E^{\text{mol}} = \frac{2E(\text{molecule}) - E(\text{bulk})}{2} \quad [75]$$

In this form,  $\Delta E^{\text{mol}}$  is the cohesive energy per molecule and the factor 2 in the formula is caused by the presence of two molecules in the unit cell. The cohesive energies computed with different Hamiltonians and Pople's 6-31G(d,p) basis set are reported in Table 9. At variance with the case of ionic crystals, molecular-devised basis sets can generally be used for molecular crystals as such, without any exponent reoptimization. As shown in Table 9,  $\Delta E_{\text{conformation}}$  accounts for 5–8 kJ/mol. Table 9 also shows that DFT-based

**Table 9** Cohesive Energy ( $\Delta E^{\text{mol}}$ ) per Molecule of Crystalline Urea Calculated with the Experimental Lattice Parameters (in kJ/mol)

Hamiltonian	$\Delta E_{\text{condensation}}$	$\Delta E_{\text{conformation}}$	$\Delta E^{\text{mol}}$	$\Delta E^{\text{CP}}$
HF	−80.2	4.9	75.3	54.4
LDA	−177.0	4.9	172.1	135.2
PW91	−124.5	7.5	117.0	79.1
B3LYP	−105.2	6.7	98.5	63.6

$\Delta E^{\text{CP}}$  includes the correction for the basis set superposition error, estimated via the counterpoise method (CP).

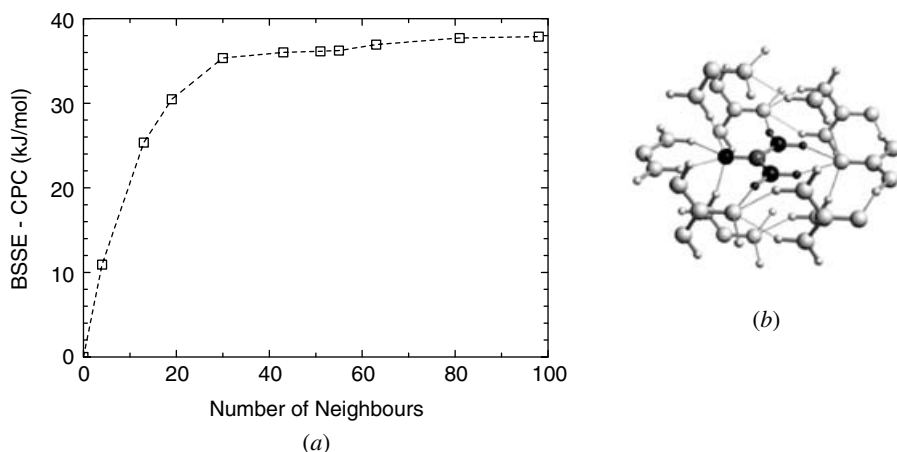


methods give larger  $\Delta E^{\text{mol}}$  than HF, as the correlation energy, somehow included in DFT, makes the crystalline structure more stable.

When the local basis set is incomplete at some extent, the basis set superposition error (BSSE)<sup>114,115</sup> affects the cohesive energy more extensively than the binding energy in molecules. In fact, in calculating the wave function and total energy of a molecular crystal with a finite basis set, the description of molecule A in the crystal will be improved by the variational freedom provided by the functions of the adjacent molecule B, and vice versa. As a consequence, the energy content of A and B in the crystalline environment turns out to be overestimated, as if an extra binding occurred between A and B. This error is commonly corrected via the counterpoise (CP) method, as proposed by Boys and Bernardi,<sup>116</sup> by supplementing the basis set of an isolated molecule with the functions of an increasing number of atoms (ghost atoms) belonging to the surrounding array of molecules that would be present in the crystal. An introductory tutorial about the theory and the practice of BSSE has been published in this book series.<sup>115</sup> An interesting discussion of the BSSE problem in molecular crystals can be found in a recent paper by Spackman and Mitchell.<sup>117</sup>

The dependence of the BSSE on the number of the accounted neighboring atoms in crystalline urea is shown in Figure 24a. The correction converges to a limiting value of about 38.0 kJ/mol, with the largest calculation including 98 ghost atoms. Nevertheless, 30 neighbors are enough to reach a value of 35.3 kJ/mol that represents about 93% of the entire CP correction.

Thus, proper consideration of the BSSE correction decreases  $\Delta E^{\text{mol}}$  considerably at all levels of theory considered here. On the contrary, basis set



**Figure 24** (a) Dependence of the BSSE on the number of neighboring atoms included in the CP correction for crystalline urea, with the LDA approximation (SVWN) and a 6-31G(d,p) basis set. (b) Urea molecule surrounded by a star of 63 neighboring ghost atoms.

improvements, which have been shown to be so relevant in the calculation of the cohesive energy of ionic systems, and particularly with respect to the atoms, are much less important in computing  $\Delta E^{\text{mol}}$  of a molecular crystal like urea. The BSSE-corrected B3LYP cohesive energies obtained with two standard Pople's basis sets, namely 6-31G(d,p) and 6-311G(d,p), and a double- $\zeta$  plus polarization basis proposed by Thakkar et al.<sup>118</sup> are 63.6, 61.1 and 64.1 kJ/mol, respectively, thus showing an almost negligible dependence on the basis set used, although the BSSE correction decreases with the basis set size.

The experimental sublimation energy of crystalline urea is 90.0 kJ/mol<sup>119</sup> and 97.7 kJ/mol after correction for the zero point energy (ab initio estimate), which is more properly compared with the results from calculations. Unfortunately, none of the calculated  $\Delta E^{\text{CP}}$  in Table 9 compares well with this latter value, with the minimum found error being about 20% of the cohesive energy. However, the computed cohesive energies follow closely the trend obtained for hydrogen bonded molecular adducts (see, for instance, Civalleri et al.<sup>120</sup>), where LDA functionals, like SVWN, tend to greatly overestimate the interaction energy, whereas gradient-corrected and hybrid functionals represent definite improvements with respect to LDA. On the other hand, as thermal effects are expected to account only for a few kJ/mol, the difference with respect to experiment must then be traced back to other effects, probably to the lack of dispersion interactions, which are not taken into account<sup>121</sup> at one-electron levels of theory.

Other examples of applications of the CRYSTAL code to molecular crystals include ice polymorphs,<sup>122–124</sup> orthoboric acid,<sup>125</sup> vitamin C,<sup>126</sup> oxalic acid dihydrate,<sup>127</sup> and *p*-benzoquinone.<sup>128</sup>

## Polymorphism

Silica is of great interest in solid state chemistry, as it exists in many different crystalline forms, from high-density polytypes (e.g., quartz, cristobalite, tridymite) to low-density microporous all-silica zeotypes. Despite their enormous structural complexity, silica polymorphs show similar stability. From calorimetric measurements,<sup>129–132</sup> it is known that  $\alpha$ -quartz is the most stable polymorph at room temperature and pressure, with all other polymorphs being confined within a range of just 15 kJ/mol. This small range makes an accurate simulation of the relative stability of silica polymorphs a delicate but challenging task. In fact, quartz and all-silica polytypes have been the subject of many theoretical studies.<sup>133–135</sup> Here, we consider the relative stability of  $\beta$ -quartz ( $\beta$ -Q),  $\alpha$ -cristobalite ( $\alpha$ -C),  $\alpha$ -tridymite ( $\alpha$ -T), sodalite (SOD), chabazite (CHA), faujasite (FAU), and edingtonite (EDI) with respect to  $\alpha$ -quartz ( $\alpha$ -Q). Frameworks are shown in Figure 25.

In Table 10, CRYSTAL all-electron (AE) calculations at both HF and DFT level of theory (LDA and B3LYP) are compared with shell model results and ab initio PW calculations.<sup>135</sup> G(HF)<sup>136</sup> and G(B3LYP)<sup>137</sup> refer to

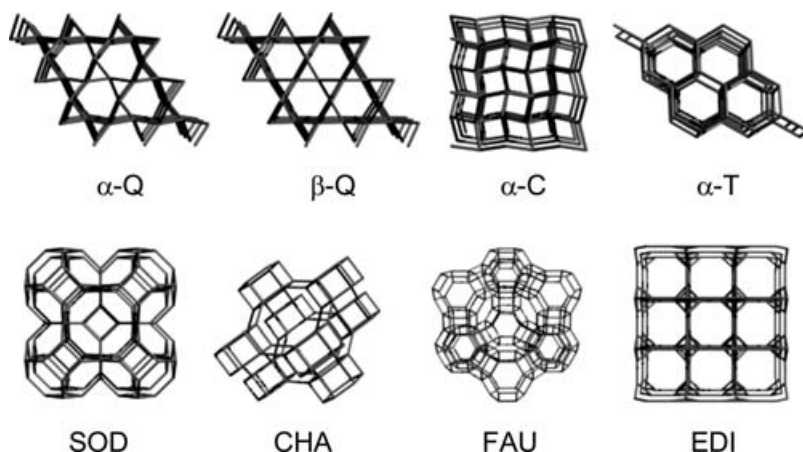


Figure 25 Frameworks of the studied all-silica polymorphs.

a semiclassical approach with model interatomic potentials that were fitted to ab initio calculations on molecular clusters simulating fragments of microporous all-silica frameworks.

The range and the order of stability computed with a local AE basis set are in good agreement with the experimental measurements available, particularly as concerns the B3LYP results (third row in Table 10). Also the G(B3LYP) parameterization of the semiclassical method provides results that are in reasonably good agreement with the experimental evidence. However, the ab initio approach is still to be preferred in the case of more complicated systems, such as Ti-substituted zeolites, for example, where the parameterization procedure may become critical.

Perhaps the most surprising feature in Table 10 is the large difference between AE and PW results, obtained at the same levels of theory. The origin of this inconsistency is probably caused by the different computational

Table 10 Relative Stability (in kJ/mol per  $\text{SiO}_2$  unit) of Silica Polymorphs with Respect to  $\alpha$ -Q

	$\beta$ -Q	$\alpha$ -C	$\alpha$ -T	SOD	CHA	FAU	EDI
AE-HF	3.1	0.0	1.7	4.6	6.3	7.6	11.6
AE-LDA		9.1	12.9				
AE-B3LYP		1.2	4.3	5.7	7.0	7.8	11.6
G(HF)	1.6	-3.8	-2.8	0.5	4.4	8.1	9.5
G(B3LYP)		1.4	7.7	7.7	8.6	9.8	13.0
PW-LDA	2.5	2.4	3.3				
PW-GGA	0.2	-3.1	-2.2				
Exp. ( $\Delta H^{298}$ )		2.8	3.2		11.4	13.6	

conditions used by scientists in the different implementations of the method, and especially in the different representations of the wavefunction. These results are a good indication of how much computational and methodological aspects can affect results in solid state calculations, where full standardization has not yet been achieved.

## Magnetic Phases

Another context in which simple energy differences provide useful information to chemists is magnetism. The magnetic properties of transition metal insulators such as  $\text{KMnF}_3$  perovskites ( $M = \text{Mn, Fe, Co, Ni, Cu}$ )<sup>138–140</sup>, simple MO oxides ( $M = \text{Ni, Mn}$ )<sup>140,141</sup>,  $\text{M}_2\text{O}_3$  sesquioxides ( $M = \text{Cr, Fe}$ )<sup>142,143</sup>, complex oxides such as  $\text{Mn}_3\text{O}_4$  (hausmannite),<sup>144</sup> rutile type fluorides  $\text{MF}_2$  ( $M = \text{Mn, Fe, Co, Ni, Cu}$ )<sup>145–147</sup> and high- $T_c$  superconductor parent compounds<sup>148</sup> have been investigated extensively.

As an example, we discuss the relative stability of some magnetic phases of  $\text{KMnF}_3$ , an ionic perovskite, with cubic lattice, where each Mn ion is at the center of a regular octahedron of fluorine ions. Different magnetic structures can be envisaged, as shown in Figure 26. The structure on the left represents a ferromagnetic (FM) phase with the spin at the Mn ions all parallel ( $S_z = 5/2$ ). In the middle, an antiferromagnetic (AFM) phase is shown, where the Mn ions, carrying five unpaired electrons each, are arranged in alternating stacked spin-up and spin-down (111) planes. In the AFM' structure on the right, each Mn atom is surrounded by two Mn nearest neighbors with opposite spin and four with the same spin, thus forming a sequence of alternating spin-up and spin-down Mn (001) planes. Because of the presence of alternating antiparallel spin Mn planes, the two different antiferromagnetic phases imply the use of double supercells.

Information about the basis sets and other computational details can be found in Dovesi et al.,<sup>138</sup> Harrison et al.,<sup>139</sup> and Mallia et al.<sup>140</sup> The hypothetical spin arrangement shown in Figure 26 is a schematic representation of a point charge lattice, where the unpaired electrons are fully assigned to the

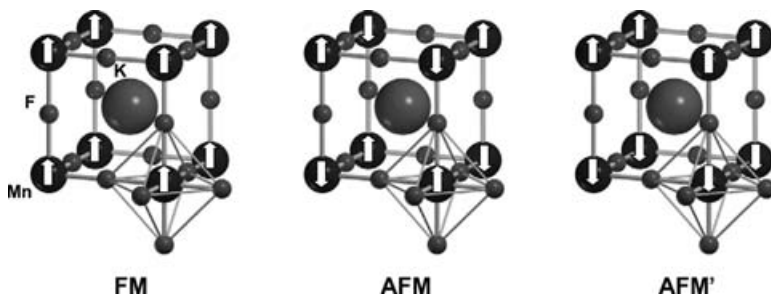


Figure 26 Magnetic phases of  $\text{KMnF}_3$ . Arrows denote spin-up and spin-down Mn ions.

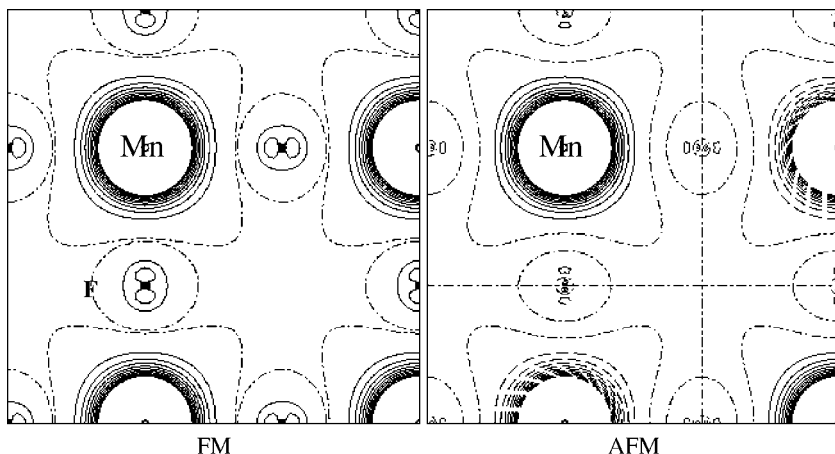
**Table 11** Net Atomic Charges ( $q$ ) and Spin Moments ( $\mu$ ) in  $\text{KMnF}_3$  Evaluated According to Mulliken Partition of Spin Density

	Mn		F		K	
	$q$	$\mu$	$q$	$\mu$	$q$	$\mu$
FM	1.77	4.94	-0.92	0.02	0.99	0.00
AFM	1.77	4.94	-0.92	0.02	0.99	0.00

Values of  $q$  and  $\mu$  are given in electrons.

transition metal ions. This picture is actually not far from quantum mechanical reality, at least at the unrestricted Hartree–Fock (UHF) level of theory, as documented in Table 11, in spite of the looser constraint the spin density distribution has to comply with in a variational calculation, i.e., that the total spin projection along one direction ( $\Sigma_z$ ) be assigned to a constant given value in every cell ( $\Sigma_z$  is 5/2 and 0 in the ferromagnetic and antiferromagnetic phases of  $\text{KMnF}_3$ , respectively). In particular, the ion net charges are close to their formal values (+2, -1, and +1 for Mn, F, and K) and the Mn spin moment is close to 5, with only small polarization of the F anion (see also Figure 27).

The different stability of FM and AFM phases results from the superexchange interaction<sup>138</sup> along the Mn–F–Mn path and is essentially caused by the different spin polarization of F in the two cases, as shown in Figure 27. Consequently, the corresponding energy differences are expected to be so small that they need to be determined with very high numerical accuracy. In



**Figure 27** UHF spin density map for the FM and AFM solutions of  $\text{KMnF}_3$  on the (001) plane through the Mn and F atoms. The separation between contiguous isodensity lines is 0.01 bohr; the function is truncated in the core region at  $\pm 0.1$  bohr. Continuous, dashed, and dot-dashed lines correspond to up-, down-, and zero-spin density, respectively.

**Table 12** Total Energy (in Hartree) of the FM and AFM Phases of  $\text{KMnF}_3$  as Obtained with Various Hamiltonians

Hamiltonian	$E_{\text{FM}}$	$E_{\text{AFM}}$	$\Delta E$	$J$
HF	-4095.286164	-4095.286758	1.56	2.50
LDA	-4089.532113	-4089.540105	21.00	33.64
B-LYP	-4101.523522	-4101.531254	20.30	32.56
PBE	-4100.020100	-4100.026538	16.90	27.10
B3LYP	-4101.056636	-4101.060127	9.17	14.70
Exp. <sup>#</sup>				7.30
				7.40

$\Delta E = E_{\text{FM}} - E_{\text{AFM}}$  in kJ/mol.  $J$  is the super-exchange coupling constant (in K).

<sup>#</sup>Experimental values from reference 149.

fact, the relative stability of the considered FM phase with respect to the anti-ferromagnetic structures is 0.780 and 0.260 kJ/mol for AFM and AFM', respectively, at the UHF level of theory. The results obtained with different Hamiltonians are compared in Table 12.

$\Delta E$  can be used to evaluate the superexchange coupling constant  $J$ , which is a measure of the superexchange interaction. A general introduction to the superexchange interaction can be found in Kahn<sup>150</sup> and Yosida.<sup>151</sup> According to the Ising model,  $\Delta E(\text{FM-AFM})$  and  $J$  are related through the following expression:

$$\Delta E = \frac{2zJS_z^2}{k} \quad [76]$$

where  $z$  corresponds to the number of nearest Mn transition metal ion neighbors with opposite spin,  $S_z$  is the formal  $S_z$  of Mn (5/2) and  $k = 120.27 \text{ K} \cdot \text{mol} \cdot \text{kJ}^{-1}$  is the conversion factor from kJ/mol to K. The values of  $J$  can be compared with the experimental observations, usually fitted to the same Ising model. The UHF value of  $J$  is about one third of the experimental value, whereas all DFT functionals provide values that are larger than the experiment<sup>149</sup> and, in the case of LDA,  $J$  is overestimated even by a factor of 4.5.

Spin contamination is the main source of error in the evaluation of  $J$ , but deviations from the Ising model may also account for part of it. However, despite the large disagreement between the calculated and experimental values of  $J$ , the prediction of the relative stability of different magnetic phases is correct. Moreover, investigation of the same properties with the other systems previously mentioned always reproduced phase stabilities correctly and  $J$  values were calculated approximately within the same error bar.

## Positional Isomorphous Phases

The relative stability of different cation sites in a zeolite framework is one more interesting example of energy difference calculation, in this case

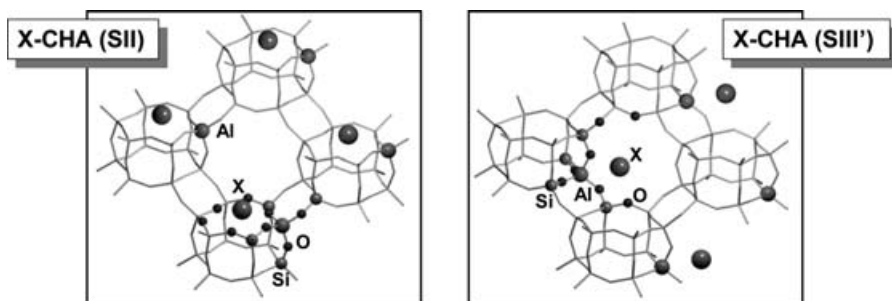


Figure 28 Cation sites in chabazite.

concerning isomorphous crystalline structures with the same composition but a positional difference.

The cation location in a high-silica Al-substituted chabazite<sup>152</sup> has been studied where two possible sites are to be considered (Figure 28) on the basis of the experimental evidence. The SII site corresponds to a cation at the top of a hexagonal prism, whereas for SIII', it is in the eight-member ring. The relative stability of these two sites was computed for a proton and three alkali metal ions: Li, Na, and K at both HF and B3LYP levels of calculation (see Civalleri et al.<sup>152</sup> for details).

In this case, the level of theory is almost irrelevant for evaluating the relative stability of the different sites (Table 13). The proton is preferably predicted at the SIII' site. Li and Na cations are more stable when coordinated on the hexagonal prism (SII), whereas K definitely prefers the eight-member ring site (SIII') for steric reasons.

## Energy Derivatives

The interest by computational chemists in the total energy originates from the possibility of comparing different systems or different phases of a given system. The information content of the total energy dependence on the crystal structure permits a wider analysis of the properties of a system. Some examples of observables related to first- and second-order derivatives are listed in Table 14. Beside the search for energy minima, investigating the energy hypersurface is, for example, also a means to go beyond the static

**Table 13** Relative Stability (in kJ/mol) of Site SIII' with Respect to the SII Site

	H	Li	Na	K
HF	-11.2	7.9	10.7	-18.3
B3LYP	-6.9	9.2	12.5	-18.1

**Table 14** First and Second Order Derivatives of the Total Energy

Differentiating variable	Total energy derivative	Observable
nuclear coordinate	$\left(\frac{\partial E}{\partial \mathbf{r}_i}\right)_T = 0$	equilibrium nuclear coordinates
	$\left(\frac{\partial^2 E}{\partial \mathbf{u}_i \partial \mathbf{u}_j}\right)_{\text{eq}} = k_{ij}$	force constants
lattice parameter	$\left(\frac{\partial E}{\partial \mathbf{a}_i}\right)_T = 0$	equilibrium unit cell
	$\left(\frac{\partial^2 E}{\partial \varepsilon_i \partial \varepsilon_j}\right)_{\text{eq}} = c_{ij}$	elastic tensor
unit cell volume	$\left(\frac{\partial E}{\partial V}\right)_S = -P$	internal pressure
	$\left(\frac{\partial^2 E}{\partial V^2}\right)_{\text{eq}} = B$	bulk modulus

$\mathbf{u}_i$  denotes a nuclear displacement from the equilibrium position,  $\varepsilon_i$  is a component of the strain tensor,  $\mathbf{a}_i$  is a lattice basis vector.

lattice model implicit in the adiabatic approximation and open to lattice dynamical and thermodynamical properties.

We first consider first-order energy derivatives. The calculation of the cohesive energy and of the relative stability of different compounds or phases relies on the hypothesis that the systems involved are all in their respective equilibrium geometries, and that these geometries are to be determined *ab initio*. Thus, the accuracy problem in a calculation of  $\Delta E$  (see Table 4) cannot be separated from the reliability of a geometry determined with the same method. The analytic calculation of energy gradients with respect to the cell parameters and nuclear coordinates is the most efficient method of finding minima on the total energy hypersurface. As an example, we report the equilibrium lattice parameters (Table 15) that were used to compute the cohesive energies reported in Table 8. The structure of all crystals sampled in the table is simple, with all atoms in a special position (see the subsection on the direct lattice), so that the lattice parameters are the only geometrical variables to be taken into account.

Table 15 shows that in most cases, the percentage deviation from the experimental lattice parameter is below 2%. In general, HF tends to overestimate lattice parameters in nonmetallic systems, whereas LDA shows the opposite trend. The alkali-halide series is more varied in results. For instance, the HF error for the lattice parameter of LiF, NaCl, and KBr increases from about 0% to 8%. This result is a consequence of the increasing importance of



**Table 15** Computed Lattice Parameters (in Å) of LiF, NaCl, KBr, MgO, Si, and Be with Different Hamiltonians

	HF	LDA	PW91	B3LYP	CCSD(T)	Exp.
LiF	4.02 (+0.7)	3.93 (−1.5)	4.09 (+2.5)	4.05 (+1.5)	3.99 ( 0.0)	3.99
NaCl	5.80 (+4.1)	5.50 (−1.3)	5.72 (+2.7)	5.73 (+2.9)	5.63 (+1.1)	5.57
KBr	7.05 (+8.0)	6.60 (+1.1)	6.92 (+6.0)	6.94 (+6.3)	6.65 (+1.8)	6.53
MgO	4.21 (+0.2)	4.18 (−0.4)	4.26 (+1.4)	4.24 (+1.0)	4.18 (−0.5)	4.20
Si	5.52 (+1.7)	5.42 (−0.2)	5.48 (+0.9)	5.50 (+1.3)	5.42 (−0.2)	5.43
Be	2.28 (−0.4)	2.21 (−3.5)	2.24 (−2.2)	2.24 (−2.2)		2.29
	3.55 (−1.1)	3.49 (−2.8)	3.53 (−1.7)	3.52 (−1.9)		3.59

Percentage difference between calculated and experimental data is given in parentheses. All the crystals are cubic, with the exception of Be, which is hexagonal and whose cell is defined in terms of two lattice parameters,  $a$  and  $c$ .

omitted correlation effects with increasing atomic number. For MgO, all Hamiltonians give reasonable lattice parameters, because the electrostatic interaction is about four times larger than for alkali halides (approximately +2 instead of +1 charges on the cation).

Thermodynamic implications also exist in energy derivatives. For example, one of the basic equations of thermodynamics relates pressure to the rate of energy change with the unit cell volume at constant temperature:

$$P = -\left(\frac{\partial E}{\partial V}\right)_T \quad [77]$$

Pressure is an important variable in condensed matter, because the structural modifications crystals undergo by the action of pressure are usually much larger than the modifications from thermal expansions or contractions. Moreover, the structure and stability of high-pressure phases is of particular interest in Earth science. This is an application where ab initio modeling can play an important role and help experimentalists to understand the behavior of minerals in the Earth's mantle. The knowledge of pressure enables investigation of phase stability and transitions. In fact, enthalpy is immediately obtained from the total energy by

$$H = E + PV \quad [78]$$

At  $T = 0$  K, where any transformation of a pure substance tends to be isentropic, phase stability can be related to the enthalpy and a phase transition occurs at those points in the phase diagram where two phases have equal enthalpy. From the computational point of view, it is possible to explore a range of crystalline volumes by isometric lattice deformations and obtain the corresponding values of pressure and, consequently, of enthalpy. It is intended that nuclei are allowed to relax to their equilibrium geometry after

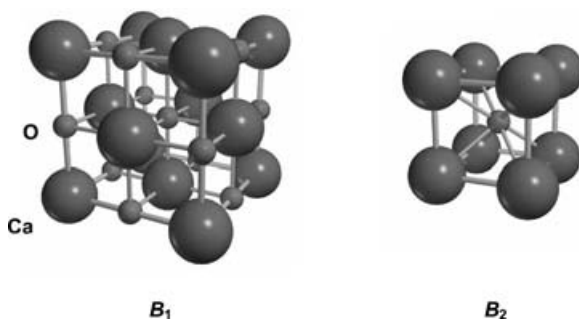


Figure 29  $B_1(Fm\bar{3}m)$  and  $B_2(Pm\bar{3}m)$  crystalline phases of CaO.

every lattice deformation. We illustrate this point with a simple example: calcium oxide. CaO presents two polymorphic cubic phases (Figure 29).  $B_1$  corresponds to a face-centered cubic lattice with calcium in a six-fold coordination environment, and  $B_2$  corresponds to a primitive cubic lattice with every Ca coordinated to eight oxygen ions. We consider the following phase transition of CaO:<sup>153</sup>



As occurs with other observables, the range of the calculated phase transition pressure ( $P_t$ ) and volumes (Table 16) also depends on the choice of the Hamiltonian to some extent. However, the overall agreement with experimental measurements is fairly good.

By interpolating the results of a series of total energy calculations in a range of different lattice volumes, the curves of enthalpy (Figure 30) and the  $V$  vs  $P$  isothermal (Figure 31) can be plotted in any pressure range easily.

**Table 16** Phase Transition Pressure,  $P_t$  (GPa), and Volumes,  $V_{B1}$  and  $V_{B2}$  ( $\text{\AA}^3$ ) Calculated with Different Hamiltonians

Hamiltonian	$P_t$	$V_{B1}$	$V_{B2}$
HF	69.2	21.2	19.0
LDA	57.2	20.6	18.5
PW91	66.1	20.8	18.8
B3LYP	72.7	20.6	18.6
Exp.	60.0 <sup>a</sup>	—	—
	65.0 <sup>b</sup>	—	—
	63.0 <sup>c</sup>	20.7 <sup>c</sup>	18.7 <sup>c</sup>

<sup>a</sup>Taken from Richet et al.<sup>154</sup>

<sup>b</sup>Taken from Jeanloz et al.<sup>155</sup>

<sup>c</sup>Taken from Mammone et al.<sup>156</sup>

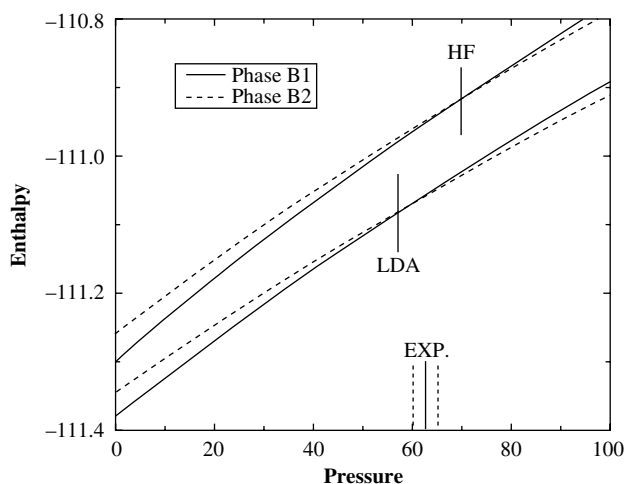


Figure 30 Hartree-Fock and LDA enthalpy of the  $B_1$  and  $B_2$  phases of CaO.

For an application on the phase stability of various  $\text{TiO}_2$  polymorphs, see references 157,158.

Other observables also depend on the total energy derivatives, in particular on second-order derivatives, such as the bulk modulus, the elastic constants, and lattice vibration frequencies.

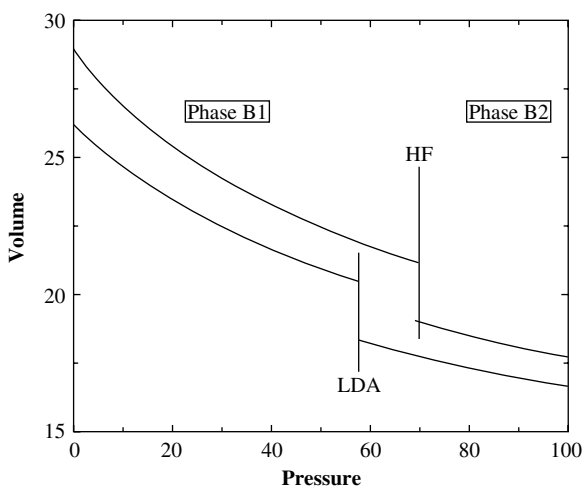


Figure 31 Isothermal phase diagram of CaO (volume,  $\text{\AA}^3$ , versus pressure, GPa), as obtained with HF and LDA approximations. Vertical lines represent the transition pressure.

The bulk modulus,  $B$ , which measures the response of a crystal to isotropic lattice expansion or compression, is related to the second-order derivative of the total energy with respect to the volume,  $V$ , evaluated at the equilibrium volume  $V_0$ :

$$B = -V \left. \frac{\partial^2 E}{\partial V^2} \right|_{V_0} \quad [80]$$

Conversely, the anisotropic response of a crystal to a mechanical force can be described by the elastic constants,  $C_{ij}$ , which are defined as the second derivatives of the total energy with respect to the components  $i$  and  $j$  of the strain tensor,  $\varepsilon$ :

$$C_{ij} = \frac{1}{V} \frac{\partial^2 E}{\partial \varepsilon_i \partial \varepsilon_j} \quad [81]$$

They provide a full description of the mechanical properties of crystalline materials.  $B$  is related to the elastic tensor.<sup>159</sup> In the case of a cubic system, where only three independent components of the elastic tensor differ from zero,  $B$  can be obtained from  $C_{11}$  and  $C_{12}$  as

$$B = \frac{C_{11} + 2 C_{12}}{3} \quad [82]$$

In many computational codes, second derivatives are evaluated numerically (this is also the case of CRYSTAL). This evaluation requires high numerical accuracy in the determination of the total energy. In particular, the lattice deformations involved in calculating lattice constants, which generally reduce the local symmetry, make basis set flexibility (additional sp shells, polarization functions, and so on) necessary. Consequently, a good basis set for the determination of total energy and lattice parameters may be inadequate for the more demanding estimation of these second derivatives.

The effect of the basis set on the bulk properties of MgO is documented in Table 17. Basis set (c), as containing three valence sp shells at oxygen, two

**Table 17** Basis Set Effects on Bulk Properties of MgO at the HF Level of Theory

Case	Basis Set	$a$	$B$	$C_{11}$	$C_{12}$	$C_{44}$	$\Delta E^{\text{atoms}}$
a	8-61/8-51	4.190	200	391	103	201	-715.4
b	8-511/8-411	4.205	181	352	95	188	-699.5
c	8-511*/8-411*	4.194	184	334	108	186	-706.2
Exp.	—	4.195	167	314	94	160	-989.8

$a$  (in Å) is the lattice parameter;  $B$ ,  $C_{11}$ ,  $C_{12}$ , and  $C_{44}$  (in GPa) denote bulk modulus and elastic constants. The cohesive energy  $\Delta E^{\text{atoms}}$  is in kJ/mol.

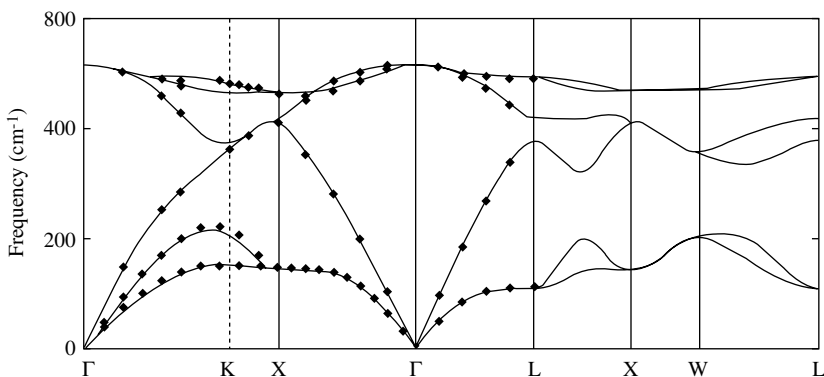
at magnesium, and a set of polarization functions at both ions, can be considered as nearly complete. The table shows that the basis set quality, which has a relatively small influence on  $\Delta E^{\text{atoms}}$ , becomes more important in the calculation of  $B$  and the elastic constants. The addition of d polarization functions is particularly important when computing elastic constants, because the related lattice deformation corresponds to lowering point symmetry. The ions undergo dipolar relaxation, for which the combination of p- and d-type orbitals is required.

Second-order derivatives of energy with respect to the nuclear coordinates are involved in lattice dynamics. Many interesting physical properties<sup>105</sup> are related to lattice dynamics such as vibration spectra (e.g., infrared, Raman, neutron-diffraction), specific heat, thermal expansion, heat conduction, electron-phonon interaction (e.g., superconductivity), and interaction of radiation with matter (e.g., reflectivity of ionic crystals, inelastic scattering of light). The decoupling of the nuclear from the electronic motion through the adiabatic approximation of Born and Oppenheimer<sup>160</sup> and the hypothesis of harmonicity are usually the basic assumptions made when computing lattice dynamics. Within these approximations, force constants relative to all pairs of nuclei in the lattice can be computed as second derivatives of the total energy with respect to small oscillations  $\mathbf{u}$  of the nuclei about their equilibrium positions

$$H_{ij}^{\mathbf{g}} = \left( \frac{\partial^2 E}{\partial \mathbf{u}_i^0 \partial \mathbf{u}_j^{\mathbf{g}}} \right)_{\text{eq}} \quad [83]$$

Equation [83] defines one element of the Hessian matrix relative to the oscillations along the  $i$ -th coordinate of atom A in the 0-cell and along the  $j$ -th coordinate of atom B in the  $\mathbf{g}$ -cell (one of the atoms can always be considered in the 0-cell because of the lattice translation invariance). The Hessian matrix of a crystal obviously has infinite size. However, energy derivatives have the same periodicity as does the potential energy, i.e., the periodicity of the lattice. Therefore, the nuclear wave functions describing nuclear oscillations in the lattice must also obey the Bloch theorem, and application of the periodic boundary conditions allows us to represent the Hessian matrix in the reciprocal space, just as in the case of the electronic Schrödinger equation. Thus, the problem is again reformulated by calculating an infinite set of finite-sized square matrices of dimension  $3N$ , with  $N$  being the number of atoms per unit cell, each corresponding to a particular pattern of displacements of the  $N$  nuclei in any direction in space. Each matrix is associated with a particular  $\mathbf{k}$  point in the reciprocal space and is obtained from Fourier transforming Eq. [83]

$$H_{ij}(\mathbf{k}) = \sum_{\mathbf{g}} e^{i\mathbf{k} \cdot \mathbf{g}} H_{ij}^{\mathbf{g}} \quad [84]$$



**Figure 32** Phonon dispersion of Si from ab initio calculations. Experimental data are denoted by diamonds. Reproduced with permission from reference 161.

Harmonic vibration frequencies are then obtained from diagonalizing the Hessian matrix scaled by the nuclear masses for a convenient sample of reciprocal space points:

$$\Omega_{ij}(\mathbf{k}) = \frac{H_{ij}(\mathbf{k})}{\sqrt{M_i M_j}} \quad [85]$$

Also vibration frequencies, like one-electron energies, then depend on the wave vector  $\mathbf{k}$  and  $3N$  modes at each  $\mathbf{k}$  exist, forming branches. Branches  $\omega_1(\mathbf{k}), \omega_2(\mathbf{k}), \dots, \omega_{3N}(\mathbf{k})$  are called *phonon frequencies*, and the relationship between  $\omega$  and  $\mathbf{k}$  determines the *phonon dispersion*. The computed phonon dispersion of Si is presented in Figure 32 as an example.

Three branches have zero frequency at the  $\Gamma$  point and are associated with the translation of the entire crystal along any direction in space. These branches are called *acoustic* modes as the corresponding vibrations behave as acoustic waves. All other branches show finite nonzero frequencies at  $\Gamma$  and are known as *optical* modes, because they correspond to unit-cell dipole moment oscillations that can interact with an electromagnetic radiation. Acoustic and optical modes can be identified clearly in Figure 32.

The dispersion relation contains the most important information concerning vibration normal modes in a crystal. Lattice vibrations can be measured experimentally by means of classical vibration spectroscopic techniques (infrared and Raman) or neutron inelastic scattering. However, only the latter technique allows one to measure the full spectrum in a range of  $\mathbf{k}$  vectors, whereas with infrared and Raman spectroscopy, only lattice vibrations at  $\Gamma$  can be detected. This limitation for measuring phonon dispersions is serious, because neutron scattering experiments are demanding.

Once the dispersion relation is known, thermodynamic functions can be calculated on the basis of statistical mechanics equations.<sup>57</sup> As an example, the Helmholtz free energy,  $F$ , can be obtained as:

$$\langle F_{\text{vib}} \rangle = \sum_{n,k} \left\{ \frac{1}{2} \hbar \omega_{nk} + k_B T \ln \left[ 1 - \exp \left( -\frac{\hbar \omega_{nk}}{k_B T} \right) \right] \right\} \quad [86]$$

where the sum is extended to all lattice vibrations,  $\omega_{nk}$ , and  $k_B$  is the Boltzmann's constant. Another way of computing thermodynamic functions is based on the use of the phonon density of states. The evolution of the crystal structure as a function of temperature and pressure can also be simulated by minimizing  $G = F + pV$ . The procedure requires a sequence of geometry optimizations, and lattice vibration calculations.

Although the ab initio calculation of vibrational frequencies of molecular systems is a well-known practice, it is not so common in the case of crystalline systems. However, quantum-mechanical calculation of lattice vibrations and phonon spectra has become a subject of increasing interest and effective methods have been developed and implemented. In this respect, a recent review by Baroni et al.<sup>162</sup> gives a detailed overview of the state of the art of ab initio calculation of vibrations and related properties for crystalline materials. Most of the current implementations are based on DFPT and use either plane waves<sup>162,163</sup> or localized functions as a basis set.<sup>164</sup> As an example, calculated and experimental vibration frequencies at  $\Gamma$  are reported in Table 18 for

**Table 18** Calculated and Experimental Vibration Frequencies ( $\text{cm}^{-1}$ ) of  $\alpha$ -quartz at  $\Gamma$

Symmetry	HF <sup>a</sup>	LDA <sup>a</sup>	LDA <sup>b</sup>	LDA <sup>c</sup>	B3LYP	Expt. <sup>d</sup>
$A_1$	216.7	261.6	193.7	238.9	216.0	219.0
	381.3	332.3	355.0	339.3	350.4	358.0
	504.9	482.1	460.1	461.7	465.1	469.0
	1144.4	1089.1	1123.3	1061.0	1085.4	1082.0
$A_{2T}$	395.4	326.3	366.4	341.3	352.3	361.3
	544.1	504.6	489.3	493.4	500.9	499.0
	823.4	791.1	792.2	762.4	783.8	778.0
	1132.4	1086.4	1115.4	1056.5	1076.4	1072.0
$A_T$	138.8	143.4	120.9	133.3	132.5	133.0
	286.5	263.5	257.3	261.3	263.6	269.0
	427.4	376.9	390.0	377.6	391.3	394.5
	490.6	443.8	448.0	443.8	447.0	453.5
	740.9	721.7	703.3	690.8	702.9	698.0
	847.7	835.0	809.6	791.7	810.5	799.0
	1125.2	1070.3	1108.7	1045.0	1068.2	1066.0
	1235.8	1141.7	1190.8	1128.1	1163.1	1158.0

<sup>a</sup> Taken from reference 165.

<sup>b</sup> Taken from reference 60.

<sup>c</sup> Taken from reference 166.

<sup>d</sup> Taken from reference 167.

$\alpha$ -quartz. The average error evaluated with respect to the experimental frequencies is 39.9, 16.6, 16.5, 12.6, and 4.7  $\text{cm}^{-1}$  for the HF, the three LDA, and the B3LYP calculations reported in the table, respectively. B3LYP harmonic frequencies exhibit a fairly good agreement with the corresponding absorptions observed in the experimental spectrum. On the contrary, the HF energy hypersurface curvature is known to be incorrect, and this leads to a regular overestimation of vibration frequencies. LDA results are better than HF but worse than B3LYP. Large differences, however, exist among the three sets of LDA data in Table 18, which represent an interesting example of how results may depend on the particular implementation of a method.

---

## MODELING SURFACES AND INTERFACES

In nature, crystals are not infinite but finite macroscopic three-dimensional (3-D) objects terminated by surfaces. Many phenomena and processes occur at the interface between a condensed phase and the environment. Modeling surfaces is then of great theoretical and practical interest.

A surface can be created by cutting a crystal, which we simulate as an infinite object, through a crystalline plane. Two semi-infinite crystals are then generated containing an infinite number of atoms in the direction orthogonal to the surface, where periodicity, which is essential for applying the Bloch theorem, is lost. We then need further approximations to be able to treat this problem, for which alternative methods have been proposed such as those based on (2-D) clusters, embedded clusters, or slabs. We will focus here on the two-dimensional (2-D) *slab model*.

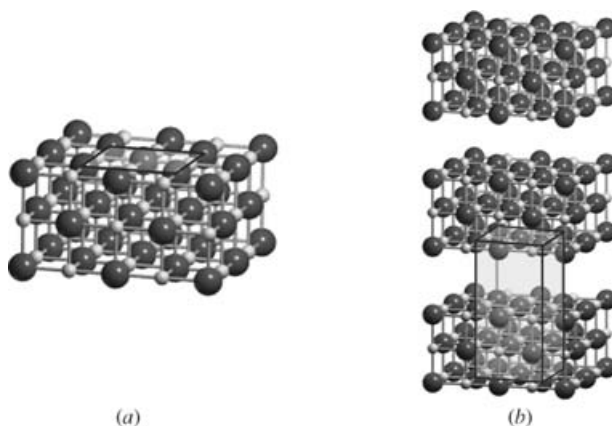
### The Slab Model

The slab model consists of a film formed by a few atomic layers parallel to the (*hkl*) crystalline plane of interest. The film, of finite thickness, is limited by two surface planes, possibly related by symmetry. For sufficiently thick slabs, this kind of model can provide a faithful description of the ideal surface. The adequacy of the model must be checked by considering convergence of geometry, energy, and electronic properties with an increasing number of atomic layers included in the slab.

In actual calculations, two different schemes can be envisaged to deal with a slab model:

1. By imposing 2-D periodic boundary conditions. The slab model is really two-dimensional, with a 2-D unit cell (Figure 33*a*).
2. By forcing a 3-D periodicity (3-D slab model). The three-dimensional system consists of an array of slabs of selected thickness along one direction, separated by vacuum zones, as shown in Figure 33*b*. The vacuum





**Figure 33** Three-layer slab models of the MgO (100) surface. (a) With 2-D periodic boundary conditions. (b) 3-D supercell approximation of the slab model as adopted in plane wave calculations.

zones must be large enough for the fictitious interactions between slabs to be negligible.

When we use a plane wave basis set, which requires a 3-D Fourier representation of many intermediate quantities, such as the charge density, only model (b) can be adopted. On the contrary, when a local basis set is adopted, no problems occur in the implementation of both schemes.

In a recent paper,<sup>168</sup> 2-D and 3-D slab models have been compared. The (110) surface of rutile  $\text{TiO}_2$  served as a case study. Calculations were carried out with CRYSTAL at the HF level with a Gaussian basis set. The convergence of the calculated surface energy and Fermi level was investigated as a function of the slab thickness and interslab vacuum gap. It was found that 2-D and 3-D slabs provide similar convergence with the slab thickness when the vacuum gap is larger than  $6.0\text{\AA}$ . However, model (a) is more general and is to be preferred, for example, in the simulation of an adsorption process, where attention must be paid to spurious interactions among periodic replicas along the direction perpendicular to the slab.

### Specifying the Surface Plane—Miller Indices

The surface is identified by three integers ( $hkl$ )—the so-called Miller indices. The three indices specify a plane of atoms in the crystal by means of the components of a vector perpendicular to that plane. Planes parallel to crystallographic axes  $YZ$ ,  $XZ$ , and  $XY$  are indicated as ( $h00$ ), ( $0k0$ ), and ( $00l$ ), respectively. The planes closest to the origin are then identified by the normal vector with the smallest indices: (100), (010), and (001). Planes

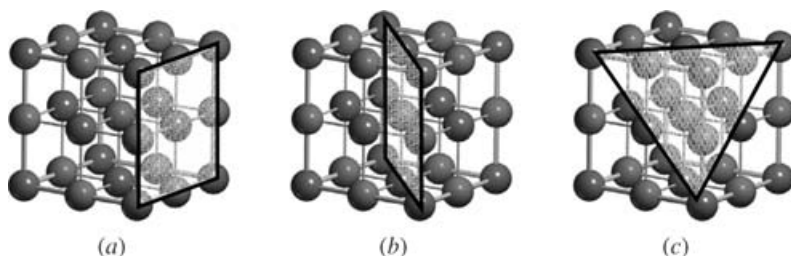


Figure 34 Three types of planes in the MgO crystal: (a) 100, (b) 110, and (c) 111.

parallel to one of the three axes  $X$ ,  $Y$ , or  $Z$  are defined by  $(0kl)$ ,  $(h0l)$ , or  $(hk0)$ , and so on. Some examples for MgO are shown in Figure 34.

### Choosing the Surface Termination

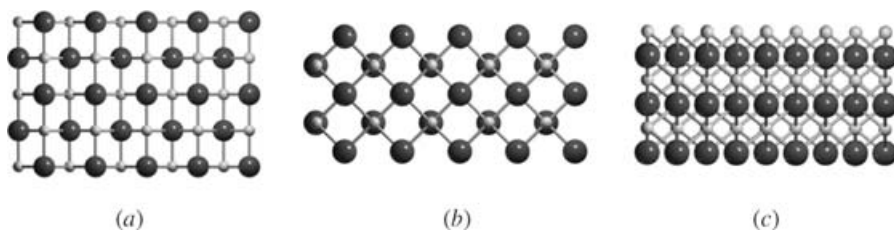
The most interesting surfaces are generally those with low indices, because their creation requires, as a rule, a smaller amount of energy and are, therefore, the most commonly observed. However, not all crystalline surfaces are physically stable or worthy of investigation.

For ionic and semi-ionic crystals, a careful analysis of the possible surface terminations has been carried out by Tasker.<sup>169</sup> Given a slab composed of a certain number of repeated units, which are in turn made up of atomic layers parallel to the selected plane, the resulting structure can be classified in one of the three following categories (Tasker's types):

- Type 1: the slab consists of neutral layers with the same stoichiometry of the host crystal.
- Type 2: the slab consists of charged layers arranged symmetrically so that the repeated unit presents no net dipole perpendicular to the surface.
- Type 3: the slab consists of charged layers alternating in such a way that the repeated unit presents a net dipole normal to the surface.

Although type 1 and 2 surfaces may exist, those of type 3, also referred to as dipolar surfaces, are unstable and can only be stabilized through some mechanism to remove the macroscopic field (i.e., by reconstruction, molecular adsorption, and so on). In the MgO case (see Figure 35), the (100) and (110) surfaces correspond to type 1, whereas the (111) surface is type 3.

In covalent solids, the creation of a surface requires cutting covalent bonds, which means that dangling bonds would be present at the surface. The resulting instability is partly reduced either by creating new bonds, giving rise to a reconstruction of the surface, or chemisorbing atoms from the environment (e.g., H, Cl). The saturation of dangling bonds by chemisorption is important, for example, in silicates. When a surface is cut out from the

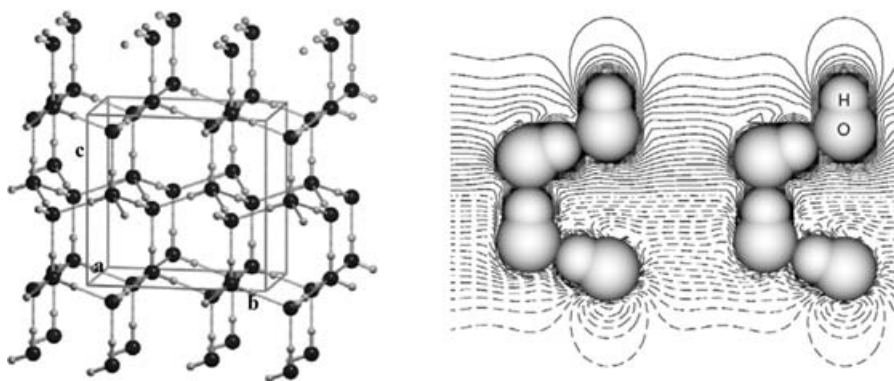


**Figure 35** Five-layer slab model of MgO surfaces (lateral view): (a) (100) surface—type 1, (b) (110) surface—type 1, and (c) (111) surface—type 3.

bulk, unstable  $\text{Si-O}\cdot$  radicals at the surface react readily with water to give a fully hydroxylated surface with hydrophilic character.

When cutting molecular crystals, the molecular topology must be preserved and only intermolecular bonds are cut. If the molecule has a dipole moment, attention must be paid to the surface termination, because the slab can possess a net dipole perpendicular to the surface. For instance, ice XI, a proton-ordered phase of ice, is ferroelectric because its basic repeating unit, consisting of four water molecules, has a net dipole along the *c*-axis (see Figure 36 on the left).

Thus, by cutting a slab parallel to the (001) face, a dipolar surface is created that would be highly unstable,<sup>170</sup> according to Tasker's classification. Figure 36, on the right, shows that the electrostatic potential difference between the two surfaces is large, which explains their instability. A stable slab of C-ice can be obtained only by cutting the crystal in such a way that the ferro-electric axis is parallel to the surface, which is equivalent to selecting the (010) surface.



**Figure 36** On the left: Structure of ice XI; On the right: Electrostatic potential at the (001) surface of ice XI. Consecutive isodensity lines differ by 0.01 a.u.; continuous, dashed, and dot-dashed curves correspond to positive, negative, and zero potential, respectively. Isopotential lines corresponding to potential values larger than 0.2 a.u. in module are not plotted.

Even if it may appear unrealistic, because of the electric strain parallel to the surface, it comes out that this structure is particularly stable.<sup>170</sup>

## Surface Formation Energy and Stability

Within the slab model approach, the surface formation energy is computed as

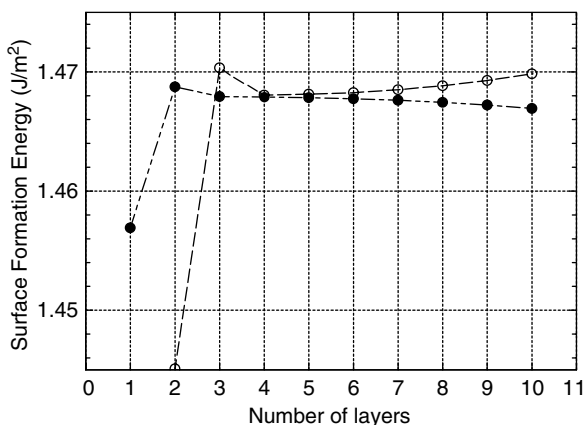
$$\Delta E_{\text{surf}}^n = \frac{(E^n - nE_{\text{bulk}})}{2A} \quad [87]$$

where  $E^n$  is the energy of an  $n$ -layer slab,  $E_{\text{bulk}}$  is the energy of a single layer's worth of bulk material, and  $A$  is the area of the primitive surface unit cell. The factor  $1/2$  accounts for the existence of two limiting surfaces.  $\Delta E_{\text{surf}}^n$  is then the energy per unit area required to form the surface from the bulk, and it is intrinsically a positive quantity (if not, the bulk would exfoliate). As more layers are added in the calculation by increasing the slab thickness ( $n \rightarrow \infty$ ),  $\Delta E_{\text{surf}}^n$  will converge to the surface formation energy per unit area. This important check should be performed, when studying surfaces.

For MgO, which is a wide band gap insulator, the computed HF  $\Delta E_{\text{surf}}^n$  converges rapidly, as shown in Figure 37.

In metals or small band gap semiconductors, convergence can be slower and numerical noise larger. In Eq. [87], total energies from 3-D and 2-D systems are used, and this in principle might create problems of “equivalent” accuracy in algorithms that are specific for 2-D and 3-D. As a cross check, we can use the following definition for the surface energy:

$$\Delta E_{\text{surf}}^n = \frac{E^n - n(E^n - E^{n-1})}{2A} \quad [88]$$



**Figure 37** Dependence of the surface energy on the number of layers for a MgO (100) slab model. Filled circles from Eq. [87]; open circles from Eq. [88].

In this expression  $E_{\text{bulk}}$  has been replaced by  $E^n - E^{n-1}$ . So the surface formation energy is determined from a series of 2-D calculations. If each additional layer in the slab is seen as the central layer, it is clear that  $E^n - E^{n-1}$  should converge to the energy of a single layer in the bulk crystal. In Figure 37, the MgO(100) surface energy computed with Eq. [88] is also reported. A more extensive discussion on the use of Eqs. [87] and [88] when computing the surface energy of metallic lithium, as a case study, can be found in Doll et al.<sup>171</sup>

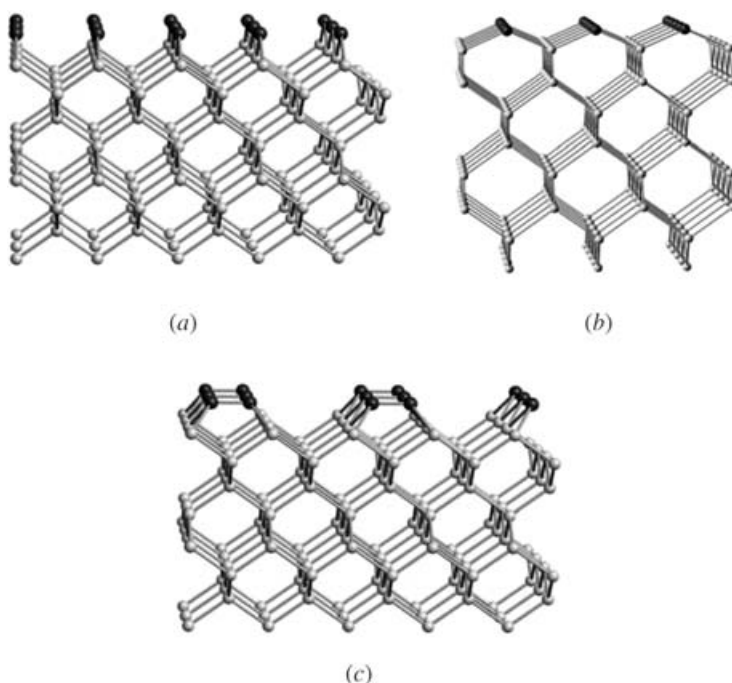
When the surface formation energy of different surfaces is available, the relative stability can be evaluated. Surface stability has relevance in determining the crystal morphology, although kinetic effects in many cases can also play an important role. As an example, we can compare the stability of the (100) and (110) MgO surfaces<sup>172</sup> (see Figure 34). When a five-layer slab model is adopted, the computed HF surface energies are 1.47 J/m<sup>2</sup> and 5.24 J/m<sup>2</sup>. The difference in stability between these two surfaces is easily explained on the basis of the different environment of the surface ions: at the (100) surface Mg and O are fivefold-coordinated, whereas at the (110) surface, coordination decreases to four.

## Surface Relaxation and Reconstruction

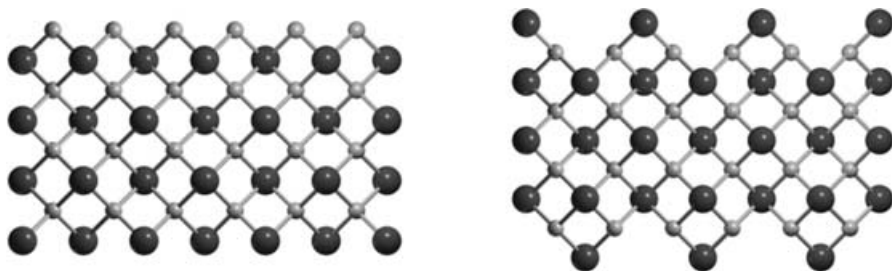
In general, when a surface is cut out of a perfect crystal, the atoms near the surface will move away from their bulk positions to minimize the surface energy. When the atomic displacements do not change the symmetry of the slab, this is referred to as surface *relaxation*. However, in some systems, the rearrangement is deeper and the surface has a tendency to *reconstruct*; that is, the periodicity of the surface layer changes from that implied by pure bulk termination. If the primitive cell of the surface is defined by lattice vectors **a** and **b**, then a reconstruction introducing a new periodicity, involving two steps in **a** and three steps in **b**, is called a (2 × 3) reconstruction. To model such a phenomenon, the slab model can be combined with the supercell approach by creating a 2-D cell and then enlarging it to introduce the new periodicity. A typical example of surface reconstruction is the (111) and the (100) surfaces of silicon.

At the (100) surface, the presence of highly unstable dangling bonds at the top of the fully unrelaxed surface (indicated in dark in Figure 38a) is partly reduced by the formation of new bonds leading to a (2 × 1) reconstruction (Figure 38c). A 2-D cell of double size is necessary to model such a reconstruction.

In ionic crystals, reconstruction effects can also be involved in the stabilization of polar surfaces (Tasker's type 3). For instance, the (100) surface of the fluorite-type crystal of Li<sub>2</sub>O becomes stable if half of the Li atoms are moved from the bottom face of the slab to the top face above the oxygen atoms to produce a zero-dipole structure (Figure 39). In fact, this kind of surface has been observed experimentally.<sup>173</sup>

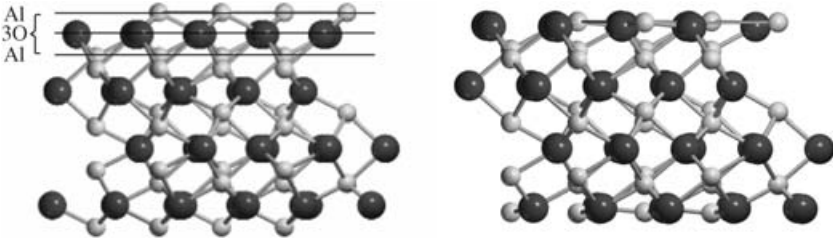


**Figure 38** Pictorial view of the (100) surface of silicon. (a) unrelaxed surface; (b) relaxed surface and (c)  $(2 \times 1)$  reconstructed surface.



**Figure 39** Lateral view of the (100) surface of  $\text{Li}_2\text{O}$ . On the left: unreconstructed dipolar surface. On the right: reconstructed zero-dipole surface.

Relaxation effects are particularly significant in partially-ionic oxides like  $\alpha\text{-Al}_2\text{O}_3$ . Alumina is of enormous technological importance and is widely used by chemists as a support material in most heterogeneous catalysts. The simplest surface is the Al-terminated (0001) one, which corresponds to the basal plane. As shown in Figure 40 on the left, the basic repeated unit perpendicular to the (0001) surface is a three-layer slab consisting of Al-3O-Al atomic layers, where 3O indicates a layer containing three oxygens per 2-D cell. In the



**Figure 40** A 12-layer slab model of the Al-terminated (0001) surface of  $\alpha$ - $\text{Al}_2\text{O}_3$ . Lateral view of the unrelaxed structure (left) and the relaxed structure (right).

unrelaxed slab model, the topmost layer consists of undercoordinated Al ions, as is shown in Figure 40. Consequently, this is an unfavorable structure and the Al atom undergoes a large relaxation (compare left and right pictures in Figure 40). In particular, the first-to-second interlayer spacing largely contracts to lower the surface energy. For a 15-layer slab model, the contraction is as large as  $-78.8\%$  at HF, and it increases further when DFT methods are adopted:  $-87.1\%$  and  $-79.2\%$ , with SVWN and BLYP methods, respectively.

Structural relaxation has also significant effects on the relative stability of different surfaces. For  $\alpha$ - $\text{Al}_2\text{O}_3$ , five different low-index nonpolar faces are usually believed to be competitive for stability.<sup>174</sup> At the HF level of theory,<sup>175</sup> the stability order is as follows:

Unrelaxed: face	(01 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}1$ )	<	(10 $\bar{1}0$ )	<	(0001)
$\Delta E_{\text{surf}}^n (\text{J/m}^2)$	2.70		3.27		4.18		4.50		4.85
Relaxed: face	(0001)	<	(10 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}0$ )	<	(10 $\bar{1}1$ )
$\Delta E_{\text{surf}}^n (\text{J/m}^2)$	1.90		2.00		2.37		2.42		2.47

These data show that the inclusion of structural relaxation has a dramatic effect: The stability order is almost completely reversed, and the surface formation energy spans a much more narrow range of values with respect to the unrelaxed data. Relaxation is then important. For example, the (0001) surface, which is the most unstable unrelaxed face, becomes the most stable when relaxation is taken into account. It is worth noting that slab models with more than 20 atomic layers are required to reach full convergence on stability order. The computed trend for the unrelaxed surfaces is close to that obtained from classical simulations within a fully ionic model by Tasker<sup>174</sup> and Mackrodt,<sup>176</sup> whereas the three sets differ after relaxation:

Unrelaxed: Tasker	(01 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}1$ )	<	(0001)	<	(10 $\bar{1}0$ )
Mackrodt	(01 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}1$ )	<	(0001)	<	(10 $\bar{1}0$ )
Relaxed: Tasker	(01 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}0$ )	<	(0001)	<	(10 $\bar{1}1$ )
Mackrodt	(0001)	<	(10 $\bar{1}0$ )	<	(01 $\bar{1}2$ )	<	(11 $\bar{2}0$ )	<	(10 $\bar{1}1$ )

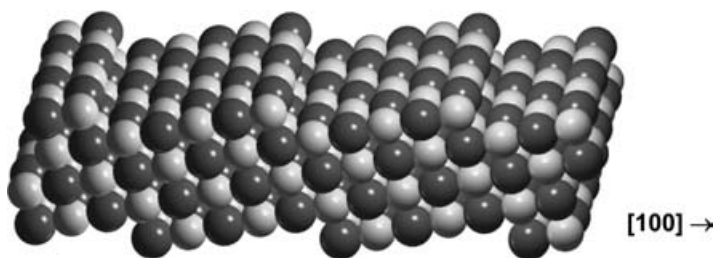


Figure 41 MgO (501) slab model with steps along the [100] direction.

The difference between classical and quantum-mechanical calculations is caused by the important role played by electronic rearrangement.

### Vicinal Surfaces—Modeling Steps and Kinks

Real surfaces are rarely atomically flat. Typically, a surface will be covered in plateaus, with edges and corners. Within the slab model, the use of *vicinal* surfaces (those cut at an angle slightly different to that of the low-energy surface) is a useful trick to model steps and kinks. For example, Figure 41 shows a few unit cells of the (501) surface of MgO, which contains a sequence of steps in the [100] direction. These faces may mimic situations occurring at defective ionic surfaces and can therefore be used to investigate the reactivity of different kinds of defects. For instance, it has been found<sup>177</sup> that heterolytic splitting of hydrogen [ $\text{H}_{2(\text{g})} \rightarrow \text{H}_{(\text{ads})}^+ + \text{H}_{(\text{ads})}^-$ ] may take place with low activation energy and favorable energy balance at a (*n*01) face of MgO with  $n \geq 3$ , which simulates a regular step at a (100) surface.

Ab initio calculations on bare surfaces are now sufficiently accurate and efficient; we propose a (certainly not exhaustive) list of surface studies performed with the CRYSTAL code: halides ( $\text{LiF}$ ,<sup>178,179</sup>  $\text{NaCl}$ ,<sup>180</sup>  $\text{CaF}_2$ <sup>181,182</sup>), oxides ( $\alpha\text{-Al}_2\text{O}_3$ ,<sup>183–186</sup>  $\text{CeO}_2$ ,<sup>187</sup>  $\alpha\text{-Cr}_2\text{O}_3$ ,<sup>188</sup>  $\text{Li}_2\text{O}$ ,<sup>189</sup>  $\text{MgO}$ ,<sup>172,190</sup>  $\text{SnO}_2$ ,<sup>191</sup>  $\text{TiO}_2$ ,<sup>192,193</sup>  $\text{ZnO}$ ,<sup>194–196</sup>  $\text{ZrO}_2$ <sup>184,187</sup>), sulphides ( $\text{Li}_2\text{S}$ ,<sup>197</sup>  $\text{FeS}_2$ ,<sup>198</sup>  $\text{PbS}$ ,<sup>199</sup>  $\text{RuS}_2$ <sup>200</sup>), metals ( $\text{Li}$ ,<sup>171</sup>  $\text{Cu}$ ,<sup>201,202</sup>  $\text{Ag}$ ,<sup>203,204</sup>  $\text{Ni}$ ,<sup>205</sup>  $\text{Pt}$ <sup>206</sup>), perovskites ( $\text{BaTiO}_3$ ,<sup>207</sup>  $\text{LaMnO}_3$ ,<sup>208</sup>  $\text{SrTiO}_3$ <sup>209,210</sup>), molecular crystals (water ice<sup>211</sup>), and silicates.<sup>212</sup>

### Adsorption on Surfaces

An appropriate slab model of a surface is also useful to study the adsorption of atoms or molecules. Both physisorption and chemisorption processes can be modeled easily.

The binding energy between surface and adsorbate,  $\Delta E$ , is a key observable. It corresponds to the process in which the molecules move from an ideal



gas state onto the surface, and it is defined as

$$\Delta E = E(\text{slab}) + N \cdot E(\text{mol}) - E(\text{slab/ads}) \quad [89]$$

where  $E(\text{mol})$  is the energy of one isolated adsorbed molecule and  $N$  is the number of adsorbed molecules per unit cell. These energies are defined per unit cell and are negative.

$\Delta E$  can also be written as the sum of two contributions:

$$\Delta E = \Delta E_{\text{ads}} + \Delta E_{\text{L}} \quad [90]$$

The first contribution is the binding energy per unit cell per adsorbed molecule and is defined as

$$\Delta E_{\text{ads}} = E(\text{slab}) + E(\text{ads}) - E(\text{slab/ads}) \quad [91]$$

where  $E(\text{slab/ads})$  is the total energy of the slab in interaction with the periodic array of adsorbed molecules,  $E(\text{slab})$  is the energy of the slab alone, and  $E(\text{ads})$  is the energy of the periodic array of adsorbed molecules without the underneath solid surface.  $\Delta E_{\text{L}}$  is the lateral interaction energy, per unit cell, among the adsorbate molecules, i.e., without the underneath surface, and can be either positive (repulsion) or negative (attraction), depending on the nature of the ad-molecules. In the limit of low coverage, i.e., large distances between molecules,  $\Delta E_{\text{L}}$  tends to zero, so that  $\Delta E_{\text{ads}} \sim \Delta E$ .

Modeling different coverages is of interest in adsorption processes and can be achieved easily by enlarging the underlying surface unit cell (i.e., within a supercell approach), so that the density of adsorbed molecules can be increased or reduced. In the limit of low coverage, lateral interactions tend to vanish and adsorbed molecules can be considered as isolated.

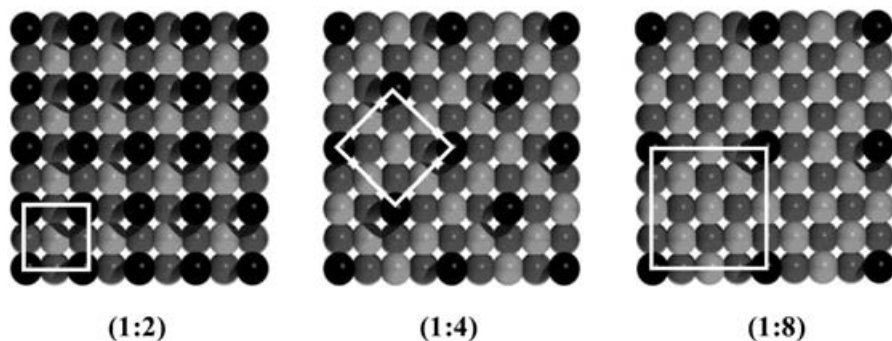
If the basis set is not complete, the computed binding energy per adsorbed molecule should be corrected for the BSSE following the CP method as proposed for molecular complexes. Therefore, the BSSE corrected binding energy  $\Delta E_{\text{ads}}^{\text{C}}$  becomes

$$\Delta E_{\text{ads}}^{\text{C}} = E(\text{slab}/[\text{ads}]) + E([\text{slab}]/\text{ads}) - E(\text{slab/ads}) \quad [92]$$

in which  $E(\text{slab}/[\text{ads}])$  and  $E([\text{slab}]/\text{ads})$  are, respectively, the energy of the slab with the basis functions of the adsorbate only, and vice versa.

Let us consider, as an example, the interaction of carbon monoxide with the MgO (100) surface. This system has been the subject of several simulations where cluster and periodic models have been used (see Damini et al.<sup>213</sup> and references therein). Here we focus on a periodic approach through a slab model.

The MgO (100) surface is modeled by a three-layer slab. As shown above (Figure 37), the adopted slab is thick enough to ensure good convergence of



**Figure 42** Pictorial view of the CO/Mg (100) system at different CO coverages. Light and dark gray spheres are Mg and O, respectively; black spheres represent the CO molecules.

the surface energy. Three different CO coverages have been considered, namely the (1:2), (1:4), and (1:8) (here (1: $n$ ) means that adsorption involves one CO molecule per  $n$   $\text{Mg}^{2+}$  ions), in order to study the effect of the lateral interactions on  $\Delta E$ . Preliminary calculations indicate that the adsorption through the oxygen atom is disfavored in comparison with that via the carbon atom, so it will not be considered in the following discussion. The resulting 2-D unit cell are shown in Figure 42.

The CO molecules were adsorbed at both faces of the MgO (100) slab model. Calculations were carried out at the HF and DFT levels. For the latter, the B3LYP method was adopted. Two basis sets were employed, hereafter indicated as A and B, with B being more accurate and costly than A. For more details concerning the all-electron basis sets and optimized geometry, see Damini et al.<sup>213</sup>

Table 19 shows the HF and B3LYP computed binding energies.  $\Delta E$  is small at both levels of theory, with a value of about 10 kJ/mol at the best B3LYP/B level. The BSSE correction is large:  $\Delta E^C$  is nearly null at HF/A, around 2.0 kJ/mol at B3LYP/A and 3.7 kJ/mol at B3LYP/B for the (1:4)

**Table 19**  $\Delta E$  (kJ/mol) of CO Adsorbed on the MgO (100) Surface as a Function of CO Coverage, Hamiltonian and Basis Set

Method	HF/A			B3LYP/A			B3LYP/B	
	Coverage (1:2)	(1:4)	(1:8)	(1:2)	(1:4)	(1:8)	(1:2)	(1:4)
$\Delta E$	6.2	6.4	6.4	13.6	14.0	14.0	9.1	10.3
$\Delta E^C$	0.9	0.8	0.7	2.8	3.0	2.0	3.0	3.7
$\Delta E_L$	2.4	0.6	0.1	1.6	0.4	0.1	2.6	0.4
$\Delta E_N$	-1.5	0.2	0.6	1.2	2.6	1.9	0.4	3.3

$\Delta E^C$  is corrected for BSSE.  $\Delta E_L$  is the lateral interaction and  $\Delta E_N = \Delta E_C - \Delta E_L$  is the net binding energy.

coverage. The lateral interaction energy,  $\Delta E_L$ , is important only when the smallest cell (1:2) is used. The final binding energy at the B3LYP/B level for the (1:4) coverage is 3.3 kJ/mol. This result is to be compared with the best experimental measurement of 13.5 kJ/mol.<sup>214,215</sup>

The disagreement with experiment suggests that dispersive contribution to the binding energy, not accounted for by both HF and DFT methods, may play an important role. To cope with this flaw, Ugliengo and Damin<sup>216</sup> proposed an interesting approach to include correlation contributions at the MP2 level of theory, through a kind of cluster-in-crystal embedding technique. Those authors were able to obtain a final extrapolated MP2 binding energy of 12.7 kJ/mol, in good agreement with the experimental value, which shows that dispersive contributions account for about 7 kJ/mol.

Further examples of adsorption systems investigated with the CRYSTAL code include H<sub>2</sub>O on NaCl<sup>180</sup>; CO, N<sub>2</sub>, and O<sub>2</sub> on LiF(100)<sup>178,179</sup>; Na,<sup>217</sup> K,<sup>218</sup> and noble metals<sup>219</sup> on TiO<sub>2</sub> surfaces; CH<sub>3</sub>OH,<sup>220</sup> CO,<sup>221</sup> and CO<sub>2</sub><sup>222</sup> on SnO<sub>2</sub>; CO on Cu<sub>2</sub>O (111)<sup>223</sup>; formic acid<sup>224</sup> and hydrogen<sup>225</sup> on ZnO (10-10); the interaction of pyrite (100) surfaces with O<sub>2</sub> and H<sub>2</sub>O<sup>226</sup>; NH<sub>3</sub> on a model of a silica surface<sup>227</sup>; HCl on water ice<sup>211</sup>; Cl on Cu (111)<sup>201</sup>; and Pd on  $\alpha$ -Al<sub>2</sub>O<sub>3</sub> (0001).<sup>228</sup>

## Interfaces

Apart from the simulation of ideal surfaces, increasing interest in “real” 2-D crystals now exists, which are quasi-periodic structures in two dimensions but only a few atomic layers thick, and which may present new and useful properties precisely because of their limited thickness. This branch of nanoscience is then an ideal ground for application of the slab model.

The study of epitaxial interfaces between crystals of different nature is an example of this flexible technique. The interface is modeled by creating two slabs and letting them interact to form a sort of supra-slab model (i.e., a wafer), as shown schematically in Figure 43. Care must be paid to the lattice mismatch because the 2-D unit cells, with different size, must match at the interface.

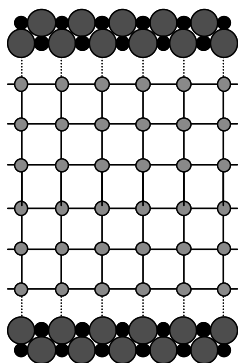
The interaction energy at the interface or adhesion energy is computed as

$$\Delta E = \frac{E(\text{interface}) - E(\text{slab1}) - E(\text{slab2})}{2} \quad [93]$$

which is the difference between the total energy of the interface model and the energies of the isolated slabs, divided by a factor two to account for the existence of two outer surfaces.

By taking the BSSE into account, we get:

$$\Delta E = \frac{E(\text{interface}) - E(\text{slab1}/[\text{slab2}]) - E([\text{slab1}]/\text{slab2})}{2} \quad [94]$$



**Figure 43** Schematic representation of an interface as simulated through a slab model.

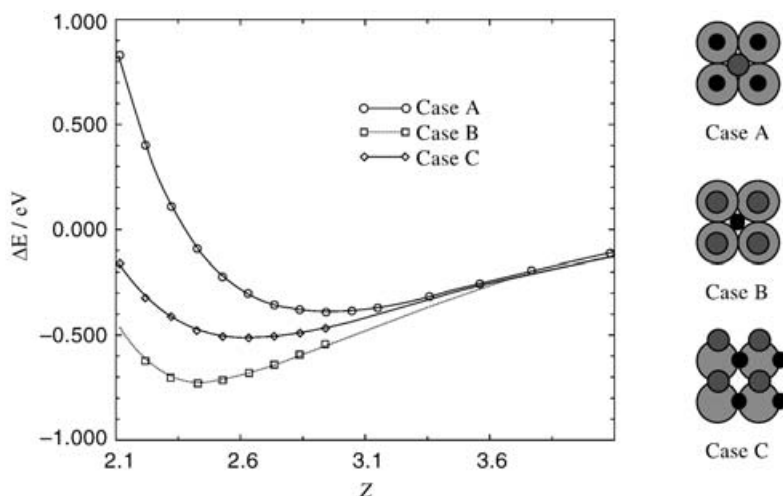
where  $E(\text{slab1}/[\text{slab2}])$  and  $E([\text{slab1}]/\text{slab2})$  are, respectively, the energy of the first slab obtained in the presence of the basis functions of the second one, and vice versa.

An interesting example of the application of the slab model approach to the study of interfaces is the modeling of ultrathin oxide films on metallic substrate, which has been the subject of recently published papers by Pisani et al.<sup>229,230</sup> It deals with a model of the epitaxially grown MgO (100) thin film on Ag (100). An advantage of studying the MgO/Ag system is that it has a small mismatch between MgO and Ag lattice parameters ( $\sim 3\%$ ) and allows a cube-on-cube epitaxy.

This system was studied at the DFT level with both LDA (SPZ) and GGA (PW91) methods. A six-layer slab model of Ag parallel to the (100) surface was adopted, covered with one or two MgO monolayers on both surfaces (as in Figure 43). Three high symmetry configurations were considered: (A)  $\text{Mg}^{2+}$  upon a surface Ag atom, (B)  $\text{O}^{2-}$  upon a silver atom, and (C) both ions bridged.

The interaction energy computed with respect to the distance between the MgO overlayer, at first considered as a rigid entity, and the silver substrate is represented in Figure 44. The most stable configuration was found to be (B). Such a configuration was further analyzed, allowing an independent optimization of Mg and O positions to check the possibility of rumpling. Calculated MgO/Ag distances and BSSE-corrected adhesion energies are reported in Table 20.

Both DFT methods show that Mg relaxes toward the silver surface, although the rumpling is partly recovered when a second MgO layer is present, probably caused by the electrostatic attraction of the oxygen ions in the upper plane. Accordingly, the additional MgO layer also reduces the adhesion energy. It is worth noting that LDA gives an interaction energy that is larger (more than double) than GGA. This confirms the general observation that LDA



**Figure 44** Interaction energy  $\Delta E$  per MgO unit between a six-layer-thick silver slab and two monolayers of MgO at distance  $z$ . The reported data refer to the SPZ method.

methods tend to overestimate binding energies with respect to gradient-corrected methods, as already pointed out here for different kinds of system.

The interface model can be further complicated by considering the possible adsorption of molecules. Obviously, the substrate will modulate the interaction of the surface with adsorbates. For instance, the interaction of the Ag-supported MgO with water was simulated and compared with that on a pure MgO<sup>229</sup> surface.

Interfaces are a field of growing interest, and some other applications have been carried out, which include ultra-thin adlayers of Ag on

**Table 20** Calculated Properties of the MgO/Ag System in the (B) Configuration

	SPZ		PW91	
	No Rumpling	Rumpling	No Rumpling	Rumpling
One layer				
$d_{\text{Mg}}$	2.45	2.32	2.55	2.39
$d_{\text{O}}$	2.45	2.47	2.55	2.55
$\Delta E$	-41.8	-53.1	-17.7	-29.6
Two layers				
$d_{\text{Mg}}$	2.45	2.38	2.55	2.46
$d_{\text{O}}$	2.45	2.47	2.55	2.55
$\Delta E$	-41.1	-42.6	-18.2	-19.2

$d_X$  is the distance of ion X from the surface.  $\Delta E$  is the BSSE-corrected adhesion energy of the MgO overlayer on the Ag substrate (in kJ/mol).

MgO(100), MgO(110),<sup>231–234</sup> and Al<sub>2</sub>O<sub>3</sub><sup>235</sup>; the interface between a transition metal oxide, NiO, and Ag(100),<sup>236</sup> and the study of oxides on oxides, like the MgO/NiO films.<sup>237,238</sup>

## MODELING DEFECTIVE SYSTEMS

### Defects in Solids

Defects in solids are ubiquitous and can be found both in the bulk and at the surface of materials.<sup>239,240</sup> Two classes can be distinguished: point defects and extended defects. The former, also called local defects, produce a modification of the site environment of an otherwise perfect lattice: for instance, the absence of an atom in a lattice position (vacancy), the presence of an atom in an interstitial position (interstitial defect), or the substitution of an atom for another atom of a different chemical species at a regular lattice site (substitutional defect). Figure 45 shows typical examples of local defects in an ionic solid.

Extended defects, on the other hand, correspond to structural imperfections in the assembly of either lattice planes (planar defects), as stacking faults in layer structures, or lattice directions (linear defects), as dislocations.

At their surface, along with point (e.g., vacancy on a terrace) and extended defects (observed in crystal growth processes), solids present other typical defects such as vertices, edges, kinks, and steps, as shown in Figure 46. Defects play an important role in determining the surface reactivity, as briefly mentioned in the previous section for a stepped MgO surface model. In a

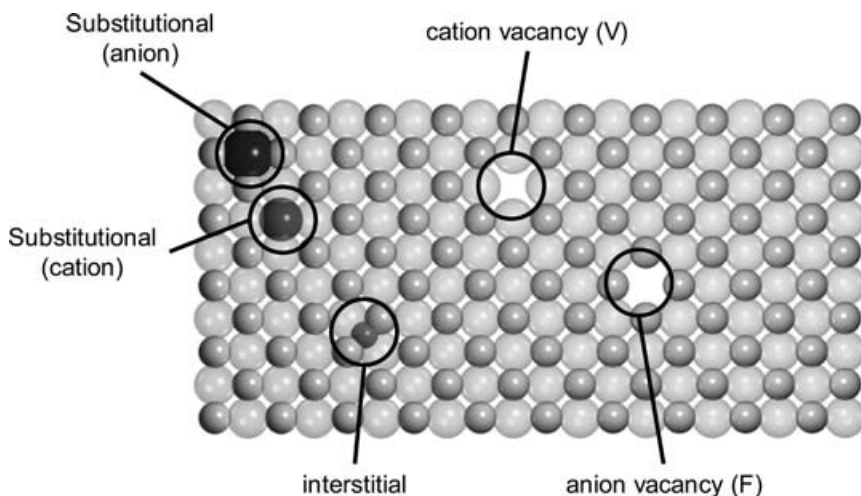


Figure 45 Examples of local defects in ionic solids.

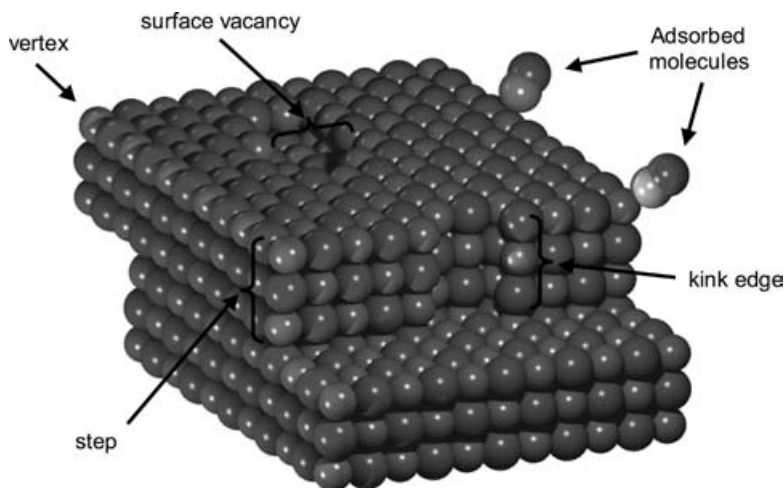


Figure 46 Examples of surface defects.

broad classification of defects, a surface with adsorbed atoms or molecules can also be considered as a defective system with respect to the bare surface.

## How to Model a Defect

Here, we mainly focus on the modeling of point defects in a perfect crystal. Because of their local nature, we can assume that the perturbation in the crystal is small. Therefore, both structural and electronic effects can be considered as confined in the finite region close to the defect (defect zone). The rest of the system can then be treated as a perfect host with a well-defined periodic structure. On the basis of these assumptions, two strategies can be envisaged in the modeling of point defects:

1. The host as an environment
2. The defect as an impurity

### *The Host as an Environment*

Both perfect and defective systems are simulated as a *finite cluster*, i.e., a relatively small cluster is cut out of the bulk structure containing the defect. This kind of strategy is also called the *cluster approach*.<sup>241</sup> The main advantage of this approach is its flexibility: first, because it does not assume any translational symmetry and allows us to investigate complex structures as defects in amorphous or disordered solids. Second, high quality standard molecular ab initio codes can be used, which thus allow a high-level quantum-mechanical treatment of the defect zone.

Nevertheless, the abrupt termination of an otherwise infinite solid gives rise to spurious effects, such as nonphysical electronic states localized at the boundary, levels in the gap, finite size quantum effects, and neglect of Coulomb and exchange interactions with the environment. Such effects depend on the nature of the chemical bond in the studied solids: ionic, semi-ionic, covalent, metallic, and molecular. Therefore, additional manipulations may be needed to take into account the environment where the cluster would be contained. For ionic systems, the most important contribution from the missing ions is the Coulomb field, which may be approximated by introducing a finite/infinite array of point charges or polarizable semiclassical ions, described with a shell model.<sup>241</sup> For metals and covalent solids, border quantum effects are important and not easily simulated. In covalent solids, the atoms at the surface present dangling bonds that have to be saturated with capping atoms, like H, F, and so on. An ambiguity remains, however, related to the field modification originating from the presence of the capping atoms, especially for mixed ionic-covalent situations.

It is worth mentioning here a variant of the cluster model that has been recently proposed to embed the cluster in its environment, based on a cluster-in-cluster scheme. The method, proposed by Morukuma et al., is the so-called our N-layer integrated molecular orbitals—molecular mechanics (ONIOM) scheme.<sup>242–244</sup> The system is modeled through a large cluster (real system) that is then partitioned into two or more regions (model cluster). Each region is described with a hierarchical level of theory, higher for the model cluster and lower for the real system. The total energy is obtained from that of the real system, i.e., the largest cluster, including a series of correcting terms. This scheme, originally developed in a molecular context, has also been reformulated for defects-in-solids.<sup>245,246</sup>

Finally, size and shape of the cluster are critically important and results should converge with the cluster size. Unfortunately, when enlarging the cluster, the number of atoms grows fast, which makes the calculation rapidly unaffordable. For solids with complex frameworks and containing a large unit cell, such as zeolites, it is also difficult to preserve the memory of the original crystalline structure in the cluster.

### *The Defect as an Impurity*

The host system is treated as a perfect crystalline structure, and the exploitation of periodicity or quasi-periodicity is an essential ingredient when treating the defect as an impurity. From a quantum-mechanical point of view, the defect is treated as a perturbation to the electronic structure of the perfect crystal environment.

Two different approaches can be adopted:

1. The embedded cluster approach
2. The supercell approach



The embedded cluster approach is based on schemes that smoothly link the quantum-mechanical solution of the cluster to the perfect crystal. Different methods have been developed based on either Green function<sup>247</sup> or group-function (localized crystalline orbitals)<sup>248</sup> techniques.

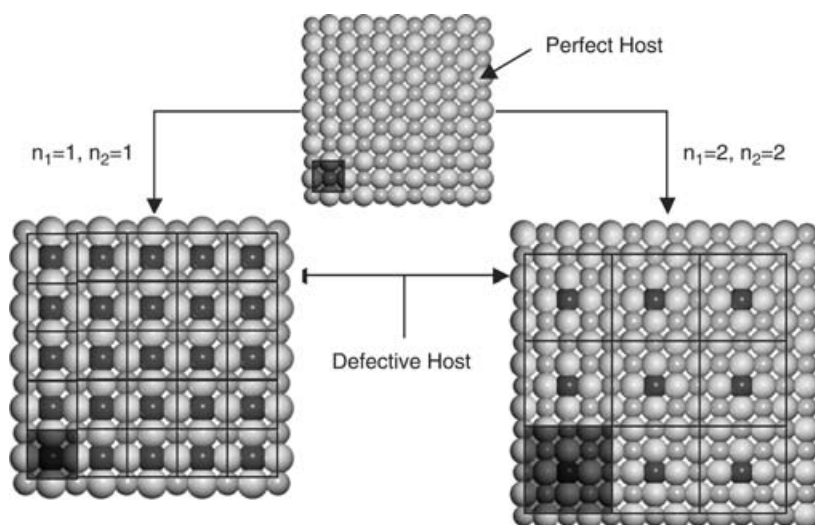
In the supercell approach, the defect is instead enclosed in a sufficiently large unit cell and periodically repeated throughout space. A common problem with both approaches is the availability of high-level quantum-mechanical periodic solutions, because, as already mentioned, it is difficult to go beyond the one-electron Hamiltonian approximations (HF and DFT), at present.

The supercell scheme is the most widely adopted approach because it is easily implemented in all periodic *ab initio* codes. Embedding approaches, on the other hand, may require specific and not widely disseminated softwares, which make their development slow, and their accuracy relatively low. A discussion of limits and merits of the embedding techniques can be found in Pisani.<sup>249</sup>

In the following pages, we illustrate in more detail the supercell approach and discuss a few examples.

## The Supercell Approach

The supercell approach consists of a periodic replica of the defect, which is enclosed in a large nonprimitive unit cell. A pictorial view (in 2-D) of the supercell approach is given in Figure 47, where, by starting from the perfect



**Figure 47** Schematic example of a local defect in a two-dimensional lattice as modeled by the supercell approach.

host, a substitutional defect is created through two bi-dimensional supercell models with different size.

The supercell scheme has some attractive features:

1. It is of wide applicability, and it may be adopted, in principle, to model both bulk and surface defects of ionic, covalent, metallic, and molecular systems.
2. It is conceptually simple. The size of the supercell depends on an expansion matrix that consists of integers. The matrix is  $2 \times 2$  or  $3 \times 3$  according to the dimensionality of the periodic system.
3. It allows for a proper definition of the defect formation energy, as will be discussed in the next section.
4. Properties of the defective solid can be calculated easily.

Obviously, computed properties are required to converge with the supercell size. This internal consistency check is important for estimating the interaction between defects in neighboring cells. In fact, two kinds of limitations in the model exist, which correspond to two different levels of complexity:

1. The supercell size must be such as to contain the defect zone, which includes all atoms involved in the structural and electronic relaxation.
2. The distance among defects must be large enough to reduce their electrostatic interaction to negligible values.

In the latter case, it must be distinguished whether the defect is neutral or charged. For neutral defects, the supercell scheme is expected to converge quickly to the isolated defect limit (we will see later that “quickly” can imply large supercells).

For charged defects, the electrostatic energy of the supercell diverges, and approximations must be adopted to neutralize the unit cell to cancel the interaction between neutralized defects. The treatment of charged defects within periodic boundary conditions is still a partially unsolved problem, and *ad hoc* solutions have been proposed like, for instance, the corrective schemes proposed in Leslie and Gillan,<sup>250</sup> Makov and Payne,<sup>251</sup> and Gerstmann et al.<sup>252</sup>

Finally, the supercell shape should be such as to exploit the point symmetry of the defect as far as possible.

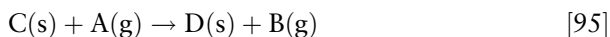
When a suitable supercell model of the defective system is devised, three major issues have to be faced:

1. The determination of the equilibrium geometry, that is the evaluation of structural effects on the surroundings because of the presence of the defect (relaxation/reconstruction).
2. The description of the electronic structure in the defect region.
3. The estimate of the defect formation energy.

Points 1 and 2 will be highlighted with a number of examples, whereas point 3 deserves some more comments.

## Defect Formation Energy

The formation process of a point defect can be described by the reaction:



where C is the perfect crystal and D is the defective system. A and B are the reactants and the products, respectively, and are typically atomic or molecular species, usually in the gas (g) phase. For instance, in a substitutional defect, A is the substituting atom for atom B. In the process of creating a vacancy, A is absent and B is the leaving moiety. The formation energy is then given by the following expression:

$$\Delta E^n = E^n(D) - nE(C) + E(B) - E(A) \quad [96]$$

i.e., the energy difference between the infinite system with the defect ( $E(D)$ ) and without ( $E(C)$ ), plus the energy difference of the atomic reactants,  $E(A)$ , and products,  $E(B)$ , with  $n$  being the ratio between the volume of the supercell  $S_n$  and that of the perfect crystal primitive cell.

The defect formation energy can also be considered as the sum of two contributions: a term originating from the creation of the defect in the perfect crystal,  $\Delta E^n(\text{dc})$ , and a term that accounts for the relaxation effects because of the perturbation caused by the presence of the defect,  $\Delta E^n(\text{rel})$ . That is,

$$\Delta E^n = \Delta E^n(\text{dc}) + \Delta E^n(\text{rel}) \quad [97]$$

Within the supercell approach,  $\Delta E^n$  should tend to a well-defined limit with increasing the supercell size:

$$\Delta E = \lim_{n \rightarrow \infty} \Delta E^n \quad [98]$$

For a defect to be considered as isolated, sufficiently large supercells must be adopted to avoid spurious interactions among neighboring defects because of both relaxation/reconstruction and long-range electrostatic effects.

To allow a consistent definition of the defect formation energy, the computational method must satisfy a size-extensivity criterion: Supercells of different size, for the perfect crystal, must provide the same value per formula unit for all properties.

## Examples

The applicability of ab initio periodic methods to the study of defects will be illustrated through a few examples, the first of which is a neutral vacancy in LiF. This defect is typical in ionic solids, called F-center, because its presence produces a color change (F comes from the German: *Farbe*) in the bulk material.

The second example concerns trapped hole centers in alkaline earth oxides (BeO, MgO, CaO, and SrO). These neutral defects essentially consist of the substitution of a monovalent cation (H/D, Li, Na, K) for one of the bivalent cations (Be, Mg, Ca, Sr). Thus, one electron is missing, so that an electron hole is expected to be localized and “trapped” at the substitutional cation. In both cases, the defect is paramagnetic, and in the second one, part of the original lattice symmetry is lost.

A third example, carbon substitution in bulk silicon will compare the cluster and supercell approaches.

### *F-center in LiF*

An F-center consists of an electron trapped at a negative-ion vacancy within the crystal. It is a paramagnetic defect, and its presence in ionic compounds has been the object of extensive and systematic experimental investigations, mainly by means of EPR and ENDOR techniques. Here we consider an F-center in LiF.<sup>140,253</sup> Computational details and references to experimental data are reported elsewhere in the original papers.<sup>140,253</sup>

This defect is a good case for size-extensivity checks, because it is simple, with the full cubic symmetry of the unperturbed lattice. Moreover, a relatively small basis set can be adopted because light atoms are involved, so it has been possible to consider supercells containing up to 256 atoms (or 128 primitive cells,  $S_{128}$ ).

Energy data are given in Table 21. The fourth column indicates that the defect formation energy converges rapidly even with relatively small supercells. The second column is reported just to show that the total energy per LiF pair of the perfect crystal is independent of the size of the supercell.

Relaxation effects are small as shown in Table 22 for  $S_{32}$ . Partial relaxation was allowed by including progressively up to the third nearest-neighbors

**Table 21** Effect of the Supercell Size on the Defect Formation Energy,  $\Delta E^n$  (in kJ/mol)

$n$	Lattice	$E^n(\text{LiF})^a$	$E^2(\text{LiF})^b$	$E^{n-1}(\text{F-c})^c$	$\Delta E^{n-d}$
4	P <sup>e</sup>	-428.221800	-107.055450	-328.590048	674.0
8	F <sup>f</sup>	-856.443601	-107.055450	-756.811638	674.6
16	I <sup>g</sup>	-1712.887207	-107.055450	-1613.255239	674.6
27	F	-2890.497155	-107.055450	-2790.865051	674.9
32	P	-3425.774381	-107.055449	-3326.142400	674.6
64	F	-6851.548839	-107.055451	-6751.916825	674.7
108	P	-11561.988593	-107.055450	-11462.356596	674.6
125	F	-13381.931279	-107.055450	-13282.299274	674.7
128	I	-13703.097661	-107.055450	-13603.465658	674.7

<sup>a</sup>Total energy (hartree) of a perfect LiF supercell.

<sup>b</sup>Total energy (hartree) per LiF formula unit.

<sup>c</sup>Total energy (hartree) of the F-center defective system.

<sup>d</sup>Referred to the unrelaxed defect geometry.

<sup>e</sup>Primitive lattice.

<sup>f</sup>Face-centered lattice.

<sup>g</sup>Body-centered lattice.

**Table 22** Relaxation Effects in F-Centred LiF for a  $S_{32}$  Supercell Model

	$N^a$	$\Delta E^c$	$I^d$	$II^d$	$III^d$
First nearest neighbours	6 (6) <sup>b</sup>	-1.7	0.028		
Second nearest neighbours	18 (12)	-2.3	0.035	0.011	
Third nearest neighbours	26 (8)	-2.3	0.035	0.011	-0.001
Fully relaxed	63	-2.6	0.040	0.013	-0.003

<sup>a</sup>Number of relaxed atoms.<sup>b</sup>Number of atoms in each *star* (set of atoms equidistant from the defect) of neighbors.<sup>c</sup>Gain in energy (in kJ/mol) with respect to the unrelaxed structure.<sup>d</sup>Displacements (in Å) of the stars of neighbors with respect to their position in the unrelaxed geometry.

of the F-center, in order to show the trend to the fully relaxed geometry. Data indicate that nearest and next-nearest neighbors move away from the defect center by a small amount, with the largest relaxation involving the nearest neighbors, whereas relaxation of the third nearest neighbors is negligible. Accordingly, the gain in energy caused by relaxation is just a few kJ/mol. Relaxation effects are thus negligible and die down quickly, so that the unrelaxed structure could safely be considered as the reference geometry.

One of the interesting features of LiF is that it has been a sort of model system in the interpretation of EPR and ENDOR data.<sup>254</sup> Experimental spectra have been fitted to model Hamiltonians, and hyperfine coupling constants up to the eighth nearest neighbors have been proposed.<sup>255</sup> In the calculation of the hyperfine coupling, it is then important to check the convergence of the spin density  $\rho^{\alpha-\beta}$  with the supercell size, not only at the center of the defect, but also at a relatively large distance from it. The spin density at the nuclear position for various supercells up to nine stars of neighbors is given in Table 23.

**Table 23** Effect of the Supercell Size on the Spin Density  $\rho^{\alpha-\beta}$  (in units of  $10^{-2}$  bohr<sup>-3</sup>) at the Nuclei of the Indicated Atoms for the F-Centre in LiF

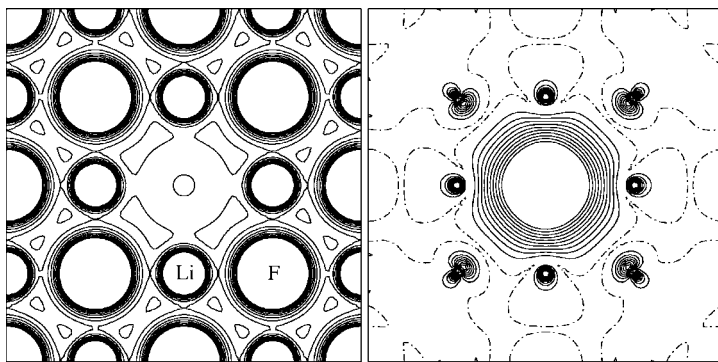
$n$	$F_C^a$	I $Li_{100}^b$	II $F_{110}$	III $Li_{111}$	IV $F_{200}$	V $Li_{210}$	VI $F_{211}$	VII $F_{220}$	VIII $Li_{221}$	IX $Li_{300}$
4	2.074	4.431	7.204	0.054						
8	2.080	2.272	3.641	0.021	0.143					
16	2.083	2.254	1.835	0.013	0.116	0.017				
27	2.082	2.255	1.837	0.006	0.022	0.008	0.034			0.000
32	2.083	2.255	1.829	0.007	0.044	0.009	0.038	0.096	0.000	
64	2.083	2.255	1.816	0.006	0.021	0.004	0.013	0.039	0.000	0.000
108	2.083	2.255	1.819	0.006	0.021	0.004	0.012	0.021	0.000	0.000
125	2.083	2.255	1.816	0.006	0.021	0.004	0.011	0.020	0.000	0.000
128	2.083	2.255	1.816	0.006	0.021	0.004	0.011	0.020	0.000	0.000

Nine sets of neighbors are considered (I-IX).

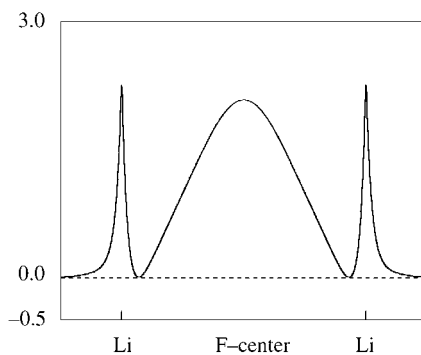
<sup>a</sup> $F_C$  is the anion vacancy.<sup>b</sup>Subscripts give the Cartesian coordinates of the vacancy neighbors in units of the cation-anion distance.

$S_{32}$  is the smallest supercell containing all atoms of interest. When  $S_{108}$  is considered, all  $\rho^{\alpha-\beta}$  at the nuclei up to the eighth nearest neighbors are numerically stable with respect to larger cells. The table shows that after the second neighbors, the spin density drops by two orders of magnitude.

Insights on the nature of the paramagnetic defect can be obtained by analysis of the electronic structure. Figure 48(a) shows the total charge (left) and the spin density (right) of the F-center in LiF ( $S_{16}$ ) obtained at the UHF level. The spin density map shows that the unpaired electron is localized at the vacancy site, whereas the spin density profile [Figure 48(b)] gives an indication



(a)



(b)

**Figure 48** (a) Total charge and spin density maps for an F-center in LiF as obtained at the UHF level of theory. The section is parallel to the (100) plane through the defect. The separation between contiguous isodensity curves is 0.01 and 0.001 bohr<sup>-3</sup> for the electron charge and spin density, respectively. The density range is 0.01–0.1(charge) and –0.01–0.01(spín). Continuous, dashed, and dot-dashed lines denote positive, negative, and zero values. (b) Spin density profile (in 10<sup>-2</sup> bohr<sup>-3</sup>) along a line connecting the F-center to opposite nearest neighboring Li ions. Ticks indicate nuclear positions.

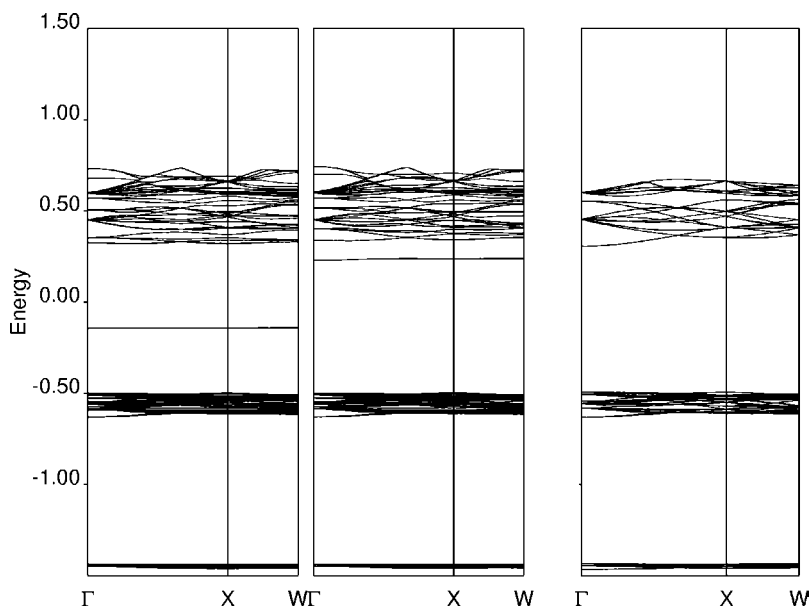
**Table 24** Mulliken Charge  $q$  (in electrons) and Spin Density  $\rho^{\alpha-\beta}$  (in units of  $10^{-2}$  bohr $^{-3}$ ) at the Vacancy Site (F-center) Obtained with Various Hamiltonians ( $S_{16}$ )

	UHF	SPZ	BLYP	PBE	B3LYP
$q$	1.05	0.89	1.06	0.93	1.03
$\rho^{\alpha-\beta}$	2.08	1.66	1.86	1.73	1.89

of the spread of the spin density within the vacancy. The Mulliken charge of the F-center (first row of Table 24) confirms that the amount of electronic charge that can be attributed to the vacancy is close to one.

One of the possible effects of the presence of a defect is the appearance of localized states in the band gap of the perfect host. Trapping and releasing electrons to and from these states requires less energy than exciting electrons from the top of the valence band to the bottom of the conduction band of the perfect crystal.

Figure 49 shows the band structure of the F-center in LiF for the  $S_{16}$ -supercell model obtained at the UHF level along with the band structure of bulk LiF. Alpha and beta electrons are described by different sets of orbitals. Two band structures are obtained for the  $\alpha$ - and  $\beta$ -spin states. The shape of the bands is similar to those of the perfect system, but a new band appears in the



**Figure 49** Band structure for  $\alpha$ - (left) and  $\beta$ -spin (middle) of the LiF F-center and perfect LiF (right). Plotted data refer to a UHF calculation with a  $S_{16}$  supercell.

band gap, because of the state associated with the F-center. The analysis of that band through the projected density of states shows that it is essentially associated with a hydrogen s-like state.

When other Hamiltonians are considered, the qualitative picture of the defect remains essentially unaltered (we will see that this is not the case for the trapped hole centers in alkaline earth oxides), as shown in Table 24, where Mulliken charge  $q$  and spin density  $\rho^{\alpha-\beta}$  at the vacancy site obtained with five different Hamiltonians is reported.<sup>140</sup>

The general picture emerging from the UHF and UDFT data confirm that the F-center is nearly totally localized at the anion vacancy and that the various functionals provide similar descriptions. Note that UDFT tends to spread the spin density onto the nearest neighbors, in particular with the local density approximation. This picture is also confirmed by spin density maps and profiles (not reported here, see Mallia et al.<sup>253</sup>).

In summary, when we consider the applicability of the supercell model to the LiF defect, we see that in this case, we are in a favorable position, because:

1. Nuclear relaxation is small.
2. Electronic perturbation is confined within nearest neighbors, and actually from the next-nearest neighbors the electronic density is indistinguishable from the perfect bulk on a “normal” scale.
3. The defect is neutral and conforms to the cubic symmetry, as no long-range electrostatic defect-defect and defect-bulk effects take place.

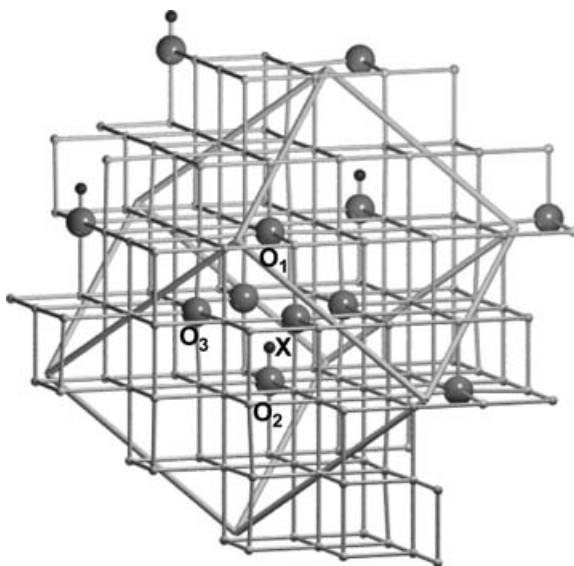
Because of these conditions,  $S_8$  or  $S_{16}$  would be large enough, if we were interested only in the properties of the F-center itself. As, however, we are also investigating the hyperfine coupling of the unpaired electron with up to the eighth nearest neighbors, a much larger supercell is required. Thus, the point is then raised concerning the definition of the extent of the “perturbed zone”: Because the hyperfine coupling is detected up to the seventh nearest neighbors by EPR/ENDOR experiments, the perturbed zone is obviously large enough to include up to the seventh nearest neighbors. However, the amount of spin density present at the nuclei farther than the next-nearest neighbors is small (see Table 23), and examining these interactions becomes interesting only in this special case, where extremely accurate and sensitive experiments are available. As other properties are concerned, such as the defect formation energy (Table 21) or the nuclear relaxation (Table 22), the perturbed zone is much smaller, and a  $S_8$ - $S_{16}$  supercell is large enough for most purposes, as previously mentioned.

Let us now consider a slightly more complicated defect, which is common in all alkaline earth oxides.

### *Trapped Hole Centers in Alkaline Earth Oxides*

*The Cubic Oxides: Li in MgO.* Ionizing radiation produces a variety of trapped hole centers in alkaline earth oxides at low temperature. In the cubic





**Figure 50** Schematic picture of a  $S_{16}$  supercell model of the trapped hole center in cubic alkaline earth oxide:  $\text{MO}:[\text{X}]^0$  (where  $\text{M} = \text{Mg}, \text{Ca}, \text{Sr}$  and  $\text{X} = \text{H}, \text{Li}, \text{Na}$ ).  $\text{O}_1$  is the oxygen ion at which the hole is trapped.

case ( $\text{MgO}$ ,  $\text{CaO}$ ,  $\text{SrO}$ , excluding  $\text{BeO}$ , which has a lower symmetry), the defect has axial symmetry along one of the main directions of the cubic lattice. The presence of various impurity ions occupying different positions along the hole-vacancy axis produces neutral defects indicated as  $[\text{X}]^0$ , where  $\text{X}$  represents  $\text{H}$ ,  $\text{Li}$ ,  $\text{Na}$ , or  $\text{K}$  (see Figure 50). When  $\text{X}$  stands for  $\text{H}$  ( $\text{D}$ ), the defects are also denoted as  $\text{V}_{\text{OH}}$  ( $\text{V}_{\text{OD}}$ ) centers. These trapped holes have been the subject of theoretical investigations at the UHF level of theory.<sup>140,175,256–259</sup> Here we only discuss the  $\text{MgO}:[\text{Li}]^0$  case and refer to the cited papers for additional information on experimental evidences and computational details not included here.

The convergence of the defect data with the supercell size must be checked again. Table 25 reports the defect formation energies of the  $S_8$ ,  $S_{16}$ ,  $S_{32}$ , and  $S_{64}$  supercells, with respect to the atomic energies of the species involved in the substitution ( $\text{Li}$ :  $-7.429609$  hartree;  $\text{Mg}$ :  $-199.602732$

**Table 25** Dependence of the  $\text{MgO}:[\text{Li}]^0$  Defect Formation Energy (in  $\text{kJ/mol}$ ) with Respect to the Supercell Size and Geometry Relaxation

	$S_8$	$S_{16}$	$S_{32}$	$S_{64}$
Li only	465.3	481.3	488.0	492.4
Full relaxation	403.8	400.8	401.1	400.4

hartree), when only the Li atom is allowed to relax (first row) and for the fully relaxed defective structure (second row).

We can compare the first row of Table 25, where the defect formation energy increases by about 40 kJ/mol in going from the smallest to the largest supercell, with the LiF case (Table 21), where convergence is already reached in the unrelaxed defective structure.

The main difference here is that the substitution of a Mg atom with lithium lowers the crystalline symmetry from cubic to tetragonal and generates a dipole moment within the cell, as a consequence of the Li displacement. A long-range dipolar defect-defect interaction originates then among defects in neighboring cells, and larger supercells are needed to compensate for it.

However, when all atoms in the cell are allowed to relax, a dramatic change is observed in the defect formation energy, which decreases by 61.5, 80.5, 86.9, and 92.0 kJ/mol for the four supercells. Again, this behavior can be compared with the case of the F-center in LiF, where relaxation effects are essentially absent (2–3 kJ/mol). In the fully relaxed case, the convergence of the defect formation energy is much faster (3 kJ/mol in going from  $S_8$  to  $S_{64}$ ), which shows that structural relaxation is an effective mechanism to screen and minimize long-range electrostatic interactions induced by the dipolar nature of the defect center.

Thus, geometry optimization plays a crucial role, which has been analyzed through partial optimizations within the  $S_{64}$  supercell, by including an increasing number of neighbors in the process. Results are reported in Table 26.

When only the atomic position of lithium is optimized, it migrates by  $\Delta z = 0.279$  Å away from  $O_1$ , with a corresponding relaxation energy of 31.3 kJ/mol; when the first nearest neighbors of Li along the  $z$  axis ( $O_1$ ,  $O_2$ )

**Table 26** Convergence of the  $\text{MgO}:[\text{Li}]^0$  Defect Formation Energy (in kJ/mol) in the  $S_{64}$  Supercell with Respect to Structural Relaxation (in Å) Allowed up to the Fourth Nearest Neighbors of the Defect in the Otherwise Unrelaxed Structure (changes in the geometry around the defect are also reported)

	Unrelaxed	Li	Li + $O_1^a$ + $O_2^a$	Li + O(6) <sup>c</sup>	Li + O(6) + Mg(12)	Li + O(6) + Mg(12) + O(8)	Li + O(6) + Mg(12) + O(8) + Mg(6)	Fully relaxed
$\Delta E$	523.7	492.4	484.2	477.2	428.3	426.3	411.2	400.4
Li $\Delta z$		0.279	0.312	0.291	0.251	0.251	0.245	0.236
$O_1^a$ $\Delta z$			−0.032	−0.031	0.003	0.002	−0.017	0.007
$O_2^a$ $\Delta z$			0.095	0.075	0.059	0.060	0.060	0.059
$O_3^b$ $\Delta z$				−0.035	−0.017	−0.015	−0.013	−0.012
$\Delta r$				0.035	0.040	0.040	0.040	−0.046

<sup>a</sup>Oxygen atoms above and below the Li atom along the defect axis (see Figure 50).

<sup>b</sup>Equatorial oxygen atom (see Figure 50).

<sup>c</sup>O(6) includes also the four equatorial oxygen atoms around the defect in addition to  $O_1$  and  $O_2$ .

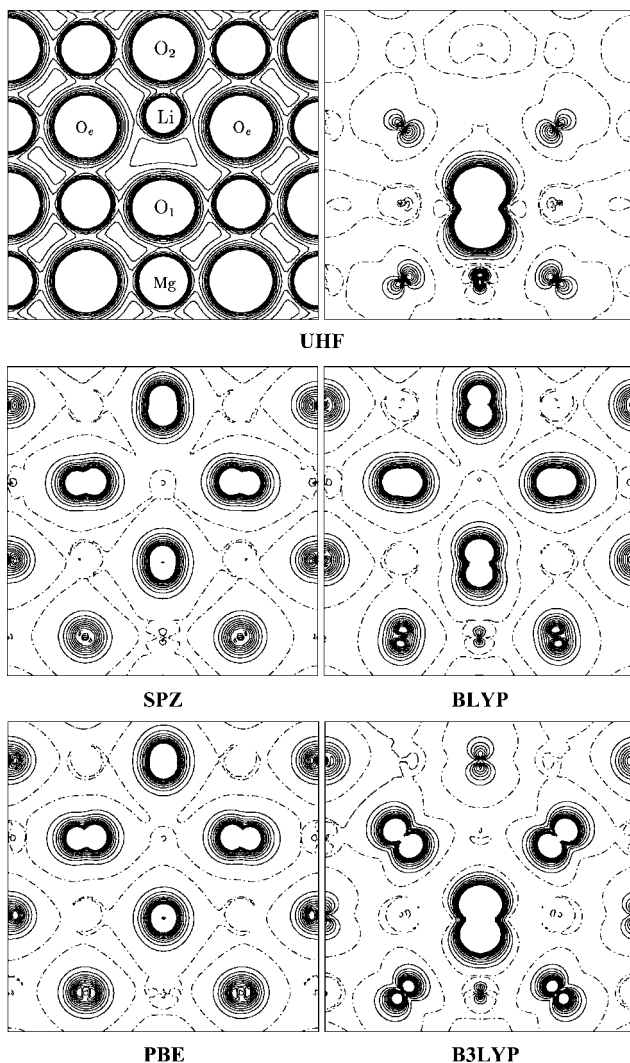
are also relaxed, the energy gain is 8.2 kJ/mol and displacements are  $\Delta z(\text{O}_1) = -0.032 \text{ \AA}$ ,  $\Delta z(\text{O}_2) = 0.095 \text{ \AA}$ . When also including the four equatorial  $\text{O}_3$  atoms, the formation energy decreases further by 7.0 kJ/mol and the displacements of  $\text{O}_3$  are  $\Delta z = -0.035 \text{ \AA}$  and  $\Delta r = 0.035 \text{ \AA}$  (in the direction perpendicular to the axial defect). Including in the geometry optimization the next star of neighbors (12 atoms of magnesium) reduces the formation energy by 48.9 kJ/mol. Full relaxation leads to a total gain of 123.3 kJ/mol. Interestingly, a large fraction of the energy gain is caused by the relaxation of the Mg ions (48.9+15.1 kJ/mol, to be compared with 15.2 + 2.0 kJ/mol for the first and third oxygen neighbors), the reason being that  $\text{Mg}^{2+}$  is smaller than  $\text{O}^{2-}$  and mobile in the cage of its six oxygen neighbors, which is rigid, because the oxygen ions are in contact. The smaller the cation, the larger this effect, as in the case of  $\text{Be}^{2+}$  (see next example).

The relaxation mechanism is simple and similar in all alkaline earth oxides: When the alkali metal ion replaces an alkaline earth cation, it relaxes from the perfect lattice position toward the oxygen ion ( $\text{O}_2$ ) along the axial direction ( $z$  axis), which brings a formal +2 charge; the electron hole localizes at the opposite oxygen ( $\text{O}_1$ ), which in turn relaxes away from the X monovalent ion.

The  $\text{S}_{32}$  supercell is certainly adequate for describing this defect, and it will be used for the analysis of its electronic and magnetic features as follows. The electron charge and spin density, computed at the UHF level of theory and shown on the top of Figure 51, illustrates the two main effects of substituting a Li ion for one Mg in  $\text{MgO}$ . Li binds to one of the neighboring O ions ( $\text{O}_2$ ) and, at the same time, like all other monovalent cations, acts as a dopant, which causes the formation of a trapped electron hole, well localized at the opposite ion ( $\text{O}_1$ ). The spin density map permits us to appreciate the localization of the unpaired electron at  $\text{O}_1$ , in a  $p_z$ -type state, with minor spin polarization on the neighboring atoms. This analysis is supported by the Mulliken population data, having a net charge of about +1 electrons for Li, -1 for  $\text{O}_1$ , and -2 for  $\text{O}_2$ . According to the spin density Mulliken analysis, the spin moment of  $\text{O}_1$  is close to 1 electron, whereas it is almost null on Li and on the other neighboring oxygen ions.

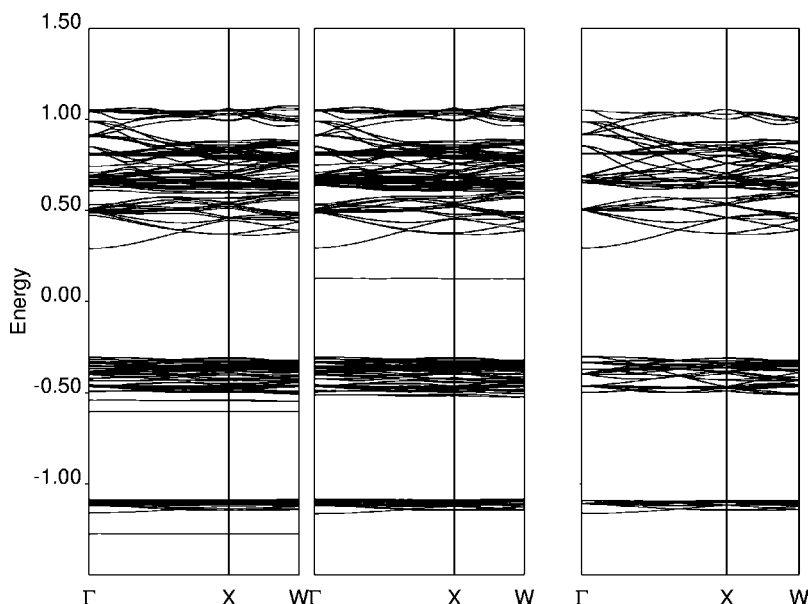
The different role of  $\text{O}_1$  and  $\text{O}_2$  among the O ions of the crystal is reflected in the band structure. In particular, the  $p$  states of both atoms split off the valence band, which is essentially contributed to by the  $p$  orbitals of the bulk oxygen atoms, as shown in Figure 52. A detailed analysis of the band structure shows that the most stable  $\alpha$  states are associated with the  $p_z$  (lower in energy) and the  $p_x, p_y$  orbitals of  $\text{O}_1$ ; the stabilization is a consequence of the lack of interelectronic repulsion with the corresponding  $\beta$  electron. The empty  $\beta - p_z$  state corresponds to the hole level and lies in the band gap.

As stated above, the UHF electronic structure of the trapped hole defect indicates the localization of the unpaired electron at the  $\text{O}_1$  atom. However, this picture changes significantly when different Hamiltonians are



**Figure 51** Total charge and spin density maps for  $\text{MgO}:[\text{Li}]^0$  as obtained at the UHF level and with various DFT Hamiltonians. The section illustrated is parallel to the (100) plane through the defect. The separation between contiguous isodensity curves is 0.01 and 0.001  $\text{bohr}^{-3}$  for the electron charge and spin density, respectively. The density range is 0.01–0.1(charge) and –0.01–0.01(spin map). Continuous, dashed, and dot-dashed lines denote positive, negative, and zero values.

considered.<sup>140</sup> Table 27 reports the net atomic charges and spin moments on Li, O<sub>1</sub>, O<sub>2</sub>, and O<sub>3</sub>, in  $\text{MgO}:[\text{Li}]^0$ , as computed with UHF and four different DFT methods. The degree of localization as quantified by the spin moment of O<sub>1</sub> (Table 27) decreases from 0.98 at the UHF level, to 0.41 for B3LYP, to



**Figure 52** Band structure for  $\alpha$ - (left) and  $\beta$ -spin (middle) of  $\text{MgO}:[\text{Li}]^0$  and perfect  $\text{MgO}$  (right). Plotted data refer to a UHF calculation with a  $S_{16}$  supercell.

about 0.1 with other DFT methods. With DFT methods, this spread of the unpaired electron over the nearest and next-nearest oxygen neighbors is evident in the spin density maps (Figure 51) and causes the  $\text{O}_1$  ion to have a net charge similar to the other oxygen ions in the lattice, i.e., about  $-1.7$  electrons.

The different degree of localization produced by DFT methods also has important consequences on the atomic relaxation, with the Li ion being less strongly attracted to  $\text{O}_2$  than at the UHF level. These effects are less pronounced when the hybrid B3LYP method is adopted.

**Table 27** Net Atomic Charges ( $q$ ) and Spin Moments ( $\mu$ ) in  $\text{MgO}:[\text{Li}]^0$  Evaluated According to a Mulliken Partition of Charge and Spin Densities

Method	Li		$\text{O}_1$		$\text{O}_2$		$\text{O}_3$	
	$q$	$\mu$	$q$	$\mu$	$Q$	$\mu$	$q$	$\mu$
UHF	+0.99	0.00	-1.03	0.97	-1.91	0.00	-1.88	0.01
SPZ	+0.98	0.00	-1.71	0.09	-1.69	0.12	-1.69	0.11
BLYP	+0.97	0.00	-1.66	0.14	-1.70	0.10	-1.68	0.11
PBE	+0.98	0.00	-1.74	0.08	-1.71	0.10	-1.71	0.10
B3LYP	+0.98	0.00	-1.48	0.41	-1.80	0.01	-1.72	0.10

Label of atoms as in Figure 50. Data in electrons.

**Table 28** Calculated and Experimental Hyperfine Isotropic (*a*) and Anisotropic (*b*) Coupling Constants and Nuclear Quadrupole (*P*) Coupling Constant for  $\text{MgO}:[\text{Li}]^0$ 

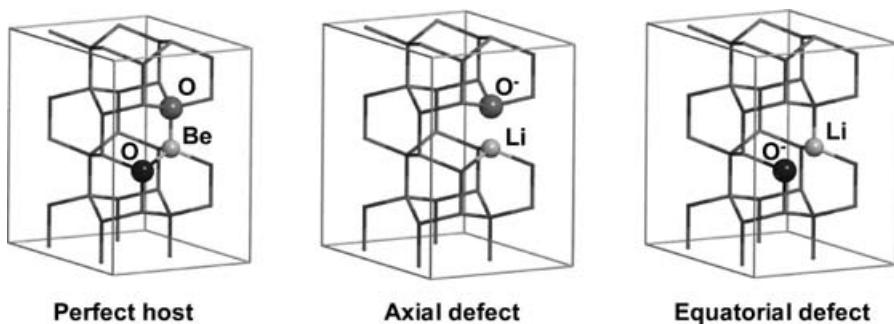
Method	<i>a</i>	<i>b</i>	<i>P</i>
UHF	-2.393	2.258	-0.017
SPZ	-3.691	-0.129	-0.003
BLYP	-3.078	0.092	-0.003
PBE	-3.979	-0.202	-0.010
B3LYP	-3.339	0.480	-0.004
Exp. <sup>a</sup>	-4.539	2.313	-0.014

<sup>a</sup>Experimental data taken from Abraham et al.<sup>260</sup> and Chen and Abraham.<sup>261</sup>

Magnetic coupling constants determined by EPR and ENDOR techniques permit a direct comparison with experimental data. Table 28 shows that, in the particular case of the Li defect, the agreement is reasonable for the UHF result, where the hole is localized at  $\text{O}_1$ . For the other Hamiltonians, the disagreement increases in parallel with the delocalization of the hole.

*The BeO Case.* The Li-trapped hole center in BeO ( $\text{BeO}:[\text{Li}]^0$ ), is slightly more complicated than the corresponding defect in MgO, and it presents some new features we must now consider. At variance with respect to the other alkaline earth oxides, BeO has a wurzite-like crystalline structure. The cation is fourfold coordinated, with one Be-O distance (the axial or vertical one) slightly different from the other three that are equivalent. Substitution of  $\text{Li}^+$  for one  $\text{Be}^{2+}$  ion in the hexagonal structure generates an electron hole that can be localized either at the axial oxygen or at one of the three equatorial oxygen ions, which is indicated in Figure 53 as  $\text{O}^-$ .

These features make studying the relative stability of the two point defects interesting. Experiment<sup>262</sup> indicates that, when the crystal is irradiated at low temperature, the hole is trapped at an axial oxygen; electron holes at



**Figure 53** Schematic view of a  $S_{16}$  supercell model of the two different trapped hole centers in BeO.  $\text{O}^-$  indicates the oxygen ion where the electron hole is localized. Axial oxygen in light-gray, and equatorial oxygen in black.

**Table 29** Defect Formation Energy,  $\Delta E$ , and Relative Stability,  $\delta E$ , (in kJ/mol) of the BeO:[Li]<sup>0</sup> Axial and Equatorial Centers as Obtained with Supercells of Different Size

Supercell <sup>a</sup>	BeO $E_{tot}$	Axial BeO:[Li] <sup>0</sup> $E_{tot}$	Equatorial BeO:[Li] <sup>0</sup> $\Delta E$	Ax-Eq $\delta E$
S <sub>16</sub> (2 2 2)	-1435.23633	-1427.83826	677.9	46.7
S <sub>54</sub> (3 3 3)	-4843.92262	-4836.53959	638.4	7.4
S <sub>128</sub> (4 4 4)	-11481.89066	-11474.50982	632.6	3.8

Total energies in hartrees.

<sup>a</sup>(i j k) are the expansion coefficients of the primitive lattice basis vectors **a**, **b** and **c** to obtain  $S_n$ .

equatorial oxygens can be obtained via thermal excitation. The energy difference between the two configurations is then expected to be small. At variance with the cases previously discussed, here we are not interested in the absolute defect formation energy, but in the relative value for the two positions. How much this relative value depends on the supercell size must be checked carefully.

Calculations have been performed at the UHF level with full optimization of the positions of all atoms in the supercell. As for cubic alkaline earth oxides, at the UHF level of the theory, the electron hole is fully localized at O<sup>-</sup>.

In both the axial and the equatorial configurations, relaxation of Li<sup>+</sup> and O<sup>-</sup> is qualitatively similar to that observed for the other members of the series. However, the extent of the O<sup>-</sup> relaxation in the case of the axial defect,  $\Delta r = -0.23$  Å, is four times as large as in MgO:[Li]<sup>0</sup>, where  $\Delta r = -0.05$  Å. Relaxation of the Be<sup>2+</sup> ions is larger than that of the Mg<sup>2+</sup> ions, because the former can migrate more easily in the oxygen cage.

The formation energy obtained with three supercells of increasing size is reported in Table 29, which shows that the convergence of the formation energy with the supercell size is much slower than the one reported for MgO:[Li]<sup>0</sup> (Table 25), where the formation energy is stable already for small supercells.

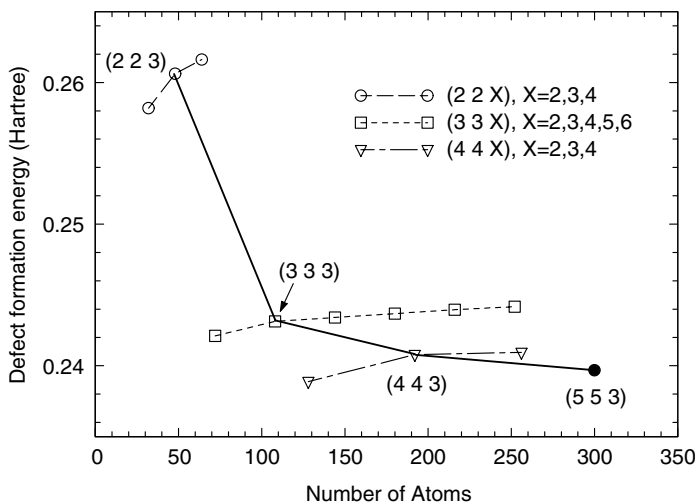
Many reasons exist for this different behavior: In MgO, the point symmetry remains high (tetragonal); the mobility of ions (in particular of the cations) is low, so that relaxation effects are not large; the site coordination of the ions remains essentially octahedral. In the present case, Be is almost free to move within the oxygen cage, the point symmetry of the defect is lower than in the MgO case, and the site symmetry of each ion is also lower so that nonzero low-order multipoles are generated by the defect. As a consequence, larger supercells are required.

Indeed, the defect formation energy of the axial defect is still changing by more than 5 kJ/mol when going from S<sub>54</sub> to S<sub>128</sub> (Table 29). The screening mechanism is more effective in the equatorial direction along which the defect is better accommodated. Table 29 shows that the equatorial defect is more

stable than the axial one for the supercells considered here, in disagreement with experimental evidence,<sup>262</sup> with the energy difference between the two configurations, however, small and decreasing when increasing the supercell size. In supercells as large as  $S_{250}(5\ 5\ 5)$  or  $S_{432}(6\ 6\ 6)$ , the stability order might be reversed. Unfortunately, they correspond to large unit cells containing 500 and 864 atoms, respectively, that would make the calculation demanding.

To check the mutual interactions of defects in different supercells, the formation energy of the axial center has been studied in more detail. A series of 13 supercells with increasing size (from 32 to 300 atoms) and different shape has been considered. In fact, given the hexagonal unit cell, the supercell may be increased by enlarging it in either the axial or the equatorial direction, or both. Figure 54 shows, graphically, the dependence of the formation energy on the number of atoms in the supercell. Each  $\text{Li}^+ - \text{O}^-$  pair can be viewed as a dipole oriented along the  $c$  axis. Making the supercell larger in the equatorial directions (see the  $S_n(i\ i\ 3)$  series connected by the continuous line in the figure) reduces the dipole lateral repulsion and decreases the formation energy. On the contrary, when separating the dipole along the axial direction, the cooperative interaction between dipoles decreases and  $\Delta E$  grows slightly (see, for example, the  $S_n(3\ 3\ i)$  series). By growing the supercell equatorially,  $\Delta E$  converges more rapidly with respect to axial growth, as can be seen comparing the  $S_n(3\ 3\ i)$  and  $S_m(4\ 4\ i)$  series.

The present example shows that it is not difficult to find situations where the supercell size must be large, and the accuracy of the calculations, referring to different supercells, must be high.



**Figure 54** Dependence of the formation energy of the axial  $\text{BeO}:[\text{Li}]^0$  defect with size and shape of the adopted supercell. (i j k) are the expansion coefficients of the primitive lattice basis vectors.



It is hard to believe that a cluster model would give reliable results when relaxation, polarization, and long-range electrostatic interactions play an important and competitive role, because border effects are expected to alter the relative energies by an amount that is orders of magnitude larger than the relative stabilities being investigated.

### *Carbon Substitution in Silicon: A Supercell Versus Cluster Investigation*

As a last example, we consider another simple defect: the carbon substitution in bulk silicon.<sup>263,264</sup> In this case, however, we will not only consider the convergence properties of the supercell approach but also compare the results of the cluster and supercell schemes. Calculations were performed at the HF level with a 6-21G basis set plus polarization functions for C and Si and a 2-1G basis set for H (the latter was used in the cluster calculations).

We consider first the convergence of the results with respect to the supercell size. Four supercells with 8 ( $S_4$ ), 16 ( $S_8$ ), 32 ( $S_{16}$ ), and 64 ( $S_{32}$ ) atoms (the last is shown in Figure 55) are considered, all with the cubic symmetry. As in previous examples, to investigate how far the perturbation propagates, an increasing number of defect neighbors has subsequently been allowed to relax in each supercell.

The stars of neighbors completely contained in the unit cell, and the number of atoms belonging to a star (in parentheses) for the four supercells are I (4), I (4), I + II (4 + 12), and I + II + III + V (4 + 12 + 12 + 12), respectively. Notice that in  $S_8$ , next-nearest neighbors are shared with defects in neighboring cells. In the  $S_{32}$  supercell, the atoms of star IV are at special positions and cannot relax.

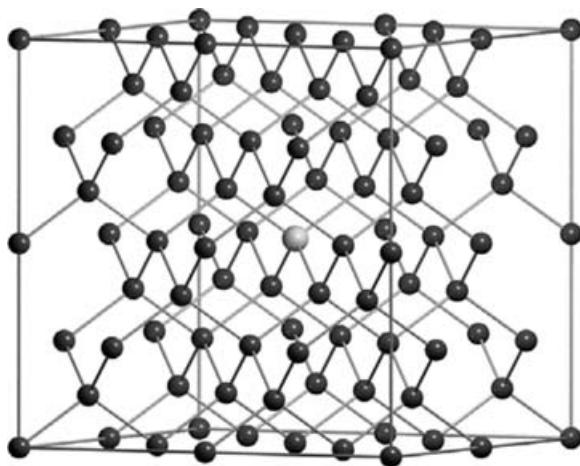


Figure 55 A 64-atom supercell ( $S_{32}$ ) model of a carbon impurity in bulk silicon.

**Table 30** Defect Formation Energies (in kJ/mol) for a Carbon Substitutional Impurity in Silicon, as a Function of Supercell and Cluster Size

	Unrelaxed	I	II	III	IV	V	Fully Relaxed
Supercell							
S <sub>4</sub>	222.6						105.1
S <sub>8</sub>	223.6	100.4					92.3
S <sub>16</sub>	223.2	59.0	24.4				9.7
S <sub>32</sub>	223.0	60.1	28.5	25.8		12.5	12.4
Cluster							
CSi <sub>4</sub> H <sub>12</sub>	220.5	56.3					
CSi <sub>34</sub> H <sub>36</sub>	230.1	97.2	49.1	50.6			
CSi <sub>86</sub> H <sub>76</sub>	225.4	68.1	39.8	40.1	41.5	6.5	

Relaxation effects are taken into account by including increasing stars of neighbors of the carbon impurity.

The substitutional defect formation energy,  $\Delta E$ , computed with respect to atomic energies according to Eq. [96], is reported in Table 30. The Si and C atomic energies are  $-288.812622$  and  $-37.654208$  hartree, respectively.

Carbon has both a higher electronegativity and a smaller covalent radius than does silicon, so a large charge transfer and atomic relaxation is foreseen. The former feature is expected to polarize the charge distribution in the cell, with a charge alternation between subsequent stars of neighbors and relatively strong electrostatic interactions among defects in neighboring cells, at least for small supercells. Table 31 shows that indeed this polarization occurs, with a huge difference in the Mulliken net charge of the central carbon ( $-0.7$  electrons) atom with respect to its four nearest neighbors. Charge oscillations, however, damp down rapidly; in the S<sub>32</sub> supercell, the Mulliken net charge of the second neighbors is as small as  $-0.018$  electrons, and third and fourth

**Table 31** Mulliken Net Charges for a Carbon Substitutional Impurity in Silicon, as a Function of Supercell and Cluster Size

	C	Si <sup>I</sup>	Si <sup>II</sup>	Si <sup>III</sup>	Si <sup>IV</sup>	Si <sup>V</sup>
Supercell						
S <sub>4</sub>	-0.644	0.210	-0.065 <sup>a</sup>			
S <sub>8</sub>	-0.661	0.219	-0.035 <sup>a</sup>			
S <sub>16</sub>	-0.758	0.237	-0.017	0.001	0.002 <sup>a</sup>	
S <sub>32</sub>	-0.755	0.237	-0.018	-0.003	0.001	0.002
Cluster						
CSi <sub>4</sub> H <sub>12</sub>	-0.738	0.575 <sup>b</sup>				
CSi <sub>34</sub> H <sub>36</sub>	-0.723	0.166	0.102 <sup>b</sup>	0.118 <sup>b</sup>		
CSi <sub>86</sub> H <sub>76</sub>	-0.746	0.235	-0.038	-0.023	-0.027	0.109 <sup>b</sup>

Net charges refer to fully relaxed supercell structures whereas clusters are partially relaxed up to first, third, and fifth nearest-neighbors of CSi<sub>4</sub>H<sub>12</sub>, CSi<sub>34</sub>H<sub>36</sub>, and CSi<sub>86</sub>H<sub>76</sub>, respectively.

<sup>a</sup>Incomplete star of neighbors.

<sup>b</sup>Silicon atoms bonded to one or more hydrogens.

nearest neighbors are almost neutral. These data refer to the relaxed solution, but results are similar in the unrelaxed calculations, which helps us understand why the unrelaxed defect formation energy is essentially the same for the various supercells (first column in Table 30: numbers differ by less than 1 kJ/mol): The defect is already screened effectively at the nearest neighbor level, so that the electrostatic defect–defect interaction is essentially null even with small cells (in the  $S_4$  supercell, for example, the defect–defect distance is as small as 5.52 Å). This nearly perfect screening is also a consequence of symmetry: The defect zone has zero dipole and quadrupole.

When relaxation is taken into account, however, the defect perturbation propagates farther away, and convergence of the defect substitutional energy is slower.

The structural relaxation is large, and the process is dominated by modifications in the covalent network with no dependence on electrostatic effects. The largest structural change involves nearest neighbors as is shown in Table 32. The C–Si<sup>I</sup> bond length reduces by 0.312 Å with respect to the bulk Si–Si

**Table 32** Relaxation Effects in a Carbon Substitutional Impurity in Bulk Silicon as a Function of Supercell and Cluster Size

	N	I	II	III	IV	V
	R	2.390	3.903	4.577	5.520	6.015
Supercell						
$S_4$	1	–0.205				
$S_8$	1	–0.213				
	fr <sup>a</sup>	–0.217				
$S_{16}$	1	–0.261				
	2	–0.300	–0.066			
	fr <sup>a</sup>	–0.315	–0.078	0.018		
$S_{32}$	1	–0.259				
	2	–0.300	–0.065			
	3	–0.300	–0.065	–0.003		
	5	–0.313	–0.081	–0.002		–0.038
	fr <sup>a</sup>	–0.312	–0.082	–0.002		–0.040
Cluster						
CSi <sub>4</sub> H <sub>12</sub>	1	–0.208				
CSi <sub>34</sub> H <sub>36</sub>	1	–0.237				
	2	–0.279	–0.058			
	3	–0.278	–0.062	–0.016		
CSi <sub>86</sub> H <sub>76</sub>	1	–0.254				
	2	–0.292	–0.058			
	3	–0.291	–0.056	0.007		
	4	–0.291	–0.056	0.009	0.014	
	5	–0.306	–0.072	0.010	0.014	–0.038

This table gives the variation (in Å) of the distance  $R$  between the defect and its neighbors ( $N$  in Roman numerals). Relaxation effects are taken into account by including increasing stars of neighbors (Arabic numerals) of the carbon impurity.

<sup>a</sup>Full relaxation in the supercell.

distance, from 2.390 to 2.078 Å, whereas the Si<sup>I</sup>-Si<sup>II</sup> bond elongates to 2.446 Å. When only the nearest neighbors are allowed to relax, the energy gain is 120 kJ/mol for  $S_4$  and  $S_8$ , which increases to 163 kJ/mol in the  $S_{16}$  and  $S_{32}$  supercell. A further gain of about 30 kJ/mol is obtained by relaxing the next-nearest neighbors (this is possible only in  $S_{16}$  and  $S_{32}$ ). Relaxation of the third nearest neighbors is small and results in essentially no energy gain, whereas a further 13 kJ/mol is obtained by relaxing the fifth nearest neighbors. Full relaxation yields a further small energy gain.

Full relaxation of  $S_{16}$  reduces  $\Delta E$  to 9.7 kJ/mol, close to the value obtained with the fully relaxed  $S_{32}$  supercell. Thus, from the point of view of the convergence of defect properties, the present example is similar to the F-center in alkali halides previously discussed, where a relatively small supercell ( $S_{16}$ ) is large enough to “contain” all the defect perturbation.

Let us now investigate the same problem with a cluster model.<sup>264</sup> Clusters containing 5, 35, and 87 silicon atoms have been considered; they are referred to as small, medium, and large in Figure 56. Dangling bonds were saturated with hydrogen atoms along the Si-Si directions of the perfect crystal. The Si-H distance was kept constant at 1.46 Å during the optimization of all cluster atoms (for this reason clusters in the tables are not “fully relaxed”). The medium and large clusters were constructed subject to the constrain that silicon atoms at the cluster surface are connected to no more than two hydrogen atoms. The total number of atoms in the three clusters is then 17, 71, and 163. Computational conditions are the same for cluster and supercell models. The defect formation energies, Mulliken net charges, and structural relaxation effects as evaluated with the cluster models are given in the lower part of Tables 30, 31, and 32, respectively.

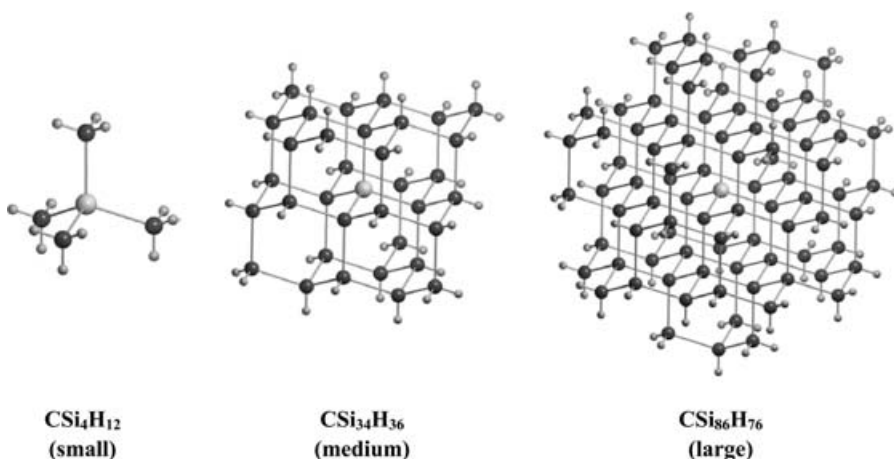


Figure 56 Adopted cluster models of the carbon impurity in bulk silicon.

The defect formation energies for the unrelaxed clusters are similar to those calculated by the supercell approach, even in the case of small clusters. The same screening mechanism, already discussed for the supercell models, is active also in the cluster calculations.

When relaxation effects are taken into account, the defect perturbation extends farther from the impurity into regions where the screening of border effects (H saturation) is only partial. The cluster defect substitutional energy differs from the  $S_{32}$  datum by -3.8 (small), 37.1 (medium), and 8.1 (large) kJ/mol when only the nearest neighbors are allowed to relax, and by 20.6 (medium) and 11.3 (large) kJ/mol when next-nearest neighbors can relax. Atomic charges in Table 31 confirm that the positive charge on the first Si neighbors compensates the negative charge on carbon, which makes the electrostatic contribution short-ranged. However, closer inspection of Table 31 shows that only for the largest cluster do net charges damp down as in the supercell. In fact, silicon atoms close to the border H atoms are positively charged, and this is expected to perturb their interaction to move toward the cluster center. Eventually, as regards the defect formation energies, the largest supercell and cluster models data are reasonably close to each other.

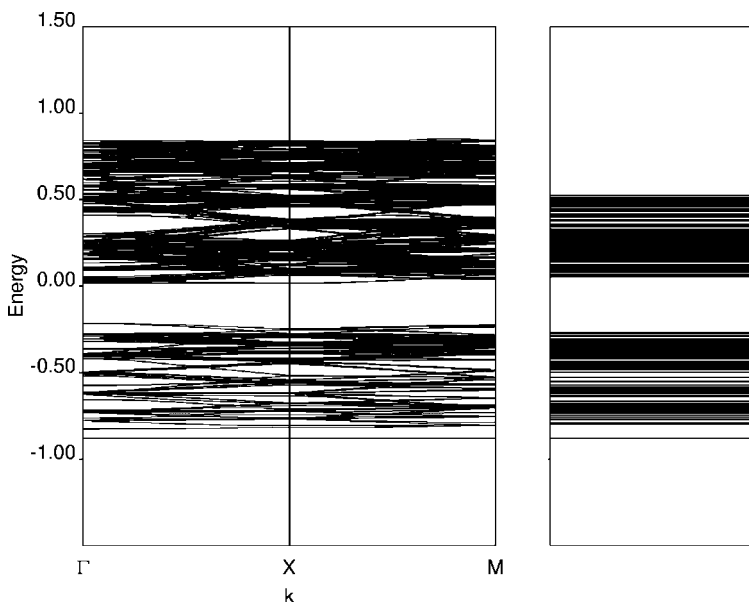
It is worth noting that if the hydrogen atoms are allowed to relax, a dramatic change in the structure originates, and the formation energy becomes negative. Fixed, saturating hydrogen atoms force the cluster to retain the memory of the bulk geometry.

The comparison between supercell and cluster approaches shows that the cluster size is crucial to derive accurate results and care must be taken to account for border effects. Only the largest cluster considered here (86 silicon atoms) is adequate for describing relaxation and defect formation energy properly, whereas a  $S_{16}$  (32-atom) supercell is already large enough.

Larger differences between the supercell and the cluster scheme may appear when other properties related to the infinite nature of the perturbed crystalline system are analyzed. In Figure 57, on the left, the band structure of defective  $S_{32}$  supercell is shown and, on the right, the energy levels of the  $CSi_{86}H_{76}$  cluster are shown (the energy scale is shifted in such a way that the 1s level of carbon of the two systems coincide). In the present example, no defect states are present in the gap; the computed band gap is 6.38 eV for the defective  $S_{32}$  supercell, to be compared with the 8.8 eV HOMO-LUMO gap of the cluster.

These examples show that the supercell approach is an accurate and, in many cases, relatively cheap tool for the study of neutral defects in crystalline systems, once properly gauged with respect to supercell size.

The supercell approach, as implemented in CRYSTAL, has been applied to the study of many different bulk and surface defective systems. These include Ca and Be substitution in bulk  $MgO$ ,<sup>265</sup> F-center in  $CaF_2$ ,<sup>266</sup> Fe doped  $NiO$ ,<sup>267</sup> Li doped  $NiO$ ,<sup>268</sup> V doped  $TiO_2$ ,<sup>269</sup> and Ti substitution in an all-silica Chabazite.<sup>270,271</sup> An example of reactivity of a surface defect has been



**Figure 57** Band structure of the fully relaxed carbon doped  $S_{32}$  silicon supercell (left) and energy levels of the  $CSi_{86}H_{76}$  cluster (right). In the cluster, atomic positions were relaxed up to the fifth nearest neighbors of the central carbon atom. The cluster energy levels are shifted so that the lowest level of both systems coincide.

reported by Orlando et al.,<sup>272</sup> where they studied the hydrogen abstraction from methane by Li doped MgO. An interesting application in the modeling of extended defects has been reported recently by Gruen et al.<sup>273</sup> in which the growth of crystalline diamond through planar defects was investigated. Convergence properties of the cluster model in the study of local perturbations has also been studied in ionic systems as in the case of bulk defects in MgO.<sup>274</sup>

---

## ACKNOWLEDGMENTS

We would like to acknowledge many fruitful discussions and interactions with Cesare Pisani. The authors wish also to thank Piero Ugliengo for providing some figures and assisting with the production of this chapter and for many helpful discussions.

---

## APPENDIX 1: AVAILABLE PERIODIC PROGRAMS

Many periodic codes are mentioned in the literature, mostly based on either a plane waves basis set and pseudopotentials or projector-augmented waves. Some of this software can be purchased, and the rest is available through collaboration with the main development team or downloadable

**Table A1.1** List of Available Solid State ab initio Computer Programs

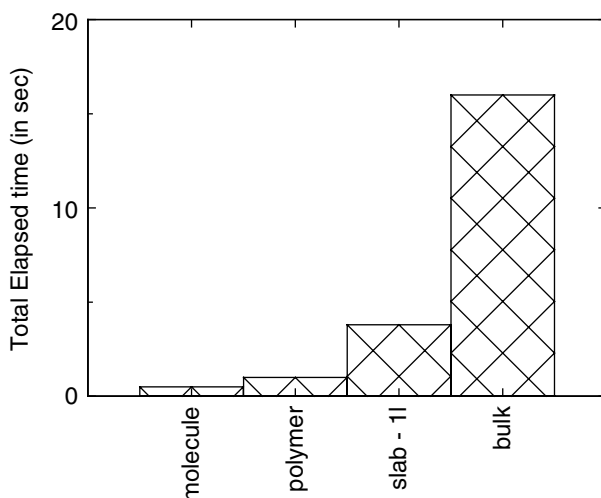
Program	Description	Web site	Refs
ABINIT	DFT(LDA,GGA); TD-DFT; NCPP; PW; PAW; TE; GO; PH	www.abinit.org	163
CASTEP	DFT(LDA,GGA); PP; PW; TE; GO; CP-MD	www.accelrys.com/mstudio/ ms_modeling/castep.html	275
CPMD	DFT(LDA,GGA); NCPP; PW; TE; GO; CP-MD	www.cpmc.org	276
Dacapo	DFT(LDA,GGA); PP; PW; TE; GO	www.fysik.dtu.dk/CAMP/daca- po.html	
DoD-Planewave	DFT(LDA,GGA); PP; PW; TE; GO	cst-www.nrl.navy.mil/people/ singh/planewave/	
FHI98md	DFT(LDA,GGA); NCPP; PW; TE; GO; BO-MD	www.fhi-berlin.mpg.de/th/ fhi98md/index.html	277
PARATEC	DFT(LDA,GGA); TD-DFT; NCPP; PW; TE; GO	www.nersc.gov/projects/para- tec/	278
PWSCF	DFT (LDA,GGA); DF-PT; NCPP; PW; TE; GO; PH; CP- MD	www.pwscf.org	279
VASP	DFT (LDA,GGA); USPP; PW; PAW; TE; GO; PH; CP-MD	cms.mpi.univie.ac.at/vasp/	280
CP-PAW	DFT(LDA,GGA); PAW; TE; GO; CP-MD	www.pt.tu-clausthal.de/~paw/ index.html	281
PWPAW	DFT(LDA); PAW; TE; GO	www.wfu.edu/~natalie/papers/ pwpaw/man.html	282, 283
QUICKSTEP/ CP2K	DFT(LDA,GGA); PP; hybrid GTO/PW	cp2k.berlios.de	284
SIESTA	DFT(LDA,GGA); PP; NTO; TE; GO; MD	www.uam.es/departamentos/ ciencias/fismateriac/siesta/	285
DMol <sup>3</sup>	DFT(LDA,GGA); AE; NTO; TE; GO	www.accelrys.com/mstudio/ ms_modeling/dmol3.html	286
LmtART	DFT; AE; LMTO; TE; GO; PH	physics.njit.edu/~savrasov/Pro- grams/index_lmtart.htm	164
FLEUR	DFT; AE; FLAPW; TE	www.flapw.de	287
WIEN2K	DFT(LDA,GGA); AE; FLAPW; TE; GO; PH	www.wien2k.at/	288
MOPAC2002	SE; TE; GO	www.schrodinger.com/ Products/mopac.html	
ADF2002 (BAND)	DFT(LDA,GGA); TD-DFT; AE; STO; TE	www.scm.com	289
Gaussian03	HF; DFT; AE; GTO; TE; GO; PH	www.gaussian.com	290
CRYSTAL03	HF; DFT; AE; PP; GTO; TE; GO	www.crystal.unito.it	
AE	All-electron basis set	BO	Born–Oppenheimer approximation
CP	Car–Parrinello method	DFT	Density functional theory
FLAPW	Fully linearized augmented plane wave	GGA	Generalized gradient approximation
GO	Geometry optimization	GTO	Gaussian-type orbitals
HF	Hartree–Fock	LDA	Local density approximation
MD	Molecular dynamics	NCPP	Norm-conserving pseudopotentials
NTD	Numerical type orbitals	PAW	Projector-augmented wave method
PH	Phonons	PP	Pseudopotentials
PT	Perturbation theory	PW	Plane waves
SE	Semi-empirical methods	STO	Slater-type orbitals
TD	Time dependent	TE	Total energy
USPP	Ultra-soft pseudopotentials		

with no restriction. A tentative list of public codes for solid state ab initio calculations is reported in Table A1.1, along with a list of the main functionalities and indication of the corresponding homepage.

## APPENDIX 2: PERFORMANCE OF THE PERIODIC PROGRAM CRYSTAL

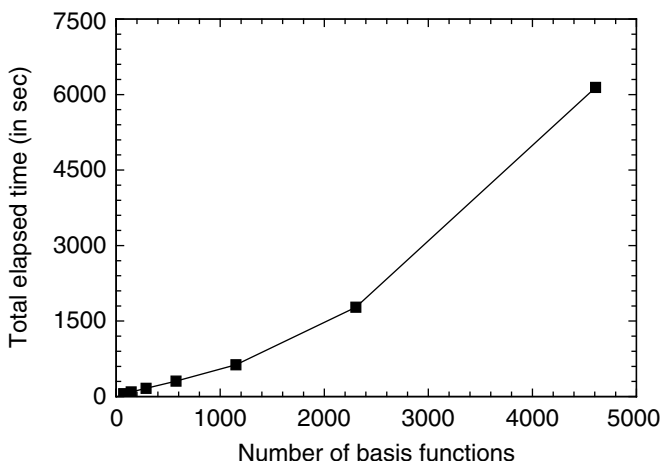
The cost of a quantum mechanical calculation depends on a large number of variables, of which the hardware available is certainly the most volatile, owing to rapid technological evolution. Also the compiler and compilation options that generate the binary have an influence on the performance of a program, as well as the more specific choice of the computational parameters controlling the accuracy of the computer program. For these reasons, the data reported in this appendix are intended to give only a rough indication of the cost of a periodic calculation and show how computational time scales with the dimensionality of the system, the approximation used, the unit cell, and the basis set size. To achieve this goal, we reconsider some of the examples illustrated before from the point of view of performance using CRYSTAL.<sup>56</sup>

We begin the analysis by comparing periodic calculations of increasing dimensionality with a single molecule calculation of a molecule consisting of the atoms in the unit cell. Data for MgO are represented in Figure A2.1, where the time required for the calculation of the HF total energy and wave function of a MgO molecule (0-D) is compared with the time required for the



**Figure A2.1** Computer time of a CRYSTAL calculation of the HF total energy and wave function for a MgO molecule, polymer, one-layer (001) slab, and bulk.





**Figure A2.2** Total elapsed time for HF total energy calculations of MgO cubic supercells with 8, 16, 32, 64, 128, 256, and 512 ions. The CRYSTAL program was compiled using the Intel Fortran Compiler IFC7.0 with the `-O2 -tp p7` options. Calculations were carried out with a Pentium Xeon 2.4-GHz single-processor computer, with 2-GB RAM, 512-KB cache, wide SCSI disks.

corresponding polymer (1-D), one-layer (001) slab (2-D), and bulk (3-D). We have used the same basis set, computational conditions, and geometry in all cases.

Cost increases exponentially with the dimensionality of the system, following approximately the progression 1:2:8:30 from 0-D to 3-D. In this simple case, even the bulk calculation takes only a few seconds on a small PC.

The cost of a bulk calculation is primarily a function of the unit cell size. Figure A2.2 shows the total time required for the calculation of the total energy with the MgO supercells that we used in the study of trapped-hole centers. We are considering the supercells before creating the defect. Therefore, the system possesses the full symmetry of the perfect crystal (48 point-symmetry operations). Calculations refer to the  $S_4$ ,  $S_8$ ,  $S_{16}$ ,  $S_{32}$ ,  $S_{64}$ ,  $S_{128}$ , and  $S_{256}$  supercells; nine AOs are used for every ion, so that the largest cell contains 4608 basis functions.

A full energy and gradient calculation with 512 ions in the supercell takes less than 2 hours on a low-level PC. Comparing the cost of calculations with supercells of increasing size, we see in the left part of the curve (Figure A2.2) that the elapsed time scales linearly with the number of the basis functions in the supercell up to  $S_{64}$ . The departure from linearity at that point corresponds to some nonlinear processes becoming important, in particular the diagonalization of large Fock (or KS) matrices.

However, the MgO supercells are a favorable case for their high point symmetry, of which CRYSTAL takes full advantage. For example, the

computer time spent for a calculation with the  $S_8$  supercell increases by about 20 times when symmetry is neglected, and full geometry optimization of large supercells ( $S_{64}$  or  $S_{128}$ ) containing a trapped-hole center, where only 8 of 48 original point symmetry operations exist, takes a few days on a single PC.

The cost of a calculation can increase also because of basis set enlargement. In particular, the truncation of the infinite series in CRYSTAL is effective for relatively sharp Gaussian functions. In ionic and covalent compounds, the exponents of the most diffuse Gaussians in a double-zeta-type basis set are relatively high, typically between 0.13 and 0.3 bohr<sup>-2</sup>, and the corresponding computational effort is not too large. On the contrary, the exponent of the outermost Gaussian in a metal can be smaller than 0.1 bohr<sup>-2</sup>, with a dramatic increase in computer time, as a huge amount of one- and two-electron integrals have to be computed, which is also the case for any large molecular basis sets containing diffuse functions with a high angular quantum number used for describing crystalline systems.<sup>89</sup> As an example, we consider the crystal of urea, with two molecules in the unit cell (16 atoms, 8-point symmetry operations), and we compare the elapsed time for the integral evaluation, SCF iteration, and gradient calculation as a function of the basis set. Six different molecular basis sets have been considered, from 3-21G to 6-311G(d,p). A few details on the basis sets (number of AOs and primitive functions per cell, exponents of the outermost atomic Gaussian) are given in Table A2.1, and results are shown on the right side of Table A2.2. The cost of a typical HF optimization of atomic positions (lattice parameters were fixed at their experimental values)<sup>291</sup> is also reported.

The 3-21G and 6-21G rows in Table A2.2 are similarly cheap, because the additional core Gaussian function exponents in 6-21G are too large to affect the total cost of the calculation. The higher cost of the 6-31G basis set is caused by the increased number of primitives in the valence shell, and to the smaller exponent of the outermost uncontracted valence Gaussian functions, as shown in Table A2.1. For similar reasons, 6-311G is more expensive than 6-31G. Although d-type functions are usually not diffuse, computing integrals involving d-type AOs is more demanding than with p-type AOs

**Table A2.1** Description of the Basis Sets Used for Crystalline Urea

Basis set	$N_{AO}^a$	$N_{\gamma}^b$	$\alpha_{sp}(H)^c$	$\alpha_{sp}(C)^c$	$\alpha_{sp}(N)^c$	$\alpha_{sp}(O)^c$
3-21G	88	144	0.183	0.196	0.283	0.374
6-21G	88	166	0.183	0.196	0.283	0.374
6-31G	88	208	0.161	0.169	0.212	0.270
6-311G	128	248	0.103	0.146	0.201	0.256
6-31G(d,p)	152	280	0.161 (1.100)	0.169 (0.800)	0.212 (0.800)	0.270 (0.800)
6-311G(d,p)	192	320	0.103 (0.750)	0.146 (0.626)	0.201 (0.913)	0.256 (1.292)

<sup>a</sup>Number of basis functions.

<sup>b</sup>Number of primitive Gaussians.

<sup>c</sup>Exponent of the outermost sp Gaussian (exponent of polarization functions in parentheses).

**Table A2.2** Elapsed Time (in sec) for Hartree–Fock Calculations of Urea Molecular Dimer and Crystal

Basis set	Molecular			Crystal				
	$t_{\text{INT}}^a$	$t_{\text{SCF}}^b$	$t_{\text{GRAD}}^c$	$t_{\text{INT}}^a$	$t_{\text{SCF}}^b$	$t_{\text{GRAD}}^c$	$t_{\text{OPT}}(N_{\text{step}})^d$	$t_{\text{OPT}}/N_{\text{step}}^e$
3-21G	3	5	11	13	13	43	853 (13)	65
6-21G	4	5	12	14	13	45	833 (12)	69
6-31G	5	7	24	30	20	124	2765 (15)	184
6-311G	15	21	58	135	66	466	8398 (13)	645
6-31G(d,p)	24	32	103	72	67	318	5673 (13)	457
6-311G(d,p)	53	79	199	263	251	973	18267 (13)	1405

With an AMD Athlon 2.8-GHz single-processor computer, 1-GB RAM, 512-KB cache, EIDE disks. Compilation of the program source with release 4.2 of the Portland Group PGF90 compiler, with -O1 -tp athlon options.

<sup>a</sup>Elapsed time (sec) for the calculation of integrals.

<sup>b</sup>Elapsed time (sec) for the full self-consistent field cycle.

<sup>c</sup>Elapsed time (sec) for the calculation of energy gradients.

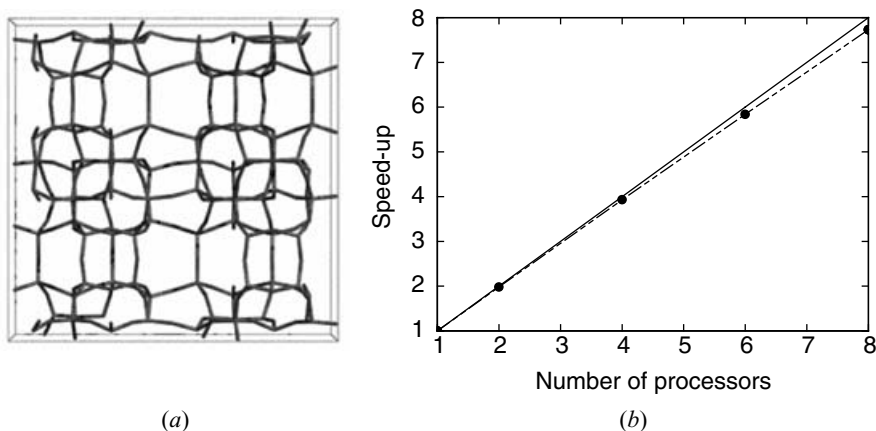
<sup>d</sup>Elapsed time (sec) for geometry optimization (number of optimization steps in parentheses).

<sup>e</sup>Average elapsed time (sec) per geometry optimization step.

with the same exponent. In summary, full geometry optimization with the largest basis set is 20 times slower than with the smallest basis set considered. For comparison, the computer time necessary to run similar calculations for a dimer of urea molecules (0-D), with the same geometry as in the crystal unit cell, is also reported on the left side of Table A2.2.

Beside the composition of the basis set, the choice of the Hamiltonian is also important in determining the length of a quantum-mechanical calculation. For example, the time spent in a single point energy plus gradient calculation, i.e., 667 seconds at the Hartree–Fock level (see Table A2.2), becomes 548, 898, and 1535 seconds, respectively, when we use LDA, GGA, and B3LYP, with a pruned grid for numerical integration of the exchange-correlation functional defined by 75 radial points and 494 angular points (the total number of points in the cell is 27,364; such a grid is labeled as LGRID in the CRYSTAL manual).<sup>56</sup>

Because parallel machines are becoming common in routine calculations, we also provide a few examples of parallel calculations. Two versions of parallel CRYSTAL exist: The first is based on a replicated data scheme (PCRYSTAL), whereas the second implements a distributed data algorithm (MPP-CRYSTAL). The most important difference between these versions concerns the Fock (or KS) matrix diagonalization step. Every matrix associated with one  $\mathbf{k}$ -point is diagonalized by a single processor in PCRYSTAL, and maximum efficiency is obtained when the number of  $\mathbf{k}$ -points in the reciprocal unit cell ( $n_{\mathbf{k}}$ ) is an exact multiple of the number of processors. When  $n_{\mathbf{k}}$  is small and the matrices to be diagonalized are large, as for example with large unit cell systems, parallelization with PCRYSTAL may become inefficient with



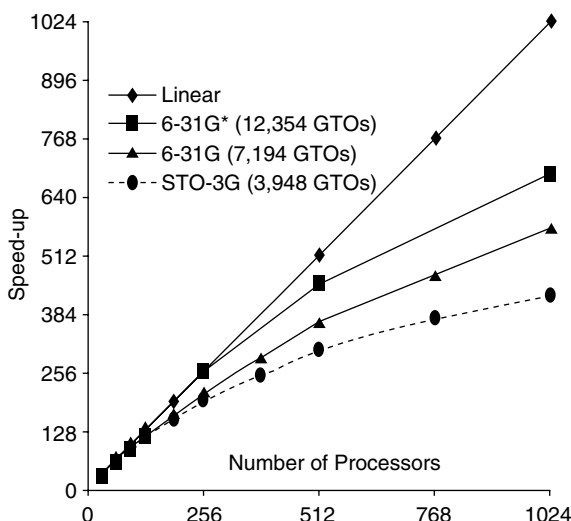
**Figure A2.3** (a) Structure of an all-silica zeolite silicalite, MFI framework; (b) Speed-up of the total energy and wavefunction calculation with the number of processors in a cluster of AMD Athlon 1.9-GHz single-processor, 1-GB RAM, 512-KB cache, EIDE disks. The code was compiled with PGF90 v4.2 and -O1 -tp athlon options.

any number of processors greater than  $n_k$ . MPP-CRYSTAL obviates this problem in massive-parallel computers by distributing every diagonalization task across a large number of processors.

We consider the scaling of PCRYSTAL with the number of active processors for silicalite (MFI framework; see Figure A2.3a), an all silica zeolite with 288 atoms in the unit cell (4416 AOs, 8-point symmetry operations; 8 k-points). Figure A2.3b shows that an almost linear correspondence exists between the number of active processors and execution time. The CPU time decreases from about 23 to 3 hours when distributing the job to eight processors.

MPP-CRYSTAL has been used recently for the calculation of the HF total energy of a small structural protein that has been characterized by X-ray diffraction studies (0.52 Å) to a very high precision: crambin,<sup>292</sup> which has  $P2_1$  symmetry with two chains per unit cell, 46 aminoacidic residues per chain, and 1284 atoms per cell.

The scaling of such a calculation with the number of active processors of IBM p-series 690 RS6000/P4 1.3 GHz (1240 processors) has been tested on three basis sets of increasing size: STO-3G (3948 AOs), 6-31G (7194 AOs), and 6-31G(d,p) (12354 AOs) at the HPCx Supercomputing Center in Daresbury, (U.K.). Figure A2.4 shows that scalability increases with increasing the basis set size. A total of 1024 active processors speed-up the calculation with the largest basis set by a factor of 700, with the calculation taking about 3 hours to be completed. Almost linear scaling is observed up to 256 processors.



**Figure A2.4** Comparison of the CPU time speed-up with the number of processors for three different basis sets of benchmark calculations on crambin (by courtesy of I. J. Bush).

## APPENDIX 3: ACRONYMS

AE	All-electron
AFM	Antiferromagnetic
AO	Atomic orbital
APW	Augmented plane waves
B3LYP	Becke 3-parameter exchange-correlation functional
BLYP	Becke and Lee–Yang–Parr exchange-correlation functional
BSSE	Basis set superposition error
BZ	Brillouin zone
CCA	Coupled-cluster approximation
CCSD	CC truncated to singles and doubles substitutions
CCSD(T)	CC truncated to singles, doubles, and (approximated) triples substitutions
CHA	Chabazite
CI	Configuration-interaction
CO	Crystalline orbital
CP	Counterpoise method
DFT	Density functional theory
EDI	Edingtonite
ENDOR	Electron-nuclear double-resonance

EPR	Electron paramagnetic resonance
FAU	Faujasite
FLAPW	Fully linearized augmented plane waves
FM	Ferromagnetic
GGA	Generalized gradient approximation
HF	Hartree–Fock
HOMO	Highest occupied molecular orbital
IR	Irreducible representation
KKR	Korringa–Kohn–Rostoker
KS	Kohn–Sham
LAPW	Linearized augmented-Plane-Waves
LCAO	Linear combination of atomic orbitals
LDA	Local density approximation
LSDA	Local spin-density approximation
LUMO	Lowest unoccupied molecular orbital
MP2	Møller–Plesset second-order perturbation expansion
ONIOM	Our N-layer integrated molecular orbitals—molecular mechanics
OPW	Orthogonalized plane waves
PP	Pseudopotentials
PW	Plane waves
PW91	Perdew–Wang 91 exchange-correlation functional
SOD	Sodalite
SPZ	LDA functional formulation (Slater exchange, Perdew–Zunger correlation)
STO	Slater type orbital
SVWN	LDA functional formulation (Slater exchange, Vosko–Wilk–Nusair correlation)
UHF	Unrestricted Hatree–Fock

---

## REFERENCES

1. R. Colle and O. Salvetti, *Theor. Chim. Acta*, **37**, 329 (1975). Approximate Calculation of Correlation Energy for Closed Shells.
2. E. Wigner, *Phys. Rev.*, **46**, 1002 (1934). On the Interaction of Electrons in Metals.
3. P. Hohenberg and W. Kohn, *Phys. Rev.*, **136**, B864 (1964). Inhomogeneous Electron Gas.
4. W. Kohn and L. J. Sham, *Phys. Rev.*, **140**, A1133 (1965). Self-Consistent Equations Including Exchange and Correlation Effects.
5. L. Hedin and S. Lundqvist, in *Solid State Physics*, Vol. 23, F. Seitz, D. Turnbull, and H. Ehrenreich, Eds., Academic Press, New York, 1969, p. 1. Effects of Electron–Electron and Electron–Phonon Interactions on the One-Electron States of Solids.
6. S. H. Vosko, L. Wilk, and M. Nusair, *Can. J. Phys.*, **58**, 1200 (1980). Accurate Spin-dependent Electron Liquid Correlation Energies for Local Spin Density Calculations: A Critical Analysis.
7. J. P. Perdew and A. Zunger, *Phys. Rev.*, **B23**, 5048 (1981). Self-Interaction Correction to Density Functional Approximations for Many-Electron Systems.

8. V. R. Saunders, in *Computational Techniques in Quantum Chemistry and Molecular Physics*, Vol. 15, G. H. F. Diercksen, B. T. Sutcliffe, and A. Veillard, Eds., Reidel, Dordrecht, the Netherlands, NATO ASI Series C, 1975, p. 347. An Introduction to Molecular Integral Evaluation.
9. V. R. Saunders, in *Methods in Computational Molecular Physics*, Vol. 113, G. H. F. Diercksen and S. Wilson, Eds., Reidel, Dordrecht, the Netherlands, NATO ASI Series C, 1983, p. 1. Molecular Integrals for Gaussian Type Functions.
10. C. F. Bender, *Phys. Rev.*, **183**, 23 (1969). Studies in Configuration Interaction: The First-Row Diatomic Hydrides.
11. W. Meyer, *Int. J. Quantum Chem. Symp.*, **5**, 341 (1971). Ionization Energies of Water from PNO-CI Calculations.
12. P. J. Hay and I. Shavitt, *J. Chem. Phys.*, **60**, 2865 (1974). Ab Initio Configuration Interaction Studies of the  $\pi$ -electron States of Benzene.
13. R. Ahlrichs, H. Lischka, V. Stämmler, and W. Kutzelnigg, *J. Chem. Phys.*, **62**, 1225 (1975). PNO-CI (Pair Natural Orbital Configuration Interaction) and CEPA-PNO (Coupled Electron Pair Approximation with Pair Natural Orbitals) Calculations of Molecular Systems. I. Outline of the Method for Closed-Shell States.
14. V. R. Saunders and J. H. Van Lenthe, *Mol. Phys.*, **48**, 923 (1983). The Direct CI. A Detailed Analysis.
15. P. Pulay, *Mol. Phys.*, **17**, 197 (1969). Ab Initio Calculation of Force Constants and Equilibrium Geometries in Polyatomic Molecules.
16. R. Krishnan, H. B. Schlegel, and J. A. Pople, *J. Chem. Phys.*, **72**, 4654 (1980). Derivative Studies in Configuration-Interaction Theory.
17. E. Clementi and A. Veillard, IBMOL, Computation of Wavefunction for Molecules of General Geometry Using the LGCO-MO-SCF Method. QCPE program number 92, 1966, Quantum Chemistry Program Exchange, Indiana University, Bloomington, Indiana.
18. W. J. Hehre, W. A. Lathan, R. Ditchfield, M. D. Newton, and J. A. Pople, GAUSSIAN 70, QCPE program number 237, 1970, Quantum Chemistry Program Exchange, Indiana University, Bloomington, Indiana.
19. F. Bassani, G. Pastori Parravicini, and R. A. Ballinger, *Electronic States and Optical Transitions in Solids*, Pergamon Press, Oxford, 1975.
20. V. L. Moruzzi, J. F. Janak, and A. R. Williams, *Calculated Electronic Properties of Metals*, Pergamon Press, New York, 1978.
21. D. R. Hamann, *Phys. Rev. Lett.*, **42**, 662 (1979). Semiconductor Charge Densities with Hard-Core and Soft-Core Pseudopotentials.
22. D. R. Hamann, M. Schlüter, and C. Chang, *Phys. Rev. Lett.*, **43**, 1494 (1979). Norm-Conserving Pseudopotentials.
23. A. Zunger, *Phys. Rev.*, **B21**, 4785 (1980). Ground-State Properties of Crystalline Silicon in a Density-Functional Pseudopotential Approach.
24. G. B. Bachelet, D. R. Hamann, and M. Schlüter, *Phys. Rev.*, **B26**, 4199 (1982). Pseudopotentials that work: From H to Pu.
25. J. C. Slater, in *Computational Methods in Band Theory*, P. Marcus, J. F. Janak, and A. R. Williams, Eds., Plenum Press, New York, 1971, p. 447. The Self-Consistent Field Method for Crystals.
26. J. Korringa, *Physica*, **13**, 392 (1947). On the Calculation of the Energy of a Bloch Wave in a Metal.
27. W. Kohn and N. Rostoker, *Phys. Rev.*, **94**, 1111 (1954). Solution of the Schrödinger Equation in Periodic Lattices with an Application to Metallic Lithium.
28. C. Herring, *Phys. Rev.*, **57**, 1169 (1940). A New Method for Calculating Wave Functions in Crystals.
29. J. C. Slater, *Phys. Rev.*, **92**, 603 (1953). An Augmented Plane Wave Method for Periodic Potential Problem.

30. M. M. Saffren and J. C. Slater, *Phys. Rev.*, **92**, 1126 (1953). An Augmented Plane-Wave Method for the Periodic Potential Problem. II.
31. O. K. Andersen, *Phys. Rev.*, **B12**, 3060 (1975). Linear Methods in Band Theory.
32. J. R. Leite, B. I. Bennett, and F. Herman, *Phys. Rev.*, **B12**, 1466 (1975). Electronic Structure of the Diamond Crystal Based on an Improved Cellular Calculation.
33. A. Zunger and A. J. Freeman, *Int. J. Quantum Chem. Symp.*, **10**, 383 (1976). Combined Fourier Transform and Discrete Variational Method Approach to the Self-Consistent Solution of the Electronic Band Structure Problem within the Local Density Formalism.
34. C. R. A. Catlow and W. C. Mackrodt, *Computer Simulation of Solids*, Springer-Verlag, Berlin, Germany, 1982.
35. J. M. André, L. Gouverneur, and G. Leroy, *Int. J. Quantum Chem.*, **1**, 427 (1967). L'Etude Théorique des Systèmes Périodiques. I. La Méthode LCAO-HCO.
36. J. M. André, L. Gouverneur, and G. Leroy, *Int. J. Quantum Chem.*, **1**, 451 (1967). L'Etude Théorique des Systèmes Périodiques. II. La Méthode LCAO-SCF-CO.
37. F. E. Harris and H. J. Monkhorst, *Phys. Rev. Lett.*, **23**, 1026 (1969). Complete Calculations of the Electronic Energies of Solids.
38. G. Del Re, J. Ladik, and G. Biczó, *Phys. Rev.*, **155**, 997 (1967). Self-Consistent-Field Tight-Binding Treatment of Polymers. I. Infinite Three-Dimensional Case.
39. H. Stoll and H. Preuss, *Phys. Status Solidi*, **60**, 185 (1973). Convergence of Lattice Sums in Hartree-Fock LCAO Calculations.
40. H. Stoll and H. Preuss, *Int. J. Quantum Chem.*, **IX**, 775 (1975). Hartree-Fock Calculation of Cohesive Energies and Equilibrium Lattice Constants for Solid Li and Be.
41. F. E. Harris, in *Advances and Perspectives in Theoretical Chemistry*, Vol. 2, H. Eyring and D. Henderson, Eds., Academic Press, New York, 1975, p. 147. Hartree-Fock Studies of Electronic Structures of Crystalline Solids.
42. R. N. Ewema, *Phys. Rev.*, **B7**, 818 (1972). General Crystalline Hartree-Fock Formalism: Diamond Results.
43. R. N. Ewema, G. T. Surratt, D. L. Wilhite, and G. G. Wepfer, *Philos. Mag.*, **29**, 1033 (1974). Hartree-Fock Electron Distribution of Cubic BN.
44. C. Pisani and R. Dovesi, *Int. J. Quantum Chem.*, **17**, 501 (1980). Exact Exchange Hartree-Fock Calculations for Periodic Systems. I. Illustration of the Method.
45. R. Dovesi, C. Pisani, C. Roetti, M. Causà, and V. R. Saunders, CRYSTAL88, an Ab Initio All-Electron LCAO Hartree-Fock Program for Periodic Systems. QCPE program number 577, 1989, Quantum Chemistry Program Exchange, Indiana University, Bloomington, Indiana.
46. L. Bellaiche and D. Vanderbilt, *Phys. Rev.*, **B61**, 7877 (2000). Virtual Crystal Approximation Revisited: Application to Dielectric and Piezoelectric Properties of Perovskites.
47. N. J. Ramer and A. M. Rappe, *J. Phys. Chem. Solids*, **61**, 315 (2000). Application of a New Virtual Crystal Approach for the Study of Disordered Perovskites.
48. S. W. De Leeuw, J. W. Perram, and E. R. Smith, *Proc. R. Soc. Lond. A*, **373**, 27 (1980). Simulation of Electrostatic Systems in Periodic Boundary Conditions. I. Lattice Sums and Dielectric Constants.
49. S. W. De Leeuw, J. W. Perram, and E. R. Smith, *Proc. R. Soc. Lond. A*, **373**, 57 (1980). Simulation of Electrostatic Systems in Periodic Boundary Conditions. II. Equivalence of Boundary Conditions.
50. E. R. Smith, *Proc. R. Soc. Lond. A*, **375**, 475 (1981). Electrostatic Energy in Ionic Crystals.
51. E. R. Smith, *Proc. R. Soc. Lond. A*, **381**, 241 (1982). Effects of Surface Charge on the Electrostatic Energy of an Ionic Crystal.
52. P. P. Ewald, *Ann. Phys. (Leipzig)*, **64**, 253 (1921). Die Berechnung Optischer und Elektrostatischer Gitterpotentiale.



- 
53. F. E. Harris and H. J. Monkhorst, *Phys. Rev.*, **B2**, 4400 (1970). Electronic-Structure Studies of Solids. I. Fourier Representation Method for Madelung Sums.
  54. V. R. Saunders, C. Freyria-Fava, R. Dovesi, L. Salasco, and C. Roetti, *Mol. Phys.*, **77**, 629 (1992). On the Electrostatic Potential in Crystalline Systems where the Charge Density is Expanded in Gaussian Functions.
  55. V. R. Saunders, C. Freyria-Fava, R. Dovesi, and C. Roetti, *Comput. Phys. Commun.*, **84**, 156 (1994). On the Electrostatic Potential in Linear Periodic Polymers.
  56. V. R. Saunders, R. Dovesi, C. Roetti, R. Orlando, C. M. Zicovich-Wilson, N. M. Harrison, K. Doll, B. Civalleri, I. J. Bush, P. D'Arco, and M. Llunell, *CRYSTAL03 User's Manual*, 2003, Università di Torino, Torino, Italy.
  57. A. R. Oganov, J. P. Brodholt, and D. Price, in *Energy Modelling in Minerals*, Vol. 4, C. M. Gramaccioli, Ed., Eötvös University Press, Budapest, EMU Notes in Mineralogy, 2002, p. 83. Ab Initio Theory of Phase Transitions and Thermoelasticity of Minerals.
  58. R. Car and M. Parrinello, *Phys. Rev. Lett.*, **55**, 2471 (1985). Unified Approach for Molecular Dynamics and Density-Functional Theory.
  59. F. Mauri, B. G. Pfrommer, and S. G. Louie, *Phys. Rev. Lett.*, **77**, 5300 (1996). Ab Initio Theory of NMR Chemical Shifts in Solids and Liquids.
  60. P. Umari, A. Pasquarello, and A. Dal Corso, *Phys. Rev.*, **B63**, 094305 (2001). Raman Scattering Intensities in Alpha-Quartz: A First-Principles Investigation.
  61. M. Veithen, X. Gonze, and P. Ghosez, *Phys. Rev. Lett.*, **93**, 187401 (2004). First-Principles Study of the Electrooptic Effect in Ferroelectric Oxides.
  62. M. Schütz, G. Hetzer, and H. J. Werner, *J. Chem. Phys.*, **111**, 5691 (1999). Low-Order Scaling Local Electron Correlation Methods. I. Linear Scaling Local MP2.
  63. M. Schütz, *Phys. Chem. Chem. Phys.*, **4**, 3941 (2002). A New, Fast, Semi-Direct Implementation of Linear Scaling Local Coupled Cluster Theory.
  64. N. Marzari and D. Vanderbilt, *Phys. Rev.*, **B56**, 12847 (1997). Maximally Localized Generalized Wannier Functions for Composite Energy Bands.
  65. C. M. Zicovich-Wilson, R. Dovesi, and V. R. Saunders, *J. Chem. Phys.*, **115**, 9708 (2001). A General Method to obtain Well Localized Wannier Functions for Composite Energy Bands in Linear Combination of Atomic Orbital Periodic Calculations.
  66. C. Pisani, M. Busso, G. Capecchi, S. Casassa, R. Dovesi, L. Maschio, V. R. Saunders, C. M. Zicovich-Wilson, and M. Schütz, *in preparation*, Local-MP2 Electron Correlation Method for Non Conducting Crystals.
  67. A. Lüchow and J. B. Anderson, *Annu. Rev. Phys. Chem.*, **51**, 501 (2000). Monte Carlo Methods in Electronic Structures for Large Systems.
  68. S. Suhai, *Phys. Rev.*, **B27**, 3506 (1983). Quasiparticle Energy-Band Structures in Semiconducting Polymers: Correlation Effects on the Band Gap in Polyacetylene.
  69. J. Q. Sun and R. J. Bartlett, *J. Chem. Phys.*, **104**, 8553 (1996). Second-Order Many-Body Perturbation-Theory Calculations in Extended Systems.
  70. Q. Y. Ayala, K. N. Kudin, and G. E. Scuseria, *J. Chem. Phys.*, **115**, 9698 (2001). Atomic Orbital Laplace-Transformed Second-Order Møller-Plesset Theory for Periodic Systems.
  71. H. Stoll, *Phys. Rev.*, **B46**, 6700 (1992). Correlation Energy of Diamond.
  72. H. Stoll, *Chem. Phys. Lett.*, **191**, 548 (1992). The Correlation Energy of Crystalline Silicon.
  73. K. Doll, M. Dolg, P. Fulde, and H. Stoll, *Phys. Rev.*, **B52**, 4842 (1995). Correlation Effects in Ionic Crystals: The Cohesive Energy of MgO.
  74. K. Doll and H. Stoll, *Phys. Rev.*, **B56**, 10121 (1997). Cohesive Properties of Alkali Halides.
  75. A. Shukla, M. Dolg, P. Fulde, and H. Stoll, *Phys. Rev.*, **B60**, 5211 (1999). Wave-Function-Based Correlated Ab Initio Calculations on Crystalline Solids.
  76. E. Runge and E. K. U. Gross, *Phys. Rev. Lett.*, **52**, 997 (1984). Density-Functional Theory for Time-Dependent Systems.

77. M. E. Casida, in *Recent Developments and Applications of Modern Density Functional Theory*, Vol. 4, J. Seminario, Ed., Elsevier, Amsterdam, Theoretical and Computational Chemistry, 1996, p. 391. Time-dependent Density Functional Response Theory of Molecular Systems: Theory, Computational Methods, and Functionals.
78. M. Petersilka and E. K. U. Gross, *Int. J. Quantum Chem.*, **60**, 181 (1996). Spin-Multiplet Energies from Time-Dependent Density Functional Theory.
79. L. Bernasconi, M. Sprik, and J. Hutter, *Chem. Phys. Lett.*, **394**, 141 (2004). Hartree-Fock Exchange in Time Dependent Density Functional Theory: Application to Charge Transfer Excitations in Solvated Molecular Systems.
80. M. S. Hybertsen and S. G. Louie, *Phys. Rev. Lett.*, **55**, 1418 (1985). First-Principle Theory of Quasiparticles: Calculations of Band Gaps in Semiconductors and Insulators.
81. M. Rohlfing, P. Krüger, and J. Pollmann, *Phys. Rev.*, **B52**, 1905 (1995). Efficient Scheme for GW Quasiparticle Band-Structure Calculations.
82. S. Albrecht, L. Reining, R. Del Sole, and G. Onida, *Phys. Rev. Lett.*, **80**, 4510 (1998). Ab Initio Calculation of Excitonic Effects in the Optical Spectra of Semiconductors.
83. B. Arnaud and M. Alouani, *Phys. Rev.*, **B62**, 4464 (2000). All-Electron Projector-Augmented-Wave GW Approximation: Application to the Electronic Properties of Semiconductors.
84. S. Lebègue, B. Arnaud, M. Alouani, and P. E. Blöchl, *Phys. Rev.*, **B67**, 155208 (2003). Implementation of an All-Electron GW Approximation Based on the Projector Augmented Wave Method without Plasmon Pole Approximation to Si, SiC, AlAs, InAs, NaH and KH.
85. A. W. Hewat and C. Riekel, *Acta Crystallogr. Sec. A*, **35**, 569 (1979). The Crystal Structure of Deuteroammonia between 2 and 180 K by Neutron Powder Profile Refinement.
86. M. F. C. Ladd, *Symmetry in Molecules and Crystals*, Ellis Horwood, New York, 1989.
87. T. Hahn, Ed., *International Tables of Crystallography*, Vol. A, 5th ed., Kluwer Academic Publishers, Dordrecht, the Netherlands, 2002.
88. F. Bloch, *Z. Phys.*, **52**, 555 (1928). Über die Quantenmechanik der Electonen in Kristallgittern.
89. C. Pisani, R. Dovesi, and C. Roetti, *Hartree-Fock Ab Initio Treatment of Crystalline Systems*, Springer, Berlin, 1988.
90. P. A. M. Dirac, *Proc. Cambridge Philos. Soc.*, **26**, 376 (1930). Note on Exchange Phenomena in the Thomas Atom.
91. A. D. Becke, *J. Chem. Phys.*, **98**, 5648 (1993). Density-Functional Thermochemistry. III. The Role of Exact Exchange.
92. H. Landolt and R. Börnstein, *Numerical Data and Functional Relationships in Science and Technology*, Springer, Berlin, 1982.
93. T. Koopmans, *Physica*, **1**, 104 (1934). Über die Zuordnung von Wellen Funktionen und Eigenwerten zu den Einzelnen Elektronen eines Atom.
94. A. Savin, C. J. Umrigar, and X. Gonze, *Chem. Phys. Lett.*, **288**, 391 (1998). Relationship of Kohn-Sham Eigenvalues to Excitation Energies.
95. A. I. Al-Sharif, R. Resta, and C. J. Umrigar, *Phys. Rev.*, **A57**, 2466 (1998). Evidence of Physical Reality in the Kohn-Sham Potential: The Case of Atomic Ne.
96. A. Görling and M. Levy, *Phys. Rev.*, **A50**, 196 (1994). Exact Kohn-Sham Scheme Based on Perturbation-Theory.
97. A. Görling and M. Levy, *Int. J. Quantum Chem. Symp.*, **29**, 93 (1995). DFT Ionization Formulas and a DFT Perturbation-Theory for Exchange and Correlation, through Adiabatic Connection.
98. A. D. Becke, *Phys. Rev.*, **A38**, 3098 (1988). Density-Functional Exchange-Energy Approximation with Correct Asymptotic Behavior.
99. C. Lee, W. Yang, and R. G. Parr, *Phys. Rev.*, **B37**, 785 (1988). Development of the Colle-Salvetti Correlation-Energy Formula into a Functional of the Electron Density.

100. J. D. Pack and H. J. Monkhorst, *Phys. Rev.*, **B16**, 1748 (1977). Special Points for Brillouin-zone Integration - A Reply.
101. G. Gilat and P. Raubenheimer, *Phys. Rev.*, **144**, 390 (1966). Accurate Numerical Method for Calculating Frequency-Distribution Functions in Solids.
102. G. Gilat, *J. Comput. Phys.*, **10**, 432 (1972). Analysis of Methods for Calculating Spectral Properties in Solids.
103. M. Methfessel and A. T. Paxton, *Phys. Rev.*, **B40**, 3616 (1989). High-precision Sampling for Brillouin-zone Integration in Metals.
104. R. Dovesi, *Int. J. Quantum Chem.*, **29**, 1755 (1986). On the Role of Symmetry in the Ab Initio Hartree-Fock Linear-Combination-of-Atomic-Orbitals Treatment of Periodic Systems.
105. N. W. Ashcroft and N. D. Mermin, *Solid State Physics*, Saunders College Publishing, Fort Worth, Texas, 1976.
106. M. Lax, *Symmetry Principles in Solid State and Molecular Physics*, Wiley, New York, 1974.
107. C. M. Zicovich-Wilson and R. Dovesi, *Int. J. Quantum Chem.*, **67**, 299 (1998). On the Use of Symmetry Adapted Crystalline Orbitals in SCF-LCAO Periodic Calculations. I. The Construction of the Symmetrized Orbitals.
108. C. M. Zicovich-Wilson and R. Dovesi, *Int. J. Quantum Chem.*, **67**, 311 (1998). On the Use of Symmetry-Adapted Crystalline Orbitals in SCF-LCAO Periodic Calculations. II. Implementation of the Self-Consistent-Field Scheme and Examples.
109. L. A. Curtiss, K. Raghavachari, G. W. Trucks, and J. A. Pople, *J. Chem. Phys.*, **94**, 7221 (1991). Gaussian-2 Theory for Molecular-Energies of 1st-Row and 2nd-Row Compounds.
110. M. Prencipe, A. Zupan, R. Dovesi, E. Aprà, and V. R. Saunders, *Phys. Rev.*, **B51**, 3391 (1995). Ab Initio Study of the Structural Properties of LiF, NaF, KF, LiCl, NaCl, and KCl.
111. D. R. Lide, Ed., *The Handbook of Chemistry and Physics*, 72nd edition, CRC Press, Boca Raton, Florida, 1991.
112. T. L. Hill, *An Introduction to Statistical Thermodynamics*, Addison-Wesley, Reading, Massachusetts, 1960.
113. A. Masunov and J. J. Dannenberg, *J. Phys. Chem. A*, **103**, 178 (1999). Theoretical Study of Urea. I. Monomers and Dimers.
114. E. R. Davidson and D. Feller, *Chem. Rev.*, **86**, 681 (1986). Basis Set Selection for Molecular Calculations.
115. N. R. Kestner and J. E. Combariza, in *Reviews in Computational Chemistry*, Vol. 13, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 1999, p. 99. Basis Set Superposition Errors: Theory and Practice.
116. S. F. Boys and F. Bernardi, *Mol. Phys.*, **19**, 553 (1970). The Calculation of Small Molecular Interaction by the Difference of Separate Total Energies. Some Procedures with Reduced Errors.
117. M. A. Spackman and A. S. Mitchell, *Phys. Chem. Chem. Phys.*, **3**, 1518 (2001). Basis Set Choice and Basis Set Superposition Error (BSSE) in Periodic Hartree-Fock Calculations on Molecular Crystals.
118. A. J. Thakkar, T. Koga, M. Saito, and R. E. Hoffmeyer, *Int. J. Quantum Chem. Symp.*, **27**, 343 (1993). Double and Quadruple Zeta Contracted Gaussian Basis Sets for Hydrogen through Neon.
119. R. Dovesi, M. Causà, R. Orlando, C. Roetti, and V. R. Saunders, *J. Chem. Phys.*, **92**, 7402 (1990). Ab Initio Approach to Molecular Crystals: A Periodic Hartree-Fock Study of Crystalline Urea.
120. B. Civalieri, E. Garrone, and P. Ugliengo, *J. Mol. Struct. (THEOCHEM)*, **419**, 227 (1997). Density Functional Study of Hydrogen-Bonded Systems: Energetic and Vibrational Features of Some Gas-Phase Adducts of Hydrogen Fluoride.
121. J. F. Dobson, K. McLennan, A. Rubio, J. Wang, T. Gould, H. M. Le, and B. P. Dinte, *Aust. J. Chem.*, **54**, 513 (2001). Prediction of Dispersion Forces: Is There a Problem?

122. L. Ojamäe, K. Hermansson, C. Pisani, M. Causà, and C. Roetti, *J. Chem. Phys.*, **100**, 2128 (1994). Mechanical and Molecular Properties of Ice VIII from Crystal-Orbital Ab Initio Hartree-Fock Calculations.
123. B. Silvi, A. Beltrán, and J. Andrés, *J. Mol. Struct.*, **436-437**, 443 (1997). Periodic Hartree-Fock Calculation of the  $A_{1g}$  ( $T_z$ ) and  $E_g$  ( $T_x$ ,  $T_y$ ) Phonon Modes in Ice VIII.
124. C. Pisani, S. Casassa, and P. Ugliengo, *Chem. Phys. Lett.*, **253**, 201 (1996). Proton-Ordered Ice Structures at Zero Pressure: A Quantum-Mechanical Investigation.
125. P. Zapol, L. A. Curtiss, and A. Erdemir, *J. Chem. Phys.*, **113**, 3338 (2000). Periodic Ab Initio Calculations of Orthoboric Acid.
126. M. Milanesio, R. Bianchi, P. Ugliengo, C. Roetti, and D. Viterbo, *J. Mol. Struct. (THEO-CHEM)*, **419**, 139 (1997). Vitamin C at 120 K: Experimental and Theoretical Study of the Charge Density.
127. S. Camus, K. D. M. Harris, and R. L. Johnson, *Chem. Phys. Lett.*, **276**, 186 (1997). Ab Initio Calculation of H-2 Quadrupole Coupling Constants in Molecular Crystals: Application to Polymorphs of Oxalic Acid Dihydrate.
128. G. L. Cárdenas-Jirón, A. Masunov, and J. J. Dannenberg, *J. Phys. Chem. A*, **103**, 7042 (1999). Molecular Orbital Study of Crystalline *p*-Benzoquinone.
129. I. Petrovic, A. Navrotsky, M. E. Davis, and S. I. Zones, *Chem. Mater.*, **5**, 1805 (1993). Thermochemical Study of the Stability of Frameworks in High-Silica Zeolites.
130. I. Petrovic, P. J. Heaney, and A. Navrotsky, *Phys. Chem. Miner.*, **23**, 119 (1996). Thermochemistry of the New Silica Polymorph Moganite.
131. S. K. Saxena, N. Chatterjee, Y. Fei, and G. Shen, *Thermodynamic Data on Oxides and Silicates*, Springer, Berlin, 1993.
132. P. M. Piccione, C. Laberty, S. Y. Yang, M. A. Camblor, A. Navrotsky, and M. E. Davis, *J. Phys. Chem. B*, **104**, 10001 (2000). Thermochemistry of Pure-Silica Zeolites.
133. B. Civalleri, C. M. Zicovich-Wilson, P. Ugliengo, V. R. Saunders, and R. Dovesi, *Chem. Phys. Lett.*, **292**, 394 (1998). A Periodic Ab-Initio Study of the Structure and Relative Stability of Silica Polymorphs.
134. M. Catti, B. Civalleri, and P. Ugliengo, *J. Phys. Chem.*, **B104**, 7259 (2000). Structure and Energetics of SiO<sub>2</sub> Polymorphs by Quantum-Mechanical and Semiclassical Approaches.
135. T. Demuth, Y. Jeanvoine, J. Hafner, and J. G. Ángyán, *J. Phys.-Condens. Matter*, **11**, 3833 (1999). Polymorphism in Silica Studied in the Local Density and Generalized-Gradient Approximations.
136. K. P. Schröder and J. Sauer, *J. Phys. Chem.*, **100**, 11043 (1996). Potential Functions for Silica and Zeolite Catalysts Based on Ab Initio Calculations. 3. A Shell Model Ion Pair Potential for Silica and Aluminosilicates.
137. M. Sierka and J. Sauer, *Faraday Discuss.*, **106**, 41 (1997). Structure and Reactivity of Silica and Zeolite Catalysts by Combined Quantum Mechanics Shell-Model Potential Approach Based on DFT.
138. R. Dovesi, F. Freyria-Fava, C. Roetti, and V. R. Saunders, *Faraday Discuss.*, **106**, 173 (1997). Structural, Electronic and Magnetic Properties of KMF<sub>3</sub> (M = Mn, Fe, Co, Ni).
139. N. M. Harrison, V. R. Saunders, R. Dovesi, and W. C. Mackrodt, *Philos. Trans. R. Soc. Lond., Ser. A*, **356**, 75 (1998). Transition Metal Materials: A First Principles Approach to the Electronic Structure of the Insulating Phase.
140. G. Mallia, R. Orlando, M. Llunell, and R. Dovesi, in *Computational Materials Science*, Vol. 187, C. R. A. Catlow and E. Kotomin, Eds., IOS Press, Amsterdam, NATO Science Series III, 2003, p. 102. On the Performance of Various Hamiltonians in the Study of Crystalline Compounds. The Case of Open Shell Systems.
141. M. D. Towler, N. L. Allan, N. M. Harrison, V. R. Saunders, W. C. Mackrodt, and E. Aprà, *Phys. Rev.*, **B50**, 5041 (1994). Ab Initio Study of MnO and NiO.
142. M. Catti, G. Valerio, and R. Dovesi, *Phys. Rev.*, **B51**, 7441 (1995). Theoretical Study of Electronic, Magnetic, and Structural Properties of  $\alpha$ -Fe<sub>2</sub>O<sub>3</sub> (Hematite).

- 
143. M. Catti, G. Sandrone, G. Valerio, and R. Dovesi, *J. Phys. Chem. Solids*, **57**, 1735 (1996). Electronic, Magnetic and Crystal Structure of  $\text{Cr}_2\text{O}_3$  by Theoretical Methods.
  144. A. Chartier, P. D'Arco, R. Dovesi, and V. R. Saunders, *Phys. Rev.*, **B60**, 14042 (1999). An Ab Initio Hartree-Fock Investigation of the Structural, Electronic and Magnetic Properties of  $\text{Mn}_3\text{O}_4$ -Hausmannite.
  145. G. Valerio, M. Catti, R. Dovesi, and R. Orlando, *Phys. Rev.*, **B52**, 2422 (1995). Ab Initio Study of Antiferromagnetic Rutile-Type  $\text{FeF}_2$ .
  146. I. d. P. R. Moreira, R. Dovesi, C. Roetti, V. R. Saunders, and R. Orlando, *Phys. Rev.*, **B72**, 7816 (2000). Ab Initio Study of  $\text{MF}_2$  ( $\text{M} = \text{Mn, Fe, Co, Ni}$ ) Rutile-Type Compounds Using the Periodic UHF Approach.
  147. P. Reinhardt, I. d. P. R. Moreira, R. Dovesi, C. de Graaf, and F. Illas, *Chem. Phys. Lett.*, **319**, 625 (2000). Detailed Ab-Initio Analysis of the Magnetic Coupling in  $\text{CuF}_2$ .
  148. I. d. P. R. Moreira and R. Dovesi, *Phys. Rev.*, **B67**, 134513 (2003). Ab Initio Periodic Approach to Electronic Structure and Magnetic Exchange in  $\text{A}_2\text{CuO}_2\text{X}_2$  ( $\text{A} = \text{Ca, Sr}$  and  $\text{X} = \text{F, Cl}$ ) High- $T_c$  Superconductor Parent Compounds.
  149. L. J. De Jong and R. Block, *Physica B*, **79**, 568 (1975). On the Exchange Interactions in Some 3D-Metal Ionic Compounds.
  150. O. Kahn, *Molecular Magnetism*, VCH, New York, 1993.
  151. K. Yosida, *Theory of Magnetism*, Springer, Heidelberg, Germany, 1996.
  152. B. Civalieri, A. M. Ferrari, M. Llunell, R. Orlando, M. Merawa, and P. Uglierio, *Chem. Mater.*, **15**, 3996 (2003). Cation Selectivity in Alkali-Exchanged Chabazite: an Ab-Initio Periodic Study.
  153. M. P. Habas, R. Dovesi, and A. Lichanot, *J. Phys.-Condens. Matter*, **10**, 6897 (1998). B1-B2 Phase Transition in Alkaline-earth Oxides: A Comparison of Ab Initio Hartree-Fock and Density Functional Calculations.
  154. P. Richet, H. K. Mao, and P. M. Bell, *J. Geophys. Res.*, **93**, 15279 (1986). Static Compression and Equation of State of  $\text{CaO}$  to 1.35 Mbar.
  155. R. Jeanloz, H. K. Ahrens, H. K. Mao, and P. M. Bell, *Science*, **206**, 829 (1979). B1-B2 Transition in Calcium Oxide from Shock-Wave and Diamond-Cell Experiments.
  156. J. F. Mammone, H. K. Mao, and P. M. Bell, *Geophys. Res. Lett.*, **8**, 140 (1981). Equation of State of  $\text{CaO}$  Under Static Pressure Conditions.
  157. J. Muscat, V. Swamy, and N. M. Harrison, *Phys. Rev.*, **B65**, 224112 (2002). First-Principles Calculations of the Phase Stability of  $\text{TiO}_2$ .
  158. L. S. Dubrovinsky, N. A. Dubrovinskaia, V. Swamy, J. Muscat, N. M. Harrison, R. Ahuja, B. Holm, and B. Johansson, *Nature*, **410**, 653 (2001). Cotunnite-Structured Titanium Dioxide: The Hardest Known Oxide.
  159. J. F. Nye, *Physical Properties of Crystals*, Oxford University Press, London, 1957.
  160. M. Born and J. R. Oppenheimer, *Ann. Phys. (Leipzig)*, **84**, 457 (1927). On the Quantum Theory of Molecules.
  161. P. Giannozzi, S. de Gironcoli, P. Pavone, and S. Baroni, *Phys. Rev.*, **B43**, 7231 (1991). Ab Initio Calculation of Phonon Dispersions in Semiconductors.
  162. S. Baroni, S. de Gironcoli, A. Dal Corso, and P. Giannozzi, *Rev. Mod. Phys.*, **73**, 515 (2001). Phonons and Related Crystal Properties from Density-Functional Perturbation Theory.
  163. X. Gonze, J.-M. Beuken, R. Caracas, F. Detraux, M. Fuchs, G.-M. Rignanese, L. Sindic, M. Verstraete, G. Zerah, F. Jollet, M. Torrent, A. Roy, M. Mikami, P. Ghosez, J.-Y. Raty, and D. C. Allan, *Computational Materials Science*, **25**, 478 (2002). First-Principles Computation of Material Properties: The ABINIT Software Project.
  164. S. Y. Savrasov, *Phys. Rev.*, **B54**, 16470 (1996). Linear-Response Theory and Lattice Dynamics: A Muffin-Tin-Orbital Approach.

165. F. Pascale, C. M. Zicovich-Wilson, F. Lopez Gejo, B. Civalleri, R. Orlando, and R. Dovesi, *J. Comput. Chem.*, **25**, 888 (2004). The Calculation of the Vibrational Frequencies of the Crystalline Compounds and Its Implementation in the CRYSTAL Code.
166. X. Gonze, D. C. Allan, and M. P. Teter, *Phys. Rev. Lett.*, **68**, 3603 (1992). Interatomic Force-Constants from First Principles - The Case of Alpha-Quartz.
167. F. Gervais and B. Piriou, *Phys. Rev.*, **B11**, 3944 (1975). Temperature Dependence of Transverse and Longitudinal Optic Modes in the *alpha* and *beta* Phases of Quartz.
168. R. A. Evarestov and A. V. Bandura, *Int. J. Quantum Chem.*, **96**, 282 (2004). Hartree-Fock Calculations of Electronic Structure of (110)-Surface of Rutile TiO<sub>2</sub>: Comparison of Single (2D) and Periodic (3D) Slab Models.
169. P. W. Tasker, *J. Phys. C*, **12**, 4977 (1979). The Stability of Ionic Crystal Surfaces.
170. S. Casassa, P. Ugliengo, and C. Pisani, *J. Chem. Phys.*, **106**, 8030 (1997). Proton-Ordered Models of Ordinary Ice for Quantum-Mechanical Studies.
171. K. Doll, N. M. Harrison, and V. R. Saunders, *J. Phys.-Condens. Matter*, **11**, 5007 (1999). A Density Functional Study of Lithium Bulk and Surfaces.
172. M. Causà, R. Dovesi, E. Kotomin, and C. Pisani, *J. Phys. C*, **20**, 4983 (1987). MgO (110) Surface and CO Adsorption Thereon. I. Clean (110) Surface.
173. L. G. M. Pettersson, M. Nyberg, J.-L. Pascual, and M. A. Nygren, in *Chemisorption and Reactivity on Supported Clusters and Thin Films*, Vol. 331, G. Pacchioni and R. M. Lambert, Eds., Kluwer Academic Press, Dordrecht, the Netherlands, 1997, p. 425. Theoretical Modeling of Chemisorption and Reactions at Metal-Oxide Surfaces.
174. P. W. Tasker, *Adv. Ceram.*, **10**, 176 (1984). Surface of Magnesia and Alumina.
175. G. Mallia, PhD thesis, Università di Torino, Torino, Italy, 2002.
176. W. C. Mackrodt, *Philos. Trans. R. Soc. Lond., Ser. A*, **341**, 301 (1992). Classical and Quantum Simulation of the Surface Properties of  $\alpha$ -Al<sub>2</sub>O<sub>3</sub>.
177. A. D'Ercole, A. M. Ferrari, and C. Pisani, *J. Chem. Phys.*, **115**, 509 (2001). On the Role of Electrostatics in the Heterolytic Splitting of Covalent Bonds at Defective Oxide Surfaces.
178. M. Causà, R. Dovesi, and F. Ricca, *Surf. Sci.*, **280**, 1 (1993). Regular Adsorption of CO Molecules on LiF (001).
179. K. Hermansson and M. Alfredsson, *Surf. Sci.*, **411**, 23 (1998). N<sub>2</sub> and HF Vibrations on LiF(001): the Effect of Surface Coverage.
180. D. P. Taylor, W. P. Hess, and M. I. McCarthy, *J. Phys. Chem. B*, **101**, 7455 (1997). Structure and Energetics of the Water/NaCl(100) Interface.
181. A. V. Puchina, V. E. Puchin, M. Huisinga, R. Bennewitz, and M. Reichling, *Surf. Sci.*, **404**, 687 (1998). Theoretical Modelling of Steps and Surface Oxidation on CaF<sub>2</sub>(111).
182. V. E. Puchin, A. V. Puchina, M. Huisinga, and M. Reichling, *J. Phys.-Condens. Matter*, **13**, 2081 (2001). Theoretical Modelling of Steps on the CaF<sub>2</sub>(111) Surface.
183. C. Pisani, M. Causà, R. Dovesi, and C. Roetti, *Prog. Surf. Sci.*, **25**, 119 (1987). Hartree-Fock Ab Initio Characterization of Ionic Crystal Surfaces with a Slab Model. The (0001) Face of  $\alpha$ -Al<sub>2</sub>O<sub>3</sub>.
184. M. Causà, R. Dovesi, C. Pisani, and C. Roetti, *Surf. Sci.*, **215**, 259 (1988). Ab Initio Characterization of the (0001) and (10-10) Crystal Faces of  $\alpha$ -Alumina.
185. V. E. Puchin, J. D. Gale, A. L. Shluger, E. A. Kotomin, J. Gunster, M. Brause, and V. Kempter, *Surf. Sci.*, **370**, 190 (1997). Atomic and Electronic Structure of the Corundum (0001) Surface: Comparison with Surface Spectroscopies.
186. A. Wander, B. Searle, and N. M. Harrison, *Surf. Sci.*, **458**, 25 (2000). An Ab Initio Study of  $\alpha$ -Al<sub>2</sub>O<sub>3</sub> (0001): the Effects of Exchange and Correlation Functionals.
187. S. Gennard, F. Corà, and C. R. A. Catlow, *J. Phys. Chem. B*, **103**, 10158 (1999). Comparison of the Bulk and Surface Properties of Ceria and Zirconia by Ab Initio Investigations.
188. C. Rehbein, F. Michel, N. M. Harrison, and A. Wander, *Surf. Rev. Lett.*, **5**, 337 (1998). Ab Initio Total Energy Studies of the Alpha-Cr<sub>2</sub>O<sub>3</sub> (0001) and (01-1)2 Surfaces.

189. T. Ouazzani, A. Lichanot, and C. Pisani, *J. Phys. Chem. Solids*, **56**, 915 (1995). Effect of the Quality of the Atomic Orbitals Basis Set about the Relaxation and Electronic Structure of the (110) Surface of Lithium Oxide.
190. M. Causà, R. Dovesi, C. Pisani, and C. Roetti, *Surf. Sci.*, **175**, 551 (1986). Ab initio Hartree-Fock Study of the MgO (001) Surface.
191. F. R. Sensato, R. Custodio, M. Calatayud, A. Beltrán, J. Andrés, J. R. Sambrano, and E. Longo, *Surf. Sci.*, **511**, 408 (2002). Periodic Study on the Structural and Electronic Properties of Bulk, Oxidized and Reduced SnO<sub>2</sub> (110) Surfaces and the Interaction with O<sup>2-</sup>.
192. J. Muscat, N. M. Harrison, and G. Thornton, *Phys. Rev.*, **B59**, 2320 (1999). Effects of Exchange, Correlation, and Numerical Approximations on the Computed Properties of the Rutile TiO<sub>2</sub> (100) Surface.
193. J. Muscat and N. M. Harrison, *Surf. Sci.*, **446**, 119 (2000). The Physical and Electronic Structure of the Rutile (001) Surface.
194. A. Wander and N. M. Harrison, *Surf. Sci.*, **457**, L342 (2000). Ab Initio Study of ZnO (10-10).
195. A. Wander and N. M. Harrison, *Surf. Sci.*, **468**, L851 (2000). Ab initio study of ZnO (11-20).
196. A. Beltrán, J. Andrés, M. Calatayud, and J. B. L. Martins, *Chem. Phys. Lett.*, **338**, 224 (2001). Theoretical Study of ZnO (10-10) and Cu/ZnO (10-10) Surfaces.
197. T. Ouazzani, A. Lichanot, C. Pisani, and C. Roetti, *J. Phys. Chem. Solids*, **54**, 1603 (1993). Relaxation and Electronic Structure of Surfaces in Lithium Sulphide: a Hartree-Fock Ab Initio Approach.
198. K. M. Rosso, U. Becker, and M. F. Hochella, *Am. Mineral.*, **84**, 1535 (1999). Atomically Resolved Electronic Structure of Pyrite (100) Surfaces: An Experimental and Theoretical Investigation with Implications for Reactivity.
199. J. Muscat and J. D. Gale, *Geochim. Cosmochim. Acta*, **67**, 799 (2003). First Principle Studies of the Surface of Galena PbS.
200. F. Frechard and P. Sautet, *Surf. Sci.*, **336**, 146 (1995). Hartree-Fock Ab-Initio Study of the Geometric and Electronic Structure of RuS<sub>2</sub> and its (100) and (111) Surfaces.
201. K. Doll and N. M. Harrison, *Chem. Phys. Lett.*, **317**, 282 (2000). Chlorine Adsorption on the Cu (111) Surface.
202. K. Doll, *Eur. Phys. J. B*, **22**, 389 (2001). Density Functional Study of the Adsorption of K on the Cu (111) Surface.
203. K. Doll and N. M. Harrison, *Phys. Rev.*, **B63**, 165410 (2001). Theoretical Study of Chlorine Adsorption on the Ag(111) Surface.
204. K. Doll, *Phys. Rev.*, **B66**, 155421 (2002). Density-Functional Study of the Adsorption of K on the Ag (111) Surface.
205. K. Doll, *Surf. Sci.*, **544**, 103 (2003). Density Functional Study of Ni Bulk, Surfaces and the Adsorbate Systems Ni (111)( $\sqrt{3} \times \sqrt{3}$ )R30 degrees-Cl, and Ni (111)(2 × 2)-K.
206. A. Kokalj and M. Causà, *J. Phys.-Condens. Matter*, **11**, 7463 (1999). Periodic Density Functional Theory Study of Pt (111): Surface Features of Slabs of Different Thickness.
207. L. Fu, E. Yashenko, L. Resca, and R. Resta, *Phys. Rev.*, **B60**, 2697 (1999). Hartree-Fock Studies of Surface Properties of BaTiO<sub>3</sub>.
208. R. A. Evarestov, E. Kotomin, E. Heifets, J. Maier, and G. Borstel, *Sol. St. Commun.*, **127**, 367 (2003). Ab Initio Hartree-Fock Calculations of LaMnO<sub>3</sub> (110) Surfaces.
209. E. Heifets, R. I. Eglitis, E. A. Kotomin, J. Maier, and G. Borstel, *Phys. Rev.*, **B64**, 235417 (2001). Ab Initio Modeling of Surface Structure for SrTiO<sub>3</sub> Perovskite.
210. E. Heifets, R. I. Eglitis, E. A. Kotomin, J. Maier, and G. Borstel, *Surf. Sci.*, **513**, 211 (2002). First-Principles Calculations for SrTiO<sub>3</sub> Surface Structure.
211. G. Bussolin, S. Casassa, C. Pisani, and P. Ugliengo, *J. Chem. Phys.*, **108**, 9516 (1998). Ab Initio of HCl and HF Interaction with Crystalline Ice. I. Physical Adsorption.

212. B. Civalleri, S. Casassa, E. Garrone, C. Pisani, and P. Ugliengo, *J. Phys. Chem. B*, **103**, 2165 (1999). Quantum Mechanical Ab Initio Characterization of a Simple Periodic Model of the Silica Surface.
213. A. Damin, R. Dovesi, A. Zecchina, and P. Ugliengo, *Surf. Sci.*, **479**, 253 (2001). CO/MgO (001) at Different CO Coverages: A Periodic Ab Initio Hartree-Fock and B3-LYP Study.
214. R. Wichtendahl, M. Rodriguez-Rodrigo, U. Härtel, H. Kuhlenbeck, and H. J. Freund, *Phys. Status Solidi A*, **173**, 93 (1999). Thermodesorption of CO and NO from Vacuum-Cleaved NiO (100) and MgO (100).
215. R. Wichtendahl, M. Rodriguez-Rodrigo, U. Härtel, H. Kuhlenbeck, and H. J. Freund, *Surf. Sci.*, **423**, 90 (1999). TDS Study of the Bonding of CO and NO to Vacuum-Cleaved NiO (100).
216. P. Ugliengo and A. Damin, *Chem. Phys. Lett.*, **366**, 683 (2002). Are Dispersive Forces Relevant for CO Adsorption on the MgO (001) Surface?
217. J. F. Sanz and C. M. Zicovich-Wilson, *Chem. Phys. Lett.*, **303**, 111 (1999). A Periodic Hartree-Fock Study of Na Adsorption on the TiO<sub>2</sub> (110) Rutile Surface.
218. J. Muscat and N. M. Harrison, *Phys. Rev.*, **B59**, 15457 (1999). First-Principles Study of Potassium Adsorption on TiO<sub>2</sub> Surfaces.
219. L. Giordano, G. Pacchioni, T. Bredow, and J. F. Sanz, *Surf. Sci.*, **471**, 21 (2001). Cu, Ag, and Au Atoms Adsorbed on TiO<sub>2</sub> (110): Cluster and Periodic Calculations.
220. M. Calatayud, J. Andrés, and A. Beltrán, *Surf. Sci.*, **430**, 213 (1999). A Theoretical Analysis of Adsorption and Dissociation of CH<sub>3</sub>OH on the Stoichiometric SnO<sub>2</sub> Surface.
221. M. Melle-Franco and G. Pacchioni, *Surf. Sci.*, **461**, 54 (2000). CO Adsorption on the SnO<sub>2</sub> (110): Cluster and Periodic Ab Initio Calculations.
222. M. Melle-Franco, G. Pacchioni, and A. V. Chadwick, *Surf. Sci.*, **478**, 25 (2001). Cluster and Periodic Ab Initio Calculations on the Adsorption of CO<sub>2</sub> on the SnO<sub>2</sub> (110) Surface.
223. T. Bredow and G. Pacchioni, *Surf. Sci.*, **373**, 21 (1997). Comparative Periodic and Cluster Ab Initio Study on Cu<sub>2</sub>O (111)/CO.
224. P. Persson and L. Ojamäe, *Chem. Phys. Lett.*, **321**, 302 (2000). Periodic Hartree-Fock Study of the Adsorption of Formic Acid on ZnO (10-10).
225. A. Wander and N. M. Harrison, *J. Phys. Chem. B*, **105**, 6191 (2001). Ab Initio Study of Hydrogen Adsorption on ZnO (10-10).
226. K. M. Rosso, U. Becker, and M. F. Hochella, *Am. Mineral.*, **84**, 1549 (1999). The Interaction of Pyrite (100) Surfaces with O<sub>2</sub> and H<sub>2</sub>O: Fundamental Oxidation Mechanisms.
227. B. Civalleri and P. Ugliengo, *J. Phys. Chem. B*, **104**, 9491 (2000). First Principle Calculations of the Adsorption of NH<sub>3</sub> on a Periodic Model of Silica Surface.
228. J. R. B. Gomes, F. Illas, N. Cruz Hernandez, J. F. Sanz, A. Wander, and N. M. Harrison, *J. Chem. Phys.*, **116**, 1684 (2002). Surface Model and Exchange-Correlation Functional Effects on the Description of Pd/ $\alpha$ -Al<sub>2</sub>O<sub>3</sub> (0001).
229. M. Sgroi, C. Pisani, and M. Busso, *Thin Solid Films*, **400**, 64 (2001). Ab Initio Density Functional Simulation of Structural and Electronic Properties of MgO Ultra-Thin Adlayers on the (001) Ag Surface.
230. C. Giovanardi, A. di Bona, T. S. Moia, S. Valeri, C. Pisani, M. Sgroi, and M. Busso, *Surf. Sci. Lett.*, **505**, L209 (2002). Experimental and Theoretical Study of the MgO/Ag (001) Interface.
231. E. Heifets, E. A. Kotomin, and R. Orlando, *J. Phys.-Condens. Matter*, **8**, 6577 (1996). Periodic Hartree-Fock Simulation of the Ag/MgO Interface Structure.
232. E. Heifets, Y. F. Zhukovskii, E. A. Kotomin, and M. Causà, *Chem. Phys. Lett.*, **238**, 395 (1998). The Adhesion Nature of the Ag/MgO (100) Interface: An Ab Initio Study.
233. Y. F. Zhukovskii, E. A. Kotomin, P. W. M. Jacobs, A. M. Stoneham, and J. H. Harding, *Surf. Sci.*, **441**, 373 (1999). Comparative Theoretical Study of the Ag-MgO (100) and (110) Interfaces.



- 
234. Y. F. Zhukovskii, E. A. Kotomin, D. Fuks, S. Dorfman, and A. Gordon, *Surf. Sci.*, **482**, 66 (2001). Hartree-Fock Study of Adhesion and Charge Redistribution on the Ag/MgO (001) Interface.
235. E. Kotomin, J. Maier, Y. F. Zhukovskii, D. Fuks, and S. Dorfman, *Mater. Sci. Eng. C*, **23**, 247 (2003). Ab Initio Modelling of Silver Adhesion on the Corundum (0001) Surface.
236. S. Casassa, A. M. Ferrari, M. Busso, and C. Pisani, *J. Phys. Chem. B*, **106**, 12978 (2002). Structural, Electronic and Magnetic Properties of the NiO Monolayer Epitaxially Grown on the (001) Ag Surface: an Ab-Initio Density Functional Study.
237. M. D. Towler, N. M. Harrison, and M. I. McCarthy, *Phys. Rev.*, **B52**, 5375 (1995). Ab Initio Study of the Surface and Interfacial Properties of a Layered MgO/NiO Film.
238. W. C. Mackrodt, C. Noguera, and N. L. Allan, *Faraday Discuss.*, **114**, 105 (1999). A Study of the Electronic, Magnetic, Structural and Dynamic Properties of Low-Dimensional NiO on MgO (100) Surfaces.
239. B. Henderson, *Defects in Crystalline Solids*, Arnold, London, 1972.
240. W. Haynes and A. M. Stoneham, *Defects and Defect Processes in Non Metallic Solids*, Wiley, New York, 1984.
241. G. Pacchioni, P. S. Bagus, and F. Parmigiani, Eds., *Cluster Models for Surface and Bulk Phenomena*, Vol. 283, Plenum Press, New York, 1992.
242. F. Maseras and K. Morokuma, *J. Comput. Chem.*, **16**, 1170 (1995). IMOMM: A New Integrated Ab Initio + Molecular Mechanics Geometry Optimization Scheme of Equilibrium Structures and Transition States.
243. S. Dapprich, I. Komáromi, K. S. Byun, K. Morokuma, and M. J. Frisch, *J. Mol. Struct. (THEOCHEM)*, **462**, 1 (1999). A New ONIOM Implementation in Gaussian98. Part I. The Calculation of Energies, Gradients, Vibrational Frequencies and Electric Field Derivatives.
244. T. Vreven and K. Morokuma, *J. Comput. Chem.*, **21**, 1419 (2000). On the Application of the IMOMO (Integrated Molecular Orbital Plus Molecular Orbital) Method.
245. U. Eichler, C. M. Kölmel, and J. Sauer, *J. Comput. Chem.*, **18**, 463 (1997). Combining Ab-Initio Techniques with Analytical Potential Functions for Structure Predictions of Large Systems: Method and Applications to Crystalline Silica Polymorphs.
246. J. R. Shoemaker, L. W. Burggraf, and M. S. Gordon, *J. Phys. Chem. A*, **103**, 3245 (1999). SIMOMM: An Integrated Molecular Orbital / Molecular Mechanics Optimization Scheme for Surfaces.
247. C. Pisani, R. Dovesi, C. Roetti, M. Causà, R. Orlando, S. Casassa, and V. R. Saunders, *Int. J. Quantum Chem.*, **77**, 1032 (2000). CRYSTAL and EMBED, Two Computational Tools for the Ab Initio Study of the Electronic Properties of Crystals.
248. C. Pisani, *J. Mol. Struct. (THEOCHEM)*, **621**, 141 (2003). Local Techniques for the Ab-Initio Quantum Mechanical Treatment of the Chemical Properties of Crystalline Materials.
249. C. Pisani, in *Handbook of Molecular Physics and Quantum Chemistry*, Vol. 3, S. Wilson, Ed., Wiley: Chichester, United Kingdom, 2003, p. 339. Molecules at Crystalline Surfaces and in Crystal Cages.
250. M. Leslie and M. J. Gillan, *J. Phys. C*, **18**, 973 (1985). The Energy and Elastic Dipole Tensor of Defects in Ionic Crystals Calculated by the Supercell Method.
251. G. Makov and M. C. Payne, *Phys. Rev.*, **B51**, 4014 (1995). Periodic Boundary-Conditions in Ab-Initio Calculations.
252. U. Gerstmann, P. Deák, R. Rurali, B. Aradi, T. Frauenheim, and H. Overhof, *Phys. Rev. Lett.*, submitted (2004). Accurate Treatment of Charged Defects within Periodic Boundary Conditions.
253. G. Mallia, R. Orlando, C. Roetti, P. Ugliengo, and R. Dovesi, *Phys. Rev.*, **B63**, 235102 (2001). F Center in LiF: A Quantum Mechanical Ab Initio Investigation of the Hyperfine Interaction Between the Unpaired Electron and the Vacancy and its First Seven Neighbors.
254. H. Seidel and H. C. Wolf, *The Physics of Colour Centers*, Fowler, New York, 1968.

255. W. C. Holton and H. Blum, *Phys. Rev.*, **125**, 89 (1962). Paramagnetic Resonance of F Centers in Alkali Halides.
256. A. Lichanot, C. Larrieu, R. Orlando, and R. Dovesi, *J. Phys. Chem. Solids*, **59**, 7 (1998). Lithium Trapped-Hole Centre in Magnesium Oxide. An Ab-Initio Supercell Study.
257. A. Lichanot, C. Larrieu, C. M. Zicovich-Wilson, C. Roetti, R. Orlando, and R. Dovesi, *J. Phys. Chem. Solids*, **59**, 1119 (1998). Trapped-Hole Centres Containing Lithium and Sodium in MgO, CaO and SrO. An Ab Initio Supercell Study.
258. A. Lichanot, R. Orlando, G. Mallia, M. Merawa, and R. Dovesi, *Chem. Phys. Lett.*, **318**, 240 (2000).  $V_{OH}$  Center in Magnesium Oxide: An Ab Initio Supercell Study.
259. A. Lichanot, P. Baranek, M. Merawa, R. Orlando, and R. Dovesi, *Phys. Rev.*, **B62**, 12812 (2000).  $V_{OH}$  and  $V_{OD}$  Centers in Alkaline-Earth Oxides: An Ab Initio Supercell Study.
260. M. M. Abraham, W. P. Unruh, and Y. Chen, *Phys. Rev.*, **B10**, 3540 (1974). Electron-Nuclear-Double-Resonance Investigations of  $[Li]^0$  and  $[Na]^0$  Centers in MgO, CaO, and SrO.
261. Y. Chen and M. M. Abraham, *J. Phys. Chem. Solids*, **51**, 747 (1990). Trapped-Hole Centers in Alkaline-Earth Oxides.
262. O. F. Schirmer, *J. Phys. Chem. Solids*, **29**, 1407 (1968). The Structure of the Paramagnetic Lithium Center in Zinc Oxide and Beryllium Oxide.
263. R. Orlando, R. Dovesi, P. Azavant, N. M. Harrison, and V. R. Saunders, *J. Phys.-Condens. Matter*, **6**, 8573 (1994). A Super-Cell Approach for the Study of Localized Defects in Solids: Carbon Substitution in Bulk Silicon.
264. R. Orlando, P. Azavant, M. D. Towler, R. Dovesi, and C. Roetti, *J. Phys.-Condens. Matter*, **8**, 1123 (1996). Cluster and Supercell Calculations for Carbon-doped Silicon.
265. C. Freyria-Fava, R. Dovesi, V. R. Saunders, M. Leslie, and C. Roetti, *J. Phys.-Condens. Matter*, **5**, 4793 (1993). Ca and Be Substitution in Bulk MgO: Ab Initio Hartree-Fock and Ionic Model Supercell Calculation.
266. A. V. Puchina, V. E. Puchin, E. A. Kotomin, and M. Reichling, *Solid State Commun.*, **106**, 285 (1998). Ab Initio Study of the F Centers in  $CaF_2$ : Calculations of the Optical Absorption, Diffusion and Binding Energies.
267. W. C. Mackrodt, *Ber. Bunsenges. Phys. Chem.*, **101**, 169 (1997). The Nature of Valence Band Holes in Pure and Fe-Doped NiO: An Ab Initio Hartree-Fock Study.
268. W. C. Mackrodt, N. M. Harrison, V. R. Saunders, N. L. Allan, and M. D. Towler, *Chem. Phys. Lett.*, **250**, 66 (1996). Direct Evidence of O(P) Holes in Li-Doped NiO from Hartree-Fock Calculations.
269. R. A. Evarestov and V. P. Smirnov, *Phys. St. Sol.*, **B215**, 949 (1999). Supercell Model of V-Doped  $TiO_2$ : Unrestricted Hartree-Fock Calculations.
270. C. M. Zicovich-Wilson and R. Dovesi, *J. Phys. Chem.*, **102**, 1411 (1998). Titanium Containing Zeolites. A Periodic Ab-Initio Hartree Fock Characterization.
271. A. Damin, S. Bordiga, A. Zecchina, K. Doll, and C. Lamberti, *J. Chem. Phys.*, **118**, 10183 (2003). Ti-Chabazite as a Model System of Ti(IV) in Ti-Zeolites: A Periodic Approach.
272. R. Orlando, F. Corà, R. Millini, G. Perego, and R. Dovesi, *J. Chem. Phys.*, **105**, 8937 (1996). Hydrogen Abstraction from Methane by Li Doped MgO: A Periodic Quantum Mechanical Study.
273. D. M. Gruen, P. C. Redfern, D. A. Horner, P. Zapol, and L. A. Curtiss, *J. Phys. Chem. B*, **103**, 5459 (1999). Theoretical Studies on Nanocrystalline Diamond: Nucleation by Dicarbon and Electronic Structure of Planar Defects.
274. R. Orlando, R. Dovesi, C. Roetti, and V. R. Saunders, *Chem. Phys. Lett.*, **228**, 225 (1994). Convergence Properties of the Cluster Model in the Study of Local Perturbations in Ionic Systems. The Case of Bulk Defects in MgO.
275. M. D. Segall, P. J. D. Lindan, M. J. Probert, C. J. Pickard, P. J. Hasnip, S. J. Clark, and M. C. Payne, *J. Phys.-Condens. Matter*, **14**, 2717 (2002). First-Principles Simulation: Ideas, Illustrations and the CASTEP Code.

- 
276. W. Andreoni and A. Curioni, *Parallel Computing*, **26**, 819 (2000). New Advances in Chemistry and Material Science with CPMD and Parallel Computing.
277. M. Bockstedte, A. Kley, J. Neugebauer, and M. Scheffler, *Comput. Phys. Commun.*, **107**, 187 (1997). Density-Functional Theory Calculations for Polyatomic Systems: Electronic Structure, Static and Elastic Properties and Ab Initio Molecular Dynamics.
278. B. G. Pfrommer, J. Demmel, and H. Simon, *J. Comput. Phys.*, **150**, 287 (1999). Unconstrained Energy Functionals for Electronic Structure Calculations.
279. S. Baroni, A. Dal Corso, S. de Gironcoli, and P. Giannozzi. PWSCFV.2.1, 2004. Available: <http://www.pwscf.org>.
280. G. Y. Sun, J. Kurti, P. Rajczy, M. Kertesz, J. Hafner, and G. Kresse, *J. Mol. Struct. (THEOCHEM)*, **624**, 37 (2003). Performance of the Vienna Ab Initio Simulation Package (VASP) in Chemical Applications.
281. P. E. Blöchl, C. J. Först, and J. Schimpl, *Bull. Mater. Sci.*, **26**, 33 (2003). Projector Augmented Wave Method: Ab-Initio Molecular Dynamics with Full Wave Functions.
282. A. R. Tackett, N. A. W. Holzwarth, and G. E. Matthews, *Comput. Phys. Commun.*, **135**, 329 (2001). A Projector Augmented Wave (PAW) Code for Electronic Structure Calculations, Part I: Atompaw for Generating Atom-Centered Functions.
283. A. R. Tackett, N. A. W. Holzwarth, and G. E. Matthews, *Comput. Phys. Commun.*, **135**, 348 (2001). A Projector Augmented Wave (PAW) Code for Electronic Structure Calculations, Part II: Pwpaw for Periodic Solids in a Plane Wave Basis.
284. G. Lippert, J. Hutter, and M. Parrinello, *Mol. Phys.*, **92**, 477 (1997). A Hybrid Gaussian and Plane Wave Density Functional Scheme.
285. J. M. Soler, E. Artacho, J. D. Gale, A. García, J. Junquera, P. Ordejón, and D. Sánchez-Portal, *J. Phys.-Condens. Matter*, **14**, 2745 (2002). The Siesta Method for Ab Initio Order-N Materials Simulation.
286. B. Delley, *J. Chem. Phys.*, **113**, 7756 (2000). From Molecules to Solids with the DMol3 Approach.
287. G. Bihlmayer and S. Blügel, *Comput. Phys. Commun.*, in press (2005). FLEUR: a Parallelized Electronic Structure Code for Surfaces and Bulk.
288. K. Schwarz and P. Blaha, *Comput. Mat. Sci.*, **28**, 259 (2003). Solid State Calculations Using WIEN2k.
289. G. te Velde, F. M. Bickelhaupt, E. J. Baerends, C. Fonseca Guerra, S. J. A. van Gisbergen, J. G. Snijders, and T. Ziegler, *J. Comput. Chem.*, **22**, 931 (2001). Chemistry with ADF.
290. K. N. Kudin and G. E. Scuseria, *Phys. Rev.*, **B61**, 16440 (2000). Linear-Scaling Density-Functional Theory with Gaussian Orbitals and Periodic Boundary Conditions: Efficient Evaluation of Energy and Forces via the Fast Multipole Method.
291. S. Swaminathan, B. N. Craven, and R. K. McMullan, *Acta Crystallogr. Sec. B*, **40**, 300 (1984). The Crystal Structure and Molecular Thermal Motion of Urea at 12, 60 and 123 K from Neutron Diffraction.
292. I. J. Bush, *Capability Computing*, **1**, 10 (2003). CRYSTAL on HPCx: Ab Initio Study of Proteins.



# Molecular Quantum Similarity: Theory and Applications

Patrick Bultinck,<sup>\*</sup> Xavier Gironés,<sup>†</sup> and  
Ramon Carbó-Dorca<sup>‡</sup>

<sup>\*</sup>*Department of Inorganic and Physical Chemistry,  
Ghent University, Krijgslaan 281 (S-3), B-9000 Gent, Belgium*

<sup>†</sup>*Computational Chemistry, Medicinal Chemistry Department,  
AstraZeneca R&D, Pepparedsleden 1, S-43183 Mölndal, Sweden*

<sup>‡</sup>*Institute of Computational Chemistry, University of Girona,  
Campus de Montilivi, 17005 Girona, Spain*

---

---

## INTRODUCTION

Similarity is an impressively broad topic and ubiquitous in virtually any field of human activity, including natural sciences. In this context, it seems to be a paradox that such an important and widely used concept by chemists is often so poorly defined. A paradox is a statement that, although true, seems false and even selfcontradictory.<sup>1</sup> As such, it is clearly applicable to the often observed lack of a clear and general definition of similarity. Precisely because of the widespread use of the similarity concept by chemists, ranging from exact up to human sciences, giving an exact all-covering definition of similarity is a difficult task. Moreover, similarity and its perception depend on the observer. A clear example is found in the field of human psychology, where it was shown that the degree of similarity perceived by observers between different facial expressions, depends on their mental state.<sup>2,3</sup> This paradox and the

dependence of its perception on observers makes it clear that similarity is often considered as a qualitative concept.

The ubiquitous nature of the similarity concept was well expressed by Rouvray, who stated that all issues of comparison, and thus of classification, are in fact related to similarity.<sup>4</sup> It is well known that both procedures play a crucial role in chemistry. Aspects such as basicity and acidity are based on classification, the periodic system involves comparison and classification, and so on. As a consequence, similarity appears to be an important subject in chemistry.

This chapter deals with quantum similarity. Contrary to what is stated previously, the concept of similarity here is clearly and well defined from basic aspects of quantum chemistry and the quantum mechanical postulates. The combination of the ubiquitous nature of similarity, especially in chemistry, with the possibility offered by quantum chemistry to provide a neat definition for similarity, is the reason why quantum similarity is a major topic of scientific interest. This review in quantum similarity aims both at providing a clear introduction for nonexperts as well as at discussing the latest developments in the field. We will also indicate some shortcomings of the theory and highlight some caveats that may not immediately be clear and will refer to some sources of arbitrariness that may influence the perception of the degree of similarity. Examples of the latter include choices in similarity descriptor type, distance definitions, and alignment issues. Although several reviews on this topic have been published earlier,<sup>5-12</sup> this chapter gives the latest report on advances in quantum similarity, starting from an easily accessible level, and thus fulfilling the pedagogical theme of this book series.

The next section, entitled Basic Aspects of Molecular Similarity, gives a general overview of molecular similarity and the usual vocabulary used by chemists in this field. In the section entitled The Electron Density as Molecular Descriptor, some elements of quantum chemistry will be presented. These sections should not be expected by the reader to offer a rigorous discussion of all aspects of such broad fields, and therefore, only the most important definitions and concepts will be introduced. The following sections will then address extensively the subject of molecular quantum similarity in both theoretical and practical aspects. An ample references list will help the interested reader to look up the more specialized literature.

---

## BASIC ELEMENTS OF MOLECULAR SIMILARITY

It is worthwhile to introduce molecular similarity as a broader field before starting the discussion on molecular quantum similarity. Molecular similarity is a well-studied area and continues to be a major topic in modern chemical research. It is that area in which we look for methods to identify a degree of similarity between molecules in a dataset or in which we apply these

methods to identify groups of more or less similar molecules. A question we may raise is: Where does this interest in molecular similarity come from? Several plausible answers exist, but these may be mostly retraced to the central chemical idea that *similar molecules have similar properties*. The huge interest in this field within the pharmaceutical context of structure-activity relationships can thus be easily understood (see, for example, Dean<sup>13</sup>). If a certain medicinally active molecule or even a drug-like molecule is known, but exhibits some specific drawbacks, the molecular similarity concept can find similar molecules in a molecular dataset. It is hoped that those molecules, having some relevant degree of similarity, will have similar biological activity but, perhaps, without the unwanted drawbacks. This example is only one possibility; applications of similarity are naturally also found in many other fields of chemical research.

An issue that immediately develops is how to express the degree of similarity between molecules. In what computationally implementable way can molecular similarity be expressed? The fact is that many different ways exist. Only a few of them can be discussed here. The following statement by Herndon and Bertz<sup>14</sup> describes well a basic property of similarity: “Similarity, like beauty, lies in the eyes of the beholder.” Although this statement might seem true for more human applications of the similarity concept, it is also true to some extent for its chemical application. Although emotional aspects are largely absent in chemical applications, the beholder introduces some personal character into his or her perception of similarity. A clear example of this introduction occurs when chemists are asked to, even qualitatively, discern the similarity between different molecules. An organic chemist may use the reactivity of the molecules in a certain type of reaction, whereas a quantum chemist may look at other molecular features, and a physical chemist will even use other properties to discern the similarity between molecules. Clearly then, different chemists use different molecular descriptors to express degrees of molecular similarity. It is appropriate in this context to introduce the concept of a molecular descriptor and to present some typical examples, which will make it easier to introduce quantum chemistry as an alternative path to defining similarity with the electron density  $\rho$  as a descriptor. A molecular descriptor is a feature of the molecule under study. Examples may be molecular weight, boiling point, density, and many other such properties. These features or descriptors can then express degrees of similarity between molecules, as is often done in clustering, combinatorial library design, data mining, and so on. They also play an important role in finding quantitative structure property/activity relationships (QSP/ARs). Over time, several molecular descriptors have been used by chemists in different applications, and it would go far beyond the scope of the present chapter to discuss all of these. Reviews of molecular descriptors may be found in publications by Downs,<sup>15</sup> Brown,<sup>16</sup> Mason and Picket,<sup>17</sup> and Bajorath.<sup>18</sup> A specific class of molecular descriptors are those derived from quantum chemistry, whereby the molecular wave function or the

electron density yields molecular descriptors that are observables or even new classes of descriptors that are not observables, such as atomic charges. A good account of the application of quantum chemistry to obtain molecular descriptors may be found in the work of Karelson et al.<sup>19,20</sup>

Different ways exist to classify molecular descriptors. We will adopt the classification of Downs.<sup>15</sup> Most common descriptor types are largely atom based. These descriptors are based on the individual atoms composing the molecule, but with their description extended to incorporate information about the environment of the atom. To introduce the radically different path taken in quantum similarity as opposed to the more “classic” procedures in most similarity studies, a small presentation of different types of molecular descriptors first will be given.

- **Feature counts:** This type represents a simple class of molecular descriptors, based on simple calculation of the number of specific features in a molecule. Such a specific descriptor may be, for instance, the number of atoms of a certain element in the molecule, the number of chemical bonds of a specific type, and the number of hydrogen bond acceptors. It is clear that this is a simple descriptor set, and it might be expected that this kind of descriptor would be superseded by more informative descriptors. Yet, several of the descriptors used in the so-called Lipinski rule-of-five are simple summations or feature counts. The Lipinski rule-of-five is a quick method to classify molecules as being drug-like or non-drug-like on the basis of mostly simple feature count descriptors.<sup>21</sup> Experience has shown that although a feature count is a simple method, it performs well in distinguishing drug- and non-drug-like molecules.
- **Physicochemical parameters:** Physicochemical parameters are common molecular descriptors. The most well-known descriptor of this class in QSP/AR is logP.<sup>22</sup> Other examples include the heat of formation, and the Hammett  $\sigma$  constant.
- **Fragment descriptors:** Both two-dimensional (2-D) connection tables for the molecular structures involved, as well as three-dimensional (3-D) information for these structures may generate a variety of fragment descriptors for substructure search systems (for a review of substructure search systems, see Chen<sup>23</sup>). Active research continues to be carried out in this field, and the interested reader can peruse the review by Downs<sup>15</sup> for more information.
- **Topological and topographical indices:** In QSP/AR, this class of descriptors is extensively used by chemists. Extensive literature on the subject exists, and many novel indices are still being introduced. Among the classic ones, the Wiener index is best known.<sup>24</sup> This index is



calculated as half the sum of all path lengths (= the number of bonds) between two atoms *i* and *j* in the molecule. To name but a few other well-known descriptors of this type, one can quote the Balaban index,<sup>25</sup> the indices introduced by Randic,<sup>26–30</sup> the Zagreb index,<sup>31,32</sup> and the Hosoya index.<sup>33</sup>

- **Field-based descriptors:** The previous indices are mostly atom-based, and although they have proven useful to chemists, more and more work has been done and is currently being done on so-called field-based descriptors. Examples of some fields are the electron density, steric fields, electrostatic potentials, and hydrophobic fields. This definition is the proper starting point for the discussion of quantum similarity, because the electron density plays an important role in these fields as we will discuss extensively. This class of field-based descriptors also is used in 3-D QSAR.<sup>34</sup>

Although we have mentioned the electron density as being only one of a large set of possible molecular descriptors, several reasons exist for why this should be regarded as the ultimate entity to study molecular similarity. Although the electron density is obtained from a quantum chemical calculation, it is important to distinguish two different applications of quantum chemistry in molecular similarity. Often, the descriptors needed for a set of molecules may be obtained from tabulated physicochemical data, or they may be obtained from tabulated fragment contributions. The latter fragment-based approach has been successful in many applications such as rationalizing observed differences in properties for sets of molecules differing in their substitution pattern. Sometimes this fragment-based approach is inadequate, as for instance, when some fragment contribution is not available for a molecule, because it has not yet been synthesized or because an experimental determination of the descriptor is not possible. Needless to say, this situation is common when looking for drug-lead molecules. In these cases, quantum chemistry can supply the needed data computationally. Another possible scenario for quantum chemistry is to obtain new kinds of molecular descriptors that are not observables, i.e., for which no unique operator can be drawn from the classic-quantum correspondence principle for operators. Good examples of such descriptors are atomic charges<sup>35,36</sup> or several interesting quantities from conceptual density functional theory<sup>37</sup> (DFT). In both cases, quantum chemistry is often a black box, capable of filling holes in the existing dataset for the problem at hand. In molecular quantum similarity, explicit reference is made to the fundamental postulates of quantum mechanics, rather than to quantum chemistry as a black box. As such, molecular quantum similarity is a research domain where an original theoretical concept has been developed and where well-defined computational tools have been constructed for its practical application.

## THE ELECTRON DENSITY AS MOLECULAR DESCRIPTOR

It is not the aim of this chapter to give a detailed account of quantum chemistry, but it is important to introduce some of its elements that will be needed to comprehend the rest of the chapter. One of the basic postulates of quantum mechanics applied to chemistry states that for every molecule or state a wave function exists that is all-determining. This means that, once the wave function is known, every so-called observable property for an  $N$ -electron molecule may be obtained by straightforward integration as in Eq. [1]:

$$\Omega[\Psi] \equiv \int \int \cdots \int \Psi^*(\mathbf{x}_1, \dots, \mathbf{x}_N) \hat{\Omega} \Psi(\mathbf{x}_1, \dots, \mathbf{x}_N) d\mathbf{x}_1 \cdots d\mathbf{x}_N \quad [1]$$

where  $\hat{\Omega}$  is the appropriate operator for that specific observable property. In Eq. [1], the remaining symbols have the well-known meanings from standard textbook notations (see, e.g., Szabo and Ostlund<sup>38</sup>). The operator  $\hat{\Omega}$  is obtained with the correspondence principle between classical physics and quantum mechanics (for details, see the books by Levine<sup>39</sup> or Pilar<sup>40</sup>).

Despite the fundamental role played by the wave function in quantum chemistry, it has no physical meaning in itself. The square of the wave function given in Eq. [2], in contrast, is interpreted as the probability of finding an electron in the volume  $d\mathbf{r}$  around  $\mathbf{r}$ . For ease of notation, consider a wave function. The aforementioned probability density becomes

$$\Psi^*(\mathbf{r})\Psi(\mathbf{r})d\mathbf{r} = \rho(\mathbf{r}) \quad [2]$$

It is this concept that plays the most important role in molecular quantum similarity. Consider now the wave function shown in Eq. [1]. It is well known how to reduce the dimensionality of the density function  $\rho$  in Eq. [2] by integration. The density function  $\rho$  is generated from the wave function, which belongs to  $H(\mathbf{C})$ , a Hilbert space over the complex field. The density function belongs to  $H(\mathbf{R}^+)$ , a Hilbert space over the positive real field only. In other words, generating the density function occurs through the following rule:

$$\forall \Psi \in H(\mathbf{C}) \rightarrow \exists \rho = \Psi^* \Psi = |\Psi|^2 \in H(\mathbf{R}^+) \quad [3]$$

If a spinless  $n$ -th order density function is desired for an  $N$ -electron molecule, integration should be performed over  $(N-n)$  spatial coordinate sets and all  $N$  spin variables. The first-order electron density in a point  $\mathbf{r}$  is then:

$$\rho(\mathbf{r}_1) = N \int \int \cdots \int \Psi^*(\mathbf{x}_1, \dots, \mathbf{x}_N) \Psi(\mathbf{x}_1, \dots, \mathbf{x}_N) ds_1 d\mathbf{x}_2 \cdots d\mathbf{x}_N \quad [4]$$

Note that in Eq. [4], the integration is carried out over all coordinates (spatial and spin) of the electrons 2 to  $N$ , and over the spin only for the first electron,

where  $s$  has denoted the spin variable. Naturally one can also introduce the joint probability of finding simultaneously an electron at points  $\mathbf{r}_1$  and  $\mathbf{r}_2$ . This second-order electron density has a similar expression as in Eq. [4]:

$$\rho(\mathbf{r}_1, \mathbf{r}_2) = N(N-1) \int \cdots \int \Psi^*(\mathbf{x}_1, \dots, \mathbf{x}_N) \Psi(\mathbf{x}_1, \dots, \mathbf{x}_N) ds_1 ds_2 d\mathbf{x}_3 \dots d\mathbf{x}_N \quad [5]$$

Although one can extend this discussion to even higher order electron densities,  $\rho(\mathbf{r}_1)$  and  $\rho(\mathbf{r}_1, \mathbf{r}_2)$  are the most commonly used quantities in molecular quantum similarity. It is worth noting that often no order is mentioned in publications when considering electron density. It is then commonly accepted that one then refers to the first-order density. In the remainder of this chapter, we assume that the first-order density is used, except when an order is mentioned explicitly. An extremely important property of  $\rho(\mathbf{r}_1)$  is its positive definite nature, which may seem like a simple consequence of its probability nature, but this point can hardly be overemphasized for its application in quantum similarity, as will be shown.

From these straightforward equations, we might be tempted to think that the calculation of the electron density is a relatively trivial task; the truth is, however, that the wave function  $\Psi(\mathbf{x}_1, \dots, \mathbf{x}_N)$  is required first. Obtaining this wave function is far from a trivial task and is the main purpose of most wave function-based quantum chemical methods. A plethora of methods exists, which range from semi-empirical methods<sup>41,42</sup> up to highly advanced correlated methods.<sup>43,44</sup>

A somewhat different route to obtaining the electron density is used in DFT. In DFT, the central concept is not the wave function but directly the electron density. The electron density is considered to be the all-determining property of a molecule. DFT has an important advantage over wave function-based methods because the wave function for an  $N$ -electron molecule contains  $4N$  variables, whereas  $\rho(\mathbf{r}_1)$  contains only three variables.

As stated, the most important concept in molecular quantum similarity is the electron density. The idea of the electron density as the ultimate molecular descriptor is founded on the basic elements of quantum mechanics. It is the all-determining quantity in DFT, and it holds a close relation to the wave function. It is therefore appropriate in this context to raise the question of whether the electron density can really be considered as the all-determining entity in quantum similarity studies. Clear indications of this conclusion were described by Handy and are attributed to Wilson,<sup>45</sup> although initial ideas can also be traced back to Born<sup>46</sup> and von Neumann.<sup>47</sup> The electron density  $\rho(\mathbf{r})$  has several important features. First, integrated over all space, it gives the number of electrons:

$$\int \rho(\mathbf{r}) d\mathbf{r} = N \quad [6]$$

Second, information is obtained on the nature of the nuclei in the molecule from the cusp condition.<sup>48</sup> Third, the Hohenberg–Kohn theorems point out that besides determining the number of electrons, the density also determines the external potential that is present in the molecular Hamiltonian.<sup>49–51</sup> Once the number of electrons is known from Eq. [6] and because the external potential is determined by the electron density, the Hamiltonian is completely determined. Once the electronic Hamiltonian is determined, we can solve Schrödinger’s equation for the wave function, subsequently determining all observable properties of the system.

A convincing argument is then formed that the electron density is the most fundamental molecular descriptor. As stated by Dean,<sup>52</sup> similarity in the activity or properties of molecules will occur when the molecules show similar electron densities. Importantly, this process does not require a similar molecular skeleton or bonding pattern. With the electron density as the main molecular descriptor, we can study the similarity of molecules that are structurally dissimilar. Moreover, if a link can be made between some specific feature of the electron density and the biological activity, one could deduce a scaffold of electron density into which a new molecule should fit, which in turn opens the way for drug design to surpass the classic techniques that include limitations like variation of substituents on a known active compound or staying within a congeneric series.

---

## MOLECULAR QUANTUM SIMILARITY

Having established that the electron density is the basic molecular descriptor, and that a theoretical justification exists for its selection, the theory of molecular quantum similarity can now be developed.

In 1980, Carbó, Arnau and Leyda<sup>53</sup> were the first to use molecular quantum similarity. As an anecdote, in the submitted version of the manuscript, the title was “How far is one molecule from another?” After a reviewer’s comment, this title was changed to “How similar is one molecule to another?” The revised title has a much more obvious reference to similarity. In a sense, both titles are descriptive, because in that manuscript, the first degree of molecular similarity with a distance measure was presented. More precisely, a distance measure was introduced as

$$d_{AB}^2 = \int [\rho_A(\mathbf{r}) - \rho_B(\mathbf{r})]^2 d\mathbf{r} \quad [7]$$

In this equation, A and B denote two different molecules for which the electron densities are represented as  $\rho_A(\mathbf{r})$  and  $\rho_B(\mathbf{r})$  respectively. Equation [7] is an Euclidean distance between the electron densities of both molecules. Euclidean

distance is a specific kind of so-called Minkowski distances where  $k = 2$ . In general, Minkowski distances are<sup>54,55</sup>

$$d_{AB,k} = \sqrt[k]{\sum_{\alpha} |\alpha_A - \alpha_B|^k} \quad [8]$$

In the case of  $k = 2$ , Eq. [8] corresponds to the well-known Euclidean distance of which Eq. [7] is an integral version. In the case of  $k = 1$ , we find the Manhattan or city-block distance. The choice of Euclidean distance is a computationally interesting choice, but it is by no means the only one possible. In this context, it is appropriate to mention the four requirements that should be associated with a true distance measure:

$$d_{AB} \in \mathbf{R}_0^+ \quad [9]$$

$$d_{AA} = 0 \quad [10]$$

$$d_{AB} = d_{BA} \quad [11]$$

$$d_{AB} \leq d_{AC} + d_{CB} \quad [12]$$

When working with Euclidean distances between density functions, we can also work out Eq. [7], which yields

$$d_{AB}^2 = \int \rho_A(\mathbf{r})\rho_A(\mathbf{r})d\mathbf{r} + \int \rho_B(\mathbf{r})\rho_B(\mathbf{r})d\mathbf{r} - 2 \int \rho_A(\mathbf{r})\rho_B(\mathbf{r})d\mathbf{r} \quad [13]$$

The following notation is common in molecular quantum similarity:

$$Z_{AB} = \int \rho_A(\mathbf{r})\rho_B(\mathbf{r})d\mathbf{r} \quad [14]$$

It is called a molecular quantum similarity measure (MQSM), and it corresponds to the similarity integral between molecules A and B. In the first integral in Eq. [13], a so-called molecular quantum self-similarity measure (MQSSM) is obtained:

$$Z_{AA} = \int \rho_A(\mathbf{r})\rho_A(\mathbf{r})d\mathbf{r} \quad [15]$$

As we will describe, these self-similarity measures can have several relationships with different physicochemical properties.

The more similar two molecules A and B are, the smaller the (squared) Euclidean distance will be. It is immediately clear that the overlap of electron

densities of A and B is larger when the electron densities are more similar. In this sense, the distance is also a dissimilarity measure. Perfect similarity is characterized by  $d_{AB} = 0$ ; decreasing similarity is characterized by an increase in the distance and a decrease of  $Z_{AB}$ .

Following the distance measure of molecular (dis)similarity, another important similarity index was introduced.<sup>53</sup> Consider two vectors  $\mathbf{a}$  and  $\mathbf{b}$ . The classic scalar product between two vectors is calculated as

$$\begin{aligned}\vec{a} \cdot \vec{b} &= \|\mathbf{a}\| \|\mathbf{b}\| \cos \vartheta \\ \vec{a} \cdot \vec{b} &= x_a x_b + y_a y_b + z_a z_b\end{aligned}\quad [16]$$

where  $\|\mathbf{a}\|$  is the norm of vector  $\mathbf{a}$ . We can generalize this scalar product to much higher dimensions than simple Cartesian space, obtaining the following generalized cosine function:

$$C_{AB} = \frac{\int \rho_A(\mathbf{r}) \rho_B(\mathbf{r}) d\mathbf{r}}{\sqrt{\int \rho_A(\mathbf{r}) \rho_A(\mathbf{r}) d\mathbf{r} \int \rho_B(\mathbf{r}) \rho_B(\mathbf{r}) d\mathbf{r}}} \quad [17]$$

or

$$C_{AB} = \frac{Z_{AB}}{\sqrt{Z_{AA} Z_{BB}}} \quad [18]$$

This generalized cosine index is often called the Carbó index. Naturally the Carbó index is limited to the range (0,1), where  $C_{AB} = 1$  means perfect similarity. Still more (dis)similarity measures have been introduced and will be discussed further in this chapter. The range of the Carbó index naturally agrees with the Schwartz integral inequality:<sup>54,55</sup>

$$\left[ \int \rho_A(\mathbf{r}) \rho_B(\mathbf{r}) d\mathbf{r} \right]^2 \leq \int \rho_A^2(\mathbf{r}) d\mathbf{r} \int \rho_B^2(\mathbf{r}) d\mathbf{r} \quad [19]$$

When A equals B, the similarity measure of Eq. [14] is called the molecular quantum self-similarity measure, which corresponds to the square of the Euclidean norm of the density function. In other words,

$$Z_{AA} = \|\rho_A\|^2 \quad [20]$$

The work mentioned so far relies on quantum similarity indices that in turn rely on the overlap measure of electron densities. Simple overlap between electron density functions is, however, not the only possible choice. Based on the developed theory of vector semispaces and tagged sets, Carbó et al. have

developed the mathematical basis of molecular quantum similarity. A review of these mathematical aspects is delayed to a later section, but for the moment, it suffices that we can extend the ideas of quantum similarity to other operators, which give rise to similarity measures:

$$Z_{AB}[\Omega] = \iint \rho_A(\mathbf{r}_1) \Omega(\mathbf{r}_1, \mathbf{r}_2) \rho_B(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [21]$$

It is recognized that the previous development as in Eq. [13] is the special case when the operator in Eq. [21] is the Dirac delta function

$$\Omega(\mathbf{r}_1, \mathbf{r}_2) = \delta(\mathbf{r}_1 - \mathbf{r}_2) \quad [22]$$

Before proceeding to other operator types, it should be mentioned that the density functions introduced here are not the only possible way to carry out such analyses in molecular quantum similarity. It has been shown how extended wave functions may be derived that also hold, e.g., the gradient of the wave functions, which means that a new class of wave functions may be derived that are vector-like, much like what is found in relativistic quantum theory. However, it is beyond the scope of the present chapter to discuss this entire field, and the interested reader is referred to the literature, especially Carbó-Dorca et al.<sup>56</sup> where a clear discussion is given.

---

## EXTENSION TO OTHER OPERATORS

---

Until now, we have only used the Dirac delta function to yield molecular quantum similarity measures. The key equations in this regard are

$$Z_{AB} = \iint \rho_A(\mathbf{r}_1) \delta(\mathbf{r}_1 - \mathbf{r}_2) \rho_B(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [23]$$

or

$$Z_{AB} = \int \rho_A(\mathbf{r}_1) \rho_B(\mathbf{r}_1) d\mathbf{r}_1 \quad [24]$$

which define the quantitative expression of the overlap of two molecular electron density functions. As shown, this MQSM appears naturally when we introduce the concept of Euclidean distance between two molecular electron densities. The Dirac delta function is a point-by-point comparison with no reference to any other points.

The idea of molecular quantum similarity can be extended to other operators, provided they are positive definite. In this sense, they will lead to real, positive definite values for the MQSM evaluated over the density functions of the involved quantum objects.

The Dirac delta operator is an operator that does not introduce any weighting of the surrounding points in the overlap MQSM. A way of weighting the similarity measure, which does include the surrounding points, is to use as an operator  $|r_1 - r_2|^{-1}$ , which gives rise to a Coulomb style MQSM:

$$Z_{AB} = \iint \rho_A(\mathbf{r}_1) |\mathbf{r}_1 - \mathbf{r}_2|^{-1} \rho_B(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [25]$$

Naturally the name of this MQSM stems from the coincidence of Eq. [25] with Coulomb's law for two continuous electron distributions.

A third positive definite operator is a kinetic energy-like expression,<sup>57</sup> where the following MQSM is introduced:

$$Z_{AB} = - \iint \rho_A(\mathbf{r}_1) \nabla^2 \rho_B(\mathbf{r}_1) d\mathbf{r}_1 \quad [26]$$

The actual calculation of this integral is then performed using Green's identity theorem,<sup>54</sup> which turns Eq. [26] into

$$Z_{AB} = \iint \nabla \rho_A(\mathbf{r}_1) \nabla \rho_B(\mathbf{r}_1) d\mathbf{r}_1 \quad [27]$$

This equation is an overlap measure between the gradients of the electron density. Note that in passing from Eq. [21] to Eq. [27], a Dirac delta function is also implied.

As argued, positive definite values for the MQSM are obtained with positive definite operators only, because the involved density functions are positive definite. It is then straightforward for the density function to also be used a possible operator. For  $Z_{AB}$ , we can think of the MQSM in Eq. [28] as a similarity measure weighted by molecule C as a template molecule:

$$Z_{AB} = \iint \rho_A(\mathbf{r}_1) \rho_C(\mathbf{r}_1) \rho_B(\mathbf{r}_1) d\mathbf{r}_1 \quad [28]$$

Such MQSM have also been reported previously in the literature.<sup>58,59</sup> It is also immediately clear that we can generalize this MQSM to include more templates by

$$Z_{AB} = \iint \rho_A(\mathbf{r}_1) [\rho_C(\mathbf{r}_1) \rho_D(\mathbf{r}_1) \cdots] \rho_B(\mathbf{r}_1) d\mathbf{r}_1 \quad [29]$$

MQSM as in Eq. [28] evaluated over all molecules A and B in a set, again constitute a matrix in which one should also specify the molecular density function that acts as an operator:  $\mathbf{Z}_C = \{Z_{AB,C}\}$ . This matrix can be seen as



the representation of the operator  $\rho_C(\mathbf{r}_1)$  in the basis set of the density functions of all molecules in the set, including C. Naturally the entire collection of matrices may be considered, differing in the molecule used as an operator.

Another MQSM is obtained through the generalization of the self-similarity expression to any order  $n$  as follows:<sup>60</sup>

$$Z_{AA}^{(n)} = \int \Omega(\mathbf{r}) \rho_A^n(\mathbf{r}) d\mathbf{r} \quad [30]$$

When the operator is given by

$$\Omega(\mathbf{R}) = \sum_i \sum_{j>i} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \text{ and } \mathbf{R} = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n] \quad [31]$$

the expression

$$Z_{AA}^{(1)} = \int \sum_i \sum_{j>i} \frac{\rho_A(\mathbf{R})}{|\mathbf{r}_i - \mathbf{r}_j|} d\mathbf{R} \quad [32]$$

becomes a well-known quantum chemical quantity, namely, the electronic repulsion energy  $\langle V_{ee} \rangle$ , which justifies this quantity as a MQS-based molecular descriptor.<sup>61</sup>

Once an operator has been chosen for the calculation of the MQSM for a set of  $N$  molecules, one can calculate all MQSMs between every two molecules, which gives rise to the whole  $N \times N$  array of MQSM. This symmetrical matrix is called the molecular quantum similarity measure matrix (MQSMM), denoted  $\mathbf{Z}$ .

$$\mathbf{Z} = \begin{bmatrix} Z_{11} & \cdots & Z_{1N} \\ \vdots & \ddots & \vdots \\ Z_{N1} & \cdots & Z_{NN} \end{bmatrix} \quad [33]$$

Each of the columns or the rows of this matrix can be a representation of the corresponding molecular density function in the subspace formed by the density functions of the  $N$  molecules. Then, the  $I$ th column of  $\mathbf{Z}$ , denoted  $z_I$ , can be the representation of molecule  $I$  in the space spanned by the density functions of the  $N$  molecules. In this sense, these vectors reduce the infinite dimensional representations of the molecular density functions into an  $N$ -dimensional vector representation, where all numbers in the vector are real, positive definite values. These molecular vector representations have two important special properties. Once the operator has been chosen to

evaluate the elements, the vector descriptors are universal in the sense that they can be obtained for any molecule of the set and from any molecular set. Furthermore, they are unbiased except in the stage of the selection of the operator that evaluates the MQSM. Another way of looking at the MQSMM is to state that it is an  $N$ -dimensional matrix representation of the operator within the set of  $N$  density functions. Also, in this way, each molecule corresponds to a point in the  $N$ -dimensional similarity space. For the collection of all points in the set of molecules, one can describe the point cloud. Such molecular point clouds have been used previously to investigate possible relations between molecules and their properties.<sup>9,62-66</sup>

Another possibility not yet fully exploited in MQS corresponds to the construction of  $N \times M$  representations of an MQSMM, which can be done whenever two molecular sets A and B of  $N$  and  $M$  cardinality respectively, provide the corresponding density functions  $M_A = \{\rho_I^A(I = 1, N)\}$  and  $M_B = \{\rho_I^B(I = 1, M)\}$ . Then, the following  $N \times M$  MQSMM can be obtained:

$$Z^{AB} = \{Z_{IJ}^{AB} = \langle \rho_I^A | \Omega | \rho_J^B \rangle\} \quad [34]$$

This equation in fact corresponds to one set of molecules being a template for a second set.

---

## STOCHASTIC MANIPULATIONS AND GRAPHICAL REPRESENTATIONS OF THE SIMILARITY MATRIX

---

The similarity matrix  $Z$  obtained here, can naturally be further manipulated in a number of ways. One option that has been explored in some detail<sup>67-69</sup> is the stochastic transformation, which implies that all elements in a certain column of the matrix are divided by the sum of all elements in this column. Denoting the sum of the elements of a column  $I$  as  $\langle z_I \rangle$ , the stochastic matrix is given by

$$S = \begin{bmatrix} \frac{Z_{11}}{\langle z_1 \rangle} & \dots & \frac{Z_{1N}}{\langle z_N \rangle} \\ \vdots & \ddots & \vdots \\ \frac{Z_{N1}}{\langle z_1 \rangle} & \dots & \frac{Z_{NN}}{\langle z_N \rangle} \end{bmatrix} \quad [35]$$

It is immediately clear that the resulting matrix  $S$  is no longer a symmetrical matrix. On the other hand, the sum of all elements in a column of  $S$  belong to a unit shell vector semispace (see below), meaning that  $\forall I : \langle s_I \rangle = 1$ . It is clear that the columns of the stochastic matrix become a new descriptor set for the

molecules in an  $N$ -dimensional semispace, instead of in the infinite space of the density functions. The stochastic matrices could be symmetry transformed again, which is most easily done with any of the classic techniques or with inward matrix products.<sup>68</sup> For a detailed review of inward matrix products, the reader is referred to the section entitled Mathematical Aspects of Quantum Similarity. We can then use these new descriptors in so-called quantum QSAR.<sup>67-69</sup>

We have already used similarity matrices to cluster molecules. As such, they provide the necessary data to investigate the construction of a molecular set taxonomy. The most common techniques to do so include the molecular point clouds previously described. There, the columns of the molecular quantum similarity matrix yielded coordinates of the molecules in the  $N$ -dimensional space. Often, the  $N$  dimensionality cannot yet be used for a graphical representation. However, several techniques exist to reduce the dimensionality of the data, which allow it to be represented graphically on common devices like a computer screen or a plotter.<sup>70</sup> In addition to these plots, several instances have involved Kruskal trees and other algorithms.<sup>9,62-66,70</sup>

A second technique for clustering, used on several occasions,<sup>59,70,71</sup> consists of the construction of tree graphs or dendrograms. There, we manipulate the similarity matrices in such a way as to show graphically the proximity relations between different molecules. Bultinck and Carbó-Dorca<sup>71</sup> have shown how such dendrograms can be derived. All elements of a MQSMM  $\mathbf{Z}$  can then be transformed into a number lying in the interval (0,1) through the generalized cosine index or Carbó index as in Eq. [18], which means that the diagonal elements of the matrix  $\mathbf{C}$ , obtained from the diagonal elements of  $\mathbf{Z}$ , are transformed into one, whereas all other elements will lie within the interval (0,1). The extent of dissimilarity between molecules  $A$  and  $B$  is reflected in the extent of the deviation from 1 of  $C_{AB}$ . The construction of, for instance, a sequential agglomerative hierarchical nonoverlapping (SAHN) dendrogram proceeds by searching the sequences of largest off-diagonal elements of  $\mathbf{C} = \{C_{AB}\}$ . However, once such a first step is taken, and two quantum objects have been gathered into a new cluster, the matrix  $\mathbf{Z}$  in its current state is no longer applicable as it is. An exact definition of a quantum object will be given later in the section on the mathematical aspects of quantum similarity. For the present discussion, it suffices to look at a quantum object as a labeled object with a density function attached to it. Consider that, in a first step, molecules  $A$  and  $B$  have been gathered, and then all matrix elements  $Z_{KA}$  and  $Z_{KB}$  ( $K = 1, N$ ) have lost their ordering meaning. To obviate such a problem, a new object is introduced. This new object is, in fact, the cluster consisting of elements  $A$  and  $B$ . It is not a physical object like a molecule but instead corresponds to an artificial object. We can, however, use it as an object in the same way as a physical object, because we can easily obtain its representative vector  $\mathbf{z}_X$ . The elements in this column can be obtained from the MQSMM  $\mathbf{Z}^0$ , where the superscript<sup>o</sup> denotes the original startup similarity

matrix. For each of the  $N$  discrete descriptors in the column vector for the new quantum object, the average is taken of the two quantum objects composing this new quantum object. Denoting  $A$  and  $B$  as the quantum objects gathered in the new quantum object  $X$ , one has

$$\forall K \in \{1, \dots, N\} : Z_{KX} = \frac{Z_{KA}^0 + Z_{KB}^0}{2} \quad [36]$$

Note that for the consistency of the  $\mathbf{Z}$  and  $\mathbf{C}$  matrix definitions, we stress that averages are made in the elements of  $\mathbf{Z}$ , and *not* in the elements of  $\mathbf{C}$ , which corresponds to an unweighted pair-group average method. With this averaging method, the matrix  $\mathbf{Z}$  remains symmetric throughout the averaging steps, because averages are convex combinations of the  $\mathbf{Z}^0$  matrix columns and as such the result can always be a new quantum object.

Once the first cluster, holding elements  $A$  and  $B$ , has been constructed, the new matrix  $\mathbf{Z}$  is constructed. The process of constructing new  $\mathbf{Z}$  matrices and searching for the biggest off-diagonal element of  $\mathbf{C}$  may be repeated unchanged at every stage in the construction of the dendrogram. Because every new agglomeration into a new averaged quantum object involves two quantum objects  $A$  and  $B$ , we could opt to calculate an average of the elements of the column vectors of both quantum objects  $A$  and  $B$ , which would, however, continuously reduce the weight of the first two associated objects in the subsequential clustering steps, which is clearly not desired. Therefore, after every stage in the dendrogram construction, all elements of the  $\mathbf{Z}$  matrix are reconstructed as

$$Z_{XY} = \frac{1}{N_X} \frac{1}{N_Y} \sum_{I \in X} \sum_{J \in Y}^{N_X} Z_{IJ}^0 \quad [37]$$

where  $N_X$  is the number of parent quantum objects present in the cluster  $X$  and  $N_Y$  is the number of parent objects present in cluster  $Y$ . The so-called parent quantum objects are hereby defined as those individual quantum objects that form the matrix  $\mathbf{Z}^0$ .

To illustrate these techniques, we consider a set of structural isomers. Such isomers have the same stoichiometric composition, but they differ in their bonding pattern. The set of molecules in this example is shown in Figure 1.

A typical dendrogram obtained from overlap MQSM using atomic shell approximation electron densities is shown in Figure 2.

Such dendrograms can reveal interesting additional information about the similarity not only between two molecules in pair-wise comparisons, but also between clusters of molecules.

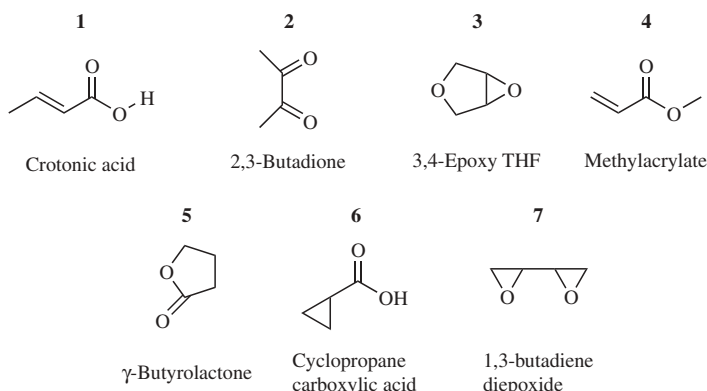


Figure 1 Isomers of  $C_4H_6O_2$  applied in dendrogram construction.

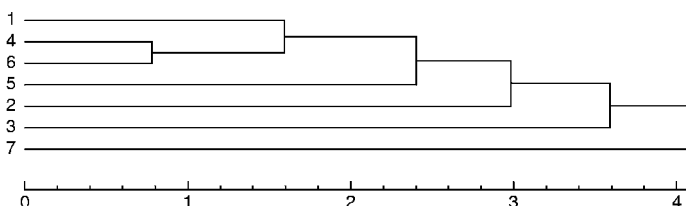


Figure 2 Dendrogram obtained for the isomers of  $C_4H_6O_2$  of Figure 1.

## ELECTRON DENSITIES FOR MOLECULAR QUANTUM SIMILARITY

Many techniques to obtain electron densities exist. The most common techniques rely on the wave function, from which the electron density is easily obtainable. Methods of different complexity for obtaining the wave function exist, ranging from the simple Hartree–Fock method up to those like Møller–Plesset perturbation, coupled-cluster, and configuration interaction. A detailed account of these methods may be found in books like that of Helgaker et al.<sup>44</sup> The range of applicability of these methods in medicinal chemistry was discussed by Barden and Schaefer.<sup>72</sup> A different path for obtaining the electron density is taken in DFT, where the electron density is used directly as the all-determining entity. A clear introduction to DFT may be found in Parr and Yang,<sup>48</sup> Koch and Holthausen,<sup>73</sup> or Ayers and Yang.<sup>74</sup>

Clearly, the electron density will depend on the method that obtains it. Nonetheless, even the simple Hartree–Fock (HF) electron density performs well to obtain first-order electron densities. The electron density in this MO LCAO framework is given by

$$\rho(\mathbf{r}) = \sum_v \sum_\mu D_{v\mu} \chi_v^*(\mathbf{r}) \chi_\mu(\mathbf{r}) \quad [38]$$

Here, this usual notation may be found in Szabo and Ostlund.<sup>38</sup> The  $\chi_v(\mathbf{r})$  and  $\chi_\mu(\mathbf{r})$  are the basis functions and  $D_{v\mu}$  is the density matrix, also called the charge density-bond order matrix. It follows that the MQSM in this framework is given by

$$Z_{AB} = \sum_{v \in A} \sum_{\mu \in A} \sum_{\delta \in B} \sum_{\lambda \in B} D_{v\mu} D_{\delta\lambda} \iint \chi_v^*(\mathbf{r}_1) \chi_\mu(\mathbf{r}_1) \Omega(\mathbf{r}_1, \mathbf{r}_2) \chi_\delta^*(\mathbf{r}_2) \chi_\lambda(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [39]$$

It is immediately clear that evaluating the MQSM is not such an easy computational task. Even though the use of contractions of Gaussian-type primitives is well-established and good algorithms are available for the evaluation of the integrals in Eq. [39], the number of integrals makes calculating the MQSM time-consuming, especially if we use many higher angular momentum basis functions. Furthermore, if we work with ab initio densities, we should also take into account the time needed to perform the ab initio calculations. On the other hand, this approach gives the ab initio MQSM, which can be considered as a reference value. Programs to calculate similarity measures according to Eq. [39] have been implemented for different operators, e.g., with the similarity routines of the BRABO ab initio program.<sup>75,76</sup>

Naturally, great interest exists in methods to enhance the speed of the evaluation of the MQSM, especially when large molecules or large sets of molecules are involved. It would be beneficial to avoid time-consuming ab initio calculations and to skip the evaluation of computationally demanding integrations as appear in Eq. [39]. A common technique consists of replacing the difficult expression of the electron density by a simpler expression, thus yielding simpler or computationally less time-consuming integrals. One such possibility involves s-type Gaussians only, which is an especially economical way of handling the integrals when molecular superposition is performed (see the section on molecular alignment) because in this case, the MQSM have to be calculated many times.

It is worth locating where the electron density is highest in a molecule. It is found that the electron density values are highest in the volumes defined by the subvalence regions of every atom. It is also well known that in those regions, chemical bonding produces only small density changes compared with the isolated atoms, which suggests the possibility of composing an approximate molecular electron density as a sum of isolated atom densities. Admittedly, molecule formation will change the electron densities from the summed-up atomic electron densities, but we may expect that this will have only a relatively small effect on the MQSM values. The idea of an approximate electron density composed of superposing isolated atom electron densities organized in the same geometrical way as in the true molecule has a long history. The so-called promolecular electron density was introduced in 1977 and still plays an important role in the Hirshfeld, X-ray based, idea of population analysis.<sup>77</sup> Carbó-Dorca et al. have derived the atomic shell

approximation to obtain approximate electron densities that correspond to promolecular densities, but with some special characteristics, as will be shown below. The Hirshfeld and ASA approaches both expand the promolecular electron density as a sum of atomic densities. For a molecule A with  $N$  atoms  $\alpha$ , the promolecular electron density  $\rho_A^{PROM}(\mathbf{r})$  is given as

$$\rho_A^{PROM}(\mathbf{r}) = \sum_{\alpha=1}^N \rho_{\alpha}^0(\mathbf{r}) \quad [40]$$

Clearly, the Hirshfeld promolecular electron density is not likely to simplify the integrals in Eq. [39]. The essential difference between the Hirshfeld and ASA promolecular densities is that in the Hirshfeld method, the isolated atom electron densities  $\{\rho_{\alpha}^0(\mathbf{r})\}$  are obtained in the same basis set as the one in the ab initio calculation of the true molecular electron density, whereas in the ASA approach, the isolated atom densities are obtained in the way as described below. In the ASA method, we use a slightly different promolecular atomic shell approximation (PASA) electron density, where the number of electrons  $P_{\alpha}$  attached to each atom  $\alpha$  is introduced. The total promolecular electron density for an  $N$ -atom molecule is given by

$$\rho_A^{PROM}(\mathbf{r}) = \sum_{\alpha=1}^N \rho_{\alpha}^0(\mathbf{r}) = \sum_{\alpha=1}^N P_{\alpha} \rho_{\alpha}^{ASA}(\mathbf{r}) \quad [41]$$

In most applications, this number  $P_{\alpha}$  is set equal to the atomic number of the element involved,  $Z_{\alpha}$ . In some applications,  $P_{\alpha}$  is set to the number of electrons attributed to the atom according to the results of some type of population analysis. The isolated atom electron densities are optimal linear combinations of s-type Gaussians. That is,

$$\rho_{\alpha}^{ASA}(\mathbf{r}) = \sum_{i=1}^M w_i |s_i(\mathbf{r})|^2 \quad [42]$$

The  $w_i$  are the expansion coefficients for the  $M$  s-type Gaussians, and we can see immediately the link between Eq. [42] and the wave function quadrature. So, for the calculation of ASA-based promolecular electron densities, we first need to develop a scheme for the fitting of the atomic densities. The exponents of the Gaussians may be chosen from, e.g., a well-tempered series.<sup>36</sup> The coefficients may then be fitted against the true atomic ab initio electron density. Once these exponents and coefficients are set, these Gaussian exponents and coefficients are universally applicable. Promolecular densities  $\rho_A^{PROM}(\mathbf{r})$  can then be obtained quickly from Eq. [41].

Fitting electron densities is by no means an unexplored terrain and has been studied in many ways for many different uses. Providing a complete overview of this subject is well beyond the scope of this chapter, but examples include fitting electron densities under constraints such as conserving the electric field or the electrostatic potential, among others.<sup>78–81</sup> The most related application of density fitting, however, goes back to the 1970s where it was introduced mainly for DFT, more precisely for the calculation of Coulomb terms.<sup>82–84</sup> In its first applications in molecular quantum similarity, ASA applied similar techniques. The idea of this ASA implementation is to fit an s-type GTO expansion to the atomic ab initio density, obtained through, e.g., a Hartree–Fock calculation. The fitting is done in a least squares sense by minimizing

$$\Delta^2 = \int [\rho_\alpha^{AI}(\mathbf{r}) - \rho_\alpha^{ASA}(\mathbf{r})]^2 d\mathbf{r} \quad [43]$$

where  $\rho_\alpha^{AI}(\mathbf{r})$  denotes the ab initio atomic electron density. Working out the square in Eq. [43], one derives the following expression:

$$\Delta^2 = \int [\rho_\alpha^{AI}(\mathbf{r})]^2 d\mathbf{r} - 2 \int \rho_\alpha^{AI}(\mathbf{r}) \rho_\alpha^{ASA}(\mathbf{r}) d\mathbf{r} + \int [\rho_\alpha^{ASA}(\mathbf{r})]^2 d\mathbf{r} \quad [44]$$

Using the symbol  $Z_{\alpha\alpha}^{AI}$  to denote the ab initio self-similarity of atom  $\alpha$ , and introducing Eq. [42], we find:

$$\Delta^2 = Z_{\alpha\alpha}^{AI} - 2 \sum_{i=1}^M w_i \int \rho_\alpha^{AI}(\mathbf{r}) |s_i(\mathbf{r})|^2 d\mathbf{r} + \sum_{i=1}^M \sum_{j=1}^M w_i w_j \int |s_i(\mathbf{r})|^2 |s_j(\mathbf{r})|^2 d\mathbf{r} \quad [45]$$

Naturally, we must take into account the normalization constraint:

$$\int \rho_\alpha^{ASA}(\mathbf{r}) d\mathbf{r} = 1 \quad [46]$$

This constraint automatically makes any promolecular density for some molecule, calculated with Eq. [41], fulfill the normalization constraint for that molecule. Such a constraint is most easily handled by a Lagrange multiplier. In ASA, normalized s-type Gaussians are applied, so the following constraint is introduced:

$$\sum_{i=1}^M w_i = 1 \quad [47]$$



After some straightforward manipulations, one finds the following linear equation:

$$\mathbf{w} = \mathbf{S}^{-1}(\mathbf{t} + \lambda |1\rangle) \quad [48]$$

and for the Lagrange multiplier

$$\lambda = \frac{1 - \langle 1 | \mathbf{S}^{-1} \mathbf{t} \rangle}{\langle 1 | \mathbf{S}^{-1} | 1 \rangle} \quad [49]$$

where the following matrix  $\mathbf{S} = \{S_{ij}\}$  and vector  $\mathbf{t} = \{t_i\}$  have been introduced:

$$S_{ij} = \int |s_i(\mathbf{r})|^2 |s_j(\mathbf{r})|^2 d\mathbf{r} \quad [50]$$

$$t_i = \int \rho_\alpha^{AI}(\mathbf{r}) |s_i(\mathbf{r})|^2 d\mathbf{r} \quad [51]$$

In these equations,  $|1\rangle$  stands for a unity vector, a vector containing only elements 1. These results are the common expressions found in DFT<sup>73</sup> and in early applications of molecular quantum similarity.<sup>85,86</sup>

It has been argued by Constans and Carbó-Dorca that important, dubious aspects of this linear fitting exist.<sup>87</sup> Because the electron density is always a positive definite quantity, the expansion coefficients should be positive definite. If not, no absolute guarantee exists that the electron density will be positive throughout the entire space. Therefore, we should also add as a constraint to the minimization that

$$\forall i : w_i \in \mathbf{R}^+ \wedge w_i > 0 \quad [52]$$

Imposing these constraints may reduce the goodness-of-fit between the ab initio and ASA densities. The elimination of a possible negative electron density outweighs this smaller error by far, however. Imposing the constraints in Eq. [52] on the Lagrangian and finding a way to minimize this function may be done in different ways. In a first approach, Constans and Carbó-Dorca used a positive restricted fitting technique.<sup>87</sup> In other implementations,<sup>88-93</sup> we can use a fairly simple idea for assuring that the constraints of Eq. [52] are obeyed. This idea consists of considering all coefficients as being the squares of a real number:

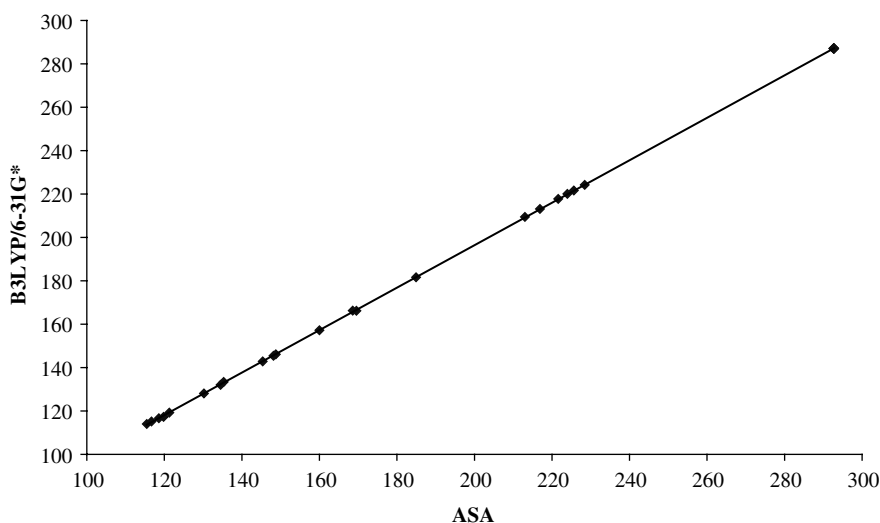
$$\forall i : w_i = x_i^2, x_i \in \mathbf{R} \wedge \sum_{i=1}^M x_i^2 = 1 \quad [53]$$

The optimization is then no longer a linear problem and becomes a difficult optimization with the usual caveats such as the existence of local minima. Amat et al.<sup>88–93</sup> used the elementary Jacobi rotation technique as an optimization technique to find the proper coefficients for the ASA s-type Gaussians. Other techniques can be applied as well, such as the Newton–Raphson technique.<sup>94</sup> Lists of s-type Gaussian coefficients and exponents may be found on the Internet for different basis sets.<sup>95</sup>

The ASA method should be critically tested in its application to molecular quantum similarity, and this has been done systematically for the promolecular electron densities obtained via Eq. [41]. Gironés et al. have found these PASA densities to be valuable for large molecules when comparing isodensity surfaces. Based on a small set of molecules, they found the ASA MQSM to be close to the *ab initio* MQSM.<sup>96,97</sup> Furthermore, Amat et al. also found that the squared errors of Eq. [43] are small.<sup>89–93</sup> On the other hand, it was also found that the calculation of quantum chemical expectation values with the promolecular density does not give good results, e.g., for the electrostatic potential.<sup>91</sup> This result is not unexpected because we only used an s-type basis set in ASA. Amat et al. also used promolecular densities to obtain starting vectors for the SCF iterations in *ab initio* calculations. Only minor improvement over other classic methods for obtaining such starting vectors<sup>93</sup> was found for molecules with first and second period atoms only. For molecules with heavier atoms, this approach performs better, and it provides a general and uniform procedure applicable to any molecule. The fact that expectation values from ASA promolecular densities are not very good, and that the quality of the starting vectors for SCF calculations does not always improve stand in great contrast with the reliable ASA and PASA densities in molecular quantum similarity studies. A clear illustration of the utility of ASA and PASA techniques may be found with the isomers shown in Figure 1. DFT-B3LYP/6-31G\* calculations yielded MQSM that were compared with ASA-based MQSM. It was found that the ASA overlap MQSM agree well with the *ab initio* MQSM, which is shown graphically in Figure 3.

The agreement is clearly superb, which is important because ASA promolecular electron densities are, in fact, just a sum of atomic densities, and we might have feared that for isomeric molecules, ASA would not be able to discriminate between different MQSM. Figure 3 shows that this fear is unfounded, and that ASA reproduces the B3LYP/6-31G\* results well. The same is true for larger sets of less congeneric molecules.<sup>98</sup> It has recently been found that PASA also give very good Coulomb MQSM.<sup>99</sup>

Two reasons exist for the good performance of ASA and PASA densities in studying molecular quantum similarity. First, because the highest electron density is concentrated near the nucleus, it is easy to understand that if we use functions that give a fair or good representation of that area, the contribution from that area of the molecule in the total similarity will be approximated well. Second, these subvalence regions are usually reasonably transferable



**Figure 3** Agreement between ASA and B3LYP/6-31G\* overlap MQSM for the structural isomers in Figure 1. *Figure reprinted with permission from Bultinck et al.*<sup>98</sup> Copyright 2003 American Chemical Society.

from one molecule to another. These two reasons are the basis why the PASA quantum similarity measures agree so well with ab initio measures. These two points are in accordance with the observation that the improvement of the correlation is relatively small when we fit the ab initio density of the entire molecule to obtain molecule-specific ASA coefficients.<sup>100</sup> Nevertheless, as was found by Amat and Carbó-Dorca<sup>91</sup> and Van Damme,<sup>100</sup> molecular fitting of the ASA coefficients always improves the agreement between the ab initio self-similarities and the ones obtained from fitted ASA densities.

When doing the molecular fitting of ASA coefficients compared with promolecular ASA, the ASA coefficients no longer sum to 1 for all atoms in the molecule, because the sum of the coefficients reflects the effect of polarization. In other words, rather than having  $P_\alpha = Z_\alpha$  in promolecular ASA, the values of  $P_\alpha \neq Z_\alpha$  in fitted ASA reflect the internal electron density redistribution. In such a way, fitted ASA coefficients are reminiscent of the “Stewart atoms” idea, where atoms in molecules are considered as radially distorted atoms and  $P_\alpha$  are connected with so-called Stewart charges.<sup>101</sup> Contrary to common promolecular ASA where the atomic densities of Eq. [41] integrate to  $Z_\alpha$ :

$$\int \rho_\alpha^0(\mathbf{r}) d\mathbf{r} = Z_\alpha \quad [54]$$

in molecular ab initio density fitted ASA (ASA-Fit), the atomic densities integrate to a different value, which reflects the molecular polarization:

$$\int \rho_{\alpha}^0(\mathbf{r})d\mathbf{r} = P_{\alpha} \quad [55]$$

The global molecular normalization naturally remains valid. For a neutral molecule A with  $N$  electrons, this means that

$$\int \rho_A^{PASA}(\mathbf{r})d\mathbf{r} = N = \int \rho_A^{ASA-Fit}(\mathbf{r})d\mathbf{r} \quad [56]$$

The interest in ASA-Fit techniques stems from those cases where high accuracy is wanted or for charged molecules. For some atoms in charged molecules, ASA-Fit densities give a better picture. A concern is whether it is cheaper to simply calculate the ab initio MQSM rather than to introduce the fitting step. In this context, one should realize that during a superposition procedure, the required integrals are evaluated many times. This evaluation is much faster with the simple ASA s-type Gaussians than when with a large general basis set. Thus, we first perform an ab initio calculation and then fit the ASA coefficients to this density. After that, we use the ASA method with the fitted coefficients in molecular superposition, thus saving computer time.

A possible criticism against the total electron density, whether it is based on the true ab initio density or the core density-containing ASA method previously described, is that the emphasis is put on the chemically less interesting regions. In other words, we may question to what extent chemical information is lost because the MQSM is dominated by the core electron density. It is well known that much of chemistry is dictated by the fuzzy lower electron density regions involved in chemical bonding. As will be described, this topic is closely related to molecular superposition. To sketch the problem briefly, consider two molecules where in each, one heavier element exists (e.g., a P atom and a Si atom) and the remaining atoms are lighter atoms like C, O, and H. The largest value for the MQSM over both molecules will be obtained when the Si and P atoms are superposed, which will happen even if the other atoms are not aligned. In contrast, chemists would say that from their knowledge, the alignment of the C, O, and H atoms might be more relevant in a chemical context. These chemists probably base their opinions on their perception of the chemical bonds in the molecules along with related features such as lone electron pairs. Both the bonds and possibly lone electron pairs are valence density effects and do not contribute to a large extent in the MQSM when we are using total electron densities.

Several attempts have been made to alleviate the problems of the atom cores dominating the MQSM by emphasizing the role of the chemically more interesting outer electron density. Because no physical ground exists to

designate some space volumes as containing core density and some as containing valence density, the idea should be considered a qualitative reasoning to emphasize the problems encountered. Nonetheless, several suggestions have been made to alleviate the overwhelming influence of the core density.

A first route lies in the use of pseudopotential methods. Pseudopotentials, also known as effective core potentials (ECPs), are techniques that help reduce the computational cost of *ab initio* calculations<sup>103</sup> by replacing the core electrons with a fixed potential. In that way, no effort needs to be spent on determining the orbitals for the subvalence electrons. The idea of chemists using pseudopotentials is based on the fact that the electron density nearest to the nucleus of the atoms does not undergo major changes upon molecular formation or changes in conformation. Accordingly, a transferable potential can be applied, which mimics the effects of this subvalence electron density on the outer electrons.<sup>103</sup> When *ab initio* calculations are performed with pseudopotentials, the electron density obtained is a valence density, which could be strictly the valence density as understood in a chemical context, or it may extend up to one subvalence shell depending on the pseudopotential applied.<sup>104</sup> We can then use the resulting valence densities in the evaluation of the MQSM in Eq. [21]. Another way to obtain valence electron densities is to apply the densities of only those orbitals that can be characterized as valence orbitals, which is straightforward in an MO-LCAO approach.

Although we have addressed up to now mostly the electron density, several studies have reported the use of some derived entity in lieu of the electron density. In fact, such a case has already been introduced in the so-called kinetic MQSM of Eq. [27], where we use the gradient of the electron density. This process is also true for those situations in which extended density functions are used in which several new density functions were derived from extended Hilbert space.<sup>56</sup>

Yet another series of MQSM can be derived from the field of study called conceptual density functional theory. In this field, many concepts are obtained that may be written as derivatives of the electron density. For a recent overview of this field, one may consult the review by Geerlings et al.<sup>37</sup> or the classic textbook by Parr and Yang.<sup>48</sup> One example of such a derivative is the Fukui function, defined as

$$f(\mathbf{r}) = \left( \frac{\partial \rho(\mathbf{r})}{\partial N} \right)_{V_{\text{ext}}} \quad [57]$$

The Fukui function at point  $\mathbf{r}$  is the derivative of the electron density over the total number of electrons in the molecule. Another density derivative-based property is the local softness, which is defined as the density derivative with respect to the chemical potential:

$$s(\mathbf{r}) = \left( \frac{\partial \rho(\mathbf{r})}{\partial \mu} \right)_{V_{\text{ext}}} \quad [58]$$

Both expressions assume a constant external potential, which usually means that the geometry of the molecule must be left unchanged. Boon et al.<sup>105–107</sup> have used different conceptual DFT properties to calculate MQSM. The basic argument for these quantities versus the electron density is that quantities such as the local softness are more related to chemical reactivity than the electron density. The main conceptual DFT quantity used by these authors has been the local softness and the atom condensed local softness. The molecular quantum similarity indices for pairs of molecules differ, depending on the molecular descriptor.<sup>105</sup> In a different paper, Boon et al.<sup>106</sup> examined autocorrelation functions to remove molecular alignment problems associated with obtaining the MQSM. This alignment-free algorithm will be discussed in the section on molecular alignment.

Still another idea is introduced in the contributions by Cooper and Allan.<sup>108</sup> They removed the core electron density dominance problem by using momentum transformations. A similar expression exists for the electron density in momentum space:

$$\rho(\mathbf{p}) = \Psi^*(\mathbf{p})\Psi(\mathbf{p}) \quad [59]$$

The wave functions in momentum space are obtained from those in position space by the usual momentum space Fourier transformation (see introductory quantum mechanics books such as that by Bransden and Joachain<sup>109</sup>). Obtaining the wave functions in momentum space is not always straightforward computationally, although numerical techniques have been developed to do this,<sup>110</sup> and analytical forms have been derived for specific functions.<sup>111</sup> It is worth describing why the momentum space density approach is appealing for molecular quantum similarity. When we consider the nature of density in momentum space, and the decay of the basis functions in momentum space, clearly the largest contributions come from the slow-moving electrons. These electrons naturally are moving furthest from the core region of the atoms. The electrons nearest to the core move at great speeds, approaching in heavy elements a substantial fraction of the speed of light. Electrons with such large momenta contribute little to the momentum space density, whereas the slow ones contribute substantially more. The latter electrons are the outer electrons. The momentum space density thus emphasizes the chemically most interesting parts of the molecule, such as where chemical bonding takes place, without relying on the chemical topology. To summarize, position space electron density is concentrated in the core regions, and momentum space density is concentrated in the valence region.

Cooper and Allan<sup>108</sup> have used momentum density in several studies. A problem remains in obtaining the momentum space densities because most calculations are performed with position space wave functions. In a sense, working in momentum space is yet another way to reduce the overweighting of the core electron density. Most of the following discussions on, e.g., molecular alignment and quantum similarity indices, remain valid when we

use momentum space densities, replacing where necessary the position space concepts by momentum space concepts. Some concerns about the proper convergence of the involved similarity integrals in momentum space remain.

As discussed earlier, we cannot only derive first-order electron densities, but also we can extend them to higher order electron densities. We have used the second-order electron density  $\rho(\mathbf{r}_1, \mathbf{r}_2)$  in lieu of the first-order electron density on several occasions in molecular quantum similarity, because the second-order electron density is in fact the lowest order density where electron correlation becomes apparent. It has been used extensively by Ponec et al.<sup>112–118</sup> in the study of similarity in pericyclic reactions where the second-order electron density offers important advantages over the first-order electron density. In another contribution, Ponec et al. went to the third-order electron density.<sup>119</sup> Again, most of the discussion relating to molecular quantum similarity indices and molecular alignment is also applicable to higher order electron densities, replacing where necessary the first-order electron density by, for example, the second-order electron density.

In a somewhat different area of similarity research, chemists use not the electron density, but rather the electrostatic potential. It is somewhat different in the sense that the entire electrostatic potential also contains a contribution not related directly to the electron density and is therefore not positive definite. The total electrostatic potential contains two contributions. One originates from the nuclei, and one from the electron density. For a system with  $M$  nuclei with nuclear charges  $\{Z_\alpha\}$  positioned at locations  $\{\mathbf{R}_\alpha\}$ , we obtain:

$$V^{pot}(\mathbf{r}') = \sum_{\alpha=1}^M \frac{Z_\alpha}{|\mathbf{r}' - \mathbf{R}_\alpha|} - \int \frac{\rho(\mathbf{r})}{|\mathbf{r}' - \mathbf{r}|} d\mathbf{r} \quad [60]$$

The electrostatic potential thus has as many discontinuities as the molecule has nuclei. When  $\mathbf{r}' = \mathbf{R}_\alpha$ , the electrostatic potential exhibits an infinite value creating large problems when an integral similarity measure of two electrostatic potentials is calculated as in

$$\int V_A^{pot}(\mathbf{r}_1) \Omega(\mathbf{r}_1, \mathbf{r}_2) V_B^{pot}(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [61]$$

A more direct link with molecular quantum similarity as introduced up to now occurs when we consider only the electronic contribution to the electrostatic potential. In this scenario, no discontinuities can develop. If the total electrostatic potential is considered instead, special techniques should remove the discontinuity problems. An example of such a technique was presented by Carbó-Dorca et al.,<sup>120</sup> where an exponential function served as a weight to avoid the divergence of the integrals of Eq. [61]. One problem that develops is that by removing the discontinuity problem, the numerical measures of the degree of quantum similarity are influenced. To compensate for this problem,

an alternative way to compare two electrostatic potentials has been proposed in the form of a similarity measure defined as:

$$Z_{AB}(r^{-2}) = \iint \frac{\rho_A(\mathbf{r}_1)\rho_B(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|^2} d\mathbf{r}_1 d\mathbf{r}_2 \quad [62]$$

Equation [62] can be taken as an approximation to the overlap of the electronic parts of the electrostatic potentials. Experience indicates that the behavior of the MQSMM with elements as in Eq. [62] is comparable with the simple density overlap as introduced already in Eq. [14]. In fact, a more general measure can thus be defined as

$$Z_{AB}(r^{-n}) = \iint \frac{\rho_A(\mathbf{r}_1)\rho_B(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|^n} d\mathbf{r}_1 d\mathbf{r}_2 \quad [63]$$

A more in-depth study of similarity indices will, however, be presented in the section on Quantum Molecular Similarity Indices. Modeling the potential of a charge distribution by an approximate charge distribution was previously studied in detail by Gill et al.<sup>121</sup>

Another application of molecular quantum similarity involves the work by Leherter et al.<sup>122–127</sup> Their low-resolution electron density maps are generated with different techniques such as ASA, wavelet transforms, or crystallography-based formalisms. Once these maps are available, the authors use critical points in the electron density to express both the degree of similarity and to perform molecular superposing. The similarity measure is no longer an integral measure, but it is turned into a root-mean-square measure with summations. Full details on this method may be found in Leherter et al.<sup>122–127</sup> An interesting feature of ASA densities, described by Leherter, et al. is that the critical points of the ASA densities agree well with those obtained from smoothed ab initio electron densities for both the values and the positions of these points in space.<sup>126–128</sup>

A special note should be made concerning the use of atom-in-molecule densities. The basic aspects of atom-density-based similarity continue to be valid. The reader is referred to the section on atoms-in-molecules similarity and chirality for an in-depth discussion of these domains of research.

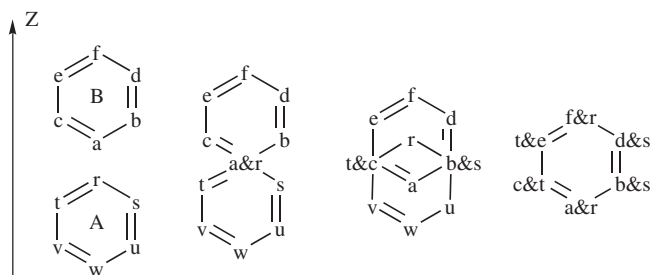
---

## THE ALIGNMENT ISSUE IN MOLECULAR QUANTUM SIMILARITY

### Statement of the Problem

If we consider the molecular quantum similarity measures in, e.g., Eq. [21], we should be aware of the important fact that the MQSM depends

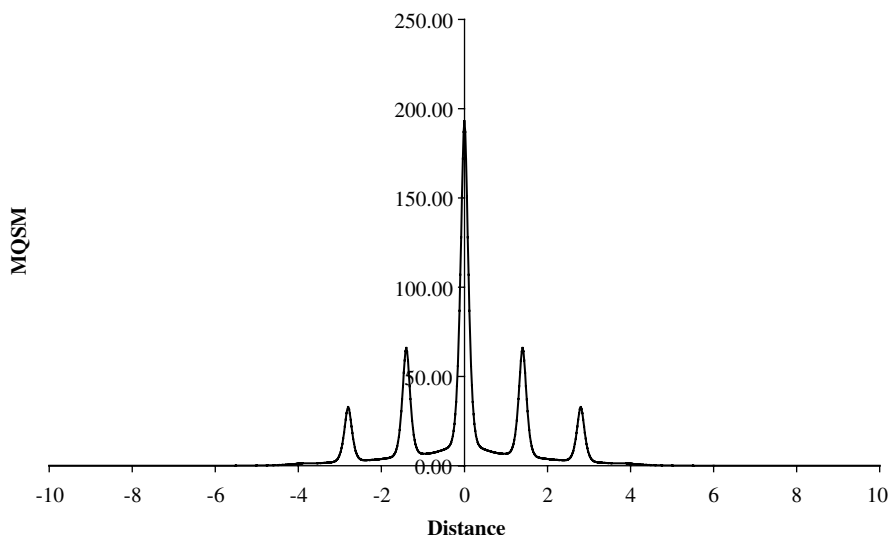




**Figure 4** Two benzene molecules A and B. A is moved along the Z-axis producing higher spatial overlap with B.

on the mutual position of the molecules A and B in space. To illustrate this dependence, consider two benzene molecules positioned as in Figure 4. Molecule A is moved along the Z-axis from a distance between the two centers of mass of  $-10 \text{ \AA}$  to  $+10 \text{ \AA}$ . At regular intervals, the overlap quantum similarity measure is calculated. The results are plotted in Figure 5.

Figure 5 clearly shows a dependence of the overlap MQSM on the mutual position of both molecules in space. Molecular alignment thus has important influence on the MQSM. One of the applications of molecular quantum similarity is the possibility to quantitatively express the similarity and especially to compare and order the degree of similarity over pairs of



**Figure 5** Overlap MQSM between benzene molecules A and B as a function of the distance (in  $\text{\AA}$ ) between the centers of mass of A and B.

molecules. As a consequence of the dependence of the MQSM on the mutual position of the two molecules involved in the MQSM, two molecules AB may be found to incorrectly have a larger similarity than a pair CD depending on the mutual position of A and B used in calculating  $Z_{AB}$  and the mutual position of C and D in evaluating  $Z_{CD}$ . Clearly, an efficient algorithm must be developed to avoid this problem and to ensure that  $Z_{AB}$  and  $Z_{CD}$  both correspond to a chemically meaningful similarity measure such that a meaningful similarity order is obtained.

Different important characteristics of MQSM may be discerned from Figure 5. Looking only at the carbon atoms in the two benzene molecules, the first maximum occurs when the first two atoms (a and r in Figure 4) coincide. Additional maxima occur when b and c coincide with s and t, and eventually the global maximum is found when all atoms of molecule A coincide with those of B. The other maxima in the example are obviously symmetry related to these three. The first important characteristic of MQSM is that several maxima are possible. The second important characteristic is that a high MQSM is found whenever two atoms coincide. Third, to illustrate the dependence on the core electron density of the MQSM, consider that the similarity of two coinciding hydrogen atoms yields a small MQSM, which means that any alignment where two heavy elements coincide will contribute more to the molecular MQSM, even if the topologically more similar molecular regions do not coincide at all. Table 1 shows atomic self-similarities calculated with ASA for some common elements.

The MQSM can be obtained roughly by summing the atomic self-similarities of the coinciding atoms. Table 1 shows that two molecules, each containing one sulfur atom and five carbon atoms, will give an optimal alignment when the two sulfur atoms coincide, even if this means that no other atoms are superimposed. As a rough measure,  $Z_{AB}$  will be 815.91 au when both molecules have superposed sulfur atoms, but no overlapping carbon atoms. An alternative alignment might have the five carbon atoms aligned and not the sulfur atom. The similarity will then be only 159.25 au, far from the maximum value.

**Table 1** ASA Overlap MQSSM for a Selection of Atoms

Atom	ASA Self-Similarity
H	0.0400
C	31.8500
N	52.7330
O	81.7261
F	120.5609
S	815.9067

## Quantum Similarity Maximization—MaxiSim and QSSA

From a mathematical point of view, we can propose maximization of the MQSM as a function of a set of molecular alignment parameters in such a way as to obtain a superposition procedure with a general scope. Several algorithms aimed directly at maximizing the MQSM have been published. The first example is the MaxiSim algorithm developed by Constans et al.<sup>129</sup> A second one is the quantum similarity superposition algorithm (QSSA) by Bultinck et al.<sup>130</sup> An algorithm developed by Stefanov and Cioslowski<sup>131,132</sup> is similar to the MaxiSim and QSSA ideas, and it will therefore not be described explicitly because it invokes similar points of view.

Whatever algorithm we use, we must ensure that the maximizing process does not get trapped into a local maximum. It is mainly there that the MaxiSim and QSSA algorithms differ. We should note that in both procedures, the optimization of the alignment involves repeated calculation of the MQSM. As described earlier, the evaluation of this similarity measure with general basis sets is demanding of computer time, and it is therefore an untenable option, even with the current generation of computers. Therefore, in both the MaxiSim and QSSA algorithms, the ASA method, previously described, is usually used by chemists. In the QSSA, however, we can also invoke calls to subroutines for the repeated calculation of the *ab initio* MQSM, even though ASA has been found to be appropriate for the calculation and maximization of the MQSM. This latter finding was established by calculating both the ASA and the *ab initio* MQSM at different regions on the optimization hypersurface formed by the alignment parameters.<sup>98</sup>

The molecular alignment problem could be dealt with by the many existing techniques available for, e.g., conformational analysis, where an extremum is sought as a function of several structural parameters. In MQS, the objective is to maximize the MQSM as a function of different alignment parameters. These parameters correspond to three translation parameters and three rotation parameters when no molecular flexibility is taken into account, i.e., for rigid bodies. The translations occur through the X, Y, and Z axes of a common coordinate system for molecules A and B. The rotations about the Cartesian axes can be described by Euler angles, although it is sometimes more efficient for chemists to use quaternions.<sup>54,55</sup> The alignment problem consists of finding a maximum for six alignment parameters. At first glance, this problem seems fairly simple but the presence of multiple maxima may be problematic for some systems. One approach is to consider a systematic search, which would seem feasible, because the well-known combinatorial explosion is not expected to be too large when only six parameters are present. However, experience shows that the grids in which the alignment parameters must be developed are so fine that this technique becomes impracticable for routine alignment, especially when large, nonlinear molecules are considered. An alternative is a random search procedure. The drawback of this method is

that the number of random points needed to ensure success is so overwhelmingly large that it too becomes an impracticable method.

To overcome this computational dilemma, Constans et al.<sup>129</sup> created an efficient algorithm that proceeds in different steps. Using the ASA method, the authors succeed in replacing the true integral measure of Eq. [21] by a summation, which was accomplished by introducing a function in the ASA density expression that effectively compresses the atomic shells. In the limit, the atomic electron densities are infinitely compacted into the nuclei. To illustrate the procedure, consider first expression [42]. A compression parameter  $t$  is introduced in the different shells such that

$$\rho_{\alpha}^{ASA}(\mathbf{r}, t) = \sum_{i=1}^M w_i |s_i(\mathbf{r}, t)|^2 \quad [64]$$

The molecular electron density then becomes in the limit of  $t$  approaching infinity:

$$\rho_A^{PASA}(\mathbf{r}) = \sum_{\alpha=1}^N P_{\alpha} \delta(\mathbf{r} - \mathbf{r}_{\alpha}) \quad [65]$$

where again  $\alpha$  refers to a specific atom, with its nucleus located at  $\mathbf{r}_{\alpha}$ . The efficiency of this approach is the consequence of the fact that the MQSM can be rewritten as a simple double summation:

$$\tilde{Z}_{AB} = \sum_{\alpha=1}^{N_A} \sum_{\beta=1}^{N_B} P_{\alpha} P_{\beta} \delta(\mathbf{r}_{\alpha} - \mathbf{r}_{\beta}) \quad [66]$$

$\alpha$  and  $\beta$  are, respectively, atoms of molecules A and B, holding, respectively,  $N_A$  and  $N_B$  atoms.

This deformed MQSM  $\tilde{Z}_{AB}$  was found to be a reliable and efficient way of scanning the true MQSM  $Z_{AB}$  as a function of the alignment parameters. The reader will note from the previous description that the MaxiSim algorithm seeks alignments that will superimpose atomic nuclei and will rely on finding alignments where interatomic distances in molecule A are equal to, or close to, interatomic distances in molecule B. Consider two sets of three atoms, one belonging to molecule A and a second set belonging to molecule B. The first step in the algorithm will superpose the first atom of A with the first atom of B by simple translation of B. In the following step, molecule B is oriented in such a way that the first atom pair of A coincides with that in B. A rotation is then performed such that the planes of the two sets of three atoms coincide. Finally, the MQSM is calculated. By cycling through the different possible collections of sets of three atoms in A and B, we can locate the alignment giving the highest MQSM. The premise of this alignment algorithm inherently is the idea that the MQSM is dominated by coinciding atomic core electron densities

in both molecules. More details and the exact implementation of the MaxiSim algorithm can be found in the publication by Constans et al.<sup>129</sup> The method has the important feature that the optimization is fast.

A different approach to the alignment problem was published by Bultinck et al.,<sup>130</sup> who developed the QSSA algorithm. This technique involves a genetic algorithm in which a local optimizer is interlaced, thereby making it a Lamarckian genetic algorithm. The starting point for this evolutionary procedure is to have an algorithm that can function for different electron densities, varying from ASA densities to *ab initio* total electron densities and valence-only densities. Whereas the MaxiSim algorithm will work efficiently if the maxima for the MQSM coincide with those alignments where a maximal amount of nuclei coincide, QSSA is more general in scope. MaxiSim in fact applies in an efficient manner the fact that the MQSM is dominated to a large extent by overlapping core electron densities. When this core electron density dominance disappears, MaxiSim can no longer be guaranteed to locate the maximum values for the MQSM and the global maximum in particular, unless if it can be found through the local optimization starting from a MaxiSim regular alignment. QSSA, in contrast, makes no assumption about the nature of the electron density and does not base the search on the assumption that the maxima will be located by superposing atoms of both molecules. When chemists use ASA electron densities, which are total electron densities, the QSSA algorithm usually yields the same global maximum as does MaxiSim, although in rare cases QSSA generates a slightly higher MQSM.

The QSSA algorithm is implemented in the following way. First the electron density is generated for both molecules A and B. We denote these  $\rho_A^X$  and  $\rho_B^X$ . The X refers to a specific type of density, which may be one of the many types previously described or they may be AIM densities that are described later. The most common densities are promolecular ASA densities (X = PASA) or *ab initio* total (X = TOT) and valence densities (X = VAL). Once these densities are available, a random mutual orientation in space is generated for both molecules. A common feature of most alignment algorithms in molecular quantum similarity is the use of three translations and the Euler angles, although quaternions have also been applied in QSSA. Denoting the coordinates of atoms b in molecule B in the coordinate system of B as  $\{x_b^0, y_b^0, z_b^0\}$ , the coordinates  $\{x_b, y_b, z_b\}$  of these atoms with respect to molecule A become

$$\begin{bmatrix} x_b \\ y_b \\ z_b \end{bmatrix} = \begin{bmatrix} T_X \\ T_Y \\ T_Z \end{bmatrix} + \begin{bmatrix} \cos\alpha\cos\beta\cos\gamma - \sin\alpha\sin\gamma & \sin\alpha\cos\beta\cos\gamma + \cos\alpha\sin\gamma & -\sin\beta\cos\gamma \\ -\cos\alpha\cos\beta\sin\gamma - \sin\alpha\cos\gamma & -\sin\alpha\cos\beta\sin\gamma + \cos\alpha\cos\gamma & \sin\beta\sin\gamma \\ \cos\alpha\sin\beta & \sin\alpha\sin\beta & \cos\beta \end{bmatrix} \begin{bmatrix} x_b^0 \\ y_b^0 \\ z_b^0 \end{bmatrix} \quad [67]$$

The set of parameters  $\{T_X, T_Y, T_Z, \alpha, \beta, \gamma\}$  are, respectively, the translations along the X, Y, and Z axes of molecule B with respect to molecule A and the Euler angles, which describe the rotation. Both molecules are considered rigid, and molecule A is kept on the same place in space. Molecule B is overlaid on molecule A with a random choice of these six parameters. The MQSM is evaluated at this random orientation and then maximized for the six translation and rotation parameters. Referring to Figure 5, we anticipate that the random alignment will be one in which the MQSM is far from maximal. Maximization can be done with different techniques, although the computationally most robust procedure was found to be the simplex method.<sup>133–134</sup> It is well known that the simplex algorithm is slow, but it does perform well, especially in the first steps of the maximization. After several steps, one can switch to a better algorithm such as a quasi-Newton–Raphson scheme like the BFGS<sup>133</sup> algorithm. When a maximum is located that was not found previously, the set of parameters corresponding to this maximum are recorded for future reference. After this run, a new attempt is initiated to locate another maximum, or possibly a previously located one. As previously mentioned, pure random searching is not a computationally interesting option, and the QSSA algorithm implements a genetic algorithm<sup>135</sup> to refine and accelerate the search. Genetic algorithms are powerful techniques to find a global extremum when numerous local extrema exist. Genetic algorithms work with populations consisting of a certain number of members. Each member holds a set of alignment parameters, which will be called the *chromosomes*, along with an MQSM value. The members of the population are ranked according to their *fitness*, where in the current application, fitness refers to the magnitude of the MQSM. A higher *fitness* infers a higher MQSM. Using ideas loosely based on genetics, the most fit members of a population have a higher chance of producing even fitter offspring than do less fit members. Therefore, we use a selection scheme for selecting parent members of a population biased toward the fitter members. The parent members then produce offspring by exchanging their *chromosomes*, thereby forming a new generation in the population. The key feature of a genetic algorithm lies in knowing that usually, with every new generation, the average fitness increases along with the fitness of the fittest member, unless the global minimum was already located. Techniques such as *crossovers*, *mutation*, and *elitism* further improve the genetic algorithm. More details of the genetic algorithm used in QSSA may be found in the literature.<sup>130,135</sup> In most applications of genetic algorithms, no local optimizer is being used. That is, when a new offspring is created, the fitness is calculated and no local optimization is performed. In so-called Lamarckian genetic algorithms, a local optimizer is included in the algorithm, performing a local maximization of the offspring fitness, before they in turn generate the next generation. In QSSA, the simplex or quasi-Newton–Raphson techniques are used as local optimizers.

Maximizing the Carbó index has been described in several previous scientific reports. McMahon and King<sup>136</sup> described the use of gradient methods in 1997, and in the same year, Parretti et al.<sup>137</sup> described the use of Monte Carlo techniques. In many studies, including the two just mentioned, these maximizations do not refer to quantum similarity, but instead they refer to maximizing the similarity in molecular electrostatic potentials, which is different.

One of the interesting features of QSSA is that no reference is made to coinciding atoms nor to locating similar interatomic distances in molecules A and B. As such, when valence densities are used, we can use the genetic algorithm to locate the globally maximal MQSM, even if this means that the heaviest atoms do not coincide at all or if only a few atoms of A coincide with atoms of molecule B. A special remark should be made about the alignment issue in the context of similarity between atoms in molecules. As will be detailed in a later section, this process involves the first exact alignment of the atoms under consideration. Naturally, then, the translational parameters are immediately determined, and the alignment reduces to a maximization problem in only three variables, which is a much easier problem.

## Structural Alignment

It is beyond the scope of this chapter to discuss the range of structure-based methods that chemists can use for molecular alignment. This field of research has been, and continues to be, very active. One algorithm, called TGSA,<sup>138</sup> will be presented here in some detail, however, because of its popularity in molecular quantum similarity studies. Structure-based techniques differ from the aforementioned techniques in several respects. First, they not attempt to maximize the MQSM for a pair of molecules. Second, they do not make a specific reference to molecular quantum similarity; as such, they are aimed at a wider range of applications. Third, they are not based on electron density in a formal way, but instead they take a more familiar approach based on chemical topology. Consequently, they apply well-known concepts such as chemical bonds and try to overlap the most similar and largest common structure elements in both molecules.

The first question we could ask is whether it is meaningful for us to use any other alignment than the one that gives the globally maximal MQSM value. As depicted in Figure 5, one of several values could be found depending on how we implement the topo-geometrical algorithm. In the MaxiSim and QSSA algorithms, we use total electron densities that are aimed at producing the maximal value of the MQSM, which comes down to aligning the heaviest atoms in both molecules. The example given at the beginning of this section, with the data from Table 1, is instructive in this regard. The overlap of only

the two heavy sulfur atoms gives the highest MQSM, even though overlapping more atoms, except the sulfur atoms, might be more meaningful chemically. Chemistry is mainly based on many changes involving the valence region where chemical bonding occurs. A chemist would then naturally be inclined to align the molecules in such a way that the common substructure coincides, even if this means that in our example the sulfur atoms will not be aligned. The valence electron density is well reflected in the bonding topology and structural representation of a molecule, and as such, a topo-geometrical alignment could be attempted to superimpose the molecule in such a way as to give the chemically more relevant alignment. Although somewhat exaggerated, we could say that rather than maximizing the MQSM, the chemical relevance of the alignment has been maximized this way. At this chemically relevant alignment, the MQSM could then be evaluated for the sake of comparison. Even though the MQSM might be far from its maximal value, it is expected to contain a higher chemical information content. To compare and contrast these different alignment philosophies, some aspects of the topo-geometrical method by Gironés et al.<sup>138</sup> will be described, followed by a discussion of its relation to QSSA.

The topo-geometrical superposition approach (TGSA) by Gironés, Robert and Carbó-Dorca is based on the recognition of the largest common substructure in the aligned molecules. Common substructure searching is an active field of research; for a lucid review of the different aspects in substructure searching, the reader is referred to the work by Chen.<sup>23</sup>

For TGSA, all that is required for both molecules A and B is the list of atoms with their atomic number and Cartesian coordinates. In both molecules, all chemical bonds are defined and the interatomic distances are stored. Hydrogen atoms are not considered in these calculations. The chemical bonds between the heavy atoms are referred to as dyads. Each dyad of the first molecule is then compared with each dyad of the second molecule. The two dyads are considered similar if the dyad distances in both molecules are the same or smaller than a given threshold. The dyads meeting this requirement then serve as a starting point to construct a database of atomic triads. A triad is made by attaching to one of the dyads a third atom, connected to at least one of the dyad atoms. All triads obtained for the first molecule are then compared with those of the second molecule in a similar way as the dyad comparison. From the pair of common triads, a set of translation and rotation parameters is derived that superpose these two triads. These translations and rotations are then applied to all atoms of molecule B. The different triad combinations are ranked according to the number of atoms that have been superimposed. If several triads exist that give the same number of superposed atoms, a scoring function assesses the quality of the alignment. The pair of triads yielding the best score is the optimal topo-geometrical alignment. If at first no common triad can be identified, the procedure is repeated with higher thresholds for the identification of similar dyads or triads.



The TGSA algorithm has recently been extended to include conformational degrees of freedom.<sup>139</sup> TGSA performs well in aligning molecules topo-geometrically. The alignments produced by TGSA agree better with common chemical intuition than with those produced by MaxiSim or QSSA based on total electron densities. The main limitation of TGSA is that to work efficiently, it requires at least some structural similarity between the molecules involved. Topo-geometrical techniques are naturally not limited to TGSA, and many algorithms have been described. For an overview, the reader is referred to the review by Lemmen and Lengauer.<sup>140</sup>

## Comparison of Alignment Techniques

Until now two different classes of alignment techniques have been described, each with its own perspective of what is important for alignment. On the one hand are the MaxiSim and QSSA methods that usually rely on total electron density, and on the other hand are the chemically more intuitive topological techniques like TGSA. Compared with QSSA and MaxiSim, the TGSA philosophy lacks a strong connection to molecular quantum similarity. However, it has the advantage that chemically more reasonable alignments are produced. It is therefore interesting to see if a connection between the simple topological-based methods can be made to a quantum chemistry-based technique, like the QSSA algorithm. When we use QSSA with valence electron densities, interestingly enough, an alignment is produced that is comparable with the TGSA alignment, which puts topo-geometrical alignments on a firmer basis within the framework of molecular quantum similarity.<sup>141</sup> Thus, the TGSA alignment corresponds qualitatively to the maximal MQSM when we use valence electron densities. The QSSA algorithm can thus yield both types of alignment by simply varying the type of electron density. QSSA has, however, a still wider scope because it can also be applied to perform alignments based on the HOMO density, Fukui functions, or in fact any partial density function.

Another matter is the consistency of the molecular quantum similarity matrix  $\mathbf{Z}$ . The MQSM produced by a specific alignment technique for a given molecular pair of the set of molecules that construct  $\mathbf{Z}$  should not be contradictory with the computed MQSM for the other pairs of molecules. To illustrate this point, consider the Euclidean distance, as defined by the square root of Eq. [13]:

$$d_{AB} = \sqrt{Z_{AA} + Z_{BB} - 2Z_{AB}} \quad [68]$$

So, if an MQSM  $Z_{AB}$  has been obtained with one of the alignment techniques, a Euclidean distance can be calculated. Naturally, the fundamental

requirements for a distance metric, including the triangular inequality, have to hold. In other words,

$$d_{AB} \leq d_{AC} + d_{CB} \quad [69]$$

A scheme to check whether all MQSMs over a set of molecules are internally consistent can be derived. If inconsistencies exist, the alignment should be checked carefully. Maximizing the MQSM with programs such as MaxiSim and QSSA present a typical case in which different local maxima occur. Many of these local maxima are likely to infer inconsistencies in  $\mathbf{Z}$ . We may then use this situation to signal cases in which the global maximum has not yet been located. Such consistency checking has been used by Bultinck et al.<sup>130</sup> to establish whether the genetic algorithm works. The issue of consistency checking has also been described in detail by Bultinck, Carbó-Dorca and Van Alsenoy.<sup>98</sup> It was found that MaxiSim matrices are usually internally consistent. TGSA-derived similarity matrices sometimes contain some, usually minor, inconsistencies. Consistency is always present with QSSA because it is an inherent part of the algorithm.

As a final note, on several occasions, alignment-free methods have been used to quantify molecular similarity; in the field of molecular quantum similarity, these methods have not yet found extensive application. One method to obtain molecular quantum similarity measures without the need for molecular alignment was published by Boon et al.<sup>106</sup> They use statistical techniques, more specifically, the autocorrelation function. This technique offers an interesting alternative method for similarity studies by removing completely the important obstacle of molecular alignment.

---

## QUANTUM SIMILARITY INDICES

Until now the molecular quantum similarity measure has primarily been the integral measure  $Z_{AB}$ . The direct comparison of two values  $Z_{AB}$  and  $Z_{CD}$  does not directly yield an idea of the degree of similarity. Consequently, a numerical transformation must be established, which allows the comparison of the similarity degree between different pairs of molecules.

It is almost impossible to enumerate and discuss all of the similarity indices that have been published. Furthermore, new indices continue to be published. Therefore, this discussion will be limited to those indices that have been most common in the application in molecular quantum similarity. For a general review of similarity indices and their application, the reader is referred to the work by Willett, Barnard and Downs.<sup>142</sup>

The most common index in molecular quantum similarity is the generalized cosine, introduced by Carbó et al.<sup>53</sup> in their first paper on quantum

similarity, published in 1980. The index has also become known as the Carbó index, and it is given by

$$C_{AB} = \frac{\int \rho_A(\mathbf{r})\rho_B(\mathbf{r})d\mathbf{r}}{\sqrt{\int \rho_A(\mathbf{r})\rho_A(\mathbf{r})d\mathbf{r} \int \rho_B(\mathbf{r})\rho_B(\mathbf{r})d\mathbf{r}}} \quad [70]$$

or using the elements of  $\mathbf{Z}$ :

$$C_{AB} = \frac{Z_{AB}}{\sqrt{Z_{AA}Z_{BB}}} \quad [71]$$

This is an example of the so-called C-class descriptors, which give a value in the interval  $[0,1]$ , where a higher value signifies a greater degree of similarity. One of the nice features of these indices is that they can be extended beyond the overlap similarity given in Eq. [70]. So, with any operator  $\Omega$ , a similarity measure can be defined as

$$C_{AB} = \frac{Z_{AB}(\Omega)}{\sqrt{Z_{AA}(\Omega)Z_{BB}(\Omega)}} \quad [72]$$

In the Carbó index, the denominator is the geometric mean of the self-similarities. Naturally, we could also think about using the arithmetic mean, which gives rise to the Hodgkin–Richards index<sup>143,144</sup> originally developed for comparisons of electrostatic potentials:

$$H_{AB} = \frac{2Z_{AB}(\Omega)}{Z_{AA}(\Omega) + Z_{BB}(\Omega)} \quad [73]$$

This index, however, has no direct geometrical interpretation, as has the Carbó index.

A third index that is common is the so-called Petke index,<sup>145</sup> given by

$$P_{AB} = \frac{Z_{AB}(\Omega)}{\max(Z_{AA}(\Omega), Z_{BB}(\Omega))} \quad [74]$$

Finally we introduce the Tanimoto index,<sup>146</sup> given by

$$T_{AB} = \frac{Z_{AB}(\Omega)}{Z_{AA}(\Omega) + Z_{BB}(\Omega) - Z_{AB}} \quad [75]$$

These C-class indices vary from 0 to 1, 1 denoting perfect similarity, which does not infer equality between the molecules involved. Maggiora, Petke

and Mestres<sup>147</sup> have investigated the relations between different indices and found that the following inequalities hold:

$$0 \leq P_{AB} \leq H_{AB} \leq C_{AB} \leq 1 \quad [76]$$

These inequalities agree with what is found in practice when applying the different indices for pairs of molecules. It has been shown by Carbó et al.<sup>148,149</sup> that many of these indices can be generalized in some way. For example, the Hodgkin–Richards and Tanimoto indices belong to the same class of indices. A general formula, or so-called Girona index, is given by

$$G_{AB} = (k - x) \frac{Z_{AB}}{\left[\frac{k}{2}(Z_{AA} + Z_{BB}) - xZ_{AB}\right]} \quad [77]$$

The Hodgkin–Richards index is recovered when  $k = 2$  and  $x = 0$ , whereas the choice  $k = 2$  and  $x = 1$  gives the Tanimoto index. Robert and Carbó-Dorca<sup>150</sup> presented an in-depth analysis of the relations between the different C-class indices using dendrograms and statistical analysis. Their aim was to find the extent of redundancies between different indices. Within the C-class indices, they found two subclasses, the first consisting of the Carbó, Hodgkin–Richards, and Petke indices and the second containing the Tanimoto index. It has also been shown that the Hodgkin–Richards index can be related to the Carbó index.<sup>148,149</sup>

It is interesting to note that another, generalized C index can be derived:

$$C_{AB}^{(n)} = \frac{Z_{AB}^{(n)}(\Omega)}{\sqrt{Z_{AA}^{(n)}(\Omega)Z_{BB}^{(n)}(\Omega)}} \quad [78]$$

where we have in general

$$Z_{IJ}^{(n)}(\Omega) = \iint \rho_I^{n/2}(\mathbf{r}_1)\Omega(\mathbf{r}_1 - \mathbf{r}_2)\rho_J^{n/2}(\mathbf{r}_2)d\mathbf{r}_1d\mathbf{r}_2 \quad [79]$$

in which we use the positive nature of  $\rho(\mathbf{r})$ . It is clear that the Carbó index is recovered when  $n = 2$ . It is interesting to note that when  $n$  equals 1 and we use the Dirac delta operator  $\delta$ , the following equations are found:

$$C_{AB}^{(1)} = \frac{Z_{AB}^{(1)}(\delta)}{\sqrt{Z_{AA}^{(1)}(\delta)Z_{BB}^{(1)}(\delta)}} = \frac{\int \rho_A^{1/2}(\mathbf{r})\rho_B^{1/2}(\mathbf{r})d\mathbf{r}}{\sqrt{\int \rho_A(\mathbf{r})d\mathbf{r} \int \rho_B(\mathbf{r})d\mathbf{r}}} = \frac{\int \rho_A^{1/2}(\mathbf{r})\rho_B^{1/2}(\mathbf{r})d\mathbf{r}}{\sqrt{N_A N_B}} \quad [80]$$

$$C_{AB}^{(1)} = \int \sigma_A^{1/2}(\mathbf{r})\sigma_B^{1/2}(\mathbf{r})d\mathbf{r} \quad [81]$$

which means that the first-order Carbó index is in fact the overlap of the square roots of the shape functions of the two molecules involved. Also, given the nature of shape functions,  $C_{AA}^{(1)} = 1$ . Shape functions will be mentioned in more detail later in this chapter. More information on these functions may be found in the paper by Parr and Bortolotti,<sup>151</sup> and the review by Geerlings, De Proft and Langenaeker.<sup>37</sup> The connection of shape functions with similarity was recently described by Bultinck et al. in the context of vector semispaces and molecular quantum similarity.<sup>152</sup> The scalar products and norms involved in Eqs. [78]–[81] are well founded within the vector semispace metric structure, as will be shown in the section on the mathematical foundations of quantum similarity.

A different class of indices are those that behave like a distance. With these indices, a perfect similarity corresponds to a value 0, and no strict upper limit exists. The best-known of these D-class indices is the Euclidean distance introduced previously in Eqs. [7]–[13]. This Euclidean distance index has been extended to give an index for any operator used by chemists in evaluating the MQSMs by using

$$D_{AB}(\Omega) = Z_{AA}(\Omega) + Z_{BB}(\Omega) - 2Z_{AB}(\Omega) \quad [82]$$

and has been generalized to include two parameters  $k$  and  $x$  to give:

$$D_{AB}(\Omega, k, x) = \frac{k}{2} [Z_{AA}(\Omega) + Z_{BB}(\Omega)] - xZ_{AB}(\Omega) \quad [83]$$

This result yields the Euclidean distance when  $k = x = 2$ .

Many more indices can be derived than those described here. Moreover, we can imagine many ways to obtain a D-class index by a mathematical manipulation of a C-class index and vice versa. Several examples of such transformations have been published by Carbó-Dorca et al.<sup>148,149</sup>

---

## QUANTUM ATOMS-IN-MOLECULES SIMILARITY

---

A field that has attracted special attention on different occasions is that of the similarity between two atoms located in two different molecules. For example, we can imagine calculating the similarity between two carbonyl carbon atoms, one in molecule A and one in molecule B, or calculating the similarity between a carbon atom in a molecule and the isolated atom, or between two different carbon atoms in the same molecule. From the perspective of molecular quantum similarity, calculating the similarity measure will require atomic electron densities within the molecule. Therefore, before turning to the similarity calculation, two methods for obtaining such densities will be briefly presented.

The methods differ in a basic aspect. The first method, used by Boon et al.<sup>105–107</sup> involves the Hirshfeld concept, previously described to some extent and in more detail later. The second method involves the atoms-in-molecules (AIM) theory developed by Bader.<sup>153</sup> The Hirshfeld method involves atomic densities that are distributed all over space, reminiscent of molecular electron densities that are also spread all over space. The atomic densities spread out to infinity. The AIM method involves atomic basins that are finite, or often in molecules semifinite, extending to infinity in one radial direction. A third method based on restricted summations over basis functions will also be presented.

### The Hirshfeld Approach

The Hirshfeld idea, already developed in 1977,<sup>77</sup> calculates the so-called stockholder charges and is a popular method in conceptual DFT.<sup>37</sup> It consists of the following rationale. First an electron density, represented as  $\rho_A^{ai}$ , is obtained for a molecule A with some Hamiltonian and basis set. For every atom  $\alpha$ , an isolated electron density  $\rho_\alpha^0$  is calculated within the same model. With the isolated atom electron densities for all N atoms comprising the molecule, a Hirshfeld promolecular density is obtained as

$$\rho_A^{PROM}(\mathbf{r}) = \sum_{i=1}^N \rho_i^0(\mathbf{r}) \quad [84]$$

The idea of Hirshfeld was to obtain atoms-in-molecules densities by defining the “stock-amount” or weight of an atom  $\alpha$  in the electron density at  $\mathbf{r}$  as

$$w_\alpha^A(\mathbf{r}) = \frac{\rho_\alpha^0(\mathbf{r})}{\sum_{i=1}^N \rho_i^0(\mathbf{r})} = \frac{\rho_\alpha^0(\mathbf{r})}{\rho_A^{PROM}(\mathbf{r})} \quad [85]$$

If it is furthermore assumed that these weight coefficients remain valid for the true, ab initio density of the molecule. The Hirshfeld atomic electron density of the atom  $\alpha$  in the molecule, denoted  $\rho_\alpha^{Hirsh}(\mathbf{r})$ , can be calculated as

$$\rho_\alpha^{Hirsh}(\mathbf{r}) = w_\alpha^A(\mathbf{r}) \rho_A^{ai}(\mathbf{r}) \quad [86]$$

Application of Eq. [86] in the general definition of a molecular quantum similarity measure then gives:

$$Z_{\alpha\beta}(\Omega) = \iint w_\alpha^A(\mathbf{r}_1) \rho_A^{ai}(\mathbf{r}_1) \Omega(\mathbf{r}_1, \mathbf{r}_2) w_\beta^B(\mathbf{r}_2) \rho_B^{ai}(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [87]$$

Although this is the most obvious expression for defining the similarity between two atomic densities in the Hirshfeld way, a different expression was used in the study of molecular chirality.<sup>107</sup> The latter approach will be discussed in some more detail in the application of molecular similarity to study molecular chirality.

## AIM-Based Methods

The AIM theory by Bader yields so-called atom basins in a molecule. For a complete and thorough discussion of this theory, the reader is referred to the book by Bader.<sup>153</sup> To make clear the application of this theory in quantum similarity, some key elements will be presented briefly. Bader's theory starts from the well-known postulates of quantum mechanics and adds just one additional postulate, which allows for discerning atoms in molecules. Regions in Cartesian space, bordered by surfaces of zero-flux in the field of the electron density gradient, are called attractor basins. They exhibit all properties of quantum chemical systems. Each of these regions is spanned by gradient paths terminating at a single point, which is called an attractor. When an attractor coincides with a nucleus, the basin and the nucleus together form the atom in the molecule. Atomic basins are denoted as  $\Omega_\alpha^A$ , which means the basin of atom  $\alpha$  in molecule A.

The atoms-in-molecules theory was used by Cioslowski et al.<sup>131,132</sup> for studies in quantum similarity. To obtain a measure of the similarity between an atom  $\alpha$  and an atom  $\beta$  in molecules A and B, respectively, Cioslowski and Nanayakkara<sup>154</sup> developed the following index:

$$C_{\alpha\beta} = \left[ \frac{\int_{\Omega_{\alpha\beta}} \rho_A(\mathbf{r}) d\mathbf{r}}{\int_{\Omega_\alpha^A} \rho_A(\mathbf{r}) d\mathbf{r}} \right] \left[ \frac{\int_{\Omega_{\beta\alpha}} \rho_B(\mathbf{r}) d\mathbf{r}}{\int_{\Omega_\beta^B} \rho_B(\mathbf{r}) d\mathbf{r}} \right] \quad [88]$$

where the notation  $\Omega_{\alpha\beta}$  means the intersection of the atomic basins  $\Omega_\alpha^A$  and  $\Omega_\beta^B$  and  $\Omega_{\alpha\beta} = \Omega_{\beta\alpha}$ .

Examples of the application of this theory may be found in the work of Cioslowski et al.<sup>131,132</sup> It was found that when applying the similarity index described here, all oxygen atoms in a series of aldehydes, show similarities between 98.91% and 99.89%.

Popelier et al.<sup>155–159</sup> proposed a new similarity measure operating in an abstract space spanned by properties evaluated at bond critical points defined by the theory of AIM. This measure gives rise to a field known as quantum topological molecular similarity, and is used, among others, to develop QSP/AR models such as toxicity predictors.

### Atom-Centered Basis Function Approach

A simple technique that could assess the degree of similarity between atoms in molecules can be based on the technique used by Mezey et al.<sup>160</sup> in their discussion on molecular chirality. There the aim was to investigate the degree of similarity between enantiomers, which can naturally be done with the total molecular electron density. Mezey et al.<sup>38</sup> provided an interesting way to calculate the similarity between fragments in both enantiomers. It is based on the well-known fact that the electron density in LCAO-MO theory may be written as

$$\rho(\mathbf{r}) = \sum_{\nu\mu} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) \quad [89]$$

where  $\nu$  and  $\mu$  refer to the basis functions applied in the calculations. Mezey et al.<sup>160</sup> used a restriction on the summation to derive fragment densities. We can obtain the electron density for a fragment  $F$ , for example, by limiting the summation over basis functions centered on the atoms of this fragment only:

$$\rho_F(\mathbf{r}) = \sum_{\nu\mu \in F} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) \quad [90]$$

If the electron density for only one atom is desired, we could reduce the fragment to only a single atom  $\alpha$ , which would give

$$\rho_{\alpha}(\mathbf{r}) = \sum_{\nu\mu \in \alpha} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) \quad [91]$$

An expression for the quantum similarity measure can then be obtained easily.

This method has been applied to molecular fragments containing three atoms as well as to the case of a single atom by Mezey et al.<sup>160</sup> A criticism of this approach is that the fragment or atom densities do not sum to the total electron density of an  $M$ -atom molecule, as is shown in Eq. [92]:

$$\sum_{\alpha}^M \rho_{\alpha}(\mathbf{r}) = \sum_{\alpha}^M \sum_{\nu\mu \in \alpha} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) \neq \sum_{\nu\mu} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) = \rho(\mathbf{r}) \quad [92]$$

This inequality is caused by the absence of the off-diagonal blocks

$$\sum_{\substack{\nu \in \alpha \\ \mu \notin \alpha}} P_{\nu\mu} \chi_{\nu}^*(\mathbf{r}) \chi_{\mu}(\mathbf{r}) \quad [93]$$

and in further research, these terms should be properly taken into account.



In all three of these approaches, the similarity indices should be maximized. Both the Carbó index and the Hodgkin–Richards index were used by Boon et al.<sup>107</sup>, whereas in the applications of AIM, the Cioslowski similarity measure was used.<sup>131,132,154</sup> For the atom-centered approach, we can use all similarity indices. In all cases, maximal similarity is obtained if the two atoms being compared are coincident. Looking back at the QSSA algorithm, and more specifically at Eq. [67], this means that the translation serves to make the atoms under consideration coincident. The maximization then comes down to finding the Euler angles that maximize the molecular similarity. These optimizations are not always analytically straightforward. AIM expressions for gradient and Hessian elements were given by Cioslowski et al.,<sup>131,132</sup> in the Hirshfeld case, numerical integration and maximization is needed.

---

## PHYSICAL CONNOTATIONS OF (SELF) SIMILARITY MEASURES

---

What physical meaning can be attached to the molecular quantum similarity indices calculated with some positive definite operator? No direct indication exists that any meaning should be attached in general. For the self-similarity measures given as

$$Z_{AA}(\Omega) = \iint \rho_A(\mathbf{r}_1) \Omega(\mathbf{r}_1, \mathbf{r}_2) \rho_A(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad [94]$$

such a connection has been established, however. It will be shown that the self-similarity measures are related to physical observables, such as  $\log P$ .  $\log P$  is one of the most successful molecular descriptors in quantitative structure-activity relationships, and it is a simple measure of the hydrophobic/lipophilic character of a molecule. It corresponds to the octanol/water partition coefficient. The utility of  $\log P$  as a potent molecular descriptor for QSAR in drug design has as its basis the fact that drug-ligands interaction are often strongly related to the hydrophobicity of the ligands. The importance of  $\log P$  is also evident from the fact that it is one of the descriptors contained in the so-called Lipinski rules.<sup>21</sup> These rules allow us to assess the drug-likeness of molecules, with only a few descriptors, which are mostly atom based (see earlier in this chapter). Values for  $\log P$  have already been determined for many molecules, but on many occasions, the values are unknown. Even in the latter case, an approximate value can usually be obtained with additive group contributions.<sup>162–168</sup> For a comparison of different methods to obtain these approximate values from group contributions, the reader may consult the test by Moriguchi et al.<sup>169</sup> Turning to the discussion of the relationship between  $\log P$  and the self-similarity measures, consider molecule X dissolved in water and the same molecule, but dissolved in octanol. Denoting this

molecule as X, and with the superscripts “w” and “o” for water and octanol, respectively, one finds as MQSM:

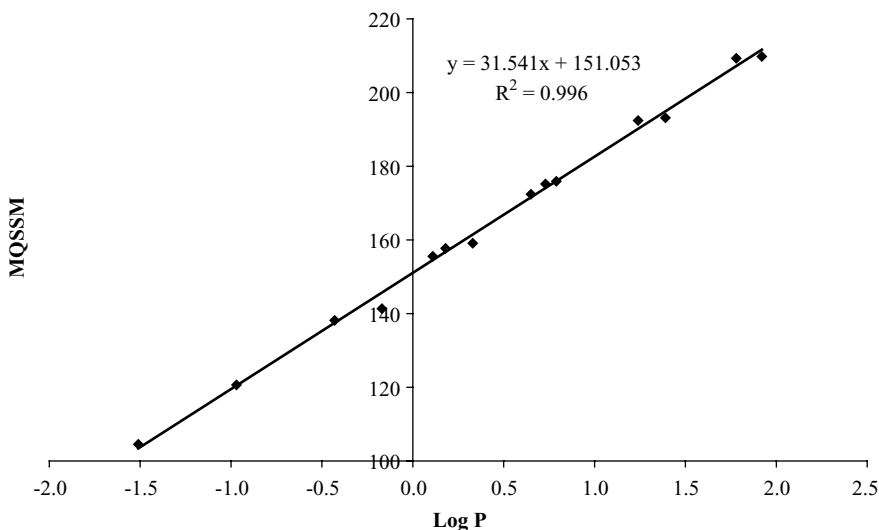
$$Z_{XX}^{wo} = \int \rho_X^w(\mathbf{r}) \rho_X^o(\mathbf{r}) d\mathbf{r} \quad [95]$$

Substituting Eq. [41] in Eq. [95], one obtains

$$Z_{XX}^{wo} = \sum_{\alpha}^{N_A} \sum_{\alpha'}^{N_A} P_{\alpha}^w P_{\alpha'}^o \int \rho_{\alpha}^w(\mathbf{r}) \rho_{\alpha'}^o(\mathbf{r}) d\mathbf{r} \quad [96]$$

Amat et al.<sup>170–174</sup> have investigated the correlation between these MQSM and log P, with ASA electron densities. As population characteristics  $P_{\alpha}^w$  and  $P_{\alpha}^o$ , they used the population numbers obtained from ab initio calculations including solvent models. Although ASA means that  $\rho_{\alpha}^w(\mathbf{r}) = \rho_{\alpha}^o(\mathbf{r}) = \rho_{\alpha}^{ASA}(\mathbf{r})$ , it was found by these authors that correlations obtained for log P versus the MQSM are good, as is shown in Figure 6.

It was later found that good correlations are also obtained with gas phase only data, therefore, negating the extra work associated with doing separate calculations for the same molecule immersed in the two different solvents. When we use ab initio calculated gas phase electron populations as



**Figure 6** Correlation between observed Log P values and self-similarity measures for a series of acetic acid esters, carboxylic acids, and amides. Reprinted by permission of John Wiley and Sons, Inc. from Amat et al.<sup>170</sup>

$P_\alpha$  in Eq. [41], correlations above 99% could again be obtained, which means that for predictive purposes,  $\log P$  for molecule A can be substituted by the MQSSM. Thus, the Lipinski rules can be modified taking into account easily and universally computed self-similarities.

In another study, Ponec et al.<sup>172</sup> looked for correlations between the self-similarity and the Hammett sigma constant  $\sigma$ , which plays an important role in defining the relationships between biological activity and molecular structure (QSAR). It was shown that molecular fragment quantum self-similarity measures can describe the substituent effect in a reaction series. The Hammett sigma constant was found to be related to the self-similarity measure for a well-chosen fragment in the reacting molecule. In fact, relationships between the molecular self-similarity measure and the Hammett parameter also exist, but more in the spirit of Hammett's idea, it is worthwhile to seek relationships with molecular fragments. The underlying idea is to choose the fragment where the reactive center for a given reaction is found, and to link the self-similarity measure for this fragment to the observed value for the  $\sigma$  constant. Ponec et al. used a series of substituted benzoic acids and investigated the possible relationship between the self-similarity measure of the  $-\text{COOH}$  fragment and the Hammett  $\sigma$  constants. Over a series of five groups of benzoic acid derivatives, each containing seven to eight molecules, a correlation could be established with a regression coefficient over 96%.<sup>174</sup> Similar results were obtained by Carbó-Dorca et al.<sup>10</sup> for yet another series of 12 benzoic acid derivatives. In all cases, the fragment self-similarity was calculated for a fragment  $F$ , containing the atoms  $\alpha$  as

$$Z_{XX}^{wo} = \sum_{\alpha \in F} \sum_{\alpha' \in F}^{N_A} P_\alpha P_{\alpha'} \int \rho_\alpha(\mathbf{r}) \rho_{\alpha'}(\mathbf{r}) d\mathbf{r} \quad [97]$$

Gross atomic populations were again used by the authors for the  $P_\alpha$ , as obtained from quantum chemical calculations.

In a recent contribution by Gironès et al.,<sup>175</sup> MQSSM were calculated from domain-averaged Fermi holes. The correlation between Hammett sigma constants and these MQSSM values for a set of meta- and para-substituted benzene molecules was investigated, and a good linear correlation between both entities was found. The predictive power of the relationship was examined for a test set of molecules, and it was found that the predictive capacity is competitive with the approach described by Sullivan et al.,<sup>176</sup> but having the advantage that only one parameter is used instead of five as in the latter study.

In a similar context, fragment self-similarities have been used to derive QSAR models and identified active molecular sites in molecules.

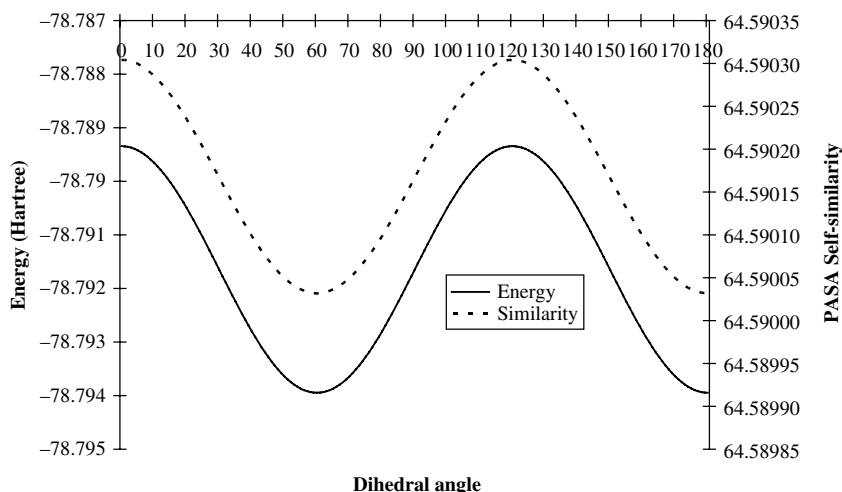
It was found by Solà et al.<sup>177</sup> that a good linear relationship exists between the overlap self-similarity measure and the molar volume for isoelectronic molecules when the atoms of the systems being compared belong to the same row of the periodic table. In the same context, it was shown that the

overlap self-similarity measure is a good descriptor for estimating the electronic charge concentration in molecules.

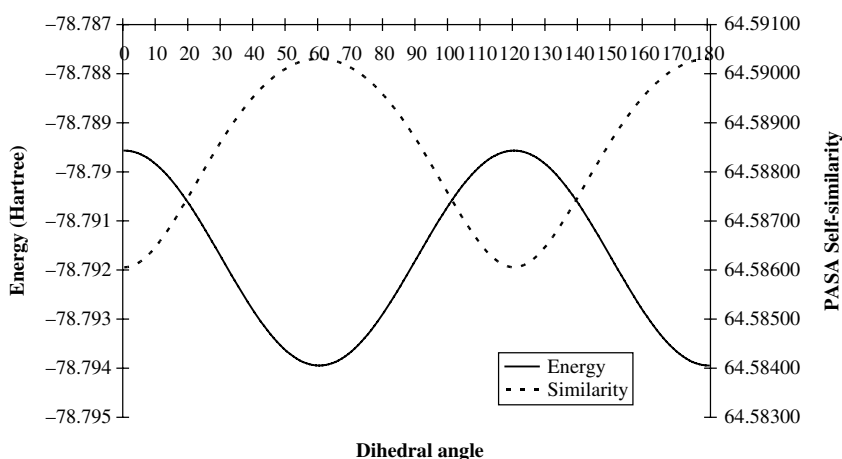
In an interesting study by Oliva et al.<sup>178</sup> the charge density redistribution during torsional rotations was examined for a series of alkanes. The supposition that electron density will redistribute with changes in molecular geometry is obvious, and the point of that investigation was to determine how these changes behave. To that end, ethane and propane served as probe molecules, and two approaches were followed. In an adiabatic approach they used the structural parameters of the relaxed staggered geometry of ethane and propane and they changed the dihedral angles without relaxation of the other structural parameters. In a non-adiabatic approach, all structural parameters were allowed to relax except the fixed dihedral angles. In the eclipsed ethane structure in the non-adiabatic approach versus the adiabatic one, we observe larger C—C bond lengths, and in propane also wider C—C—C angles, which is expected to cause a lowering of the charge density concentration. This result is indeed observed for both molecules; when going from the adiabatic staggered structure to the relaxed, non-adiabatic staggered conformation, the self-similarity measure decreases. From an energetic perspective, the formation of repulsive interactions is seen when going from staggered to eclipsed geometries. The evolution of the molecular quantum self-similarity measures depends on the type of approach. If we use a non-adiabatic approach, the relaxation of the nuclei causes a depletion of the electron density distribution, thereby lowering the self-similarity measure. Note that this result is consistent with the findings by Solà et al.,<sup>177</sup> who state that the self-similarity is a measure of the charge density concentration. If, however, we use an adiabatic approach, a reversed trend is observed where the total electron density gets more concentrated, which gives rise to higher self-similarity measures. The two trends are thus found to work in opposite directions. Figure 7 shows as an example the *ab initio* self-similarity for ethane as a function of the dihedral angle for both the adiabatic and the nonadiabatic approach.

Note that point-wise calculations of the staggered and eclipsed structures confirm these trends with higher levels of calculation of the electron density. It was argued by Oliva et al.<sup>178</sup> that one can perform fast conformational analysis from the self-similarities by qualitatively locating regions where energetic local minima could be found. This entails the assumption that approximate electron densities exhibit the same trend, which enables us to calculate self-similarity measures in a fast way without need for time-consuming *ab initio* calculations. When using ASA, we can split the contributions to the total overlap self-similarity measure into diagonal terms and off-diagonal terms in the following way:

$$Z_{AA} = \sum_{\alpha}^{N_A} P_{\alpha}^2 \int \rho_{\alpha}^2(\mathbf{r}) d\mathbf{r} + 2 \sum_{\alpha}^{N_A} \sum_{\beta > \alpha}^{N_A} P_{\alpha} P_{\beta} \int \rho_{\alpha}(\mathbf{r}) \rho_{\beta}(\mathbf{r}) d\mathbf{r} \quad [98]$$



a. Adiabatic Bond Rotation



b. Nonadiabatic Bond Rotation

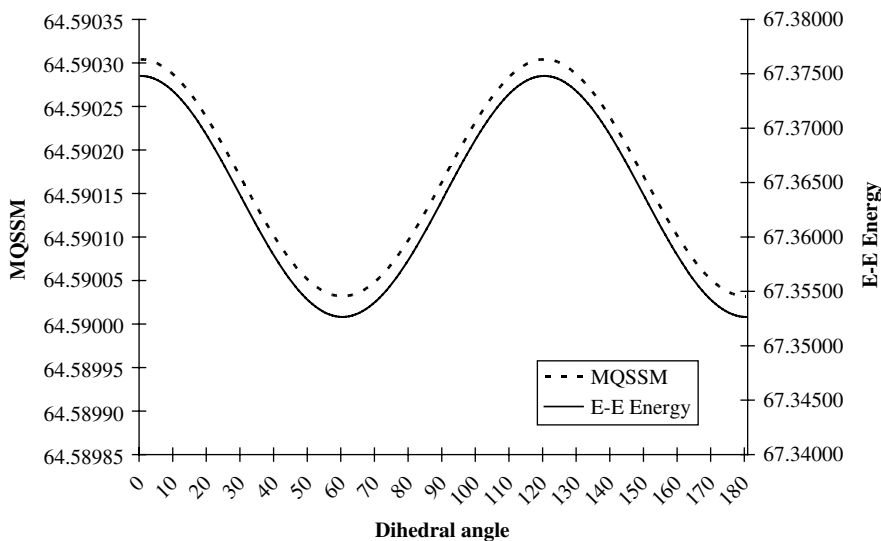
**Figure 7** Evolution of the ab initio (HF/3-21G) energy (in Hartree) and self-similarity measure versus dihedral angle in ethane in the adiabatic (a) and the non-adiabatic (b) approach.

The first term is constant; the second term was found to account for the observed trends. Oliva et al.<sup>178</sup> pursued this study further to examine the applicability of this method for conformational analysis of larger molecules. The approach described was found to be successful for larger molecules,

although several concerns were raised concerning the quality of the electron densities.

Although the trends of the total energy and the PASA MQSSM are shown on the same graphical size, we must be aware that in the adiabatic approach, the units on the MQSSM scale are substantially smaller than in the non-adiabatic approach, which illustrates the shape of the curve in the adiabatic approach. If in the adiabatic case, the MQSSM was to be plotted with the same scale as in the non-adiabatic approach, the MQSSM variation form would look like a near-constant valued line. The constant MQSSM behavior on adiabatic conformational changes is an interesting characteristic of this molecular descriptor, as the active conformation will have, within the precision required in QSAR models, the same value. However, not only is this remarkable MQSSM overlap density trend well suited for QSAR purposes, but also in Figure 8, we show another particular function trend along with the overlap MQSSM, namely, the total Coulomb electron–electron repulsion energy changes within the adiabatic model.

As was shown previously, the Coulomb bielectronic repulsion energy can be considered as a kind of MQSSM, and that their behaviors are similar. Interestingly, in the adiabatic case, the variation of the bielectronic energy is much larger than that in the overlap MQSSM. Although not shown to avoid repetitive graphical information, the trend of kinetic energy is the same as in Figure 8. Again, this result is not unexpected because kinetic energy can



**Figure 8** Evolution of the ab initio (HF/3-21G) Coulomb electron–electron repulsion energy (E-E energy, in Hartree) and the overlap self-similarity measure (in atomic units) versus dihedral angle in ethane during adiabatic bond rotation.

also be considered a first-order similarity measure. The variation scale of the three functions, within the adiabatic scheme, is also comparable: a slowly varying form under rotation angle progression.

---

## CHIRALITY AND THE HOLOGRAPHIC ELECTRON DENSITY THEOREM

---

Before discussing molecular quantum similarity as a basis for so-called quantum QSP/AR, a typical example of the usage of molecular quantum similarity will be discussed.

Chirality plays a key role in many fields of chemistry. Examples include the area of catalysis and especially medicinal chemistry, where two molecules that differ solely in their absolute configuration can exhibit substantially different biological activities. Consider the case of enantiomers. One enantiomer can exhibit the desired medicinal activity, whereas the other may have other characteristics. The latter can simply lack this biological activity, or it can have an undesirable effect, and in extreme cases, it may be poisonous and even deadly. Knowledge of the absolute configuration of a molecule is therefore of prime importance. Chirality is often considered to be a black/white topic, which means that a molecule is either chiral or achiral. As will be shown below, from a molecular quantum similarity point of view, this opinion is not the case, because we can express a degree of chirality. To highlight the utility of molecular quantum similarity, two case studies of molecular chirality will be discussed in some detail.

The first study concentrating on chirality from a molecular quantum similarity point of view was published by Mezey et al.<sup>160</sup> Their study was limited to the case of enantiomeric structures and did not consider diastereomers. Although the usefulness of molecular quantum similarity ideas in chirality had been proposed previously,<sup>161</sup> the publication by Mezey et al. is the first to give actual numerical results for the study of a degree of chirality. With molecular quantum similarity, the binary degree of chirality, which means that a molecule is either chiral or not, is replaced by a continuous degree of chirality. As a metric for chirality, we could use one of the molecular quantum similarity indices already described. In the study by Mezey et al.,<sup>160</sup> they used the overlap similarity measure and calculated the Euclidean distances. Mezey et al. calculated the self-similarities of both enantiomers for a set of eight amino acids and the overlap similarity measure between the two enantiomers. They used ASA promolecular densities as electron densities. No in-depth discussion of the correlation between the Euclidean distances and experimentally observed specific rotations  $[\alpha]_D$  in water with sodium light is given, but an analysis from the reported numerical results reveals that no correlation exists ( $r^2 < 50\%$ ). As a source of error, Mezey et al. refer to the molecular superposition, which may

influence the value for the overlap similarity measure between the two enantiomers. Molecular superposition was performed with the MaxiSim method.

Boon et al.<sup>107</sup> also studied several chiral molecules, which included again two amino acids (Ala and Leu) and  $\text{CHFCIBr}$ , a prototype of chiral molecules. Ab initio total molecular electron densities yielded both Euclidean distances and Carbó indices between the enantiomers of these molecules. Molecular superposition was performed with, on the one hand, a manual alignment based on chemical intuition and the QSSA method, on the other hand. When analyzing the tables of the work by Boon et al. and comparing the results to the work by Mezey et al., similar values for the Euclidean distances between the two enantiomers appear for Ala and Leu, which once again illustrates the power of the ASA promolecular densities to yield quantum similarity measures in good agreement with those obtained from ab initio calculations.

If we abandon the application of the total electron density and use only fragments of the enantiomers involved, it is possible to check whether the holographic electron density theorem holds.<sup>179</sup> According to this theorem, all information on any property of the total molecule is not only contained in the total molecular electron density, but also in any nonzero volume fragment of the highly fuzzy electron density of that molecule. Accordingly, the degree of chirality of a molecule should also be reflected in such fragments. Examples of fragments can be certain functional groups in a molecule, just a single atom, or in fact any volume element of that molecule. Both the studies by Mezey et al.<sup>160</sup> and by Boon et al.<sup>107</sup> try to establish the validity of this theorem by using atom-in-molecules electron densities. As discussed in the context of atom-in-molecule densities, Mezey et al. use a technique based on the classic expression for the electron density and limit the summation to only those basis functions centered within the fragment (see Eqs. [89]–[91]). Mezey et al. consider both the  $\text{NC}_\alpha\text{C}$  fragment and the chiral atom as fragments. When applying their technique, a qualitative correlation is found between the quantum similarity measure  $Z_{RS}^C$  and the absolute value of  $[\alpha]_D$  for the asymmetric carbon in the R and S enantiomers of the eight amino acids considered. This agreement is depicted in Figure 9.

The correlation between  $[\alpha]_D$  and the Euclidean distance between the carbon atoms in the R and S enantiomers remains only qualitative, possibly because of the approximations used in obtaining the fragment density. Besides this, diverse experimental sources and possible errors in  $[\alpha]_D$  may play a role as well as the fact that  $[\alpha]_D$  is a colligative property, depending on concentration and temperature. Another feature in Figure 9 is the presence of an outlier, namely, Serine. Leaving this point out, the correlation increases to 90%.

Boon et al.<sup>107</sup> rely on the Hirshfeld picture of an atom in a molecule. As we mentioned, probably the most obvious way to proceed is via a product of two Hirshfeld coefficients. Boon et al., however, opted to use only one coefficient to retain the spirit of the stockholder idea. First a superposed electron



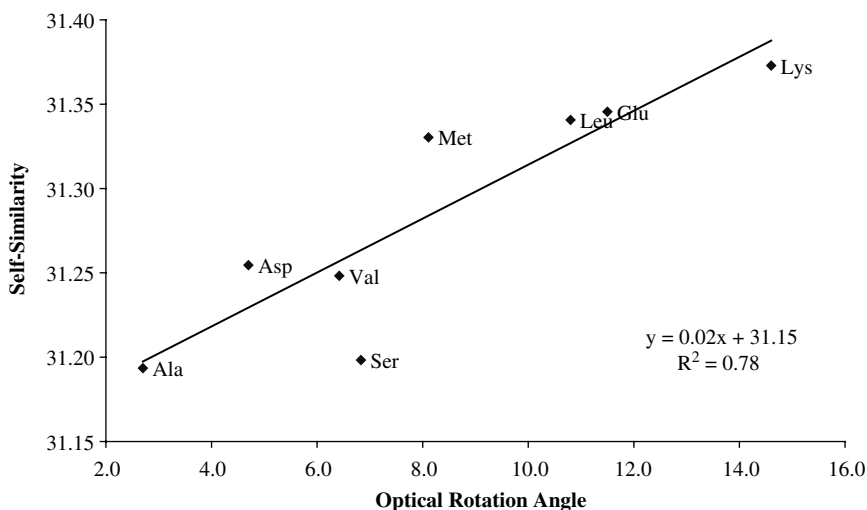


Figure 9 Correlation between  $[\alpha]_D$  and  $Z_{RS}^C$  for eight amino acids.<sup>160</sup>

density over the two enantiomers is constructed, and from that a stockholder coefficient  $w_{C,R+S}(\vec{r})$  is derived as

$$w_{C,R+S}(\vec{r}) = \frac{\rho_{C,R}^0(\vec{r}) + \rho_{C,S}^0(\vec{r})}{\sum_X \rho_{X,R}^0(\vec{r}) + \sum_X \rho_{X,S}^0(\vec{r})} \quad [99]$$

where  $\sum_X \rho_{X,R}^0(\vec{r}) + \sum_X \rho_{X,S}^0(\vec{r})$  is the total promolecular density of the two enantiomers with their asymmetric carbon atom put together. The summations run over all atoms X in both enantiomers.  $Z_{RS}^C$  then becomes

$$\begin{aligned} Z_{RS}^C &= \int w_{C,R+S}(\vec{r}) \rho_R(\vec{r}) \rho_S(\vec{r}) d\vec{r} \\ &= \int \left( \frac{\rho_{C,R}^0(\vec{r}) + \rho_{C,S}^0(\vec{r})}{\sum_X \rho_{X,R}^0(\vec{r}) + \sum_Y \rho_{Y,S}^0(\vec{r})} \right) \rho_R(\vec{r}) \rho_S(\vec{r}) d\vec{r} \end{aligned} \quad [100]$$

Both in the case of Mezey et al. and Boon et al. the calculation of  $Z_{RS}^C$  involves coalescence of the chiral atom in both molecules, similar to the atoms-in-molecules work by Cioslowski et al.<sup>131,132</sup> mentioned earlier. This coalescence provides a drastic reduction of the work needed to obtain the maximal quantum similarity superposition of the two molecules involved. All three translation parameters align the two atoms, so only three parameters remain to be

optimized, which reduces the dimensionality of the QSSA approach by half and improves computational efficiency dramatically.

---

## MATHEMATICAL ASPECTS OF QUANTUM SIMILARITY

A continuing effort has been maintained by the scientific community to provide a firm theoretical and mathematical basis for molecular quantum similarity. According to Carbó-Dorca et al., the basis of molecular quantum similarity and of quantum QSAR (as described later) is founded in the concepts of tagged sets and vector semispaces. To make quantum similarity understandable to the novice, the following explanatory paragraphs provide first the required mathematical basis and second an extensive discussion of some useful aspects of vector semispaces.

As mentioned, the application of molecular descriptors with quantum mechanical origins was proposed several years ago,<sup>180</sup> but the first ideas about quantum similarity (QS) and QS measures (QSM) were published around 1980.<sup>53</sup> However, it has been not until recently that the mathematical and physical foundations of QS have been developed in a series of publications.<sup>5,8-10,65,68,148-149,181-186</sup> In those articles, several new theoretical definitions, related to old concepts, were described. The first idea described was the so-called *tagged set concept*.<sup>187</sup> Tagged sets not only seemed to be essential to QS theory, but they also constituted a generalization of the well-known fuzzy set theoretical setup.<sup>181</sup>

A tagged set  $Z$  is defined as the Cartesian product of two sets:  $Z = O \times T$ , where  $O$ , the *object set*, contains as elements the so-called *objects* and  $T$ , the *tag set*, contains as elements the so-called *tags*. Thus, any element of  $Z$  is constructed by an ordered pair made by an object and a tag. That is,  $\forall \theta \in Z = O \times T \rightarrow \exists \omega \in O \wedge \exists \tau \in T : \theta = (\omega; \tau)$ .

After the seminal structure building of the QS formalism, several additional studies appeared over time, which developed new theoretical details. Especially noteworthy is the concept of *vector semispace* (VSS).<sup>181</sup> This mathematical construction will be shown to be the main platform on which several QS ideas are built, related in turn, to probability distributions and hence to quantum mechanical probability density functions. Such quantum mechanical density distributions<sup>188-196</sup> form a characteristic functional set, which can be easily connected to VSS properties. Construction of the so-called *quantum objects* (QO) and their collections: the *QO sets* (QOS) (see, for example, Carbó-Dorca<sup>181,187</sup>), easily permit the interpretation of the nature of quantum similarity measures for relationships between such quantum mechanically originated elements. Within quantum similarity context, QOS appear as a particular kind of tagged sets, where objects are submicroscopic systems and their density functions become tags.

This section of the chapter deals with the structure and properties of VSS and its mathematical framework. Because VSS is an interesting conceptual formalism, which supports *inward matrix product* (IMP) algebra,<sup>184,186</sup> the applications of such a product in quantum chemistry, particularly in QS theory, will be discussed.

The IMP is based on the structure of the *Hadamard product*,<sup>54</sup> which is related to the result of multiplying two sums, retaining the diagonal cross-terms only. In this way, the Hadamard (or inward) product of two sums can be specified by the algorithm:  $\left(\sum_I^N a_I\right) * \left(\sum_I^N b_I\right) = \sum_I^N a_I b_I$ . The feasibility of a Hadamard product is restricted by the fact that the sums entering the product must possess the same number of terms  $N$ .

Defined in this way, IMP becomes an internal composition law, which can be defined within a matrix (or hypermatrix) vector space  $M_{(m \times n)}(\mathbf{K})$  of arbitrary dimension  $(m \times n)$  and defined over a field  $\mathbf{K}$ , which produces a matrix whose elements are products made in turn by the elements of the matrices appearing in the IMP,<sup>186</sup> according to the straightforward algorithm:

$$\begin{aligned} \forall \mathbf{A} = \{a_{ij}\}, \mathbf{B} = \{b_{ij}\} \in M(\mathbf{K}) : \\ \mathbf{P} = \mathbf{A} * \mathbf{B} \rightarrow \mathbf{P} = \{p_{ij}\} \in M(\mathbf{K}) \wedge p_{ij} = a_{ij} b_{ij}; \forall i, j. \end{aligned} \quad [101]$$

IMP and classic matrix products coincide within the diagonal matrix subspaces. From now on, the IMP and Hadamard products will be synonyms of an operation, which can be applied not only to matrix spaces but also to a wide variety of mathematical objects. Keep in mind that the IMP main characteristic is the result, defined as producing another mathematical object of the same kind as the objects involved in the operation.

The IMP is equivalent to a feature, included in high-level computer languages such as Fortran 95,<sup>197</sup> where it has been implemented in an easy useable manner, so the practical programming of the following IMP properties and characteristics can be carried out rapidly. IMP is commutative, associative, and distributive with respect to the matrix sum. Moreover, it has a multiplicative neutral element, the unity matrix, which is customarily represented by a bold real unit symbol and formally defined as  $\mathbf{1} = \{1_{ij} = 1\}$ .

IMP inverse matrices exist whenever we take into account matrices without zeroes as elements. Any nonsingular IMP matrix is defined as

$$\mathbf{A} = \{a_{ij}\} \in M \wedge \forall i, j : a_{ij} \neq 0 \rightarrow \mathbf{A}^{[-1]} = \{a_{ij}^{-1}\}, \quad [102]$$

where the following equation holds:  $\mathbf{A} * \mathbf{A}^{[-1]} = \mathbf{A}^{[-1]} * \mathbf{A} = \mathbf{1}$ . With respect to both the usual sum and the IMP, matrix sets of the same dimension behave

almost like scalars. Any IMP function may be expressed in the same way as the inward operation; that is,

$$\varphi[\mathbf{A}] = \mathbf{F} = \{f_{ij} = \varphi(a_{ij})\} \quad [103]$$

However, in computational practice, the IMP inverse of a given matrix having elements with zeros can be performed, admitting the definition of a particular unity element having values of ones in the appropriate places, i.e., associated with a matrix having ones instead of zeroes in the same positions. That is, defining the matrix:

$$\mathbf{1}_0(\mathbf{A}) = \{1_{0;ij} \equiv \delta(a_{ij} = 0)\} \quad [104]$$

with a *logical Kronecker delta*,<sup>198,199</sup> to inform which elements are zeroes in the matrix  $\mathbf{1}_0$ , and which are ones, coincident with zero and nonzero elements, respectively, in the associated matrix  $\mathbf{A}$ . Then, the IMP pseudo-inverse of such a matrix can be defined as

$$\mathbf{A}_0^{[-1]} = (\mathbf{A} + \mathbf{1}_0(\mathbf{A}))^{[-1]} - \mathbf{1}_0(\mathbf{A}) \quad [105]$$

where the following relationships will hold:

$$\mathbf{A}_0^{[-1]} * \mathbf{A} = \mathbf{A} * \mathbf{A}_0^{[-1]} = \mathbf{1} - \mathbf{1}_0(\mathbf{A}) \quad [106]$$

IMP powers will from here on be defined as successive IMP over the same matrix, as many times as needed. IMP powers will be depicted with square brackets to distinguish them from the usual matrix powers, for example,

$$\mathbf{A}^{[2]} = \mathbf{A} * \mathbf{A} \quad [107]$$

For more details on IMP features, see Carbó-Dorca.<sup>184,186</sup>

IMP composite operations can also be applied to define scalar products. For this purpose, we can define first the *total sum* of the elements of an arbitrary matrix:

$$\mathbf{A} = \{a_{ij}\} \in M \quad [108]$$

by means of the symbol:

$$\langle \mathbf{A} \rangle = \sum_i \sum_j a_{ij} \quad [109]$$

Connecting this definition with both the Hadamard and IMP definitions, we can easily write:

$$\langle \mathbf{A} * \mathbf{B} \rangle = \langle \mathbf{A} \rangle * \langle \mathbf{B} \rangle \quad [110]$$

Then, it is a matter of simple reasoning to construct the scalar product of two matrices, which will be symbolized here as  $\langle \mathbf{A} | \mathbf{B} \rangle$ , by means of the IMP structure. Thus,

$$\langle \mathbf{A} | \mathbf{B} \rangle = \sum_i \sum_j a_{ij} b_{ij} = \langle \mathbf{A} * \mathbf{B} \rangle \quad [111]$$

In this way, the further definition of distances and cosines of the angle between two matrices can be outlined straightforwardly. For instance, the cosine of the angle subtended by two matrices can be written, according to Eq. [111], as

$$\cos(\alpha) = (\langle \mathbf{A} * \mathbf{A} \rangle \langle \mathbf{B} * \mathbf{B} \rangle)^{-\frac{1}{2}} \langle \mathbf{A} * \mathbf{B} \rangle \quad [112]$$

Similar statements and definitions can be made for other kinds of norms, as will be discussed.

The most important concept in the theoretical development of molecular quantum similarity, lies in the VSS definition. A VSS<sup>181</sup> is a vector space in which, instead of the additive group, an additive *semigroup* has been chosen. Although the VSS is not defined in the *Encyclopaedia of Mathematics*,<sup>54</sup> the semigroup structure is well documented. An additive semigroup is an additive group without reciprocal elements, which is the same as saying that negative elements are not present in the VSS. Construction of a specific matrix VSS, for example, will be made below by matrices whose elements are positive definite or semidefinite, although this latter case will not be deeply discussed further on, as QS structures matrices are of the first kind only. Thus, the matrix elements forming a VSS are constructed by positive definite real numbers extracted from  $\mathbf{R}^+$ . In this sense, all elements of a matrix VSS are nonsingular matrices from the IMP point of view, whereas the additive neutral element or any matrix with a zero element will be nonexistent in a VSS, if this strict sense is adopted. A function-made VSS can be constructed by positive definite functions over a given domain and lacking of the null function element to comply with the strict VSS characteristics.

Because of the positive definite structure of the components of a VSS, the easiest way to construct a norm within such a mathematical configuration is, perhaps, the Minkowski definition. That is, in a matrix VSS, we can easily write:

$$\forall \mathbf{A} \in M(\mathbf{R}^+) \rightarrow \langle \mathbf{A} \rangle \in \mathbf{R}^+ \quad [113]$$

where the symbol:  $\langle \mathbf{A} \rangle$  refers to the complete sum of the matrix elements as defined in Eq. [109].

In any general functional VSS, an equivalent form can also be defined:

$$\forall \rho(\mathbf{r}) \in F(\mathbf{R}^+) \rightarrow \langle \rho \rangle = \int_D \rho(\mathbf{r}) d\mathbf{r} \in \mathbf{R}^+ \quad [114]$$

Finally, we point out that a Minkowski norm in  $M(\mathbf{R}^+)$  and, thus, the complete sum of a matrix elements can be considered as a linear operator. That is,

$$\langle \alpha \mathbf{x} + \beta \mathbf{y} \rangle = \alpha \langle \mathbf{x} \rangle + \beta \langle \mathbf{y} \rangle \quad [115]$$

the Minkowski-like norms classify the VSS elements in subsets, which will be called here  $\sigma$ -shells,  $S(\sigma)$ , whose elements are defined by the value of such a norm. That is,

$$\forall x \in S(\sigma) \subset V(\mathbf{R}^+) \rightarrow \langle x \rangle = \sigma \in \mathbf{R}^+ \quad [116]$$

The *unit shell*  $S(1)$  becomes in this way a VSS subset, which can generate all other VSS shells. The following evidence of this property can be shown:

$$\forall z \in S(\sigma) \rightarrow \exists x \in S(1) : z = \sigma x \quad [117]$$

This simple VSS substructure permits us to consider roughly that some ordering principle is induced in VSS subsets, because of the VSS present  $\sigma$ -shell structure. That is, we can write the following ordering relationship, between two vectors belonging to a pair of different  $\sigma$ -shells:

$$x \in S(\sigma) \wedge y \in S(\varsigma) \rightarrow x \prec y \text{ iff } \sigma < \varsigma \quad [118]$$

In this case, we can say that the vector  $x$  *precedes* the vector  $y$  in the  $\sigma$ -shell sense.

It is simple to extend the idea of shell structure into any normed space. Suppose a VS over a field  $\mathbf{K}$ , provided with some positive definite norm,  $N(x)$ , that is,

$$\forall x \in V(\mathbf{K}) \rightarrow N(x) \in \mathbf{R}^+ \quad [119]$$

Thus, a  $\sigma$ -shell in  $V(\mathbf{K})$  is defined in the same fashion as in the previous Minkowski norm case:

$$\forall x \in S(\sigma) \subset V(\mathbf{K}) \rightarrow N(x) = \sigma \in \mathbf{R}^+ \quad [120]$$

Equivalent considerations to the discussed  $\sigma$ -shell properties can be also admitted and new ones added, as follows.

Operations involving vectors belonging to a VSS  $\sigma$ - shell can be defined as closed with the appropriate modifications. Otherwise, the vector sum and IMP cannot bear such closure property. However, defining a *modified vector sum* in the following terms:

$$\forall x, y \in S(\sigma) : x \oplus y = \sigma \langle x + y \rangle^{-1} (x + y) = \frac{1}{2}(x + y) \quad [121]$$

with the Minkowski norm acting as a linear operator. In the same manner, a closed summation symbol within a  $\sigma$ - shell can be defined whenever the vector sum is computed with the algorithm:

$$X = \{x_I\} \subseteq S(\sigma) \wedge \#X = N : \oplus_I [x_I] = \frac{1}{N} \sum_I x_I \quad [122]$$

Moreover, a closed IMP product can be defined with a similar technique:

$$\forall x, y \in S(\sigma) : x \odot y = \sigma \langle x * y \rangle^{-1} (x * y) \quad [123]$$

The scalar product of the involved vectors corresponds, after inversion, to a Minkowski normalization factor, which transforms the IMP resultant vector into an element of the unit shell. In fact, Eq. [124] involves the already demonstrated property (see Eq. [111]), meaning that the Minkowski norm of an IMP, involving two matrices, coincides with their scalar product.

Another important concept in molecular quantum similarity is associated with convex conditions. The idea underlying *convex conditions*, associated with a numerical set, a vector, a matrix, or a function, has been described previously in the initial work on VSS and related issues.<sup>181</sup> Convex conditions correspond to several properties of some mathematical objects. The symbol  $K(\mathbf{x})$  means that the conditions  $\langle \mathbf{x} \rangle = 1 \wedge \mathbf{x} \in V(\mathbf{R}^+)$  hold simultaneously for a given mathematical object  $\mathbf{x}$ , which is present as an argument in the *convex conditions symbol*. Convex conditions become the same as considering the object as a vector belonging to the unit shell of some VSS. For such kind of elements,

$$K(\mathbf{x}) = \{ \langle \mathbf{x} \rangle = 1 \wedge \mathbf{x} \in V(\mathbf{R}^+) \} \equiv \{ \mathbf{x} \in S(1) \} \quad [125]$$

Alternatively, the following property holds over any element of the unit shell:

$$\forall \mathbf{x} \in S(1) \rightarrow K(\mathbf{x}) \quad [126]$$

Given an arbitrary  $\sigma$ - shell  $S(\sigma) \subset V(\mathbf{R}^+)$ , of some VSS, convex combinations of the elements of the  $\sigma$ - shell then produce a new element of  $S(\sigma)$ .

That is, suppose that the convex conditions

$$K(\{\gamma_I\}) = \left\{ \sum_I \gamma_I = 1 \wedge \forall I : \gamma_I \in \mathbf{R}^+ \right\} \quad [127]$$

hold on a known set of scalars  $\{\gamma_I\}$ . Then, the following property will be fulfilled for any arbitrary subset of elements belonging to the chosen  $\sigma$ - shell:

$$\{\mathbf{x}_I\} \subset S(\sigma) \wedge K(\{\gamma_I\}) \rightarrow \mathbf{x} = \sum_I \gamma_I \mathbf{x}_I \in S(\sigma) \quad [128]$$

owing to the fact that the summation symbol, associated here with a Minkowski norm, can be considered as a linear operator. Then the following property holds:

$$\langle \mathbf{x} \rangle = \left\langle \sum_I \gamma_I \mathbf{x}_I \right\rangle = \sum_I \gamma_I \langle \mathbf{x}_I \rangle = \sum_I \gamma_I \sigma = \sigma \sum_I \gamma_I = \sigma \rightarrow \mathbf{x} \in S(\sigma) \quad [129]$$

Such a property is related to the possibility of constructing approximate atomic and molecular *density functions*, by means of convex combinations with a basis of structurally simpler functions, which belong to the same VSS  $\sigma$ - shell. The ASA described extensively earlier and in Amat et al.,<sup>87-95</sup> Mestres et al.,<sup>200</sup> and Carbó-Dorca and Gironés<sup>201</sup> is a good illustration of these ideas. Moreover the VSS shell structure can transform density functions into a homothetic construct, a characteristic that is discussed elsewhere.<sup>202</sup> The need to take into account the convex conditions [127] to construct approximate density functions has not been properly performed in the literature, as discussed.

Any VSS  $\sigma$ - shell structure can be generated from a conventional *vector space* (VS). A VS can be defined over the complex or real fields and can be provided for convenience with a positive definite metric structure. Indeed, suppose such a VS, defined for the sake of generality over the complex field,  $V(\mathbf{C})$ . Then, from a general point of view, the following algorithm can be envisaged:

$$\begin{aligned} \forall v \in V(\mathbf{C}) \wedge v \neq \mathbf{0} \rightarrow \langle v | v \rangle &= \sigma \in \mathbf{R}^+ \\ \Rightarrow \exists x = v^* * v \in S(\sigma) &\subset V(\mathbf{R}^+) \end{aligned} \quad [130]$$

where the IMP or Hadamard product  $x = v^* * v$  has constructed the VSS vector. Then the following sequence:

$$\langle x \rangle = \langle v^* * v \rangle = \langle v | v \rangle = \sigma \quad [131]$$

holds and has set the form of Eq. [130].



In addition, from a complementary point of view, a symbol to briefly summarize Eq. [130] could be easily constructed. We can say, in this way, that vector  $v$  *generates* vector  $x$ , whereas writing accordingly a *generating symbol*  $R(v \rightarrow x)$ , whenever the sequence of relationships in Eq. [130] holds.<sup>181</sup> The quantum mechanical image of the density function construction appears as a particular case of the definition attached to Eq. [130]. Thus, it is interesting now that from a quantum mechanical point of view, when a wave function  $\Psi(\mathbf{r})$  is known, the attached density function  $\rho(\mathbf{r})$  is generated in the following way:  $R(\Psi \rightarrow \rho)$ , according to the mathematical details of the previous section related to VSS generation.<sup>181</sup>

In this generating sense, we can consider the *dimension of a VSS*  $V(\mathbf{R}^+)$  as bearing the dimension of the generating VS  $V(\mathbf{R})$ . Alternatively, we can formally write:  $\text{Dim}[V(\mathbf{R}^+)] \equiv \text{Dim}[V_N(\mathbf{R})] = N$ , which may be based on the fact that we can write the generating implication:

$$R(V_N(\mathbf{R}) \rightarrow V_N(\mathbf{R}^+)) \quad [132]$$

We can use a similar set of symbolic forms when the generating space is defined over the complex field.

At this stage, the problem of the basis set in VSS develops. As the solution is not as obvious as in VS, because of the positive definite structure of the VSS elements, some details are discussed briefly here. When considering VSS made of  $N$ -dimensional column matrices as chosen elements, many simple VS basis sets can be used to construct linearly independent vectors in VSS. For instance, when choosing the *canonical* basis set in a column matrix VS, each element of the basis set, the  $I$ -th, say, is made by the  $I$ -th column of the corresponding unit matrix,  $\mathbf{I}_N$ . Thus, such a canonical basis set element is made by a unit in the  $I$ -th position and zeros in the rest. Therefore, the canonical basis set  $\{|\mathbf{e}_I\rangle\}$  could be defined using the Kronecker's delta symbol,  $\delta_{JI}$ , as

$$|\mathbf{e}_I\rangle = \{e_{JI} = \delta_{JI}\} \quad [133]$$

A collection of linearly independent sets in VSS can be obtained from the unit matrix by the IMP multiplicative unit forming the scalar matrix  $\alpha \mathbf{1}$ . So, the resultant matrix:

$$\mathbf{J}_N(\alpha) = \mathbf{I}_N + \alpha \mathbf{1}_N \quad [134]$$

has a simple determinant value:

$$\text{Det}[\mathbf{J}_N(\alpha)] = 1 + N\alpha \quad [135]$$

and in this way generates a continuous set of nontrivial linearly independent vectors in VSS, whenever  $\alpha > 0$ . Also, the columns of the matrix

$\mathbf{J}_N(\alpha) = \{\mathbf{j}_I(\alpha)\}$  belong by construction to the  $\sigma$ - shell  $S(1 + N\alpha)$ . Thus, convex combinations of the column vectors of a chosen matrix  $\mathbf{J}_N(\alpha)$ , with a fixed parameter value,  $\alpha$ , generate the vectors:

$$\begin{aligned} \forall K(\{\gamma_I\}) \wedge \alpha > 0 : \mathbf{v} &= \sum_I \gamma_I \mathbf{j}_I(\alpha) \rightarrow \langle \mathbf{v} \rangle = \sum_I \gamma_I \langle \mathbf{j}_I(\alpha) \rangle \\ &= \sum_I \gamma_I (1 + N\alpha) = (1 + N\alpha) \\ &\rightarrow \mathbf{v} \in S(1 + N\alpha) \end{aligned} \quad [136]$$

and the resultant vectors belong to the same  $\sigma$ - shell as the  $\mathbf{J}_N(\alpha)$  columns, as expected. Thus, with the components of the vectors of the unit shell, we can construct the vectors of the  $\sigma$ - shell where the columns of  $\mathbf{J}_N(\alpha)$  belong.

In the VSS, however, the back transformation from a chosen basis, like that represented by the collection attachable to the  $\mathbf{J}_N(\alpha)$  matrix, is not well defined, because the matrix inverses can no longer furnish the VSS basis set elements; their columns belong to a VS in general, because the inverse matrix elements are no longer positive definite. For instance, the inverse of the  $\mathbf{J}_N(\alpha)$  matrix can be expressible as

$$\mathbf{J}_N^{-1}(\alpha) = \mathbf{I}_N - \left( \frac{\alpha}{1 + N\alpha} \right) \mathbf{1}_N \quad [137]$$

The minus sign forbids the inverse matrix columns in Eq. [137] from belonging in general to any VSS. Other obvious examples can be presented, but this one is sufficient to show the difficulties that develop when defining isomorphism rules or endomorphism transformations in the VSS.

Another interesting point concerns the connection between the VS structure and the VSS form, as well as with the peculiar structure of the VSS. Any probability distribution, discrete or continuous, belongs to some VSS unit shell. Probability distributions can be generated by the normalized elements of some normed or metric space, so that the resultant VSS element belongs to the unit shell,  $S(\sigma)$ , and in this way potentially to any other VSS  $\sigma$ - shell.

Thus, the unit shell represents every other  $\sigma$ - shell in the VSS. Probability distributions of any kind can be transformed into any other element of the associated VSS in this manner. Any VSS  $\sigma$ - shell,  $S(\sigma)$ , can be considered like a homothetic construct of the unit shell,<sup>202</sup>  $S(1)$ . Because the elements of the unit shell comply with the form of a probability distribution, they then also fulfill the convenient convex conditions as shown in Eq. [126]. In other words, any probability distribution can be considered an element of the unit shell forming part of a VSS provided with the appropriate dimension.

Because of these connections to probability vectors, scalar products of two distinct compatible probability distributions are always positive definite, so we have:

$$\forall x, y \in S(1) \rightarrow \langle x | y \rangle = \langle x * y \rangle \in \mathbf{R}^+ \quad [138]$$

Therefore, the cosine of the angle subtended by two probability distributions has to be contained in the open interval:  $(0, 1]$ , because of Eq. [138]. This requirement is so because the cosine of the subtended angle between two elements can be computed as

$$\forall x, y \in S(1) : \cos(\alpha) = (\langle x * x \rangle \langle y * y \rangle)^{-\frac{1}{2}} \langle x * y \rangle \in (0, 1] \quad [139]$$

Furthermore, Eq. [138] shows that the scalar product between two, or more, elements of the unit shell of a given metric VSS constitute a *measure*. In this way, such a scalar product can be considered a generalized volume.

However, these conclusions are not the only ones that we can reach from the VSS metric properties: A *generalized scalar product* of the sort commented on before can be also constructed with the aid of the following algorithm associated with the IMP:

$$\forall \{x_I\} \subset S(1) \rightarrow \langle x_1 * x_2 * \dots x_P * \dots \rangle \in \mathbf{R}^+ \quad [140]$$

This algorithm ensures in any circumstance that the algorithm of the products constructed by a rule with Eq. [140] will always yield a positive definite result. Such a generalized scalar product can be considered as a measure constructed over probability distributions.

A similarity measure over the unit shell of any VSS can be defined through the description of the mathematical elements, which have been described so far. In the simplest way, an MQSM can be defined knowing the appropriate density function tags of two quantum objects  $\{\rho_A; \rho_B\}$ , adapted to some shells of the corresponding VSS, and with as a weight some positive definite operator  $\Omega$ . In that case, the integral measure

$$z_{AB}(\Omega) = \iint_D \rho(\mathbf{r}_1) \Omega(\mathbf{r}_1; \mathbf{r}_2) \rho(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 = \langle \rho_A * \Omega * \rho_B \rangle \in \mathbf{R}^+ \quad [141]$$

will correspond to a weighted scalar product, defined over the shell elements, made in turn by the compatible quantum density functions. The MQSM of Eq. [141] can be associated with a property comparable with the one encountered in Eqs. [138] and [140], and it always must have a positive definite nature.

Consider now a set of quantum systems  $S = \{s_I\}$  in a well-defined set of states. Suppose that to every quantum system, a known state density function is attached, which forms the set  $P = \{\rho_I\}$ , belonging to the unit shell of some functional VSS. A tagged set can be constructed with the Cartesian product  $T = S \times P$ , where each element,  $\tau_I \in T$ , is constructed by the ordered pair composition rule  $\tau_I = (s_I; \rho_I)$ , which forms in this way a QO. The tagged set  $T$  constitutes a QOS; that is,  $T = \{\tau_I\}$ . The MQSM defined earlier in Eq. [141] can be interpreted, in turn, as a tensor product of the tag part of the QOS.

Collecting all MQSMs computed between the element pairs of a given QOS, a so-called *quantum similarity matrix* is obtained, having been constructed according to the definition [141] by means of  $\mathbf{Z} = \{z_{IJ}\}$ . Because of the structure of the quantum similarity matrix elements, the matrix can be considered as an element of a VSS of some appropriate dimension. The similarity matrix  $\mathbf{Z}$  is a symmetric matrix with positive definite elements, whose columns  $\{z_I\}$  (or rows) are also elements of some  $N$ -dimensional VSS. As such, a real symmetric matrix  $\mathbf{X}$  exists that, in general, generates  $\mathbf{Z}$  as

$$\mathbf{Z} = \mathbf{X} * \mathbf{X} = \mathbf{X}^{[2]} \vee \mathbf{X} = \mathbf{Z} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad [142]$$

or  $R(\mathbf{X} \rightarrow \mathbf{Z})$ , which is the same. As a consequence, any similarity matrix belongs to a precise  $\sigma$ -shell of some VSS. That is,

$$\forall \mathbf{Z} : \langle \mathbf{Z} \rangle = \sum_i \sum_j z_{ij} = \sigma \rightarrow \mathbf{Z} \in S(\sigma) \subseteq M(\mathbf{R}^+) \quad [143]$$

Even if the columns or rows of the similarity matrix  $\mathbf{Z}$  belong to different  $\sigma$ -shells of some VSS, all of them can be brought to the unit shell easily by a set of simple homothetic transformations, which involve a product by a diagonal matrix, with elements constructed by the Minkowski norms of the columns (or rows) of the similarity matrix. That is, the diagonal matrix:

$$\mathbf{D} = \text{Diag}(\langle z_1 \rangle; \langle z_2 \rangle; \dots \langle z_I \rangle; \dots) \quad [144]$$

can transform the quantum similarity matrix  $\mathbf{Z}$  into a column (or row) *stochastic matrix*, as described. This transformation simply involves multiplying on the right (or the left) of  $\mathbf{Z}$  by the inverse of  $\mathbf{D}$ . For instance, the stochastic column matrix associated with the quantum similarity matrix is

$$\mathbf{S} = \mathbf{ZD}^{-1} \quad [145]$$

In this way, the columns  $s_I$  of the stochastic matrix,  $\mathbf{S}$ , belong to the unit shell of the column vector VSS of the appropriate dimension. That is,

$$\mathbf{S} = \{s_I\} \rightarrow \forall I : \langle s_I \rangle = \left\langle \langle z_I \rangle^{-1} z_I \right\rangle = \langle z_I \rangle^{-1} \langle z_I \rangle = 1 \rightarrow s_I \in S(1) \quad [146]$$

However, the column stochastic similarity matrix [145] appears to be no longer symmetric as is its originating quantum similarity matrix **Z**. The columns of the stochastic matrix [145], as defined in Eq. [146], can substitute the QOS density function tag elements, which have previously generated them. A new kind of discrete QOS can be constructed in this way.

This brief mathematical background seems sufficient to prove how the quantum similarity application structure can be well founded on a new set of theoretical concepts, which encompasses basic quantum mechanical ideas.

---

## THE CRAMER STEROID SET—A WORKED OUT EXAMPLE OF MQS

---

It is appropriate to illustrate these concepts with a worked out example of a typical application of molecular quantum similarity ideas. The example addressed here involves the set of globulin bindings steroids used by Cramer et al.<sup>203</sup> and subsequently in other studies to develop QSAR models.<sup>204</sup> This dataset has also been used by chemists in molecular quantum similarity studies and to develop quantum QSAR models.<sup>205,206</sup>

The 31 molecules included in this set are shown in Figure 10. From these simple graphical representations of the molecular structures, the degree of quantum similarity cannot be derived, so 3-D structures were generated first. In the application of molecular quantum similarity and quantum QSAR, cheap and relatively simple methods like AM1<sup>41,42</sup> are used often to obtain 3-D molecular geometries. Naturally, more sophisticated methods like Hartree–Fock and post-Hartree–Fock methods and DFT can also be applied, and nothing opposes the application of experimental geometries either. Once the molecular geometries are available, the density functions are computed and the MQSM is calculated.

For the set of steroid molecules illustrated in Figure 10, the geometries were obtained with the semi-empirical AM1 method. Electron densities were then obtained with both B3LYP/6-31G\* single-point calculations and promolecular ASA densities. The latter are generated much faster than the ab initio densities. Once the electron density is known, the MQSM can be calculated with any of the positive definite operators mentioned earlier.

A complete discussion of the molecular quantum similarity matrix would be repetitive. Therefore, discussion is limited here to a 5-by-5 submatrix. The molecules included from Figure 10 are molecules 3, 5, 6, 17, and 24. These molecules were chosen to show, among other aspects of quantum similarity, the influence of molecular alignment for evaluating the overlap MQSM. Molecular alignment was performed with both TGSA and QSSA. The latter technique provides results similar to MaxiSim, and it seeks the alignment that provides the maximum MQSM. TGSA is a typical structural alignment technique that does not attempt to maximize the MQSM. In Table 2 some values

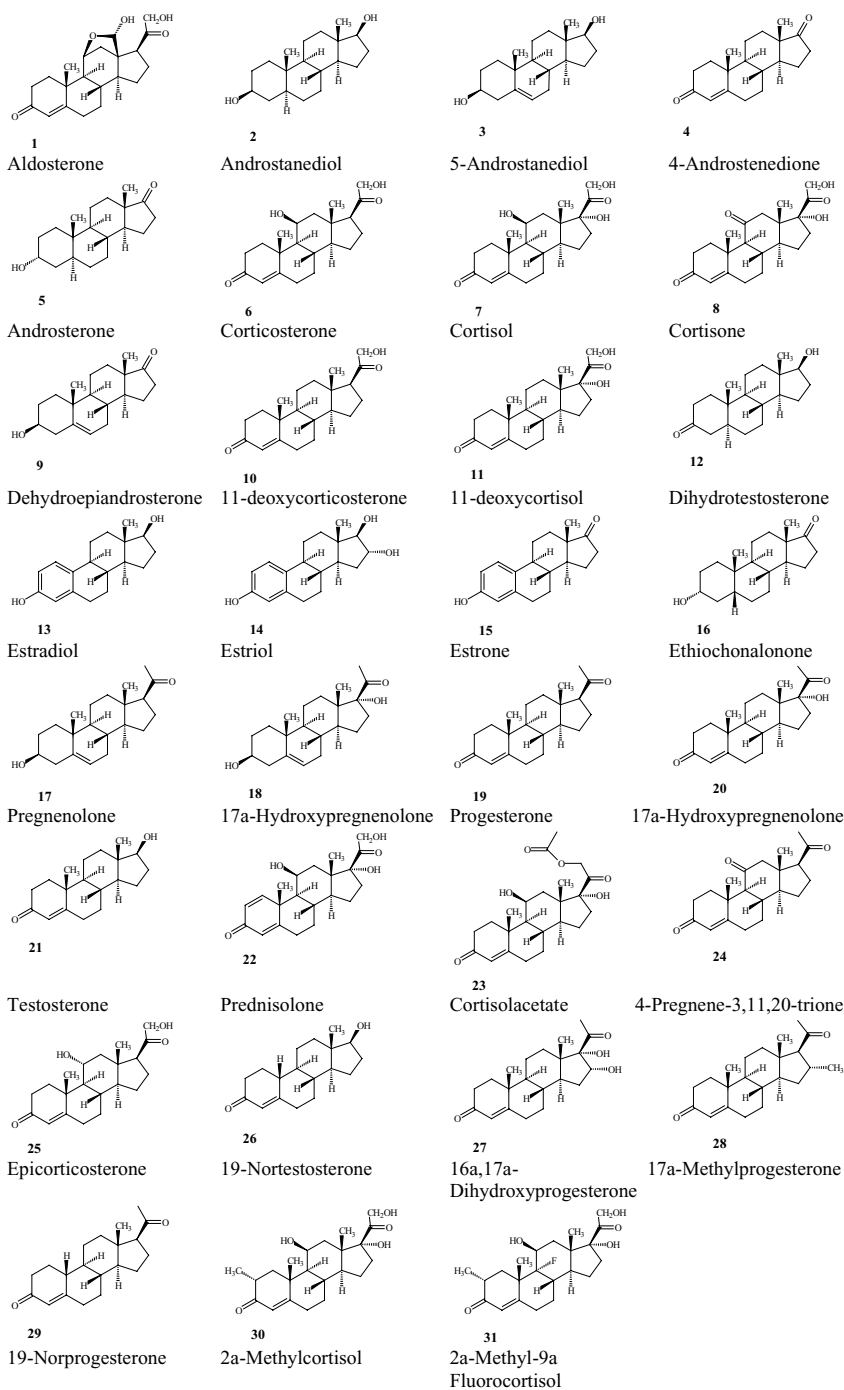


Figure 10 The set of 31 steroids considered in the present example application.

**Table 2** Some Overlap MQSM Between Several Sets of Two Steroid Molecules A and B of Figure 10

A	B	$Z_{AB}$
3	18	723.2 <i>519.4</i>
3	4	488.3 <i>439.4</i>
14	15	657.2 <i>530.4</i>
22	26	466.1 <i>345.4</i>
22	31	826.5 <i>535.3</i>

The plain numbers refer to QSSA alignment, the ones in italics to TGSA alignments. *Numerical data reprinted with permission from Bultinck et al.*<sup>98</sup> Copyright 2003 American Chemical Society.

are listed for the overlap MQSM for combinations of two steroid molecules A and B selected from Figure 10.

Table 2 clearly shows that alignment influences the MQSM obtained, even in the present case in which the molecules are congeneric. A graphical inspection of the alignments produced reveals that with TGSA, most atoms are aligned approximately, but they do not always overlap exactly. With QSSA, the heaviest atoms coincide better than with TGSA but at the cost of some other atoms that are then slightly more separated. Given the behavior of the MQSM as shown in Figure 5, a small difference in alignment can significantly affect the value of the MQSM. That the differences in MQSM with two different alignments are indeed substantial is revealed when looking at the Carbó indices for each pair of molecules. Both resulting matrices are shown here:

$$C_{TGSA} = \begin{bmatrix} 1.0000 & 0.4421 & 0.5320 & 0.6435 & 0.4221 \\ & 1.0000 & 0.4583 & 0.3954 & 0.3785 \\ & & 1.0000 & 0.6691 & 0.4957 \\ & & & 1.0000 & 0.5476 \\ & & & & 1.0000 \end{bmatrix} \quad [147]$$

$$C_{QSSA} = \begin{bmatrix} 1.0000 & 0.4952 & 0.5868 & 0.8954 & 0.5238 \\ & 1.0000 & 0.4708 & 0.4747 & 0.4491 \\ & & 1.0000 & 0.6878 & 0.4957 \\ & & & 1.0000 & 0.5880 \\ & & & & 1.0000 \end{bmatrix} \quad [148]$$

As shown, in some extreme cases, differences up to nearly 25% occur, which is naturally an appreciable amount. Moreover, the difference between the TGSA and QSSA alignments is not constant, and consequently, the similarity ordering of the pairs of molecules differs. Both approaches have their limitations and their merits. Whereas QSSA (and MaxiSim) gives higher MQSM values, TGSA often gives results in better accordance with chemical intuition. It is therefore up to users to carefully decide which approach is best suited for their needs.

We can use several other indices beside the Carbó index. One of the most interesting indices is the Euclidean distance, which constitutes a real metric, in the sense that it obeys the triangular inequality.<sup>54,55</sup> So, a further requirement that must be fulfilled is that the MQSM matrix should not contain any inconsistencies. Inconsistency checking should always be carried out to increase the reliability of the MQSM in the similarity matrix.<sup>98</sup>

Once the MQSM is available, and as a consequence the Carbó indices, we can also obtain dendrograms that graphically depict the relation between different molecules and sets of molecules. The dendrogram obtained from the overlap MQSM for the complete set of steroids is shown in Figure 11.

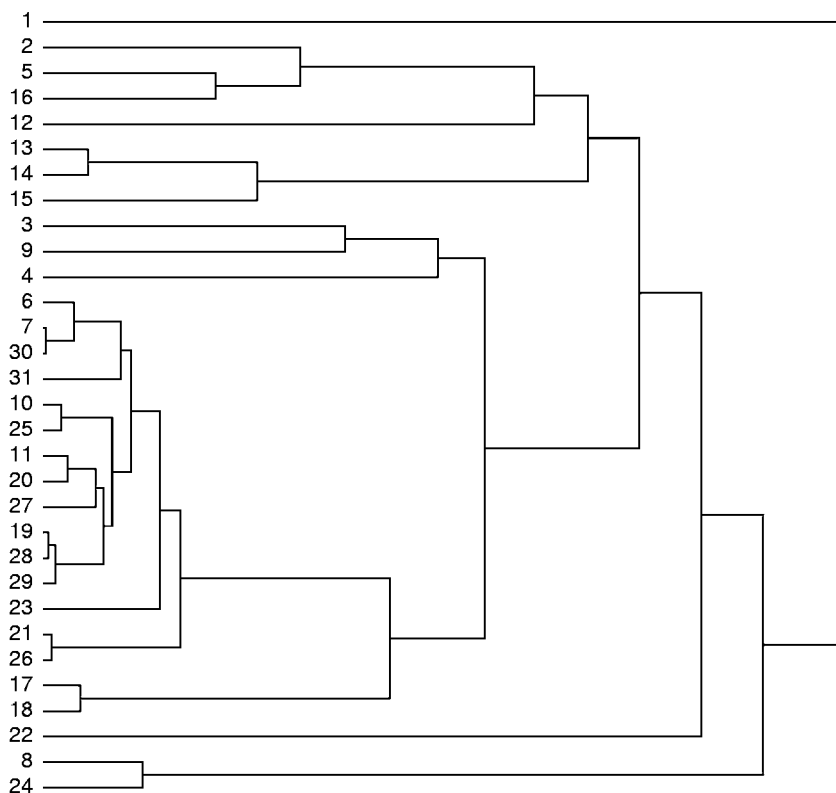
Such dendrograms provide interesting additional information, especially where links appear between previously clustered sets of molecules.

A noticeable influence on the results of the similarity analysis is exerted by the type of alignment algorithm. Another matter is the choice of the electron density. It is worthwhile to examine how the results change when we use the approximate promolecular ASA density instead of the *ab initio* density. For steroid molecules 1-9 of Figure 10, B3LYP/6-31G\* electron densities were calculated with the Gaussian-98 program. This level of calculation is generally considered as a sufficiently high in computational medicinal chemistry. The *ab initio* MQSMs were then calculated for the  $9 \times 9$  similarity matrix with the BRABO *ab initio* program. The agreement between the B3LYP/6-31G\*- and the ASA-derived similarity measures is depicted graphically in Figure 12.

The application of the promolecular ASA densities is certainly warranted, considering the good agreement between the DFT and the ASA MQSM, which does not mean that expectation values obtained with the promolecular ASA density would also be good. These results confirm previous findings for a set of binding isomers, in which promolecular ASA densities were also found to be applicable to the calculation of MQSM.

We can also use the MQSM matrices in a field known as quantum QSAR, where QSAR is performed with similarity matrices instead of the more classic molecular descriptors. The Cramer set of steroids have been studied extensively in this regard, and it was found that good QSAR models can be obtained with only the similarity matrices as obtained through molecular quantum similarity. Quantum QSAR is, however, outside the scope of this chapter because of its involved nature, which would require a lengthy and

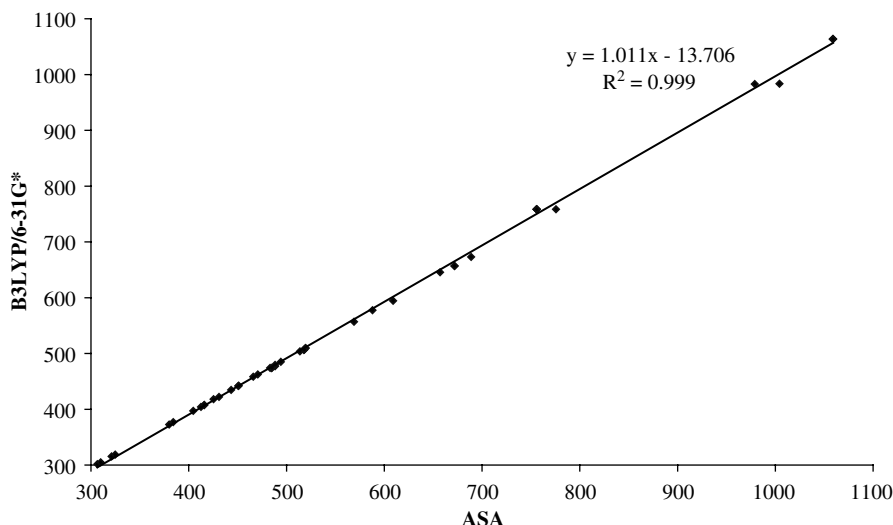




**Figure 11** Dendrogram obtained with the Carbó similarity index for the set of 31 steroids. *Figure reprinted with permission from Bultinck and Carbó-Dorca.<sup>71</sup> Copyright 2003 American Chemical Society.*

profound discussion on both the theoretical aspects as well as the applications. Reviews of quantum QSAR may be found in Carbó-Dorca et al.,<sup>9–12</sup> and more background in the mathematical aspects may be found in Carbó-Dorca et al.<sup>183–186</sup>

Although a limited number of examples of molecular quantum similarity applications have been described in this chapter, including chirality, the holographic electron density theorem, molecular clustering and visualization, calculation of log P and the Hammett sigma constant, molecular quantum similarity has many other applications. We can mention the optimization of parameters in hybrid density functionals, in which parameters are optimized by pursuing maximal similarity between the DFT densities and high-level ab initio electron densities<sup>207,208</sup> or by investigating the influence of solvation on electron densities<sup>209</sup> as examples.



**Figure 12** Agreement between B3LYP/6-31G\*- and ASA-derived quantum similarity measures for molecules 1–9 of the steroid set. *Reprinted with permission from Bultinck et al.*<sup>98</sup> Copyright 2003 American Chemical Society.

---

## CONCLUSIONS

This chapter has clearly shown molecular quantum similarity to be a worthy field of research with a large range of applications. It allows quantification of the similarity between any pair of atoms or molecules, and it can be used by chemists for all quantum entities for which a density probability function is available. Consequently, molecular similarity can be extended far beyond the traditional approaches that address mainly congeneric molecules, or molecules differing in substituent composition.

A firm theoretical basis has been established for molecular quantum similarity, and many computational tools have been developed that allow for the evaluation and quantification of molecular quantum similarity measures among sets of molecules or atoms. Molecular quantum similarity is also the basis of quantum QSAR, another active field of research.

---

## ACKNOWLEDGMENTS

The Fund for Scientific Research-Flanders (Belgium) is thanked for their grants to the Computational Chemistry group at Ghent University. We acknowledge the European Community—Access to Research Infrastructure action of the Improving Human Potential Programme, which

allows the applications of the CEPBA infrastructure at the PolyTechnical University of Catalonia (Spain) and the fellowships with the Institute of Computational Chemistry at the University of Girona (Catalonia, Spain). This work has been partially funded by a grant from the Spanish Ministerio de Ciencia y Tecnología: #BQU2003-07420-C05-01 as well as from a visiting fellowship granted by the Research Council of Ghent University. The authors are grateful to Prof. Dr. Jan P. Tollenaere and Dr. Wilfried Langenaeker for much appreciated, helpful comments.

---

## REFERENCES

1. *The New Lexicon, Webster's Dictionary of the English Language*, Lexicon Publications, Inc., New York, 1991.
2. A. Tversky, *Psychol. Rev.*, **84**, 327 (1997). Features of Similarity.
3. D. H. Rouvray, in *Concepts and Applications of Molecular Similarity*, M. A. Johnson and G. M. Maggiora, Eds., Wiley-Interscience, New York, 1990, pp. 15–43. The Evolution of the Concept of Molecular Similarity.
4. D. H. Rouvray, *Top. Curr. Chem.*, **173**, 1 (1995). Similarity in Chemistry—Past, Present and Future.
5. R. Carbó, and B. Calabuig, in *Concepts and Applications of Molecular Similarity*, M. A. Johnson and G. M. Maggiora, Eds., Wiley-Interscience, New York, 1990, pp. 147–171. Molecular Similarity and Quantum Chemistry.
6. E. Besalú, R. Carbó, J. Mestres, and M. Solà, *Top. Curr. Chem.*, **173**, 31 (1995). Foundations and Recent Developments on Molecular Quantum Similarity.
7. R. Carbó-Dorca and E. Besalú, *J. Mol. Struct. (THEOCHEM)*, **451**, 11 (1998). A General Survey of Molecular Quantum Similarity.
8. R. Carbó-Dorca, Ll. Amat, E. Besalú, and M. Lobato, in *Advances in Molecular Similarity*, Vol. 2, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1998, pp. 1–42. Quantum Similarity.
9. R. Carbó-Dorca, Ll. Amat, E. Besalú, X. Gironés, and D. Robert, *J. Mol. Struct. (THEOCHEM)*, **504**, 181 (1999). Quantum Mechanical Origin of QSAR: Theory and Applications.
10. R. Carbó-Dorca, D. Robert, Ll. Amat, X. Gironés, and E. Besalú, *Lecture Notes in Chemistry*, **73**, 1 (2000). Molecular Quantum Similarity in QSAR and Drug Design.
11. E. Besalú, X. Girones, Ll. Amat, and R. Carbó-Dorca, *Acc. Chem. Res.*, **35**, 289 (2002). Molecular Quantum Similarity and the Fundamentals of QSAR.
12. R. Carbó-Dorca and X. Gironés, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 364–386. Quantum Similarity and Quantitative Structure-Activity Relationships.
13. P. M. Dean, *Molecular Similarity in Drug Design*, Blackie Academic & Professional, Chapman & Hall, London, 1995.
14. W. C. Herndon and S. H. Bertz, *J. Comput. Chem.*, **8**, 367 (1987). Linear Notations and Molecular Graph Similarity.
15. G. M. Downs, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 364–386. Molecular Descriptors.
16. R. D. Brown, *Perspect. Drug Discov. Design*, **7/8**, 31 (1997). Descriptors for Diversity Analysis.
17. J. S. Mason and S. D. Pickett, *Perspect. Drug Discov. Design*, **7/8**, 85 (1997). Partition-Based Selection.

18. J. Bajorath, *J. Chem. Inf. Comput. Sci.*, **41**, 233 (2001). Selected Concepts and Investigations in Compound Classification, Molecular Descriptor Analysis, and Virtual Screening.
19. M. Karelson, V. S. Lobanov, and A. R. Katritzky, *Chem. Rev.*, **96**, 1027 (1996). Quantum-Chemical Descriptors in QSAR/QSPR Studies.
20. M. Karelson, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 641–667. Quantum-Chemical Descriptors in QSAR.
21. C. A. Lipinski, F. Lombardo, B. W. Dominy, and P. J. Feenay, *Adv. Drug Deliv. Rev.*, **23**, 3 (1997). Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings.
22. R. Mannhold and H. van de Waterbeemd, *J. Comp.-Aided Mol. Design*, **15**, 337 (2001). Substructure and Whole Molecule Approaches for Calculating logP.
23. L. Chen, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 483–513. Substructure and Maximal Common Substructure Searching.
24. H. Wiener, *J. Am. Chem. Soc.*, **69**, 17 (1947). Structural Determination of Paraffin Boiling Point.
25. A. T. Balaban, Ed., *Chemical Applications of Graph Theory*, Academic Press, London, 1976.
26. M. Randic and C. L. Wilkins, *J. Chem. Inf. Comput. Sci.*, **19**, 31 (1979). Graph Theoretical Approach to Recognition of Structural Similarity in Molecules.
27. G. Ruecker and C. Ruecker, *J. Chem. Inf. Comput. Sci.*, **33**, 683 (1993). Counts of All Walks as Atomic and Molecular Descriptors.
28. M. Randic, in *Concepts and Applications of Molecular Similarity*, M. A. Johnson and G. M. Maggiora, Eds., Wiley-Interscience, New York, 1990, pp. 77–145. Design of Molecules With Desired Properties.
29. M. Randic, *J. Am. Chem. Soc.*, **97**, 6609 (1975). On Characterization of Molecular Branching.
30. L. B. Kier and L. H. Hall, *Molecular Connectivity in Structure-Activity Analysis*, Wiley, New York, 1986.
31. I. Gutman and N. Trinajstić, *Chem. Phys. Lett.*, **17**, 535 (1972). Graph Theory and Molecular-Orbitals—Total Pi-Electron Energy of Alternant Hydrocarbons.
32. S. Nikolic, G. Kovacevic, A. Milicevic, and N. Trinajstić, *Croat. Chem. Acta*, **76**, 113 (2003). The Zagreb Indices 30 Years After.
33. H. Hosoya, *Bull. Chem. Soc. Japan*, **44**, 2332 (1971). Topological Index: A Newly Proposed Quantity Characterizing the Topological Nature of Structural Isomers of Saturated Hydrocarbons.
34. T. I. Oprea, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 571–616. 3D QSAR Modeling in Drug Design.
35. S. M. Bachrach, in *Reviews in Computational Chemistry*, Vol. 5, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley-VCH, New York, 1994, pp. 171–227. Population Analysis and Electron Densities from Quantum Mechanics.
36. F. Jensen, *Introduction to Computational Chemistry*, Wiley, New York, 1999.
37. P. Geerlings, F. De Proft, and W. Langenaeker, *Chem. Rev.*, **103**, 1793 (2003). Conceptual Density Functional Theory.
38. A. Szabo and N. S. Ostlund, *Modern Quantum Chemistry-Introduction to Advanced Electronic Structure Theory*, Dover Publications, New York, 1996.
39. I. N. Levine, *Quantum Chemistry*, Prentice-Hall, New York, 1999.
40. F. L. Pilar, *Elementary Quantum Chemistry*, Dover Publications, New York, 2001.
41. J. J. P. Stewart, in *Reviews in Computational Chemistry*, Vol. 1, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 1990, pp. 45–81. Semiempirical Molecular Orbital Methods.

42. T. Bredow, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 29–55. Semi-empirical Methods.
43. T. Helgaker, P. Jorgensen, J. Olsen, and W. Klopper, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 57–88. Wave-Function Based Quantum Chemistry.
44. T. Helgaker, P. Jorgensen, and J. Olsen, *Molecular Electronic-Structure Theory*, Wiley, Chichester, United Kingdom, 2000.
45. N. C. Handy, in *Lecture Notes in Quantum Chemistry II*, B. O. Roos, Ed., Springer, Heidelberg, Germany, 1994, pp. 91–124. Density Functional Theory.
46. M. Born, *Atomic Physics*, Blackie and Son, London, 1945.
47. J. von Neumann, *Mathematical Foundations of Quantum Mechanics*, Princeton University Press, Princeton, New Jersey, 1955.
48. R. G. Parr and W. Yang, *Density Functional Theory of Atoms and Molecules*, Oxford University Press, Oxford, United Kingdom, 1989.
49. P. Hohenberg and W. Kohn, *Phys. Rev.*, **136B**, 864 (1964). Inhomogeneous Electron Gas.
50. K. D. Sen, E. Besalú, and R. Carbó-Dorca, *J. Math. Chem.*, **25**, 253 (1999). A Naive Look on the Hohenberg-Kohn Theorem.
51. R. Carbó-Dorca and P. Bultinck, *J. Math. Chem.*, **34**, 75 (2003). Analysis of a General Theorem Concerning Two Non-Commuting Hermitian Matrices: Quantum Mechanical Implications for Ground and Excited States.
52. P. M. Dean, in *Concepts and Applications of Molecular Similarity*, M. A. Johnson and G. M. Maggiora, Wiley-Interscience, New York, 1990. Molecular Recognition: The Measurement and Search for Molecular Similarity in Ligand-Receptor Interaction.
53. R. Carbó, J. Arnau, and L. Leyda, *Int. J. Quantum Chem.*, **17**, 1185 (1980). How Similar Is a Molecule to Another? An Electron Density Measure of Similarity Between Two Molecular Structures.
54. I. M. Vinogradov, *Encyclopaedia of Mathematics*, Kluwer Academic Publishers, Dordrecht, the Netherlands, 1989.
55. E. W. Weisstein, *CRC Concise Encyclopedia of Mathematics (CD-ROM)*, Chapman & Hall/CRC, London, 2003. The contents of this Encyclopedia may be consulted through the internet at <http://mathworld.wolfram.com>.
56. R. Carbó-Dorca, E. Besalú, and X. Gironés, *Advances in Quantum Chemistry*, **38**, 1 (2000). Extended Density Functions.
57. X. Gironés, A. Gallegos, and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **40**, 1400 (2000). Modeling Antimalarial Activity: Application of Kinetic Energy Density Quantum Similarity Measures as Descriptors in QSAR.
58. R. Carbó-Dorca, B. Calabuig, E. Besalú, and A. Martínez, *Molec. Eng.*, **2**, 43 (1992). Triple Density Molecular Quantum Similarity Measures: A General Connection Between Theoretical Calculations and Experimental Results.
59. D. Robert and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **38**, 620 (1998). Analyzing the Triple Density Molecular Quantum Similarity Measures With the INDSCAL Model.
60. D. Robert, X. Gironés, and R. Carbó-Dorca, *Polycyc. Arom. Comp.*, **19**, 51 (2000). Molecular Quantum Similarity Measures as Descriptors for Quantum QSAR.
61. X. Gironés, Ll. Amat, D. Robert, and R. Carbó-Dorca, *J. Comp.-Aided Molec. Des.*, **14**, 477 (2000). Use of Electron-Electron Repulsion Energy as a Molecular Descriptor in QSAR and QSPR Studies.
62. R. Carbó and B. Calabuig, *Int. J. Quant. Chem.*, **42**, 1695 (1992). Molecular Quantum Similarity Measures and N-Dimensional Representation of Quantum Objects. II. Practical Applications.

63. R. Carbó and B. Calabuig, *J. Chem. Inf. Comput. Sci.*, **32**, 600 (1992). Quantum Similarity Measures, Molecular Cloud Description and Structure-Property Relationships.
64. R. Carbó and B. Calabuig, *J. Mol. Struct. (THEOCHEM)*, **254**, 517 (1992). Quantum Molecular Similarity Measures and the N-Dimensional Representation of a Molecular Set: Phenylidimethyldiazines.
65. R. Carbó, B. Calabuig, L. Vera, and E. Besalú, *Adv. Quant. Chem.*, **25**, 253 (1994). Molecular Quantum Similarity: Theoretical Framework, Ordering Principles and Visualization Techniques.
66. R. Carbó and E. Besalú, in *Molecular Similarity and Reactivity: From Quantum Chemical to Phenomenological Approaches*, R. Carbó, Ed., Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995, pp. 3–30. Theoretical Foundation of Quantum Molecular Similarity.
67. R. Carbó-Dorca, *J. Math. Chem.*, **27**, 357 (2000). Quantum QSAR and the Eigensystems of Stochastic Similarity Matrices.
68. R. Carbó-Dorca, *Int. J. Quant. Chem.*, **79**, 163 (2000). Stochastic Transformation of Quantum Similarity Matrices and Their Use in Quantum QSAR (QQSAR) Models.
69. X. Gironés and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **42**, 317 (2002). Using Molecular Quantum Similarity Matrices to Describe Physical Properties of Molecular Systems.
70. R. Carbó and B. Calabuig, in *Computational Chemistry: Structure, Interactions and Reactivity*, S. Fraga, Ed., Elsevier, Amsterdam, 1992, pp. 300–324. Quantum Similarity: Definitions, Computational Details and Applications.
71. P. Bultinck and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **43**, 170 (2003). Molecular Quantum Similarity Matrix Based Clustering of Molecules Using Dendrograms.
72. C. J. Barden and H. F. Schaefer III, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 133–149. Accuracy and Applicability of Quantum Chemical Methods in Computational Medicinal Chemistry.
73. W. Koch and M. C. Holthausen, *A Chemist's Guide to Density Functional Theory*, Wiley-VCH, Weinheim, Germany, 2002.
74. P. W. Ayers and W. Yang, in *Computational Medicinal Chemistry for Drug Discovery*, P. Bultinck, H. De Winter, W. Langenaeker, and J. P. Tollenaere, Eds., Dekker, Inc., New York, 2003, pp. 89–188. Density-Functional Theory.
75. C. Van Alsenoy, *J. Comput. Chem.*, **9**, 620 (1988). Ab Initio Calculations on Large Molecules: The Multiplicative Integral Approximation.
76. C. Van Alsenoy and A. Peeters, *J. Mol. Struct. (THEOCHEM)*, **286**, 19 (1993). BRABO: A Program for Ab Initio Studies on Large Molecular Systems.
77. F. L. Hirshfeld, *Theoret. Chim. Acta*, **44**, 129 (1977). Bonded-Atom Fragments for Describing Molecular Charge Densities.
78. G. G. Hall and D. Martin, *Isr. J. Chem.*, **19**, 255 (1980). Approximate Electron Densities for Atoms and Molecules.
79. G. G. Hall and C. M. Smith, *Int. J. Quant. Chem.*, **25**, 881 (1984). Fitting Electron Densities of Molecules.
80. C. M. Smith and G. G. Hall, *Theor. Chim. Acta*, **69**, 63 (1986). The Approximation of Electron Densities.
81. G. G. Hall and C. M. Smith, *J. Mol. Struct.*, **179**, 293 (1988). Electric Fields around Molecules.
82. E. J. Baerends, D. E. Ellis, and P. Ros, *Chem. Phys.*, **2**, 41 (1973). Self-Consistent Molecular Hartree-Fock-Slater Calculations. I. The Computational Procedure.
83. H. Sambe and R. H. Felton, *J. Chem. Phys.*, **62**, 1122 (1975). A New Computational Approach to Slater's  $X_\alpha$  Equation.
84. B. I. Dunlap, J. W. Connolly, and J. R. Sabin, *J. Chem. Phys.*, **71**, 3396 (1979). On Some Approximations in Applications of  $X_\alpha$  Theory.

85. J. Mestres, M. Solà, M. Duran, and R. Carbó, *J. Comput. Chem.*, **15**, 1113 (1994). On the Calculation of Ab Initio Quantum Molecular Similarities for Large Systems: Fitting the Electron Density.
86. J. Mestres, M. Solà, M. Duran, and R. Carbó, in *Molecular Similarity and Reactivity: From Quantum Chemical to Phenomenological Approaches*, R. Carbó, Ed., Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995, pp. 89–111. General Suggestions and Applications of Quantum Molecular Similarity Measures from Ab Initio Fitted Electron Densities.
87. P. Constans and R. Carbó, *J. Chem. Inf. Comput. Sci.*, **35**, 1046 (1995). Atomic Shell Approximation: Electron Density Fitting Algorithm Restricting Coefficients to Positive Values.
88. P. Constans, Ll. Amat, X. Fradera, and R. Carbó-Dorca, in *Advances in Molecular Similarity*, Vol. 1, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1996, pp. 187–211. Quantum Molecular Similarity Measures (QMSM) and the Atomic Shell Approximation (ASA).
89. Ll. Amat and R. Carbó-Dorca, *J. Comput. Chem.*, **18**, 2023 (1997). Quantum Similarity Measures Under Atomic Shell Approximation: First Order Density Fitting Using Elementary Jacobi Rotations.
90. Ll. Amat and R. Carbó-Dorca, *J. Comput. Chem.*, **20**, 911 (1999). Fitted Electron Density Functions from H to Rn for Use in Quantum Similarity Measures: Cis-Diammine, Dichloroplatinum(II) Complexes as an Application Example.
91. Ll. Amat and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **40**, 1118 (2000). Molecular Electron Density Fitting Using Elementary Jacobi Rotations under Atomic Shell Approximation.
92. X. Gironés, Ll. Amat, and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **42**, 847 (2002). Modeling Large Macromolecular Structures Using Promolecular Densities.
93. Ll. Amat and R. Carbó-Dorca, *Int. J. Quant. Chem.*, **87**, 59 (2002). Use of Promolecular ASA Density Functions to Obtain Starting MO in SCF Calculations.
94. P. Bultinck (Ghent University), Unpublished Work, 2002.
95. ASA coefficients and exponents. Available: <http://iqc.udg.es/cat/similarity/ASA/funcset.html>.
96. X. Gironés, Ll. Amat, and R. Carbó-Dorca, *J. Mol. Graph. and Model.*, **16**, 190 (1998). A Comparison of Isodensity Surfaces Using Ab Initio and ASA Density Functions.
97. X. Gironés, R. Carbó-Dorca, and P. G. Mezey, *J. Mol. Graph. and Model.*, **19**, 343 (2001). Application of Promolecular ASA Densities to Graphical Representation of Density Functions of Macromolecular Systems.
98. P. Bultinck, R. Carbó-Dorca, and C. Van Alsenoy, *J. Chem. Inf. Comput. Sci.*, **43**, 1208 (2003). Quality of Approximate Electron Densities and Internal Consistency of Molecular Alignment Algorithms in Molecular Quantum Similarity.
99. P. Bultinck (Ghent University), Unpublished Work, 2003.
100. S. Van Damme (Ghent University), *Graduation Thesis*, Ghent Quantum Chemistry Group, 2003.
101. R. F. Stewart, E. R. Davidson, and W. J. Simpson, *J. Chem. Phys.*, **42**, 3175 (1965). Coherent X-Ray Scattering for the Hydrogen Atom in the Hydrogen Molecule.
102. R. F. Stewart, *Isr. J. Chem.*, **16**, 124 (1977). One-Electron Density Functions and Many-Centered Finite Multipole Expansions.
103. T. R. Cundari, M. T. Benson, M. Leigh Lutz, and S. O. Sommerer, in *Reviews in Computational Chemistry*, Vol. 8, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 1995, pp. 145–202. Effective Core Potential Approaches to the Chemistry of the Heavier Elements.
104. P. J. Hay and W. R. Wadt, *J. Chem. Phys.*, **82**, 299 (1985). Ab Initio Effective Core Potentials for Molecular Calculations. Potentials for K to Au Including the Outermost Orbitals.

105. G. Boon, F. De Proft, W. Langenaeker, and P. Geerlings, *Chem. Phys. Lett.*, **295**, 122 (1998). The Use of Density Functional Theory-Based Reactivity Descriptors in Molecular Similarity Calculations.
106. G. Boon, W. Langenaeker, F. De Proft, H. De Winter, J. P. Tollenaere, and P. Geerlings, *J. Phys. Chem. A*, **105**, 8805 (2001). Systematic Study of the Quality of Various Quantum Similarity Descriptors. Use of the Autocorrelation Function and Principal Component Analysis.
107. G. Boon, C. Van Alsenoy, F. De Proft, P. Bultinck, and P. Geerlings, *J. Phys. Chem. A*, **107**, 11120 (2003). Similarity and Chirality: Quantum Chemical Study of Dissimilarity of Enantiomers.
108. D. L. Cooper and N. L. Allan, in *Molecular Similarity and Reactivity: From Quantum Chemical to Phenomenological Approaches*, R. Carbó, Ed., Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995, pp. 31–55. Molecular Similarity and Momentum Space.
109. B. H. Bransden and C. J. Joachain, *Introduction to Quantum Mechanics*, Longman Scientific and Technical, Harlow, United Kingdom, 1995.
110. F. Rioux, *J. Chem. Educ.*, **74**, 605 (1997). Numerical Methods for Finding Momentum Space Distributions.
111. P. Kaijser and V. H. Smith, Jr., *Adv. Quant. Chem.*, **10**, 37 (1977). Evaluation of Momentum Distributions and Compton Profiles for Atomic and Molecular Systems.
112. R. Ponec, *Collect. Czech. Chem. Commun.*, **52**, 555 (1987). Topological Aspects of Chemical Reactivity. On the Similarity of Molecular Structures.
113. R. Ponec and M. Strnad, *Collect. Czech. Chem. Commun.*, **55**, 2363 (1990). Topological Aspects of Chemical Reactivity. On Physical Meaning of the Similarity Index.
114. R. Ponec and M. Strnad, *Int. J. Quant. Chem.*, **42**, 501 (1992). Electron Correlation in Pericyclic Reactivity: A Similarity Approach.
115. R. Ponec and M. Strnad, *J. Phys. Org. Chem.*, **5**, 764 (1992). Structure of Transition States in Forbidden Pericyclic Reactions. The Second-Order Similarity Approach.
116. M. Strnad and R. Ponec, *Int. J. Quant. Chem.*, **49**, 35 (1994). Novel Approach to Molecular Similarity: Second-Order Similarity Indices from Geminal Expansion of Pair-Densities.
117. R. Ponec, *Top. Curr. Chem.*, **174**, 1 (1995). Similarity Models in the Theory of Pericyclic Reaction.
118. R. Ponec, in *Advances in Molecular Similarity*, Vol. 1, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1996, pp. 121–133. Electron Correlation in Allowed and Forbidden Pericyclic Reactions from Geminal Expansion of Pair Densities.
119. R. Ponec and M. Strnad, *Collect. Czech. Chem. Commun.*, **55**, 2583 (1990). Similarity Approach to Chemical Reactivity. Specificity of Multibond Processes.
120. R. Carbó, F. Lapeña, and E. Suñé, *Afinidad*, **43**, 483 (1986). Similarity Measures on Electrostatic Molecular Potentials.
121. P. M. W. Gill, B. G. Johnson, J. A. Pople, and S. W. Taylor, *J. Chem. Phys.*, **96**, 7178 (1992). Modeling the Potential of a Charge Distribution.
122. N. Meurice, L. Leherter, and D. P. Vercauteren, *SAR & QSAR in Environ. Res.*, **8**, 195 (1998). Comparison of Benzodiazepine-Like Compounds Using Topological Analysis and Genetic Algorithms.
123. L. Leherter, N. Meurice, and D. P. Vercauteren, *J. Chem. Inf. Comput. Sci.*, **40**, 816 (2000). Critical Point Analysis of Electron Density Maps for the Comparison of Benzodiazepine-Type Ligands.
124. L. Leherter, *J. Math. Chem.*, **29**, 47 (2001). Application of Multiresolution Analyses to Electron Density Maps of Small Molecules: Critical Point Representations for Molecular Superposition.
125. L. Leherter, L. Dury, and D. P. Vercauteren, *J. Phys. Chem.*, **107**, 9875 (2003). Structural Identification of Local Maxima in Low-Resolution Promolecular Electron Density Distributions.



126. L. Leherter and D. P. Vercauteren, *J. Molec. Model.*, **3**, 156 (1995). Critical Point Analysis of Calculated Electron Density Maps at Medium Resolution: Application to Shape Analysis of Zeolite-Like Systems.
127. L. Leherter and F. H. Allen, *J. Comput.-Aided Mol. Des.*, **8**, 257 (1994). Shape Information from a Critical Point Analysis of Calculated Electron Density Maps; Application to DNA-Drug Systems.
128. L. Leherter (University of Namur, Belgium), Personal Communication, 2003.
129. P. Constans, Ll. Amat, and R. Carbó-Dorca, *J. Comput. Chem.*, **18**, 826 (1997). Toward a Global Maximization of the Molecular Similarity Function: Superposition of Two Molecules.
130. P. Bultinck, T. Kuppens, X. Gironés, and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **43**, 1143 (2003). QSSA: A Consistent Scheme for Molecular Alignment and Molecular Similarity Based on Quantum Chemistry.
131. B. B. Stefanov and J. Cioslowski, in *Advances in Molecular Similarity*, Vol. 1, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1996, pp. 43–59. Similarity of Atoms in Molecules.
132. J. Cioslowski, B. Stefanov, and P. Constans, *J. Comput. Chem.*, **17**, 1352 (1996). Efficient Algorithm for Quantitative Assessment of Similarities Among Atoms in Molecules.
133. F. H. Walters, L. R. Parker, Jr., S. L. Morgan, and S. N. Deming, *Sequential Simplex Optimization*, CRC Press, New York, 1991. An electronic reprint of this book may be obtained from [www.multisimplex.com](http://www.multisimplex.com).
134. W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes in Fortran*, Cambridge University Press, Cambridge, United Kingdom, 1992.
135. R. Judson, in *Reviews in Computational Chemistry*, Vol. 10, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 1997, pp. 1–73. Genetic Algorithms and Their Use in Chemistry.
136. A. J. McMahon and P. J. King, *J. Comput. Chem.*, **18**, 151 (1997). Optimization of the Carbó Molecular Similarity Using Gradient Methods.
137. M. F. Parretti, R. T. Kroemer, J. H. Rothman, and W. G. Richards, *J. Comput. Chem.*, **18**, 1344 (1997). Alignment of Molecules by the Monte Carlo Optimization of Molecular Similarity Indices.
138. X. Gironés, D. Robert, and R. Carbó-Dorca, *J. Comput. Chem.*, **22**, 255 (2001). TGSA: A Molecular Superposition Program Based on Topo-Geometrical Considerations.
139. X. Gironés and R. Carbó-Dorca, *J. Comput. Chem.*, **25**, 153 (2004). TGSA-Flex: Extending the Capabilities of the Topo-Geometrical Superposition Algorithm to Handle Rotary Bonds.
140. C. Lemmen and T. Lengauer, *J. Comp.-Aided Mol. Des.*, **14**, 215 (2000). Computational Methods for the Structural Alignment of Molecules.
141. P. Bultinck (Ghent University), Unpublished Work, 2002.
142. P. Willett, J. M. Barnard, and G. M. Downs, *J. Chem. Inf. Comput. Sci.*, **38**, 983 (1998). Chemical Similarity Searching.
143. E. E. Hodgkin and W. G. Richards, *Int. J. Quant. Chem. Quantum. Biol. Symp.*, **14**, 105 (1987). Molecular Similarity Based on Electrostatic Potential and Electric Field.
144. E. E. Hodgkin and W. G. Richards, *Chem. Ber.*, **24**, 1141 (1988). Molecular Similarity.
145. J. D. Petke, *J. Comput. Chem.*, **14**, 928 (1993). Cumulative and Discrete Similarity Analysis of Electrostatic Potentials and Fields.
146. J. T. Tou and R. C. González, *Pattern Recognition Principles*, Addison-Wesley, Reading, Massachusetts, 1974.
147. G. M. Maggiora, J. D. Petke, and J. Mestres, *J. Math. Chem.*, **31**, 251 (2002). A General Analysis of Field-Based Molecular Similarity Indices.
148. R. Carbó, E. Besalú, Ll. Amat, and X. Fradera, *J. Math. Chem.*, **19**, 47 (1996). On Quantum Molecular Similarity Measures (QMSM) and Indices (QMSI).

149. R. Carbó-Dorca, E. Besalú, Ll. Amat, and X. Fradera, in *Advances in Molecular Similarity*, Vol. 1, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1996, pp. 1–42. Quantum Molecular Similarity Measures: Concepts, Definitions, and Applications to Quantitative Structure-Property Relationships.
150. D. Robert and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **38**, 469 (1998). A Formal Comparison Between Molecular Quantum Similarity Measures and Indices.
151. R. G. Parr and L. J. Bartolotti, *J. Phys. Chem.*, **87**, 2810 (1983). Some Remarks on the Density Functional Theory of Few-Electron Systems.
152. P. Bultinck and R. Carbó-Dorca, *J. Math. Chem.*, **36**, 191 (2004). A Mathematical Discussion on Density and Shape Functions, Vector Semispaces and Related Questions.
153. R. F. W. Bader, *Atoms in Molecules: A Quantum Theory*, Clarendon Press, Oxford, United Kingdom, 1990.
154. J. Cioslowski and A. Nanayakkara, *J. Am. Chem. Soc.*, **115**, 11213 (1993). Similarity of Atoms in Molecules.
155. P. L. A. Popelier, *J. Phys. Chem. A*, **103**, 2883 (1999). Quantum Molecular Similarity. 1. BCP Space.
156. S. E. O'Brien and P. L. A. Popelier, *Can. J. Chem.*, **77**, 28 (1999). Quantum Molecular Similarity. Part 2: The Relation Between Properties in BCP Space and Bond Length.
157. S. E. O'Brien and P. L. A. Popelier, *J. Chem. Inf. Comput. Sci.*, **41**, 764 (2001). Quantum Molecular Similarity. 3. QTMS Descriptors.
158. S. E. O'Brien and P. L. A. Popelier, *J. Chem. Soc. Perk. Trans.*, **2**, 478 (2002). Topological Molecular Similarity. Part 4. A QSAR Study of Cell Growth Inhibitory Properties of Substituted (E)-1-Phenylbut-1-en-3-ones.
159. P. L. A. Popelier, U. A. Chaudry, and P. J. Smith, *J. Chem. Soc. Perk. Trans.*, **2**, 1231 (2002). Quantum Topological Molecular Similarity. Part 5. Further Development With An Application to the Toxicity of Polychlorinated Dibenzo-P-Dioxins (PCDDS).
160. P. G. Mezey, R. Ponec, Ll. Amat, and R. Carbó-Dorca, *Enantiomer*, **4**, 371 (1999). Quantum Similarity Approach to the Characterization of Molecular Chirality.
161. R. Carbó and B. Calabuig, *Int. J. Quant. Chem.*, **42**, 1681 (1992). Molecular Quantum Similarity Measures and N-Dimensional Representation of Quantum Objects. I. Theoretical Foundations.
162. T. Fujita, J. Iwasa, and C. Hansch, *J. Am. Chem. Soc.*, **86**, 5175 (1964). A New Substituent Constant,  $\Pi$ , Derived from Partition Coefficients.
163. R. F. Rekker, *Hydrophobic Fragmental Constants. Its Derivation and Applications. A Means of Characterizing Membrane Systems*, Elsevier, New York, 1977.
164. R. F. Rekker and H. M. Kort, *Eur. J. Med. Chem.*, **14**, 479 (1979). Hydrophobic Fragmental Constant—Extension to a 1000 Data Point Set.
165. C. Hansch and A. Leo, *Substituent Constants for Correlation Analysis in Chemistry and Biology*, Wiley, New York, 1979.
166. G. Klopman and L. D. Iroff, *J. Comput. Chem.*, **2**, 157 (1986). Calculation of Partition-Coefficients by the Charge-Density Method.
167. G. Klopman, K. Namboodiri, and M. Schochet, *J. Comput. Chem.*, **1**, 28 (1985). Simple Method of Computing the Partition-Coefficient.
168. A. K. Ghose and G. M. Crippen, *J. Comput. Chem.*, **7**, 565 (1986). Atomic Physicochemical Parameters for 3-Dimensional Structure-Directed Quantitative Structure-Activity-Relationships. 1. Partition-Coefficients as a Measure of Hydrophobicity.
169. I. Moriguchi, S. Hirono, I. Nakagome, and H. Hirano, *Chem. Pharm. Bull.*, **42**, 976 (1994). Comparison of Reliability of LogP Values for Drugs Calculated by Several Methods.
170. Ll. Amat, R. Carbó-Dorca, and R. Ponec, *J. Comput. Chem.*, **14**, 1575 (1998). Molecular Quantum Similarity Measures as an Alternative to LogP Values in QSAR Studies.

171. X. Gironés, Ll. Amat, and R. Carbó-Dorca, *SAR and QSAR in Environ. Res.*, **10**, 545 (1999). Using Molecular Quantum Similarity Measures as Descriptors in Quantitative Structure-Toxicity Relationships.
172. R. Ponec, Ll. Amat, and R. Carbó-Dorca, *J. Comput.-Aided Mol. Des.*, **13**, 259 (1999). Molecular Basis of Quantitative Structure-Properties Relationships (QSPR): A Quantum Similarity Approach.
173. R. Ponec, Ll. Amat, and R. Carbó-Dorca, *J. Phys. Org. Chem.*, **12**, 447 (1999). Quantum Similarity Approach to LFER: Substituent and Solvent Effects on the Acidities of Carboxylic Acids.
174. Ll. Amat, R. Carbó-Dorca, and R. Ponec, *J. Med. Chem.*, **42**, 5169 (1999). Simple Linear QSAR Models Based on Quantum Similarity Measures.
175. X. Gironés, R. Carbó-Dorca, and R. Ponec, *J. Chem. Inf. Comput. Sci.*, **43**, 2033 (2003). Molecular Basis of LFER. Modeling of the Electronic Substituent Effect Using Fragment Quantum Self-Similarity Measures.
176. J. J. Sullivan, A. D. Jones, and K. K. Tanji, *J. Chem. Inf. Comput. Sci.*, **40**, 1113 (2000). QSAR Treatment of Electronic Substituent Effects Using Frontier Orbital Theory and Topological Parameters.
177. M. Solà, J. Mestres, J. M. Oliva, M. Duran, and R. Carbó, *Int. J. Quant. Chem.*, **58**, 361 (1996). The Use of Ab Initio Quantum Molecular Self-Similarity Measures to Analyze Electron Charge Density Distributions.
178. J. M. Oliva, R. Carbó-Dorca, and J. Mestres, in *Advances in Molecular Similarity*, Vol. 1, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1996, pp. 135–165. Conformational Analysis from the Viewpoint of Molecular Similarity.
179. P. G. Mezey, *Mol. Phys.*, **96**, 169 (1999). The Holographic Electron Density Theorem and Quantum Similarity Measures.
180. R. Carbó, M. Martín, and V. Pons, *Afinidad*, **34**, 348 (1977). Application of Quantum Mechanical Parameters in Quantitative Structure-Activity Relationships.
181. R. Carbó-Dorca, in *Advances in Molecular Similarity*, Vol. 2, R. Carbó-Dorca, and P. G. Mezey, Eds., JAI Press, London, 1998, pp. 43–72. Fuzzy Sets and Boolean Tagged Sets; Vector Semispaces and Convex Sets; Quantum Similarity Measures and ASA Density Functions; Diagonal Vector Spaces and Quantum Chemistry.
182. R. Carbó-Dorca, Ll. Amat, E. Besalú, X. Gironés, and D. Robert, in *Fundamentals of Molecular Similarity*, R. Carbó-Dorca, X. Gironés, and P. G. Mezey, Eds., Kluwer Academic /Plenum Press, New York, 2001, pp. 187–320. Quantum Molecular Similarity: Theory and Applications to the Evaluation of Molecular Properties, Biological Activities and Toxicity.
183. R. Carbó-Dorca and E. Besalú, *Int. J. Quantum Chem.*, **88**, 167 (2002). Fundamental Quantum QSAR (Q<sup>2</sup>SAR) Equation: Extensions, Non-Linear Terms and Generalizations Within Extended Hilbert-Sobolev Spaces.
184. R. Carbó-Dorca, in *Proceedings of the European Congress on Computational Methods in Applied Sciences and Engineering (ECCOMAS 2000)*, CDROM edited by Facultat d'Informàtica de Barcelona (FIB)—Universitat Politècnica de Catalunya (UPC)—International Centre for Numerical Methods in Engineering (CIMNE): Barcelona, 2000, Computational Chemistry Section, Chapter 12. Quantum Quantitative Structure-Activity Relationships (QQSAR): A Comprehensive Discussion Based on Inward Matrix Products, Employed as a Tool to Find Approximate Solutions of Strictly Positive Linear Systems and Providing QSAR-Quantum Similarity Measures Connections.
185. R. Carbó, E. Besalú, Ll. Amat, and X. Fradera, *J. Math. Chem.*, **18**, 237 (1995). Quantum Molecular Similarity Measures (QMSM) as a Natural Way Leading Towards a Theoretical Foundation of Quantitative Structure-Activity Relationships (QSAR).
186. R. Carbó-Dorca, *J. Mol. Struct. (THEOCHEM)*, **537**, 41 (2001). Inward Matrix Products: Extensions and Applications to Quantum Mechanical Foundations of QSAR.

187. R. Carbó-Dorca, *J. Math. Chem.*, **23**, 353 (1998). Tagged Sets, Convex Sets and QS Measures.
188. H. Eyring, H. J. Walter, and G. E. Kimball, *Quantum Chemistry*, Wiley, New York, 1944.
189. P. O. Löwdin, *Phys. Rev.*, **97**, 1474 (1955). Quantum Theory of Many-Particle Systems. I. Physical Interpretations By Means of Density Matrices, Natural Spin-Orbitals, and Convergence Problems in the Method of Configuration Interaction.
190. P. O. Löwdin, *Phys. Rev.*, **97**, 1490 (1955). Quantum Theory of Many-Particle Systems. II. Study of the Ordinary Hartree-Fock Approximation.
191. P. O. Löwdin, *Phys. Rev.*, **97**, 1509 (1955). Quantum Theory of Many-Particle Systems. III. Extension of the Hartree-Fock Scheme to Include Degenerate Systems and Correlation Effects.
192. R. McWeeny, *Rev. Mod. Phys.*, **32**, 335 (1960). Some Recent Advances in Density Matrix Theory.
193. R. E. Moss, *Advanced Molecular Quantum Mechanics*, Chapman and Hall, London, 1973.
194. P. A. M. Dirac, *Principles of Quantum Mechanics*, Clarendon Press, Oxford, United Kingdom, 1983.
195. R. McWeeny and B. T. Sutcliffe, *Methods of Molecular Quantum Mechanics*, Academic Press, London, 1969.
196. R. Carbó-Dorca, *Institute of Computational Chemistry Technical Report*, IT-IQC-00-03, 2000. Mathematical Elements of Quantum Electronic Density Functions.
197. LF 95 Language Reference, Lahey Computer Systems, Incline Village, Nevada, 1998. Available: <http://www.lahey.com>.
198. R. Carbó and E. Besalú, in *Strategies and Applications in Quantum Chemistry*, Y. Ellinger, and M. Defranceschi, Eds., Kluwer Academic Publishers, Dordrecht, The Netherlands, 1996, p. 229. Applications of Nested Summation Symbols to Quantum Chemistry: Formalism and Programming Techniques.
199. R. Carbó and E. Besalú, *J. Math. Chem.*, **18**, 37 (1995). Definition and Quantum Chemical Applications of Nested Summation Symbols and Logical Functions: Pedagogical Artificial Intelligence Devices for Formulae Writing, Sequential Programming and Automatic Parallel Implementation.
200. J. Mestres, M. Solà, E. Besalú, M. Duran, and R. Carbó, in *Molecular Similarity and Reactivity: From Quantum Chemical to Phenomenological Approaches*, R. Carbó, Ed., Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995, pp. 77–85. Electron Density Approximations for the Fast Evaluation of Quantum Molecular Similarity Measures.
201. R. Carbó-Dorca and X. Gironés, *Institute of Computational Chemistry Technical Report*, IT-IQC-02-17, 2002. Brief Theoretical Description, With Appropriate Application Examples, of Density Functions Structure and Approximations, Leading to the Foundation of Quantum Similarity Measures and Conducting Towards Quantum Quantitative Structure-Properties Relationships.
202. R. Carbó-Dorca, *Institute of Computational Chemistry Technical Report*, IT-IQC-02-18, 2002. Shell Partition and Metric Semispaces: Minkowski Norms, Root Scalar Products, Distances and Cosines of Arbitrary Order.
203. R. D. Cramer III, D. E. Patterson, and J. D. Bunce, *J. Am. Chem. Soc.*, **10**, 5969 (1988). Comparative Molecular Field Analysis (CoMFA). 1. Effect on Shape on Binding of Steroids to Carrier Proteins.
204. M. Wagener, J. Sadowski, and J. Gasteiger, *J. Am. Chem. Soc.*, **117**, 7769 (1995). Auto-correlation of Molecular Surface Properties for Modeling Corticosteroid Binding Globulin and Cytosolic Ah Receptor Activity by Neural Networks.
205. D. Robert, Ll. Amat, and R. Carbó-Dorca, *J. Chem. Inf. Comput. Sci.*, **39**, 333 (1999). 3D QSAR from Tuned Molecular Quantum Similarity Measures: Prediction of the CBG Binding Affinity for a Steroid Family.

206. M. Lobato, Ll. Amat, E. Besalú, and R. Carbó-Dorca, *Quant. Struct.- Act. Relat.*, **16**, 465 (1997). Structure-Activity Relationships of a Steroid Family Using QSM and Topological QS Indices.
207. M. Solà, J. Mestres, R. Carbó, and M. Duran, *J. Chem. Phys.*, **104**, 1 (1996). A Comparative Analysis By Means of Quantum Molecular Similarity Measures of Density Distributions Derived from Conventional Ab Initio and Density Functional Methods.
208. M. Solà, M. Forés, and M. Duran, in *Advances in Molecular Similarity*, Vol. 2, R. Carbó-Dorca and P. G. Mezey, Eds., JAI Press, London, 1998, pp. 1–42. Optimizing Hybrid Density Functionals by Means of Quantum Molecular Similarity Techniques.
209. J. Mestres, M. Solà, R. Carbó, F. J. Luque, and M. Orozco, *J. Phys. Chem.*, **100**, 606 (1996). Effect of Solvation on the Charge Distribution of a Series of Anionic, Neutral, and Cationic Species. A Quantum Molecular Similarity Study.



## Enumerating Molecules

Jean-Loup Faulon,<sup>\*</sup> Donald P. Visco, Jr.,<sup>†</sup> and  
Diana Roe<sup>‡</sup>

<sup>\*</sup>*Sandia National Laboratories, Computational Biology  
Department, P.O. Box 969, MS 9951, Livermore, CA, 94551*

<sup>†</sup>*Tennessee Technological University, Department of Chemical  
Engineering, Box 5013, Cookeville, TN, 38502*

<sup>‡</sup>*Sandia National Laboratories, Biosystems Research  
Department, P.O. Box 969, MS 9951, Livermore, CA, 94551*

---

---

### ENUMERATING MOLECULES: WHY

Enumerating molecules is a mind-boggling problem that has fascinated chemists and mathematicians alike for more than a century. Taking the definition from various dictionaries, to enumerate means (1) “to name things separately, one by one”, and (2) “to determine the number of, to count.” Interestingly enough, both definitions have been taken when enumerating molecules. Historically, the latter definition was first used by chemists, and mathematical solutions were devised to count molecules. Some of the solutions developed were not only valuable to chemists but to mathematicians as well. Indeed, as we shall see in this chapter, while trying to solve the problem of counting the isomers of paraffin structures<sup>1</sup> or of counting substituted aromatic compounds,<sup>2</sup> important concepts in graph theory and combinatorics were developed. The terms *graph* and *tree* were even coined in a chemistry context.<sup>3</sup>

About four decades ago, with the advance of computer science, researchers started to look at the former definition of enumeration and devised computer codes to explicitly list molecules. Again, while studying this challenging problem, important concepts in computer science were developed. Artificial intelligence textbooks<sup>4</sup> generally quote DENDRAL, a code to enumerate molecules, as the first expert system.

Historically, molecular enumeration has brought a fertile ground of research among chemistry, mathematics, and computer science. Still today new concepts and techniques are being developed at the interstice of these fields.<sup>5</sup>

Enumerating molecules is not only an interesting academic exercise, but it has practical applications as well. The foremost application of enumeration is structure elucidation. Ideally, the wishful bench chemist collects experimental data (NMR, MS, IR, ...) for an unknown compound, the data are fed to a code, and the resulting unique structure is given back. Although such a streamlined picture is not yet fully automated, and may never be, commercial codes are available that can, for instance, list all structures matching a given molecular formula, an infrared (IR) spectrum, or a nuclear magnetic resonance (NMR) spectrum. Another important application is in molecular design. Here the problem is to design compounds (drugs, for example) that optimize some physical, chemical, or biological property or activity. Although not as prolific as structure elucidation, molecular design has introduced some novel stochastic solutions to molecular enumeration. Finally, with the advent of combinatorial chemistry, molecular enumeration takes a central role as it allows computational chemists to construct virtual libraries, test hypotheses, and provide guidance to design optimal combinatorial experiments.

Our primary goal in this chapter is to explain how molecules are enumerated. This explanation is the objective of the first four sections. We start with the problem of counting molecules, then describe how molecules are explicitly enumerated, and finish with a review of stochastic techniques to sample molecules. Our discussion is directed toward structure elucidation and molecular design. However, these applications involve nearly all aspects of counting, enumerating, and sampling. Before understanding how molecules can be elucidated and designed, important theoretical concepts and interesting results relevant to chemistry have to first be assimilated.

The purpose of the final section of this chapter is to review the practical applications of molecular enumeration and to give the reader interested in any of these applications pointers to relevant codes and techniques. In particular, the numbers of isomers for a specific molecular series are given, popular structure elucidation codes are reviewed, computed-aided structure elucidation successes are surveyed, and the connections between structure enumeration and combinatorial library design are established. The field of molecular design with inverse quantitative structure activity relationship is also reviewed. We conclude the chapter outlining future research directions.



Before we start, we want to point out that this chapter is limited to structural [i.e., two-dimensional (2-D)] enumeration and does not cover conformational [i.e., three-dimensional (3-D)] enumeration. This latter topic has already been discussed in the book series for small- and medium-sized molecules<sup>6</sup> and peptides.<sup>7</sup>

---

## ENUMERATING MOLECULES: HOW

The term *enumerating* is found in the literature for both listing molecules one by one and determining the number of molecules corresponding to a given set of constraints. In this chapter, we use the term *counting* for the latter case, and we use the term *enumerating* only when molecules are explicitly listed. Starting with some elementary definitions from graph theory, we then describe how molecules are counted, enumerated, and finally stochastically sampled. The counting, enumerating, and sampling sections can be read separately. Although counting is mostly solved through mathematical treatments, enumerating and sampling are essentially algorithmic problems. In each of the following sections, theoretical results are first explained and illustrated with examples relevant to chemistry. Second, chemical applications are surveyed. To illustrate the problem being studied, a question is attached to each section. The answers to the questions can be found in the text.

### From Graph Theory to Chemistry

We provide here elementary definitions that we later use to count, enumerate, and sample molecules. Rather than a formal mathematical presentation, examples and illustrations are given.

A *simple graph*  $G$  is defined as an ordered pair  $G = (V(G), E(G))$ , where  $V = V(G)$  is a nonempty set of elements called *vertices* and  $E = E(G)$  is a set of unordered pairs of distinct element of  $V$  called *edges*. In most cases of chemical interest, the sets  $V$  and  $E$  are finite. An example of a simple graph is given in Figure 1(a). Of course, a relationship exists between graphs and chemical structures. Sylvester<sup>3</sup> proposed the term *graph* in 1878 on the basis of the structural formulas of molecules. Figure 1(a) can, for instance, be viewed as a representation of cyclohexane. But molecules exist that do not fit the simple graph picture. A *multigraph* is a graph in which the edge set is not necessarily composed of distinct pair of vertices; in other words, *multiple edges* are allowed in a multigraph. A multigraph is without a loop when vertices are not allowed to be paired with themselves. Figure 1(b) is a representation of benzene. In a simple graph or a multigraph, the *degree* of a vertex is the number of edges attached to it, and the *multiplicity* of an edge is the number of

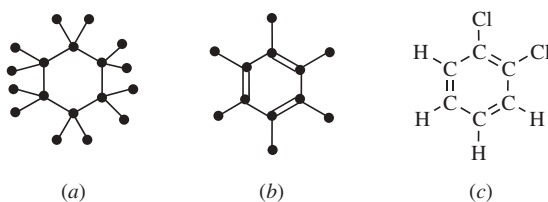
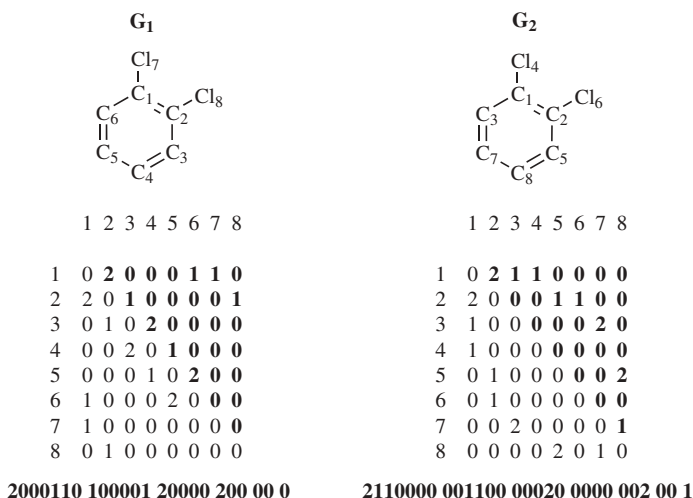


Figure 1 (a) Simple graph, (b) multigraph, and (c) molecular graph.

times that edge occurs in the graph. In Figure 1 diagram (a) contains vertices of degree 1 and 4, and all edges have multiplicity 1; in (b), the vertices have degrees 1 and 4 and the edges have multiplicities 1 and 2. The *degree sequence* of a graph or a multigraph is the sequence of numbers of vertices having a given degree starting with degree 0 and ending with the maximum degree for all vertices. Figure 1(a) has no vertices of degree 0, 12 vertices of degree 1, no vertices of degree 2 and degree 3, and 6 vertices of degree 4; the degree sequence is (0,12,0,0,6). Figure 1(b) has the degree sequence (0,6,0,0,6).

Although Figure 1(b) could correspond uniquely to benzene, we cannot distinguish 1,2-dichlorobenzene from 1,4-dichlorobenzene with this representation. To make the distinction between the two compounds, we have to attach to each vertex a label, or color, that is unique to each element of the periodic table (for instance, the atomic symbol). Finally, in a molecular structure, atoms are always connected through some bonds; in other words, a molecular structure is in one piece. A *molecular graph* is thus defined as a connected multigraph with vertices colored by the atomic symbols of the periodic table. We use the term *color* instead of *label* because, as we shall see next, labeled graphs have a specific definition in graph theory. Figure 1(c) is the molecular graph of 1,2-dichlorobenzene. Clearly, in a molecular graph, each vertex is an atom and each edge is a bond. The terms *atom valence* replace the terms *vertex degree*, and *bond order* replace *edge multiplicity*. Note that with the exception of rare gases, a molecular graph comprises more than one atom. Because molecular graphs are connected, their valence sequences start with valence 1 and usually end with valences 4 or 5 for most organic compounds. The valence sequence of benzene is (6,0,0,6).

Now that we have defined molecular graphs, we need to find an appropriate representation for computer manipulation and storage. Assuming our molecular graph  $G$  has  $n$  atoms, we first start to label each atom with numbers 1 through  $n$ . We then create a vector of  $n$  entries where each entry  $i$ ,  $1 \leq i \leq n$ , is the symbol of atom  $i$ . We also create an  $n \times n$  matrix called the *adjacency matrix*, where each entry  $i,j$ ,  $1 \leq i,j \leq n$  is set to the order of the bond between atom  $i$  and atom  $j$ . The maximum bond order is 3, and the order is set to 0 when the two atoms are not bonded. Examples of adjacency matrices are given in Figure 2. Note that the diagonals of the adjacency matrices are filled with 0s, as atoms are not bonded to themselves. Adjacency matrices are symmetric



**Figure 2** Two hydrogen-suppressed molecular graphs with corresponding adjacency matrices and connectivity stacks.

matrices, because when atom  $i$  is bonded to  $j$ , atom  $j$  is also bonded to  $i$ . A convenient way to store adjacency matrices into a compact code was introduced by Kudo and Sasaki.<sup>8</sup> This code, called the connectivity stack, is obtained by reading the upper triangle of the adjacency matrix row by row from left to right. Examples of connectivity stacks are given in Figure 2. Connectivity stacks can be compared. Let  $A = a_1a_2 \dots a_i \dots$  and  $B = b_1b_2 \dots b_i \dots$  be two connectivity stacks where  $a_i$  and  $b_i$  take the values 0, 1, 2, or 3. We then write  $A \geq B$  if an index  $i$  exists such that  $a_i \geq b_i$ ,  $a_{i-1} = b_{i-1}, \dots, a_1 = b_1$ . Taking the example of Figure 2, the connectivity stack of graph  $G_2$  is greater than the connectivity stack of graph  $G_1$ .

To code a molecular graph, we have to label all atoms of our graph and transform the graph into what is called in graph theory a *labeled graph*, i.e., a graph for which each vertex has a distinct label. This process, of course, does not mean that a one-to-one correspondence exists between molecules and labeled graphs, as the two different labeled graphs shown in Figure 2 correspond to the same molecule. Some instances exist, however, in which graphs used by chemists can be appropriately represented by labeled graphs. For instance, in linear reaction networks and protein and gene networks, all vertices have a unique label (e.g., a compound name). Another example is combinatorial libraries obtained by attaching reactants to a scaffold having no symmetry. In this case, the reactants have a unique name, and the reacting sites on the scaffold being different can be labeled uniquely. In general, molecules should not be considered as labeled graphs, and as we shall see in this chapter, the techniques used by chemists to count, enumerate, and sample molecules are all derived from combinatorial results obtained with *unlabeled* graphs.

Although molecules are not appropriately represented by labeled graphs, they are stored and manipulated (by computers) as such. We thus need to find a way to detect when two labeled representations correspond to the same molecular structures. Two labeled representations correspond to the same (unlabeled) graph if a one-to-one mapping between the two sets of labels can be found to also map the edges of the graphs. This mapping is called *isomorphism*. Formally, two labeled graphs  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$  are isomorphic if a one-to-one mapping  $\pi$  from  $V_1$  to  $V_2$  can be found such that  $\pi(E_1) = E_2$ . Note that with molecular graphs, we have to restrict the mapping  $\pi$  between atoms having the same atomic symbol. With the notation  $(i\ j)$  to specify that atom  $i$  from graph  $G_1$  is mapped to atom  $j$  in graph  $G_2$ , the isomorphism  $(1\ 1)(2\ 2)(3\ 5)(4\ 8)(5\ 7)(6\ 3)(7\ 4)(8\ 6)$  maps graph  $G_1$  to graph  $G_2$  in Figure 2. Note that  $(1\ 2)(2\ 1)(3\ 3)(4\ 7)(5\ 8)(6\ 5)(7\ 6)(8\ 4)$  is also an isomorphism from  $G_1$  to  $G_2$ .

Because several labelings of the same molecular graph can occur, it is important to distinguish one of them. We shall call this labeling the *canonical* one. Several ways of obtaining a canonical labeling exist. The one we chose in this chapter is the one leading to the maximal connectivity stack. Taking all possible  $8! = 40320$  labeling of the graphs in Figure 2, we can verify (with the help of a computer code) that no other labeling of the vertices has a connectivity stack greater than the one of graph  $G_2$ . Of course, better ways of canonizing a connectivity stack exist than checking all possible labelings. Algorithms capable of performing this checking can be found in the literature,<sup>9</sup> but reviewing them is not the purpose of this chapter. The computational complexity of canonizing a general graph is unknown; i.e., no fast algorithm has yet been found. However, molecular graphs can theoretically be canonized efficiently.<sup>9</sup> Furthermore, some fast graph canonizers such as Brendan McKay's code Nauty<sup>10</sup> can easily be adapted to canonize molecular graphs.

From the definition of isomorphism given earlier, two atoms  $x_1$  of graph  $G_1$  and  $x_2$  of graph  $G_2$  are isomorphic if an isomorphism can be found matching  $x_1$  to  $x_2$  and matching the bonds of  $G_1$  to those of  $G_2$ . For instance in Figure 2, atom 4 in  $G_1$  is isomorphic to atom 8 in  $G_2$ . Now, if  $G_1 = G_2$ , we say that the two atoms are *equivalent* and instead of isomorphism we use the term *automorphism*. Taking again the example of graph  $G_1$  in Figure 2, and taking again the notation  $(i\ j)$  to specify that atom  $i$  is mapped to atom  $j$ , the mapping  $(1\ 1)(2\ 2)(3\ 6)(4\ 5)(5\ 4)(6\ 3)(7\ 7)(8\ 8)$  of the vertices of  $G_1$  leads to a graph identical to  $G_1$ . That mapping, also called a *permutation*, is an automorphism. The permutation notation can be simplified by noting that when  $(i\ j)$  occurs,  $(j\ i)$  also occurs, and writing  $(i)$  instead of  $(i\ i)$ . Thus, the permutation  $(1\ 1)(2\ 2)(3\ 6)(4\ 5)(5\ 4)(6\ 3)(7\ 7)(8\ 8)$  reduces to  $(1)(2)(36)(4\ 5)(7)(8)$ .

Several isomorphisms may exist between two graphs (cf.  $G_1$  and  $G_2$  in Figure 2). Similarly, a given graph may have several automorphisms. The *automorphism group* of a graph is the set of all of its automorphisms. The automorphism group of the hydrogen-suppressed molecular graph of benzene is given in Figure 3. This group is the dihedral  $D_{6h}$  group.

Because two atoms are equivalent, if they can be mapped by an automorphism, we can partition the atoms into *equivalent classes* using the automorphism group of a graph. In graph theory, the atom equivalent classes are named the orbits of the automorphism group. Chemically, atoms that belong to the same equivalent class are symmetrical and among other properties have the same chemical shift in NMR spectra. Many algorithms in the chemistry literature compute the atom equivalent classes of molecular graphs.<sup>9</sup>

Another term we use in this chapter is the *subgraph* of a graph. A subgraph of a graph is obtained by selecting any subset of vertices of the graph and by selecting any subset of edges of the graph that are attached to the selected vertices. In chemistry, subgraphs are molecular fragments. As depicted in Figure 4, the fragments of a molecular graph may or may not overlap.

We finish this subsection with a few additional definitions. A molecular graph is a *tree* if it does not contain cycles. Multiple bonds are allowed in molecular trees, and alkanes and alkenes are examples of molecular trees. A rooted tree is a tree where one vertex (the root) is distinguished from the others. Isotopically monolabeled alkanes are a rooted tree. Of course, not all molecules are trees, but all molecules are *bonded valence* graphs. Bonded valence graphs are graphs for which the degree of any vertex is below some threshold. For most organic molecules, the maximum valence of any atom is 4 (sometimes 5), and for any molecule, the number of bonds attached to any atom is always limited because of the three-dimensional space limitations surrounding an atom. Thus, molecular graphs are bonded valence graphs. This property is important because bonded valence graphs can usually be treated more easily than general graphs. For instance, isomorphism can be solved efficiently for that class of graph.<sup>11</sup> The term *efficient* has a specific meaning here. An algorithm is said to be efficient if the time taken and the space allocated to complete the job is a polynomial of the size of the problem. With a molecular graph, the size of the problem is usually the number of atoms. An example of an efficient algorithm is Kudo-Sasaki's quadratic  $O(n^2)$  time and space algorithm that computes the connectivity stack from an adjacency matrix.<sup>8</sup> Problems that cannot be solved by a polynomial time and space algorithm are said to be *intractable*. Searching all occurrences of a fragment (subgraph) in a molecular graph is an intractable problem.<sup>12</sup>

---

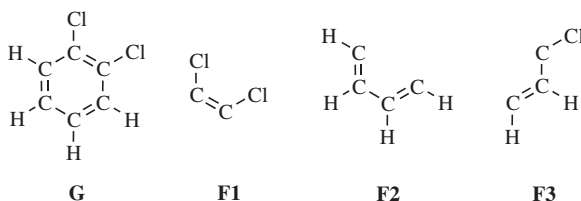
## COUNTING STRUCTURES: HOW MANY ISOMERS HAS DECANE?

### Counting Labeled and Unlabeled Graphs

In this subsection, we briefly summarize results relevant to labeled graphs. We then survey the work on counting series by Cayley, Pólya, Harary,

Symmetry operation	Permutation	Symmetry operation	Permutation
E			$(1)(4)(2\ 6)(3\ 5)$
$C_6^-$			$(3)(6)(1\ 5)(2\ 4)$
$C_6^+$			$(2)(5)(1\ 3)(4\ 6)$
$C_3^-$			$(1\ 6)(2\ 5)(3\ 4)$
$C_3^+$			$(1\ 2)(3\ 6)(4\ 5)$
$C_2$			$(1\ 4)(2\ 3)(5\ 6)$
		$\sigma_v^{(1)}$	$\sigma_v^{(1)}$
		$\sigma_v^{(2)}$	$\sigma_v^{(2)}$
		$\sigma_v^{(3)}$	$\sigma_v^{(3)}$
		$\sigma_v^{(4)}$	$\sigma_v^{(4)}$
		$\sigma_v^{(5)}$	$\sigma_v^{(5)}$
		$\sigma_v^{(6)}$	$\sigma_v^{(6)}$

**Figure 3** List of all (12) automorphisms for hydrogen-suppressed benzene. In the permutation notation, (1 2) reads “1 goes to 2 and 2 goes to 1” and (1 3 5) reads “1 goes to 3, 3 goes to 5, and 5 goes to 1.”



**Figure 4** Fragments of a molecular graph *G*. *F1* and *F2* are nonoverlapping fragments, and *F3* overlaps with both *F1* and *F2*.

and Read. Particular attention is given to Pólya's work because it has led to many applications in chemistry.

Although labeled graphs are not the appropriate objects to describe molecular graphs, we recall that combinatorial libraries, protein and gene networks, and linear reaction networks can be represented by labeled graphs. Furthermore, we will use some of the results given here to count unlabeled graphs and molecular graphs.

We first note that  $\binom{n}{2} = n(n-1)/2$  possible distinct edges exist between  $n$  vertices, and  $\binom{n(n-1)/2}{k}$  ways exist of choosing  $k$  edges in a set of  $n(n-1)/2$  edges. Summing for all possible  $k$  values gives the number of labeled graphs of  $n$  vertices:

$$L_n = \sum_{k=0}^{n(n-1)/2} \binom{n(n-1)/2}{k} = 2^{n(n-1)/2} \quad [1]$$

As the objects we are interested in this chapter are connected, let  $C_k$  be the number of connected labeled graphs of  $k$  vertices.  $kC_k$  rooted connected labeled graphs exist because  $k$  ways exist of choosing a root. The number of rooted, labeled graphs of  $n$  vertices in which the root is in a connected component containing  $k$  vertices is  $kC_k \binom{n}{k} L_{n-k}$ . This expression, summed from  $k=1$  to  $n$ , is equal to  $nL_n$ , which is the number of rooted labeled graphs. Thus,  $nL_n = \sum_{k=1}^n kC_k \binom{n}{k} L_{n-k}$ , from which we derive the recursive formula for counting the number of connected labeled graphs of  $n$  vertices:

$$C_n = 2^{n(n-1)/2} - \frac{1}{2} \sum_{k=0}^{n-1} k \binom{n}{k} 2^{(n-k)(n-k-1)/2} C_k \quad [2]$$

Interestingly enough, investigations related to counting unlabeled graphs started with the pragmatic problem of calculating the number of paraffin

structures. Cayley was the first to propose a solution,<sup>1</sup> and in doing so, he introduced the notion of trees.<sup>13</sup> Applications of Cayley's counting formula are tedious and prone to errors. Cayley made several errors in his work, some of which were later corrected by Hermann.<sup>14</sup> Almost 75 years after Cayley's initial work, Henze and Blair<sup>15</sup> proposed a recursive formula much easier to apply than the Cayley formulation. The next significant advance came with the work of Pólya<sup>16</sup> and his famous theorem. Most of today's counting techniques apply Pólya's theory, and we shall first describe the theory before using it to count objects relevant to chemistry.

*Pólya Theory of Counting.* In 1935, Pólya proposed a counting theory that is probably the most powerful counting technique available to chemists. In fact, in his original paper, Pólya already described his theory by counting the number of structural isomers when hydrogen atoms in benzene are successively substituted with monovalent atoms or groups.<sup>16</sup> The theory relies on the concept of the cycle index,  $Z(A)$ , where  $A$  is a permutation group with object set  $X = \{1, 2, \dots, n\}$  and  $Z$  stands for the German word *Zyklenzeiger* (meaning cycle index). Applied to a graph,  $X$  is the set of vertices and  $A$  is the automorphism group of the graph as defined in subsection "From Graph Theory to Chemistry". Keeping in mind that the automorphism group of a graph takes into account the permutations, or symmetry operations, namely, the proper rotation axes of the structure, the cycle index of a given permutation,  $\alpha$ , is obtained by decomposing  $\alpha$  into disjoint cycles, formally,

$$Z(\alpha; s_1, s_2, \dots, s_n) = \prod_{k=1}^n s_k^{j_k(\alpha)} \quad [3]$$

where  $s_k$  is a variable representing cycles of length  $k$  and  $j_k(\alpha)$  is the number of cycles  $s_k$  in  $\alpha$ . To show cycle decomposition, let  $\alpha = \sigma_v^{(1)}$  be the reflection of benzene perpendicular to the axis going through atoms 1 and 4. As depicted in Figure 3,  $\alpha = (1)(4)(2\ 6)(3\ 5)$  is composed of two cycles of length 1: (1) and (4), and two cycles of length 2: (2 6) and (3 5). Its cycle index is  $Z(\alpha) = s_1^2 s_2^2$ . The cycle index of an automorphism group  $A$  is simply obtained by summing the cycle decompositions for all permutations in the group and dividing by the size of the group:

$$Z(A; s_1, s_2, \dots, s_n) = \frac{1}{|A|} \sum_{\alpha \in A} Z(\alpha; s_1, s_2, \dots, s_n) = \frac{1}{|A|} \sum_{\alpha \in A} \prod_{k=1}^n s_k^{j_k(\alpha)} \quad [4]$$

From Figure 3 and Table 1, it is easy to verify that the cycle index of benzene is

$$Z = \frac{1}{12} (s_1^6 + 4s_2^3 + 2s_3^2 + 2s_6^1 + 3s_1^2 s_2^2) \quad [5]$$



Table 1 Cycle Index for Benzene

Symmetry Operation	Permutation	Cycle Index
E	(1)(2)(3)(4)(5)(6)	$s_1^6$
$C_6^-$	(1 2 3 4 5 6)	$s_6^1$
$C_6^+$	(6 5 4 3 2 1)	$s_6^1$
$C_3^-$	(1 3 5)(2 4 6)	$s_3^2$
$C_3^+$	(5 3 1)(6 4 2)	$s_3^2$
$C_2$	(1 4)(2 5)(3 6)	$s_2^3$
$\sigma_v^{(1)}$	(1)(4)(2 6)(3 5)	$s_1^2 s_2^2$
$\sigma_v^{(2)}$	(3)(6)(1 5)(2 4)	$s_1^2 s_2^2$
$\sigma_v^{(3)}$	(2)(5)(1 3)(2 4)	$s_1^2 s_2^2$
$\sigma_v^{(4)}$	(1 6)(2 5)(3 4)	$s_2^3$
$\sigma_v^{(5)}$	(1 2)(3 6)(4 5)	$s_2^3$
$\sigma_v^{(6)}$	(1 4)(2 3)(5 6)	$s_2^3$

Automorphism permutations are listed in Figure 3.

To introduce Pólya's theorem, we take the example of counting the numbers of isomers obtained when substituting hydrogen atoms in benzene with chlorine atoms. Pólya's theorem applied to this problem states that the number of isomers, when  $k$  hydrogen atoms are substituted by  $k$  chlorine atoms, is the coefficient  $C_k$  of the generating function,  $C(x)$ , obtained by substituting in the cycle index each variable  $s_k$  by  $(1 + x^k)$ . Although the proof of this theorem can be found in Pólya's paper,<sup>16</sup> or more recent sources such as the book of Harary and Palmer,<sup>49</sup> intuitively the substitution comes from the fact that each hydrogen atom can be either replaced or not by a chlorine atom. These two possibilities (0 or 1) are expressed by the function  $x^0 + x^1 = 1 + x$ , which Pólya calls the *figure* generating function. Now, observing that more ways of coloring an object with no symmetry exist than ways of coloring an object in which all vertices are symmetrical, we realize that the automorphism group of the studied object has to play a role in counting the number of configurations. The exact relationship between the number of configurations and the automorphism group is given by Pólya's theorem.

*Theorem (Pólya).* The configuration generating function, or counting series,  $C(x)$ , is obtained by substituting the figure generating function,  $c(x)$ , in the cycle index, by replacing every occurrence of  $s_k$  in the cycle index by  $c(x^k)$ . Thus,

$$C(x) = Z(A; c(x), c(x^2), c(x^3), \dots) \quad [6]$$

A corollary of Pólya's theorem is that the total number of configurations,  $N$ , obtained after coloring an object with permutation group  $A$  with  $n$  colors is

obtained by replacing every occurrence of  $s_k$  in the cycle index of  $A$  by  $n$ . Formally,

$$N(A; s_1, s_2, \dots, s_n) = \frac{1}{|A|} \sum_{\alpha \in A} \prod_{k=1}^n n^{j_k(\alpha)} = \frac{1}{|A|} \sum_{\alpha \in A} n^{\sum_{k=1}^n j_k(\alpha)} = \frac{1}{|A|} \sum_{\alpha \in A} n^{|\alpha|} \quad [7]$$

where  $|\alpha|$  is the number of orbits of the permutation  $\alpha$ . For a graph,  $|\alpha|$  is the number of vertex equivalent classes induced by  $\alpha$ .

As an illustration, the generating function  $c(x) = 1 + x$  substituted in the benzene cycle index of Eq. [5] gives the counting series:

$$\begin{aligned} C(x) &= 1/12[(1+x)^6 + 4(1+x^2)^3 + 2(1+x^6) + 2(1+x)^2(1+x^2)^2] \\ &= 1 + x + 3x^2 + 3x^3 + 3x^4 + x^5 + x^6 \end{aligned} \quad [8]$$

and the total number of configurations is  $1/12 [(2)^6 + 4(2)^3 + 2(2) + 2(2)^2(2)^2] = 13$ . The coefficients of Eq. [8] represent the number of isomers of benzene (1), chlorobenzene (1), dichlorobenzene (3), trichlorobenzene (3), and soon, up to hexachlorobenzene (1). The various structural isomers counted in Eq. [8] are listed in Figure 5.

### Counting Molecules

Although we have already seen examples on how Pólya's theorem can enumerate chemical compounds, we consider this idea in greater detail in

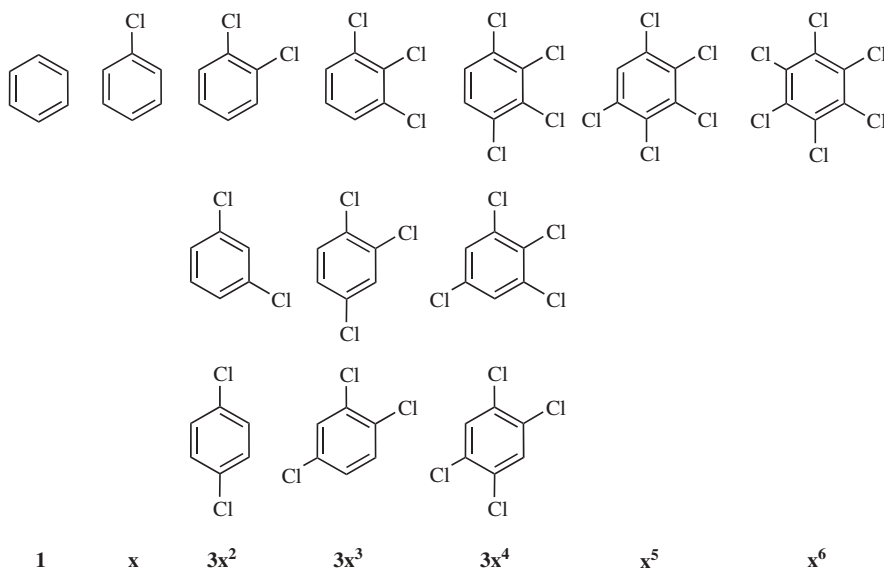


Figure 5 Count of  $k$ -chloro-benzene isomers with Pólya's theorem.

this subsection. The most general problem of this kind is to determine the number of isomers given a molecular formula. Although this problem can be solved with explicit enumeration techniques (cf. “Enumerating Structures” section), currently no counting series provides the number of isomers of a given molecular formula. However, if we lower our expectations and confine our attention to some restricted class of compounds, a mathematical treatment of the problem then becomes possible.

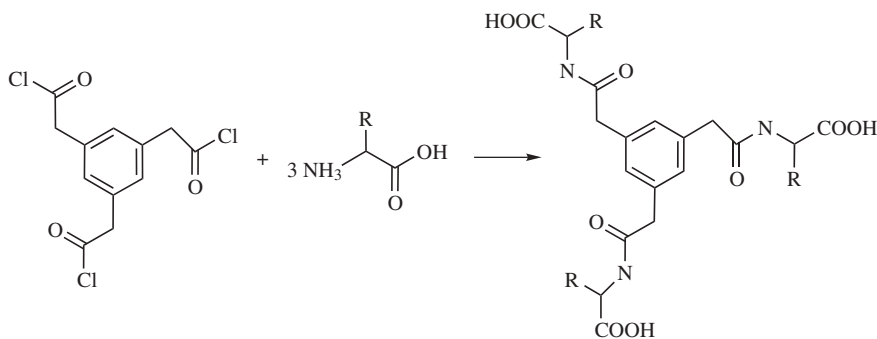
The most straightforward application of Pólya’s theorem is with substituted or labeled hydrocarbons. Indeed, we have already seen that the count of structural isomers obtained after substituting hydrogen atoms in benzene with chlorine atoms is derived directly by plugging the figure generating function  $c(x) = 1 + x$  into the cycle index of benzene. In Table 2, the same exercise is carried out for other benzenoid hydrocarbons, and the number of isomers obtained after substituting  $k$  hydrogen atoms is the  $x^k$  coefficient of the corresponding counting series. This type of calculation can also be performed to count substituted fullerenes,<sup>18–20</sup> polyhedral cages,<sup>20</sup> and substituted

**Table 2** Cycle Indices and Counting Series for Some Substituted Benzenoids and Hydrocarbons Cages

Benzenoids and Cages	Symmetry Group	Cycle Index	Counting Series
Benzene	$D_{6h}$	$1/12(s_1^6 + 4s_2^3 + 2s_3^2 + 2s_6^1 + 3s_1^2s_2^2)$	$1 + x + 3x^2 + 3x^3 + 3x^4 + x^5 + x^6$
Naphthalene	$D_{2h}$	$1/4(s_1^8 + 3s_2^4)$	$1 + 2x + 10x^2 + 14x^3 + 22x^4 + 14x^5 + 10x^6 + 2x^7 + x^8$
Anthracene	$D_{2h}$	$1/12(s_1^{10} + s_1^2s_2^4 + 2s_2^5)$	$1 + 3x + 15x^2 + 32x^3 + 60x^4 + 66x^5 + 60x^6 + 32x^7 + 15x^8 + 3x^9 + x^{10}$
Phenanthrene	$C_{2v}$	$1/2(s_1^{10} + s_2^5)$	$1 + 5x + 25x^2 + 60x^3 + 110x^4 + 126x^5 + 110x^6 + 60x^7 + 25x^8 + 5x^9 + x^{10}$
Tetracene	$D_{2h}$	$1/4(s_1^{12} + 2s_2^6)$	$1 + 3x + 21x^2 + 55x^3 + 135x^4 + 198x^5 + 236x^6 + 198x^7 + 125x^8 + 55x^9 + 21x^{10} + 3x^{11} + x^{12}$
Triphenylene	$D_{3h}$	$1/6(s_1^{12} + 2s_2^6 + 2s_3^4)$	$1 + 2x + 14x^2 + 38x^3 + 90x^4 + 132x^5 + 166x^6 + 132x^7 + 90x^8 + 38x^9 + 14x^{10} + 2x^{11} + x^{12}$
$C_{60}$	$I_h$	$1/120(24s_1^{10} + 20s_6^{10} + 24s_5^{12} + 20s_3^{20} + 16s_2^{30} + 15s_1^4s_2^{28} + s_1^{60})$	$1 + x + 23x^2 + 303x^3 + 4190x^4 + 45718x^5 + 418470x^6 + 3220218x^7 + 21330558x^8 + 123204921x^9 + 628330629x^{10} + \dots$

cycloalkanes.<sup>21</sup> A general approach to counting substituted isomers based on their symmetries can be found in Baraldi and Vanossi.<sup>22</sup>

Another direct application of Pólya's theorem is to compute the sizes of combinatorial libraries that can be generated from a scaffold and a set of reactants. We recall that the total number of configurations obtained after coloring an object with  $k$  colors is obtained by substituting  $k$  in the cycle index of the object. Library sizes are computed by taking the scaffold as the object and the reactant as the colors. As an example, consider the following reaction scheme:



Three reacting sites exist on the benzene ring; the cycle index of benzene reduced to these three sites is  $\frac{1}{6}(s_1^3 + 3s_1s_2 + 2s_3)$  (cf. Figure 3 for a list of the permutations involved). According to Eq. [7] attaching  $n$  different R reactants to tri-acidchloride will result in a library of size  $\frac{1}{6}(n^3 + 3n^2 + 2n)$ . Now,  $n = 2$  different reactants will give 4 compounds, and if the reactants are all possible ( $n = 20$ ) amino acids, the library will be composed of 1540 compounds. Scaffolds in library design often have no symmetry, and if such a scaffold is composed of  $r$  reacting site, its cycle index is  $s_1^r$ . The size of a library composed of  $n$  reactants for a scaffold with no symmetry and  $r$  reacting sites is  $n^r$ .

A bit more challenging is the application of Pólya's theorem to count alkyl groups. An alkyl group has the formula  $-\text{C}_n\text{H}_{2n+1}$ , and it contains one free bond or bonding site. An alkyl group is a rooted tree because the carbon atom carrying the bonding site can be distinguished from the other. Let  $A_n(x)$  be the counting series for alkyl groups having  $n$  atoms. The remarkable idea in counting alkyl groups with Pólya's theorem is to apply a figure generating function that is the counting series itself. It is thus a recursive process in which the number of alkyl group of  $n$  atoms is counted from the number of alkyl groups of  $n-1$  atoms. To apply Pólya's theorem, we must first determine the group of permutations attached to atom number  $n$ . The permutations attached to any carbon atom in an alkyl chain are listed in Figure 6.

Clearly, all alkyl groups attached to a carbon atom are interchangeable. The permutation group is called the *symmetric* group  $S_3$  with cycle index  $\frac{1}{6}(s_1^3 + 3s_1s_2 + 2s_3)$ . Substituting  $A_{n-1}(x)$  into the cycle index of  $S_3$  gives a

symmetry operation		permutation	cycle index
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array}$		$(\text{R1})(\text{R2})(\text{R3})$	$s_1^3$
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R1} \\   \\ \text{R2}-\text{C}-\text{R3} \\   \\ \text{X} \end{array}$		$(\text{R1})(\text{R2 R3})$	$s_1 s_2$
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R3} \\   \\ \text{R1}-\text{C}-\text{R2} \\   \\ \text{X} \end{array}$		$(\text{R2})(\text{R1 R3})$	$s_1 s_2$
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R2} \\   \\ \text{R3}-\text{C}-\text{R1} \\   \\ \text{X} \end{array}$		$(\text{R3})(\text{R1 R2})$	$s_1 s_2$
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R2} \\   \\ \text{R1}-\text{C}-\text{R3} \\   \\ \text{X} \end{array}$		$(\text{R1 R2 R3})$	$s_3$
$\begin{array}{c} \text{R1} \\   \\ \text{R3}-\text{C}-\text{R2} \\   \\ \text{X} \end{array} \longrightarrow \begin{array}{c} \text{R3} \\   \\ \text{R2}-\text{C}-\text{R1} \\   \\ \text{X} \end{array}$		$(\text{R3 R2 R1})$	$s_3$

Figure 6 Permutation group and cycle index for carbon atoms in alkyl groups.

counting series representing the number of ways of attaching three alkyl groups to an additional atom. Multiplying the resulting series by  $x$ , that is, adding the additional atom number  $n$ , leads to the following counting series for alkyl groups. Starting with  $A_0(x) = 1$ :

$$A_n(x) = 1 + \frac{1}{6}x[A_{n-1}^3(x) + 3A_{n-1}(x)A_{n-1}(x^2) + 2A_{n-1}(x^3)] \quad [9]$$

The 1 on the right-hand side must be added to ensure that the term  $A_0$  corresponding to a hydrogen is properly counted. In Eq. [9] the coefficient of  $A_n(x)$  has to be computed only up to  $x^n$ . To avoid this restriction, it is customary to write Eq. [9] up to  $n = \infty$ :

$$A(x) = \sum_{n=0}^{\infty} A_n x^n = 1 + \frac{1}{6}x[A^3(x) + 3A(x)A(x^2) + 2A(x^3)] \quad [10]$$

The basic operations in Eq. [10] are summations and products of polynomials and a scalar multiplication. For polynomials of order  $n$ , summations

and scalar multiplications are performed with no more than  $n$  integer arithmetic operations, whereas polynomial products necessitate at most  $O(n^2)$  integer operations. The total cost of computing the counting series is, therefore,  $O(n^3)$ . The first elements of the series are  $A(x) = 1 + x + x^2 + 2x^3 + 4x^4 + 8x^5 + 17x^6 + 39x^7 + 89x^8 + 211x^9 + 507x^{10} + \dots$

### *The Number of Isomers of Acyclic Compounds*

We can use the coefficients  $A_0, A_1, \dots, A_n$  of Eq. [10] to evaluate the number of isomers for several families of acyclic compounds comprising up to  $n$  carbon atoms. Computationally, these numbers can be obtained by summations, products, and scalar multiplications of the polynomial  $A(x)$ . The results that follow have been derived by Read.<sup>23</sup>

*Primary alcohols.* The primary alcohols are of the form  $R\text{-CH}_2\text{-OH}$ , where  $R$  is an alkyl group with  $n-1$  carbons atoms. To maintain the correct number of carbon atoms, the counting series for primary alcohols becomes

$$xA(x) \quad [11]$$

*Secondary alcohols.* The secondary alcohols are of the form  $R_1\text{-CH(R}_2\text{)-OH}$ , where  $R_1$  and  $R_2$  are alkyl groups. To count these isomers, we apply Pólya's theorem with the figure counting series  $A(x) - 1$  because  $R_1$  and  $R_2$  are not hydrogen atoms. The permutation group is the symmetric group  $S_2$  because  $R_1$  and  $R_2$  are interchangeable, and the counting series for secondary alcohols is

$$xZ(S_2; A(x) - 1) = 1/2x[A^2(x) - 2A(x) + A(x^2)] \quad [12]$$

*Tertiary alcohols.* The formula for a tertiary alcohol is  $\text{HO-C(R}_1\text{)-(R}_2\text{)(R}_3\text{)}$ . The counting series is obtained with the same arguments as for secondary alcohols but with the permutation group  $S_3$  instead of  $S_2$ :

$$xZ(S_3; A(x) - 1) = 1/6x[A^3(x) - 3A^2(x) + 3A(x)A(x^2) - 3A(x^2) + 2A(x^3)] \quad [13]$$

*Aldehydes and ketones.* These compounds have the form  $R_1\text{-C=O-R}_2$ , where  $R_1$  and  $R_2$  are alkyl groups and, possibly, hydrogen atoms. As hydrogen atoms are included, the counting series is

$$xZ(S_2; A(x)) = 1/2x[A^2(x) + A(x^2)] \quad [14]$$

*Alkynes.* The formula for acetylene compounds takes the form  $R_1\text{-C}\equiv\text{C-R}_2$ , and because two additional carbon atoms are available when no terminal hydrogens exists, the counting series is

$$xZ(S_2; A(x)) = 1/2x^2[A^2(x) + A(x^2)] \quad [15]$$

*Esters.* The general ester formula is  $R_1-C(=O)-OR_2$ , with  $R_1$  and  $R_2$  being alkyl groups.  $R_1$  can be a hydrogen atom but not  $R_2$ ; otherwise the compound would be an acid. Consequently, the counting series is  $A(x)$  for  $R_1$  and  $A(x) - 1$  for  $R_2$ . One more carbon must be added in the ester formula. The final counting series becomes

$$xA(x)[A(x) - 1] \quad [16]$$

*Isotopically labeled alkanes.* This class of alkanes has one carbon atom labeled with, for instance, a C-13 isotope. The general formula for these compounds is  $C(R_1)(R_2)(R_3)(R_4)$ , where  $C$  is the labeled carbon atom and  $R_i$ ,  $i = 1, \dots, 4$  are alkyl groups whose counting series is  $A(x)$ . As all alkyl groups can be exchanged with one another around the labeled carbon atom, the permutation group is the symmetric group  $S_4$  with cycle index  $1/24(s_1^4 + 6s_1^2s_2 + 3s_2^2 + 8s_1s_3 + 6s_4)$  and the counting series for labeled alkanes is

$$\begin{aligned} P(x) &= \sum_{n=1}^{\infty} P_n x^n = xZ(S_4; A(x)) \\ &= \frac{1}{24}x[A^4(x) + 6A^2(x)A(x^2) + 3A^2(x^2) + 8A(x)A(x^3) + 6A(x^4)] \end{aligned} \quad [17]$$

We now turn our attention to a class of compounds that is not of primary importance in chemistry, but the results are used later to derive the counting series for alkanes. These structures are of the type  $R_1-R_2$ , where  $R_1$  and  $R_2$  are non-hydrogen alkyl groups with counting series  $A(x) - 1$ . The permutation group is  $S_2$ , and the counting series is

$$Q(x) = \sum_{n=1}^{\infty} Q_n x^n = Z(S_2; A(x) - 1) = \frac{1}{2}[(A(x) - 1)^2 + A(x^2) - 1] \quad [18]$$

*Alkanes.* We may think that counting alkanes would be less difficult than counting alcohols, ketones, esters, and other substituted or labeled structures, especially because these compounds are all derived from alkanes, but this is not the case. Actually, Cayley, Henze and Blair, and even Pólya had a great deal of difficulty finding the alkane counting series. Their solutions are complex involving tree centers and bicenters. For instance, in the case of Cayley, the solution for alkanes was developed in 1875,<sup>1</sup> 18 years after finding the counting series for rooted trees.<sup>13</sup> It was only in 1948 that a simple formula for unlabeled trees was found by Otter.<sup>24</sup> The solution we review next was first given by Read,<sup>23</sup> and it is an application of Otter's formula to alkanes.

Let us first consider an arbitrary unlabeled alkane. We want to find  $p^*$ , the number of different atom-labeled alkanes obtained after labeling all carbon atoms one after another. Two carbon atoms once labeled will produce

the same labeled structure if they are symmetrical. Thus,  $p^*$  is the number of equivalent classes among atoms, formally the number of orbits in the automorphism group as defined in the subsection “From Graph Theory to Chemistry”. Using the same arguments, we find that  $q^*$ , the number of bond-labeled alkanes, is the number of bond’s equivalent classes. Otter<sup>24</sup> and Harary and Norman<sup>17</sup> have shown that for any unlabeled tree,  $p^* - q^* + s = 1$ , where  $s = 1$  if the tree has a symmetric bond (e.g., a bond between two identical subtrees) and  $s = 0$  otherwise. Now, if we sum the previous equation over all alkanes having  $n$  carbon atoms, we obtain  $P_n - Q_n + \sum s = a_n$ , where  $P_n = \sum p^*$  and  $Q_n = \sum q^*$  and  $a_n$  is the number of alkanes having  $n$  carbon atoms. Clearly,  $P_n$  is the number of atom-labeled alkanes having  $n$  carbon atoms and is thus the  $n$ th coefficient of the counting series in Eq. [17]. Similarly,  $Q_n$  is the  $n$ th coefficient in Eq. [18]. To compute  $a_n$ , we have to evaluate  $s$ . As already mentioned,  $s = 1$  for those alkanes having a bond splitting the structure into two identical alkyl groups. These alkanes must, therefore, have an even number  $n$  of carbon atoms, and their count is simply equal to the number  $A_{n/2}$  of alkyl groups having  $n/2$  carbon atoms. The number of alkanes having  $n$  carbon atoms is thus  $a_n = P_n - Q_n + A_{n/2}$ , with  $A_{n/2} = 0$  when  $n$  is odd. The corresponding counting series is obtained by multiplying  $a_n$  by  $x^n$  and summing starting with  $n = 1$ :

$$a(x) = \sum_{n=1}^{\infty} a_n x^n = P(x) - Q(x) + A(x^2) - 1$$

$$a(x) = \frac{1}{24} x [A^4(x) + 6A^2(x)A(x^2) + 3A^2(x^2) + 8A(x)A(x^3) + 6A(x^4)] \quad [19]$$

$$- \frac{1}{2} [(A(x) - 1)^2 - A(x^2) + 1]$$

The first elements of the series are  $A(x) = 1 + x + x^2 + x^3 + 2x^4 + 3x^5 + 5x^6 + 9x^7 + 18x^8 + 35x^9 + 75x^{10} + \dots$ , and to answer the section’s question, 75 structural isomers exist for decane. Applying Eq. [19] the number of alkane isomers up to 25 carbon atoms is given in Table 4 in the “Chemical Information” subsection appearing later in this chapter. Note that  $a(x)$  can be evaluated computationally, with the product, sum, and scalar multiplication operator on the polynomial  $A(x)$  representing the alkyl group counting series. Considering the computational cost to evaluate  $A(x)$ , alkanes up to  $n$  carbon atoms can be counted with no more than  $n^3$  elementary arithmetic operations.

*Hydroxyl ethers.* In a recent development, Wang, Li and Wang<sup>25</sup> proposed a counting series for compounds of the form  $C_i H_{2i+2} O_j$ , which is a first step toward a general counting series for molecular formulas. Their technique involves Pólya’s cycle index and two generating functions for alkyl groups



R(I), where the root is a carbon atom, and alkoxyl groups R(II) are rooted on an oxygen atom.

### The Number of Stereoisomers of Acyclic Compounds

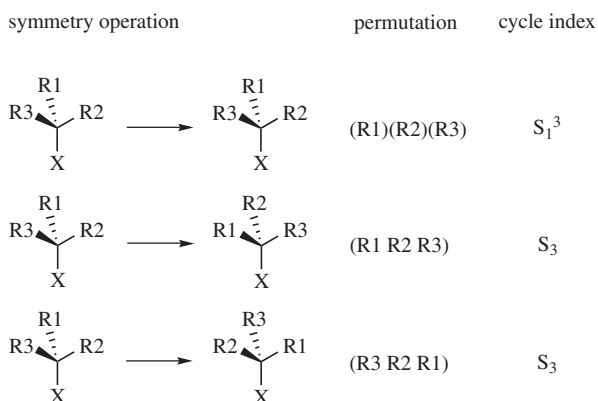
Stereoisomers of acyclic compounds are derived the same way as structural isomers, but the permutation group included in Pólya's cycle index is no longer the symmetric group  $S_3$  with structural isomers, the three alkyl groups attached to any carbon atom are interchangeable, with stereoisomers, the alkyl groups can be arranged in two distinct enantiomeric forms, rectus (*R*) and sinister (*S*). Consequently, in the permutation group attached to carbon atoms, all permutations mapping an *R* form onto an *S* form must be discarded. The remaining permutations are listed in Figure 7, and the permutation group is the cyclic group  $C_3$  with cycle index  $Z(C_3) = 1/3 [s_1^3 + 2s_3]$ .

Now that we have determined the permutation group we can count stereoisomers using the formulas obtained with structural isomers but by replacing the group  $S_3$  by  $C_3$ . For instance, from Eq. [10], the counting series for alkyl groups becomes

$$A'(x) = 1 + xZ(C_3; A'(x)) = 1 + 1/3x[A'^3(x) + 2A'(x^3)] \quad [20]$$

The counting series for functionalized stereoalkanes are summarized in Table 3. All results have been derived by Read.<sup>23</sup> The first elements of the counting series for stereoalkanes are  $a'(x) = 1 + x + x^2 + x^3 + 2x^4 + 3x^5 + 5x^6 + 11x^7 + 24x^8 + 55x^9 + 136x^{10} + \dots$ , and decane thus has 136 stereoisomers.

All isomer counts we have given so far are derived from Pólya's theorem and the alkyl group counting series. Our intention was to illustrate the power of Pólya's counting theory and to make things easier to follow because all



**Figure 7** Permutation group and cycle index for stereo carbon atoms in alkyl groups. All permutations maintain the *R* or *S* stereocenter.

**Table 3** Counting Series for the Stereoisomers of Functionalized Alkanes

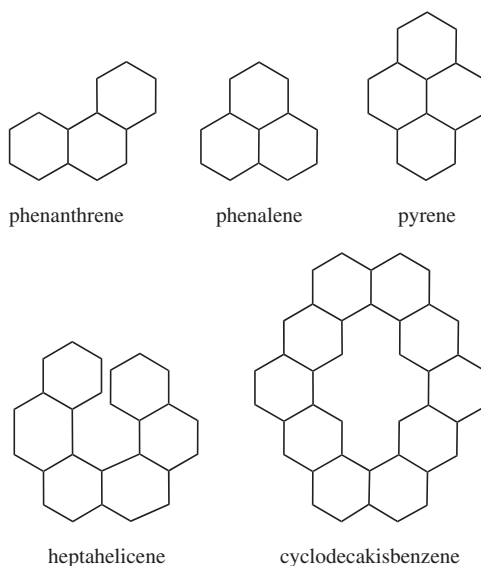
Compound	Formula	Counting Series
Alkyl groups	—R	$A'(x) = 1 + \frac{1}{3}x[A'^3(x) + 2A'(x^3)]$
Primary alcohols	R—CH <sub>2</sub> —OH	$x A'(x)$
Secondary alcohols	R <sub>1</sub> —CH(R <sub>2</sub> )—OH	$x[A'(x) - 1]^2$
Tertiary alcohols	HO—C(R <sub>1</sub> )(R <sub>2</sub> )(R <sub>3</sub> )	$\frac{1}{3}x\{[A'(x) - 1]^3 + 2[A'(x^3) - 1]\}$
Aldehydes, ketones	R <sub>1</sub> —C=O—R <sub>2</sub>	$\frac{1}{2}x[A'^2(x) + A'(x^2)]$
Alkynes	R <sub>1</sub> —C≡C—R <sub>2</sub>	$\frac{1}{2}x^2[A'^2(x) + A'(x^2)]$
Esters	R <sub>1</sub> —C=O—OR <sub>2</sub>	$A'(x)[A'(x) - 1]$
Stereoalkanes	C <sub>n</sub> H <sub>2n+2</sub>	$\frac{1}{12}x[A'^4(x) + 3A'^2(x^2) + 8A'(x)A'(x^3)]$ $- \frac{1}{2}[(A'(x) - 1)^2 - A'(x^2) + 1]$

formulas are derived with the same technique. The reader interested in further details on the applications of Pólya's theory to chiral and achiral compounds and to reaction processes is referred to the book of Fujita.<sup>26</sup> Also, Pólya's theory has been applied to count staggered conformers of alkanes and monocyclic cycloalkanes.<sup>27</sup> Staggered conformers of alkanes are represented by systems that can be embedded in the diamond lattice. Beyond Pólya, few other methods have been proposed in the literature to count acyclic hydrocarbons. In particular, in a series of papers, Yeh gives counting series for alkanes,<sup>28</sup> polyenoids,<sup>29</sup> alkenes,<sup>30</sup> and structures excluding steric strain<sup>31,32</sup> based on Cayley's counting series. Bytautas and Klein<sup>33</sup> have more recently derived a new alkane counting series using a graph's diameter instead of Otter's formula.

### *The Number of Benzenoids and Polyhex Hydrocarbons*

This particular class of hydrocarbons has lead to numerous investigations and probably deserves an entire chapter to be properly reviewed. Here we summarize only the major findings related to counting. The reader further interested by polyhexes and benzenoids can consult the books of Gutman and Cyvin<sup>34–37</sup> as well as the books of Dias.<sup>38,39</sup> These books, as well as that by Trinajstić,<sup>40</sup> provide valuable information regarding the counting and enumeration of Kekulé structures and the conjugated-circuit model, neither of which is reviewed here because of space limitations.

As illustrated in Figure 8, a *polyhex* is a connected system of congruent regular hexagons such that two hexagons either share exactly one edge or are disjoint. Among polyhex hydrocarbons are *helicenes* such as heptahelicene, which are nonplanar, and *coronoids*, such as cyclodecakisbenzene, which are systems with holes. The most heavily studied class of polyhexes has been, by far, *benzenoid* hydrocarbons, which are planar and simply connected. In other words, benzenoid hydrocarbons are condensed polycyclic unsaturated fully conjugated hydrocarbons composed of six-membered rings. The class of benzenoid hydrocarbon is further divided into two subsets: *catacondensed* and *pericondensed*. Catacondensed benzenoids, such as phenanthrene, are systems

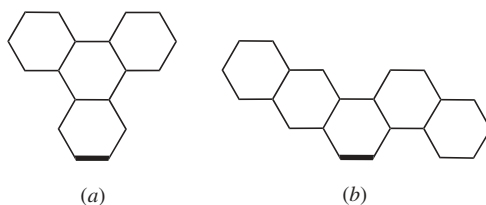


**Figure 8** Some polyhex hydrocarbons.

in which all carbon atoms are lying on the perimeter of the structure. Pericondensed benzenoids are structures having  $n_i \neq 0$  internal atoms, e.g., atoms that do not belong to the perimeter. Phenalene ( $n_i = 1$ ) and pyrene ( $n_i = 2$ ) are examples of pericondensed benzenoids. Finally, all polyhexes are either *Kekuléan* (cf. pyrene) or *non-Kekuléan* (cf. phenalene) depending on whether they possess Kekulé structures. While we are discussing nomenclature, it is worth outlining the distinctions between benzenoid hydrocarbons and polycyclic aromatic hydrocarbons (PAHs). PAHs possess features that are not shared with benzenoid hydrocarbons; they may contain rings with sizes different from six, and they may also comprise  $sp^3$  carbon atoms and side groups.

Essentially, two types of approaches exist to count polyhexes. One is to apply a counting series and Pólya's theorem, and the other is an algorithmic approach based on explicit enumeration. The algorithmic approach is reviewed in the "Enumerating Structures" section as counting is performed through enumeration and each solution is actually generated. It is nonetheless worth mentioning that the algorithmic approach can count planar benzenoid systems, whereas the former approach cannot, as helicenes are included in the counting series. Additionally, further limitations with counting series exist. Polyhex hydrocarbons that cannot be represented by tree-like structures, such as, for instance, pericondensed benzenoids with many internal atoms, cannot be counted.

The first serious attempts to count polyhexes are due to Balaban and Harary<sup>41</sup> and Harary and Read.<sup>42</sup> While Balaban and Harary proposed a



**Figure 9** Bond-rooted catafusenes. (a) *S*-catafusene. Only one hexagon adjoins the hexagon with the root bond (thick line). (b) *D*-catafusene. Two hexagons adjoin the hexagon with the root bond.

nomenclature and simple counting formulas for some benzenoid systems, Harary and Read derived the first counting series for catacondensed polyhexes. The catacondensed systems counted by Harary and Read include helicenes. These systems are also named catafusenes and, strictly speaking, are not benzenoids because they can be nonplanar.

To count catafusenes, like with alkyl groups and alkanes, we first derive a counting formula for bond-rooted catafusenes. A bond-rooted catafusene is a catafusene in which one peripheral bond (the root) has been labeled. We can distinguish two kinds of bond-rooted catafusenes according to whether one or two hexagons are attached to the hexagon containing the root bond (cf. Figure 9). Note that these possibilities are the only ones available if perifusenes are to be avoided. We call them *S*-catafusenes and *D*-catafusenes, respectively.

Let  $S_n$  and  $D_n$  denote the numbers of *S*-catafusenes and *D*-catafusenes having  $n$  hexagons, and let  $U_n = S_n + D_n$  be the total of bond-rooted catafusenes with  $n$  hexagons. From Figure 9, it is easy to be convinced that

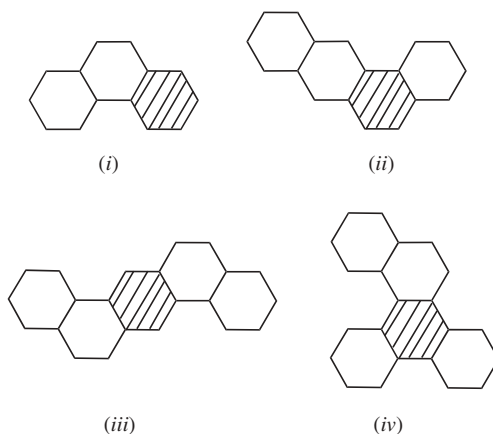
$$\begin{aligned} S_n &= 3U_n \\ D_{n+1} &= \sum_{k=1}^{n-1} U_k U_{n-k} \end{aligned} \quad [21]$$

We now define the three generating functions  $S(x) = \sum_{i=1}^{\infty} S_i x^i$ ,  $D(x) = \sum_{i=1}^{\infty} D_i x^i$ , and  $U(x) = \sum_{i=1}^{\infty} U_i x^i$ . As  $U_n = S_n + D_n$ , we have

$$U(x) = S(x) + D(x) + x \quad [22]$$

The  $x$  on the right-hand side come from the fact that  $U_1 = 1$ , whereas  $S_1 = D_1 = 0$ . Substituting Eq. [21] into Eq. [22] we derive the counting series for bond-rooted catafusene hydrocarbons:

$$U(x) = 3xU(x) + xU^2(x) + x \quad [23]$$



**Figure 10** The four types of rooted catafusenes. The root hexagon is the shaded one. (i) Only one bond-rooted catafusene is attached; (ii) two bond-rooted catafusenes are attached in “meta” position; (iii) two bond-rooted catafusenes are attached in “para” position; and (iv) three bond-rooted catafusenes are attached.

We now wish to count catafusenes in which one hexagon (the root) has been distinguished from the other. Such a rooted catafusene is obtained by taking the root hexagon and attaching one, two, or three of its bonds to a bond-rooted catafusene. As depicted in Figure 10, this process can be performed in four ways.

With the four cases depicted in Figure 10, the number of rooted catafusene of type (i) having  $n + 1$  hexagons is the number  $U_n$  of bond-rooted catafusenes, and the counting series for type (i) rooted catafusenes is  $xU(x)$ . Note that to count the root hexagon in the counting series, we have to multiply  $U(x)$  by  $x$ . To count rooted catafusenes of type (ii) comprising  $n + 1$  hexagons, we have to choose two bond-rooted catafusenes having, respectively,  $k$  and  $n - k$  hexagons. This procedure is similar to the calculation of  $D_{n+1}$  in Eq. [21]. Thus, the counting series for rooted catafusenes of type (ii) is  $xU^2(x)$ . With the rooted catafusenes of type (iii), we have the possibility of building catafusenes, which are invariant under a rotation of  $180^\circ$ , such as in Figure 10(iii). The permutation group attached to the root hexagon in case (iii) is the symmetric group  $S_2$ , and applying Pólya’s theorem, we find the counting series for type (iii) rooted catafusenes to be  $xZ(S_2, U(x)) = x/2[U^2(x) + U(x^2)]$ . Finally, to count rooted catafusenes of type (iv), we first observe that this time we have a possibility of symmetry under rotations of  $120^\circ$ . The permutation group is, therefore, the cyclic group  $C_3$  (already encountered when counting stereoisomers). With Pólya’s theorem, the counting series for type (iv) rooted catafusenes is  $xZ(C_3, U(x)) = x/3[U^3(x) + 2U(x^3)]$ . Summing all terms

corresponding to cases (i) through (iv), the counting series for rooted catafusenes becomes

$$F(x) = x + xU(x) + \frac{3}{2}xU^2(x) + \frac{1}{2}xU(x^2) + \frac{1}{3}xU^3(x) + \frac{2}{3}xU(x^3) \quad [24]$$

The derivation of the counting series for unlabeled catafusenes can be found in Harary and Read.<sup>42</sup> Their solution makes use of Otter's formula<sup>24</sup> the same way the counting series for alkanes was derived with counting series for labeled alkanes and alkyl groups. The counting series for unlabeled catafusenes is

$$H(x) = F(x) - \frac{1}{2}[U^2(x) - U(x^2)]$$

$$H(x) = x + xU(x) + \frac{1}{2}(3x - 1)U^2(x) + \frac{1}{2}(1 + x)U(x^2) + \frac{1}{3}xU^3(x) + \frac{2}{3}xU(x^3) \quad [25]$$

So far we have regarded a catafusene and its mirror image as distinct, provided that the catafusene has no symmetry that would allow it to be rotated into its mirror image. The counting series in Eq. [25] was corrected by Harary and Read<sup>42</sup> to count only once catafusenes and their mirror images. The series is

$$h(x) = \frac{1}{12}(1 + 9x) - \frac{1}{12}(1 - x)(1 - 5x)U(x) + \frac{1}{4}(3 + 5x)U(x^2) + \frac{1}{3}xU(x^3) \quad [26]$$

The first terms of this counting series are  $h(x) = x + x^2 + 2x^3 + 5x^4 + 12x^5 + 37x^6 + 123x^7 + 446x^8 + 1689x^9 + 6693x^{10} + \dots$

Other counting series have been developed, expanding on the initial work of Harary, Balaban, and Read. Harary–Read numbers have been classified and deconvoluted according to symmetries.<sup>43,44</sup> Counting series have been developed for fluorantenoids and fluorenoids,<sup>45</sup> annelated catafusenes,<sup>45</sup> catacondensed monohiptafusenes,<sup>45</sup> and catacondensed octagonal systems.<sup>45</sup> Cyvin et al. have developed a combinatorial summation method that does not invoke counting series and explicit reference to Pólya's theorem. The method has counted perifusenes with one<sup>46</sup> and two internal vertices.<sup>47</sup>

Finally, we should mention the work on conjugated polyene hydrocarbons, which are not polyhexes, but have been counted<sup>45</sup> with a treatment similar to the one we just described for catafusene. The counting series for polyene hydrocarbons is

$$p(x) = \frac{1}{12}[4U(x^3) + (6 + \frac{9}{x})U(x^2) + \frac{4}{x}U(x) - \frac{1}{x^2}U(x)] \quad [27]$$

where  $U(x)$  is the number of bond-rooted polyenes with a counting series similar to Eq. [23]:

$$U(x) = 2xU(x) + xU^2(x) + x \quad [28]$$

### *The Number of Molecular Cages (Fullerenes and Nanotubes)*

To the best of our knowledge, isomers for fullerenes, nanotubes, spher-oalkanes, and other molecular cages have so far been counted only through explicit enumeration (cf. “Enumerating Structures” section). In other words, we are not aware of any formula, counting series, or applications of Pólya’s theorem from which we could compute the number of isomers for these compounds. In fact, molecular cages present a challenge for Pólya’s theory of counting. Looking back, all compounds we have treated so far are either acyclic or have acyclic representations (cf. Balaban and Harary’s paper<sup>41</sup> to see how catafusenes can be represented by trees). Although a solution to enumerate general graphs, including cyclic graphs, applying Pólya’s theorem appeared in 1955,<sup>17</sup> difficulties develop with the class of locally restricted graphs.<sup>23</sup> A locally restricted graph is a graph in which the degrees of its vertices are predefined. Molecular cages are regular graphs in which all atoms have the same degree (for instance, three for fullerenes); they thus belong the class of locally restricted graphs.

In conclusion, we have seen how Pólya’s theory of counting is a powerful and efficient tool to count chemical objects. All counting series derived in this review can be computed with no more than  $O(n^3)$  elementary arithmetic operations for compounds comprising up to  $n$  carbon atoms or  $n$  hexagons. Yet, difficulties exist in deriving counting series for locally restricted graphs, especially if these graphs cannot be represented by trees. A substantial number of chemical compounds unfortunately belong to that difficult class of graphs. Each atom in a molecular graph has a specific degree given by the valence of the atom. Thus, molecules are *always* locally restricted graphs, and unless they have acyclic representations, we cannot easily deal with molecules using counting series. To overcome these difficulties, an alternative is for us to use the explicit enumerations. We review this approach next.

---

## **Enumerating Structures: Are there any isomers of decane having seven methyl groups?**

### **Enumerating Labeled and Unlabeled Graphs**

We begin with the enumeration of labeled graphs because, as with counting, they are easier to deal with. The algorithm we outline next for enumerating labeled graphs we will use and modify to enumerate unlabeled graphs.

Our goal here is to enumerate all possible graphs that can be constructed with a set of vertices labeled 1 through  $n$ . The algorithm given in Scheme I is recursive. At each step of the recursion, we augment the graph by one edge. We start with a graph containing no edges; this is our first labeled graph. Next, we add one edge between any pair of vertices  $[i, j]$ ,  $1 \leq i \leq n, j > i$ . Clearly,  $n(n-1)/2$  of such edges exist. Each of  $n(n-1)/2$  possibilities is a different labeled graph containing one edge. For each of these graphs, a second edge is then added in all possible ways. To avoid generating the same labeled graph, the second edge  $[k, l]$  must be lexicographically greater than the first; i.e.,  $[k, l] > [i, j]$  ( $k > i$  or  $k = i$  and  $l > j$ ). To be convinced the requirement is necessary, consider the graphs  $U_1$  and  $V_1$  having, respectively,  $[1, 2]$  and  $[3, 4]$  as the first edge. Without lexicographic ordering, we can add edge  $[3, 4]$  to  $U_1$  and edge  $[1, 2]$  to  $V_1$ . The two resulting graphs are identical, both being composed of edges  $[1, 2]$  and  $[3, 4]$ . Now, the lexicographic requirement is sufficient because the edges of any labeled graph can be sorted lexicographically. The process of adding edges is repeated until no more can be added; i.e., edge  $[n-1, n]$  already belongs to the graph. Running the algorithm given in Scheme I without constraints, we generate  $m = n(n-1)/2$  labeled  $(n, 1)$ -graphs having  $n$  vertices and one edge, and  $\binom{m}{2}$   $(n, 2)$ -graphs having two edges, which is the number of ways of selecting two edges in a set of  $m$  edges. In general, the algorithm produces  $\binom{m}{q}$   $(n, q)$ -graphs with  $q$  edges. Summing all contributions, the total number of labeled graphs is  $2^m$  in agreement with Eq. [1].

Scheme I: Label-Enumeration(G)

1. IF graph G is completed
2.     PRINT G
3.     ELSE
4.         FOR all edge e lexicographically greater than  
           the edges of G DO
5.             IF constraints are not violated for the graph  
               GU e
6.             Label-Enumeration(GU e)
7.             FI
8.         DONE
9.     FI

The algorithm given in Scheme I can also be run with constraints (cf. step 5) such as degree sequence, specific ranges for the number of edges, number of connected components, and cycle sizes. Additionally, some edges between specific labels may be forbidden, and the presence or absence of specific subgraphs may also be imposed. This algorithm has actually counted and enumerated



gene regulatory networks matching gene expression profiles (i.e., mRNA concentrations).<sup>48</sup> The algorithm was run with two constraints: a list of forbidden edges compiled from the expression profiles and a maximum degree (2 and 3). Scheme I can also generate combinatorial libraries when the scaffold has no symmetry. In such a case, the number of edges is at most the number of reacting sites on the scaffold, and the only edges authorized are between scaffold and reactants.

To use Scheme I to enumerate unlabeled graphs, we need to remove duplicates, i.e., isomorphic graphs. Of course, this removal can be performed after generating all labeled graphs with  $n$  vertices, but this becomes lengthy (i.e.,  $2^{n(n-1)/2}$ ) even for modest  $n$ . A better strategy is to build unlabeled  $(n,q)$ -graphs from unlabeled  $(n,q-1)$ -graphs, which can be carried out by augmenting all unlabeled  $(n,q-1)$ -graphs by one edge. But again, we have to remove duplicates. Observing that  $n(n-1)/2 - (q-1)$  edges can augment any unlabeled  $(n,q-1)$ -graph, and letting  $N_{n,q-1}$  be the number of  $(n,q-1)$ -graphs, we have to test isomorphism between  $[n(n-1)/2 - (q-1)]^2 N_{n,q-1}^2$  pairs of graphs. The problem is that  $N_{n,q}$  scales exponentially with  $n$  and  $q$ .<sup>49</sup> The ideal solution would be to augment each unlabeled  $(n,q-1)$ -graph by one edge without having to be concerned with isomorphism. Fortunately, this solution is possible as Read<sup>50</sup> has shown that the canonical representation of any  $(n,q)$ -graph is an augmentation of the canonical representation of exactly one  $(n,q-1)$  graph. Recall from the subsection "From Graph Theory to Chemistry" that the canonical representation of a graph is a unique ordering of its vertices, such as the one, for instance, that maximizes its connectivity stack. Using Read's results, we can easily modify Scheme I to produce unlabeled graphs. The modified algorithm is given in Scheme II and is named *orderly generation*.

Scheme II: Orderly-Generation-Read-Faradzev(G)

```

1.   IF graph G is completed
2.     PRINT G
3.   ELSE
4.     FOR all edge e lexicographically greater than
       the edges of G DO
5.       IF constraints are not violated for the graph
          G U e
6.       AND CANON(G U e) = G U e
7.       Orderly-Generation-Read-Faradzev(G U e)
8.     FI
9.   DONE
10.  FI
```

The orderly algorithm is to enumeration what Pólya's theorem is to counting. Orderly generation is generally attributed to Read,<sup>50</sup> although

Faradzev<sup>51</sup> independently published an orderly technique. Both Read and Faradzev use the fact that a graph is legitimate if it is identical to its canonical representation (cf. step 6,  $\text{CANON}(G \cup e) = G \cup e$ ). To this end, an artificial ordering must be imposed on the set of graphs that are generated such that a canonical representative always contains a subgraph that is also canonical. A more general orderly algorithm proposed by McKay<sup>52</sup> does not require artificial ordering of graphs and is thus independent of the way the canonical code is constructed. The only requirement is that the canonization procedure induces an ordering of the edges of the graph being canonized. An example of the McKay algorithm is given in Scheme III. This algorithm produces all canonical edge augmentations of a given graph  $G$  having  $q-1$  edges (steps 4-9), which results in a set  $S$  of labeled graphs  $G'$  with  $q$  edges. Identical graphs are removed from the set  $S$  (step 10). Then, in steps 11-16, for every  $(n,q)$ -graph  $G'$  in  $S$ , the algorithm explicitly searches the  $(n,q-1)$ -graph it came from. In other words, the algorithm searches the parent of every child produced. The parent is obtained removing the last edge  $e'$  in  $\text{CANON}(G')$  (step 12). If the parent (e.g., graph  $G'-e'$ ) is the one that was just augmented (i.e., graph  $G$ ), then the child is legitimate (step 13), and the algorithm is recursively run with  $G'$  (step 14); otherwise, graph  $G'$  is ignored.

Scheme III: Orderly-Generation-McKay( $G$ )

```
1. IF graph  $G$  is completed
2.   PRINT  $G$ 
3. ELSE
4.    $S = \emptyset$ 
5.   FOR all edges  $e$  not already in  $G$  DO
6.     IF constraints are not violated for the graph
        $G' = G \cup e$ 
7.        $S = S \cup G'$ 
8.     FI
9.   DONE
10.  Remove duplicates from the set  $S$ 
11.  FOR all graph  $G'$  of  $S$  DO
12.    let  $e'$  be the last edge of  $\text{CANON}(G')$ 
13.    IF  $\text{CANON}(G'-e') = \text{CANON}(G)$ 
14.      Orderly-Generation-McKay( $G'$ )
15.    FI
16.  DONE
17. FI
```

One issue we have not yet addressed with orderly generation is computational complexity. Although orderly generation is certainly faster than labeled enumeration followed by a removal of the duplicated structures, is it the optimum solution? First we have to ask what optimum means when

dealing with enumeration. We certainly cannot hope for a polynomial time algorithm because the number of solutions may be exponentially large, and it already takes an exponential time just to write the solutions. The best we can hope for is an algorithm that runs in polynomial time per output. Such an algorithm indeed exists at least theoretically, as was shown by Goldberg.<sup>53</sup> Goldberg proved that an orderly algorithm can be designed to generate all graphs of  $n$  vertices adding one vertex at a time (not an edge) such that the time delay between two outputs is polynomial. In the proof, Goldberg uses the fact that more graphs of  $n$  vertices than  $n - 1$  vertices always exist and that canonization can be performed in polynomial time for more than half of the graphs of  $n$  vertices, which implies that the enumeration tree always grows, that is, to every  $n - 1$  vertex graph corresponds at least one  $n$  vertex graph. Unfortunately, we cannot use that proof directly when growing graphs by adding edges, because the number of  $(n, q)$ -graphs is not necessarily greater than the number of  $(n, q - 1)$ -graphs. For example, only one  $(n, n(n - 1)/2)$ -graph exists, which is the complete graph (each vertex is connected to all others). Also, one  $(n, n(n - 1)/2 - 1)$ -graph exists, a complete graph without one edge. However, several ways of removing a second edge exist, and thus, more than one  $(n, n(n - 1)/2 - 2)$ -graph exists. Goldberg's result is thus not directly applicable to Schemes II or III. More generally, no guarantee exists that locally restricted graphs, such as molecular graphs restricted by valence sequences, can be constructed in an iterative process such that the number of graphs at a given iteration is always greater than the number of graphs of the previous iteration. Although the theoretical complexity of enumerating molecular graphs is still an open problem, in practice, as we shall see next, fast algorithms exist to enumerate molecules.

As far as general graphs are concerned, some codes are available for their enumeration. In particular, two codes to enumerate small graphs and bipartite graphs can be downloaded along with Nauty, a graph canonizer we mentioned earlier.<sup>10</sup>

## Enumerating Molecules

Enumerating molecules is not only the main subject of this chapter, but it has also been a prolific field of research for decades. Rather than reviewing every approach that has so far been taken, we have chosen to present examples of orderly generation. Our reasons are many. First, as discussed earlier, orderly generation is the most elegant technique to enumerate graphs. Second, no other technique has had as many applications in chemistry than has orderly generation. Finally, focusing on one technique will help the reader understand how molecules are enumerated. As we shall see in all subsections that follow, the main problem in applying orderly generation to a specific class of molecules is to find the appropriate canonical code. That is, a code that uniquely

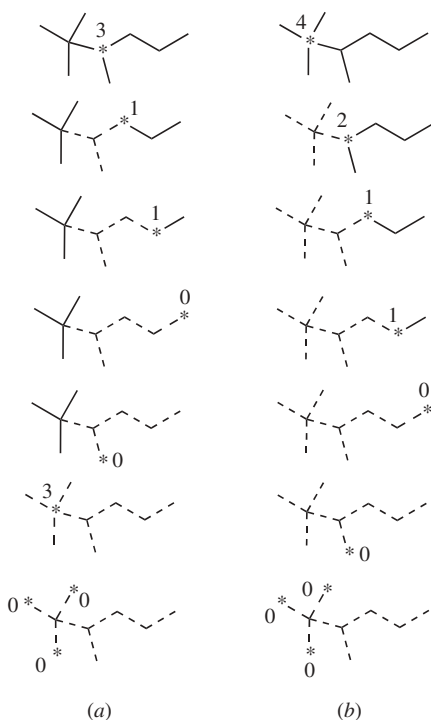
represents the class of molecules we want to enumerate, and a code that is easily computable, ideally, in polynomial time.

### *Acyclic Molecular Graph Enumeration*

As with counting, it is more simple to enumerate acyclic structures than cyclic ones. For this reason, the field of molecular structure enumeration started with acyclic hydrocarbons with an algorithm published by Nobel Laureate J. Lederberg.<sup>54</sup> The algorithm was later integrated into a code named DENDRAL and enumerated the isomers for a variety of acyclic compounds containing C, H, O, and N atoms.<sup>55</sup> Much could be said about the DENDRAL project, which is described in many computer science textbooks as the first expert system. The reader further interested by DENDRAL is referred to the books by Lindsay et al. and Gray,<sup>56,57</sup> where the history of the project is reviewed. A decade after the initial DENDRAL effort, a powerful approach appeared based on the  $n$ -tuple code developed by Knop et al.<sup>58</sup> We present this technique in the context of an orderly algorithm.

The  $n$ -tuple code is a set of non-negative integers smaller than  $n$ , the number of atoms of an acyclic molecular structure. Each number in the  $n$ -tuple represents the degree of an atom in the structure or in one of its substructures. To compute the  $n$ -tuple of a structure, we first choose a starting atom (a root) as illustrated in Figure 11(a). For the purpose of this example, any atom will do, but as we shall see later, the root atom must be the atom with the highest degree if we are to construct a canonical represent of the  $n$ -tuple. The first element of the tuple is  $k$ , the degree of the root. Next, the root and all bonds attached to it are removed from the structure, thus creating  $k$  disconnected substructures. The process is repeated for each of the  $k$  substructures where the new roots are the atoms that were bonded to the initial root. The process stops when all atoms have been removed.

Looking at Figure 11(a), it is obvious the  $n$ -tuple code is nothing else but a list of atom degrees obtained by reading the structure in a depth-first order. All degrees are reduced by one except for the initial root. Now, for any given rooted structure, a canonical  $n$ -tuple (cf. Figure 11(b)) is computed with the procedure described here, but at each step, the tuples associated with the substructures are sorted and read in decreasing lexicographic order. Finally, to compute a canonical  $n$ -tuple for an unrooted structure, we compute the canonical  $n$ -tuples for all structures rooted at atoms with the highest degree while keeping the lexicographically maximal tuple as the canonical representation for the structure. Note that no need exists to compute  $n$ -tuples rooted on atoms with degrees smaller than the maximum one, as these rooted structures produce lexicographically smaller  $n$ -tuples. The code corresponding to Figure 11(b) is the canonical  $n$ -tuple of 2,2,3-trimethylhexane because only one quaternary carbon is in the structure. As shown by Hopcroft and Tarjan,<sup>59</sup> this canonization procedure can be implemented with an  $O(n)$  time complexity. Finally it is worth mentioning that modifications of the  $n$ -tuple code have

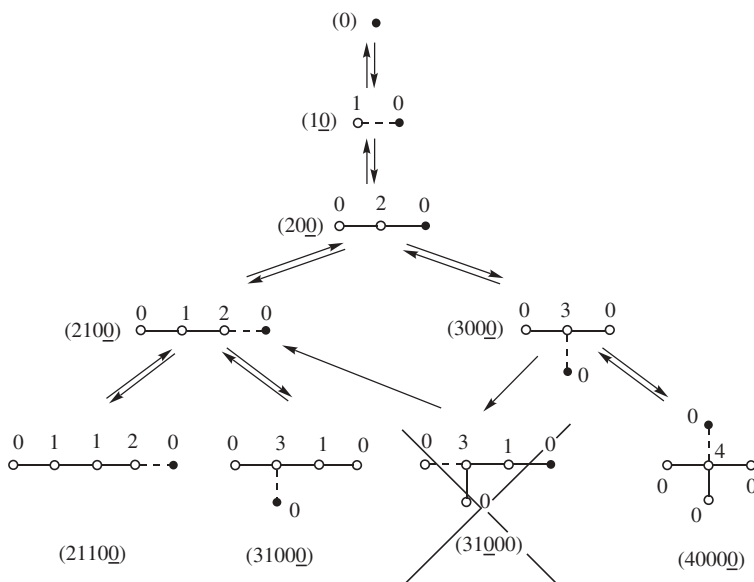


**Figure 11** Some  $n$ -tuple codes for 2,2,3-trimethylhexane. Successive roots are indicated with a “\*” symbol. (a) The code is 311003000. (b) The code is 421100000, and this code is canonical.

been proposed to take into account atom and bond types.<sup>60</sup> Instead of just writing the degree of the atoms in the  $n$ -tuple, we also include atom types and bond orders.

Now that we have a code to canonize acyclic structures we can use an orderly algorithm. Next, we illustrate how the  $n$ -tuple code enumerates alkanes up to  $n$  carbon atoms with a McKay-type orderly generation (Scheme III). For simplicity, all hydrogen atoms are ignored, and carbon atoms may thus have a number of bonds ranging between 1 and 4. As depicted in Figure 12, the initial graph contains one atom and no bond, so its canonical  $n$ -tuple is (0).

Following Scheme III, we first verify that the construction process is completed for structure  $G$  (step 1). In this case, we must check that the required number of atoms  $n$  is met and that the number of bonds attached to every atom ranges between 1 and 4. If structure  $G$  passes the completion test, it is printed (step 2); otherwise, we augment  $G$  in all possible ways by adding a bond  $e$  and a new atom (step 5). Augmentations violating the maximum valence requirement are rejected (step 6). For all other  $G \cup e$  structures, a



**Figure 12** The three pentane isomers obtained with McKay's orderly algorithm and the  $n$ -tuple code. Hydrogen atoms are not represented. All atoms are carbons and can have up to four bonds. Parent-child and child-parent relationships are indicated with arrows. Canonical  $n$ -tuples are written in parentheses. At each layer, a bond and a new atom are added. The added atom is represented by a solid node. The last bond/atom in the canonical  $n$ -tuple is represented by a dashed line and is underlined in the canonical  $n$ -tuple. A graph is rejected when its legitimate parent is not the graph it came from. This case develops when the added bond/atom is not the last digit of the canonical  $n$ -tuple (the dashed line is not linked to the solid node).

canonical  $n$ -tuple  $G'$  is constructed and  $G'$  is added to the set of  $n$ -tuples  $S$  (step 7). Duplicated  $n$ -tuples are removed (step 10). For each resulting  $n$ -tuple  $G'$  in  $S$ , McKay's algorithm removes the last edge of  $G'$  (step 12), which in this case is the last digit of the  $n$ -tuple. If the resulting  $n$ -tuple equals the  $n$ -tuple of the initial graph  $G$  (step 13), then  $G'$  is a legitimate child of  $G$ , and the process repeats with  $G'$  (step 14); otherwise  $G'$  is an illegitimate child and is ignored.

The application of Scheme III to generate alkane structures up to pentane is illustrated in Figure 12, in which examples of legitimate and illegitimate parent-child relationships are depicted. Of course, Figure 12 could be expanded up to decane, and we could then answer the section question "Are there any isomers of decane having seven methyl groups?" As we shall see later, more efficient ways are available to enumerate all decane isomers having seven methyl groups.

The  $n$ -tuple technique has lead to numerous implementations and extensions. In particular, Contras et al. extended the  $n$ -tuple enumeration algorithm

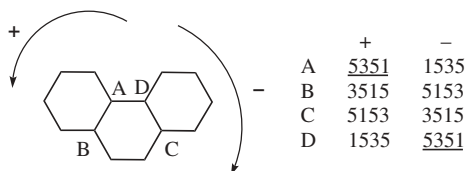
in a series of papers to acyclic compounds with heteroelements and multiple bonds,<sup>60</sup> cyclic structures,<sup>61</sup> mixed compounds,<sup>62</sup> acyclic stereoisomers,<sup>63</sup> and unsaturated stereoisomers.<sup>64,65</sup> We should also mention the tree enumeration technique proposed by Lukovits.<sup>66</sup> Instead of an  $n$ -tuple, Lukovits uses a compressed adjacency matrix (CAM). The CAM is a vector in which each element  $e_i$  represents a column  $i$  of the adjacency matrix ( $a_{ij}$ ). The value of element  $e_i$  is the row number  $j < i$  for which a bond appears, i.e.,  $a_{ij} > 1$ . Lukovits proposes a set of rules to generate all trees having a maximal CAM.<sup>67</sup> The technique may not be as efficient as the  $n$ -tuple code because during the construction process, many structures do not meet the rules and are thus rejected.

### *Benzenoids and Polyhex Hydrocarbons Enumeration*

The reader is referred to the "Number of Benzenoids and Polyhex Hydrocarbons" subsection for the definition and classification of benzenoids and polyhex hydrocarbons, as well as for additional references for this class of compounds, which is only partially reviewed because of space limitations. Let us recall that the direct counting approach has difficulties with molecules that cannot be represented by tree-like structures, such as pericondensed polyhexes. Furthermore, the counting approach cannot separate nonplanar polyhexes (helicenes) from planar benzenoids. Consequently, for benzenoids and polyhexes, enumeration is not only a valuable tool that provides a concise description of the structures being enumerated, but also enumeration computes isomer numbers that cannot be derived otherwise.

The first algorithm to enumerate polyhexes was proposed by Balasubramanian et al.<sup>68</sup> The enumeration of planar simply connected polyhexes to  $h = 10$  hexagons,<sup>69</sup>  $h = 11$ ,<sup>70</sup> and  $h = 12$ <sup>71</sup> applied this algorithm. The next advance in polyhex enumeration came from a code based on the dual graph associated with every polyhex.<sup>72</sup> This code allowed enumeration of all polyhexes for  $h = 13$ ,<sup>73</sup>  $h = 14$ ,<sup>74</sup>  $h = 15$ ,<sup>72</sup> and  $h = 16$ .<sup>75</sup> The next progress was made by Tosić et al.,<sup>76</sup> who proposed a lattice-based approach with a "cage," within which the polyhexes are placed. This method led to enumeration of all polyhexes with  $h = 17$ .<sup>76</sup> Three years later, Caporossi and Hansen<sup>77</sup> developed a McKay-type orderly algorithm and enumerated polyhexes up to  $h = 21$  and  $h = 24$ .<sup>78</sup> Finally, in 2002, another lattice-based method was proposed and polyhexes were enumerated up to  $h = 35$ .<sup>79</sup> Next, we briefly describe the orderly generation and the lattice enumeration approaches.

*Orderly generation of polyhexes.* As usual with orderly generation algorithms, polyhexes comprising  $h$  hexagons are constructed from polyhexes having  $h-1$  hexagons. To avoid repetitions, each polyhex with  $h$  hexagons is generated from one and only one parent, i.e., a polyhex with  $h-1$  hexagons. As we observed with alkanes in Figure 12, once a structure is generated from a potential parent, its canonical code must be scanned to verify if the parent is legitimate. To apply Scheme III to polyhexes, we only have to find the appropriate canonical code. One possible code for this purpose is the



**Figure 13** BEC code. Canonical codes are underlined. Starting at vertex A and turning clockwise, we first encounter one edge from the center face, then we find five edges belonging to the right face, next are the three edges from the center face, and finally are the five edges belonging to the left face. Turning clockwise, the BEC code starting at A is 1535.

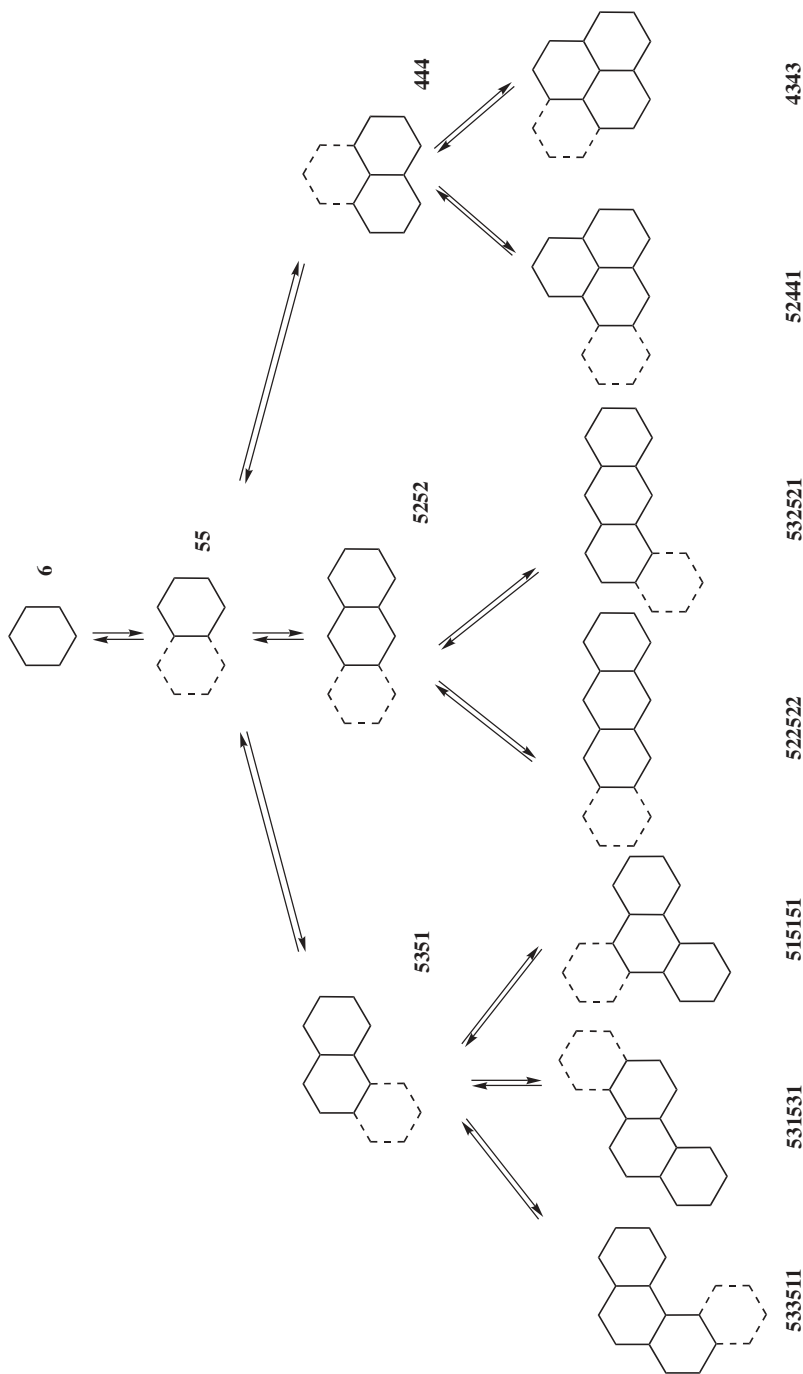
boundary edges code (BEC).<sup>77</sup> This code is outlined next and illustrated in Figure 13.

Beginning at any external vertex of degree three, which thus belongs to only two hexagons, follow the boundary of the polyhex noting by a digit the number of edges on the boundary for each successive hexagon. The procedure is repeated clockwise and counterclockwise, the canonical code is the lexicographically maximum code. In Figure 13, we observe that the code is unique but may be obtained in several ways in case of symmetry of the polyhex. The high efficiency of the BEC code is caused by an alternative way to check whether a polyhex must be considered as being legitimate. To this end, Caporossi and Hansen<sup>77</sup> established the following rule: A polyhex is legitimate if and only if the first digit of its BEC code corresponds to the last added hexagon. This simple rule induces the enumeration tree illustrated in Figure 14 up to  $h = 4$ . Note that the cost of determining whether a polyhex is legitimate equals the cost of computing the BEC code,  $O(h^2)$ . Caporossi and Hansen<sup>77</sup> assessed the computational time per output of their algorithm, and it appears to increase quadratically with the system size.

*Lattice enumeration of benzenoids.* With lattice enumeration techniques, we make use of the fact that only eight symmetry groups are associated with benzenoids.<sup>35</sup> These groups are (1)  $C_s$  for benzenoids of  $h$  hexagons with no rotational or reflection symmetry, (2)  $C_{2v}$  for those with one axis of reflection symmetry, (3)  $C_{2b}$  for those invariant with respect to rotations through  $\pi$ , (4)  $D_{2b}$  for those with two axes of reflection symmetry and invariant with respect to rotations through  $\pi$ , (5)  $C_{3b}$  for those invariant with respect to rotations through  $2\pi/3$ , (6)  $D_{3b}$  for those with three axes of reflection symmetry and invariant with respect to rotations through  $2\pi/3$ , (7)  $C_{6b}$  for those invariant with respect to rotations through  $\pi/3$ , and finally, (8)  $D_{6b}$  for those with six axes of reflection symmetry and invariant with respect to rotations through  $\pi/3$ . For these, the number of benzenoids  $b_h$  comprising  $h$  hexagons may be written as

$$b_h = C_s^{(h)} + C_{2v}^{(h)} + C_{2b}^{(h)} + D_{2b}^{(h)} + C_{3b}^{(h)} + D_{3b}^{(h)} + C_{6b}^{(h)} + D_{6b}^{(h)} \quad [29]$$





**Figure 14** The seven polyhexes with four hexagons obtained with orderly generation and BEC code. At each layer, the last added hexagon (dashed lines) corresponds to the first digit in the BEC code.

where, for instance,  $C_s^{(h)}$  is the number of benzenoids of  $h$  hexagons with symmetry  $C_s$ . Now, let  $B_h$  be the number of fixed hexagonal systems. Fixed hexagonal systems are simply all possible benzenoids we can construct on a hexagonal lattice disregarding rotational and reflection symmetries. From the definitions of symmetry groups, provided here it is easy to verify that

$$B_h = 12C_s^{(h)} + 6C_{2v}^{(h)} + 6C_{2h}^{(h)} + 3D_{2h}^{(h)} + 4C_{3h}^{(h)} + 2D_{3h}^{(h)} + 2C_{6h}^{(h)} + D_{6h}^{(h)} \quad [30]$$

Eliminating  $C_s^{(h)}$ , we arrive at

$$b_h = \frac{1}{12}(B_h + 6C_{2v}^{(h)} + 6C_{2h}^{(h)} + 9D_{2h}^{(h)} + 8C_{3h}^{(h)} + 10D_{3h}^{(h)} + 10C_{6h}^{(h)} + 11D_{6h}^{(h)}) \quad [31]$$

The lattice enumeration technique consists of generating and counting all hexagonal systems that appear on the right-hand side of Eq. [31] to evaluate  $b_h$ . Let us start with  $B_h$ , the number of fixed hexagonal systems of size of  $h$ . Generating fixed polygonal systems on lattices can be solved by enumerating self-avoiding polygons on lattices. This problem has been studied in the physics literature and will not be reviewed here. The reader interested by this particular problem is referred to the work of Enting and Guttmann.<sup>80</sup> To enumerate benzenoids, Vöge et al.<sup>79</sup> use the Enting and Guttmann technique, whereas Tosic et al.<sup>76</sup> use an original algorithm based on a brute force approach enumerating all fixed hexagonal systems on a lattice.

Once  $B_h$  has been computed, the other terms of Eq. [31] are derived as follows. We first consider the elements of  $C_{2v}^{(h)}$ . Each element of  $C_{2v}^{(h)}$  can be decomposed into two identical  $h/2$  hexagonal systems, joined together at the symmetry axis. Thus, the elements of  $C_{2v}^{(h)}$  can be generated from the elements of  $B_{h/2}$ . Similar arguments apply to the elements of  $C_{2h}^{(h)}$ . From the definitions of the symmetry groups given previously, it is easy to verify that the elements of  $D_{2h}^{(h)}$  can be generated from the fixed hexagonal systems of  $B_{h/4}$ , the elements of  $C_{3h}^{(h)}$  from  $B_{h/3}$ , the elements of  $D_{3h}^{(h)}$  and  $C_{6h}^{(h)}$  from  $B_{h/6}$ , and the elements of  $D_{6h}^{(h)}$  from  $B_{h/12}$ . Thus, all elements in Eq. [31] can be computed from  $B_h$ . In other words, benzenoids can be counted and enumerated from the enumeration of fixed hexagonal systems.

Results obtained with this approach as well as the orderly generation technique have been compiled in Table 7 in the “Chemical Information” subsection appearing later in this chapter.

### *Molecular Cages Enumeration (Fullerenes and Nanotubes)*

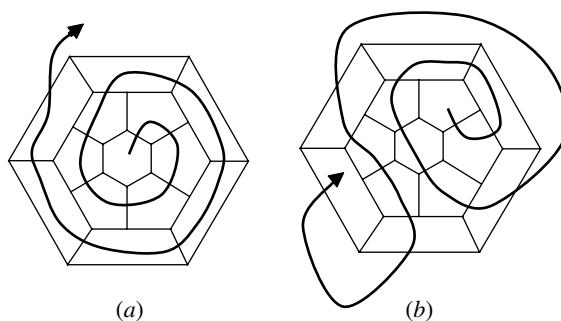
Fullerenes, nanotubes, spheroalkanes, and other molecular cages belong to the class of regular graphs. A regular graph is a graph in which all vertices have the same degree. Among the class of regular graphs of interest in

chemistry are  $(k,g)$ -cages in which all atoms have the same valence  $k$  and all rings are at least of size  $g$ . We first review the literature for regular graphs and cages and then describe algorithms specifically designed for fullerenes.

*Regular graphs and cages.* Enumerating regular graphs is one of the oldest problem in combinatorics. In the nineteenth century, Jan de Vries<sup>81</sup> enumerated all 3-regular graphs, also named cubic graphs, up to ten vertices. The first computational approach is due to Balaban,<sup>82</sup> who in 1966, enumerated all cubic regular graphs up to 10, and later 12 vertices.<sup>82</sup> In 1976, Bussemaker et al.<sup>83</sup> computed all cubic graphs up to 14 vertices. About the same period, Faradzev<sup>51</sup> worked out the case for 18 vertices when he suggested the general orderly algorithm presented in Scheme II. In 1986, McKay and Royle settled the case for 20 vertices,<sup>84</sup> whereas in 1996, Brinkmann<sup>85</sup> enumerated all 24 vertices cubic graphs and  $(3,8)$  cages up to 40 vertices. Finally in 1999, based of the Brinkmann technique, Meringer enumerated all  $k$ -regular graphs up to  $k=6$  and a number of vertices ranging between 15 and 24.<sup>86</sup> Meringer's orderly algorithm is an integral part of the latest version of the MOLGEN isomer generator.<sup>87</sup> Next, we describe this algorithm, which a classic example of the Read-Faradzev orderly generation.

Meringer's algorithm generates all  $k$ -regular graphs of  $n$  vertices. The process starts with an initial graph,  $G$ , composed of  $n$  vertices labeled 1 through  $n$  and no edges. Meringer's algorithm is recursive, thus, following scheme II, in steps (1) and (2), the graph is printed if it is fully constructed. That is, if all  $n$  vertices have  $k$  neighbors. When the graph is not fully constructed, in step (4), all edges,  $e$ , are enumerated only when they are lexicographically greater than the edges built so far. In steps (5) and (6), the algorithm checks if the graph  $G \cup e$  obtained for each enumerated edge,  $e$ , is identical to its canonical representation, i.e.,  $G \cup e = \text{CAN}(G \cup e)$ . When the graph is canonical and the additional constraints of step (5) are verified, the same process is repeated; the algorithm backtracks otherwise. The main constraint in step (5) is regularity. All vertices must have at most  $k$  neighbors, and supplementary constraints such as connectivity and minimum cycle size (girth) may also be added. According to its author, the most time-consuming part of the algorithm is the canonization step. To reduce the number of times graphs are canonized, not all possible edges are enumerated in step (4), but only the edges attached to the lexicographically smallest vertex having less than  $k$  neighbors.

*Fullerenes and nanotubes.* A fullerene is a spherically shaped carbon molecule composed exclusively of five- and six-membered rings. In the language of graph theory, a fullerene is a 3-regular spherical map having pentagonal and hexagonal faces only. Furthermore, by definition, any fullerene  $C_n$ ,  $n \geq 20$ , has exactly 12 pentagons and  $n/2 - 10$  hexagons. Because of these restrictions, we cannot use here the polyhexes, benzenoids, and regular cages generators presented earlier. For instance, the BEC canonical code cannot be applied because fullerenes do not have edges on their boundaries. In the early algorithms that enumerate fullerenes,<sup>88-90</sup> orderly generation was not used.



**Figure 15** Spiral codes for  $C_{24}$ . (a) Starting at a hexagonal face the code is  $65555555555556 = 65^{12}6$ . (b) Starting at a pentagonal face the code is  $55555655655555 = 5^565^265^5$ . Code (b) is canonical.

Yet, no reasons exist why orderly generation could not potentially be applied, provided that a canonical code exists to uniquely identify fullerenes. Next we describe the spiral canonical code for fullerenes,<sup>91–93</sup> and we then propose a sketch of a Read–Faradzev orderly generation taken from the algorithms of Fowler and Manolopoulos<sup>93</sup> and Brinkmann.<sup>85</sup>

The spiral canonical code for a  $C_{24}$  fullerene is illustrated in Figure 15. Starting at one face, choose a first neighboring face and an orientation (clockwise or counterclockwise). Visit all faces of the fullerene by recursively choosing a new face as the next one to be visited. The next face must not have already been visited, and it must be adjacent to the last face visited. Additionally, the next face is the first one encountered running around the last face in a clockwise (counterclockwise) direction from the intersection with the next-to-last face. The code is simply the sequence of face sizes in the order they are visited. The process is repeated choosing all faces one after another as the starting one, choosing all possible first neighbors, and choosing the two possible directions. The lexicographically minimum code is the canonical one. The major pitfall of the spiral code is that not all fullerenes admit ring spirals;<sup>91</sup> however, this problem can be overcome by identifying the edges adjacent to consecutive faces and adding these identifiers to the spiral code.<sup>92</sup>

Now that we have a way to canonize fullerenes, we construct fullerenes adding pentagonal or hexagonal faces one at a time starting with a pentagonal face; otherwise the final spiral codes would not be canonical. In other words,  $n$  digits spirals (i.e.,  $n$  faces fullerenes) are constructed from  $n - 1$  digits spirals (i.e.,  $n - 1$  faces fullerenes) by appending to the code either a 5 or a 6. Let  $s_{n-1}$  be a  $n - 1$  digits spiral code; the child  $s_n = s_{n-1}5$  ( $s_{n-1}6$ ) is legitimate if the canonical spiral code of that child is indeed  $s_{n-1}5$  ( $s_{n-1}6$ ), and if not, the child is rejected. It is important to realize some spiral codes do not lead to final fullerenes at all. For instance, starting with eleven 5's in the code, i.e., 11 pentagonal faces, we can see that this code cannot lead to a fullerene unless we have

no hexagon and the structure to be constructed is  $C_{20}$ . Consequently, the orderly generation applied to fullerene creates unproductive branches in the enumeration tree.

Faster than this algorithm is the technique proposed by Brinkmann et al.<sup>94–96</sup> Instead of building fullerenes from the ground up, this algorithm generates structures by gluing together “benzenoid” patches composed of five- and six-membered rings. This approach was taken because fullerenes can be decomposed into either two or three patches following a Petrie path. Petrie paths are constructed as follows: Start at any edge  $e_1$  in the fullerene and with a scissor cut that edge. Next cut edge  $e_2$  on the right side of  $e_1$ , cut edge  $e_3$  on the left side of the  $e_2$ , and repeat the process turning alternately right and left until you reach an edge  $e_k$  that has already been cut. If  $e_k = e_1$ , you have separated the fullerene into two patches. Now, if  $e_k \neq e_1$ , the fullerene is also separated into two parts, but the job is not completed because one part is partially cut, i.e., the part containing edges  $e_1, e_2, \dots, e_{k-1}$ . Take that part and start again at  $e_1$ , but now cut in the opposite direction; you will eventually split the part into two patches, and create a total of three patches. Because any fullerenes can be decomposed into at most three patches, from a given number  $h$  of hexagons, all fullerenes can be constructed by attaching in all possible ways a catalogue of all patches composed of at most  $h$  hexagons and 12 pentagons. Results obtained with this algorithm can be found in Table 8 in the “Chemical Information” subsection appearing later in the chapter.

Before closing this subsection, we should also mention a simple algorithm that enumerates the isomers of a toroidal polyhex.<sup>97</sup> Toroidal polyhexes are fullerenes embedded on the surface of a torus. The word “fullerene” is not appropriate here because only structures having six-membered rings (not five) are enumerated. This limitation greatly reduces the number of solutions. The number of isomers is found to increase at only a modest rate that does not exceed 30% of the number of atoms.

### *General Structural Isomer Enumeration*

By general structural isomer enumeration, we mean the enumeration of all molecular graphs corresponding to a molecular formula. We do not include here solutions that construct molecular structures from additional constraints, such as the presence or the absence of substructural fragments. Enumeration with constraints is reviewed in the next subsection.

Techniques to enumerate molecules (including cyclic ones) from a molecular formula appeared in the 1970s. The first algorithm to do so, CONGEN,<sup>98</sup> was a product of the DENDRAL project. The solution consisted of decomposing the molecular formula into cyclic substructures, which were combined by bridges to get molecules. The cyclic substructures were built from a database of 3000 elementary cycles. A second approach, simpler in principle, has been the technique chosen by the researchers involved in the CHEMICS project.<sup>99</sup> In this approach, only canonical structures are

generated. However, orderly generation was not applied in the earlier version of CHEMICS. Instead, all labeled structures were generated and noncanonical ones were rejected. A similar approach was also taken by the authors who developed the ASSEMBLE generator,<sup>100</sup> although this code was designed to combine fragments. Since these initial developments, CONGEN, CHEMICS, and ASSEMBLE have lead to numerous improvements, most of which involve enumeration with constraints.

Another development to enumerate isomers has been a method based on an atom's equivalent classes. In this method pioneered by Bangov,<sup>101</sup> and generalized by Faulon,<sup>102</sup> the atoms corresponding to the molecular formula are partitioned into equivalent classes. Next, a class of atom is selected and all atoms of the class are saturated; that is, bonds are added until each selected atom has a number of bond equals to its valence. Atom saturation is performed in all possible ways, and to avoid generating isomorphic structures, noncanonical graphs are rejected. For each resulting graph, equivalent classes are computed again, a new unsaturated class is chosen, and the process is repeated until all atoms are saturated. With the equivalent-classes technique, we can choose the atoms to be saturated. Thus, we can drive the process to first build tree-like structures, choosing classes of atoms that do not create cycles when being saturated, and then create cycles adding bonds to the unsaturated atoms of the trees. The advantage of building tree-like structures first is that we can canonize them efficiently using, for instance, the  $n$ -tuple code mentioned earlier. For acyclic isomers, the equivalent-classes algorithm is efficient because canonization can be performed in linear time. However, for all other compounds, the cost of canonization has to be factored in.

The next approach to enumerate isomers is orderly generation. One of the first algorithms is due to Kvasnicka and Pospichal.<sup>103</sup> Their orderly technique is based on Faradzev's algorithm. The proposed solution constructs all molecular graphs of maximum valence matching given numbers of atoms and bonds. The technique was soon modified to enumerate all molecular graphs matching a prescribed valence sequence.<sup>104</sup> Faradzev's orderly generation was also involved in developing the SMOG program that enumerates compounds from molecular formulas with fragments.<sup>105,106</sup> The isomer generator MOLGEN<sup>87,107,108</sup> is also based on orderly generation.

The latest development with isomer enumeration is the method of homomorphisms proposed by Grüner et al.<sup>87</sup> Interestingly, the homomorphism method is a systematization of the early solution developed within the DENDRAL project. The homomorphism method has been implemented in the latest version of MOLGEN.<sup>98</sup> The enumeration relies on a strategy of determining how all molecular graphs with a given valence sequence can be built up recursively from regular graphs. Grüner et al. observe that any molecular graph  $G$  can be decomposed into two subgraphs:  $T$ , a subgraph comprising all atoms of a fixed valence, for instance, the largest valence, and  $H$ , a subgraph comprising the remaining atoms. Attached to the two subgraphs

an incidence structure,  $I$ , is constructed such that each column corresponds to an atom  $t$  of  $T$ , each row to an atom  $h$  of  $H$  and noting a bond connecting two atoms  $t$  and  $h$  by the entry 1 in the corresponding place of  $I$ . The authors then prove that all possible valence sequences for  $T$  and  $H$  and all possible numbers of entries 1 in each row and each column of  $I$  can be determined directly from the valence sequence of  $G$ . This decomposition of the valence sequence is repeated recursively until all resulting valence sequences correspond to regular graphs. The strategy obviously reduces the construction problem of molecular graphs with prescribe valence sequences to that of regular graphs and the problem of pasting the subgraphs  $T$  and  $H$  together. Regular graphs are constructed with Meringer's algorithm,<sup>86</sup> and all possible ways of pasting  $T$  and  $H$  are enumerated with an orderly algorithm. According to the authors, the resulting algorithm is fast as it has determined up to  $10^{30}$  molecular graphs (without actually constructing them) corresponding to valence sequences up to 50 atoms.

### *Molecular Graph Enumeration with Constraints*

Molecular structure enumeration subjected to constraints has practical application in structure elucidation and molecular design. Many codes have been developed to address these two applications; most of them can be found in the section entitled "Enumerating Molecules: What Are the Uses". For structure elucidation, the constraints are generally composed of fragments that must be present and/or absent in the final solutions. With molecular design, the goal is to generate all structures matching a specified property or activity. This problem, also named inverse imaging, is generally solved in a two-step procedure. First, from the target property or activity, a molecular descriptor is computed, which is usually performed through a quantitative structure-activity relationship in which the molecular descriptors are fragments or topological indices. In a second step, all structures matching the descriptor value are enumerated. We next present the methods that have been developed for structure elucidation and molecular design purposes.

*Enumerating structures with molecular fragments.* We first consider the simple case in which the molecular fragments do not overlap. Each fragment must be unsaturated and, thus, contain some free bonds or bonding sites. Then the problem consists of connecting the bonding sites in all possible ways. This process can be solved by generating all possible labeled graphs in which the vertices are bonding sites. Duplicates can be eliminated in a postprocess,<sup>109</sup> or noncanonical graphs can be rejected as they are generated such as in ASSEMBLE<sup>110</sup> and CHEMICS.<sup>111</sup> Another solution is the equivalent-classes algorithm, in which the equivalent classes are computed only for the unsaturated atoms, and of course, only these atoms are saturated.<sup>102</sup> Orderly generation can and has enumerated structures from fragments.<sup>105,112</sup> During the orderly process, the search for canonical structures is performed without permuting the elements of the adjacency matrices that correspond to the

fragments. Any of the aforementioned algorithms can answer the section question, “are there any isomers of decane having exactly seven methyl groups?” All solutions (if any) must contain seven methyl groups. As the final structure has the molecular formula  $C_{10}H_{22}$ , the additional fragments are three carbon atoms and one hydrogen atom. These 11 fragments were given as input to the equivalent-classes algorithm. The code returned two solutions: (2,2,3,4,4)-pentamethyl-pentane and (2,2,3,3,4)-pentamethyl-pentane.

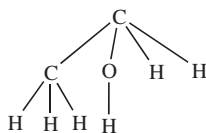
In most structure elucidation instances, fragments unfortunately overlap. For instance, consider the fragments provided by  $^{13}C$  NMR spectra. To each  $^{13}C$  NMR peak, a corresponding fragment exists (the environment of a  $^{13}C$  carbon atom) and two neighboring atoms in the probed structure have corresponding overlapping fragments. The problem of overlapping fragments can be addressed with manual intervention as in GENOA,<sup>113</sup> another product of the DENDRAL project. At first the code’s user selects one fragment as a core. The user then chooses a second fragment and the code generates all possible ways of breaking those two fragments into nonoverlapping, ever smaller fragments. The process is repeated until all fragments have been decomposed into nonoverlapping ones. Final structures are then generated assembling the nonoverlapping fragments with a technique similar to those we just presented.

More systematic is the approach taken with the EPIOS code.<sup>114,115</sup> A large database of assigned  $^{13}C$  NMR spectra is the source of a library of carbon-centered fragments to which are assigned chemical shifts and signal multiplicities. With the experimental spectrum, fragments are extracted from the database and the construction proceeds by attaching carbon atoms only if their fragments overlap. Partially assembled structures with chemical shift deviations that exceed a preset threshold are discarded. Once the structures are fully assembled, a spectrum prediction code is run and the predicted spectrum is checked against the experimental one. Structure assembly with overlapping information is also the method implemented in the SpecSolv system.<sup>116</sup>

Another method dealing with overlapping fragments was devised with the so-called *signature equation*.<sup>9</sup> The signature of an atom is a fragment comprising all atoms and bonds that are at a specified distance  $h$  from the probed atom. The fragment is written as a tree with a height equal to the specified distance, the tree is canonized, and the signature is written reading the tree in a depth first order. Examples of signatures of various heights are given in Figure 16.

The signature of a molecule or a molecular fragment is simply the sum of its atomic signatures. The signature of a bond is the difference between the signature of the structure containing the bond and the signature of the structure where the bond has been removed. Now, assuming we know the signature up to a certain height of a yet unresolved compound and assuming we know that the compound contains a number of fragments that may or not overlap, the purpose of the signature equation is to compute lists of nonoverlapping fragments matching the signature of the unresolved compound. Simply stated,





**Figure 16** The fragment centered on the carbon atom attached to the alcohol group in ethanol is depicted. The height-0 signature of this carbon atom is  ${}^0\sigma(\text{C}) = \text{C}$ , the height 1 signature is  ${}^1\sigma(\text{C}) = \text{C}(\text{COHH})$ , and the height 2 is  ${}^2\sigma(\text{C}) = \text{C}(\text{C}(\text{HHH})\text{O}(\text{H})\text{HH})$ . The height 1 signature of ethanol is obtained by summing the height 1 signatures for all atoms,  ${}^1\sigma(\text{ethanol}) = \text{C}(\text{COHH}) + \text{C}(\text{CHHH}) + \text{O}(\text{CH}) + 5\text{H}(\text{C}) + \text{H}(\text{O})$ . The height 1 signature of the bond C–O is the difference between the signature of ethanol and the signature of the structure where the bond has been removed,  ${}^1\sigma(\text{C–O}) = \text{C}(\text{COHH}) + \text{O}(\text{CH}) - \text{C}(\text{CHH}) - \text{O}(\text{H})$ .

fragments and signatures are related by the expression: signature of the fragments + signature of the interfragment bonds = signature of the unknown compound. Formally, lists of nonoverlapping fragments are computed solving the equation with unknowns  $x_i$  and  $y_j$ :

$$\sum_i x_i {}^b\sigma(\text{fragment } i) + \sum_j y_j {}^b\sigma(\text{bond } j) = {}^b\sigma(\text{unknown compound}) \quad [32]$$

The variables  $x_i$  and  $y_j$  are, respectively, the number of fragments  $i$  present in the final structure and the number of interfragment bonds  $j$ . The signature equation (Eq. [32]) is an integer equation and can be solved using integer linear programming (ILP) tools.<sup>117</sup> Note, however, that in general ILP problems are intractable.<sup>12</sup> Each solution of Eq. [32] is a list of nonoverlapping fragments and interfragment bonds. To enumerate the final structures, each list of fragments and interfragment bonds is fed to an isomer generator working with nonoverlapping fragments. In the structure elucidation instances in which the signature equation was applied<sup>118</sup> elemental analysis, NMR, and functional group analysis provided the height 0, 1, and 2 signatures of the unknown compounds, and fragments were derived from chemical degradation and pyrolysis.

An elegant approach dealing with overlapping fragments is the structure reduction method proposed by Christie and Munk.<sup>119</sup> In contrast with all enumeration algorithms we have presented so far, this method begins with a hyperstructure containing all possible bonds between unsaturated atoms. The algorithm removes inconsistent bonds until valences of atoms are respected. This results in a more efficient way to deal with overlapping fragments because all fragments are contained (i.e., are subgraphs) in the hyperstructures, and as bond deletion occurs, the resulting graphs are kept if they still contain the fragments and are rejected otherwise. Although it is not clear

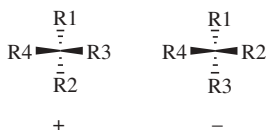
from reading the original paper on structure reduction how duplicated structures are removed, orderly generation can certainly avoid the production of duplicates. Checking that fragments occur in a given structure requires running a subgraph isomorphism routine. As already stated, general subgraph isomorphism is an intractable problem.<sup>12</sup> In a recent development, the structure reduction method was coupled with a convergent structure generation technique.<sup>120,121</sup> In this technique, instead of having a list overlapping fragments, a network of substructures is first constructed. Substructures are linked in this network when they overlap, and alternative neighborhoods are indicated when overlapping is ambiguous. The initial structure is a hyperstructure composed of all possible bonds between atoms. The reduction method determines all possible ways in which the substructures of the network can be mapped to the actual atoms of the structure being constructed.

*Enumerating structures with molecular descriptors.* Enumerating molecules matching molecular descriptors or topological descriptors is a long-standing problem. Surprisingly, few reports in the literature provide answers to the question. Most of the proposed techniques are stochastic in nature and are reviewed in the "Sampling Structures" section. In a series of five papers, Kier, et al.<sup>122–126</sup> reconstruct molecular structures from the count of paths up to length  $l = 3$ . Their technique essentially computes all possible valence sequences matching the count of paths up to length  $l = 2$ . Then, for each valence sequence, all molecular structures are generated with a classic isomer generator (cf. General Structural Isomer Enumeration subsection), and the graphs that do not match the path length  $l = 3$  count are rejected. Skvortsova et al.<sup>127</sup> use a similar technique, but from the count of paths, they derive a bond sequence in addition to the valence sequence. A bond sequence counts the number of bonds between each distinct pair of atom valences. The two sequences are then fed to an isomer generator that produces all structures matching the sequences. Regrettably, the authors do not provide details on how the isomer generator deals with the bond sequence. Another approach to enumerate molecular graphs matching a given signature has been introduced recently.<sup>128</sup> As defined earlier, the signature is the collection of all atoms' environments in a molecule (cf. Figure 16). Like other fragmental molecular descriptors, signature works well in quantitative structure-activity relationships.<sup>129</sup> The input information to the algorithm is a signature. To each atomic signature, we associate an atom in the initial graph. At first, the graph is composed of isolated atoms without any bond. The construction proceeds by adding bonds one at a time with the equivalent-classes technique (cf. General Structural Isomer Enumeration subsection). Orderly generation can also enumerate structure matching signatures. During the generation process, bonds are created only if the signatures of the bonded atoms are compatible and the resulting graph is canonical. This algorithm is capable of enumerating molecular structures up to 50 non-hydrogen atoms on a time scale of few CPU seconds.<sup>128</sup>

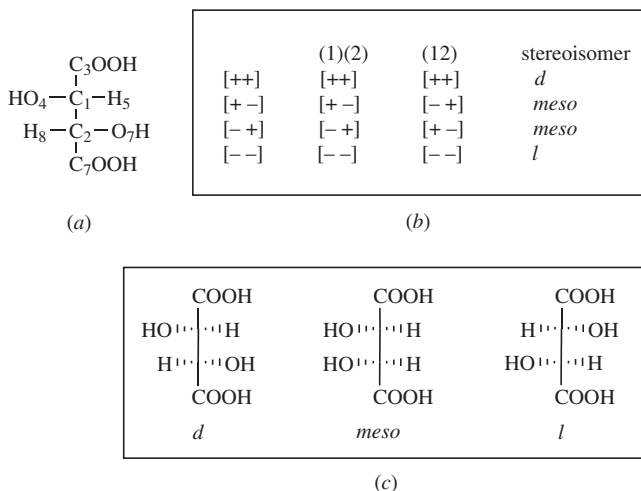
### Stereoisomer Enumeration

Few approaches have been reported to enumerate stereoisomers.<sup>63,64,130–134</sup> We describe here the technique proposed by Nourse.<sup>130</sup> This method has been developed within the CONGEN<sup>131</sup> structure generator, but it is also the method applied in MOLGEN.<sup>108</sup> Nourse's technique computes all stereoisomers of a given structural isomer. Thus, to enumerate the stereoisomers of a given molecular formula, we first generate all structural isomers with the techniques presented earlier. Then, for each structure, we apply Nourse's algorithm. Essentially three steps are involved in this algorithm. (1) All potential stereocenters are determined for the given structural isomer. (2) A permutation group called the configuration group is constructed from the automorphism group of the structure (cf. definition in "From Graph Theory to Chemistry" subsection). (3) The permutations of the configuration group are applied to all possible orientations of the stereocenters, and orientations found identical under the permutations are removed. The number of stereoisomers is the number of remaining orientations.

A stereocenter is defined to be any trivalent or tetravalent atom with at most one hydrogen that is not part of an aromatic system or cumulenes with H<sub>2</sub>-ends, and not triple bonded. A stereocenter has two possible orientations induced by the labels of the neighboring atoms. These labels are simply the atom numbers defined by the generator that produced the structural isomer, and these numbers remain unchanged during stereoisomer enumeration. Because the orientation is defined by an arbitrary labeling, we use the notation +, – instead of the *R,S* nomenclature. However *R,S* notations can be restored in a postprocess.<sup>131</sup> Let  $R1 < R2 < R3 < R4$  be four atom labels attached a given stereocenter; the two possible orientations are as follows:



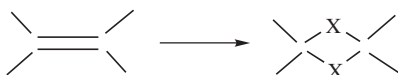
For a structure comprising  $n$  stereocenters, each having two possible orientations,  $2^n$  potential stereoisomers exist. Taking the example of tartaric acid of Figure 17(a), this structure has two stereocenters ( $C_1$  and  $C_2$ ). The potential stereoisomers are  $[++]$ ,  $[+-]$ ,  $[-+]$ , and  $[- -]$ . Some of these stereoisomers are identical because of the symmetry of the structure. The labels of Figure 17(a), show that only two permutations in the automorphism group preserve the structure of tartaric acid:  $(1)(2)(3)(4)(5)(6)(7)(8)$  and  $(12)(36)(47)(58)$ . In this case, the configuration group is simply the set of permutations of the automorphism group restricted to the stereocenters:  $(1)(2)$  and  $(12)$ . To compute the exact number of stereoisomers, we apply the configuration group to all potential stereoisomers and remove all equivalent orientations. The application of the configuration group on the four possible stereoisomers



**Figure 17** The stereoisomers of tartaric acid. (a) Tartaric acid structural isomer with atom labels 1 through 8 (only atoms attached to stereocenters are labeled). (b) Application of the configuration group  $\{(1)(2), (12)\}$  on the four possible stereoisomers. The second and third stereoisomers are identical. (c) The three resulting stereoisomers, a *meso* form and a *dl* pair.

of tartaric acid is given in Figure 17(b). The three resulting stereoisomers are depicted in Figure 17(c).

From the tartaric acid example, it may seem that the configuration group is no different than the automorphism group restricted to the stereocenters. However, more complicated cases exist in which permutations can change the orientations of stereocenters even when the stereocenters are not permuted. As an example, consider the permutation  $(1)(24)(3)$  acting on the labels of 1,2,3,4-tetrachlorocyclobutane. Stereocenter  $C_1$  is attached to  $C_2$ ,  $C_4$ , a chlorine atom, Cl, and a hydrogen atom, H. The permutation  $(1)(24)(3)$  changes this order to  $C_4$ ,  $C_2$ , Cl, and H; consequently, the orientation of  $C_1$  is reversed by  $(1)(24)(3)$ . The same observation can be made for  $C_3$ . To indicate that the orientations of  $C_1$  and  $C_3$  are reversed by the permutation  $(1)(24)(3)$ , Nourse uses the notation  $(1')(24)(3')$ . Application of  $(1')(24)(3')$  on the stereoisomer  $[++++]$  gives the correct configuration  $[-++-]$ , which differs from  $[++++]$ , the configuration given by  $(1)(24)(3)$ . Finally, a stereoisomer induced by double bonds can also be enumerated with Nourse's technique. When double bonds are involved, a special configuration group is computed. This group is the product of the atom automorphism groups and bond automorphism groups. A simpler solution was later suggested by Wieland et al.<sup>135</sup> and consists of converting double bonds into single bonds with fictitious bivalent nodes:



Expanding on Nourse's technique, Wieland<sup>133</sup> proposed an enumeration algorithm of stereoisomers in which the valence of the stereocenters can be larger than four.

To conclude this section on enumeration, it seems that enumerating structural isomers is no longer a technical challenge. The reader not convinced of this conclusion can access the web page of the journal MATCH,<sup>136</sup> enter any molecular formula, and visualize the list of corresponding isomers. The algorithm that produced this list is MOLGEN. Although not every compound family can be counted, as far as isomer enumeration is concerned, up to 50 non-hydrogen atoms, all molecular graphs, can be enumerated according to the authors of MOLGEN. Unfortunately, structural elucidation and molecular design problems do not fit this optimistic picture. The pitfall of isomer enumeration is the number of solutions produced. Of course, the number of solutions can be reduced by adding constraints, but the problem becomes computationally harder and most likely intractable, especially when dealing with overlapping fragments. The usual way to deal with intractable problems in computer science is to apply stochastic techniques in solutions that are only guaranteed up to some probability. The purpose of the next section is to review the stochastic techniques that sample molecular structures for the purpose of structure elucidation and molecular design.

---

### Sampling Structures: What is the Decane Isomer with the Highest boiling point?

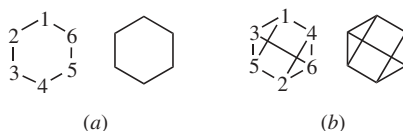
The premise of the sampling approach is the following question: Is it necessary to generate all molecular graphs corresponding to a set of constraints to design compounds having specified activities or properties? As far as structure elucidation is concerned, the question is whether the concept of a unique chemical graph has a physical or chemical significance for complex natural compounds, such as *lignin*, *coal*, *kerogen*, or *humic substances*. As far as sampling is concerned, no method of choice exists like there was for counting molecules (Pólya's theory) and enumerating structures (orderly generation). The reason perhaps is that the field is relatively new. Both in graph theory and computational chemistry, the techniques to sample graphs and chemical graphs were introduced mostly in the last decade. In the subsections that follow, we first summarize what can be learned from graph theory about sampling graphs and then we review their applications in chemistry.

## Sampling Labeled and Unlabeled Graphs

Randomly sampling labeled graphs of  $n$  vertices and  $q$  edges can easily be performed by selecting at random  $q$  pairs of vertices in the set of  $n(n-1)/2$  possible pairs. Such a random selection can be performed with or without replacement depending on whether we wish to create multiple edges.

As we have seen with counting and enumeration, unlabeled graphs are harder to deal with than labeled ones. Nijenhuis and Wilf<sup>137</sup> have shown how to sample unlabeled rooted trees. The approach was extended by Wilf,<sup>138</sup> who gave an algorithm to sample unlabeled unrooted trees. The algorithm is based on a counting series for trees. More complicated is the case of cyclic graphs. Dixon and Wilf<sup>139</sup> were the first to give an algorithm for sampling unlabeled graphs with a specified number  $n$  of vertices. First, a permutation,  $\pi$ , of  $n$  vertices is chosen in the set of all possible permutations, that is, in the symmetric group  $S_n$ . As an example, assume the selected permutation is  $\pi = (135)(246)$  (cf.  $C_3^-$  in Figure 3). Next, a graph is constructed at random from those graphs that are fixed by  $\pi$ , i.e., graphs like benzene that remain unchanged under the action of  $\pi$ . To construct this graph, the permutation  $\pi^*$  acting on the edges is computed from  $\pi$ , where for any edge  $[i,j]$ ,  $\pi^*([i,j]) = [\pi(i), \pi(j)]$ . Using our benzene example, we have  $\pi^* = (12\ 34\ 56)(13\ 35\ 15)(14\ 36\ 25)(16\ 23\ 45)(24\ 46\ 26)$ . Then, for each cycle of  $\pi^*$  independently, we choose with probability  $1/2$  whether *all* or *none* of the edges of the cycle will appear in the graph. Taking our benzene example, we may choose edges in cycles  $(12\ 34\ 56)$  and  $(16\ 23\ 45)$  to be turned on as in Figure 18(a) or edges in cycles  $(13\ 35\ 15)(14\ 36\ 25)(24\ 46\ 26)$  as in Figure 18(b). Both of the resulting graphs are drawn at random from the set of all possible unlabeled graphs of six vertices.

The Dixon and Wilf technique was later expanded by Wormald<sup>140</sup> to sample regular graphs with degrees equal or greater than 3, and by Goldberg and Jerrum<sup>244</sup> to graphs of prescribed degree sequences. The case of degree sequences is of particular interest to chemistry, and in fact, in the article published by Goldberg and Jerrum there is an extension to sample molecules. Their algorithm is a two-step procedure. First, a core structure that does not contain vertices of degree one or two is sampled with a Dixon–Wilf–Wormald's type algorithm. Then, the core is extended adding trees and chains of trees (vertices of degree one or two). Interestingly, a parallel can be drawn between Goldberg



**Figure 18** Two unlabeled graphs drawn at random and unchanged under the permutation  $\pi = (135)(246)$ .

and Jerrum's core structures and the cyclic substructures of CONGEN,<sup>98</sup> or the regular subgraphs of MOLGEN<sup>87</sup> (cf. General Structural Isomer Enumeration subsection). In these approaches, structures are enumerated or sampled by first constructing cyclic subgraphs and then either by connecting these subgraphs or by adding vertices and edges that do not create additional cycles. The main result of Goldberg and Jerrum's article is that molecules can be sampled in polynomial time, which is an interesting result considering that the computational complexity of counting and enumerating molecules are still open questions.

## Sampling Molecules

As with enumeration, we use sampling in structure elucidation and molecular design applications. With both applications in mind, three different techniques have been developed: random sampling, Monte-Carlo sampling, and genetic algorithms.

### *Sampling Molecules at Random*

The first published sampling technique is a generator that constructs linear polymers at random.<sup>141</sup> The random construction is repeated until a polymer is found matching a given set of physical properties. Note that the method is time consuming because the solutions are not refined as the sampling progresses. In the context of drug design, a random sampling technique was proposed<sup>142</sup> to generate random structures by combining fragments. Specifically, fragments are chosen from a database of known drugs with a probability proportional to some statistical weight. Bonding sites are picked randomly for the chosen fragments and for the molecule built so far, and the two are joined together. Fragments are added in such a manner until the total molecular weight exceeds some predefined threshold. The random selection of bonding sites for fusion often produces structures that are chemically unstable or unusual. These structures are eliminated during a selection process based on topological indices and quantitative structure activity relationships. The structures that survive selections are archived in a database of compounds to be considered for synthesis. As in the polymer case, this latter approach appears to be time consuming for molecular (drug) design purposes, because the solutions are not improved as the algorithm progresses. Additionally, these techniques may generate duplicated structures because they essentially sample labeled graphs. In the context of structure elucidation, a random sampling of nonidentical molecular graphs was proposed in 1994.<sup>118</sup> The method is a randomized version of a deterministic structure generation algorithm. Underlying all algorithms enumerating molecules is a construction tree (cf. examples in Figures 12 and 14), and that method selects branches at random instead of exploring all of them. Structures produced by the random selection are different if the branches of the construction tree lead to nonidentical structures,

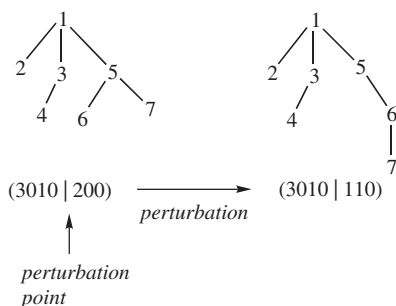
which is the case with the orderly algorithm or the equivalent-classes algorithm the sampling technique was based on. Running the algorithm, we observed that large samples of nonidentical structures could be generated efficiently. Aside from generating nonidentical structures at random, this sampling technique also provides an estimate of the number of solutions. This number is then used to carry out statistical analysis; for instance, mean values and standard deviations of some properties calculated with molecular simulations on the sample can be extrapolated to the entire population of potential structures.

### *Monte Carlo Sampling of Molecules*

Random sampling techniques are appropriate to calculate average properties of compounds matching specific constraints, but they are time consuming when we use them to search for the best compounds matching target properties or experimental data. In such an instance, optimization methods such as Monte-Carlo or genetic algorithms are best suited. Monte-Carlo (MC) and simulated annealing (SA) are simple algorithms that were initially designed to provide efficient simulations of collections of particles in condensed matter physics.<sup>143</sup> In each step of these algorithms, a particle is given a small random displacement, and the resulting change,  $\Delta E$ , in the energy of the system is computed. If  $\Delta E \leq 0$ , the displacement is accepted, and the new configuration is the starting point of the next step. The case  $\Delta E \geq 0$  is treated probabilistically: The probability that the configuration is accepted is  $\exp(-\Delta E/kT)$ , where  $k$  is the Boltzmann constant and  $T$  is the temperature. With MC, the simulations are carried out at equilibrium at a constant temperature  $T$ , whereas with SA, the temperature is decreased according to a predefined cooling program (annealing schedule). Using a cost function in place of the energy and defining configurations by a set of parameters, it is straightforward with this procedure to generate a population of configurations for a given optimization problem. For instance, SA techniques have searched for the global minimum of energy in conformational space.<sup>144</sup> For conformational isomers, the random displacement of the MC/SA algorithm consists of slightly modifying the conformation by either moving atoms or rotating bonds.

In the structural space, any MC/SA random displacement must consist of changing the connectivity between the atoms. A solution to this problem, proposed by Kvasnicka and Pospichal,<sup>104</sup> and illustrated in Figure 19, is to introduce perturbations in bonding patterns starting at a randomly chosen atom. Specifically, assuming an initial structure is constructed, a linear code is computed for this structure, and atoms are ordered according to the code. Examples of suitable codes are the  $n$ -tuple code for acyclic compounds, and the connectivity stack. Next, an atom is chosen at random and the code is randomly modified starting at the chosen atom. Not every perturbation is a valid one; for instance, we need to check that after a perturbation, the valences of the atoms and the total number of bonds are maintained.





**Figure 19** Perturbation of the  $n$ -tuple code of a hydrogen-suppressed  $C_7H_{16}$  isomer. Starting from a randomly selected point, the code is randomly modified and 2-methylhexane is obtained by bond perturbation of 1,2-dimethylpentane.

A disadvantage of this perturbation technique is that it is difficult to control to what extent the structure is changed because bonding pattern changes start at a randomly chosen atom. Ideally, in the spirit of the MC algorithm, we would like to keep the random displacement as small as possible. To this end, another solution to the random displacement came about by observing that connectivity between atoms can be changed by deleting bonds, creating bonds, or modifying bond order.<sup>9</sup> With the convention that a bond is deleted when its order is set to zero and a bond is created when its order is switched from zero to a positive value, all changes of connectivity can be performed by modifying the bond order. Because all structures must have the same total number of bonds, when a bond order is increased, another bond order must be decreased. Hence, changing the connectivity implies the selection of at least two bonds, or four atoms,  $x_1, y_1, x_2$ , and  $y_2$ . Let  $a_{11}, a_{12}, a_{21}$ , and  $a_{22}$  be the order of the bonds  $[x_1, y_1]$ ,  $[x_1, y_2]$ ,  $[x_2, y_1]$  and  $[x_2, y_2]$  in the initial structure, and let  $b_{11}, b_{12}, b_{21}$ , and  $b_{22}$  be the order of the same bonds after the random displacement occurs. The random displacement is performed by a bond order switch. Precisely, a value  $b_{11} \neq a_{11}$  is chosen at random verifying:

$$b_{11} \geq \text{MAX}(0, a_{11} - a_{22}, a_{11} + a_{12} - 3, a_{11} + a_{21} - 3) \quad [33]$$

$$b_{11} \leq \text{MIN}(3, a_{11} + a_{12}, a_{11} + a_{21}, a_{11} - a_{22} + 3) \quad [34]$$

These equations are derived based on the fact that bond orders range between 0 and 3. The orders for all other bonds are computed maintaining the valences of the atoms:

$$b_{12} = a_{11} + a_{12} - b_{11} \quad [35]$$

$$b_{21} = a_{11} + a_{21} - b_{11} \quad [36]$$

$$b_{22} = a_{22} - a_{11} + b_{11} \quad [37]$$

All possible structural isomers of a given molecular formula can be reached with this bond order switch.<sup>145</sup> Also, every structure produced by

the bond order switch is a valid one. Thus, contrary to the bond perturbation technique of Figure 19, no need exists to check for structure consistency. The bond order switch with a SA algorithm was used to search compounds having the maximum and minimum Wiener indices. The correct solutions were found up to 84 carbon atoms. Many *quantitative structure-activity/property relationships* exist between the Wiener index and the boiling point of organic compounds.<sup>40</sup> These relationships may not be linear, but as a general rule, the larger the Wiener number is, the higher is the boiling point. Thus, searching for compounds having the highest boiling point can be achieved by finding molecular graphs having the maximum Wiener index. For dodecane, the maximum Wiener index was found by this algorithm to be  $W = 286$ , and the corresponding structure is the linear dodecane isomer. The bond order switch was also integrated in the SENECA software, and it elucidated structures as large as triterpenes matching experimental 1-D and 2-D NMR spectra.<sup>146</sup>

### *Genetic Algorithms to Sample Molecules*

A genetic algorithm (GA) is a method of producing new examples from combinations of previous individuals, or parents. The algorithm has the same logical structure as inheritance in biological systems. The probability that an individual will be produced and participate as a parent in a succeeding generation must be defined by some standard. For optimization purposes, the suitability of an offspring is usually assessed with a “fitness” function, which is a direct analogy to Darwin’s evolutionary rules of natural selection and survival of the fittest.

The applications of GAs in chemistry have been reviewed in this series of books.<sup>144</sup> We focus here on the use of genetic algorithms to sample and search molecular graphs. The implementation of a GA usually invokes three data processing steps on the genetic code: mutation, crossover (recombination), and selection. The genetic codes suitable for chemical graphs are the ones we have seen with the MC/SA algorithms, the  $n$ -tuple code, and the connectivity stack. Mutations of the genetic code can be performed the same way random displacements are carried out in MC/SA, that is, bond perturbation or bond order switch. Several steps are required to crossover genetic codes. First two parents are selected. Next, a crossover point is chosen at random, the two codes are spliced into two segments, and the corresponding fragments are recombined taking a segment from each code. The crossover operation is illustrated in Figure 20 with the  $n$ -tuple code. As with bond perturbations in MC/SA, no guarantee exists that a crossover operation will maintain the valences of the atoms. Thus, all structures created by crossover must be checked for consistency, which is a disadvantage that GA has versus MC/SA when we use bond order switching. An interesting solution to avoid consistency check during crossover operation appeared in 1999.<sup>147</sup> In this solution, a bond is chosen randomly in each parent. The bond is deleted, and if the parent is not spliced into two pieces, a second bond is then removed in the shortest path linking the two atoms attached to the deleted bond. The process of

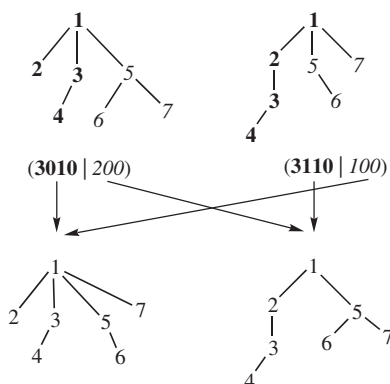


Figure 20 Crossover operation with the  $n$ -tuple code.

bond deletion is repeated until the parent is cut into two disconnected parts. The four resulting pieces (two per parent) are then recombined by saturating the atoms where bonds have been deleted.

The last genetic operation is selection. Elements of the population are selected to form the next generation with a problem-specific fitness function. Taking the simple example of searching for the structure having the highest boiling, the fitness function can be, for instance, the Wiener index of the structure.

The first application of a genetic algorithm to sample molecular structures was in the context of the design of polymers with desired properties.<sup>148</sup> Later, an article appeared to construct combinatorial libraries<sup>149</sup> and targeted libraries<sup>150</sup> for drug design purposes. These applications are limited to linear genetic codes and are thus unable to create structures by recombining parents in a cyclic manner. A general GA algorithm that includes cyclic recombination was implemented for the purpose of structure elucidation of organic compounds from  $^{13}\text{C}$  NMR spectra.<sup>151</sup> In this GA, mutations are performed with bond perturbations as in Figure 19, and crossovers are carried out as in Figure 20. The selection operator is a root-mean-square deviation between the experimental chemical shifts and the predicted chemical shifts obtained with neural network technology. Structures up to 20 heavy atoms have been elucidated with this algorithm.

---

## ENUMERATING MOLECULES: WHAT ARE THE USES?

### Chemical Information

The combination of counting series and enumerating algorithms described previously in this chapter has allowed researchers to generate isomer

lists for not only popular compounds, but for several specific compound classes as well. In this next section of the chapter, we provide a brief review of isomer lists available in the literature as well as tabulate some important and popular lists to provide the reader a quick resource for this information.

Alkanes and alkane-like substances have captured the interest of researchers in isomer enumeration for a long time because of their commercial importance. For example, Henze and Blair published the first isomer enumeration of alkanes in 1931.<sup>15</sup> Here we provide, for reference, tables that list the number of isomers of alkanes, alkenes, alkynes,<sup>152</sup> and stereoalkanes (Table 4), ketones and esters (Table 5), and primary, secondary, and tertiary alcohols (Table 6) up to 25 carbon atoms.

Although the alkane and alkane-like substances are the most important, no series of compounds has received as much interest in generating isomer series as has the polyhexes, with much being done on various classes of benzenoids.<sup>78,153–155</sup> For example, isomer series are available for benzenoids of a variety of classes, including peri-condensed,<sup>156,157</sup> cata-condensed,<sup>35,158</sup> essentially disconnected,<sup>159</sup> helicenes,<sup>158,160</sup> resonant sextets,<sup>161</sup> quinones,<sup>162</sup> coronenes,<sup>163,164</sup> and pyrenes.<sup>165</sup> Recently, with the aid of a new lattice

**Table 4** Isomers of Alkanes, Alkenes, Alkynes, and Stereoalkanes

Carbon Atoms	Alkanes	Alkenes	Alkynes	Stereoalkanes
1	1			1
2	1	1	1	1
3	1	1	1	1
4	2	3	2	2
5	3	5	3	3
6	5	13	7	5
7	9	27	14	11
8	18	66	32	24
9	35	153	72	55
10	75	377	171	136
11	159	914	405	345
12	355	2281	989	900
13	802	5690	2426	2412
14	1858	14397	6045	6553
15	4347	36564	15167	18127
16	10359	93650	38422	50699
17	24894	240916	97925	143255
18	60523	623338	251275	408429
19	148284	1619346	648061	1173770
20	366319	4224993	1679869	3396844
21	910726	11062046	4372872	9892302
22	2278658	29062341	11428365	28972080
23	5731580	76581151	29972078	85289390
24	14490245	202365823	78859809	252260276
25	36797588	536113477	208094977	749329719

**Table 5** Isomers of Ketones and Esters

Carbon Atoms	Ketone Isomers	Ester Isomers
1	1	0
2	1	1
3	2	2
4	3	4
5	7	9
6	14	20
7	32	45
8	72	105
9	171	249
10	405	599
11	989	1463
12	2426	3614
13	6045	9016
14	15167	22695
15	38422	57564
16	97925	146985
17	251275	377555
18	648061	974924
19	1679869	2529308
20	4372872	6589734
21	11428365	17234114
22	29972078	45228343
23	78859809	119069228
24	208094977	314368027
25	550603722	832193902

enumeration algorithm, the number of benzenoid hydrocarbons was calculated up to 35 hexagons.<sup>79</sup> We provide this information in Table 7.

Another class of polyhexes, which are fullerenes, has also been the subject of much interest in the field of isomer generation because of the impact that this class of compounds has made in many areas of science and technology. Although tables exist for the general and isolated-pentagon rule fullerenes,<sup>166,167</sup> we provide in Table 8 an isomer list for these classes up to a large number of vertices, courtesy of the Fullgen code.<sup>94,168</sup>

Although we have provided popular and useful tables in this section for a variety of substances, many other isomer tables exist in the literature. Several of these tables have been generated in an attempt to verify or compare codes that generate isomers. To best aid the reader, references to these listings are provided. Novak gives a small list of some halogen derivatives of a few molecules and ions and their chirality.<sup>169</sup> In a series of several papers, Contreras et al. have isomer tables on dozens of organic compounds, both cyclic and acyclic. Wieland et al. have generated configurational and constitutional isomers for a variety of hydrocarbons up to 10 carbon atoms.<sup>135</sup> Dias provides a constant isomer series for fluorenoid and fluoranthenoid hydrocarbons.<sup>170</sup>

**Table 6** Isomers of Alcohols Series

Carbon Atoms	Primary Alcohols	Secondary Alcohols	Tertiary Alcohols
1	1	0	0
2	1	0	0
3	1	1	0
4	2	1	1
5	4	3	1
6	8	6	3
7	17	15	7
8	39	33	17
9	89	82	40
10	211	194	102
11	507	482	249
12	1238	1188	631
13	3057	2988	1594
14	7639	7528	4074
15	19241	19181	10443
16	48865	49060	26981
17	124906	126369	69923
18	321198	326863	182158
19	830219	849650	476141
20	2156010	2216862	1249237
21	5622109	5806256	3287448
22	14715813	15256265	8677074
23	38649152	40210657	22962118
24	101821927	106273050	60915508
25	269010485	281593237	161962845

**Table 7** Benzenoids Isomers as a Function of the Number of Hexagons

Hexagons	Benzenoids	Hexagons	Benzenoids
1	1	18	8553649747
2	1	19	41892642772
3	3	20	205714411986
4	7	21	1012565172403
5	22	22	4994807695197
6	81	23	24687124900540
7	331	24	122238208783203
8	1435	25	606269126076178
9	6505	26	3011552839015720
10	30086	27	14980723113884739
11	141229	28	74618806326026588
12	669584	29	372132473810066270
13	3198256	30	1857997219686165624
14	15367577	31	9286641168851598974
15	74207910	32	46463218416521777176
16	359863778	33	232686119925419595108
17	1751594643	34	1166321030843201656301
		35	5851000265625801806530

**Table 8** Fullerene Isomers as a Function of the Number of Carbons

Number of Atoms	Numbers of Fullerenes	Number of Atoms	Numbers of Fullerenes
20	1	100	285914
22	0	102	341658
24	1	104	419013
26	1	106	497529
28	2	108	604217
30	3	110	713319
32	6	112	860161
34	6	114	1008444
36	15	116	1207119
38	17	118	1408553
40	40	120	1674171
42	45	122	1942929
44	89	124	2295721
46	116	126	2650866
48	199	128	3114236
50	271	130	3580637
52	437	132	4182071
54	580	134	4787715
56	924	136	5566948
58	1205	138	6344698
60	1812	140	7341204
62	2385	142	8339033
64	3465	144	9604410
66	4478	146	10867629
68	6332	148	12469092
70	8149	150	14059173
72	11190	152	16066024
74	14246	154	18060973
76	19151	156	20558765
78	24109	158	23037593
80	31924	160	26142839
82	39718	162	29202540
84	51592	164	33022572
86	63761	166	36798430
88	81738	168	41478338
90	99918	170	46088148
92	126409	172	51809018
94	153493	174	55817091
96	191839	176	64353257
98	231017		

Luinge generated dioxane isomers as well as isomer lists for a variety of CHO and CHN compounds.<sup>171</sup> CHO, CHON, and CON isomer lists were generated by Molodtsov in 1994.<sup>172</sup> A subclass of indacenoids, namely di-5-catafusenes, were studied by Cyvin et al. and isomer lists were generated.<sup>173</sup> These same authors later provided isomer lists for systems containing pentagons and

heptagons, with both one pentagon (azulenoids)<sup>174</sup> and multiple pentagons.<sup>175</sup> Dolhaine et al.<sup>176</sup> have provided some tables and formulae for the number of isomers of a variety of substituted molecules including benzene, anthracene, and fullerenes. Dolhaine and Honig<sup>177</sup> have also published a large list of inositol oligomers up to the tetramer with estimates of larger isomers for larger oligomers provided. Finally, Davidson<sup>178</sup> provides alkyl frequency distributions in alkane isomers up to 21 carbon atoms.

Before we leave this subsection, we note that a few studies comparing various isomer generation techniques have been published in an attempt to provide both a validation of an algorithm against a test case and a comparison of a variety of methods for consistency and execution time. Such studies may also contain isomer generation lists of the type mentioned in the previous paragraphs. To the reader interested in these tests, we provide references.<sup>60,171,179–181</sup>

## Structure Elucidation

A clear and important application of enumerating molecules falls within a larger framework of structure elucidation. In brief, the ultimate goal of structure elucidation is to take input information and identify *the* compound that is consistent with that information. A more pragmatic goal of this endeavor is to generate *all* candidate molecules consistent with the input information. Although we use information before or after the candidate generation to focus the solution space to a single solution, such an ideal result is not often met, and thus, lists of solutions (perhaps ranked) are produced. The reason for this result has to do with many factors, including the combinatorial explosion of isomers, the quality of the input information, and the efficiencies of the algorithms applying this information. A much more detailed assessment of this situation is provided elsewhere, and the reader is referred to those studies.<sup>110,120,182</sup>

Our goal in this section is to *highlight* some of the popular codes we can use to perform structural elucidation. Where applicable, information required on the types of input is provided. We must note that not all structural elucidation codes are equal. Some codes are considered expert systems containing large databases of initial (stored) information and apply complex algorithms that attempt to reach the ultimate goal. Others are more modest isomer generation codes with some preprocessing or postprocessing to include experimental information to assist with the arrival of a solution (or solution set) for a particular problem. Accordingly, we will describe some isomer generation codes first and finish with the expert systems.

A program to enumerate all possible saturated hydrocarbons was introduced in 1991 by Hendrickson and Parks.<sup>183</sup> This code, called SKEL\_GEN, was tested for structures containing up to 11 carbon atoms, but limited work was extended to larger ringed structures.

In 1992, Contreras introduced Computer-Assisted Molecular Generation and Counting (CAMGEC),<sup>60</sup> which is an exhaustive, selective, and



nonredundant structural generation C code under a Unix platform. The program requires just a molecular formula as input. No means to input other information or for postprocessing exists. Recent improvements to CAM-GEC<sup>61–65</sup> have been presented that improve efficiency and allow for stereoisomer generation.

To aid in the interpretation of infrared spectra, Luinge developed the structure generator Algorithm for the Exhaustive Generation of Irredundant Structures (AEGIS).<sup>171</sup> Although it is reportedly simple for chemists to use and requires only the molecular formula as input, it is written in the PROLOG language that is computationally expensive.

Although not designed to compete with other isomer generation codes, Barone et al.<sup>180</sup> designed an exhaustive method to generate organic isomers from base 2 and base 4 numbers called Generation of Isomers (GI). They have used GI in an attempt to check the consistency between other, much faster, isomer generation codes and have found some discrepancies.

A large, yet exhaustive isomer generation package called ISOGEN<sup>184</sup> produces an irredundant list of structure isomers consistent with a given empirical formula. Revisions to this code with the same name involve modified algorithms that include evolutionary approaches.<sup>185</sup>

Le Bret<sup>181</sup> put forth a novel approach to the structural elucidation problem by using a genetic algorithm to exhaustively generate isomers. The program, called GalvaStructures, has the molecular formula as input and can take various spectral information to aid in the efficiency of solution. Although stochastic, the program seems to be consistent with other generators; yet it is much slower for large problems.

The “grandfather” of knowledge-based structural elucidation codes is the DENDRAL project at Stanford University.<sup>56,57</sup> DENDRAL (DENDRitic ALgorithm) provided a recipe (plan, generation, test) to exhaustively enumerate all isomers given an input set of atoms and spectral information. The generation of structures was performed with CONstrained GENERator (CONGEN) and, ultimately, with a more-advanced structure generator called GENERation with Overlapping Atoms (GENOA).<sup>113</sup> The latter code added some automated features as well as a different way to handle overlapping substructural units.

ASSEMBLE 2.0<sup>186</sup> is a structure generator taking molecular formula and fragments as input. On output, candidates can be ranked based on fragment spectra given on input.

CHEMICS<sup>111</sup> is an automated structure elucidation system for organic compounds that applies 630 fragments in developing structures. Spectroscopic data in the form of IR, <sup>1</sup>H-NMR, and <sup>13</sup>C-NMR as well as bond correlations limit the candidate structures output.

Elucidation by Progressive Intersection of Ordered Substructures (EPIOS)<sup>187</sup> is a code that takes a database of <sup>13</sup>C NMR spectra and generates candidate structures through overlapping fragments.

The structure generator GEN applies up to 30 fragments (obtained from, say, spectral information) as input and can be given various types of constraints

during generation such as molecular formula, molecular weight, and structural considerations. The code is involved in two systems, GENSTR and GENMAS.<sup>179</sup> GENSTR is used when specific fragments can be selected and additional information can be introduced into the generation process. GENMAS is used when only molecular formula is to be input.

GENM a program written in both C and Fortran that generates all non-isomorphic molecular graphs given a set of labeled vertices with a specific valence.<sup>172</sup> Although the code lacks postprocessing, forbidden and required fragments can be input to the program.

MOLGEN, developed from MOLGRAPH,<sup>107</sup> is a structure elucidation code that has made its way into the commercial market and is, perhaps, the most widely known program of its type.<sup>108,188</sup> Upon input of a molecular formula, MOLGEN produces a complete set of redundancy-free isomers. MOLGEN can be used online and provide the user with the number and structure of isomers corresponding to a given molecular formula. An online version can be found in the MATCH journal webpage.<sup>136</sup>

When StrucEluc was first introduced, it only allowed 1-D-NMR data with other information for structural elucidation of molecules containing fewer than 25 atoms. A recent enhancement of StrucEluc now involves a variety of 2-D NMR data as well as data from IR and mass spectra.<sup>189</sup> Such improvements have allowed this code to elucidate product molecules containing more than 60 atoms. This program now forms a suite of programs offered under the name ACD/StructureElucidator.

COCON<sup>190</sup> and a web-based version (WebCocon 4J) search for compounds of known molecular formula compatible with 2-D NMR data as input. This code also can interpret heteronuclear multibond correlations as 2-, 3-, and 4-bond connectivities.

SpecSolv<sup>116</sup> is a structure elucidation system based on <sup>13</sup>C chemical shifts with additional options to apply NMR information of most any kind to aid in candidate refinement. Note that SpecSolv does not require an initial molecular formula as input.<sup>191</sup>

The Expert System for the Elucidation of the Structures of Organic Compounds (ESESOC)<sup>192,193</sup> system is used by chemists for structural elucidation and presents candidate structures consistent with a molecular formula and spectroscopic data. In addition to being a structure generator, the system can extract various information, including IR, NMR and COSY as constraints.

Faulon developed a stochastic structure generator, named SIGNATURE,<sup>118,145</sup> and used it for a variety of natural compound structure elucidation problems. The SENECA structural elucidation system<sup>146</sup> later incorporated algorithms developed by Faulon with the goal of finding the constitution of a molecule given spectroscopic data, most notably, NMR.

COstrained COmbination of Atom-centered fragments (Cocoa)<sup>119</sup> is a structural elucidation method, which rather than combining fragments (based on input information) to make structures, applies the information to remove

fragments. Such an approach makes the best use of all input information. Cocoa was later incorporated into a larger software, SESAMI, that includes a spectrum interpreter.<sup>194</sup>

The HOUDINI program, part of the SESAMI system, contains elements of both structure assembly (ASSEMBLE) and structure reduction (Cocoa).<sup>120</sup>

The Computerized Information System of Organic Chemistry-Structure Elucidation Subsystem (CISOC-SES) is another structural elucidation program to generate candidate structures given NMR information.<sup>195</sup> The algorithms in this system emphasize the use of long-range distance constraints by chemists.

Elyashberg et al.<sup>196</sup> developed a structure elucidation system called X-PERT that applies molar mass, IR, and NMR spectra in combination with a gradual growth of constraints in reaching candidate solutions.

### *Some Structure Elucidation Success Stories*

All codes and systems listed previously have, at some point, been tested to evaluate effectiveness. Once these approaches are reported in the literature and/or presented at a conference, the next step involves adoption of the system for a particular need. However, those who adopt these systems are normally scientists from companies whose work is proprietary. Hence, most of the “real world” successes of structure elucidation are not disseminated. However, cases exist in which difficult structural elucidation problems have been solved and published and we provide a few interesting examples here.

As reported by Munk,<sup>110</sup> the earliest example of a real-world structure elucidation problem with computational techniques was performed in 1967 to determine the structure of actinobolin. Degradation reactions of actinobolin were performed leading to substructures from which an early version of ASSEMBLE generated six viable candidates. Subsequent spectroscopic studies were performed on these six, and the correct structure was isolated. This procedure, in total, reportedly took several man-years to complete. By way of comparison, this problem was revisited recently with the SESAMI system with 1-D and 2-D NMR data derived from actinobolin acetate. SESAMI, in conjunction with some simple experiments and other data, resolved the correct structure, with the entire process numbering in days as opposed to years.

Lignin is an important polymer found in the cell wall of plants and plays a key role in a variety of industries including pulp and paper as well as fuel and wood science. Although many studies had been performed on lignin, a clear and compelling picture of the structure of this polymer corroborating existing experimental information had yet to be determined. Using the SIGNATURE program in conjunction with molecular simulation as well as NMR data and known fragments for lignin monomers, Faulon and Hatcher<sup>197</sup> concluded that a structure with a helical template for lignin was preferred over random structures. Such a conclusion was consistent with Raman spectroscopic information. More recent uses of SIGNATURE by chemists include the design of sample structures for humic substances<sup>198</sup> and asphaltenes.<sup>199</sup>

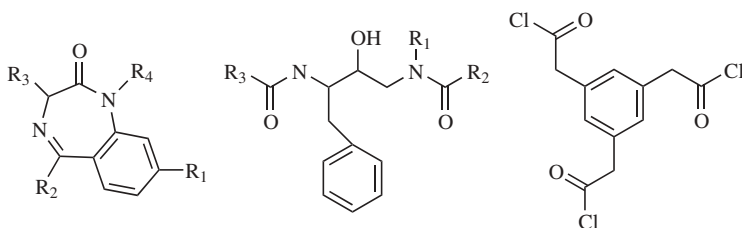
In another, recent example, the ACD/Structure Elucidator resolved 2-D NMR data on a C31 alkaloid that had several ambiguities in connectivity associated with spectral overlap. Twelve candidates were revealed, of which eleven were ultimately ruled out for violating a variety of constraints. The new compound, named quindolinocrypto-tackieine, was thus solved.<sup>200</sup>

## Combinatorial Library Design

Methods for synthesizing large combinatorial libraries of organic compounds emerged in the mid-1990s.<sup>201–203</sup> This introduction revolutionized drug discovery, as millions of candidate compounds could then be synthesized in parallel and evaluated with high throughput screening techniques.<sup>204,205</sup> However, given that the number of compounds that could be synthesized exceeds  $10^{12}$  for even a simple combinatorial library scheme based only on commercially available reagents,<sup>206</sup> and estimated to be over  $10^{30}$  in the whole accessible chemical space,<sup>207</sup> the effectiveness of these early “brute force” experimental approaches were necessarily limited. Therefore, along with the development of combinatorial chemistry came a growth in virtual chemistry and software tools to sift *in silico* through large numbers of potential compounds in combinatorial libraries to select the most promising subset for synthesis and experimental testing. The ability to enumerate molecules is a crucial step in many virtual chemistry algorithms for designing combinatorial libraries, which in turn can help discover lead pharmaceutical compounds.

Many different approaches have been taken to designing these libraries. Diversity approaches<sup>208,209</sup> design general exploratory libraries that maximize chemical diversity for initial drug discovery. Biased approaches design focused libraries when *a priori* knowledge exists of either the structure of the target (structure-based approaches<sup>210–215</sup>) or a small lead compound (similarity approaches<sup>216</sup>). Informative design,<sup>217</sup> a relatively new approach, designs a library that will provide the maximum amount of information from each experimental cycle of synthesis and testing. In addition to examining the potential drug-binding properties of the library, most library design efforts try to simultaneously maximize the ADMET (adsorption, distribution, metabolism, excretion, and toxicity) properties of the library members, with heuristics of “drug-likeness”<sup>218–221</sup> as well as other functions such as cost. The cost functions may be implemented as simple postprocessing filters or as objective functions to maximize. Although library design methods can deal in the chemical space of the reactants, product-based design has been shown to be superior, albeit more expensive computationally.<sup>222,223</sup> For a complete review of computational techniques applied to combinatorial libraries, see Lewis et al.<sup>209</sup> Most of the product-based approaches involve either full or partial enumeration of the products of the combinatorial library.

The basic chemistry for combinatorial synthesis usually involves a core group or scaffold, to which a set of reagents or R groups is systematically



**Figure 21** Benzodiazepine scaffold (left) and a statine-base peptomimetic (middle) are typical asymmetric scaffolds. Benzene triacylchloride scaffold (right) is a typical symmetric scaffold.

reacted at each substitution site (see Figure 21 for some typical scaffolds). As all combinations of reagents are synthesized, the total size of the combinatorial library can be estimated with Pólya's theorem as given by Eq. [7] in the "Counting Structures" section. For most combinatorial libraries, the scaffold is asymmetric and the size of the library is simply the product of the number of possible reagents at each substitution site. For example, if a scaffold has three variable positions R1, R2, and R3, each with 1000 possible reagents that can react at that site,  $1000^3$  or  $10^9$  possible compounds could be synthesized. Because often even more than 1000 reagents are commercially available per reactant site, the numbers are often even greater.

The first step in constructing a virtual combinatorial library for a given scaffold is to identify the pool of reagents available for each substitution site on the scaffold, which usually involves searching substructures in a database of commercially available and in-house compounds for those containing the appropriate reactivity for the synthesis protocol at each substitution site. Filters eliminate reagents with inappropriate chemistries such as functional groups predicted to cause side reactions, or those that may interfere with or cause false-positive tests in the biological assays, or those with insufficient "drug-like" properties.

Once the reagents are selected, the next step is to enumerate product compounds of the library. Most enumeration programs take into account only the simple case of an asymmetric scaffold; however, MOLGEN-COMB<sup>224</sup> is specifically designed to handle symmetry. For the asymmetric scaffold case, two basic approaches exist. The first approach is referred to as "fragment marking". Here, the reactant pools are treated with a preprocessing step where they are "marked" by removing the reacting functional group in each reagent and replacing it with a free valance. Enumeration consists of systematically placing all clipped reagents onto the scaffold<sup>225,226</sup> with an algorithm similar to Scheme I given in the "Enumerating Structures" section. A simple version of this approach has been applied in the structure-based programs CombiDOCK<sup>212</sup> and CombiBUILD.<sup>211</sup> One problem with this approach is that it cannot handle all synthetic reaction types, such as the

Diels–Alder reaction, or systems with no clear core scaffold such as oligomeric libraries of variable length.<sup>227</sup> The second approach is for chemists to use a “reaction transform” to perform the same chemical transformations *in silico* that are being performed chemically. Advantages of this approach are that it can be applied on all chemistries, does not involve any preprocessing of the reagents, and the transforms can be reused by chemists. It has the disadvantage, however, of being computationally more demanding and thus slower to perform. Many programs apply the SMARTS molecular query notation from the Daylight toolkit<sup>228</sup> to design reaction transformation tools,<sup>229</sup> an approach that has been incorporated into the ADEPT program.<sup>227</sup> Commercially available programs that perform computational enumeration include CombiLibMaker in Sybyl<sup>122</sup>, Analog Builder in Cerius<sup>2, 123</sup> PRO\_SELECT,<sup>215</sup> and the QuaSAR-CombiGen module in MOE (Montreal, Canada).

As full enumeration of large combinatorial libraries is impractical, several methods have been developed to avoid explicit enumeration of all library members. Many diversity- and similarity-based design strategies apply sampling approaches such as genetic algorithms,<sup>149,230,231</sup> simulated annealing,<sup>232,233</sup> and stochastic sampling<sup>234</sup> to optimize libraries (see Gillet et al.<sup>235</sup> for a full review). Some approaches involve descriptors that can estimate the properties of library members without explicit enumeration of the full library, either by descriptors that can be calculated roughly from the sum of the reactants<sup>206,225,226</sup> or by a neural net to estimate properties from a small sampling of enumerated products.<sup>229</sup> Combining multiple approaches can also reduce the problem to a computationally tractable number of possible solutions. For example, a diversity search can be performed to select a smaller library that can then be explicitly enumerated as a starting point for a structure-based library design of a focused library.<sup>236</sup> Enumeration can also be reduced in structure-based programs, which start with the three-dimensional structure of the target, by taking a “divide-and-conquer” strategy.<sup>211,212,215</sup> In a divide-and-conquer scheme, the scaffold is first docked to the binding site. The reagents are then evaluated *individually* at each substitution site for predicted binding, thus turning the problem from  $n^r$  enumerations to  $r \times n$ , where  $n$  is number of reagents and  $r$  is the number of substitution sites. Only top-scoring reagents are saved for full compound enumeration and evaluation. Virtual libraries designed in this manner have rapidly led to potent lead compounds.<sup>210,237</sup>

## Molecular Design with Inverse-QSAR

The forward QSAR procedure defines an equation or a set of equations that relates a variable of interest (dependent variable) with independent variables. The dependent variable is normally an activity/property of interest (binding affinity, normal boiling point, IC<sub>50</sub>, etc.), whereas the independent variables are related to the structure of the substance. Developing a QSAR

for a particular activity/property involves training the parameters of the model against a well-defined set of data (training set), with a small portion of the data held back for validation of the model (test set). Once the QSAR is effectively trained and validated, we can use this model to predict the activity/property value of a given compound by determining the values of its independent variables in a straightforward manner. On the other hand, rather than determining an activity/property value for a particular compound from the QSAR, what if we want to determine a compound from the QSAR given a *particular activity/property value*? This question is known as the inverse-QSAR problem (I-QSAR) and is the subject of this section.

Anyone reading this chapter has, undoubtedly, solved an inverse-type problem in one form or another. The key to efficient solution lies in the restriction of the solution space. If constraints are composed such that the solution space is limited, a brute-force technique (try all candidates) can guarantee a solution. In the field of molecular design, however, the solution space comes from all compounds that can be reasonably made from the various atoms in the Periodic Table. Hence, we need a way to limit this solution space to arrive at candidate solutions efficiently. We will describe some of these techniques.

As mentioned earlier, Kier and Hall published a series of papers in the early 1990s that described the inverse QSAR methodology with chi indices. The QSARs they developed had a maximum of four descriptors. Example applications included the inverse design of alkanes from molar volume<sup>122</sup> and the identification of isonarcotic agents.<sup>238</sup>

Simultaneous with the work from Kier and Hall, Zefirov et al.<sup>127</sup> developed a similar technique with the count of paths. The QSARs they used were given for three Kappa-shape descriptors and they considered three functional groups, namely, alkanes, alcohols, and small oxygen-containing compounds.

In 2001, Bruggemann et al.<sup>239</sup> demonstrated Hasse diagrams combined with a similarity measure in the generation of solutions to the inverse problem involving toxicity of algae. Their method is based on partial ordered sets and does not assume a particular model for the QSAR.

Garg and Achenie also demonstrated a reasonable approach to the solution of the I-QSAR problem in 2001.<sup>240</sup> Taking a target scaffold of an antifolate molecule for dihydrofolate reductase inhibition, these authors generated a QSAR for both activity and selectivity. They solved the I-QSAR problem to maximize selectivity through changing substituents on the scaffold, subject to a constraint of a threshold activity. Finally, a work by Skvortsova, et al.<sup>241</sup> from 2003 demonstrated that the I-QSAR problem could be solved for the Hosoya index plus constraints on the number of carbon atoms for a system of 78 hydrocarbons.

All previous methods have limitations. As has been demonstrated, we can limit the problem size by working with a QSAR derived with only a few descriptors. Hence, many solutions can be found associated with the given problem. Additionally, we can limit the solution space to contain, say, only

hydrocarbons or alcohols. A third issue on the approaches described concerns the degeneracy of the solution. It is not uncommon for a particular value of a topological index to correspond to a large number of possible compounds. A novel I-QSAR methodology has been developed recently that addresses these issues and will be described briefly.

The *signature* molecular descriptor, previously mentioned in this chapter for structural elucidation (cf. Figure 16), has found utility in the solution of the inverse-QSAR problem. The reason for this utility is that *signature* can produce meaningful QSARs<sup>129,242</sup> and is the least degenerate of dozens of topological indices tested.<sup>129</sup> Additionally, *signature* lends itself to the inversion process.<sup>128</sup> An algorithm that will enumerate and sample chemical structures corresponding to the numerical solutions from the I-QSAR problem has already been developed and tested for a variety of compounds, including alkanes, fullerenes, and HIV-1 protease inhibitors.<sup>128</sup>

The inverse-QSAR problem with *signature* has also been applied to a small set of LFA-1/ICAM-1 peptide inhibitors to assist in the search and design of more potent inhibitory compounds. After developing a QSAR, the inverse-QSAR technique with *signature* generated many novel inhibitors. Two more potent inhibitors were synthesized and tested in vivo, confirming them to be the strongest inhibiting peptides to date.<sup>243</sup>

---

## CONCLUSION AND FUTURE DIRECTIONS

We have seen in this chapter that counting, enumerating, and sampling of molecular graphs from a molecular formula are not the technical challenges they once were. Counting formulas exist for a large variety of chemical compounds, isomer generators can enumerate without construction, or count, up to  $10^{30}$  molecular graphs,<sup>87</sup> and sampling molecules can theoretically be performed efficiently.<sup>244</sup> Nonetheless, the computational complexity of counting and enumerating molecular graphs remains an unsolved problem. It is thus expected that research collaboration among mathematics, computer science, and computational chemistry will continue to devise better techniques to count and enumerate molecules.

As far as structure elucidation and molecular design are concerned, enumerating molecules from a molecular formula is only part of the problem. Indeed, as we have argued in this chapter, enumerating structures with constraints, such as including the presence or absence of overlapping fragments, is most probably an intractable problem. Alternative stochastic sampling approaches have been devised recently to overcome the difficulties of enumerating molecules with constraints. Only a few stochastic techniques have so far been published, and it is likely that the sampling approach will continue to be developed and used by chemists in the near future for practical purposes such as elucidation from NMR spectra.



Even if we knew how to efficiently enumerate or sample molecular graphs under constraints, our job would not be completed. Molecules are 3-D objects, and ultimately, structure generators should produce 3-D structures. Enumerating stereoisomers alone is not sufficient as we also need to generate the structural conformations corresponding to the problem constraints. The natural solution that comes to mind is to first enumerate all molecular graphs matching the constraints and then to explore the conformational space of each graph. Although codes exist for constructing 3-D representations of molecular graphs,<sup>245–248</sup> exploring the conformational space of each molecular graph is a cumbersome task that must be added to the already costly endeavor of structure enumeration. Such a strategy does not appear to be computationally feasible. One alternative to avoid that computational bottleneck may be to use the geometrical enumeration we have seen for benzenoids.<sup>76,79</sup> Recall that this approach consists of enumerating self-avoiding polygons on lattices. The enumeration is performed directly on a 2-D lattice space (benzenoids are planar) ignoring the underlying molecular graphs. The advantage of the geometrical approach is that energetically unfavorable structures are never constructed. Such is obviously not the case when enumerating dimensionless molecular graphs. Considering that geometrical enumeration is currently the most powerful technique to enumerate benzenoids, such a promising approach may be further explored, perhaps, for structure elucidation and molecular design purposes.

---

## ACKNOWLEDGMENTS

We thank the editors of this volume for inviting us to contribute this chapter. We also gratefully acknowledge the support of the Mathematics Information and Computer Science Program of the U.S. Department of Energy, and the Computer Science Research Foundation Program, a U.S. Department of Energy and Sandia National Laboratories program, under Grant DE-AC04-76DP00789.

---

## REFERENCES

1. A. Cayley, *Rep. Brit. Ass. Adv. Sci.*, **14**, 257 (1875). On the Analytical Forms Called Trees with Applications to the Theory of Chemical Compounds.
2. G. Pólya, *C. R. Acad. Sci. Paris*, **201**, 1167 (1935). Un Probleme Combinatoire General Sur Les Groupes Des Permutations Et Le Calcul Du Nombre Des Composes Organiques.
3. D. H. Rouvray, *J. Mol. Struct. (THEOCHEM)*, **54**, 1 (1989). The Pioneering Contributions of Cayley and Sylvester to the Mathematical Description of Chemical Structure.
4. S. J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, Prentice Hall, Upper Saddle River, New Jersey, 2002.
5. P. W. Fowler and P. Hansen, in *DIMAC Workshop Report*, Rutgers University Press, New Brunswick, New Jersey, 2001. The Working Group on Computer-Generated Conjectures from Graph Theoretic and Chemical Databases I.

6. A. R. Leach, in *Reviews in Computational Chemistry*, Vol. 2, K. B. Lipkowitz and D. B. Boyd, Eds., VCH Publishers, New York, 1991, pp. 1–55. A Survey of Methods for Searching the Conformational Space of Small and Medium-Sized Molecules.
7. H. A. Sheraga, in *Reviews in Computational Chemistry*, Vol. 3, K. B. Lipkowitz and D. B. Boyd, Eds., VCH Publishers, New York, 1992, pp. 73–142. Predicting Three-Dimensional Structures of Oligopeptides.
8. Y. Kudo and S.-I. Sasaki, *J. Chem. Doc.*, **14**, 200 (1974). The Connectivity Stack, a New Format for Representation of Organic Chemical Structures.
9. J.-L. Faulon, *J. Chem. Inf. Comput. Sci.*, **38**, 432 (1998). Isomorphism, Automorphism Partitioning, and Canonical Labeling Can Be Solved in Polynomial-Time for Molecular Graphs.
10. B. D. McKay, Nauty User's Guide, Version 2.2. 2004 Available: <http://cs.anu.edu.au/people/bdm/nauty/>.
11. E. M. Luck, *J. Comput. Sys. Sci.*, **25**, 42 (1982). Isomorphism of Graphs of Bounded Valence Can Be Tested in Polynomial Time.
12. M. R. Garey and D. S. Johnson, *Computers and Intractability. A Guide to the Theory of Np-Completeness*, W. H. Freeman & Company, New York, 1979.
13. A. Cayley, *Philos. Mag.*, **13**, 172 (1857). On the Analytical Forms Called Trees.
14. F. Hermann, *Ber. Dtsch. Chem. Ges.*, **13**, 792 (1880). On the Problem of Evaluating the Number of Isomeric Paraffins of the Formula  $C_nH_{2n+2}$ .
15. H. R. Henze and C. Blair, *J. Am. Chem. Soc.*, **53**, 3077 (1931). The Number of Isomeric Hydrocarbons of the Methane Series.
16. G. Pólya, *Acta Math.*, **68**, 145 (1937). Kombinatorische Anzahlbestimmungen Fur Gruppen, Graphen Und Chemische Verbindungen.
17. F. Harary and Z. Norman, *Proc. Am. Math. Soc.*, **11**, 134 (1960). Dissimilarity Characteristic Theorems for Graphs.
18. F. Zhang, R. Li, and G. Lin, *J. Mol. Struct. (THEOCHEM)*, **453**, 1 (1998). The Enumeration of Heterofullerenes.
19. H. Friepertinger, *MATCH Commun. Math. Comput. Chem.*, **33**, 121 (1996). The Cycle Index of the Symmetry Group of the Fullerene C<sub>60</sub>.
20. P. W. Fowler, D. B. Redmond, and J. P. B. Sandall, *J. Chem. Soc. Faraday Trans.*, **19**, 2883 (1998). Enumeration of Fullerene Derivatives C<sub>70</sub>xm of Given Symmetries.
21. R. M. Nembra and A. T. Balaban, *J. Chem. Inf. Comput. Sci.*, **38**, 1145 (1998). Algorithm for the Direct Enumeration of Chiral and Achiral Skeleton of a Homosubstituted Derivative of a Monocyclic Cycloalkane with a Large and Factorizable Ring Size N.
22. I. Baraldi and D. Vanossi, *J. Chem. Inf. Comput. Sci.*, **40**, 386 (2000). Regarding Enumeration of Molecular Isomers.
23. R. C. Read, *J. London Math. Soc.*, **35**, 344 (1960). The Enumeration of Locally Restricted Graphs II.
24. R. Otter, *Annals Math.*, **49**, 583 (1948). The Number of Trees.
25. J. Wang, R. Li, and S. Wang, *J. Math. Chem.*, **33**, 171 (2003). Enumeration of Isomers of Acyclic Saturated Hydroxyl Ethers.
26. S. Fujita, *Symmetry and Combinatorial Enumeration in Chemistry*, Springer-Verlag, Berlin, 1992.
27. S. J. Cyvin, B. N. Cyvin, J. Brunvoll, and J. Wang, *J. Mol. Struct. (THEOCHEM)*, **445**, 127 (1998). Enumeration of Staggered Conformers of Alkanes and Monocyclic Cycloalkanes.
28. C. Y. Yeh, *J. Chem. Inf. Comput. Sci.*, **35**, 912 (1995). Isomer Enumeration of Alkanes, Labeled Alkanes, and Monosubstituted Alkanes.
29. C. Y. Yeh, *J. Phys. Chem.*, **100**, 15800 (1996). Theory of Acyclic Chemical Networks and Enumeration of Polyenoids Via Two-Dimensional Chirality.

30. C. Y. Yeh, *J. Chem. Inf. Comput. Sci.*, **36**, 854 (1996). Isomer Enumeration of Alkenes, and Aliphatic Cyclopropane Derivatives.
31. C. Y. Yeh, *J. Chem. Phys.*, **105**, 9706 (1996). Counting Linear Polyenes by Excluding Structures with Steric Strain.
32. C. Y. Yeh, *J. Mol. Struct. (THEOCHEM)*, **432**, 153 (1996). Isomerism of Asymmetric Dendrimers and Stereoisomerism of Alkanes.
33. L. Bytautas and D. J. Klein, *J. Chem. Inf. Comput. Sci.*, **38**, 1063 (1998). Chemical Combinatorics for Alkane-Isomer Enumeration and More.
34. S. J. Cyvin and I. Gutman, *Kekule Structures in Benzenoid Hydrocarbons*, Springer-Verlag, Berlin, 1988.
35. I. Gutman and S. J. Cyvin, *Introduction to the Theory of Benzenoid Hydrocarbons*, Springer-Verlag, Berlin, 1989.
36. I. Gutman and S. J. Cyvin, *Advances in the Theory of Benzenoid Hydrocarbons*, Springer-Verlag, Berlin, 1990.
37. I. Gutman, S. J. Cyvin, and J. Brunvoll, *Advances in the Theory of Benzenoid Hydrocarbons II*, Springer-Verlag, Berlin, 1992.
38. J. R. Dias, *Handbook of Polycyclic Hydrocarbons: Part A, Benzenoid Hydrocarbons*, Elsevier, Amsterdam, 1987.
39. J. R. Dias, *Handbook of Polycyclic Hydrocarbons: Part B: Polycyclic Isomers and Heteroatom Analogs of Benzenoid Hydrocarbons*, Elsevier, Amsterdam, 1988.
40. N. Trinajstić, *Chemical Graph Theory*, CRC Press, Boca Raton, Florida, 1992.
41. A. T. Balaban and F. Harary, *Tetrahedron*, **24**, 2505 (1968). Enumeration and Proposed Nomenclature of Benzenoid Cata-Condensed Polycyclic Aromatic Hydrocarbons.
42. F. Harary and R. C. Read, *Proc. Edinburgh Math. Soc., Ser. II*, **17**, 1 (1970). Enumeration of Tree-Like Polyhexes.
43. S. J. Cyvin and J. Brunvoll, *J. Math. Chem.*, **9**, 33 (1992). Generating Functions for the Harary-Read Numbers Classified According to Symmetry.
44. S. J. Cyvin, J. Brunvoll, and B. N. Cyvin, *J. Math. Chem.*, **9**, 19 (1992). Harary-Read Numbers for Catafusenes: Complete Classification According to Symmetry.
45. J. Brunvoll, S. J. Cyvin, and B. N. Cyvin, *J. Math. Chem.*, **21**, 193 (1997). Enumeration of Tree-Like Octagonal Systems.
46. S. J. Cyvin, F. Zhang, and J. Brunvoll, *J. Math. Chem.*, **3**, 103 (1992). Enumeration of Perifusenes with One Internal Vertex - A Complete Mathematical Solution.
47. S. J. Cyvin, F. Zhang, B. N. Cyvin, G. Xiaofeng, and J. Brunvoll, *J. Chem. Inf. Comput. Sci.*, **32**, 532 (1992). Enumeration and Classification of Benzenoid Systems. 32. Normal Perifusenes with Two Internal Vertices.
48. T. Akutsu, S. Miyano, and S. Kuhara, *Bioinformatics*, **16**, 727 (2000). Inferring Qualitative Relations in Genetic Networks and Metabolic Pathways.
49. F. Harary and E. M. Palmer, *Graphical Enumeration*, Academic Press, New York, 1973.
50. R. C. Read, *Annals of Discrete Math.*, **2**, 107 (1978). Every One a Winner, or How to Avoid Isomorphism Search When Cataloguing Combinatorial Configurations.
51. I. A. Faradzev, in *Problemes Combinatoires et Theorie des Graphes*, University of Paris, Orsay, France, 1978, pp. 131–135. Constructive Enumeration of Combinatorial Objects.
52. B. D. McKay, *J. Algorithms*, **26**, 306 (1998). Isomorph-Free Exhaustive Generation.
53. L. A. Goldberg, *J. Algorithms*, **13**, 128 (1992). Efficient Algorithms for Listing Unlabeled Graphs.
54. J. Lederberg, in *The Mathematical Science*, R. Cosrims, Ed., MIT Press, Cambridge, Massachusetts, 1969, pp. 37–51. Topology of Molecules.
55. J. Lederberg, G. L. Sutherland, B. G. Buchanan, E. A. Feigenbaum, A. V. Robertson, A. M. Duffield, and C. Djerassi, *J. Am. Chem. Soc.*, **91**, 2973 (1969). Applications of Artificial

- Intelligence for Chemical Inference. I. The Number of Possible Organic Compounds. Acyclic Structures Containing C, H, O, and N.
56. R. K. Lindsay, B. G. Buchanan, E. A. Feigenbaum, and J. Lederberg, *Applications of Artificial Intelligence for Organic Chemistry: The Dendral Project*, McGraw-Hill, New York, 1980.
  57. N. A. B. Gray, *Computer-Assisted Structure Elucidation*, John Wiley & Sons, New York, 1986.
  58. J. V. Knop, W. R. Müller, Z. Jericevi, and N. Trinajstić, *J. Chem. Inf. Comput. Sci.*, **21**, 91 (1981). Computer Enumeration and Generation of Trees and Rooted Trees.
  59. J. E. Hopcroft and R. E. Tarjan, in *Complexity of Computer Computation*, R. Miller and E. Thatcher, Eds., Plenum Press, New York, 1972, pp. 131–152. Isomorphism of Planar Graphs.
  60. M. L. Contreras, R. Valdivia, and R. Rozas, *J. Chem. Inf. Comput. Sci.*, **32**, 323 (1992). Exhaustive Generation of Organic Isomers. 1. Acyclic Structures.
  61. M. L. Contreras, R. Valdivia, and R. Rozas, *J. Chem. Inf. Comput. Sci.*, **32**, 483 (1992). Exhaustive Generation of Organic Isomers. 2. Cyclic Structures.
  62. M. L. Contreras, R. Rozas, and R. Valdivia, *J. Chem. Inf. Comput. Sci.*, **34**, 610 (1994). Exhaustive Generation of Organic Isomers. 3. Acyclic, Cyclic and Mixed Compounds.
  63. M. L. Contreras, R. Rozas, R. Valdivia, and R. Aguero, *J. Chem. Inf. Comput. Sci.*, **35**, 752 (1994). Exhaustive Generation of Organic Isomers. 4. Acyclic Stereoisomers with One or More Chiral Carbon Atoms.
  64. M. L. Contreras, G. M. Trevisiol, J. Alvarez, G. Arias, and R. Rozas, *J. Chem. Inf. Comput. Sci.*, **35**, 475 (1999). Exhaustive Generation of Organic Isomers. 5. Unsaturated Optical and Geometrical Stereoisomers and a New CIP Subrule.
  65. M. L. Contreras, J. Alvarez, M. Riveros, G. Arias, and R. Rozas, *J. Chem. Inf. Comput. Sci.*, **41**, 964 (2001). Exhaustive Generation of Organic Isomers. 6. Stereoisomers Having Isolated and Spiro Cycles and New Extended N-Tuples.
  66. I. Lukovits, *J. Chem. Inf. Comput. Sci.*, **39**, 563 (1999). Isomer Generation: Syntactic Rules for Detection of Isomorphism.
  67. I. Lukovits, *J. Chem. Inf. Comput. Sci.*, **40**, 361 (2000). Isomer Generation: Semantic Rules for Detection of Isomorphism.
  68. K. Balasubramanian, J. J. Kaufman, W. S. Koski, and A. T. Balaban, *J. Comput. Chem.*, **1**, 149 (1980). Graph Theoretical Characterization Computer Generation of Certain Carcinogenic Benzenoid Hydrocarbons and Identification of Bay Regions.
  69. J. V. Knop, K. Szymanski, Z. Jericevi, and N. Trinajstić, *J. Comput. Chem.*, **4**, 23 (1983). Computer Enumeration and Generation of Benzenoid Hydrocarbons and Identification of Bay Regions.
  70. I. Stojmenović, R. Toi, and R. Doroslovacki, *Proceedings of the Sixth Yugoslav Seminar on Graph Theory*, Novi Sad, Dubrovnik (1985). Generating and Counting Hexagonal Systems In Graph Theory.
  71. W. J. He, W. C. He, Q. X. Wang, J. Brunvoll, and S. J. Cyvin, *Naturforsch.*, **43a**, 693 (1988). Supplement to Enumeration of Benzenoid and Coronoid Hydrocarbons.
  72. S. Nikolic, N. Trinajstić, J. V. Knop, W. R. Müller, and K. Szymanski, *J. Math. Chem.*, **4**, 357 (1990). On the Concept of the Weighted Spanning Tree of Dualist.
  73. W. R. Müller, K. Szymanski, and J. V. Knop, *Croat. Chem. Acta*, **62**, 481 (1989). On Counting Polyhex Hydrocarbons.
  74. W. R. Müller, K. Szymanski, J. V. Knop, S. Nikoli, and N. Trinajstić, *J. Comput. Chem.*, **11**, 223 (1990). On the Enumeration and Generation of Polyhex Hydrocarbons.
  75. J. V. Knop, W. R. Müller, K. Szymanski, and N. Trinajstić, *J. Chem. Inf. Comput. Sci.*, **30**, 159 (1990). Use of Small Computers for Large Computations: Enumeration of Polyhex Hydrocarbons.

- 
76. R. Tosic, D. Masulovic, I. Stojmenovi, J. Brunvoll, S. J. Cyvin, and B. N. Cyvin, *J. Chem. Inf. Comput. Sci.*, **35**, 181 (1995). Enumeration of Polyhex Hydrocarbons to  $H = 17$ .
77. G. Caporossi and P. Hansen, *J. Chem. Inf. Comput. Sci.*, **38**, 610 (1998). Enumeration of Polyhex Hydrocarbons to  $H = 21$ .
78. G. Brinkmann, G. Caporossi, and P. Hansen, *Commun. Math. Chem. (MATCH)*, **43**, 133 (2001). Numbers of Benzenoids and Fusenes.
79. M. Voge, A. J. Guttmann, and I. Jensen, *J. Chem. Inf. Comput. Sci.*, **42**, 456 (2002). On the Number of Benzenoid Hydrocarbons.
80. I. G. Enting and A. J. Guttmann, *J. Phys. A*, **22**, 1371 (1989). Polygons on the Honeycomb Lattice.
81. J. de Vries, *Rendiconti Circolo Mat. Palermo*, **5**, 221 (1891). Sur Les Configurations Planes Dont Chaque Point Supporte Des Droites.
82. A. T. Balaban, *Revue Roumaine de Chimie*, **12**, 103 (1967). Valence-Isomerism of Cyclopolynes (Erratum).
83. F. C. Bussemaker, S. Cobeljic, D. M. Cvetkovic, and J. J. Seidel, *J. Combin. Theory Ser. B.*, **23**, 234 (1977). Cubic Graphs on 14 Vertices.
84. B. D. McKay and G. F. Royle, *Ars Combinatoria*, **21a**, 129 (1986). Constructing the Cubic Graphs on up to 20 Vertices.
85. G. Brinkmann, *J. Graph Theory*, **23**, 139 (1996). Fast Generation of Cubic Graphs.
86. M. Meringer, *J. Graph Theory*, **30**, 137 (1999). Fast Generation of Regular Graphs and Construction of Cages.
87. T. Grüner, R. Laue, and M. Meringer, in *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, Rutgers University Press, New Brunswick, New Jersey, 1997, pp. 113–122. Algorithms for Group Actions: Homomorphism Principle and Orderly Generation Applied to Graphs.
88. X. Liu and D. J. Klein, *J. Comput. Chem.*, **12**, 1265 (1991). Sixty-Atom Carbon Cages.
89. C.-H. Sah, *Croatica Chemica Acta*, **66**, 105 (1993). Combinatorial Construction of Fullerene Structures.
90. D. J. Klein and X. Liu, *Int. J. Quantum Chem. Quantum Chem. Symp.*, **28**, 501 (1994). Elemental Carbon Isomerism.
91. D. E. Manolopoulos and P. W. Fowler, *Chem. Phys. Lett.*, **204**, 1 (1993). A Fullerene without a Spiral.
92. P. W. Fowler, T. Pisanski, A. Graovac, and J. Zerovnik, in *Discrete Mathematical Chemistry*, Vol. 51, P. Hansen, P. W. Fowler, and M. Zheng, Eds., American Mathematical Society, Providence, Rhode Island 2000, pp. 175–188. A Generalized Ring Spiral Algorithm for Coding Fullerenes and Other Cubic Polyhedra.
93. P. W. Fowler and D. E. Manolopoulos, *An Atlas of Fullerenes*, Oxford University Press, Oxford, United Kingdom, 1995.
94. G. Brinkmann and A. W. Dress, *J. Algorithms*, **23**, 345 (1997). A Constructive Enumeration of Fullerenes.
95. G. Brinkmann, A. W. Dress, S. W. Perrey, and J. Stove, *Math. Prog.*, **79**, 71 (1997). Two Applications of the Divide & Conquer Principle in the Molecular Sciences.
96. G. Brinkmann and A. W. Dress, *Advances Applied Math.*, **21**, 473 (1998). Penthex Puzzles. A Reliable and Efficient Top-Down Approach to Fullerene-Structure Enumeration.
97. E. C. Kirby and P. Pollack, *J. Chem. Inf. Comput. Sci.*, **38**, 66 (1998). How to Enumerate the Connectional Isomers of a Toroidal Polyhex Fullerene.
98. L. M. Masinter, N. S. Sridharan, J. Lederberg, and D. H. Smith, *J. Am. Chem. Soc.*, **96**, 7702 (1974). Applications of Artificial Intelligence for Chemical Inference. XII. Exhaustive Generation of Cyclic and Acyclic Isomers.
99. Y. Kudo and S.-I. Sasaki, *J. Chem. Inf. Comput. Sci.*, **16**, 43 (1976). Principle for Exhaustive Enumeration of Unique Structures Consistent with Structural Information.

100. C. A. Shelley, T. R. Hays, M. E. Munk, and R. V. Roman, *Analytica Chimica Acta*, **103**, 121 (1978). An Approach to Automated Partial Structure Expansion.
101. I. P. Bangov, *J. Chem. Inf. Comput. Sci.*, **30**, 277 (1990). Computer-Assisted Structure Generation For a Gross Formula. 3. Alleviation of the Combinatorial Problem.
102. J.-L. Faulon, *J. Chem. Inf. Comput. Sci.*, **32**, 338 (1992). On Using Graph-Equivalent Classes for the Structure Elucidation of Large Molecules.
103. V. Kvasnicka and J. Pospichal, *J. Chem. Inf. Comput. Sci.*, **30**, 99 (1990). Canonical Indexing and Constructive Enumeration of Molecular Graphs.
104. V. Kvasnicka and J. Pospichal, *J. Chem. Inf. Comput. Sci.*, **36**, 516 (1996). Simulated Annealing Construction of Molecular Graphs with Required Properties.
105. M. S. Molchanova and N. S. Zefirov, *J. Chem. Inf. Comput. Sci.*, **38**, 8 (1998). Irredundant Generation of Isomeric Molecular Structures with Some Known Fragments.
106. M. S. Molchanova, V. V. Shcherbukhin, and N. S. Zefirov, *J. Chem. Inf. Comput. Sci.*, **36**, 888 (1996). Computer Generation of Molecular Structures by the Smog Program.
107. A. Kerber, R. Laue, and D. A. Moser, *Analytica Chimica Acta*, **235**, 2973 (1990). Structure Generator for Molecular Graphs.
108. C. Benecke, R. Grund, R. Hohberger, R. Laue, A. Kerber, and T. Wieland, *Analytica Chimica Acta*, **141** (1995). Molgen+, a Generator of Connectivity Isomers and Stereoisomers for Molecular Structure Elucidation.
109. S. Bohanec and J. Zupan, *J. Chem. Inf. Comput. Sci.*, **31**, 531 (1991). Structure Generation of Constitutional Isomers from Structural Fragments.
110. M. E. Munk, *J. Chem. Inf. Comput. Sci.*, **38**, 997 (1998). Computer-Based Structure Determination: Then and Now.
111. K. Funatsu, N. Miyabayashi, and S.-I. Sasaki, *J. Chem. Inf. Comput. Sci.*, **28**, 18 (1988). Further Developments of Structure Generation in the Automated Structure Elucidation System Chemics.
112. S. G. Molodtsov, *MATCH Commun. Math. Comput. Chem.*, **30**, 203 (1994). Generation of Molecular Graphs with a Given Set of Nonoverlapping Fragments.
113. R. E. Carhart, D. H. Smith, N. A. B. Gray, J. G. Nourse, and C. Djerassi, *J. Org. Chem.*, **46**, 1708 (1981). Genoa: A Computer Program for Structure Elucidation Utilizing Overlapping and Alternative Substructures.
114. J. E. Dubois, M. Carabedian, and R. Ancian, *C. R. Acad. Sci. (Paris)*, **290**, 369 (1980). Automatic Structural Elucidation by C-13 NMR - DARC-EPIOS Method - Search for a Discriminant Chemical Structure Displacement Relationship.
115. J. E. Dubois, M. Carabedian, and R. Ancian, *C. R. Acad. Sci. (Paris)*, **290**, 383 (1980). Automatic Structural Elucidation by C-13 NMR - Darc-Epios Method - Description of Progressive Elucidation by Ordered Intersection of Substructures.
116. M. Will, W. Fachinger, and J. R. Richert, *J. Chem. Inf. Comput. Sci.*, **36**, 221 (1996). Fully Automated Structure Elucidation - A Spectroscopist's Dream Comes True.
117. A. Schrijver, *Theory of Linear and Integer Programming*, Wiley, New York, 1986.
118. J.-L. Faulon, *J. Chem. Inf. Comput. Sci.*, **34**, 1204 (1994). Stochastic Generator of Chemical Structure. 1. Application to the Structure Elucidation of Large Molecules.
119. B. D. Christie and M. E. Munk, *J. Chem. Inf. Comput. Sci.*, **28**, 87 (1988). Structure Generation by Reduction: A New Strategy for Computer-Assisted Structure Elucidation.
120. A. Korytko, K. P. Schulz, M. Madison, and M. E. Munk, *J. Chem. Inf. Comput. Sci.*, **43**, 1434 (2003). Houdini: A New Approach to Computer-Based Structure Generation.
121. K. P. Schulz, A. Korytko, and M. E. Munk, *J. Chem. Inf. Comput. Sci.*, **42**, 1447 (2003). Applications of a Houdini-Based Structure Elucidation System.
122. L. B. Kier, L. H. Hall, and J. W. Frazer, *J. Chem. Inf. Comput. Sci.*, **33**, 143 (1993). Design of Molecules from Quantitative Structure-Activity Relationship Models. 1. Information Transfer between Path and Vertex Degree Counts.

123. L. H. Hall, L. B. Kier, and J. W. Frazer, *J. Chem. Inf. Comput. Sci.*, **33**, 148 (1993). Design of Molecules from Quantitative Structure-Activity Relationship Models. 2. Derivation and Proof of Information Transfert Relating Equations.
124. L. H. Hall, R. S. Dailey, and L. B. Kier, *J. Chem. Inf. Comput. Sci.*, **33**, 598 (1993). Design of Molecules from Quantitative Structure-Activity Relationship Models. 3. Role of Higher Order Path Counts: Path 3.
125. L. B. Kier and L. H. Hall, *Quant. Struct.-Act. Relat.*, **12**, 383 (1994). The Generation of Molecular Structures from a Graph-Based QSAR Equation.
126. L. H. Hall and J. B. Fisk, *J. Chem. Inf. Comput. Sci.*, **34**, 1184 (1994). Degree Set Generation for Chemical Graphs.
127. M. I. Skvortsova, I. I. Baskin, O. L. Slovokhotova, V. A. Palyulin, and N. S. Zefirov, *J. Chem. Inf. Comput. Sci.*, **33**, 630 (1993). Inverse Problem in QSAR/QSPR Studies for the Case of Topological Indices Characterizing Molecular Shape (Kier Indices).
128. J.-L. Faulon, C. J. Churchwell, and D. P. Visco, Jr., *J. Chem. Inf. Comput. Sci.*, **43**, 721 (2003). The Signature Molecular Descriptor. 2. Enumerating Molecules from Their Extended Valence Sequences.
129. J.-L. Faulon, D. P. Visco, Jr., and R. S. Pophale, *J. Chem. Inf. Comput. Sci.*, **43**, 707 (2003). The Signature Molecular Descriptor. 1. Extended Valence Sequences and Topological Indices.
130. J. G. Nourse, *J. Am. Chem. Soc.*, **101**, 1210 (1979). The Configuration Symmetry Group and Its Application to Stereoisomer Generation, Specification, and Enumeration.
131. J. G. Nourse, R. E. Carhart, D. H. Smith, and C. Djerassi, *J. Am. Chem. Soc.*, **101**, 1216 (1979). Exhaustive Generation of Stereoisomers for Structure Elucidation.
132. H. Abe, H. Hayasaka, Y. Miyashita, and S.-I. Sasaki, *J. Chem. Inf. Comput. Sci.*, **24**, 216 (1984). Generation of Stereoisomeric Structures Using Topological Information Alone.
133. T. Wieland, *J. Chem. Inf. Comput. Sci.*, **35**, 220 (1995). Enumeration, Generation, and Construction of Stereoisomers of High-Valence Stereocenters.
134. L. A. Zaltina and M. E. Elyashberg, *MATCH Commun. Math. Comput. Chem.*, **27**, 191 (1992). Generation of Stereoisomers and Their Spatial Models Corresponding to the Given Molecular Structure.
135. T. Wieland, A. Kerber, and R. Laue, *J. Chem. Inf. Comput. Sci.*, **36**, 413 (1996). Principles of the Generation of Constitutional and Configurational Isomers.
136. Match-Online, 2005 Available: <http://www.mathe2.uni-bayreuth.de/match/online/links.html>.
137. A. Nijenhuis and H. S. Wilf, *Combinatorial Algorithms*, Academic Press, New York, 1978.
138. H. S. Wilf, *J. Algorithms*, **5**, 247 (1984). The Uniform Selection of Free Trees.
139. J. D. Dixon and H. S. Wilf, *J. Algorithms*, **4**, 205 (1983). The Random Selection of Unlabeled Graphs.
140. N. C. Wormald, *SIAM J. Comput.*, **16**, 717 (1987). Generating Random Unlabeled Graphs.
141. G. C. Derringer and R. L. Markham, *J. Appl. Polymer Sci.*, **30**, 4609 (1985). A Computer-Based Methodology for Matching Polymer Structure with Required Properties.
142. R. Nilakantan, N. Bauman, and R. Venkataraghavan, *J. Chem. Inf. Comput. Sci.*, **31**, 527 (1991). A Method for Automatic Generation of Novel Chemical Structures and Its Potential Applications to Drug Discovery.
143. N. Metropolis and A. W. Rosenbluth, *J. Chem. Phys.*, **21**, 1087 (1953). Equation of State Calculation by Fast Computing Machines.
144. R. Judson, in *Reviews in Computational Chemistry*, Vol. 10, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 1997, pp. 1–73. Genetic Algorithms and Their Use in Chemistry.
145. J.-L. Faulon, *J. Chem. Inf. Comput. Sci.*, **34**, 731 (1996). Stochastic Generator of Chemical Structure. 2. Using Simulated Annealing to Search the Space of Constitutional Isomers.

146. C. Steinbeck, *J. Chem. Inf. Comput. Sci.*, **41**, 1500 (2001). Seneca: A Platform-Independent, Distributed, and Parallel System for Computer-Assisted Structure Elucidation in Organic Chemistry.
147. A. Globus, J. Lawton, and T. Wipke, *Nanotechnology*, **10**, 290 (1999). Automatic Molecular Design Using Evolutionary Techniques.
148. V. Venkatasubramanian, K. Chan, and J. M. Caruthers, *J. Chem. Inf. Comput. Sci.*, **35**, 188 (1995). Evolutionary Design of Molecules with Desired Properties Using the Genetic Algorithm.
149. R. P. Sheridan and S. K. Kearsley, *J. Chem. Info. Comput. Sci.*, **35**, 310 (1995). Using a Genetic Algorithm to Suggest Combinatorial Libraries.
150. R. P. Sheridan, S. G. SanFeliciano, and S. K. Kearsley, *J. Molec. Graphics and Modelling*, **18**, 320 (2000). Designing Targeted Libraries with Genetic Algorithms.
151. J. Meiler and M. Will, *J. Chem. Inf. Comput. Sci.*, **41**, 1535 (2001). Automated Structure Elucidation of Organic Molecules from <sup>13</sup>C NMR Spectra Using Genetic Algorithms and Neural Networks.
152. C-W. Lam, *J. Math. Chem.*, **23**, 421 (1998). A Mathematical Relationship between the Number of Isomers of Alkenes and Alkynes: A Result Established from the Enumeration of Isomers of Alkenes from Alky Biradicals.
153. J. R. Dias, *J. Chem. Inf. Comput. Sci.*, **30**, 61 (1990). Benzenoid Series Having a Constant Number of Isomers.
154. S. J. Cyvin, *Chem. Phys. Lett.*, **181**, 431 (1991). Note on the Series of Fully Benzenoid Hydrocarbons with a Constant Number of Isomers.
155. S. J. Cyvin and J. Brunvoll, *Chem. Phys. Lett.*, **176**, 413 (1991). Series of Benzenoid Hydrocarbons with a Constant Number of Isomers.
156. J. R. Dias, *Chem. Phys. Lett.*, **176**, 559 (1991). Enumeration of Benzenoid Series Having a Constant Number of Isomers.
157. J. R. Dias, *MATCH Commun. Math. Comput. Chem.*, **26**, 87 (1991). Strictly Pericondensed Benzenoid Isomers.
158. S. J. Cyvin, B. N. Cyvin, and J. Brunvoll, *MATCH Commun. Math. Comput. Chem.*, **26**, 63 (1991). Isomer Enumeration of Catafusenes,  $C_{4n+2}H_{2n+4}$  Benzenoid and Helicenic Hydrocarbons.
159. J. Brunvoll, S. J. Cyvin, B. N. Cyvin, and I. Gutman, *MATCH Commun. Math. Comput. Chem.*, **24**, 51 (1989). Essentially Disconnected Benzenoids: Distribution of K, the Number of Kekule Structures, in Benzenoid Hydrocarbons. VIII.
160. B. N. Cyvin, Z. Fujii, G. Xiaofeng, J. Brunvoll, and S. J. Cyvin, *MATCH Commun. Math. Comput. Chem.*, **29**, 143 (1993). On the Total Number of Polyhexes with Ten Hexagons.
161. J. R. Dias, *J. Chem. Inf. Comput. Sci.*, **31**, 89 (1991). Benzenoid Series Having a Constant Number of Isomers. 3. Total Resonant Sextet Benzenoids and Their Topological Characteristics.
162. J. R. Dias, *J. Chem. Inf. Comput. Sci.*, **30**, 53 (1990). Isomer Enumeration and Topological Characteristics of Benzenoid Quinones.
163. B. N. Cyvin, J. Brunvoll, C. Rongsi, and S. J. Cyvin, *MATCH Commun. Math. Comput. Chem.*, **29**, 131 (1993). Coronenic Coronoids: A Course in Chemical Enumeration.
164. D. J. Klein, T. P. Zivkovic, and A. T. Balaban, *MATCH Commun. Math. Comput. Chem.*, **29**, 107 (1993). The Fractal Family of Coro-[N]-Enes.
165. S. J. Cyvin, B. N. Cyvin, and J. Brunvoll, *MATCH Commun. Math. Comput. Chem.*, **30**, 73 (1994). The Number of Pyrene Isomers Is Still Unknown.
166. P. W. Fowler, P. Hansen, and D. Stevanovic, *MATCH Commun. Math. Comput. Chem.*, **48**, 37 (2003). A Note on the Smallest Eigenvalue of Fullerenes.
167. M. Yoshida and E. Ōsawa, *Bull. Chem. Soc. Jpn.*, **68**, 2073 (1995). Formalized Drawing of Fullerene Nets. 1. Algorithm and Exhaustive Generation of Isomeric Structures.



168. Plantri, 2001 Available: <http://cs.anu.edu.au/~bdm/plantri/fullgen-guide.txt>.
169. I. Novak, *J. Chem. Educ.*, **73**, 120 (1996). Chemical Enumeration with Mathematica.
170. J. R. Dias, *Chem. Phys. Lett.*, **185**, 10 (1991). Series of Fluorenoïd/Fluoranthenoïd Hydrocarbons Having a Constant Number of Isomers.
171. H. J. Luinge, *MATCH Commun. Math. Comput. Chem.*, **27**, 175 (1992). AEGIS, a Structure Generation Program in Prolog.
172. S. G. Molodtsov, *MATCH Commun. Math. Comput. Chem.*, **30**, 213 (1994). Computer-Aided Generation of Molecular Graphs.
173. B. N. Cyvin, J. Brunvoll, and S. J. Cyvin, *MATCH Commun. Math. Comput. Chem.*, **33**, 35 (1996). Di-5-Catafusenes, a Subclass of Indacenoids.
174. J. Brunvoll, S. J. Cyvin, and B. N. Cyvin, *MATCH Commun. Math. Comput. Chem.*, **34**, 91 (1996). Azulenoïds.
175. B. N. Cyvin, J. Brunvoll, and S. J. Cyvin, *MATCH Commun. Math. Comput. Chem.*, **34**, 109 (1996). Isomer Enumeration of Unbranched Catacondensed Polygonal Systems with Pentagons and Heptagons.
176. H. Dolhaine, H. Honig, and M. van Almsick, *MATCH Commun. Math. Comput. Chem.*, **39**, 21 (1999). Sample Applications of an Algorithm for the Calculation of Isomers with More Than One Type of Achiral Substituent.
177. H. Dolhaine and H. Honig, *MATCH Commun. Math. Comput. Chem.*, **46**, 91 (2002). Full Isomer-Tables of Inositol-Oligomers up to Tetramers.
178. S. Davidson, *J. Chem. Inf. Comput. Sci.*, **42**, 147 (2002). Fast Generation of an Alkane-Series Dictionary Ordered by Side-Chain Complexity.
179. S. Bohanec and J. Zupan, *MATCH Commun. Math. Comput. Chem.*, **27**, 49 (1992). Structure Generator GEN.
180. R. Barone, F. Barberis, and M. Chanon, *MATCH Commun. Math. Comput. Chem.*, **32**, 19 (1995). Exhaustive Generation of Organic Isomers from Base 2 and Base 4 Numbers.
181. C. Le Bret, *MATCH Commun. Math. Comput. Chem.*, **41**, 79 (2000). Exhaustive Isomer Generation Using the Genetic Algorithm.
182. M. E. Elyashberg, *Russ. Chem. Rev.*, **68**, 525 (1999). Expert Systems for Structure Elucidation of Organic Molecules by Spectral Methods.
183. J. B. Hendrickson and C. A. Parks, *J. Chem. Inf. Comput. Sci.*, **31**, 101 (1991). Generation and Enumeration of Carbon Skeletons.
184. S. Y. Zhu and J. P. Zhang, *J. Chem. Inf. Comput. Sci.*, **22**, 34 (1982). Exhaustive Generation of Structural Isomers for a Given Empirical Formula - A New Algorithm.
185. X. Shao, C. Wen-sheng, and M. Zhang, *Jisuanji Yu Yingyong Huaxue*, **15**, 169 (1998). Generation of Isomers of Organic Molecules Using Genetic Algorithms.
186. M. Badertscher, A. Korytko, K. P. Schulz, M. Madison, M. E. Munk, P. Portmann, M. Junghans, P. Fontana, and E. Pretsch, *Chemometrics and Intelligent Laboratory Systems*, **51**, 73 (2002). Assemblé 2.0: A Structure Generator.
187. M. Carabedian, L. Dagane, and J. E. Dubois, *Anal. Chem.*, **60**, 2186 (1988). Elucidation by Progressive Intersection of Ordered Substructures from Carbon-13 Nuclear Magnetic Resonance.
188. R. Grund, A. Kerber, and R. Laue, *MATCH Commun. Math. Comput. Chem.*, **27**, 87 (1992). MOLGEN, Ein Computeralgebra System Fur Die Konstruktion Molekularer Graphen.
189. M. E. Elyashberg, K. A. Blinov, A. J. Williams, E. R. Martirosian, and S. G. Molodtsov, *J. Nat. Prod.*, **65**, 693 (2002). Application of a New Expert System for the Structure Elucidation of Natural Products from Their 1D and 2D NMR Data.
190. T. Lindel, J. Junker, and M. Kock, *J. Molec. Mod.*, **3**, 364 (1997). Cocon: From NMR Correlation Data to Molecular Constitutions.
191. M. Will and J. Richert, *J. Chem. Inf. Comput. Sci.*, **37**, 403 (1997). Specsolv - An Innovation at Work.

192. C. Hu and L. Xu, *Fenxi Huaxue*, **20**, 643 (1992). Computer Automatic Structure Elucidation Expert System, Esesoc.
193. J. Hao, L. Xu, and C. Hu, *Science in China, Series B: Chemistry*, **43**, 503 (2000). Expert System for Elucidation of Structures of Organic Compounds (Esesoc) - Algorithm on Stereoisomer Generation.
194. B. D. Christie and M. E. Munk, *J. Am. Chem. Soc.*, **113**, 3750 (1991). The Role of Two-Dimensional Nuclear Magnetic Resonance Spectroscopy in Computer-Enhanced Structure Elucidation.
195. C. Peng, S. Yuan, C. Zheng, Y. Hui, H. Wu, and K. Ma, *J. Chem. Inf. Comput. Sci.*, **34**, 814 (1994). Application of Expert System Cisoc-Ses to the Structure Elucidation of Complex Natural Products.
196. M. E. Elyashberg, E. R. Martirosian, Y. Z. Karasev, H. Thiele, and H. Somberg, *Analytica Chimica Acta*, **337**, 265 (1997). X-Pert: A User-Friendly Expert System for Molecular Structure Elucidation by Spectral Methods.
197. J.-L. Faulon and P. G. Hatcher, *Energy and Fuels*, **8**, 402 (1994). Is There Any Order in the Structure of Lignin?
198. M. S. Diallo, A. Simpson, P. Gassman, J.-L. Faulon, J. J. H. Johnson, W. A. Goddard, III and P. G. Hatcher, *Environ. Sci. & Technol.*, **37**, 1783 (2003). 3-D Structural Modeling of Humic Acids through Experimental Characterization, Computer Assisted Structure Elucidation and Atomistic Simulations. 1. Chelsea Soil Humic Acid.
199. M. S. Diallo, A. Strachan, J.-L. Faulon, and W. A. Goddard, III *Petroleum Science and Technology*, **22**, 877 (2004). Properties of Petroleum Geomacromolecules through Computer Assisted Structure Elucidation and Atomistic Simulations. 1. Bulk Arabian Light Asphaltenes.
200. A. Williams, G. Martin, K. A. Blinov, and M. E. Elyashberg, in *44th Annual Meeting of the American Society of Pharmacognosy*, Chapel Hill, North Carolina, 2003. All Good Things to Those Who Wait: Solving a Structure Computationally after 10 Years of Human Effort.
201. L. A. Thompson and J. A. Ellman, *Chem. Rev.*, **96**, 555 (1996). Synthesis and Applications of Small Molecule Libraries.
202. E. M. Gordon, M. A. Gallop, and D. V. Patel, *Acc. Chem. Res.*, **29**, 144 (1996). Strategy and Tactics in Combinatorial Organic Synthesis. Applications to Drug Discovery.
203. F. Balkenhohl, C. v. d. Bussche-Hunnefeld, A. Lansky, and C. Zechel, *Angew Chem. Int. Ed. Engl.*, **35**, 2289 (1996). Combinatorial Synthesis of Small Organic Molecules.
204. G. S. Sittampalam, S. D. Kahl, and W. P. Janzen, *Curr. Opin. Chem. Biol.*, **1**, 384 (1997). High-Throughput Screening: Advances in Assay Technologies.
205. K. R. Oldenburg, *Ann. Rep. Med. Chem.*, **33**, 301 (1998). Current and Future Trends in High Throughput Screening for Drug Discovery.
206. R. Cramer III, D. Patterson, R. Clark, F. Soltanshahi, and M. Lawless, *J. Chem. Inf. Comput. Sci.*, **38**, 1010 (1998). Virtual Compound Libraries: A New Approach to Decision Making in Molecular Discovery Research.
207. Y. C. Martin, *Perspect. Drug Disc. Des.*, **7**, 159 (1997). Challenges and Prospects For Computational Aids to Molecular Diversity.
208. D. C. Spellmeyer, J. M. Blaney, and E. M. Martin, in *Practical Application of Computer-Aided Drug Design*, P. S. Charifson, Ed., Dekker, New York, 1997, pp. 165–194. Computational Approaches to Chemical Libraries.
209. R. A. Lewis, S. D. Pickett, and D. E. Clark, in *Reviews in Computational Chemistry*, Vol. 16, K. B. Lipkowitz and D. B. Boyd, Eds., VCH Publishers, New York, 2000, pp. 1–51. Computer-Aided Molecular Diversity Analysis and Combinatorial Library Design.
210. E. Kick, D. C. Roe, A. Skillman, G. Liu, T. Ewing, Y. Sun, I. Kuntz, and J. Ellman, *Chemistry & Biology*, **4**, 297 (1997). Structure-Based Design and Combinatorial Chemistry Yield Low Nanomolar Inhibitors of Cathepsin D.

211. D. C. Roe, *Application and Development of Tools for Structure-Based Drug Design*, Ph.D. thesis, University of California, San Francisco, San Francisco, California, 1995.
212. Y. Sun, T. J. A. Ewing, A. G. Skillman, and I. D. Kuntz, *J. Comput.-Aided Mol. Design*, **12**, 597 (1998). Combidock: Structure-Based Combinatorial Docking and Library Design.
213. M. Rarey and T. Lengauer, *Perspect. Drug Disc. Des.*, **20**, 63 (2000). A Recursive Algorithm for Efficient Combinatorial Library Docking.
214. H. Bohm, D. Banner, and L. Weber, *J. Comput.-Aided Mol. Design*, **13**, 51 (1999). Combinatorial Docking and Combinatorial Chemistry: Design of Potent Non-Peptide Thrombin Inhibitors.
215. C. Murray, D. Clark, T. Auton, M. Firth, J. Li, R. Sykes, B. Waszkowycz, D. Westhead, and S. Young, *J. Comput.-Aided Mol. Design*, **11**, 193 (1997). Pro-Select: Combining Structure-Based Drug Design and Combinatorial Chemistry for Rapid Lead Discovery. 1. Technology.
216. P. Willett, J. Barnard, and G. Downs, *J. Chem. Inf. Comput. Sci.*, **38**, 983 (1998). Chemical Similarity Searching.
217. S. Teig, *J. Bio. Scr.*, **3**, 85 (1998). Informative Libraries Are More Useful Than Diverse Ones.
218. C. Lipinski, F. Lombardo, B. Dominy, and P. Feeney, *Advanced Drug Delivery Reviews*, **23**, 3 (1997). Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings.
219. J. Wang and K. Ramnarayan, *J. Combinatorial Chemistry*, **1**, 524 (1999). Toward Designing Drug-Like Libraries: A Novel Computational Approach for Prediction of Drug Feasibility of Compounds.
220. J. Sadowski and H. Kubinyi, *J. Med. Chem.*, **41**, 3325 (1998). A Scoring Scheme for Discriminating between Drugs and Nondrugs.
221. Ajay, W. Walters, and M. Murcko, *J. Med. Chem.*, **41**, 3314 (1998). Can We Learn to Distinguish between "Drug-Like" and "Nondrug-Like" Molecules?
222. V. Gillet, P. Willett, and J. Bradshaw, *J. Chem. Inf. Comput. Sci.*, **37**, 731 (1997). The Effectiveness Of Reactant Pools For Generating Structurally-Diverse Combinatorial Libraries.
223. E. Jamois, M. Hassan, and M. Waldman, *J. Chem. Inf. Comput. Sci.*, **40**, 63 (2000). Evaluation of Reagent-Based and Product-Based Strategies In the Design of Combinatorial Library Subsets.
224. R. Gugisch, A. Kerber, R. Laue, M. Meringer, and J. Weidinger, *MATCH Commun. Math. Chem.*, **41**, 189 (2000). Molgen-Comb, a Software Package for Combinatorial Chemistry.
225. G. Downs and J. Barnard, *J. Chem. Inf. Comput. Sci.*, **37**, 59 (1997). Techniques for Generating Descriptive Fingerprints in Combinatorial Libraries.
226. B. Leland, B. Christie, J. Nourse, D. Grier, R. Carhart, T. Maffett, S. Welford, and D. Smith, *J. Chem. Inf. Comput. Sci.*, **37**, 62 (1997). Managing the Combinatorial Explosion.
227. A. Leach, J. Bradshaw, D. Green, and M. Hann, *J. Chem. Inf. Comput. Sci.*, **39**, 1161 (1999). Implementation of a System for Reagent Selection and Library Enumeration, Profiling, and Design.
228. C. A. James, *Daylight Theory Manual*, Daylight, Chemical Information Systems Inc., Mission Viejo, California, 2004.
229. V. Lobanov and D. Agrafiotis, *Combinatorial Chemistry & High Throughput Screening*, **5**, 167 (2002). Scalable Methods for the Construction and Analysis of Virtual Combinatorial Libraries.
230. R. Brown and Y. Martin, *J. Med. Chem.*, **40**, 2304 (1997). Designing Combinatorial Library Mixtures Using a Genetic Algorithm.
231. V. Gillet, P. Willett, P. Fleming, and D. Green, *J. Molec. Graphics and Modelling*, **20**, 491 (2002). Designing Focused Libraries Using Moselect.

232. A. Good and R. A. Lewis, *J. Med. Chem.*, **40**, 3926 (1997). New Methodology For Profiling Combinatorial Libraries and Screening Sets: Cleaning Up The Design Process With HAR-PICK.
233. M. Hassan, J. Bielawski, J. Hempel, and M. Waldman, *Molecular Diversity*, **2**, 64 (1996). Optimization and Visualization of Molecular Diversity of Combinatorial Libraries.
234. D. Agrafiotis, *J. Chem. Inf. Comput. Sci.*, **37**, 841 (1997). Stochastic Algorithms for Maximizing Molecular Diversity.
235. V. Gillet, *J. Comput.-Aided Mol. Design*, **16**, 371 (2002). Reactant- and Product-Based Approaches to the Design of Combinatorial Libraries.
236. M. P. Beavers and X. Chen, *J. Molec. Graphics and Modelling*, **20**, 463 (2002). Structure-Based Combinatorial Library Design: Methodologies And Applications.
237. T. Haque, A. Skillman, C. Lee, H. Habashita, I. Gluzman, T. Ewing, D. Goldberg, I. Kuntz, and J. Ellman, *J. Med. Chem.*, **42**, 1428 (1999). Potent, Low-Molecular-Weight Non-Peptide Inhibitors of Malarial Aspartyl Protease Plasmeprin II.
238. L. B. Kier and L. H. Hall, *Quant. Struct.-Act. Relat.*, **12**, 383 (1993). The Generation of Molecular Structure for a Graph-Based Equation.
239. R. Bruggemann, S. Pudenz, L. Carlsen, P. B. Sorensen, M. Thomsen, and R. K. Mishra, *SAR and QSAR in Envir. Res.*, **11**, 473 (2001). The Use of Hasse Diagrams as a Potential Approach for Inverse QSAR.
240. S. Garg and L. E. K. Achenie, *Biotechnol. Prog.*, **17**, 412 (2001). Mathematical Programming Assisted Drug Design for Nonclassical Antifolates.
241. M. I. Skvortsova, K. S. Fedyaev, V. A. Palyulin, and N. S. Zefirov, *Internet Electron. J. Mol. Des.*, **2**, 70 (2003). Molecular Design of Chemical Compounds with Prescribed Properties from QSAR Models Containing the Hosoya Index.
242. D. P. Visco, Jr., R. S. Pophale, M. D. Rintoul, and J.-L. Faulon, *J. Molecular Graphics and Modelling*, **20**, 429 (2002). Developing a Methodology for an Inverse Quantitative Structure Activity Relationship Using the Signature Molecular Descriptor.
243. C. J. Churchwell, M. D. Rintoul, S. Martin, D. P. Visco, Jr., A. Kotu, R. S. Larson, L. O. Sillerud, D. C. Brown, and J.-L. Faulon, *J. Molecular Graphics and Modelling*, **22**, 263 (2004). The Signature Molecular Descriptor. 3. Inverse Quantitative Structure-Activity Relationship of ICAM-1 Inhibitory Peptides.
244. L. A. Goldberg and M. Jerrum, *SIAM J. Comput.*, **29**, 834 (1999). Randomly Sampling Molecules.
245. R. S. Pearlman, *Chem. Des. Auto. News*, **2**, 1 (1987). Rapid Generation of High Quality Approximate 3D Molecular Structures.
246. J. Gasteiger, C. Rudolph, and J. Sadowski, *Tetrahedron Comput. Method.*, **3**, 537 (1990). Automatic Generation of 3D-Atomic Coordinates for Organic Molecules.
247. J. Sadowski and J. Gasteiger, *Chem. Rev.*, **93**, 2567 (1993). From Atoms and Bonds to Three-Dimensional Atomic Coordinates: Automatic Model Builders.
248. J. Gasteiger, J. Sadowski, J. Schuur, P. Selzer, L. Steinhauer, and V. Steinhauer, *J. Chem. Inf. Comput. Sci.*, **36**, 1030 (1996). Chemical Information in 3D-Space.

# Variable Selection—Spoilt for Choice?

David J. Livingstone\* and David W. Salt†

\**ChemQuest, Delamere House, 1 Royal Crescent, Sandown, Isle of Wight, U.K., PO36 8LZ, and, Centre for Molecular Design, University of Portsmouth, Portsmouth, U.K., PO1 3HE*

†*Department of Mathematics, Buckingham Building, University of Portsmouth, Portsmouth, U.K., PO1 3HE.*

---

---

## INTRODUCTION

In the 1960s and 1970s, when computer-aided drug design was in its infancy, little need existed for any kind of variable selection. After Corwin Hansch had applied a physical organic chemistry approach to the description of biologically active molecules,<sup>1,2</sup> drug designers chose from a limited number of tabulated physicochemical descriptors. At that time, aromatic substituents were typically characterized by  $\pi$  and  $\sigma$ , describing hydrophobic and electronic effects, and other parameters such as  $E_s$  or MR were used by designers to account for bulk.<sup>3,4</sup> Aliphatic systems were less well treated. A major disadvantage of describing chemical structure in quantitative terms this way was the need for a “parent” structure. That, in turn, restricted the generation of quantitative structure-activity relationships (QSARs) to congeneric series. Tabulated data also posed a major problem in that “gaps” in the tabulations always existed, usually for the substituent(s) that the computational chemist was most interested in. In those circumstances, it was necessary either to try to

estimate values for the missing entries or to make experimental measurements on the appropriate chemical model systems. Missing tabulated values were often a consequence of chemical instability, difficulties in synthesis, extreme values, and so on, so filling the gaps in descriptor tables was often fraught with difficulty.

Topological parameters were the first molecular descriptors that could be calculated for any chemical structure (other than obvious parameters such as molecular weight, atom counts, etc.); all that was required for their computation was a standard two-dimensional (2-D) representation of a chemical structure. The best known topological parameters are the molecular connectivity indices first described by Randic<sup>5</sup> and investigated extensively by Hall and Kier et al.<sup>6-8</sup> Molecular connectivity descriptors have since been used by chemists in the construction of QSAR models for many types of biological properties, especially applications in the environmental area. The reason for such heavy application in environmental studies is because these data sets often contain diverse sets of compounds, and as just mentioned, topological descriptors can be computed for any structure.

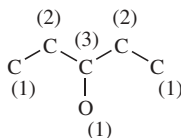
One feature of most topological descriptors, including molecular connectivity indices, is the possibility of generating many different descriptors for the same molecule. The calculation of simple molecular connectivity indices serves to illustrate this possibility. Figure 1 shows the hydrogen-suppressed graph for pentan-3-ol.

Connectivity indices are computed from this hydrogen-suppressed skeleton by the assignment of a degree of connectivity,  $\delta_i$ . This value,  $\delta_i$ , represents the number of atoms connected to the  $i$ th atom. For each bond in the structure, a bond connectivity,  $C_k$ , can be calculated by taking the reciprocal of the square root of the product of the connectivities of the atoms at either end of the bond. For example, the bond connectivity for the first carbon-carbon bond (from the left) in the structure is

$$C_1 = \frac{1}{\sqrt{(1 \times 2)}} \quad [1]$$

More generally, the bond connectivity of the  $k$ th bond is given by

$$C_k = \frac{1}{\sqrt{(\delta_i \delta_j)}} \quad [2]$$



**Figure 1** Hydrogen suppressed graph of pentan-3-ol showing bond connectivities.

where the subscripts  $i$  and  $j$  refer to the atoms at either end of the bond. The molecular connectivity index,  $\chi$ , for a molecule is found by summing the bond connectivities over all of its  $N$  bonds.

$$\chi = \sum_{k=1}^N C_k \quad [3]$$

For the pentanol shown in Figure 1, the five bond connectivities are the reciprocal square roots of  $(1 \times 2)$ ,  $(2 \times 3)$ ,  $(3 \times 1)$ ,  $(3 \times 2)$ , and  $(2 \times 1)$ , which gives a molecular connectivity of 2.808. This process is known as a first-order connectivity index because it considers only individual bonds in the structure, that is, paths of two atoms only. Higher order indices may be generated by taking longer paths and other variations including valence connectivity values, path, cluster, and chain connectivities.<sup>9</sup>

Thus, for the first time, it became possible to generate even more descriptors than molecules in a dataset, which gives rise to several problems, as discussed in the next section, but in fact, many simple molecular connectivity descriptors are not so problematic because they contain much colinearity (colinearity between a pair of descriptors means that they are highly correlated) and multicollinearity. These two properties are discussed in more detail later. High correlation between a pair of descriptors implies that they carry the same information. Although much debate has occurred about the physicochemical meaning of molecular connectivity descriptors,<sup>10-12</sup> little doubt seems to exist that molecular connectivity descriptors contain information related primarily to molecular shape.<sup>13</sup> A principal component (PC) analysis of 108 molecular connectivity descriptors for a set of  $n$ -alkanes and polychlorinated biphenyls showed<sup>14</sup> that three principal components account for 98% of the variance in the dataset. These three PCs were associated with:

- Degree of branching
- Molecular size or bulk
- Structural flexibility

The next great step forward in the quantitative description of chemical structure was caused by the rapid improvements in computing technology coupled with the availability of easy-to-use molecular modeling software. The late 1980s saw the first reports of chemists using properties calculated by molecular modeling in the construction of QSAR models<sup>15-17</sup> and, in particular, the description of the electronic properties of compounds.<sup>18,19</sup> These new computed properties describe many different aspects of chemical structure, thus confronting a “molecular designer” with problems of choice as discussed in the next section.

In the space of about 25 years, computer-aided drug design (CADD) advanced from relying almost exclusively on look-up tables of substituent constants to calculating hundreds, if not thousands, of molecular descriptors,<sup>20</sup> including experimental properties such as solubility and log *P*.<sup>21,22</sup> The excellent handbook of Todeschini and Consonni lists over 3100 separate chemical descriptors<sup>23</sup> to illustrate the depth and breadth of modern-day descriptors.

---

## THE PROBLEM

So what is the problem presented by this wealth of information for most chemical structures that can be so easily computed? Having once complained about the lack of chemical descriptors, it now seems almost churlish to say that we have an embarrassment of riches; yet, it is the large number of variables that is now the problem! With many variables, we no longer have what might be termed “complete ownership” of the data; we share it with our computer. We can no longer “eye-ball” the data in the time-honored fashion of our progenitors, seeking out obvious trends and associations. Instead, we see a “window full” of data at a time. Yes, we can calculate summary statistics of the individual variables, but those statistics tell us nothing about the interrelationships of the variables, which is what most people are interested in. We can, of course, still compute the correlation matrix, but whereas once upon a time it occupied half a page, it can now fill many pages, even with a modest number of variables (say, 50). The correlations are also of little value unless viewed alongside the corresponding scatter plots,<sup>24</sup> and although plots of the response variables against the descriptors are likely to be examined, plots of one descriptor against another are less likely to be made. What must be remembered is that multivariate methods by their nature are designed to help uncover complex relationships that are sometimes difficult to represent. Consequently, a tendency exists among some scientists to take, at face value, the results generated. Unfortunately, this tendency can be a recipe for disaster. The application of multivariate techniques require us to be even more diligent in the examination of our data as problems relating to missing values and outliers, for example, are made even worse when many variables must be considered. Many problems exist, and the type of problem that might develop will depend on what we need to do with the data. Fortunately, solutions exist to most of the problems.

If the data are to be used by chemists for comparing similarities or differences between compounds or perhaps to see if different groupings of compounds have different types of biological response, then what is needed is **Dimension Reduction**. Dimension reduction is the name given to a process that reduces the dimensionality of a multivariate dataset, while retaining most of the information that it contains. Dimension reduction is not to be



confused with variable elimination that is another issue entirely and will be considered later. The most common way in which similarities and differences are compared is visually, with so-called unsupervised learning pattern recognition methods.<sup>25</sup> In general, these methods aim to reduce the dimensionality of a dataset and are well suited to datasets containing many descriptors. Some of them, however, may fail if more variables are available than compounds; i.e., if  $p$  (variables)  $> n$  (compounds), then we refer to the data matrix as being “over square” (they are wider than they are deep). In these circumstances, the variables are linearly dependent and at least  $p - n$  of them contribute no information that is not already contained in the other variables. In this situation, the matrix product  $\mathbf{X}'\mathbf{X}$ , where  $\mathbf{X}$  is the data matrix and  $\mathbf{X}'$  is its transpose, is singular meaning that its determinant is zero. Singular matrices do not have an inverse, which is problematic in multivariate statistical methods because several important techniques depend on calculating the inverse of  $\mathbf{X}'\mathbf{X}$ . We will look at this problem in more detail in a later section.

A linear dependency between variables also presents problems when constructing models that correlate physicochemical properties and some type of biological response. In this case, and even when fewer descriptors exist than compounds, a common way to proceed is to use what is known as **Variable Elimination**. Variable elimination seeks to reduce redundancy in a dataset by identifying variables, which are correlated with one another or with combinations of other variables in the set. Finally, we point out that the aim of an analysis may be to identify the “important” variables either in their own right or for the construction of mathematical models. When the analyst is concerned with identifying useful variables to include in the model, the process is known as **Variable Selection**. The rest of this chapter sets out to illustrate the uses of these three processes with examples where possible, discussing the pros and cons of each, and indicating for the reader areas where questions and ambiguities remain.

---

## DIMENSION REDUCTION

Dimension reduction is, as the name implies, a technique for reducing the dimensionality of a dataset, which is most often applied to the columns (variables) but may also be applied to the rows (cases or compounds) and results in a reduction from  $p$  variables to  $q$  variables or dimensions where  $q$  is often 2 or 3 (for ease of display of the resulting data matrix). A common method of dimension reduction is principal component analysis (PCA). A less-frequently used but related method is factor analysis (FA). Insufficient space exists here for a complete description of these techniques, so the reader is directed to references 26 and 27 for PCA and 28 and 29 for FA. Briefly, each computes new variables

that are derived from combinations of the original variables. In the case of PCA, the new variables are called principal components (PCs); thus:

$$\begin{aligned} PC_1 &= a_{11}x_1 + a_{12}x_2 + \cdots + a_{1p}x_p \\ PC_2 &= a_{21}x_1 + a_{22}x_2 + \cdots + a_{2p}x_p \\ &\vdots \\ PC_q &= a_{q1}x_1 + a_{q2}x_2 + \cdots + a_{qp}x_p \end{aligned} \quad [4]$$

where the subscripted term,  $a_{ij}$ , shows the contribution of the original variable,  $x_j$ , to the  $i$ th PC. Equation [4] shows the generation of  $q$  PCs from a set of  $p$  variables. The generation of PCs is carried out to conform to three conditions:

- The first principal component explains the maximum variance (information) in the dataset. Subsequent components describe the maximum part of the remaining variance.
- The principal components are orthogonal to one another.
- As many PCs may be extracted as the smaller of  $n$  (data points) or  $p$  (dimensions) for a  $n \times p$  matrix (denoted by  $q$  in Eq. [4]). Actually, it is the rank of the matrix, denoted by  $r(\mathbf{A})$ , which is the maximum number of linearly independent rows (or columns) in  $\mathbf{A}$ .  $0 \leq r(\mathbf{A}) \leq \min(n, p)$ , where  $\mathbf{A}$  has  $n$  rows and  $p$  columns.

Factor analysis differs from PCA in that each variable is assumed to be made up of contributions from several common factors and a single unique factor. The common factors, as the name implies, are common to all variables in the set. The unique factors account for the variance of each variable, which is not explained by the common factors; unique factors may be noise or measurement error, which is associated with particular variables. As each variable is made up of contributions from all common factors, the argument can be turned around so that the common factors, normally just called factors, are made up of contributions from all variables in the set. Thus, each factor is composed of a weighted contribution from each of the original variables, just like the PCs in PCA. The major difference between PCA and FA is that some of the original variance in the dataset is discarded in the generation of the unique factors, which are usually ignored.

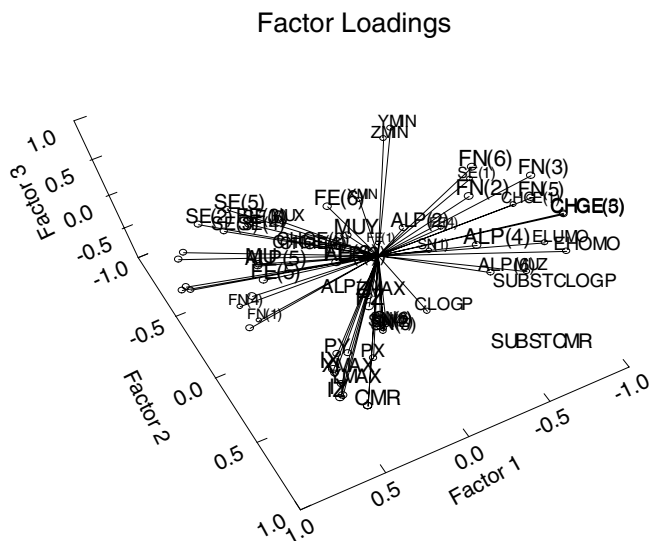
A simple chemical example may serve to show the utility of PCA as a dimension reduction method both for the display of multidimensional data and for model building. The example we give here is a chemical model system for the characterization of the intermolecular forces involved in electron-donor-acceptor complexes, also called charge-transfer complexes, which has been described previously.<sup>30</sup> Creation of this model involved the measurement of complex formation constants between monosubstituted benzenes and the

strong electron acceptor, 1,3,5 trinitrobenzene. A substituent constant,  $\kappa$ , was derived from these measurements and in the early study was shown to correlate with a combination of three “standard” substituent parameters:

$$\kappa = 0.04MR - 0.33\pi - 0.21R + 0.02 \quad [5]$$

A later examination of this system with computational chemistry to characterize the compounds resulted in a dataset of 58 computed descriptors for 40 molecules.<sup>31</sup> This dataset is wider, by 18 descriptors, than it is long so it is clearly going to present some problems in its analysis. As noted earlier, the “true” dimensionality of this set is 40 or less, so some variable elimination is necessary. Before we do this, however, it may be instructive to examine the relationships between the variables with some sort of visual display. Factor analysis carried out on this set gave 11 factors with eigenvalues greater than 1, a common test of the “significance” or importance of factors and principal components. Figure 2 shows a plot of the factor loadings of the variables on the first three factors.

The lines in this plot are vectors drawn from the origin (0,0,0) of the three factors in which the length shows the value of the loadings, with longer lines having larger loadings, and the position shows the relationship of that variable with each of the factors. Descriptors that are close to one another on the plot have similar loadings and are thus correlated with one another. This diagram is complex to interpret, because so many variables exist, but

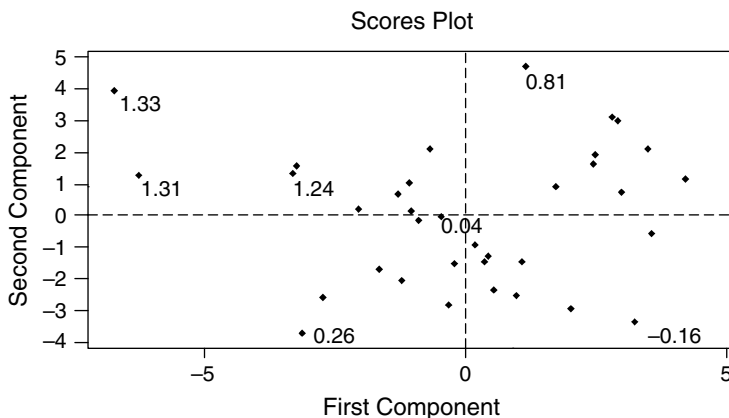


**Figure 2** Factor loadings plot of the first three factors computed from a set of 58 calculated chemical descriptors.

some features can be seen easily. A group of “size” parameters (CMR, Px, Iz, and so on) is at the bottom of the picture. ClogP appears to be separate from most other descriptors in the right of the diagram. LUMO and HOMO energies are grouped along with some self-atom polarizabilities, and so on. Viewing a factor loadings plot is a useful way to gain some insight into the inter-relationships between descriptors.

Removal of constant, or nearly constant, variables along with those having high pairwise correlations with other descriptors resulted in a reduced dataset of 31 parameters. A principal component analysis on this set yielded nine PCs with eigenvalues greater than 1, which is two less than that from the factor analysis of the starting set. Clearly some information has been removed; whether this is useful information or merely “noise” cannot, at this stage, be said. A plot of the PC scores for the compounds, calculated as shown in Eq. [4], is given in Fig. 3. Values of substituent constant  $\kappa$  have been marked on the plot for some of the compounds to show that some of the larger values are grouped in the left-hand, top quadrant. A perfect gradation of  $\kappa$  values does not occur across any part of the plot, but this scores plot does group chemically similar substituents and gives an indication that the descriptors contain information that should be useful in modeling  $\kappa$ .

This scores plot is a two-dimensional representation of the original 31 dimensional space. Points (compounds) that are close together in this plot were close together in the space described by the 31 descriptors; that is, they have similar values for most of the 31 properties. This plot, of course, is an approximation of the original 31 dimensional space but because the first two PCs explain maximum amounts of variance in the dataset, they are the “best” two-dimensional representation of that space, in terms of variance explained at least. Another form of dimension reduction that is often used



**Figure 3** Scores plot for 35 compounds with measured  $\kappa$  values. The scores are the principal component scores for the first two PC's computed from 31 descriptors.

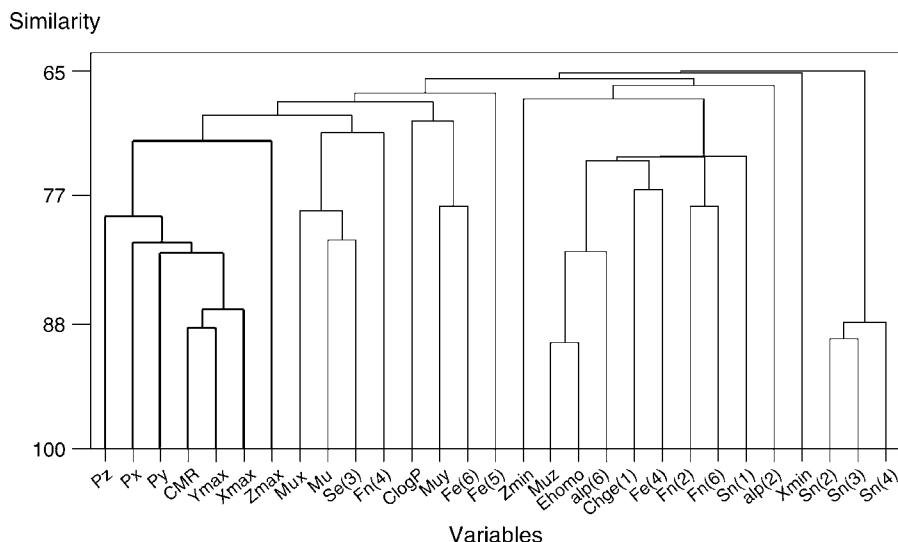


Figure 4 Dendrogram showing the associations between 31 computed descriptors.

for data inspection by chemists, is cluster analysis. Many “flavors” of cluster analysis exist (see Chatfield and Collins<sup>32</sup> and Willett<sup>33</sup> for discussion and examples), but broadly speaking, they all aim to find groupings in the data based on some measure of similarity. The similarity measure may be a measure of correlation, such as a correlation coefficient, or a measure of distance, such as Euclidean distance, in some relevant space. A common method for reporting the results of a cluster analysis is a diagram called a dendrogram in which the objects being clustered are displayed in connected groups. Figure 4 shows the relationships between the 31 variables remaining in the charge-transfer dataset.

It can be seen from the figure that a group of size descriptors, shown in bold on the left, is clustered just as was shown in the factor loadings plot in Figure 2.

Other more or less common dimension reduction techniques like artificial neural networks<sup>34</sup> exist, but space here does not permit further discussion of them.

## VARIABLE ELIMINATION

Variable elimination should not be confused with dimension reduction, which as described, may be undertaken with PCA or FA. Variable elimination, however, is a process where variables are physically removed from the dataset and are thus not applied in any model building. This topic finds little coverage

in the literature. Performing a search on the web, we inevitably find references to what are termed variable elimination or variable reduction methods but that turn out to be dimensionality reduction methods with PCA or FA. Having confirmed what may have been obvious to some that variable elimination really means what it says, the questions that need to be answered are as follows:

1. Why do we eliminate variables; aren't we just throwing information away?
2. Which variables do we eliminate; what are the criteria for deciding which variables should go?

## Why We Eliminate Variables

Variables are eliminated for two reasons. First, they are eliminated if they have a small variance, below some threshold value. It is not uncommon to find at least one variable that is constant or is nearly constant, with all but one entry being different from the others, even for moderately sized datasets. This situation can occur when a researcher includes variables (dummy variables) to record the presence or absence of particular properties of an object without checking how well the particular categories are represented. The problem is that when the data are split into training/validation groupings, it can happen that such variables take one value in one group and the other value in the other group.

The second reason for variable removal is if they are redundant.<sup>35</sup> Redundant variables develop in a dataset for two reasons: (1) because more variables ( $p$ ) exist than objects ( $n$ ) and (2) multicollinearity.

## More Variables Than Objects

The situation in which the data matrix has more variables (columns) than objects (rows) was touched on earlier, but we will now consider it in more detail as it is an important problem we encounter with large datasets. Most people are aware that with two data points, it is possible to construct a line, a one-dimensional object, and with three data points, a plane, a two-dimensional object. This process can be continued so that four data points allows a three-dimensional object, five points, four dimensions, and so on. Thus, the maximum dimensionality of an object, and hence the maximum number of dimensions, in a dataset is  $n - 1$ , where  $n$  is the number of data points. For dimensions, we can substitute "independent pieces of information," and thus, the maximum that any dataset may contain is  $n - 1$ . This, however, is a **maximum** and in reality the true dimensionality, where dimension means "information," is often much less than  $n - 1$ .

From experience, we know that some researchers need a bit more convincing that with  $p > n$ , at least  $p - n$  variables are not contributing any

additional information about the objects (compounds) in the dataset that is not already present in the others. The reason for this is that when  $p > n$  the variables (columns) of the data matrix are linearly dependent. To see what this means and to see what effect this has on our analysis, let us introduce some notation to aid the discussion. Denote the data matrix by  $\mathbf{X}$ , i.e.,

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

where  $x_{ij}$  is the observation for compound  $i(1, 2, \dots, n)$  on variable  $j(1, 2, \dots, p)$  that defines  $\mathbf{x}_j$ . The variables of  $\mathbf{X}$  are linearly dependent if it is possible to write

$$a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \cdots + a_p \mathbf{x}_p = 0 \quad [6]$$

such that the coefficients  $a_j$  are not all equal to zero. That is, if we can find values of the coefficients (not all zero) such that Eq. [6] is true, then it means that we can write some variables in our dataset as linear combinations of others. Consequently, the observations for these variables can be generated from others, thus making them redundant. To illustrate this point, let us take a simple numerical example where we have the  $2 \times 3$  ( $n \times p$ ) data matrix

$$\mathbf{X} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \text{ so that the three variables with two observations are } \mathbf{x}_1 = \begin{bmatrix} 1 \\ 4 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 2 \\ 5 \end{bmatrix}, \mathbf{x}_3 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}. \text{ It can easily be shown that } \mathbf{x}_1 - 2\mathbf{x}_2 + \mathbf{x}_3 = 0,$$

which satisfies Eq. [6], the condition for linear dependence. Rewriting this result, we have  $\mathbf{x}_3 = -\mathbf{x}_1 + 2\mathbf{x}_2$ , which shows that the third variable is simply a linear combination of the other two and we say that the variables are coplanar; i.e., they lie in the same plane, which means that if we already have  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , nothing is to be gained from the inclusion of  $\mathbf{x}_3$ , as its information is already contained in the other two. Therefore, the actual dimension of the data matrix is only 2 (the number of observations) and not, as might be thought, 3 (the number of variables). In the previous example,  $\mathbf{x}_1 - 2\mathbf{x}_2 + \mathbf{x}_3 = 0$  could have been rewritten with any one of the three variables on the left-hand side. This, however, is not always the case as the slightly more com-

plex problem with the  $3 \times 5$  matrix  $\mathbf{X} = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 5 & 4 & 3 & 2 & 1 \\ 5 & 3 & 4 & 5 & 1 \end{bmatrix}$  will illustrate.

It again can readily be shown that  $0\mathbf{x}_1 - \mathbf{x}_2 + 2\mathbf{x}_3 - \mathbf{x}_4 = 0$  and

$-5\mathbf{x}_1 + 4\mathbf{x}_2 + 4\mathbf{x}_3 - 3\mathbf{x}_5 = 0$  so that the variables  $\mathbf{x}_4$  and  $\mathbf{x}_5$  can be written as linear combinations of the other three. As in the first example, the dimension of the data matrix is less than the number of variables being equal to 3 (the number of rows of data). From these two examples, we see that if  $p > n$ , then  $p - n$  of the variables can be written as linear combinations of the remaining variables. This situation holds provided that the  $n$  variables chosen are not linearly dependent. Hence, although variables  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ , and  $\mathbf{x}_3$  can span the dimensions of this dataset in this second example, variables  $\mathbf{x}_2$ ,  $\mathbf{x}_3$ , and  $\mathbf{x}_4$  cannot because they are linearly dependent ( $-\mathbf{x}_2 + 2\mathbf{x}_3 - \mathbf{x}_4 = 0$ ).

In the two examples considered so far, the number of variables ( $p$ ) has exceeded the number of objects ( $n$ ) and  $n - p$  of the variables could be “explained” for  $n$  of them, provided that  $n$  can be found that are not linearly dependent. The question that now comes to mind is, what happens if  $n$  linearly independent variables cannot be found? The answer is, you find the maximum number that exist. This number is called the rank of the data matrix and is usually denoted by  $r$  and is such that  $r \leq \min(n, p)$ . This result is important as it places an upper bound on the number of variables that can usefully be employed in a statistical analysis of a dataset. In both examples,  $r = 3 (= n)$ , but often it is less than  $n$ . Before we leave ideas of rank and move on to the second reason for variables being redundant, let us consider the effect that mean centering of data has on the rank of the  $\mathbf{X}$  matrix. Mean centering has the effect of reducing the rank of the data matrix by unity, i.e., if the rank of the raw data matrix is  $r$ , then the rank of the mean centered matrix is  $r - 1$ , which means that for mean centered data, the rank of the data matrix  $r \leq \min(n - 1, p - 1)$ . Mean centering is the first step in variable standardization (see the section on principal component regression, partial least squares, and continuum regression) and is a device used by statisticians to simplify formulas in situations where the variance of a variable is important but not its central location, e.g., in principal component analysis. To illustrate this property of mean centering, let us use the  $3 \times 5$  data matrix in the second of the previous examples. The means of the five variables are 3.67, 3.00, 3.33, 3.67, and 2.33 and subtracting them from their respective observations produces the mean centered data matrix

$$\dot{\mathbf{X}} = \begin{bmatrix} -2.67 & -1 & -0.33 & 0.33 & 2.67 \\ 1.33 & 1 & -0.33 & -1.67 & -1.33 \\ 1.33 & 0 & 0.67 & 1.33 & -1.33 \end{bmatrix}$$

where the dot notation indicates that the entries in the matrix have been mean centered. Now as we have stated, the number of linearly independent columns is the rank of a matrix, but the rank is also the number of linearly independent rows. As can be seen, the rows of  $\dot{\mathbf{X}}$  sum to zero so the rank is less than 3 and by observation can be seen to equal 2. So mean centering the data reduces the rank of a data matrix by one.

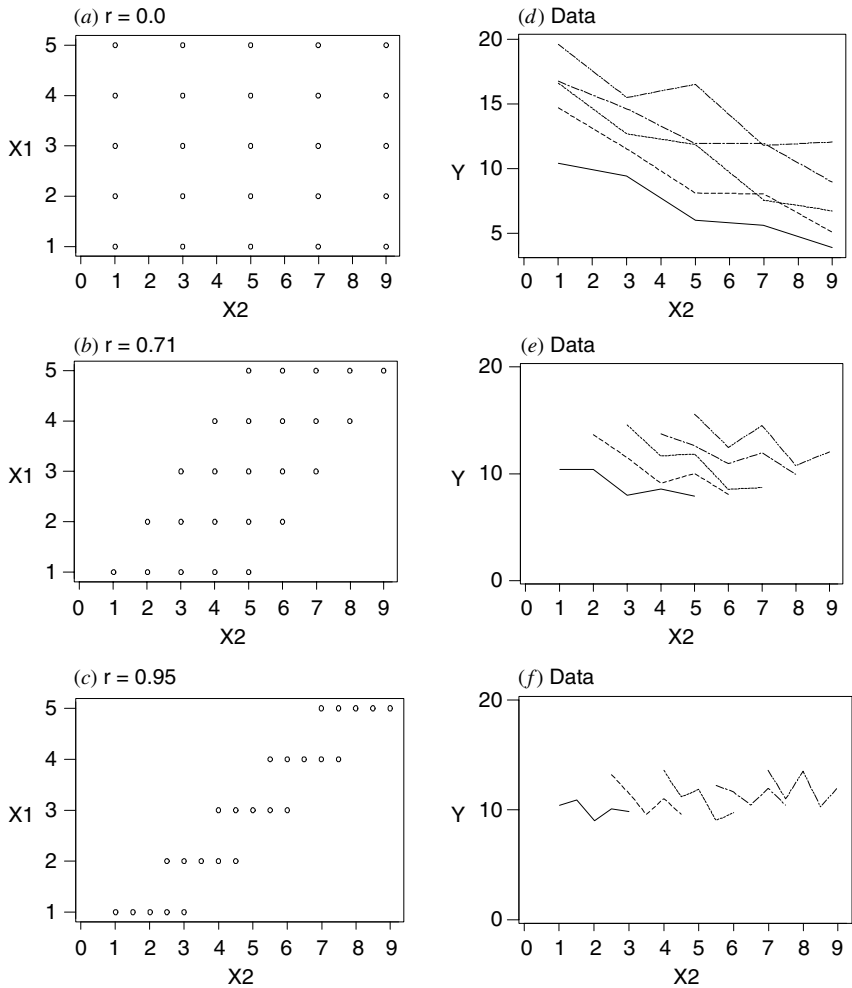


## Multicollinearity

The second reason for variables being redundant is multicollinearity. Multicollinearity is a term that denotes the presence of linear relationships (or near-linear relationships) among the variables in the  $\mathbf{X}$  matrix. With the exception of some design matrices in analysis of variance, multicollinearity exists to some degree in all datasets with  $p \leq n$ . The closer the linear combination of variables comes to satisfying the condition of Eq. [6], the stronger the multicollinearity and potentially the greater the problem if it goes undetected. The effect of perfect multicollinearity (linear dependency) has already been discussed in that it results in a singular  $\mathbf{X}'\mathbf{X}$  matrix that cannot be inverted, i.e., its inverse  $(\mathbf{X}'\mathbf{X})^{-1}$  does not exist. When the condition of Eq. [6] holds approximately, researchers performing the analysis may be totally unaware of the high degree of multicollinearity that exists in their data as the  $\mathbf{X}'\mathbf{X}$  matrix is now not singular, so that its inverse can be computed. However, any reasonable data analysis software will alert the user to the situation with the degree of multicollinearity tolerated being user defined. We will consider this last point later when we discuss methods of detecting multicollinearity.

So what exactly is the problem with multicollinearity that we have alluded to? Strong relationships between independent variables in multiple regression make it difficult, if not impossible, to identify the individual effects that each has on the response variable. Multicollinearity may result in highly unstable and inflated regression coefficient estimates, because their values may change enormously when variables are deleted or added to the regression, or when small changes are made to data values. The instability in the regression coefficients can be dramatic, which results in some instances in sign reversal,<sup>36</sup> making structure activity research somewhat problematic. Rather than embark on a detailed mathematical account of the effect of multicollinearity between the regressor variables (this can be found elsewhere<sup>37</sup>), we will illustrate its effect by means of a numerical example. To see clearly the effect of multicollinearity between the regressor variables, we have chosen the simplest situation possible, i.e., a two independent variable multiple regression. (N.B. the situation when two regressor variables are correlated is referred to by some authors as colinearity<sup>38</sup>).

To assess the effect of correlation structure on fitting the two variable model, we have generated three datasets that are shown in Figure 5. In dataset 1, the values of the explanatory variables  $x_1$  and  $x_2$  are arranged on a  $5 \times 5$  grid. For each value of  $x_2$ , the values of  $x_1$  are regularly spaced over the same range and **no correlation** exists between them. In dataset 2, the overall ranges of  $x_1$  and  $x_2$  are the same as in set 1 but now the range of values of  $x_1$  for a given value of  $x_2$  **depends** on  $x_2$ . The variables are positively correlated with  $r = 0.71$ . For dataset 3, the overall ranges are again kept the same but the range of  $x_1$  for a given value of  $x_2$  is **further restricted** to increase the correlation coefficient to 0.95.



**Figure 5** Examples to illustrate the effect of collinearity between the explanatory variables on a two variable linear model. (a)–(c) Different correlations; (d)–(f) simulated data from two variable linear model. Key: ———  $x_1 = 1$ ; - - -  $x_1 = 2$ ; — · —  $x_1 = 3$ ; — — —  $x_1 = 4$ ; . . .  $x_1 = 5$ .

To compare analyses based on these three sets of  $x$ -values, it is necessary to generate the corresponding  $y$ -values. This process is performed here by adopting the linear model in Eq. [7]

$$E(Y) = 10 + 2x_1 - x_2 \tag{7}$$

where  $E(Y)$  is the true mean response. Any two variable linear model would have sufficed, and nothing is special about the one we have used here. To

generate observed values of our response variable, we evaluate  $Y = 10 + 2x_1 - x_2 + \varepsilon$  for each pair of values of our  $x$ -variables and add on an error term  $\varepsilon$ , where  $\varepsilon \approx N(0, 1)$ , i.e., being normally distributed with a zero mean and a variance of unity. Before we look at the analyses of these three datasets, it is worth pointing out that the only thing responsible for the differences in the datasets is the  $\mathbf{X}$  matrix, since we have used the same parameter values and random errors in each case. The three datasets are depicted in Figure 5(d)–(f). We have plotted  $y$  versus  $x_2$  and have indicated the variation in the  $x_1$  direction by joining the simulated data points for the different values of this variable. Even before we look at the numerical detail of the regression analysis of these datasets, the effect of the colinearity between the two regressor variables is apparent from the plots. As the correlation between  $x_1$  and  $x_2$  increases the area of the “footprint” of the data in the  $x_1 - x_2$  plane becomes smaller and the variation of  $y$  as a result of changes in  $x_1$  decreases markedly.

In each case, the model was fitted by least squares and the results summarized for (1) sequential partitions of the regression sums of squares (Table 1); (2) the estimated regression coefficients and their standard errors (Table 2). The main points of the analyses are as follows:

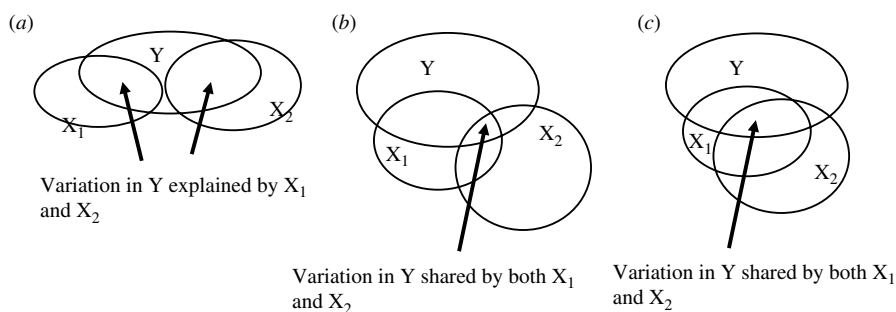
1. For dataset 1, the contribution of each variable to the regression sum of squares is the same no matter what order the variables are added when fitting the model. As a consequence, the regression sum of squares can be

**Table 1** Sequential Partition of the Regression Sum of Squares (SS)

Data	Fitting $x_1$ First		Fitting $x_2$ First	
	Source	SS	Source	SS
1 $R^2 = 94\%$	$x_1$ alone	164.14	$x_1$ after $x_2$	164.14
	$x_2$ after $x_1$	198.41	$x_2$ alone	198.41
	$x_1$ and $x_2$	362.55	$x_1$ and $x_2$	362.55
2 $R^2 = 78\%$	$x_1$ alone	32.955	$x_1$ after $x_2$	81.350
	$x_2$ after $x_1$	49.207	$x_2$ alone	0.812
	$x_1$ and $x_2$	82.162	$x_1$ and $x_2$	82.162
3 $R^2 = 43\%$	$x_1$ alone	4.863	$x_1$ after $x_2$	15.983
	$x_2$ after $x_1$	12.105	$x_2$ alone	0.984
	$x_1$ and $x_2$	16.967	$x_1$ and $x_2$	16.967

**Table 2** Estimated Regression Coefficients (standard errors)

Data	Fitting Both Variables		Fitting Single Variables	
	$x_1$	$x_2$	$x_1$	$x_2$
1	1.8118 (0.1434)	−0.9960 (0.0717)	1.8118 (0.4384)	−0.9960 (0.2015)
2	1.8039 (0.2028)	−0.9920 (0.1434)	0.8118 (0.2499)	−0.0901 (0.2126)
3	1.7880 (0.4535)	−0.9841 (0.2868)	0.3118 (0.1738)	0.0887 (0.1159)



**Figure 6** Diagram illustrating the ideas of decomposition of regression sums of squares for (a) data set 1, (b) data set 2 and (c) data set 3.

partitioned into separate contributions for each variable. For datasets 2 and 3, in contrast, the contributions depend on the order of fitting. The effect is particularly striking for set 2, where, for this case, only a small amount of variation ( $SS = 0.812$ ) is accounted for by fitting  $x_2$  alone; after fitting  $x_1$ , the SS increases to 49.207. A similar effect is observed for dataset 3. Note that the proportion of the variation explained by the fitted model decreases from 94% with two independent regressors to 43% when they are highly collinear. These results are illustrated in Figure 6 where the total sum of squares for the response variable is represented by the region inside of the oval labeled Y. The regions labeled  $x_1$  and  $x_2$  represent the information content of the two regressor variables and where these regions intersect with the Y's may be thought of as the information in Y (sum of squares) that is explained by the two X's. In dataset 1, where the X's are independent of one another, it can be seen that they make separate contributions to Y (the X regions do not overlap), but as the correlation between the two regressor variables increases, so does the overlap of the regions representing their information content. As a consequence of this, the region where the intersection of the X's overlaps that of Y becomes a major part of the information shared between the individual X's and Y. The consequence is that when these correlated descriptor variables are entered into the model, the first one takes the "lions share" of the variation and the second just picks up the "crumbs." In extreme cases, the variation left unexplained may be compatible with that associated with the error term,  $\epsilon$ , which means that the second (or subsequent) variable has nothing to explain and it is then a matter of luck, whether the regression coefficient takes a positive or a negative sign.

- When fitting  $x_1$  and  $x_2$ , the standard error of the estimated regression coefficient for  $x_1$  is smallest for set 1 (0.1434), intermediate for set 2 (0.2828), and largest for set 3 (0.4535), which is about three times that for set 1. The reduction in precision is caused by the restricted range of values

of  $x_1$  for given values of  $x_2$ . In other words, the effect of the correlation between the explanatory variables is to reduce the precision of the estimated regression coefficients. For the variable  $x_1$ , we can attribute the reduction in precision to the effect of the correlation because the values of  $x_1$  have been kept the same deliberately in each dataset (because this also effects the precision). The reduction in precision occurs with  $x_2$  although in this case the interpretation is less straightforward because the datasets differ in the distribution of their  $x_2$  values.

3. Table 2 also shows estimated regression coefficients from fitting the variables singly. Because of the independence of  $x_1$  and  $x_2$  in dataset 1, the estimated regression coefficient of each variable is the same whether the other variable is included in the fitted model. For sets 2 and 3, this similarity is not true and failing to allow for variation in both variables can give misleading estimated coefficients. For example, although the coefficient of  $x_2$  in the single variable regression for data set 3 is not significant, a sign change occurs when  $x_1$  is included.

Before leaving this example, it is instructive to see how the  $(\mathbf{X}'\mathbf{X})^{-1}$  matrix varies across the three datasets. We consider this variation because it was problems with this matrix that were first raised in the discussion of multicollinearity. The three matrices are as follows:

Data 1	Data 2	Data 3
$\begin{pmatrix} 0.020 & 0.000 & 0.000 \\ 0.000 & 0.005 & 0.000 \\ 0.000 & 0.000 & 0.040 \end{pmatrix}$	$\begin{pmatrix} 0.04 & -0.02 & 0.00 \\ -0.02 & 0.02 & 0.00 \\ 0.00 & 0.00 & 0.04 \end{pmatrix}$	$\begin{pmatrix} 0.20 & -0.12 & 0.00 \\ -0.12 & 0.08 & 0.00 \\ 0.00 & 0.00 & 0.04 \end{pmatrix}$

As can be seen, the  $(\mathbf{X}'\mathbf{X})^{-1}$  matrix is symmetrical and the submatrix formed by the first two entries in the first two columns plays an important role in what we have to discuss at this point. Also, it can be shown that the entries in this

$2 \times 2$  submatrix are proportional to  $\frac{1}{1 - r_{12}^2}$ , where  $r_{12}^2$  is the simple product

moment correlation between  $x_1$  and  $x_2$ . It can also be shown that the leading diagonal elements of this smaller matrix are proportional to the variances of the regression coefficient estimates, and the off-diagonals are proportional to

their covariance. Thus, as  $r_{12}^2$  approaches unity, the magnitude of  $\frac{1}{1 - r_{12}^2}$  will

“explode” along with the variances of the regression coefficient estimates and their covariance. Because we have used the same  $x_1$  and a set of random error values, this proportionality result can be checked easily by comparing the first entry in each of the previous three matrices. For example, dividing the first entry for dataset 1 (0.020) into those for datasets 2 (0.04) and 3 (0.20) gives 2 and 10 (approximately), respectively. As the constant of proportionality is

the same for  $x_1$ , dividing the 0.04 by the 0.02 will give the same result as calculating  $\frac{1}{1-r_{12}^2}$  with  $r_{12}^2 = (0.71)^2 = 0.5041$ , i.e.,  $\frac{1}{1-0.5041} \approx 2$ , and repeating this for dataset 3 gives  $\frac{1}{1-0.9025} \approx 10$ . These two results illustrate clearly the effect that colinearity between two regressor variables has on the variance of the estimated regression coefficients. With a correlation of 0.71 between the regressors, a doubling of the error variances occurs, and when this colinearity increases to 0.95, a 10-fold increase in error variance occurs. This argument extends to the situation where more than two regressors with the diagonal elements of  $(\mathbf{X}'\mathbf{X})^{-1}$  are proportional to  $\frac{1}{1-R_i^2}$ , where  $R_i^2$  is the proportion of the total variance in the  $i$ th regressor variable explained by its regression on the other regressor variables in the model.  $\frac{1}{1-R_i^2}$  is sometimes referred to as the variance inflation factor (for obvious reasons)<sup>39</sup> with  $1-R_i^2$  labeled as the tolerance.

---

## IDENTIFICATION OF MULTICOLLINEARITY

Several methods have been published for detecting multicollinearity. Because the focus of this chapter is on variable selection methods, we will merely list some of the earlier popular methods and discuss in more detail a method developed specifically to assess the “amount” of multicollinearity present together with the group or groups of variables involved. A more detailed account of these methods can be found elsewhere.<sup>40</sup>

The most obvious approach for detecting multicollinearity is to examine the pairwise correlations in the correlation matrix. Although high pairwise correlations are an indication of colinearity, their absence does not infer a lack of multicollinearity; it is possible for all simple pairwise correlations to be relatively low; yet a high degree of multicollinearity may exist, i.e., a near-linear relationship involving three or more variables. The variance inflation factor is another popular approach, and although it gives an overall indication of multicollinearity, it cannot identify the variables involved or indeed if several groups of multicollinear variables exist. Other methods for detecting the existence of multicollinearity involve monitoring the size, sign, and standard error of regression coefficients for variables already in a model when another variable is added or removed from the model. If one or more estimated regression coefficients change noticeably or if their standard errors increase, the most recent variable included in the model is identified as being multicollinear with existing variables in the model. One extremely straightforward method for assessing the presence of multicollinearity is to calculate the eigenvalues of the correlation matrix of the data. With  $p$  variables ( $p < n$ ),

$p$  eigenvalues  $\lambda_i, i = 1, 2, \dots, p$  will exist such that  $\lambda_1 > \lambda_2 > \dots > \lambda_p$ . In the situation where all variables are independent of one another (orthogonal), we have

$$\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \frac{1}{\lambda_i} = p. \quad [8]$$

When the variables are not orthogonal

$$\sum_{i=1}^p \frac{1}{\lambda_i} > \sum_{i=1}^p \lambda_i = p \quad [9]$$

and, accordingly, large values of  $\sum_{i=1}^p \frac{1}{\lambda_i}$  occur when high multicollinearity is present. This result is consistent with the earlier discussion on small eigenvalues, which as  $\lambda_i$  becomes small,  $\frac{1}{\lambda_i}$  becomes large and the sum of such terms explodes. Although this method has potential for discovering whether multicollinearity exists, the problem of deciding what is meant by a small eigenvalue remains. We will evaluate how well these approaches work when applied to a real dataset after we have introduced a method developed specifically to diagnose the degree of multicollinearity together with the variables involved. The method we are about to discuss consists of two stages: (1) the calculation of condition indexes and (2) the decomposition of the regression coefficient variance matrix.

The condition index is calculated as  $\frac{\mu_{\max}}{\mu_{\min}}$ , where  $\mu_{\max}$  and  $\mu_{\min}$  are, respectively, the maximum and minimum singular values of the data matrix  $\mathbf{X}$ .<sup>41</sup> To obtain the singular values, we perform what is termed a singular-value decomposition (SVD) of  $\mathbf{X}$  (see the Appendix for details). The larger the value of the condition index, the more ill-conditioned (when an inverse of a given matrix explodes) is the data matrix. If multicollinearity is present to a high degree, it is possible that the multicollinearity is a result of several separate groups of colinear or multicollinear variables. To determine how many such groups exist, we calculate  $\frac{\mu_{\max}}{\mu_j}$ , termed the  $j$ th condition index, for  $j = 1, 2, \dots, p$ , and argue that as many near-linear dependences exist as high-valued condition indexes. The decision about how large a condition index should be to be considered high can only be determined empirically. From the work of Belsley et al.,<sup>41</sup> a value around 15 or more appears to work well in practice. Therefore, the number of condition indexes greater than 15 signifies the number of near-linear dependences among the variables in the data matrix  $\mathbf{X}$ .

Although the condition index approach alerts us to the number of groups of linear dependences among the variables in  $\mathbf{X}$ , it does not tell us which variables are involved. To determine which variables are involved, it is necessary to decompose the estimated variance of each regression coefficient into a sum of terms such that each term in the sum is associated with an eigenvalue of the data matrix  $\mathbf{X}$ . The estimated variance–covariance matrix of the least squares estimates of the regression coefficients,  $C_v(\hat{\boldsymbol{\beta}})$ , is given by  $C_v(\hat{\boldsymbol{\beta}}) = s^2(\mathbf{X}'\mathbf{X})^{-1}$ , where  $s^2$  is the error mean square from the regression. As shown in the Appendix,  $C_v(\hat{\boldsymbol{\beta}})$  can be rewritten, following the SVD of  $\mathbf{X}$ , for the matrix of singular values  $\mathbf{S}$  and the matrix of eigenvectors  $\mathbf{V}$ , i.e.,

$$C_v(\hat{\boldsymbol{\beta}}) = s^2(\mathbf{X}'\mathbf{X})^{-1} = s^2\mathbf{V}\mathbf{S}^{-2}\mathbf{V}' \quad [10]$$

From Eq. [10], we can write the variance of  $\hat{\beta}_j$ , the  $j$ th component of  $\hat{\boldsymbol{\beta}}$  as

$$V(\hat{\beta}_j) = s^2 \sum_{i=1}^p \frac{v_{ji}^2}{\mu_i^2} \quad [11]$$

where the  $\mu_i$ 's are the singular values of  $\mathbf{X}$ , and  $v_{ji}$  is the  $j$ th element of  $\mathbf{V}$ . From Eq. [11], we can write the proportion of the variance of the  $j$ th regression coefficient associated with the  $i$ th component ( $\pi_{ji}$ ) of its decomposition as

$$\pi_{ji} = \frac{\omega_{ji}}{\sum_{i=1}^p \omega_j} \quad [12]$$

where  $\omega_{ji} = \frac{v_{ji}^2}{\mu_i^2}$ . The variables participating in a particular multicollinear relationship are those with high variance decomposition proportions associated with the same high condition index. Again following the work of Belsley et al.,<sup>41</sup> proportions in excess of about 0.50 identify variables that have a high degree of multicollinearity that might be detrimental to the analysis.

To illustrate this two-stage approach to identifying multicollinearity and the variables involved, a dataset with two groups of variables that have a high degree of collinearity was fabricated. The dataset contains nine random variables  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_9$  plus two others  $\mathbf{x}_{10} = \mathbf{x}_3 + \mathbf{x}_4 + \varepsilon$  and  $\mathbf{x}_{11} = \mathbf{x}_5 + \mathbf{x}_6 + \mathbf{x}_7 + \delta$ . These final two variables are a means of introducing known multicollinearity structures into the dataset where the random terms  $\varepsilon$  and  $\delta$  break the perfect multicollinearity structures that would otherwise be present.

The variance–decomposition proportions are shown in Table 3. Two condition indexes (C.I.) are greater than 15 (27.5 and 59.7). With the first of these indexes, we have the four variables  $\mathbf{x}_5$ ,  $\mathbf{x}_6$ ,  $\mathbf{x}_7$ , and  $\mathbf{x}_{11}$  with proportions in excess of 0.5 consistent with one of the two multicollinearity structures



**Table 3** Variance—Decomposition Proportions for Simulated Data

C.I.	$V(\hat{\beta}_0)$	$V(\hat{\beta}_1)$	$V(\hat{\beta}_2)$	$V(\hat{\beta}_3)$	$V(\hat{\beta}_4)$	$V(\hat{\beta}_5)$	$V(\hat{\beta}_6)$	$V(\hat{\beta}_7)$	$V(\hat{\beta}_8)$	$V(\hat{\beta}_9)$	$V(\hat{\beta}_{10})$	$V(\hat{\beta}_{11})$
1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00
1.20	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.01	0.00	0.00	0.00	0.00
1.56	0.00	0.00	0.00	0.00	0.00	0.01	0.05	0.00	0.00	0.00	0.00	0.00
2.00	0.08	0.00	0.00	0.00	0.00	0.00	0.37	0.08	0.01	0.00	0.00	0.00
2.20	0.17	0.00	0.00	0.00	0.00	0.00	0.38	0.00	0.00	0.00	0.00	0.00
2.27	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.55	0.00	0.00	0.00
2.32	0.38	0.00	0.00	0.00	0.00	0.00	0.02	0.04	0.24	0.00	0.00	0.00
2.43	0.00	0.00	0.00	0.00	0.00	0.00	0.08	0.82	0.04	0.00	0.00	0.00
2.82	0.32	0.00	0.00	0.01	0.01	0.00	0.04	0.03	0.11	0.00	0.00	0.00
27.5	0.02	0.01	0.01	0.83	0.86	0.87	0.02	0.01	0.02	0.01	0.86	0.86
59.7	0.00	<u>0.99</u>	<u>0.99</u>	<u>0.15</u>	<u>0.12</u>	<u>0.12</u>	0.01	0.01	0.03	<u>0.99</u>	<u>0.13</u>	

The underlined proportions in the same row indicate the multicollinearity.

built in to the data. The other multicollinearity structure has been picked up by the last condition number.

---

## WHICH VARIABLES DO WE ELIMINATE?

With the exception of those variables having zero variance (which pick themselves), the decision about which variables to eliminate/include and the method by which this is done depends on several factors. The two most important factors are whether the dataset consists of two blocks of variables, a response block ( $Y$ ) and a descriptor/predictor block ( $X$ ), and whether the purpose of the analysis is to predict/describe values for one or more of the response variables from a model relating the variables in the two blocks. If this result is indeed the aim of the analysis, then it seems reasonable that the choice of variables to be included should depend, to some extent, on the response variable or variables being modeled. This approach is referred to as *supervised variable selection*. On the other hand, if the variable set consists of only one block of variables, the choice of variables in any analysis will be done with what are referred to as *unsupervised variable selection*.

---

## UNSUPERVISED ELIMINATION

The need to reduce the number of variables was recognized at about the same time that quantum mechanically derived descriptors proliferated in QSAR.<sup>15,42</sup> An early procedure, designed to eliminate the smallest number

of variables while breaking the largest number of colinearities, is known as CORCHOP.<sup>43</sup> This stepwise procedure can be summarized as follows:

- Identify pairs of variables with a very high (say greater than 0.99) correlation.
- Identify 1 variable from each pair to suggest for elimination.
- Having eliminated these very highly correlated variables, count the correlations above a limit (say 0.7) for each variable.
- List the variables with the largest number of correlations and suggest additional variables for elimination.
- Recalculate and relist the variables with the largest number of correlations.

The process is cyclical and is continued until the user has reduced the dataset to a suitable size or until no pairs of variables remain with a correlation greater than the user-defined limit. One of the features of this procedure is that the algorithm suggests variables to eliminate. The criterion for judging variables for elimination is how close their distribution is to normal based on skewness and kurtosis. The less normal of the pair is suggested for removal. The eventual use of the variables in any given dataset may not depend on normality; however, in the absence of any other criteria for elimination, we believe that normality of their distribution is a reasonable criterion for selection.

Similar procedures have been proposed by others,<sup>44,45</sup> and recently, a variable elimination technique that works in the opposite sense to CORCHOP has been described.<sup>46</sup> This latter method, known as unsupervised forward selection (UFS), constructs a dataset by selecting variables having low multicollinearity. UFS begins by eliminating variables that have a standard deviation below some assigned lower limit. The algorithm then computes a correlation matrix for the remaining set of variables and chooses the pair of variables with the lowest correlation. Correlations between the rest of the variables and these two chosen descriptors are examined, and any that exceed some preset limit are eliminated. Multiple correlations between each of the remaining variables and the two selected ones are examined, and the variable with the lowest multiple correlation is chosen. The next step is to examine multiple correlations between the remaining variables and the three selected variables and to select the descriptor with the lowest multiple correlation. This process continues until some predetermined multiple correlation coefficient limit is reached.

UFS, as presently implemented, requires no user intervention other than the choice of standard deviation and correlation coefficient limits. CORCHOP, on the other hand, was designed to be an interactive process that demands user decisions, and it can be tedious for large datasets. UFS is available from the website of the Centre for Molecular Design at the University of Portsmouth ([www.cmd.port.ac.uk](http://www.cmd.port.ac.uk)).

---

## SUPERVISED ELIMINATION

Supervised variable elimination might also be regarded as variable selection. Whether we consider this to be the third major section of how to treat multivariate datasets is a matter of semantics, however. It is possible to eliminate variables in a supervised manner rather than to select them. One obvious way is to eliminate variables that have a zero or very low correlation with the response variable or variables. In the case of classified response data, this selection means those descriptors that have the same distribution (mean and standard deviation) for the two or more classes. The danger in this selection process is the possibility that a variable might have a low correlation with the response but contribute to a multivariate correlation. Although this is possible, in practice, it is unlikely.

---

## VARIABLE SELECTION

The number of papers published in the literature relating to supervised variable selection methods is vast. Its obvious association with the ubiquitous publication of regression modeling means that trying to provide an exhaustive overview of the literature is an impossible task. Although this chapter aims to review variable selection methods, all that we can hope to achieve realistically is a broad coverage of the more obvious source methods and report here techniques that are sound from a statistical point of view and for which the software is available

Baumann, Albert, and von Korff<sup>47</sup> refer to a “correct” model as “a model containing . . . all predictors that generated the data,” i.e., give rise to the observed response values. These authors allow for the possibility that a correct model “. . . may contain some redundant terms or terms not related to . . .” the response variable, which leads them to define the “true” model as the “. . . smallest correct model.” Baumann et al. then list the following “three ingredients” that a variable selection method needs to find the smallest correct model:

1. A means of modeling the data
2. A search procedure for identifying potential subsets of variables
3. A criterion function for judging the usefulness of subsets of variables identified in ingredient 2.

Unfortunately, Baumann et al. do not provide the recipe, but they do suggest that if the aim of the analysis is to predict future values of a response, then the criterion function in 3 should be one related to the prediction error of the model and give, as an example of a means of modeling the data, a regression.

The most popular modeling method for regression is multiple linear regression (MLR). Here the response variable  $y$  is linked to a block of  $p$  independent variables  $x_1, x_2, \dots, x_p$  by a model of the form

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon \quad [13]$$

where  $\varepsilon$  is the usual random error term. In its simplest form  $p = 1$ , i.e., one independent variable. In this case, the method is referred to as simple linear regression. Most undergraduates in science have attended lectures on modeling data with regression methods in courses in statistics. Those courses often cover basic principles and usually go on to assess goodness of fit and discuss the construction of confidence intervals with the  $t$  statistic. For most people, that is where their formal education in regression ends. They will not have covered diagnostic checking of assumptions and, when enquiring where the sum of squares for regression,  $t$ -values, and the like come from, will probably have been told to not worry about such issues as the computer looks after that side of things. The “black box” approach to teaching and performing statistics is easy to accept and implement; provided that we work within the constraints set by the assumptions underpinning the theory on which the software is based, all is well. It therefore comes as no surprise to experts that some researchers in subjects like biology, chemistry, and pharmacy are unsure what they should or should not do in regression modeling when faced with a plethora of descriptor variables on which to base their models. Set this against a background where rumor and conjecture concerning the “do’s” and “don’ts” of regression analysis are rife, and we have a major problem.

It is our aim in this section of the chapter to report on regression methods that are used by researchers. Where possible comparisons will be made between the approaches, highlighting issues that exist concerning the application of the techniques and suggesting a pathway through what appears to some like a statistical mine field.

Broadly speaking, two main approaches to multiple linear regression modeling exist and they are classified according to whether the goal is to predict a single response variable or whether it is to predict several response variables simultaneously. We will consider here only the single response methods because this is found in most of the literature. Single response regression methods can also be placed into two subgroups according to whether the regressions are performed (1) directly on the variables in the descriptor set or (2) on linear combinations of the variables in the descriptor set.

In the first subgroup, we have ordinary least squares (OLS), best subset methods (BSS),<sup>48,49</sup> ridge regression (RR),<sup>50</sup> and nonnegative garrote (NNG)<sup>51</sup> methods. In the second subgroup, we have principal component regression (PCR),<sup>52</sup> partial least squares (PLS),<sup>53</sup> and continuum regression (CR)<sup>54–56;57</sup> methods.

## Ordinary Least Squares

OLS is the method that most novice modelers are familiar with. The model is as given in Eq. [13].  $N$  objects measured on  $p$  variables find estimates of the regression coefficients, i.e.,  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ , where the “hats” denote values calculated from the data. The technique that obtains these coefficients is generally the method of least squares (see, e.g., p. 218 of Dillon and Goldstein<sup>37</sup>), and the solution can be written succinctly in matrix form as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \quad [14]$$

where the  $\mathbf{X}$  is an  $n \times (p + 1)$  matrix consisting of  $n$  values of  $p$  predictor variables, and a column of 1's if the constant term ( $\beta_0$ ) is being fitted.  $\mathbf{Y}$  is the corresponding  $n \times 1$  vector of observed response values, and  $\hat{\boldsymbol{\beta}}$  is the  $(p + 1) \times 1$  vector of regression coefficients estimates. Unless stated otherwise, we will assume that the rank of the matrix  $\mathbf{X}'\mathbf{X}$ ,  $r$ , equals  $p$ , which requires  $n > p$ . Notice in Eq. [14] the presence of the  $(\mathbf{X}'\mathbf{X})^{-1}$  matrix, which we earlier pointed out is the source of problems when multicollinearity is severe. These estimates are unbiased, which means that if you were able to repeat the analysis many times on datasets randomly chosen from the same population as the original sample, the average values of these coefficients would then be extremely close to the theoretical population values. Statisticians write this simply as  $E(\hat{\beta}_i) = \beta_i$  for  $i = 1, 2, \dots, p$ , where  $E(\hat{\beta}_i)$  is expressed as the “expectation of  $\hat{\beta}_i$ ” and means the “average value of  $\hat{\beta}_i$ .” However, as we have discussed at length, multicollinearity results in a whole raft of undesirable properties, not the least of which is lack of precision in its predictions. Before we move on to the other regression methods, it is important at the outset that we are all aware of the assumptions underpinning the theory giving rise to the statistical properties of estimated regression models. The assumptions are as follows:

1. The error terms in Eq. [13] are independently and identically distributed with a normal distribution having mean of zero and a variance of 1. Statisticians write this simply as  $\varepsilon_i$  are i.i.d.  $N(0,1)$ . This assumption allows, for example, the usual  $t$ -tests to be performed.
2. The  $p$  descriptor variables were chosen without recourse to using the information in the response variable; i.e., they were chosen in an unsupervised manner. This point is crucial as it allows us to use the usual  $F$ -test and corresponding  $p$ -values to assess the significance of the regression. There will be much more on this later in the chapter.
3. Measurement errors in the descriptor variables are much smaller than in the response variables.

## Ridge Regression

RR was developed by Hoerl and Kennard.<sup>50</sup> It is a technique that was specifically designed to overcome the ill-conditioning of the matrix  $\mathbf{X}'\mathbf{X}$

resulting in the undesirable properties of OLS regression coefficients. Hoerl and Kennard suggested adding a constant to the diagonal elements of the  $\mathbf{X}'\mathbf{X}$  matrix so that the estimates of the regression coefficients are calculated from

$$\hat{\boldsymbol{\beta}}_{RR} = (\mathbf{X}'\mathbf{X} + d\mathbf{I})^{-1}\mathbf{X}'\mathbf{Y} \quad [15]$$

In Eq. [15],  $d > 0$  is the ridge parameter and  $\mathbf{I}$  is the identity matrix, i.e., a matrix with 1's on the leading diagonal and zeros elsewhere. Hoerl and Kennard recommend standardizing the data to zero mean and unit variance before performing the ridge estimate calculations given by Eq. [15]. RR is, with  $d > 0$ , a biased regression method in that  $E(\hat{\beta}_{iRR}) \neq \beta_i, i = 1, 2, \dots, p$ , where  $\hat{\beta}_{iRR}$  is the RR estimate of the  $i$ th regression coefficient. Because

$$\sum \hat{\beta}_{iRR}^2 < \sum \hat{\beta}_{iOLS}^2 \quad [16]$$

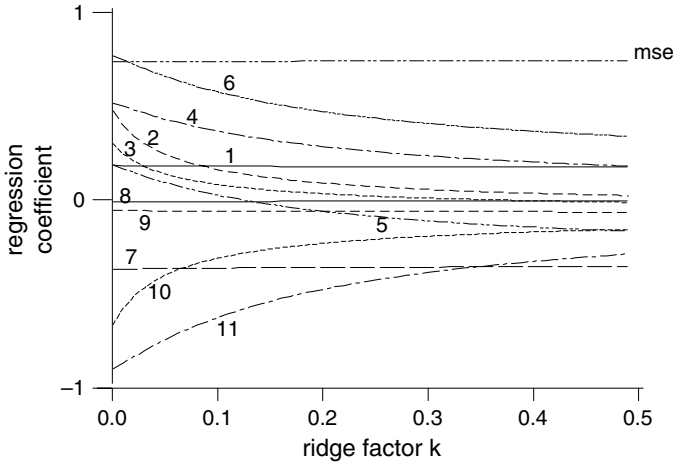
RR is also referred to as a shrinkage regression technique (see, for example, Frank and Friedman<sup>58</sup>). We will come back to this point when we have considered PCR, PLS, and CR, which are also shrinkage methods.<sup>58</sup>

Much of the literature devoted to RR centers on the mean squared error of the estimates  $\hat{\boldsymbol{\beta}}$ ,  $MSE(\hat{\boldsymbol{\beta}})$ , which may be written as

$$MSE(\hat{\boldsymbol{\beta}}) = E[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})] = \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i} \quad [17]$$

$\boldsymbol{\beta}$  is the  $p \times 1$  column vector of true or population regression coefficients, and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  are the eigenvalues of  $\mathbf{X}'\mathbf{X}$ . When multicollinearity is present to a high degree, some of the  $\lambda_i$  will approach zero, resulting in an extremely large  $MSE(\hat{\boldsymbol{\beta}})$ .  $\hat{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta}$  can be considered as two points in  $p$ -dimensional space, and the  $MSE(\hat{\boldsymbol{\beta}})$  measures the average Euclidean distance between them. Consequently, multicollinearity results in poor least squares estimates and from Eqs. [8] and [17], it can be seen that the ratio  $100 \times \sum_{i=1}^p \frac{1}{\lambda_i} / \sum_{i=1}^p \lambda_i$  gives the percentage increase in the  $MSE(\hat{\boldsymbol{\beta}})$  over what it would be if the variables had been orthogonal. Of more interest, however, is the mean square error of the predicted values,  $MSE(\mathbf{X}\hat{\boldsymbol{\beta}})$ . Lawless and Wang<sup>59</sup> have shown that although RR can result in a dramatic reduction in  $MSE(\hat{\boldsymbol{\beta}})$  for multicollinear datasets, it provides less reduction in  $MSE(\mathbf{X}\hat{\boldsymbol{\beta}})$  and occasionally produces a small increase.

Unfortunately, before RR can be applied, a value of the ridge parameter  $d$  must be found. One method that is relatively simple to apply is to perform RR for a range of values of  $d$  in the interval  $(0, 1)$  and to then select as the



**Figure 7** Ridge trace for the eleven regression coefficients for the data set generated for use in the multicollinearity detection example.

“best” value the one which the regression coefficients stabilize. This technique for selecting the value of  $d$  is clearly subjective and is not without its critics.<sup>60</sup> Figure 7 shows what are referred to as “ridge traces” for models based on the data generated earlier for the section on identification of multicollinearity. These ridge traces are the values of  $\hat{\beta}_{RR}$  evaluated with Eq. [15] plotted against the corresponding value of  $d$ . The largest eigenvalue of the correlation matrix for this data is 205.396, and the smallest is 0.058 giving a condition index (ratio of largest to smallest eigenvalue) that is extremely large. This dataset is extremely multicollinear, and we see that as  $d$  increases from zero, some of the beta coefficients change by shrinking away from the OLS values with two of them (3 and 5 changing sign). From these traces, for the most part, the coefficients that change most are those associated with multicollinearity built into their generation. Also plotted in Figure 7 is the residual mean square error. It rises only gently as expected, thus supporting that the OLS solution is the one that gives the minimum value.

Because of the subjectivity attached to the ridge trace approach for finding the “best” value for  $d$ , many automatic alternative methods have been offered in the literature. The two that appear more frequently are given in Eqs. [18] and [19].

$$d_{HKB} = \frac{ps^2}{\sum_{i=1}^p \hat{\beta}_i^2} \quad [18]$$

$$d_{LW} = \frac{ps^2}{R^2} \quad [19]$$

where  $d_{HKB}$  was suggested by Hoerl et al.<sup>61</sup> and  $d_{LW}$  by Lawless and Wang.<sup>59</sup> The denominator in Eq. [18] is the sum of squared regression, coefficient estimates for the OLS regression, and the denominator in Eq. [19] is the usual squared multiple regression correlation coefficient again for OLS. A detailed account of ridge regression is given by Draper and Van Nostrand.<sup>62</sup>

## Principal Component Regression, Partial Least Squares, and Continuum Regression

Three methods, principal components, partial least squares, and continuum regression form a convenient group for the purposes of categorization. They, like RR, produce biased estimates of the regression coefficients by shrinking the coefficient vector away from the OLS solution, and all three methods are latent variable techniques. Latent variable regression analysis is essentially a two-stage process involving the following:

1. Construction of orthogonal linear combinations of descriptor variables
2. Least squares regression of the response variable on the latent variables found in stage 1

We have already introduced orthogonal linear combinations of the descriptor variables in the sections on PCA and multicollinearity. As all three methods in this section begin with forming linear combinations (latent variables) of the X-set of variables, we will adopt the common notation in Eq. [20] to represent the  $i$ th such construction:

$$LV_i = c_{1i}\mathbf{x}_1 + c_{2i}\mathbf{x}_2 + \dots + c_{pi}\mathbf{x}_p \quad [20]$$

In this equation, where the  $\mathbf{x}$ s are assumed to have been standardized to zero mean and unit variance, i.e.,  $(x_i - \bar{x}_i)/s_i$ , where  $\bar{x}_i$  and  $s_i$  are, respectively, the mean and standard deviation of the  $i$ th descriptor variable. The only difference between the methods to be described is the way that the coefficient vectors  $\mathbf{c}'_i = (c_{1i}, c_{2i}, \dots, c_{pi})$ ,  $i = 1, 2, \dots, p$  are calculated.

### *Principal Component Regression*

In PCR, the coefficient vectors are found by performing PCA on the correlation matrix of the descriptor variables. PCA finds  $\mathbf{c}'_i$ , which maximizes the variance of  $LV_i$  i.e.,  $\text{var}(\mathbf{c}'_i\mathbf{x})$ , subject to the constraints that  $\sum_{j=1}^p c_{ji}^2 = 1$  and  $\sum_{j=1}^p c_{ji}c_{jk} = 0$ ,  $i \neq k$ . The first constraint is that the squared length of the coefficient vector equals unity. The second constraint is that the coefficient vectors are mutually orthogonal (zero correlation between any two latent variables). Statisticians sometimes write this in a type of shorthand as

$$\begin{aligned} \mathbf{c}_{iPCR} &= \arg \max \{ \text{var}(\mathbf{c}'\mathbf{x}) \} \\ \mathbf{c}'_i\mathbf{c}_k &= 0, i \neq k \\ \mathbf{c}'\mathbf{c} &= \mathbf{I} \end{aligned} \quad [21]$$



PCA constructs its latent variables (the principal components) in the order of decreasing eigenvalues, i.e.,  $\lambda_1 > \lambda_2 > \dots > \lambda_r$ , where  $r$  is the rank of the correlation matrix of the descriptor variables. The eigenvalues are the variances of the corresponding components. Because their construction is independent of the response variable, the degree of association between the response and the latent variables may not follow the same ranking as the variances, however. The second stage of the PCR process will be explained after introducing the remaining two techniques in this section.

### *Partial Least Squares*

The latent variables in a PLS analysis are constructed to maximize the squared covariance between  $Y$  and the  $\mathbf{c}'\mathbf{x}$  subject to the same two constraints as in PCA; i.e., the squared length of the coefficient vector equals unity and each of the coefficient vectors are mutually orthogonal, which can be expressed mathematically as

$$\begin{aligned} \mathbf{c}_{iPLS} &= \arg \max \{ \text{cov}(\mathbf{y}, \mathbf{c}'\mathbf{x}) \}^2 \\ \mathbf{c}'_i \mathbf{c}_k &= 0, i \neq k \\ \mathbf{c}'\mathbf{c} &= 1 \end{aligned} \quad [22]$$

where, in common with the  $\mathbf{x}$ s, the  $\mathbf{y}$  has also been standardized. Unlike PCR, because of the dependence of latent variables on the response, PLS regression will apply the latent variables in the order of their construction.

### *Continuum Regression*

Stone and Brooks<sup>54</sup> made an important step toward understanding all three of these regression methods when they introduced a method they referred to as continuum regression. These authors recognized that the coefficients of the linear combinations of descriptor variables for OLS (the  $\beta_i$ ), PCR, and PLS (the  $c_i$ ) can be found by maximizing a general criterion function with a single parameter  $\alpha$  set at three different values. With  $\alpha = 0, 0.5, 1$ , CR performs OLS, PLS, and PCR, respectively. This process is illuminating as it shows the interrelationship among these three popular methods. Before we present Stone and Brooks' formulation of CR, let us first revisit OLS and recast it in the notation introduced for PCR and PLS. Like PCR and PLS, OLS produces a linear combination of the descriptor variables (just one), which we can write as

$$\hat{\mathbf{y}} = c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_p \mathbf{x}_p = \mathbf{c}'\mathbf{x} \quad [23]$$

where the  $p \times 1$  column vector of  $\beta$ 's has now been replaced by the  $c$ 's. All that remains to bring it in line with PLS and PCR is to recognize that the  $c$ 's in the

estimated regression model (Eq. [23]) are those that give the maximum squared multiple regression coefficient  $R^2$ , i.e.,

$$\mathbf{c}_{OLS} = \arg \max_{\mathbf{c}'\mathbf{c} = 1} \{ \text{corr}(\mathbf{y}, \mathbf{c}'\mathbf{x}) \}^2 \quad [24]$$

Rewriting Eq. [24] in terms of covariance and variances gives

$$\mathbf{c}_{OLS} = \arg \max_{\mathbf{c}'\mathbf{c} = 1} \frac{[\text{cov}(\mathbf{y}, \mathbf{c}'\mathbf{x})]^2}{\text{var}(\mathbf{y})\text{var}(\mathbf{c}'\mathbf{x})} = \arg \max_{\mathbf{c}'\mathbf{c} = 1} \frac{[\text{cov}(\mathbf{y}, \mathbf{c}'\mathbf{x})]^2}{\text{var}(\mathbf{c}'\mathbf{x})} \quad [25]$$

remembering that  $\text{var}(\mathbf{y}) = 1$ . Stone and Brooks<sup>54</sup> recognized OLS, PLS, and PCR represented just three possible regressions, of a continuum of regressions, indexed by a single parameter  $\alpha$ , on the interval  $[0, 1]$ , in a generalized criterion function  $T$  given by

$$T = \{ \text{cov}(\mathbf{y}, \mathbf{c}'\mathbf{x}) \}^2 \{ \text{var}(\mathbf{c}'\mathbf{x}) \}^{\frac{\alpha}{1-\alpha}-1} \quad [26]$$

Substituting into Eq. [26],  $\alpha = 0$  gives Eq. [25], i.e., OLS and  $\alpha = 0.5$  gives Eq. [22], i.e., PLS. By letting  $\alpha$  tend toward unity, the power on the variance term in Eq. [26] approaches infinity, which makes the contribution of the squared covariance term negligible. The net effect is that when  $T$  is maximized with  $\alpha$  close to unity, the function behaves as though the variance component alone is present; i.e., the  $\mathbf{y}$  variable makes no “detectable” contribution. So, as  $\alpha$  increases from 0, the solution vectors ( $\mathbf{c}'$ ) are being “pulled” away from the direction of  $\mathbf{y}$  until eventually, at values close to unity, the only thing that matters is producing vectors that maximize the variance in the space spanned by the descriptor variables. For such values of  $\alpha$ , the criterion function given by Eq. [26] behaves as though it was Eq. [21], i.e., PCR.

Before leaving stage 1 of the two-stage regression with latent variables, we note that ridge regression can be formulated in a similar fashion to that for OLS, and the interested reader is referred to Frank and Friedman’s (re)view of regression tools in chemometrics.<sup>58</sup>

Stage 2 of LV regression moves us to one of the most contentious topics in regression: how to decide on which and how many LVs to use when building a linear regression model. Various methods are available for model construction, and they will be considered in more detail in the section on best subset selection, but generally, some form of sequential procedure is adopted. One of these, “forward stepping” as an example, builds the model equation by considering the latent variables in turn and including them in the model, provided that they pass some inclusion criteria. The model is considered

completed when variables omitted from the model cannot pass the inclusion criteria. The final model will then take the general form

$$\hat{\dot{y}} = \hat{\alpha}_1 L V_1 + \hat{\alpha}_2 L V_2 + \cdots + \hat{\alpha}_p L V_p = \sum_{i=1}^p \hat{\alpha}_i L V_i \quad [27]$$

where the dot over the  $y$  denotes that it is standardized to zero mean and unit variance, and the  $\hat{\alpha}_i$  are the OLS regression coefficient estimates found by regressing the  $\hat{\dot{y}}$  on the  $L V_i$ . Replacing the  $L V_i$  in Eq. [27] with the individual linear combinations of the  $\mathbf{x}$ 's they represent, i.e.,  $\hat{c}_{1i}\mathbf{x}_1 + \hat{c}_{2i}\mathbf{x}_2 + \cdots + \hat{c}_{pi}\mathbf{x}_p$ , where the  $\hat{c}_{ji}$  have been calculated according to whichever of these three methods is being used, we have

$$\begin{aligned} \hat{\dot{y}} = & \hat{\alpha}_1(\hat{c}_{11}\mathbf{x}_1 + \hat{c}_{21}\mathbf{x}_2 + \cdots + \hat{c}_{p1}\mathbf{x}_p) + \hat{\alpha}_2(\hat{c}_{12}\mathbf{x}_1 + \hat{c}_{22}\mathbf{x}_2 + \cdots + \hat{c}_{p2}\mathbf{x}_p) + \cdots \\ & \cdots + \hat{\alpha}_p(\hat{c}_{1p}\mathbf{x}_1 + \hat{c}_{2p}\mathbf{x}_2 + \cdots + \hat{c}_{pp}\mathbf{x}_p) \end{aligned}$$

This calculation can be simplified by removing the brackets to give

$$\begin{aligned} \hat{\dot{y}} = & (\hat{\alpha}_1\hat{c}_{11} + \hat{\alpha}_2\hat{c}_{12} + \cdots + \hat{\alpha}_p\hat{c}_{1p})\mathbf{x}_1 + (\hat{\alpha}_1\hat{c}_{21} + \hat{\alpha}_2\hat{c}_{22} + \cdots + \hat{\alpha}_p\hat{c}_{2p})\mathbf{x}_2 + \cdots \\ & \cdots + (\hat{\alpha}_1\hat{c}_{p1} + \hat{\alpha}_2\hat{c}_{p2} + \cdots + \hat{\alpha}_p\hat{c}_{pp})\mathbf{x}_p \\ = & \hat{\beta}'_1\dot{\mathbf{x}}_1 + \hat{\beta}'_2\dot{\mathbf{x}}_2 + \cdots + \hat{\beta}'_p\dot{\mathbf{x}}_p \end{aligned} \quad [28]$$

where the  $\hat{\beta}'_{iS} = \hat{\alpha}_1\hat{c}_{i1} + \hat{\alpha}_2\hat{c}_{i2} + \cdots + \hat{\alpha}_p\hat{c}_{ip}$ ,  $i = 1, 2, \dots, p$ , are standardized regression coefficients and the subscript  $S = \text{PCR, PLS, or CR}$ , depending on the techniques used. If the standardized variables  $\hat{\dot{y}} = (\hat{y} - \bar{y})/s_y$ ,  $\dot{\mathbf{x}}_i = (\mathbf{x}_i - \bar{\mathbf{x}}_i)/s_i$   $i = 1, 2, \dots, p$  are now replaced in Eq. [28], we have after further simplification

$$\hat{y} = \bar{y} + \sum_{i=1}^p \frac{s_y}{s_i} \hat{\beta}'_{iS} (\mathbf{x}_i - \bar{\mathbf{x}}_i) = \bar{y} - \sum_{i=1}^p \frac{s_y}{s_i} \hat{\beta}'_{iS} \bar{\mathbf{x}}_i + \sum_{i=1}^p \frac{s_y}{s_i} \hat{\beta}'_{iS} \mathbf{x}_i$$

This equation in turn can be further reduced if we replace the constant (intercept term)  $\bar{y} - \sum_{i=1}^p \frac{s_y}{s_i} \hat{\beta}'_{iS} \bar{\mathbf{x}}_i$  by  $\hat{\beta}_0$  and  $\frac{s_y}{s_i} \hat{\beta}'_{iS}$  by the unstandardized regression coefficients  $\hat{\beta}_{iS}$  to give the more usual form of the estimated regression model, i.e.,

$$\hat{y} = \hat{\beta}_{0S} + \sum_{i=1}^p \hat{\beta}_{iS} \mathbf{x}_i \quad [29]$$

Although Eq. [29] has the appearance of a multiple regression, remember that the parameter estimates were not calculated by OLS. Instead they were found by a biased regression method. Consequently, these parameters, which are referred to as pseudo- $\beta$ 's, will not in general equal the OLS values because they have been “shrunk,” (some more than others). However, as more components are added into the latent variable model (Eq. [27]), i.e., as  $p$  increases toward  $k$ , these pseudo- $\beta$ 's approach the values obtained by OLS. In the limit  $p = k$ , Eq. [29] will be identical to the OLS model, a result that will be illustrated later when we apply the methods to a real dataset.

So far in this section we have worked on the premise that the sample size  $n$  is larger than the number of variables in the model  $p$ . The reason for focusing on systems with  $n > p$  is because of the methods being described to extract the eigenvalues and eigenvectors. We have assumed that when the eigenvalues and vectors are being extracted they are done so all together; i.e., with  $p$  variables, we extract  $p$  components or latent variables, which is only possible if  $n > p$  so that the rank of the descriptor matrix  $r = p$ . However, algorithms exist that will extract the eigenvalues and eigenvectors one at a time and, in so doing, allow latent variables to be constructed from datasets containing many more variables than objects. One such algorithm is the NIPALS algorithm,<sup>63</sup> which serves as the basis for some applications of PLS.<sup>58</sup> We need to be convinced that regression models consisting of LVs comprising linear combinations of many descriptor variables perform any differently than ones with relatively few. We have more to say on this point later.

## Best Variable Subset Selection

At the outset, we must make it clear that despite the plethora of papers on multiple linear regressions, spanning an enormous variety of topics including drug design,<sup>31</sup> epidemiology,<sup>64</sup> psychology,<sup>65</sup> market research,<sup>66</sup> and economics,<sup>67</sup> the problems associated with variable selection we are about to discuss are little known outside the statistics fraternity. Therefore, before exploring the variable subset selection (VSS) procedures, we turn our attention briefly to several of the main issues raised by statisticians relating to the selection of a few variables from a larger pool of potential regressor variables.

Miller<sup>48</sup> identifies two forms of bias associated with subset selection: (1) omission bias and (2) competition bias, both of which will be discussed here.

### *Omission Bias*

To describe omission bias, we assume that the “true” model linking  $Y$  to the descriptors is that which consists of all  $k$  descriptor variables. To begin with, we will assume that the number of observations  $n$  exceeds the number of available descriptor variables  $k$ , so that if we wanted to include all variables

in our regression model, we would have enough degrees of freedom to do so. Let this “true” model be denoted as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + \varepsilon \quad [30]$$

which is just Eq. [13] with  $p$  replaced by  $k$  so that the least squares estimates of the unknown coefficients are given by Eq. [14]. One of the main reasons for performing regression analysis is to predict the response for a new object given a vector  $\mathbf{x}'_0 = (1, x_{01}, x_{02}, \dots, x_{0k})$  of new descriptor values. The predicted value of  $y$ ,  $\hat{y}_0$ , is obtained from Eq. [30] by replacing the estimated coefficient values to give

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \hat{\beta}_2 x_{02} + \cdots + \hat{\beta}_k x_{0k} = \mathbf{x}'_0 \hat{\boldsymbol{\beta}}$$

and it can be shown (p. 228 of Dillon and Goldstein<sup>37</sup>) that the variance is

$$\text{var}(\hat{y}_0) = \text{var}(\mathbf{x}_0 \boldsymbol{\beta}') = [\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0] \sigma^2 \quad [31]$$

Let us now consider what happens to the variance of the predicted response given by Eq. [31] when we construct a model with fewer descriptor variables  $p < k$ . It is assumed here as elsewhere in this section that the  $p$  variables are chosen independently of the  $Y$  variable, i.e., they were chosen in an unsupervised manner. We now partition our  $\mathbf{X}$  matrix into two submatrices  $\mathbf{X}_A$  and  $\mathbf{X}_B$  so that  $\mathbf{X}_A$  consists of the first  $p$  variables of  $\mathbf{X}$  plus the column of ones, and  $\mathbf{X}_B$  comprises the remaining  $k - p$  variables. Using Eq. [31], the variance of the prediction is now given by

$$\text{var}(\hat{y}_0) = \text{var}(\mathbf{x}_{0A} \boldsymbol{\beta}'_A) = [\mathbf{x}'_{0A} (\mathbf{X}'_A \mathbf{X}_A)^{-1} \mathbf{x}_{0A}] \sigma^2$$

Miller<sup>49</sup> (p. 6), using an argument based on the Cholesky factorization of  $\mathbf{X}'\mathbf{X}$ , shows that

$$\text{var}(\mathbf{x}_{0A} \boldsymbol{\beta}'_A) \leq \text{var}(\mathbf{x}_0 \boldsymbol{\beta}')$$

This result tells us that the variance of the predicted estimates reduces monotonically as the number of descriptors is decreased. The converse is also true; i.e., increasing the number of descriptors produces a monotonic increase in the variance of the prediction. This interesting result is true for linear models fitted by least squares, which at first seems to be a good argument in favor of reducing the number of descriptors, but a price will be paid. If Eq. [30] is the true model, reducing the number of variables results in biased regression coefficient estimates for the remaining variables in the model, which is what is meant by omission bias. Conversely, if we begin with fewer variables  $p(< k)$  in the

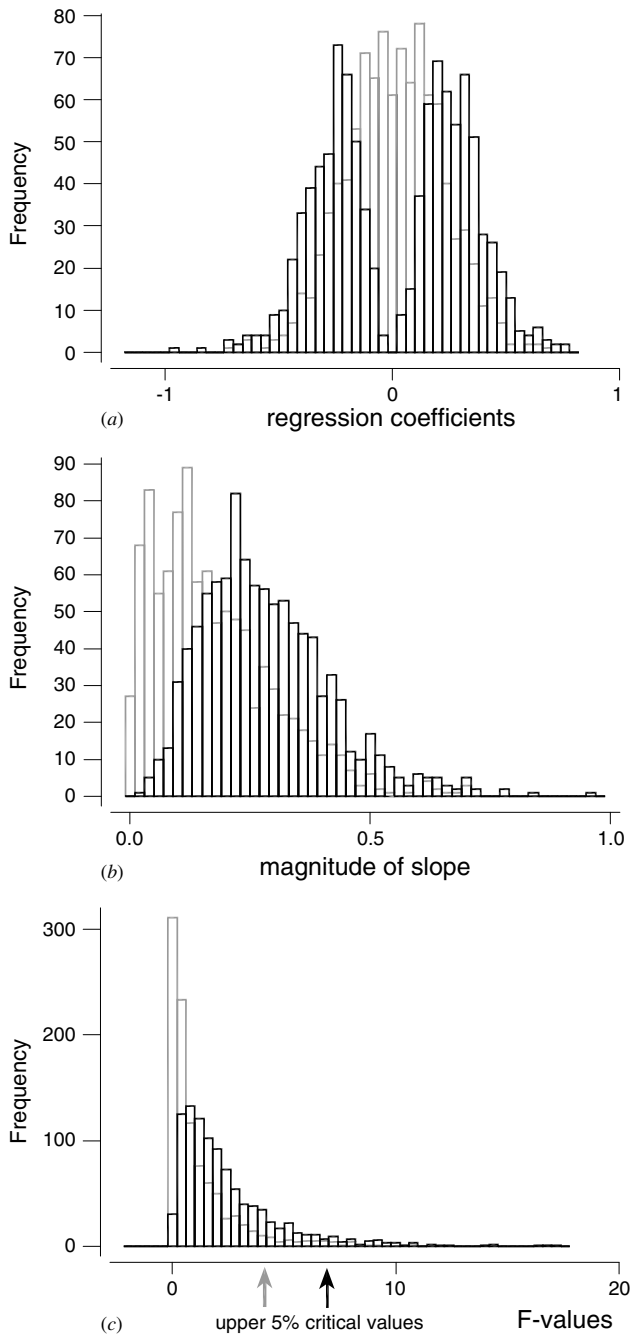
model and then add more variables, we will reduce the bias at the expense of increasing the variance of the prediction. A question that quickly comes to mind at this point is: where do we “draw” the line in this tradeoff between bias and variance? We will consider this problem in detail later on in this section, but for the moment we point out that including variables with no predictive ability just increases the variance. Some novice computational chemists reading this section on omission bias might be wondering what the point of this discussion is because to discuss it we have assumed we know the true relationship, whereas in practice we do not. Of course you are correct; the point being made here is that inclusion of unnecessary variables does little for prediction other than to increase the variance associated with it. Before leaving this discussion of omission bias, we point out the work of Thompson,<sup>68</sup> who distinguishes between omission of control variables as found in experimental design situations, and omission of observational variables over which the user has no control. In the latter situation, the biases resulting from the inclusion/exclusion of variables can be considered as part of the residual variance.

### *Competition Bias*

In all that we have said so far the choice of the set of  $p$  variables to be included in the regression has been made independently of the response variable. If this is not the case, then another type of bias is introduced. This type of bias is referred to as competition or selection bias. It occurs when the variables to be included in the model are chosen in a supervised manner to maximize some function involving the response variable. One such function is  $R^2$ , the proportion of the total sum of squares of  $Y$  explained by the regression. To understand this form of bias, consider a simple situation in which several regressions were constructed that contain only one descriptor variable each. The data and regressions were constructed as follows:

1. Four sets of 25 standard normal (i.e., zero mean and unit variance) random numbers were generated to provide data for the four variables ( $y$ ,  $x_1$ ,  $x_2$ ,  $x_3$ ).
2. One of the three  $x$ 's in step 1 was picked at random, and the  $y$  variable (considered as the response) was regressed on it and the slope and  $R^2$  values were recorded.
3. Step 2 was repeated, but instead of picking an  $x$  at random, first  $y$  was regressed on all three  $x$ 's individually and the one giving the **largest**  $R^2$  was selected.
4. This procedure was repeated 1000 times; i.e.,  $1000 \times 4 \times 25$  random numbers were produced.

The results from this simulation experiment are summarized in Figure 8 where the gray-lined histograms correspond to the randomly chosen  $x$  models and the black-lined to those obtained by selecting the  $x$  that gave the largest



**Figure 8** Histograms of regression results. Gray lines refer to randomly chosen  $x$  and the black lines for the  $x$  chosen to maximize squared correlation. See text for details.

$R^2$ . In Figure 8(a), we see the histograms for the regression coefficients where, as expected, the distribution corresponding to the randomly chosen  $x$  is fairly normal in shape and centered on zero ( $x$  and  $y$  are random and hence  $\beta_1 = 0$ ) (gray lines). In contrast is the distribution of the regression coefficient corresponding to the  $x$  selected to maximize the fit, i.e.,  $R^2$ . This distribution comes as a surprise to most people. If they had been asked beforehand the question “what difference do you think there will be between the two diagrams,” they might have responded “. . .but surely they will be about the same, won’t they?” Well clearly they are not and what we are seeing here is the result of **selection bias**. Selection bias results in regression coefficients that are generally larger in magnitude than from a purely random selection, as characterized here by the bimodal nature of the distribution. If we redraw this figure by plotting the absolute value of the observed slopes [Figure 8(b)], the bias to larger values is even more apparent.

Selection bias does not stop with inflated regression coefficients; it also has a marked effect on the significance level, which is usually quoted as a “ $p$ -value” in computer output. To see the effect selection bias has on the significance levels, we turn our attention to Figure 8(c), where the observed distribution of  $F$ -ratios for the randomly chosen  $x$  models (gray lines) and those corresponding to the maximum  $R^2$  models (black lines) are superimposed. The gray arrow marked on Figure 8(c) is the value of  $F$  (4.35), which indicates the position of the 95th percentile from the 1000 simulations. Because we are fitting  $y$  to one  $x$  with a sample size of 25, this should correspond with the upper 5% point of the  $F$ -distribution on 1 and  $25 - 1 - 1 = 23$  degrees of freedom, i.e., 4.30. The two values are in good agreement as they should be if the null hypothesis  $H_0 : \beta_1 = 0$  is true, i.e., the slope of the true regression line is zero. Because the  $y$  and  $x$ ’s have been generated at random, this is the true state of affairs. If we now look at the  $F$ -distribution from the regressions chosen to maximize the fit, and use the same critical value of 4.30, we discover that over 15% of the  $F$ -values exceed this point. Said another way, if we use the same critical value as in the case for the randomly chosen  $x$ , we have a significance level of 15% and **NOT** 5% as we may have otherwise thought. To assess the significance of the model selected when maximizing  $R^2$ , we would need to use a critical value of approximately 7.10, which is the value of  $F$  corresponding to the 95th percentile of the  $F_{max}$  values. Suffice to say that this bias gets worse the more variables there are to select from to maximize our model fit.

At one time, scientists selected variables for inclusion in a regression model based on some type of strong justification for including them. That justification would have been based on some *a priori* idea or hypothesis that the variables chosen were *the* critical factors in establishing the variability in the response variable. Although this approach may still be used by scientists, more often today researchers have no prior knowledge about the relative importance of the various independent variables. Consequently, many regression



models are fabricated by selecting variables whose observed values make “significant” contributions to explaining the variation on the observed values of the response variable; i.e., they are based on the current values in the dataset. As we have seen, such selection procedures lead to inflated regression statistics, but despite this detrimental problem, these methods are still in popular use by researchers. We finish this part of the chapter by looking at these popular yet pernicious techniques together with a look at what has been achieved in the way of overcoming some of the problems associated with variable selection.

### Forward Inclusion

In the forward inclusion procedure, the response variable is regressed on each of the  $X$  variables and the  $x$ -variable with the greatest  $F$ -to-enter,  $F_{enter}$ , is selected as a candidate for inclusion in the model. This variable will enter the model provided its  $F$ -to-enter exceeds some preselected critical value set by the user. If we assume that  $p$  variables currently exist in our model, the  $F$ -to-enter is calculated as

$$F_{enter} = \frac{ESS_p - ESS_{p+1}}{ESS_{p+1}/n - p - 1} \quad [32]$$

where  $ESS_p$  and  $ESS_{p+1}$  are the error sums of squares of the model before and after the  $(p + 1)$ th variable is included. If the maximum  $F_{enter}$  given by Eq. [32] for the variables not already in the model exceeds the  $F_{enter}$  set by the user, the variable is then included. This process continues until no variables remain or none of the remaining variables have an  $F_{enter}$  exceeding the preset value. Once a variable has entered the model, it is not removed. In most of the older statistical software packages,<sup>69</sup> the default value of this  $F_{enter}$  is 4 irrespective of the number of variables already in the model. This critical  $F$ -value is based on the upper 5% point of the  $F$ -distribution on 1 and  $(n - p - 1)$  degrees of freedom, which for  $n - p - 1$  of around about 20 or higher is approximately 4. As an alternative to using  $F_{enter}$  values, we could use the corresponding  $p$ -values with the selection process stopping when the  $p$ -value exceeded some specified level. Bendel and Afifi<sup>70</sup> have investigated various  $F_{enter}$  values for use in forward selection and suggest a value corresponding to an upper-tail  $p$ -value of 15% if the sample size is large. This  $p$ -value corresponds to an  $F_{enter}$  of approximately 2.1 for an  $F$ -distribution on 1 and infinite degrees of freedom. As can be seen, Bendel and Afifi’s minimum  $F_{enter}$  is more liberal than the above default value of 4, but it must be remembered that these values are merely notional values intended to restrict the entry of variables into the model.

### Backward Elimination

Backward elimination works along similar principals to forward inclusion, except that we start with all variables in the model and then eliminate

them one at a time. The variable eliminated at each step is the one with the smallest  $F$ -to-remove,  $F_{\text{remove}}$ , provided it is smaller than some user set value. An equation similar to Eq. [32] calculates the new  $F$ -to-remove:

$$F_{\text{remove}} = \frac{ESS_{p-1} - ESS_p}{ESS_p/n - p - 1} \quad [33]$$

where the notation is the same as that in Eq. [32]. Clearly,  $p < n$  for this procedure to be used.

### Stepwise Regression

Stepwise regression proposed by Efroymson,<sup>71</sup> is a combination of forward inclusion and backward elimination. After each variable is added (other than the first two), a test is performed to see if any of the variables entered at an earlier step can be deleted. The procedure applies both Eqs. [32] and [33] in a sequential manner. The stepping stops when no more variables satisfy either the criterion for removal or the criterion for inclusion. To prevent the procedure from unnecessarily cycling the critical values of  $F$ -to-enter and  $F$ -to-remove should be such that  $F_{\text{remove}} < F_{\text{enter}}$ .

All three of these methods that we use for variable selection are prone to entrapment in local minima, i.e., they find a combination of variables that cannot be improved on in the next step (removal or addition of one variable) for the criterion function, which can be avoided by performing either a Tabu search (TS) or the more computationally expensive “all subsets regression.” We discuss the second of these two methods in the next section and refer readers to the papers by Glover<sup>72,73</sup> for details of the TS method.

Earlier we introduced the issue of selection bias when regression models are constructed to optimize some fit criteria. The problem of selection bias for multiple descriptor models is a major problem in the scientific literature. Papers containing quoted significance values, whether in the form of statements like “... and this result is significant at the 5% level ...” or quoting a  $p$ -value of say 0.03 and remarking “... as the  $p$ -value  $0.03 < 0.05$ , the result is significant at the 5% level ...” are **incorrect**. Many regression models have severely overoptimistic statistics attached to them that can only lead to the misinterpretation of structure activity models. Solutions to this problem exist, solutions that have been known almost as long as the problem and we can only speculate why chemists, biologists, and other scientists have not picked up on them. One possible explanation is that many journals still do not have specialist statistical referees, and so overinflated significance values go unchallenged. Also, journals have a tendency to look less favorably on statistical analyses, which are not “significant” leading to another form of bias (“publication bias”). Consequently, authors are never alerted to their misuse of  $F$ -tables, and younger scientists looking to such literature for methodology pick up

and perpetuate the problem. This misuse may also result from the terminology we use in statistical packages such as *F*-to-enter and *F*-to-remove. Neither of these quantities actually has a distribution anything like an *F*-distribution and are merely numbers that the user can set to limit the number of variables selected.

Methods of overcoming inflated significance values developing from ordinary *F* tables has received much attention in the literature, albeit in more mainline statistical journals.<sup>48,74–76</sup> However, papers to be found in other areas also consider the problem.<sup>77,78</sup> Although these latter publications differ in fine detail, the main theme is use of Monte Carlo simulations to generate the empirical distribution of some fit statistic such as the coefficient of determination,  $R^2$ . These simulations are done for various combinations of  $k$  (number of variables available to select from),  $p$  (number of variables in final model),  $n$  (sample size), and method of model construction. In the earlier papers, the number of simulations was relatively small, being of the order of 200 to 500. Nowadays with advances in computing, values of 2000–5000 are common. These simulation studies are performed in a straightforward manner.

1. Values for  $k$ ,  $p$ , and  $n$  are chosen.
2.  $n$  observations are generated randomly from a normal distribution with zero mean and unit variance for  $k + 1$  variables (one  $y$  variable and  $k$  descriptors).
3.  $p$  of the  $k$  descriptors are selected with one of these stepping methods (or with an all subset approach), and the value of say  $R^2$  is recorded.  
Steps 2 and 3 are repeated 2000–5000 times to give 2000–5000 values of  $R^2$ .
5. Those  $R^2$  values are then ordered smallest to largest.
6. The  $\alpha\%$  significance level critical value of the  $R^2$  distribution is the  $(100 - \alpha)$ th percentile of the ordered values in 5.

Table 4 compiles the values of the upper 95% points of  $R^2$  as found by Rencher and Pun.<sup>74</sup> In their simulation studies with Efroymson's stepwise regression algorithm, Rencher and Pun adjusted the *F*-to-enter and *F*-to-remove values to ensure that they obtained  $p$  variables in their final model. Their results, like those from most other simulations, are based on randomly generated  $x$  variables (i.e., they are uncorrelated), which on the face of it might seem of little use because in "real" datasets, the variables will be correlated to some extent. Simulation results with correlated predictors give smaller critical values than do those for the corresponding uncorrelated ones. Consequently, if an observed  $R^2$  is greater than the critical value based on the uncorrelated  $x$  simulations, the observed  $R^2$  will be greater than that for the correlated. The use of these tables is demonstrated in the case study, which follows later in the chapter.

**Table 4** Values of the Upper 5% Points of the Distribution of  $R^2$  for Variable Subsets Selected Using the Stepwise Regression Algorithm of Efroymson

$n^{(a)}$	$k^{(b)}$	Number of Variables in Selected Subset				
		$p^{(c)} = 2$	$p = 4$	$p = 6$	$p = 8$	$p = 10$
5	5	0.990				
	10	0.997				
	20	0.999				
10	5	0.733	0.867			
	10	0.808	0.958	0.996		
	20	0.850	0.979	0.999	1.000	
	30	0.869	0.985	1.000	1.000	
20	40	0.880	0.988	1.000	1.000	
	10	0.493	0.679	0.777		
	20	0.543	0.745	0.857	0.923	
	30	0.567	0.775	0.887	0.948	0.980
30	40	0.583	0.794	0.904	0.961	0.988
	10	0.348	0.501	0.588		
	20	0.389	0.567	0.682	0.763	
	30	0.409	0.598	0.723	0.809	0.870
40	40	0.423	0.619	0.748	0.837	0.898
	10	0.268	0.394	0.468		
	20	0.302	0.453	0.558	0.634	
	30	0.319	0.482	0.598	0.686	0.753
50	40	0.331	0.501	0.624	0.718	0.790
	30	0.261	0.402	0.507	0.589	0.655
	40	0.271	0.420	0.532	0.622	0.695
60	30	0.221	0.345	0.439	0.515	0.577
	40	0.229	0.360	0.463	0.547	0.617

(a)  $n$  = sample size.(b)  $k$  = number of uncorrelated predictors.(c)  $p$  = number of variables in model.

## All Subset Regression

As might be expected from its name, all subset regression produces all possible regressions of a certain size. At one time, this would have been a way of tying up your computer for many hours, if not days, even for a moderate number of variables from which to select. This is now really not a major problem with the advent of relatively cheap and rapid computer power. Selection of the best subset of any given size usually involves picking the one that maximizes  $R^2$  (or one of the alternative procedures that will be introduced). All subset selection techniques will find the combination of variables that maximizes or minimizes a criterion function. This property is not guaranteed by any of the stepwise methods. Also, whereas forward stepping and stepwise regression may result in different subsets being selected if the order of the

variables in a dataset is changed, the all subsets method, because of its exhaustive search, results in the same answer irrespective of the variable ordering.

To conclude the list of approaches for variable selection, we briefly mention the idea of using the output from an all variable selection run as input into a latent variables regression analysis.<sup>47</sup> This technique is expensive computationally as it requires not only the subsets of variables to be found but also the number of latent variables needed to optimize some criterion function.

## Other Stopping Rules

Alternative stopping rules for variable inclusion are available, including minimizing the residual sum of squares  $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ . This quantity can be written for the squared multiple correlation coefficient  $R^2$ , as  $\sum_{i=1}^n (y_i - \bar{y})^2 (1 - R^2)$ . This sum will be minimized if  $R^2$  is maximized, but this approach is problematic because  $R^2$  always increases as new variables are included in the model. As a consequence, minimizing the residual sum of squares will always select all variables. A better way is for us to use the adjusted  $R^2$ ,  $\bar{R}^2$ , given by

$$\bar{R}^2 = R^2 - \frac{p(1 - R^2)}{(n - p - 1)}$$

where  $p$  is the number of regressors in the model. The adjusted  $R^2$  increases if added variables make a “significant” contribution to reducing the residual sum of squares; otherwise it decreases. Another criterion is Mallows’  $C_p$ ,<sup>79</sup> which can be written as

$$C_p = \frac{ESS_p}{s_k^2} - n + 2(p + 1)$$

where  $\frac{ESS_p}{s_k^2}$  is the ratio of the residual sum of squares with  $p$  variables in the model to the residual variance when all  $k$  variables are included. Variable selection with  $C_p$  is such that at each step, the variable to be included is the one giving the smallest  $C_p$ . Selection stops when  $C_p$  is at its minimum. In situations where  $n > k$ , Mallows’  $C_p$  can lead to selection of subsets, which when based on least squares estimation of the regression coefficients, could lead to a larger prediction error than if all  $k$  variables had been used by scientists and no selection had taken place.

An increasingly popular stopping rule for variable selection involves calculating a statistic referred to as *PRESS*. *PRESS* stands for predicted residual error sum of squares. To calculate *PRESS*, the data must be split into at least

two parts. In its simplest form, the data would be split at random into two sets, a **training** set and a **test** or **validation** set. Note that this terminology is not universal, and some authors refer to the whole dataset as the training data, and the two disjoint subsets as the **construction** set and the validation set. A model is first fitted to the training set. The test set  $x$  values are then applied in that model to predict their corresponding  $y$  values. The *PRESS* statistic, which indicates how well the model predicts, is then calculated as

$$PRESS = \sum_{i=1}^{n_{test}} (y_i - \hat{y}_i)^2 \quad [34]$$

where the  $y_i$  is the  $i$ th observed value in the test set,  $\hat{y}_i$  is the corresponding predicted value, and  $n_{test}$  is the size of the test set. As  $n_{test}$  increases so does the *PRESS*. So, when comparisons are being made between model predictions based on validation sets of different size, normalizing its value is required by dividing *PRESS* by, say, the validation sample size to give the mean square *PRESS*; i.e.,  $MSP = PRESS/n_{test}$ . Variable selection ceases when *PRESS* is a minimum. If further  $x$  variables are included in the model beyond this point, the model may have an improved fit to the training data, but it does so at the expense of not being able to predict “unseen” data. What is happening as variables are included is that the model is trying to predict (at least in OLS) the average value of the response variable for the given values of the  $x$  variables. Beyond the point in the variable selection process where the prediction of the average value has been optimized, the model attempts to explain the variation in individual data values. If taken to extreme, and provided sufficient data exist, the model passes through every datum leaving no residual error. Clearly the model is “overfitting” the data, and its ability to generalize the fit to predict a test set will suffer. An alternative to using *PRESS* is the validation  $R^2$ , sometimes referred to as  $q^s$ ,<sup>80</sup> which is calculated as

$$R^2 = 1 - \frac{PRESS}{PRESS_0} \quad [35]$$

Here  $PRESS_0$  is *PRESS* (Eq. [34]) with  $\bar{y}$  replacing  $y_i$  and is simply the total sum of squares for the response values in the test set. In situations where the entire dataset is small, it might not make sense to split the data into training and test/validation sets as one extra sample value placed into the test set means one less is available for constructing the model. In these situations, another form of validation undertaken is referred to as cross-validation or CV. Again in its simplest form the entire dataset is split into two equal parts but instead of fitting on one half and predicting the other, the two parts of the dataset are swapped and the half that was omitted originally fits the model and the other set validates it. The resulting *PRESS* value will then be the sum of the two

separate values. Splitting the data in this manner is only feasible if the number of variables in the model is  $p < \frac{1}{2}n$ , i.e., less than half the sample size. When smaller datasets are encountered, CV can be performed with what is termed “leave  $\nu$  out” where  $1 \leq \nu \leq n/2$ . When  $\nu = 1$ , the CV, procedure is known as leave-one-out or LOO-CV. In LOO-CV, a data point is left out of the analysis, a model is developed for the remaining dataset, and that model then predicts the response of the omitted point. This procedure is repeated  $n$  times until all data points have been left out once, and the *PRESS* is calculated as before. With  $\nu \geq 1$ , and  $n = c \times \nu$ , where  $c$  is an integer, all CV groups are disjoint and of the same size. The CV “loop” would be performed  $c$  times and is referred to as  $\nu$ -fold CV.<sup>47</sup> Yet another form of CV exists known as leave-multiple-out CV (LMO-CV). Here the dataset is split randomly into two parts of size  $d$  and  $n - d$ . The model is fitted to the  $n - d$  data points, and it then predicts the remaining  $d$  response values. This procedure is performed several times, each time randomly splitting the data into  $n - d$  and  $d$  disjoint groups.

Baumann, Albert, and von Korff<sup>47</sup> have reviewed the relative merits of this and other CV methods based on both published results and their simulation experiments. They found that for generating overfitted models,  $\nu$ -fold CV was the worst, followed by LOO-CV and then LMO-CV. These authors also report on the effect of varying the split between the construction set and the validation set and conclude that the “best” results come from models constructed from LMO-CV where the size of the construction dataset is smaller than the validation set. In fact, what Baumann et al. found was that beyond a certain minimum size construction set, which in their case was three times the number of “true”  $x$  variables (i.e., variables that generated the responses), overfitting begins to occur.

---

## CASE STUDY

To illustrate some of the techniques that have been discussed in this chapter, we have chosen the charge-transfer dataset introduced earlier in the section on dimension reduction.<sup>31</sup> This dataset comprises the response variable Kappa and the 30 descriptive variables, the definitions of which can be found in reference 31:

1. Px	2. Py	3. Pz	4. ClogP	5. CMR	6. Xmin
7. Xmax	8. Ymax	9. Zmin	10. Zmax	11. Mux	12. Muy
13. Muz	14. Mu	15. Ehomom	16. Chge(1)	17. Sn(1)	18. Fn(2)
19. alp(2)	20. Sn(2)	21. alp(3)	22. Sn(3)	23. Se(3)	24. Fn(4)
25. Fe(4)	26. alp(4)	27. Sn(4)	28. Fe(5)	29. Fn(6)	30. Fe(6)

Before beginning the regression analysis of this data, three points must be established. First, it is unlikely that just one model can be considered to be the

correct model. Second, when a model is constructed, it should be remembered that the beta coefficients in the equation are estimates subject to variability; as a consequence, we should not think of the predicted values of the response as any thing more than estimates themselves, also subject to variability. Third, if the assumption of normality of the residual terms in Eq. [13] is reasonable, and ordinary least squares was the method of model construction, then confidence intervals for these estimates can be constructed. However, if the normality assumption is invalid, alternative confidence intervals should be obtained with bootstrapping techniques.<sup>81</sup> In light of these three points, we can now develop models with the methods described earlier.

### Stepwise Regression

In this part of the tutorial, we start the data analysis with stepwise regression, which is perhaps the most popular model construction technique. Table 5

**Table 5** Stepwise Regression Results for the Charge-Transfer Data Set

Step	1	2	3	4	5	6	7	8	9
Const	-1.80	-0.47	0.53	0.33	0.15	0.31	0.34	0.22	0.25
p-val	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
CMR	0.032	0.045	0.045	0.041	0.040	0.040	0.038	0.035	0.032
T-val	4.97	16.25	20.68	16.61	16.57	17.09	16.73	11.80	9.68
p-val	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
ClogP		-0.34	-0.35	-0.34	-0.34	-0.34	-0.33	-0.33	-0.33
T-val		-13.3	-16.9	-17.7	-18.6	-19.3	-20.1	-21.1	-20.9
p-val		0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Ehomo			0.09	0.10	0.07	0.08	0.08	0.08	0.08
T-val			4.45	5.13	3.54	3.97	4.09	4.30	4.26
p-val			0.000	0.000	0.000	0.000	0.000	0.000	0.000
Pz				0.14	0.13	0.10	0.11	0.15	0.15
T-val				2.76	2.77	2.14	2.44	3.14	3.35
p-val				0.01	0.01	0.04	0.02	0.00	0.00
Mux					-0.03	-0.04	-0.04	-0.04	-0.04
T-val					-1.94	-2.26	-2.46	-2.76	-3.18
p-val					0.06	0.03	0.02	0.01	0.00
Sn(1)						-0.00	-0.00	-0.00	-0.00
T-val						-1.99	-2.44	-2.46	-2.40
p-val						0.06	0.02	0.02	0.00
Sn(2)							0.01	0.01	0.01
T-val							2.02	2.27	2.07
p-val							0.05	0.03	0.05
Px								0.03	0.04
T-val								1.94	2.49
p-val								0.06	0.02
Fe(4)									-0.26
T-val									-1.67
p-val									0.11
$R^2_{fit}$	0.4280	0.9121	0.9463	0.9572	0.9621	0.9668	0.9712	0.9748	0.9774



shows the Minitab output for stepwise regression with the recommended alpha to enter of 0.15 (Dillon and Goldstein,<sup>40</sup> p. 341). At each step, the values of the parameter estimates are shown together with their *t*-values and associated significance levels (*p*-values). Nine variables enter the model [CMR, ClogP, Ehomo, Pz, Mux, Sn(1), Sn(2), Px and Fe(4)] (no variables are removed), and although the values of the beta coefficients change, these changes are relatively small, indicative of fairly low levels of multicollinearity. The predicted regression model and associated statistics is

$$\begin{aligned} \text{Kappa} = & 0.254 + 0.0319 \text{ CMR} - 0.327 \text{ ClogP} + 0.0766 \text{ Ehomo} + 0.154 \text{ Pz} \\ & - 0.0440 \text{ Mux} - 0.00102 \text{ Sn(1)} + 0.00647 \text{ Sn(2)} + 0.0423 \text{ Px} \\ & - 0.258 \text{ Fe (4)} \end{aligned} \quad [36]$$

$$s = 0.07820 \quad \text{R-Sq} = 97.7\% \quad \text{R-Sq(adj)} = 96.9\%$$

$$\text{PRESS} = 0.361428 \quad \text{R-Sq(pred)} = 94.65\%$$

Referring to the variable numbering system given earlier, we have variables 1, 3, 4, 5, 11, 15, 17, 20, and 25 in the model. The R-Sq is the coefficient of determination telling us that 97.7% of the variation in Kappa can be explained by these variables in a least squares linear regression model. The PRESS and R-Sq(pred) are the LOO-CV statistics. The order in which the variables enter the model is largely determined by the contribution they make to explaining the variation in the response variable that is not explained by the variables already present in the model. However, if two variables exist that explain individually the same amount of residual variance, the one that goes into the model will then be determined by its position in the data file.

The analysis of variance results for this regression was obtained with Minitab by performing an ordinary regression with the variables identified in the stepwise run. These results are given in Table 6. It is the statistics in Table 6 that are usually quoted along with the *p*-value when researchers publish their regression models. The tacit assertion being made is that because the *p*-value is less than 0.05, the result is significant. In other words, at least one of the *x* variables makes a significant contribution to the variation in the response variable. Alternatively, we can use tabulated critical *F*-values, which for 9 and

**Table 6** Analysis of Variance for a Model Selected by Stepwise Regression

Source	DF <sup>a</sup>	SS <sup>b</sup>	MS <sup>c</sup>	F-ratio	<i>p</i> -value
Regression	9	6.59851	0.73317	119.90	0.000
Residual Error	25	0.15288	0.00612		
Total	34	6.75139			

<sup>a</sup>DF = Degrees of freedom.

<sup>b</sup>SS = Sums of squares.

<sup>c</sup>MS = Mean square.

25 of freedom gives a value around 3.25. The conclusions would then be that this result is significant at the 5% level because the observed value of the test statistic ( $F$ ) is clearly in excess of this critical value. Unfortunately, as we have pointed out, these significance levels are not valid because no allowance has been made for the fact that the nine variables have been selected from a pool of 30 possible regressive variables. To assess the true significance of this model, it is necessary to compare the observed statistics with critical values derived from Monte Carlo simulation studies. Rencher and Pun<sup>74</sup> compiled critical values of  $R^2$  at the 5% significance level based on Monte Carlo simulations of Efroymsen's<sup>71</sup> stepwise regression algorithm. Although their table covers a range of sample sizes  $n$  (5, 10, 20, 30, 40, 50, 60), pool size  $k$  of possible regressors (5, 10, 20, 30, 40), and model size  $p$  (2, 4, 6, 8), it is far from complete. That notwithstanding they provide the following formula, which acts as a means of calculating combinations of  $(n, k, p)$  not provided directly by their table:

$$R_\alpha = F^{-1} \left[ 1 + \frac{\log_e(1 - \alpha)}{(\log_e N)^{1.8N^{0.04}}} \right] \quad [37]$$

In Eq. [37],  $F^{-1}(\cdot)$  is the inverse of the incomplete beta function with parameters  $p/2$  and  $(n - p - 1)/2$ ,  $N = {}^kC_p$ , and  $\alpha$  is the significance level. We have used Eq. [37] to generate a set of critical  $R^2$  at the 5% significance level (Table 4). Using the relationship

$$F = \frac{R^2(n - p - 1)}{(1 - R^2)p} \quad [38]$$

we can find the corresponding critical  $F$ -values.

Applying Eq. [37] to our current model, we note that we have a sample size of  $n = 35$  and a pool of  $k = 30$  variables from which we are selecting  $p = 9$ . Substituting these values into Eq. [37] along with the calculated value for  $N = \frac{30!}{9!(30 - 9)!} = 14307150$ , we obtain a critical  $R^2$  value of 0.779 and a critical  $F$ -value of 9.80. In comparison the critical value from standard  $F$ -tables gives 2.28, which illustrates the overoptimistic nature of the significance levels when no account is made of the selection procedure. From Table 6, we have an observed  $F$ -value of 119.90, which is well in excess of 9.80 and hence is significant at the 5% level.

### Best Subset Regression

To carry out an all subset regression means that all models will be constructed of size one (of which there are 30), of size two (there are

Table 7 Best Subset Regression Models of Size 2 to 9

Model size 2	$R^2_{fit}$	Model size 3	$R^2_{fit}$	Model size 4	$R^2_{fit}$
4,5	0.91213	4,5,15	0.94634	4,5,23,29	0.96039
4,8	0.57775	4,5,29	0.94465	4,5,13,29	0.95936
5,14	0.57485	4,5,23	0.93935	3,4,5,15	0.95723
11,14	0.55195	4,5,11	0.93816	4,5,15,29	0.95669
4,6	0.55017	4,5,13	0.93575	4,5,15,17	0.95420
Model size 5	$R^2_{fit}$	Model size 6	$R^2_{fit}$	Model size 7	$R^2_{fit}$
4,5,22,23,29	0.96676	1,4,5,22,23,29	0.96962	1,3,4,5,11,13,25	0.97441
4,5,13,22,29	0.96438	4,5,22,23,24,29 29	0.96903	1,3,4,5,22,23,29	0.97162
4,5,15,22,29	0.96387	4,5,15,22,24,29	0.96819	4,5,9,10,23,25,29	0.97092
4,5,20,23,29	0.96366	4,5,9,10,23,29	0.96815	4,5,9,10,22,23,29	0.97090
4,5,21,23,29	0.96263	4,5,10,22,23,29	0.96791	1,3,4,5,20,23,29	0.96997
Model size 8	$R^2_{fit}$	Model size 9	$R^2_{fit}$		
1,3,4,5,11,13,17,25	0.97591	1,3,4,5,11,15,17,22,25	0.97817		
1,3,4,5,10,11,13,25	0.97573	1,3,4,5,11,12,21,22,25	0.97796		
1,3,4,5,11,13,20,25	0.97537	1,3,4,5,11,15,17,20,25	0.97736		
1,3,4,5,11,12,13,25	0.97524	1,3,4,5,11,15,17,25,27	0.97721		
1,3,4,5,11,13,25,28	0.97519	1,3,4,5,11,12,13,25,28	0.97661		

${}^{30}C_2 = \frac{30!}{2!28!} = 435$ ), of size three (there are  ${}^{30}C_3 = \frac{30!}{3!27!} = 4060$ ), and so on.

The top five models, in terms of  $R^2_{fit}$ , are presented in Table 7 for regressions with 2, 3, ..., 9 variables. With the exception of the two-variable model, we cannot distinguish the top five models of each size. Leave-one-out cross validation  $R^2$  has also been calculated (not shown) and like  $R^2_{fit}$ , does not change much either, ranging between 0.912 and 0.978 for the two- to nine-variable models. This plethora of models all have similar fit statistics. The novice should be cognisant when performing regression or reading the literature that although most people produce only one model, chances are that other models exist having very similar fit characteristics. Looking at Table 7, for example, we see that the nine-variable model identified by this stepwise procedure is in fact the third best model in terms of  $R^2_{fit}$ .

As with stepwise regression, the statistical significance of these best subset models must be viewed with caution. Again simulated critical values must be generated, but this time they must be based on a best subset approach to match the method of model construction. It is possible here for us to use Rencher and Pun's formula as a guide because it is based on Efroymson's stepwise approach, thus giving slightly lower values than those required for an exhaustive search. Because we are focusing on the same size model as that produced by the stepwise method, the critical values are the same, and as the observed  $R^2_{fit} > 0.779$ , we can say the result is significant at the 5% level.

---

## SUPERVISED AND UNSUPERVISED VARIABLE SELECTION

Whitley, Ford and Livingstone<sup>46</sup> recommend the following strategy for model construction:

1. Select variables that have significant correlations with the response variable above some lower bound (supervised selection)
2. Omit columns that have variance below some minimum value
3. Omit columns to reduce the multicollinearity in the variable set (unsupervised)
4. Select a subset of variables to construct the model

This procedure can be performed either with item 1, in which case it is considered to be a supervised variable selection meaning that the response variable has selected variables, or without item 1, relegating it to an unsupervised selection category.

The analysis that we now describe was performed with the software package Paragon developed at the Centre for Molecular Design at the University of Portsmouth. Paragon is available from ([www.cmd.port.ac.uk](http://www.cmd.port.ac.uk)). With this software, we can set various levels of multicollinearity (step 3) that is to

be tolerated ( $R^2_{\max}$ ) together with a range of minimum correlations between the response variable and the regressors (step 4).

Adopting the unsupervised option initially, the first two variables to be selected are those with the lowest pairwise correlation. The next variable selected has the smallest multiple squared correlation with those first two variables. This process is continued until the preset maximum level of multicollinearity (determined by the squared multiple correlation coefficient) is reached. Whitley et al. refer to this procedure as unsupervised forward selection (UFS). UFS can also be performed with a minimum variance criterion where only variables with variance above this minimum will be selected. These two criteria can be used by scientists simultaneously. With supervised variable selection, only those variables having a sufficiently high correlation with the response are considered for what effectively is UFS on this reduced set of variables. We will term this latter process, supervised forward selection (SFS). To see how these options work and to examine the effect they have on the model produced, we performed PLS on the data with both UFS and SFS configured to run with a range of  $R^2_{\max}$  values and response variable correlations (Table 8). The rows in Table 8 are numbered to aid the discussion. The first two entries are the results of setting UFS to allow a correlation of 0.999 between any  $x$  variable and the rest (this accepts a high degree of multicollinearity) but to not include variables with a standard deviation less than 0.01. The effect of this second condition results in the removal of 4 of the 30 variables. PLS is then performed with the remaining 26 variables, with the probability ( $\alpha$ ) of a component entering the model set to 0.05 and 0.15. We used this latter probability value in the stepwise regression and is included here for comparison purposes. These so-called “significance levels” need to be viewed with the

**Table 8** Influence of  $R^2_{\max}$  and  $\alpha$  on PLS Results for UFS and SFS

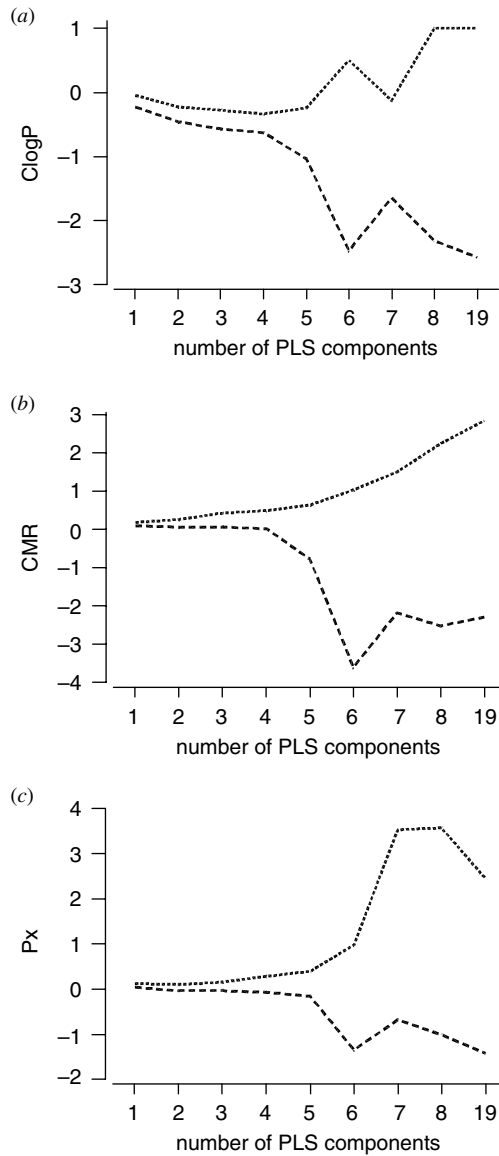
	UFS $R^2_{\max}$	Minimum Variance	SFS $R^2$	Number of Variables	$\alpha$ to Enter	Number of Comp	$R^2_{fit}$	$R^2_{CV}$
1	0.999	0.01	*	26	0.05	4	0.9727	0.9055
2	0.999	0.01	*	26	0.15	8	0.9829	0.7107
3	0.990	0.10	*	21	0.05	8	0.9693	0.7054
4	0.990	0.10	*	21	0.15	10	0.9776	0.4151
5	0.950	0.10	*	17	0.05	4	0.9547	0.8438
6	0.950	0.10	*	17	0.15	4	0.9511	0.8438
7	0.900	0.10	*	14	0.05	4	0.9508	0.8208
8	0.900	0.10	*	14	0.15	4	0.9508	0.8208
9	0.990	0.10	0.05	12	0.05	5	0.9365	0.8682
10	0.990	0.10	0.10	14	0.05	6	0.9382	0.8584
11	0.990	0.10	0.20	17	0.05	7	0.9528	0.8471
12	0.950	0.10	0.05	11	0.05	5	0.8746	0.6621
13	0.950	0.10	0.10	13	0.05	6	0.8777	0.5920
14	0.950	0.10	0.20	15	0.05	4	0.8549	0.7456
15	ols			9			0.9769	0.9470
16	bss			9			0.9774	0.9495

same critical eye as we did previously with the stepwise and best subset regressions. Because the 26 data matrix is of full rank, it is possible to extract 26 components. However, we pick only those that pass some entry test, and as a result the model is subject to the same sort of selection bias discussed earlier. With four components, the first model explains 97.27% of the variation in the response variable and produces a LOO-CV  $R^2$ ,  $R^2_{CV}$ , of 0.9055. Allowing more components to enter (row 2 of Table 8) will always increase the  $R^2_{fit}$ , but it may reduce the  $R^2_{CV}$  if the model is overfitting the data and thus is not able to generalize to unseen data.

Reducing the multicollinearity (Table 8, rows 3 to 8) clearly reduces the number of variables available for modeling and, with the exception of row 4, increasing the  $\alpha$  level of entry, has no influence on the number of components entering the model. In rows 9 to 14, we have introduced SFS by making available for model construction only those predictors that correlate (significant at 5%, 10%, or 20% level) with the response variable. This SFS has been performed at two levels of UFS (0.99 and 0.95) and, as to be expected, increasing the value of  $\alpha$  to enter, increases the number of variables available for the UFS part of the process. The only noteworthy observation of these results is that for a UFS  $R^2_{\max} = 0.99$ , more components enter the model as more variables enter the final selection pool. But, as the fit improves, the cross-validation moves in the opposite direction, albeit not very far.

The final two rows of Table 8 are the results for the stepwise (row 15) and best subset (row 16) regressions. Both of these models have the smallest number of variables (9) and have very similar fit and cross-validation statistics. The only other models approaching these values are the first two models involving 26 variables. Although it is risky for us to use the results from one data set and make generalizations, we are of the opinion that when using data sets containing a high degree of multicollinearity, we should attempt to reduce the pool of variables. The reason for promulgating this view was discussed and is restated here; two highly colinear variables with roughly the same response correlation are very unlikely to have individually any unique information about the response not present in the other.

When using these biased regression methods, we should be aware that although they are designed to cope with multicollinearity issues, as the number of components in the model increases, so does the level of multicollinearity. In fact as we continuously add components, the pseudo-beta coefficients of the original  $x$  variables approach those obtained by performing OLS on the original variables. Consequently the error in the predicted response will increase because of the reintroduction of inflated regression coefficient variances. We have illustrated this situation with a 19-variable PLS model where we have produced a series of models with 1, 2, 3, . . . , 8 components (Figure 9). Plotted in Figure 9 are the values of the bootstrapped 95% confidence intervals for the pseudo-betas for three of the variables for models of increasing size. Also plotted are the OLS values that coincide with the 19-component model. We



**Figure 9** Plots of bootstrapped pseudo-regression coefficient confidence intervals against number of PLS components for the three variables ClogP, CMR and Px.

have plotted the standardized coefficients so that all of them can be displayed on the same scale. Note that the confidence intervals remain fairly stable until component 5 or 6 is added. At that point, the multicollinearity has increased to a level that results in the instability in the regression coefficients and inflated

variances. Also, whereas the intervals have either two positive values or two negative values for the smaller component models (indicating that the actual regression coefficient may be significantly different to zero), this all changes as the multicollinearity “kicks” in when the intervals straddle zero.

These models are linear. That is, linear combinations of the descriptor variables modeled the response ( $\kappa$ ) variable. Does the situation change if we employ a nonlinear modeling method, such as an artificial neural network (ANN), and then assess variable selection? In the original report by Livingstone, Evans, and Saunders<sup>31</sup> in which a large set of descriptors was employed, a subset of 11 variables was identified as having a high correlation with the response. The first nine of these were the same variables as selected by stepwise regression. Moreover, they are in the same order, with the extra two being Mu and Sn(3). Variable selection in the neural network field is often referred to as “pruning.” It has this name because it generally aims to prune unnecessary connections between neurons and even some of the input (variable) neurons if all their connections are severed. Several different pruning algorithms exist that can be broadly classified as “magnitude based” when they depend on an analysis of the connection weights or “error based” when they take account of network error after the elimination of connections.<sup>82</sup>

Livingstone, Manallack, and Tetko<sup>83</sup> constructed a backpropagation artificial neural network<sup>84</sup> with 11 input neurons (corresponding to the 11 independent variables), five hidden layer neurons, and a single output neuron. Also, two bias neurons are present, one on the input layer and the other on the hidden layer. Because more connection weights exist than training examples, it was important for them to use a training method that would prevent overfitting so an early stopping technique was employed.<sup>85</sup> The network was initialized from random starting weights and trained 400 times, thus giving a “family” of trained networks on which to apply the pruning. Application of a magnitude-based<sup>83</sup> pruning method identified variable 11 as the least important and thus the first to be removed in the variable selection, which is not surprising as this descriptor has the lowest correlation with  $\kappa$ . The next variable to be identified for removal, however, is more surprising. That descriptor is 6 (Sn(1)), which was followed at the next step by removal of variable 7. Thus, variables 10, 9, and 8, which are identified at later steps in the stepwise regression were shown to be of greater importance in the nonlinear ANN models. The next variable to be identified for removal is perhaps even more surprising as this is variable 3 (Ehomo), the third most important variable identified by stepwise regression. The complete sequence of variables identified as being unimportant by this sensitivity analysis was 11, 6, 7, 3, 8, 10, 9, 5, 4, and 2. Clearly the most important variables in the linear models, 1, 2, 4, and 5 are also important in these nonlinear models, but the contributions of the other descriptors are ranked differently.

This demonstration was not intended to be a step-by-step approach showing how to construct a model. Instead it was meant to illustrate some



of the concepts reviewed earlier in the chapter and to demonstrate that many models exist with similar fit and CV statistics. Many other modeling approaches could have been applied to this dataset, and no doubt it is possible to build a host of models of differing complexity. What the reader should bear in mind, however, is a “common sense” approach to judging the quality of any of these models. Why do we build these models in the first place, and what are we going to do with them? The answer in most cases is either to gain understanding or to be able to make predictions, and sometimes both. Having looked at what the statistics has to offer in model building, it is now time to focus on the chemistry associated with such models.

---

## PUBLISHED VARIABLE SELECTION METHODS

---

The number of published accounts of variable selection methods in the general literature is enormous. To provide a focus, this section will concentrate just on applications to computer-aided drug design. Variable selection was identified as an important requirement at about the same time as the need for variable elimination techniques.<sup>86</sup> The simplest method of variable selection is to choose those variables that have a large correlation with the response and, for simple datasets, that method is probably not a bad choice. As we have shown in this chapter, variable selection may be an integral part of a modeling technique, but not all modeling methods lend themselves to variable selection, and in these cases, other techniques need to be applied.

The need for variable selection was pointed out when a particular dataset, the Selwood set, became something of a standard or a benchmark for scientists to assess new model building techniques. The Selwood set is named after the chemist who made most of the compounds; it consists of 53 descriptors computed for 31 compounds.<sup>87</sup> Overall,  $7.16 \times 10^{15}$  possible regression models (all equations up to 29 terms) can be built for this dataset.<sup>88</sup> If the regression equations were calculated at a rate of 1 per second, it would take 227 million years to compute all possible models. Clearly, an exhaustive search is impossible in this case so researchers turned their attention to nature for a solution to this and similar intractable problems. The “natural solution” to optimization problems in a large solution space is evolution, and J. H. Holland<sup>89</sup> provided the inspiration for scientists to use evolutionary approaches in fields of study as diverse as digital image processing and criminology. A generic term for this type of solution is *genetic algorithms* (GA), and various “flavors” of this methodology exist such as evolutionary algorithms, genetic function approximations, evolutionary programming, and evolution strategies. A good introduction to GAs and their application in drug design is given in the account of a workshop held in 1995.<sup>90</sup> GAs and their applications in chemistry have also been reviewed in a pedagogical manner in this book series by Peterson.<sup>91</sup>

The application of different types of GA strategies for variable selection in drug design all happened at about the same time. Leardi et al.<sup>92</sup> reported a GA designed specifically for variable selection and applied this to the Selwood set to generate a population of models that identify important variables.<sup>93</sup> Rogers and Hopfinger<sup>94</sup> applied a genetic function approximation to the same set and produced several sets of “best” models containing 3, 4, 5, and 6 variables. Using evolution to produce a set of “best” models is one of the features of the GA approach because they produce populations of solutions; variable selection involves counting the frequency of occurrence of variables in these best models. Kubinyi<sup>88</sup> described an evolutionary approach called Mutation and Selection Uncover Models (MUSEUM) that when applied to the Selwood set, discovered some of the same equations published by Hopfinger and some others. Because variable selection is based on the resulting population of GA-generated models, the results depend critically on the fitness function that the GA is aiming to optimize. Kubinyi went on to combine this evolutionary approach with a systematic search procedure that was claimed to be highly efficient. He demonstrated its application to regression models<sup>95</sup> and to PLS.<sup>96</sup> Before finishing the discussion of variable selection methods applied to the Selwood set, we point out that a neural network has been applied to this end,<sup>97</sup> as have a neural network coupled with a GA<sup>98</sup> and a method called cluster significance analysis<sup>99</sup> that was developed originally to judge the likelihood of clusters developing by chance. Kubinyi<sup>95</sup> has nicely summarized the results of some of these different approaches.<sup>95</sup> The Selwood set continues to be a reference standard for the development of new analytical methods, including the problem of variable selection,<sup>98,100–105</sup> but despite all this work, only a few reports have commented on the odd nature of the distribution of some of the variables in the set<sup>46,93,106</sup> and only one has suggested that it may be time for scientists to stop using it.<sup>46</sup>

GA-based approaches have optimized the production of many different types of models (artificial neural network architectures in particular) and simultaneously selected variables and optimized neural network models.<sup>98,107–109</sup> GAs coupled with well-known and less well-known modeling methods have also been used by scientists in variable selection. The combination of a GA with multiple linear regression was shown to perform well on datasets containing 15, 26, and 35 descriptors.<sup>110</sup> PLS coupled with a GA has also been shown to be useful in variable selection.<sup>111–114</sup> Spline fitting with a GA to select variables has been compared with several other classification techniques.<sup>115</sup>

A great deal of interest exists in methods for variable selection as well as for model evaluation, which are actually two sides of the same coin. Bayesian neural networks include a procedure called automatic relevance determination (ARD), allowing for the identification of important variables.<sup>116,117</sup> A *k*-nearest neighbor method for variable selection has been applied successfully to problems of biological activity<sup>118,119</sup> and metabolic stability.<sup>120</sup> Other

approaches to variable selection include the use of information theory,<sup>121</sup> artificial ant colonies,<sup>122</sup> variable weighting,<sup>123</sup> binary particle swarms,<sup>124</sup> pair correlations,<sup>125</sup> Fisher's discriminant ratio,<sup>126</sup> and aggregated, bootstrapped, regression, or classification trees.<sup>127</sup>

So, given these different approaches to variable selection, is it possible to recommend a particular procedure? The answer is no, unfortunately, because these reports involve many different datasets, model building strategies, and intended aims, thus making direct comparisons difficult. Reports on comparing variable selection strategies are beginning to appear,<sup>128</sup> however, and no doubt in the fullness of time a superior approach or sequence of approaches will be identified.

---

## CONCLUSIONS

The development of computational chemistry software and techniques, coupled with the increasing speed and decreasing costs of computing machinery, has transformed computer-aided material design, particularly drug design. The calculation of many chemical descriptors for almost any kind of molecule is now a trivial problem. Variable selection, however, is not trivial and becomes necessary when a dataset contains many variables. What constitutes as many depends both on the use that will be made of the data by the scientist and the ratio of data points (cases) to variables.

In this tutorial, we have shown that "variable selection" can be divided into three subtasks: dimension reduction, variable elimination and variable selection. The first of these tasks is reasonably well understood, and many standard and not-so-standard methods can process most types of datasets. Variable elimination is also relatively straightforward. Examination of the distributions of individual variables allows the easy identification of descriptors, containing little or no information. Distributions also allow the analyst to recognize properties that are associated with a particular compound or subset of compounds, either because of the underlying chemical rationale behind the descriptor or because of some division of the dataset into, say, training and test sets. The calculation of multicollinearity allows for the identification of redundancy and sets of variables containing essentially the same information in pairs, or as linear combinations of three or more descriptors. Algorithms for variable elimination have been described in this chapter, and software is available commercially and free from the web.

Variable selection, in contrast, is a much thornier problem; at first sight, it may seem to be solved easily in many model building approaches, but as we have pointed out, many pitfalls are possible in the process. "Established" procedures are used by many researchers, and albeit flawed in some respects, they do provide a good starting point for variable selection. Given the intense interest in the community, as judged by the numerous recent publications and

forthcoming conferences and workshops<sup>129–131</sup> we may well expect to see robust and useful methods emerging in the near future.

---

## ACKNOWLEDGMENTS

This work was partially supported by INTAS Grant 00-0363, Virtual Computational Chemistry Laboratory.

---

## APPENDIX

The singular value decomposition (SVD) exists for any matrix  $\mathbf{X}(n \times p)(n \geq p)$ . It is unique and is given by

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}' \quad [\text{A.1}]$$

where  $\mathbf{U}$  is the same size as  $\mathbf{X}$  ( $n \times p$ ) and  $\mathbf{S}$  and  $\mathbf{V}$  are square matrices of order  $p$ . The matrix  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{S}$  have the following properties:

1.  $\mathbf{U}'\mathbf{U} = \mathbf{I}$ , i.e., the columns of  $\mathbf{U}$  are orthonormal.
2.  $\mathbf{V}'\mathbf{V} = \mathbf{V}\mathbf{V}' = \mathbf{I}$ , i.e., the matrix  $\mathbf{V}$  is orthogonal.
3.  $\mathbf{S}$  is diagonal with elements  $s_1, s_2, \dots, s_p \geq 0$ , which are referred to as the singular values of  $\mathbf{X}$ .

From Eq. [10], we have  $C_\nu(\hat{\boldsymbol{\beta}}) = s^2(\mathbf{X}'\mathbf{X})^{-1}$ , which substituting from Eq. [A.1] gives

$$C_\nu(\hat{\boldsymbol{\beta}}) = s^2(\mathbf{X}'\mathbf{X})^{-1} = s^s[(\mathbf{U}\mathbf{S}\mathbf{V}')'(\mathbf{U}\mathbf{S}\mathbf{V}')]^{-1} = s^2[\mathbf{V}\mathbf{S}\mathbf{U}'\mathbf{U}\mathbf{S}\mathbf{V}']^{-1} \quad [\text{A.2}]$$

The  $\mathbf{U}'\mathbf{U}$  in the center of this last expression is equal to the identity matrix  $\mathbf{I}$  (property 1), which means that

$$C_\nu(\hat{\boldsymbol{\beta}}) = s^2[\mathbf{V}\mathbf{S}^2\mathbf{V}']^{-1} = s^2(\mathbf{V}\mathbf{S}^{-2}\mathbf{V}') \quad [\text{A.3}]$$

Comparing terms from Eqs. [A.2] and [A.3], it can be seen that

$$\mathbf{X}'\mathbf{X} = \mathbf{V}\mathbf{S}^2\mathbf{V}' \quad [\text{A.4}]$$

which has the form of the spectral decomposition of  $\mathbf{X}'\mathbf{X}$ . It is the spectral decomposition of  $\mathbf{X}'\mathbf{X}$  that is performed in PCA. It follows that the eigenvectors formed by PCA are the columns of  $\mathbf{V}$ , and the eigenvalues are the square of the singular values of  $\mathbf{X}$ . It follows, therefore, that the results of an SVD of  $\mathbf{X}$  can be obtained by performing a PCA on  $\mathbf{X}'\mathbf{X}$ .

## REFERENCES

1. C. Hansch, P. P. Maloney, T. Fujita, and R. M. Muir, *Nature*, **194**, 178 (1962). Correlation of Biological Activity of Phenoxyacetic Acids with Hammett Substituent Constants and Partition Coefficients.
2. C. Hansch and T. Fujita, *J. Am. Chem. Soc.*, **86**, 1616 (1964).  $\rho$ - $\sigma$ - $\pi$  Analysis. A Method for the Correlation of Biological Activity and Chemical Structure.
3. C. Hansch and A. J. Leo, *Substituent Constants for Correlation Analysis in Chemistry and Biology*, Wiley, New York, 1979.
4. C. Hansch and A. J. Leo, *Exploring QSAR, Fundamentals and Applications in Chemistry and Biology*, American Chemical Society, Washington, D.C., 1995.
5. M. Randic, *J. Am. Chem. Soc.*, **97**, 6609 (1975). On Characterisation of Molecular Branching.
6. L. H. Hall, L. B. Kier, and W. J. Murray, *J. Pharm. Sci.*, **64**, 1974 (1975). Molecular Connectivity II: Relationship to Water Solubility and Boiling Point.
7. L. B. Kier, W. J. Murray, and L. H. Hall, *J. Med. Chem.*, **18**, 1272 (1975). Molecular Connectivity. 4. Relationships to Biological Activities.
8. L. B. Kier and L. H. Hall, *Molecular Connectivity in Chemistry and Drug Research*, Academic Press, New York, 1976.
9. L. B. Kier and L. H. Hall, *Molecular Connectivity in Structure—Activity Analysis*, Wiley, New York, 1986.
10. A. K. Saxena, *Quant. Struct.-Act. Relat.*, **14**, 142 (1995). Physicochemical Significance of Topological Parameters: Molecular Connectivity Index and Information Content: Part 2. Correlation Studies with Molar Refractivity and Lipophilicity.
11. H. Kubinyi, *Quant. Struct.-Act. Relat.*, **14**, 149 (1995). The Physicochemical Significance of Topological Parameters. A Rebuttal.
12. A. K. Saxena, *Quant. Struct.-Act. Relat.*, **14**, 150 (1995). Reply to H. Kubinyi's Rebuttal.
13. J. C. Dearden and P. K. Mays, *J. Pharm. Pharmacol.*, **37**, 70P (1985). Can Molecular Connectivity Serve as a Steric Parameter in Quantitative Structure-Activity Relationships?
14. L. P. Burkhard, A. W. Andrew, and D. E. Armstrong, *Chemosphere*, **12**, 935 (1983). Structure Activity Relationships Using Molecular Connectivity Indices with Principal Component Analysis.
15. O. Kikuchi, *Quant. Struct.-Act. Relat.*, **6**, 179 (1987). Systematic QSAR Procedures with Quantum Chemical Descriptors.
16. R. C. Glen and V. S. Rose, *J. Mol. Graph.*, **5**, 79 (1987). Computer Program Suite for the Calculation, Storage and Manipulation of Molecular Property and Activity Descriptors.
17. R. M. Hyde and D. J. Livingstone, *J. Comp.-Aided Mol. Design*, **2**, 145 (1988). Perspectives in QSAR: Computer Chemistry and Pattern Recognition.
18. D. F. V. Lewis, *Xenobiotica*, **19**, 243 (1989). Molecular Orbital Calculations on Tumour-Inhibitory Aniline Mustards.
19. M. R. Saunders and D. J. Livingstone, in *Advances in Quantitative Structure-Property Relationships*, M. Charton, Ed., JAI Press, Greenwich, Connecticut, 1996, pp. 53–79. Electronic Structure Calculations in Quantitative Structure-Property Relationships.
20. D. J. Livingstone, *J. Chem. Inf. Comput. Sci.*, **40**, 195 (2000). The Characterisation of Chemical Structures Using Molecular Properties – a Survey.
21. J. J. Morris and P. P. Bruneau, in *Virtual Screening for Bioactive Molecules*, H. J. Bohm and G. Schneider, Eds, Wiley-VCH, Chichester, United Kingdom, 2000, pp. 33–58. Prediction of Physicochemical Properties.
22. D. J. Livingstone, *Curr. Top. Med. Chem.*, **3**, 1171 (2003). Theoretical Property Predictions.
23. R. Todeschini and V. Consonni, *Handbook of Molecular Descriptors*, Wiley-VCH, Mannheim, Germany, 2000.

24. D. J. Livingstone, *Data Analysis for Chemists*, Oxford University Press, Oxford, United Kingdom, 1995, pp. 118–121.
25. D. J. Livingstone, in *Molecular Design and Modeling: Concepts and Applications*, Vol. 203 of *Methods in Enzymology*, J. J. Langone, Ed., Academic Press, San Diego, California, 1991, pp. 613–638. Pattern Recognition Methods for use in Rational Drug Design.
26. D. J. Livingstone, *Data Analysis for Chemists*, Oxford University Press, Oxford, United Kingdom, 1995, pp. 67–81.
27. J. E. Jackson, *A User's Guide to Principal Components*, Wiley, New York, 1991.
28. D. J. Livingstone, *Data Analysis for Chemists*, Oxford University Press, Oxford, 1995, pp. 95–103.
29. E. Malinowski, *Factor Analysis in Chemistry*, 2nd Ed., Wiley, New York, 1991.
30. D. J. Livingstone, R. M. Hyde, and R. Foster, *Eur. J. Med. Chem.*, **14**, 393, (1979). Further Study of an Organic Electron-Donor-Acceptor Related Substituent Constant.
31. D. J. Livingstone, D. A. Evans, and M. R. Saunders, *Perkin Transactions II*, 1545 (1992). Investigation of a Charge-transfer Substituent Constant Using Computational Chemistry and Pattern Recognition Techniques.
32. C. Chatfield and A. J. Collins, *Introduction to Multivariate Analysis*, Chapman and Hall, London, 1980, pp. 212–229. Cluster Analysis.
33. P. Willett, *Similarity and Clustering in Chemical Information Systems*, Research Studies Press, Wiley, Chichester, United Kingdom, 1987.
34. D. J. Livingstone, in *Neural Networks in QSAR and Drug Design*, J. Devillers, Ed., Academic Press, London, 1996, pp. 157–176. Multivariate Data Display using Neural Networks.
35. B. Flury and H. Riedwyl, *Multivariate Statistics a Practical Approach*, Chapman and Hall, London, 1988.
36. D. J. Livingstone, *Data Analysis for Chemists*, Oxford University Press, Oxford, United Kingdom, 1995, pp. 126–127.
37. W. R. Dillon and M. Goldstein, *Multivariate Analysis Methods and Applications*, Wiley, New York, 1984, pp. 271–272.
38. J. F. Hair, R. E. Anderson, R. L. Tathan, and W. C. Black, *Multivariate Data Analysis*, Prentice Hall, Englewood Cliffs, New Jersey, 1998, p. 220.
39. A. A. Afifi and V. Clark, *Computer-Aided Multivariate Analysis*, Van Nostrand Reinhold, New York, 1990, p. 162.
40. W. R. Dillon and M. Goldstein, *Multivariate Analysis Methods and Applications*, Wiley, New York, 1984, pp. 272–280.
41. D. A. Belsley, E. Kuh, and R. E. Welsch, *Regression Diagnostics, Identifying Influential Data and Sources of Collinearity*, Wiley, New York, 1980, p. 104.
42. A. Cartier and J.-L. Rivail, *Chemom. Int. Lab. Sys.*, **1**, 335 (1987). Electronic Descriptors in Quantitative Structure-Activity Relationships.
43. D. J. Livingstone and E. Rahr, *Quant. Struct.-Act. Relat.*, **8**, 103–108 (1989). CORCHOP - An Interactive Routine for the Dimension Reduction of Large QSAR Data Sets.
44. B. D. Gute and S. C. Basak, *SAR QSAR Environ. Res.*, **7**, 117 (1997). Use of Topostructural, Topochemical, and Geometric Parameters in the Prediction of Vapor Pressure: A Hierarchical QSAR Approach.
45. D. T. Stanton, *J. Chem. Inf. Comput. Sci.*, **39**, 11 (1999). Evaluation and Use of BCUT Descriptors in QSAR and QSPR Studies.
46. D. C. Whitley, M. G. Ford, and D. J. Livingstone, *J. Chem. Inf. Comput. Sci.*, **40**, 1160 (2000). Unsupervised Forward Selection: A Method for Eliminating Redundant Variables.
47. K. Baumann, H. Albert, and M. von Korff, *J. Chemometrics*, **16**, 339 (2002). A Systematic Evaluation of the Benefits and Hazards of Variable Selection in Latent Variable Regression. Part I. Search Algorithms, Theory and Simulations.

48. A. J. Miller, *J. R. Statist. Soc. A*, **147**, 389 (1984). Selection of Subsets of Regression Variables
49. A. J. Miller, *Subset Selection in Regression*, Chapman and Hall, London, 1990.
50. A. E. Hoerl and R. W. Kennard, *Technometrics*, **12**, 55 (1970). Ridge Regression: Biased Estimation for Non-orthogonal Problems.
51. L. Breimann, *Technometrics*, **37**, 373 (1995). Better Subset Regression Using the Nonnegative Garrote.
52. D. C. Montgomery and E. A. Peck, *Introduction to Linear Regression Analysis*, Wiley, New York, 1982, section 8.5.5.
53. I. S. Helland, *Commun. Statist. Simul. Comput.*, **17**, 581 (1988). On the Structure of Partial Least Squares Regression.
54. M. Stone and R. J. Brooks, *J. R. Statist. Soc. B*, **52**, 237 (1990). Continuum Regression: Cross-Validation Sequentially Constructed Prediction Embracing Ordinary Least Squares, Partial Least Squares and Principal Components Regression
55. J. A. Malpass, D. W. Salt, and M. G. Ford, *Pestic. Sci.*, **46**, 282 (1996). Continuum Regression: Optimised Prediction of Biological Activity.
56. J. A. Malpass, D. W. Salt, E. W. Wynn, M. G. Ford, and D. J. Livingstone in *Trends in QSAR and Molecular Modelling*, C. G. Wermuth, Ed., Escom, Leiden, The Netherlands, 1993, pp. 314–315. Prediction of Biological Activity Using Continuum Regression.
57. J. A. Malpass, D. W. Salt, E. W. Wynn, M. G. Ford, and D. J. Livingstone, in *QSAR: Chemometric Methods in Molecular Design*, Vol. 3, H. van de Waterbeemd, Ed., VCH Publishers, Weinheim, Germany, 1995, pp. 163–189. Continuum Regression: A New Algorithm for the Prediction of Biological Activity.
58. I. E. Frank and J. H. Friedman, *Technometrics*, **35**, 109 (1993). A Statistical View of some Chemometrics Regression Tools.
59. J. F. Lawless and P. Wang, *Commun. In Statist.*, **A5**, 307 (1976). A Simulation Study of Ridge and Other Regression Estimators.
60. G. Smith, *Can. J. Stat.*, **8**, 217 (1980). An Example of Ridge Regression Difficulties.
61. A. E. Hoerl, R. W. Kennard, and K. F. Baldwin, *Commun. In Statist.*, **4**, 105 (1975). Ridge Regression: Some Simulations.
62. N. R. Draper and R. C. Van Nostrand, *Technometrics*, **21**, 451 (1979). Ridge Regression and James-Stein Estimation: Review and Comments.
63. H. Wold, in *Multivariate Analysis*, P. R. Krishnaiah, Ed., Academic Press, New York, 1966. Estimation of Principal Components and Related Models by Iterative Partial Least Squares.
64. C. Bolton-Smith, M. Woodward, W. C. S. Smith, and H. Tunstall-Pedoe, *Int. J. Epidemiol.*, **20**, 95 (1991). Dietary and Non-dietary Predictors of Serum Total and HDL-cholesterol in Men and Women: Results from Scottish Heart Health Study.
65. L. Wilkinson, *Psychological Bulletin*, **86**, 168 (1979). Tests of Significance in Stepwise Regression.
66. C. H. Mason and W. D. Perreault, *J. Market. Res.*, **XXVIII**, 268 (1991). Collinearity, Power, and Interpretation of Multiple Regression Analysis.
67. M. C. Lovell, *Rev. Econom. Statist.*, **LXV**, 1–12 (1983). Data Mining.
68. M. L. Thompson, *Internat. Statist. Rev.*, **46**, 1–19, 129–146 (1978). Selection of Variables in Multiple Regression: Part I. A Review and Evaluation. Part II. Chosen Procedures, Computations and Examples.
69. W. J. Dixon, Chief Ed., *BMDP Statistical Software Manual*, Vol. 1, University of California Press, Los Angeles, California, 1988, p. 373.
70. R. B. Bendel and A. A. Afifi, *J. Am. Stat. Assoc.*, **72**, 46 (1977). Comparison of Stopping Rules in Forward Stepwise Regression.
71. M. A. Efronson, in *Mathematical Methods for Digital Computers*, A. Ralston and H. S. Wilf, Eds., Wiley, New York, 1960. Multiple Regression Analysis.

72. F. Glover, *J. Comput.*, **1**, 4190 (1989). Tabu Search-Part I. ORSA.
73. F. Glover, *J. Comput.*, **2**, 4 (1990). Tabu Search-Part II. ORSA.
74. A. C. Rencher and F. C. Pun, *Technometrics*, **22**, 49 (1980). Inflation of  $R^2$  in Best Subset Regression.
75. L. Wilkinson and G. E. Dallal, *Technometrics*, **23**, 4 (1981). Tests of Significance in Forward Selection Regression with an F-to Enter Stopping Rule.
76. J. B. Copas, *J. R. Statist. Soc. B*, **45**, 311 (1983). Regression, Prediction, and Shrinkage.
77. S. H. McIntyre, D. B. Montgomery, V. Srinivasan, and B. A. Weitz, *J. Market. Res.*, **20**, 1 (1983). Evaluating the Statistical Significance of Models Developed by Stepwise Regression.
78. L. Wilkinson, *Psychologic. Bull.*, **86**, 168 (1979). Tests of Significance in Stepwise Regression.
79. C. L. Mallows, *Technometrics*, **15**, 661 (1973). Some Comments on Cp.
80. A. Golbraikh and A. Tropsha, *J. Molec. Graphics and Modelling*, **20**, 269 (2002). Beware of  $q^2$ .
81. B. F. J. Manly, *Randomization, Bootstrap and Monte Carlo Methods in Biology*, 2nd Ed., Chapman and Hall, London, 1998.
82. I. V. Tetko, A. E. P. Villa, and D. J. Livingstone, *J. Chem. Inf. Comput. Sci.*, **36**, 794 (1996). Neural Network Studies. 2. Variable Selection.
83. D. J. Livingstone, D. T. Manallack, and I. V. Tetko, *J. Comp.-Aided Mol. Design*, **11**, 135 (1997). Data Modeling With Neural Networks—Advantages and Limitations.
84. D. T. Manallack and D. J. Livingstone, *Eur. J. Med. Chem.*, **34**, 195 (1999). Neural Networks in Drug Discovery; Have they Lived up to Their Promise?
85. I. V. Tetko, D. J. Livingstone, and A. I. Luik, *J. Chem. Inf. Comput. Sci.*, **35**, 826 (1995). Neural Network Studies. 1. Comparison of Overfitting and Overtraining.
86. M. G. Ford and D. J. Livingstone, *Quant. Struct.-Act. Relat.*, **9**, 107 (1990). Multivariate Techniques for Parameter Selection and Data Analysis Exemplified by a Study of Pyrethroid Neurotoxicity.
87. D. L. Selwood, D. J. Livingstone, J. C. W. Comley, A. B. O'Dowd, A. T. Hudson, P. Jackson, K. S. Jandu, V. S. Rose, and J. N. Stables, *J. Med. Chem.*, **33**, 136 (1990). Structure-Activity Relationships of Antifilarial Antimycin Analogues, a Multivariate Pattern Recognition Study.
88. H. Kubinyi, *Quant. Struct.-Act. Relat.*, **13**, 285 (1994). Variable Selection in QSAR Studies. I. An Evolutionary Algorithm.
89. J. H. Holland, *Adaptation in Natural and Artificial Systems*, University of Michigan Press, Ann Arbor, Michigan, 1975.
90. J. Devillers, Ed., *Genetic Algorithms in Molecular Modeling*, Academic Press, London, 1996.
91. K. L. Peterson, in *Reviews in Computational Chemistry*, Vol. 16, K. B. Lipkowitz and D. B. Boyd, Eds., Wiley, New York, 2000, pp. 53–140. Artificial Neural Networks and Their Use in Chemistry.
92. R. Leardi, R. Boggia, and M. Terrile, *J. Chemom.*, **6**, 267 (1992). Genetic Algorithms as a Strategy for Feature Selection.
93. R. Leardi, in *Genetic Algorithms in Molecular Modeling*, J. Devillers, Ed., Academic Press, London, 1996, pp. 67–86. Genetic Algorithms in Feature Selection.
94. D. Rogers and A. J. Hopfinger, *J. Chem. Inf. Comput. Sci.*, **34**, 854 (1994). Application of Genetic Function Approximation to Quantitative Structure-Activity Relationships and Quantitative Structure-Property Relationships.
95. H. Kubinyi, *Quant. Struct.-Act. Relat.*, **13**, 285 (1994). Variable Selection in QSAR Studies. II. A Highly Efficient Combination of Systematic Search and Evolution.
96. H. Kubinyi, *J. Chemom.*, **10**, 119 (1996). Evolutionary Variable Selection in Regression and PLS Analyses.
97. J. H. Wikel and E. R. Dow, *Bioorg. Med. Chem. Lett.*, **3**, 645 (1993). The Use of Neural Networks for Variable Selection in QSAR.



98. S.-S. So and M. Karplus, *J. Med. Chem.*, **39**, 1521 (1996). Evolutionary Optimization in Quantitative Structure-Activity Relationship: An Application of Genetic Neural Networks.
99. J. W. McFarland and D. J. Gans, *Quant. Struct.-Act. Relat.*, **13**, 11 (1994). On Identifying Likely Determinants of Biological Activity in High Dimensional QSAR Problems.
100. V. V. Kovalishyn, I. V. Tetko, A. I. Luik, V. V. Kholodovych, A. E. P. Villa, and D. J. Livingstone, *J. Chem. Inf. Comput. Sci.*, **38**, 651 (1998). Neural Network Studies. 3. Variable Selection in the Cascade-Correlation Learning Architecture.
101. O. Nicolotti, V. J. Gillet, P. J. Fleming, and D. V. S. Green, *J. Med. Chem.*, **45**, 6069 (2002). Multiobjective Optimization in Quantitative Structure-Activity Relationships: Deriving Accurate and Interpretable QSARs.
102. S. J. Cho and M. A. Hermsmeier, *J. Chem. Inf. Comput. Sci.*, **42**, 927 (2002). Genetic Algorithm Guided Selection: Variable Selection and Subset Selection.
103. S.-S. Liu, H.-L. Liu, C.-S. Yin, and L.-S. Wang, *J. Chem. Inf. Comput. Sci.*, **43**, 964 (2003). VSMP: A Novel Variable Selection and Modeling Method Based on the Prediction.
104. Y. S. Prabhakar, *QSAR Comb. Sci.*, **22**, 583 (2003). A Combinatorial Approach to the Variable Selection in Multiple Linear Regression: Analysis of Selwood *et al.* Data Set—A Case Study.
105. R. Todeschini, in *Designing Drugs and Crop Protectants: Processes, Problems and Solutions*, M. Ford, D. Livingstone, J. Dearden, and H. van de Waterbeemd, Eds., Blackwell Publishing, Oxford, United Kingdom, 2003, pp. 235–242. Reality and Models. Concepts, Strategies and Tools for QSAR.
106. D. J. Livingstone, in *QSAR and Molecular Modelling: Concepts, Computational Tools and Biological Applications*, J. F. Sanz, J. Giraldo, and F. Manaut, Eds., Prous Science Publishers, Barcelona, Spain, 1995, pp. 18–26. The Trouble with Chemometrics.
107. A. Yasri and D. Hartshough, *J. Chem. Inf. Comput. Sci.*, **41**, 1218 (2001). Toward an Optimal Procedure for Variable Selection and QSAR Model Building.
108. J. K. Wegner and A. Zell, *J. Chem. Inf. Comput. Sci.*, **43**, 1077 (2003). Prediction of Aqueous Solubility and Partition Coefficient Optimized by a Genetic Algorithm Based Descriptor Selection Method.
109. P. Mazzatorta, M. Vracko, and E. Benfenati, *J. Comp.-Aided Mol. Design*, **17**, 335 (2003). ANVAS: Artificial Neural Variables Adaptation System for Descriptor Selection.
110. S.-S. Liu, C.-S. Yin, and L.-S. Wang, *J. Chem. Inf. Comput. Sci.*, **42**, 749 (2002). Combined MEDV-GA-MLR Method for QSAR of Three Panels of Steroids, Dipeptides and COX-2 Inhibitors.
111. K. Hasegawa, T. Kimura, and K. Funatsu, *Quant. Struct.-Act. Relat.*, **18**, 262 (1999). GA Strategy for Variable Selection in QSAR Studies: Enhancement of Comparative Molecular Binding Energy Analysis by GA-Based PLS Method.
112. B. T. Hoffman, T. Kopajtic, J. L. Katz, and A. H. Newman, *J. Med. Chem.*, **43**, 4151 (2000). 2D QSAR Modeling and Preliminary Database Searching for Dopamine Transporter Inhibitors Using Genetic Algorithm Variable Selection of Molconn Z Descriptors.
113. A. Tropsha and W. Zheng, *Curr. Pharm. Des.*, **7**, 599 (2001). Identification of the Descriptor Pharmacophores Using Variable Selection QSAR: Applications to Database Mining.
114. H. X. Liu, R. S. Zhang, X. J. Yao, M. C. Liu, Z. D. Hu, and B. T. Fan, *J. Chem. Inf. Comput. Sci.*, in press. Prediction of the Isoelectric Point of an Amino Acid Based on GA-PLS and SVMs.
115. J. J. Sutherland, L. A. O'Brien, and D. F. Weaver, *J. Chem. Inf. Comput. Sci.*, **43**, 1906 (2003). Spline-Fitting with a Genetic Algorithm: A Method for Developing Classification Structure-Activity Relationships.
116. F. R. Burden, M. G. Ford, D. C. Whitley, and D. A. Winkler, *J. Chem. Inf. Comput. Sci.*, **40**, 1423 (2000). Use of Automatic Relevance Determination in QSAR Studies Using Bayesian Neural Networks.

117. P. Bruneau, *J. Chem. Inf. Comput. Sci.*, **41**, 1605 (2001). Search for Predictive Generic Model of Aqueous Solubility Using Bayesian Neural Nets.
118. Z. Xiao, Y.-D. Xiao, J. Feng, A. Golbraikh, A. Tropsha, and K.-H. Lee, *J. Med. Chem.*, **45**, 2294 (2002). Antitumour Agents. 213. Modeling of Epipodophyllotoxin Derivatives Using Variable Selection  $k$  Nearest Neighbor QSAR Method.
119. M. Shen, A. LeTiran, Y.-D. Xiao, A. Golbraikh, H. Kohn, and A. Tropsha, *J. Med. Chem.*, **45**, 2811 (2002). Quantitative Structure-Activity Relationship Analysis of Functionalized Amino Acid Anticonvulsant Agents Using  $k$  Nearest Neighbor and Simulated Annealing.
120. M. Shen, Y.-D. Xiao, A. Golbraikh, V. Gombur, and A. Tropsha, *J. Med. Chem.*, **46**, 3013 (2003). Development and Validation of  $k$ -Nearest-Neighbor QSPR Models of Metabolic Stability of Drug Candidates.
121. J. W. Godden and J. Bajorath, *QSAR Comb. Sci.*, **22**, 487 (2003). An information-Theoretic Approach to Descriptor Selection for Database Profiling and QSAR Modeling.
122. S. Izrailev and D. K. Agrafiotis, *QSAR Environ. Res.*, **13**, 117 (2002). Variable Selection for QSAR by Artificial Ant Colony Systems.
123. K. Baumann, in *Designing Drugs and Crop Protectants: Processes, Problems and Solutions*, M. Ford, D. Livingstone, J. C. Dearden, and H. van de Waterbeemd, Eds., Blackwell Publishing, Oxford, United Kingdom, 2003, pp. 292–294. Variable Weighting as an Alternative to Variable Selection.
124. D. K. Agrafiotis and W. Cedeno, *J. Med. Chem.*, **45**, 1098 (2002). Feature Selection for Structure-Activity Correlation Using Binary Particle Swarms.
125. K. Heberger and R. Rajko, *QSAR Environ. Res.*, **13**, 541 (2002). Variable Selection Using Pair-correlation Method. Environmental Applications.
126. T.-H. Lin, H.-T. Li, and K.-C. Tsai, *J. Chem. Inf. Comput. Sci.*, In press. Implementing the Fisher's Discriminant Ratio in a  $k$ -Means Clustering Algorithm for Feature Selection and Data Set Trimming.
127. V. Svetnik, A. Liaw, C. Tong, J. C. Culbertson, R. P. Sheridan, and B. P. Feuston, *J. Chem. Inf. Comput. Sci.*, **43**, 1947 (2003). Random Forest: A Classification and Regression Tool for Compound Classification and QSAR Modeling.
128. K. Baumann, N. Stiefl, and M. von Korf, in *Designing Drugs and Crop Protectants: Processes, Problems and Solutions*, M. Ford, D. Livingstone, J. Dearden, and H. van de Waterbeemd, Eds., Blackwell Publishing, Oxford, United Kingdom, 2003, pp. 290–292. Validation of Variable Selection Techniques in QSAR.
129. Special issue. *J. Machine Learning Res.* **3**, 2003.
130. "Getting it Right: Variable Selection and Model Validation in (Q)SAR." Spring National ACS Meeting, Anaheim, California, March 28–April 1, 2004.
131. NIPS 2003 Workshop, Whistler, B.C., Canada, December 11–13, 2004.

# Biomolecular Applications of Poisson–Boltzmann Methods

**Nathan A. Baker**

*Department of Biochemistry and Molecular Biophysics, Center for Computational Biology, Washington University in St. Louis, School of Medicine, 700 S. Euclid Ave., Campus Box 8036, St. Louis, MO, 63110*

---

---

## INTRODUCTION TO BIOMOLECULAR ELECTROSTATICS

Throughout the 1990s, biomolecular simulation has become increasingly commonplace in biology and has gained acceptance as an important biophysical method for understanding molecular structure, dynamics, and function. The energetic properties of a biomolecule are determined by a combination of both short- and long-range forces. Short-range forces include several components, such as van der Waals, bonding forces, angular forces, and torsional interactions. Long-range forces, on the other hand, are typically dominated by electrostatic interactions. Because of their slow decay over distance, electrostatics cannot be neglected or truncated in biomolecular modeling; these forces contribute significantly to molecular interactions at all length scales. As such, methods that enable the accurate and efficient modeling of these interactions are of central importance in molecular simulation.

The exact behavior of electrostatic interactions in a simulation is generally determined by four factors: molecular charge distributions, solute atomic radii, mobile ionic species, and solvent. Molecular charge distributions are

often determined by fitting electronic distributions (as obtained from quantum mechanical calculations) to static point charges located at the atom centers.<sup>1,2</sup> However, several groups have demonstrated that static point charge distributions are not sufficient for the accurate modeling of electrostatics—polarizable models with higher order charge distribution moments are often necessary to reproduce molecular energetics.<sup>3,4</sup> Regardless of the charge model details, the molecular charge distribution is intimately connected with the steric parameters chosen for the solute atoms. Together, charges and radii determine many of the solvation properties of the molecule. Several good references discuss the various issues associated with charge and radius parameterization of a solute.<sup>5–7</sup>

This chapter focuses on the various ways the aqueous environment surrounding a molecule can be treated during a simulation. Specifically, we will examine a popular class of models for ions and solvent around a molecule and the description of electrostatic interactions in the context of these models. Treatments of the ions and solvent around a biomolecule are typically divided into two classes: *explicit* and *implicit*. As their names imply, these models differ by their treatment of solvent and mobile ions as either explicit particles or implicitly through some type of continuum model. In general, explicit models offer the greatest detail and potential for accuracy in molecular simulation. However, explicit solvent and ions often account for 90–95% of the atoms in a simulation and can, therefore, severely increase the computational time required for determining kinetic and thermodynamic properties with any precision. Implicit solvent methods sacrifice the molecular details of the solvent in return for far fewer degrees of freedom to be sampled in the simulation. The result is a substantial decrease in the computational resources required to obtain converged simulation observables. Because of this decrease in computational requirements, implicit solvent methods often enable much better sampling of larger systems than do traditional explicit solvent approaches. However, this increase in length and size of molecular simulations comes at the cost of a simplistic treatment of water and ions that does not perform well under all circumstances.

This chapter will focus on the applications of implicit solvent methods and the tools available for employing these techniques in computational chemistry and biology. Additionally, some potential pitfalls associated with these methods will be described in an effort to help users avoid problems caused by overzealous application. Given its place as a benchmark for many implicit solvent methods, particular emphasis will be placed on the Poisson–Boltzmann (PB) equation. This chapter is designed to supplement the excellent review of PB theory by Gene Lamm, which recently appeared in this series.<sup>8</sup> Readers are strongly encouraged to read Lamm’s chapter as preparation for the present discussion of the algorithms and applications associated with the PB equation.

## SIMPLIFYING THE SYSTEM: IMPLICIT SOLVENT METHODS

As mentioned, implicit solvent methods are approximations designed to simplify the description of the aqueous environment around molecules and thereby reduce the degrees of freedom in a simulation. This section outlines the various methods by which solvent-mediated polar and apolar interactions are described in an implicit solvent setting.

Before discussing the various implicit solvent methods, it is useful to briefly mention the circumstances in which an explicit description of solvent structure is desirable. In general, explicit solvent methods should be used by chemists where the detailed interactions between solvent and solute are important. Some example of such situations include solvent finite size effects in ion channels,<sup>9</sup> strong solvent-solute interactions,<sup>10</sup> strong solvent coordination of ionic species,<sup>11</sup> and saturation of solvent polarization near membranes.<sup>12</sup> Likewise, implicit descriptions of mobile ions are also inappropriate under some circumstances (cf. See Holm et al.<sup>13</sup> for an excellent review), including high ion valency or strong solvent coordination, specific ion-solute interactions, and high local ion densities. Potential users of implicit solvent methods are advised to keep in mind the underlying physics when employing these methods. Although implicit solvent techniques can substantially accelerate biomolecular modeling, the ability to quickly compute the wrong answer does little to help the simulator.

### Polar Interactions

Although numerous methods exist for computing polar interactions in an implicit solvent setting, this section gives only a brief overview of some of the most popular methods available. Interested readers should refer to the reviews by Simonson<sup>14</sup> and Roux<sup>15</sup> for a more comprehensive overview of available implicit solvent methods.

The Debye–Hückel Law:<sup>16</sup>

$$\phi(\mathbf{x}) = \frac{qe^{-\kappa\|\mathbf{x}-\mathbf{x}_0\|}}{\varepsilon\|\mathbf{x}-\mathbf{x}_0\|}$$

provides the simplest description of the electrostatic potential  $\phi(\mathbf{x})$  because of a point charge of magnitude  $q$  located at position  $\mathbf{x}_0$  in a homogeneous polarizable medium of dielectric constant  $\varepsilon$ . The ionic strength of the solution (determined by the concentration of mobile ion species) is represented by the screening parameter

$$\kappa^2 = \frac{8\pi I}{1000\varepsilon RT}$$

where  $I$  is the ionic strength (in molarity, M),  $R$  is the gas constant, and  $T$  is the temperature (in Kelvin). The Debye–Hückel law reduces to Coulomb’s law at zero ionic strength ( $\kappa \rightarrow 0$ ), providing a description of charges at infinite dilution in a polarizable continuum. Debye–Hückel or Coulombic potentials obey superposition principles, i.e., the potential caused by a sum of charges is equivalent to the sum of the potentials because of the isolated charges.

Unfortunately, most biological systems of interest cannot be described as a homogeneous dielectric medium—biomolecular interiors often have significantly lower polarizabilities than their aqueous surroundings. Therefore, biomolecular electrostatics typically cannot be modeled by Debye–Hückel or Coulomb equations. However, because of their simplicity and relative ease of evaluation, such equations are attractive bases for the modeling of biomolecular electrostatics. Therefore, it is not surprising that these equations were the starting points for early algorithms that described the inhomogeneous nature of biomolecular systems.

These simpler early models include distance-dependent dielectric functions,<sup>6,17</sup> reaction field treatments,<sup>18–21</sup> and generalized Born (GB) models.<sup>22–27</sup> Of these, GB models are arguably the most popular. GB methods were introduced by Still et al. in 1990<sup>23</sup> and have been progressively refined by several other researchers.<sup>24–27</sup> GB methods are based on the Born ion, a canonical electrostatics model describing the electrostatic potential and solvation energy of a spherical ion<sup>28</sup> (see also the later section on solvation-free energies). With the GB method, we use an analytical expression based on the Born ion model to approximate the electrostatic potential and solvation energy of small molecules. Although it fails to capture all of the details of molecular structure and ion distributions provided by more rigorous models,<sup>27,29–31</sup> such as the Poisson–Boltzmann equation, it has gained popularity and continues to be vigorously developed, as it is a very rapid method for evaluating approximate forces and energies for solvated molecules. Recently, an excellent critical comparison of GB and PB methods was presented by Feig et al.;<sup>32</sup> interested readers should refer to this review for relative speeds and accuracies of GB and PB techniques.

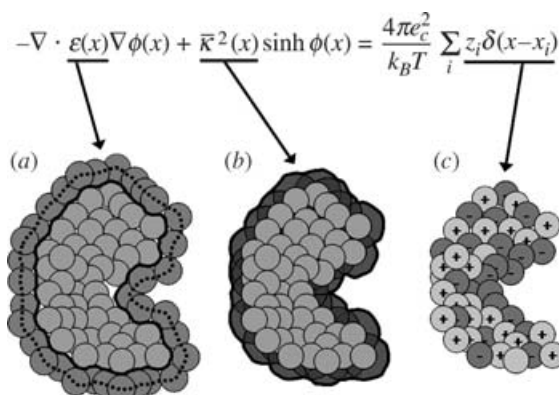
Poisson–Boltzmann methods offer a compromise between faster, but more approximate models such as GB, and more detailed explicit solvent and integral equation techniques.<sup>33</sup> The remainder of this chapter is devoted to discussion of the implementation and application of PB and related methods.

## Nonpolar Interactions

An important aspect of implicit solvent models is their ability to treat apolar energetics originating from solvent-mediated interactions. In general, apolar methods have been developed separately from their polar counterparts; none of the methods described in the previous section requires a specific apolar

treatment. The apolar term plays an important role in electrostatic force calculations; solvation energies/forces obtained from more detailed implicit solvent methods (PB and GB) calculations work to maximize the solvent–solute boundary surface area, thereby providing the maximum solvation for the biomolecule. Apolar interactions, on the other hand, tend to bias the system toward conformations with minimum surface area. The net effect of the polar and apolar solvation terms is a balance of opposing forces determined by the details of the molecular structure.

Apolar interactions are generally modeled by either scaled-particle theory methods<sup>34,35</sup> or solvent-accessible surface area techniques,<sup>36–40</sup> although more rigorous methods have been suggested.<sup>41–43</sup> Scaled-particle theories (SPTs) provide apolar interactions in terms of the work required to insert a spherical cavity into a liquid; the resulting energies are usually given as a function of the solute radius and surface area. Solvent-accessible surface area-based models (see Figure 1) are, in essence, a simpler version of scaled-particle theories. Solvent-accessible surface area (SASA) methods relate apolar energies and forces by one or more “surface tension” constants of proportionality that are either chosen globally<sup>36–38</sup> or for specific atom types.<sup>39,40</sup> However, because of the delicate balance between polar and apolar interactions and their impact on the results of the simulation, the choice of surface tension constants



**Figure 1** Description of the elements of the Poisson–Boltzmann equation for biomolecular geometry. (a) The dielectric coefficient  $\epsilon(x)$  is described in terms of the “molecular surface” (solid line) defined by the union of the surfaces of solvent-sized spherical probes (light gray circles) arranged outside the biomolecular interior (dark gray circles). The solvent-accessible surface (dashed line) determines surface-area based contributions to the apolar solvation energy; it is defined by the centers, rather than the peripheries, of the solvent probes. (b) The ion-accessibility parameter  $\bar{\kappa}^2(x)$  is proportional to the bulk ionic strength outside the ion-accessible surface (black line) defined by the union of biomolecule atoms (light gray circles) with radii increased by the radii of the counterion species (dark gray circles). (c) The biomolecular charge distribution is defined as the collection of point charges located at the centers of biomolecule atoms.

for SASA methods is still more art than science and often depends on the particular system under consideration.

---

## POISSON–BOLTZMANN THEORY: A BRIEF OVERVIEW

---

This section provides a brief overview of PB theory; however, the interested reader is urged to read the much more comprehensive treatment provided by Gene Lamm in vol. 19 of this series.<sup>8</sup> As mentioned, the Poisson–Boltzmann equation is obtained from a continuum description of the solvent and counterions surrounding a biomolecule.<sup>13,44–46</sup> Although numerous derivations of the PB equation are based on statistical mechanics,<sup>13,33</sup> the simplest begins with Poisson’s equation:<sup>47</sup>

$$-\nabla \cdot \varepsilon(x) \nabla \varphi(x) = \rho(x) \quad [1]$$

for  $x \in \Omega$ , the basic equation for describing the electrostatic potential  $\varphi(x)$  generated by a charge distribution  $\rho(x)$  in a continuum model of a polarizable solvent with dielectric constant  $\varepsilon(x)$ . This equation is generally solved in a finite domain  $\Omega$  with the potential specified as  $g(x)$  on the domain boundary  $\partial\Omega$ :

$$\varphi(x) = g(x) \quad [1b]$$

for  $x \in \partial\Omega$ . This Dirichlet boundary condition usually employs an analytic, asymptotically correct form of the potential (Coulomb’s law or Debye–Hückel) for  $g(x)$ . Therefore, to ensure the applicability of the boundary condition, the domain must be sufficiently large and the boundaries reasonably distant (usually a few Debye or Bjerrum lengths) from the biomolecule.

To obtain the PB equation from the Poisson equation, we need to consider the two types of charge distributions present in biomolecular systems. First, the partial atomic charges are usually modeled as a “fixed” charge distribution:

$$\rho_f(x) = \frac{4\pi e_c^2}{kT} \sum_{i=1}^M Q_i \delta(x - x_i) \quad [2]$$

which describes the  $M$  atomic partial charges of the biomolecule as delta functions  $\delta(x - x_i)$  located at the atom centers  $\{x_i\}$  with magnitudes  $\{Q_i\}$ . The scaling coefficients ensure the dimensionless form of the potential and include  $e_c$ , the charge of an electron, and  $kT$ , the thermal energy of the system. Additionally, the contributions of counterions are modeled as a continuous “charge



cloud” described by a Boltzmann distribution, giving rise to the “mobile” charge distribution

$$\rho_m(x) = \frac{4\pi e_c^2}{kT} \sum_{j=1}^m c_j q_j \exp[-q_j \varphi(x) - V_j(x)] \quad [3]$$

for  $m$  counterion species with charges  $\{q_j\}$ , bulk concentrations  $\{c_j\}$ , and steric potentials  $\{V_j\}$  (i.e., potentials that prevent biomolecule-counterion overlap). In the case of a monovalent electrolyte such as NaCl, Eq. [3] reduces to

$$\rho_m(x) = -\bar{\kappa}^2(x) \sinh \varphi(x) \quad [4]$$

where the coefficient  $\bar{\kappa}^2(x)$  describes both ion accessibility (via  $\exp[-V(x)]$ ) and bulk ionic strength. Combining the expressions for the charge distributions (Eqs. [2] and [4]) with Poisson’s equation (Eq. [1]) gives the Poisson–Boltzmann equation for a 1-to-1 electrolyte:

$$-\nabla \cdot \varepsilon(x) \nabla \varphi(x) + \bar{\kappa}^2(x) \sinh \varphi(x) = \frac{4\pi e_c^2}{kT} \sum_{i=1}^M Q_i \delta(x - x_i) \quad [5]$$

for  $x \in \Omega$ , where  $\varphi(x) = g(x)$  for  $x \in \partial\Omega$ . As described, details of the biomolecular structure enter into the coefficients of the PB equation (see Figure 1). The dielectric function  $\varepsilon(x)$  has been represented by a variety of models,<sup>48–52</sup> which typically involve a relatively abrupt change in the dielectric coefficient near the molecular surface.

The “full” or nonlinear form of the problem given in Eq. [5] is often simplified to the linearized PB equation by replacing the  $\sinh \varphi(x)$  term with its first-order approximation,  $\sinh \varphi(x) \approx \varphi(x)$ , to give:

$$-\nabla \cdot \varepsilon(x) \nabla \varphi(x) + \bar{\kappa}^2(x) \varphi(x) = \frac{4\pi e_c^2}{kT} \sum_{i=1}^M Q_i \delta(x - x_i) \quad [6]$$

with the same boundary conditions as Eq. [5]. However, this linearization is only appropriate at small potential values where the nonlinear contributions to Eq [5] (the  $\sinh$  term) are negligible.

Solutions to the PB equation generally calculate either energies or forces. Electrostatic free energies are obtained by integration of solutions to the PB equation over the domain of interest:<sup>53–55</sup>

$$G[\varphi] = \int_{\Omega} \left[ \rho_f \varphi - \frac{\varepsilon}{2} (\nabla \varepsilon)^2 - \bar{\kappa}^2 (\cosh \varphi - 1) \right] dx \quad [7]$$

The first term of Eq. [7] is the energy of inserting the protein charges into the electrostatic potential and can be interpreted as the interaction energy between charges. However, unlike analytic representations of charge–charge interactions, this energy also includes large “self-energy” terms associated with the interaction of a particular charge with itself. These self-energy terms are highly dependent on the discretization of the problem; as the mesh spacing increases, these terms become larger. In general, self-energies are treated as artifacts of the calculation and are removed by a reference calculation with the same discretization (see the section on applications). The second term of Eq. [7] can be interpreted as the energy of polarization for the dielectric medium. Finally, the third term is the energy of the mobile counterion distribution. Both the dielectric and mobile ion energies do not include self-energy terms and therefore do not need to be corrected by reference calculations.

As with the PB equation, the energy can be linearized by noting  $\cosh \varphi(x) - 1 \approx \varphi^2(x)/2$  for small  $\varphi(x)$  to give the simplified expression. A force expression can be obtained by differentiating the energy functional with respect to the atomic coordinates:<sup>50,56</sup>

$$\mathbf{F}_i = - \int_{\Omega} \left[ \varphi \left( \frac{\partial \rho_f}{\partial \mathbf{y}_i} \right) - \frac{(\nabla \varphi)^2}{2} \left( \frac{\partial \varepsilon}{\partial \mathbf{y}_i} \right) - (\cosh \varphi - 1) \left( \frac{\partial \bar{\kappa}^2}{\partial \mathbf{y}_i} \right) \right] d\mathbf{x} \quad [8]$$

where  $\mathbf{F}_i$  is the force on atom  $i$  and  $\partial/\partial \mathbf{y}_i$  denotes the derivative with respect to displacements of atom  $i$ . The first term of Eq. [8] represents the force density for displacements of atom  $i$  in the potential; it can also be rewritten in the classic form  $q_i \nabla \varphi$  for a charged particle in an electrostatic field. The second term is the dielectric boundary pressure; i.e., the force exerted on the biomolecule by the high dielectric solvent surrounding the low-dielectric interior. Finally, the third term is equivalent to the osmotic pressure or the force exerted on the biomolecule by the surrounding counterions.

It is worthwhile noting that the PB equation (nonlinear or linearized) is an approximate theory and therefore cannot be applied to all biomolecular systems. Specifically, the PB equation is derived from mean-field or saddle-point treatments of the electrolyte system and therefore neglects counterion correlations and fluctuations that can affect the energetics of highly charged biomolecular systems such as DNA, RNA, and some protein systems. A good review of the scope and impact of these deviations from PB theory can be found in the chapter by Lamm<sup>8</sup> and the text of Holm et al.<sup>13</sup> In short, Poisson–Boltzmann theory gives reasonable quantitative results for biomolecules with low linear charge density in monovalent symmetric salt solutions; however, PB theory can be *qualitatively incorrect* for highly charged biomolecules or more concentrated multivalent solutions. Therefore, application of PB theory and software requires some discretion on the part of the user.

## SOLVING THE POISSON–BOLTZMANN EQUATION

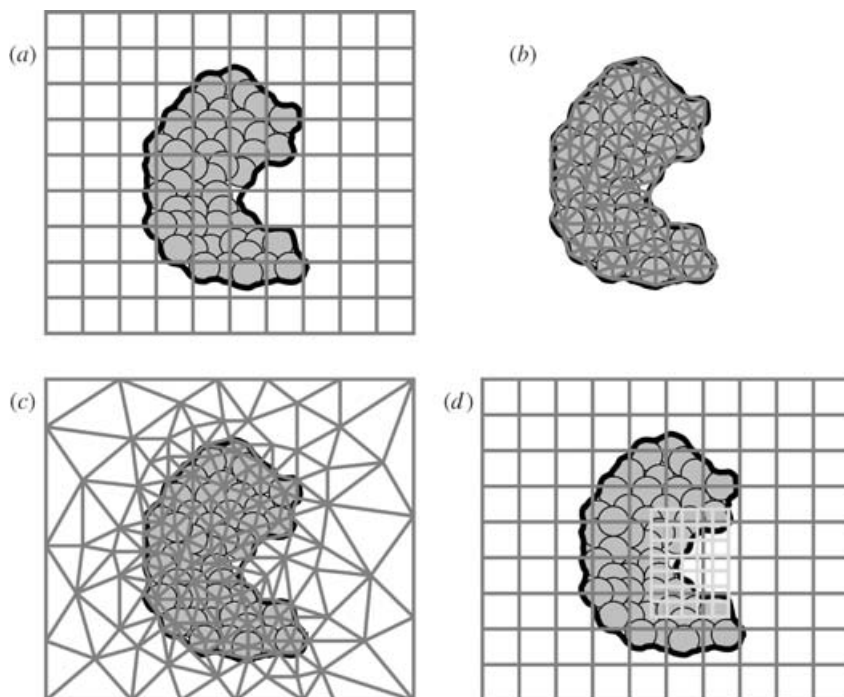
Because of the complicated nature of biomolecular geometries and charge distributions, the PB equation (PBE) is usually solved numerically by a variety of computational methods. These methods typically discretize the (exact) continuous solution to the PBE via a finite-dimensional set of basis functions. In the case of the linearized PBE, the resulting discretized equations transform the partial differential equation into a linear matrix-vector form that can be solved directly. However, the nonlinear equations obtained from the full PBE require more specialized techniques, such as Newton methods, to determine the solution to the discretized algebraic equation.<sup>57</sup>

### Discretization Methods

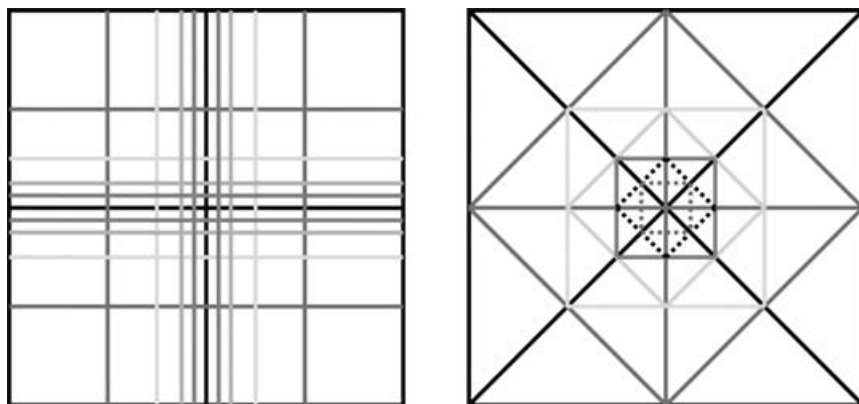
The three most popular discretization methods, finite difference, boundary element, and finite element, are shown in Figure 2(a)–(c). Some of the most popular discretization techniques employ Cartesian meshes to subdivide the domain over which the PB equation is to be solved. Of these, the finite difference (FD) method has been at the forefront of PBE solvers (see Hanig and Nicholls,<sup>46</sup> Davis et al.,<sup>58</sup> Rocchia et al.,<sup>59</sup> and Davis and McCammon<sup>60</sup> and references therein). In its most general form, the finite difference method solves the PBE on a nonuniform Cartesian mesh, as shown in Figure 2(a) for a two-dimensional domain. In this general setting, the differential operator in the PBE ( $-\nabla \cdot \epsilon \nabla$ ) is transformed into a sparse difference matrix by means of a Taylor expansion. The resulting matrix equations are then solved by a variety of matrix algebra techniques. Although FD grids offer relatively simple problem setup, they provide little control over how unknowns are placed in the solution domain. Specifically, as shown by Figure 3, the Cartesian or tensor-product nature of the mesh makes it impossible to locally increase the accuracy of the solution in a specific region without increasing the number of unknowns across the entire grid.

Finite difference methods can employ a unique method known as “electrostatic focusing”<sup>61,62</sup> to provide a limited degree of adaptivity in their calculations. As illustrated in Figure 2(d), focusing enables users to apply relatively coarse FD grids for calculations and much finer grids in regions of interest (binding or active sites, titratable residues, etc.). Specifically, a coarse grid calculation is performed over the entire problem domain, and the resulting solution provides boundary conditions for a much higher resolution calculation on the subdomain of interest. The result is a highly accurate local solution to the PB equation with a reduced amount of computational effort.

Generic simplicial discretizations offer a much more flexible alternative to Cartesian mesh finite difference techniques. Boundary element (BE)



**Figure 2** Popular discretization schemes for numerical solution of the Poisson–Boltzmann equation. The solid black line and circles denote a model protein; other lines denote the mesh on which the system is discretized. (a) Finite difference. (b) Boundary element. (c) Finite Element. (d) Focusing on finite difference grids. See color insert.



**Figure 3** Adaptive refinement for finite different (left) and finite element (right) methods. Shading denotes successive levels of refinement. See color insert.

methods [see Figure 2(b)] discretize the surface of the molecule with triangular simplices.<sup>63,64</sup> The solution is represented for an induced surface charge, which is then convolved with the Coulomb potential (Green's function) to give the desired solution. The result is a highly detailed description of the macromolecular geometry with a very small number of unknowns, thereby offering a highly efficient computational method. Unfortunately, BE methods are only applicable to the linearized PB equation, thereby limiting their overall use by scientists in biomolecular electrostatics. Like BE methods, finite element (FE) discretizations<sup>65–67</sup> also offer the ability to place computational effort in specific regions of the problem domain. Finite element meshes [see Figure 2(c)] are composed of simplices (e.g., triangles or tetrahedra) that span the entire volume in which the PB equation is to be solved. The electrostatic potential is constructed from piecewise polynomial basis functions that are associated with mesh vertices and typically are nonzero only over a small set of neighboring simplices. Solution accuracy can be increased in specific areas by locally increasing the number of vertices through simplex refinement. As shown in Figure 3, the number of unknowns (vertices) is generally increased only in the immediate vicinity of the simplex refinement and not throughout the entire problem domain, as in FD methods. This ability to locally increase the solution resolution is called “adaptivity” and is the major strength of finite element methods applied to the PB equation.<sup>68–72</sup> As with the FD method, this discretization scheme leads to sparse symmetric matrices with a few nonzero entries in each row.<sup>65–67</sup>

## Multilevel Solvers

Multilevel solvers<sup>73,74</sup> provide the most efficient solution of the algebraic equations obtained by discretization of the PBE with either finite difference or finite element techniques.<sup>44,57,68,70,75,76</sup> Most matrix equations are solved by iterative methods which start with an initial guess and repeatedly apply a set of operations to improve this guess until a solution of the desired accuracy is reached. However, the speed of traditional iterative methods has been limited by their inability to quickly reduce long-range error in the solution.<sup>73,74</sup> This problem can be overcome by projecting the discretized system onto meshes (or grids) at multiple resolutions. This projection quickly reduces the error in the slowly converging low-frequency components of the solution on the finest mesh via solutions on the coarser levels of the system. The coupling of scales gives rise to a “multilevel” solver algorithm, in which the algebraic system is solved directly on the coarsest level and then used by scientists to accelerate solutions on finer levels of the mesh.

The assembly of the multilevel hierarchy depends on the method that discretizes the PBE. For FD types of methods, we use so-called “multigrid” methods: The nature of the FD grid lends itself to the assembly of a hierarchy with little additional work.<sup>57,73,74</sup> In the case of adaptive finite element

discretizations, “algebraic multigrid” methods<sup>70,75</sup> are employed. For FE meshes, the most natural multiscale representation is constructed by refinement of an initial mesh that typically constitutes the coarsest level of the hierarchy (see Figure 3).

## Parallel Methods

Regardless of the scalability of the numerical algorithm that solves the PB equation, some systems are simply too large to be solved sequentially (i.e., on one processor). For example, although small- to medium-sized protein systems (100–1000 residues) are amenable to sequential calculations, an increasing interest exists in macromolecular assemblages with tens to hundreds of thousands of residues (e.g., microtubules, entire viral capsids, ribosomes, polymerases, etc.). Studies of these large systems are not feasible on most sequential platforms; instead, they require multiprocessor computing platforms to solve the PB equation in a parallel fashion. However, recent progress in parallel solution methods<sup>75,77</sup> has extended the applicability of PB theory to these large biomolecular systems. Parallel methods have been applied to both finite element<sup>70</sup> and finite difference<sup>76</sup> discretizations of the PB equation.

## Software for Computational Electrostatics

Table 1 presents a list of the major software that currently solves the Poisson–Boltzmann equation for biomolecular systems. A variety of such programs exist, ranging from multipurpose computational biology packages (e.g., CHARMM, Jaguar, UHBD, and MacroDox) to specialized PB solvers (e.g., APBS, MEAD, and DelPhi).

In addition to the traditional “stand-alone” software packages listed in Table 1, web-based services are becoming available for solving the PB equation. Such services have the advantage of removing the troublesome details of software installation from the user. Quite often, web services also simplify and/or automate the process of setting up calculations. One example of such web-based packages is the APBS Web Portal (<https://gridport.npaci.edu/apbs/>), a service for preparing, submitting, and organizing electrostatics calculations on supercomputing platforms. This web portal is aimed at *quantitative* electrostatics calculations and currently does not offer substantial analysis or visualization options. On the other hand, the GRASP (<http://trantor.bioc.columbia.edu/>) web service supports *qualitative* electrostatics calculations with extensive visualization capabilities. GRASP allows users to easily calculate and visualize electrostatic potentials and other properties with uploaded structures or PDB entries. Finally, the MolSurfer (<http://projects.villa-bosch.de/dbase/molsurfer/>) service also offers qualitative analysis of the electrostatic properties of biomolecules with specific attention to surface-based quantities.

**Table 1** Poisson–Boltzmann Software for Biomolecular Systems

Software Package	Description	URL	Availability
APBS <sup>76,149</sup>	Solves PBE in parallel with FD MG and FE AMG solvers.	<a href="http://agave.wustl.edu/apbs/">http://agave.wustl.edu/apbs/</a>	Windows, All Unix. Free, open source.
DelPhi <sup>59,150</sup>	Solves PBE sequentially with highly optimized FD GS solver.	<a href="http://trantor.bioc.columbia.edu/delphi/">http://trantor.bioc.columbia.edu/delphi/</a>	SGI, Linux, AIX. \$250 academic.
GRASP <sup>151</sup>	Visualization program with emphasis on graphics; offers sequential calculation of qualitative PB potentials.	<a href="http://trantor.bioc.columbia.edu/grasp/">http://trantor.bioc.columbia.edu/grasp/</a>	SGI. \$500 academic.
MEAD <sup>152</sup>	Solves PBE sequentially with FD SOR solver.	<a href="http://www.scripps.edu/bashford">http://www.scripps.edu/bashford</a>	Windows, All Unix. Free, open source.
UHBD <sup>58,153</sup>	Multipurpose program with emphasis on SD; offers sequential FD SOR PBE solver.	<a href="http://mccammon.ucsd.edu/uhbd.html">http://mccammon.ucsd.edu/uhbd.html</a>	All Unix. \$300 academic.
MacroDox	Multipurpose program with emphasis on SD; offers sequential FD SOR PBE solver.	<a href="http://pirn.chem.tntech.edu/macrodex.html">http://pirn.chem.tntech.edu/macrodex.html</a>	SGI. Free, open source.
Jaguar <sup>71,72,154</sup>	Multipurpose program with emphasis on QM; offers sequential FE MG, SOR, and CG solvers.	<a href="http://www.schrodinger.com/Products/jaguar.html">http://www.schrodinger.com/Products/jaguar.html</a>	Most Unix. Commercial.
CHARMM <sup>155,156</sup>	Multipurpose program with emphasis on MD; offers sequential FD MG solver and can be linked with APBS.	<a href="http://yuri.harvard.edu">http://yuri.harvard.edu</a>	All Unix. \$600 academic.

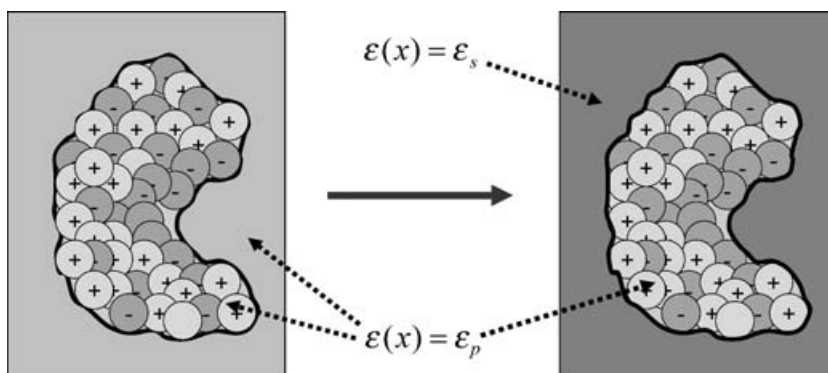
Poisson–Boltzmann equation (PBE), multigrid (MG), algebraic multigrid (AMG), finite difference (FD), finite element (FE), Gauss–Seidel (GS), conjugate gradient (CG), successive over relaxation (SOR), stochastic dynamics (SD), quantum mechanics (QM), molecular mechanics (MM), molecular dynamics (MD).

## APPLYING POISSON–BOLTZMANN METHODS

The purpose of this section is to illustrate the application of PB methods to various molecular and biomolecular problems. Where feasible, the problems are described in sufficient detail to allow readers to further investigate these systems with any of the software packages listed in Table 1. Additionally, these examples are included in the APBS software package distribution (<http://agave.wustl.edu/apbs/>)<sup>76</sup> and are freely available for download.

### Solvation Free Energies

One of the simplest PB calculations is the solvation-free energy. These types of problems consider the energy of transferring a solute from a uniform dielectric continuum (of permittivity  $\epsilon_p$ ) to an inhomogeneous medium with bulk dielectric ( $\epsilon_s$ ) equal to that of the solvent (see Figure 4). To evaluate solvation energies, two PB calculations are performed: (1) Calculate the energy  $G_1$  of the system with a constant dielectric  $\epsilon_p$ , and (2) calculate the energy  $G_2$  of the system with solute dielectric  $\epsilon_p$  and solvent dielectric  $\epsilon_s$ . The solvation energy is then the difference of these two calculations:  $\Delta G_{solv} = G_2 - G_1$ . One aspect of this procedure that should immediately be questioned is Step (1). If we are using a homogeneous dielectric, why do we need to perform a numerical PB calculation; why not just use Coulomb's law? The reasons for this apparently superfluous calculation are twofold. First, PB energies include “self-interaction” terms, i.e., the energy of a charge distribution interacting with itself. For point charges, an exact solution to the PB equation should give infinite self-interaction terms. The finite (but large!) nature of these



**Figure 4** Schematic of a solvation energy calculation. The initial state treats the solute in a homogeneous dielectric material with both solvent and solute dielectric coefficients set to solute value  $\epsilon_p$ . The final state involves an inhomogeneous dielectric coefficient with solute value  $\epsilon_p$  and bulk solvent value  $\epsilon_s$ .



components for numerical solutions is caused by the finite discretization applied in FD and FE methods. Second, these terms are extremely sensitive to the discretization scheme that solves the PB equation. In other words, small changes in grid spacing or the location of these charges on the grid gives rise to very large changes in the self-interaction terms. For these reasons, self-interactions are always removed from PB calculations—*energies obtained from a single PB calculation are meaningless*.

The Born ion<sup>28</sup> is the simplest model of solvation: It considers the solvation energy of a spherical, nonpolarizable ( $\epsilon_p = 1$ ) solute of radius  $R$  with a single point charge of magnitude  $z$  at its center. In this case, an analytical expression is available for the solvation energy:

$$\Delta G_{Born} = -\frac{z^2 e_c^2}{8\pi\epsilon_0 R} \left(1 - \frac{1}{\epsilon_e}\right) \quad [9]$$

This simple model reproduces many of the basic characteristics we would expect from ion solvation; the solvation energy decreases with the ion size, increases with the charge magnitude, and is roughly independent of the charge sign. In fact, the Born model performs very well at describing the solvation properties of low charge-density ions (large and/or monovalent) but fails for higher charge systems where dielectric saturation and solvent electrostriction become significant.

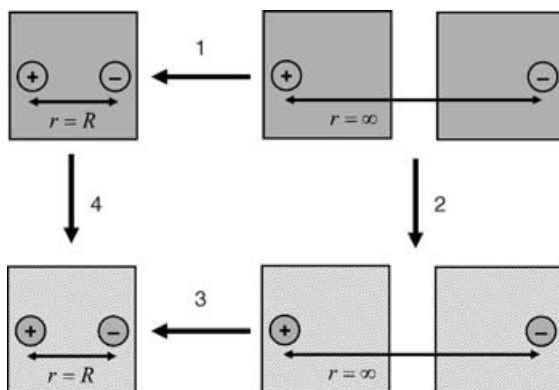
A somewhat more interesting problem is the implicit solvent potential of mean force (PMF) between two Born ions. Unlike the single ion solvation, no closed-form analytic solution is available for this system. Instead, it must be either modeled with one of the numerous published series solutions<sup>78–80</sup> or with numerical calculations. Figure 5 presents a thermodynamic cycle illustrating a typical PB calculation of an ion–ion PMF. The energy of bringing the ion to a distance  $R$  in solution is

$$\Delta G_3(R) = \Delta G_4(R) + \Delta G_1(R) - \Delta G_2 \quad [10]$$

Let the solute dielectric be  $\epsilon_p$  and the solvent dielectric be  $\epsilon_s$ ; then these quantities are obtained from the following types of Poisson calculations:

- $\Delta G_1(R)$  is the energy of bringing the ions to a distance  $R$  in a homogeneous medium of constant dielectric  $\epsilon_p$ .
- $\Delta G_2$  is the solvation energy of the isolated ions (i.e., at infinite distance). This is the energy of transferring each ion from a homogeneous dielectric  $\epsilon_p$  to an inhomogeneous dielectric with constants  $\epsilon_p$  and  $\epsilon_s$ .
- $\Delta G_4(R)$  is the solvation energy of the ion complex (i.e., at distance  $R$ ).

Therefore, Eq. [10] indicates that the energy of forming the complex at distance  $R$  in solution ( $\Delta G_3(R)$ ) is simply equal to the change in solvation

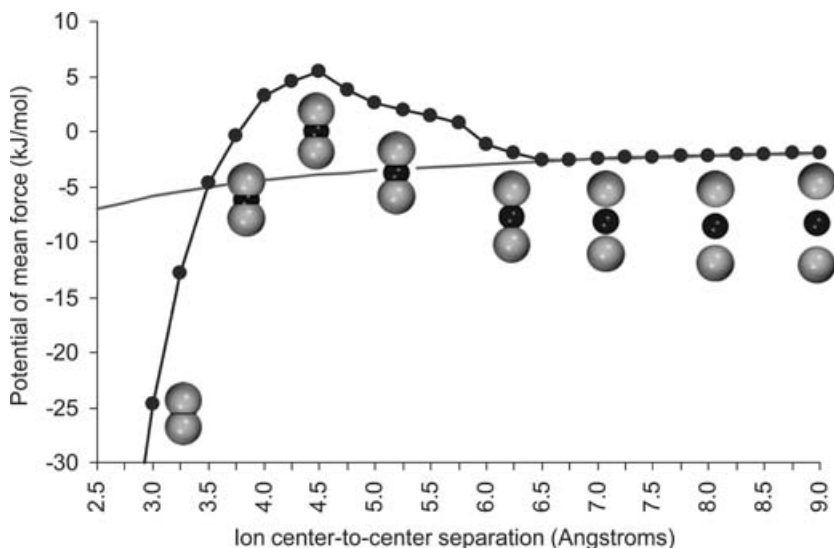


**Figure 5** Thermodynamic cycle illustrating the standard numerical procedure for calculating the energy of two solutes at distance  $R$ . When  $R$  is varied, this method can construct a potential of mean force. The steps are (1) association of solutes in a medium with a homogeneous dielectric, (2) transfer of the isolated solutes from a homogeneous dielectric into solution with an inhomogeneous dielectric, (3) association of the solutes in a medium with an inhomogeneous dielectric, and (4) transfer of the associated solutes from a homogeneous dielectric into solution with an inhomogeneous dielectric.

energy ( $\Delta G_4(R) - \Delta G_2$ )—or desolvation energy—plus the Coulombic interaction energy of the solutes in a homogeneous dielectric ( $\Delta G_1(R)$ ).

Figure 6 illustrates the application of these calculations to a pair of non-polarizable, oppositely charged, monovalent ions in a solvent of dielectric 78.54 (water). Clearly the relative energetic contributions of desolvation and simple Coulombic attraction strike a delicate balance. At large separations, desolvation is negligible and the PMF closely follows Coulomb's law for two point charges in a homogeneous dielectric of 78.54. However, on the length scale of the solvent molecule, the PMF deviates significantly from Coulomb's. For this system of oppositely charged ions, the potential even becomes net repulsive at small separations. By examining the ion–solvent–ion configuration schematics in Figure 6, we can see that the onset of the unfavorable desolvation contribution corresponds with separations that are smaller than  $(\sigma_1 + \sigma_2 + 2\sigma_s)/2$ , where  $\sigma_1$  is the radius of ion 1,  $\sigma_2$  is the radius of ion 2, and  $\sigma_s$  is the radius of the solvent molecules. In other words, solvation energies become unfavorable when solvent is “squeezed” out of charged interfaces.

The solvation and desolvation methods and phenomena described in this section are relatively generic and are applied in a variety of electrostatics calculations. In particular, desolvation plays an important role in biomolecular electrostatics and, as we will see in the following sections, must be treated with care.



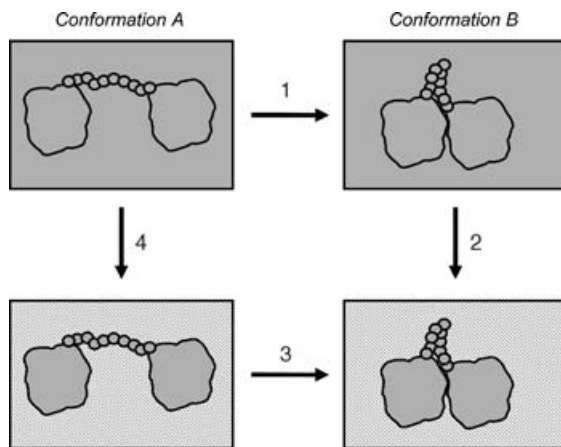
**Figure 6** The potential of mean force for two oppositely charged monovalent ions as obtained from Poisson's equation with a solvent dielectric of 78.54 and a solute dielectric of 1.0. The black line is the Coulomb's law interaction potential in a medium of constant solvent dielectric. The line with data points depicts the energies obtained from numerical solutions of Poisson's equation. The spheres depict the relative separation of the ions (large spheres) with respect to a model solvent molecule (small sphere); note the increase in energy as the separation falls below the size of the solvent probe.

## Conformational Free Energies

One common pitfall in PB applications develops when it computes electrostatic energy differences associated with conformational changes. Such applications appear in a variety of contexts, many of which are described in a later section. However, the particular caveats for conformational free energy evaluation warrant their own section.

As mentioned, energies computed from the PB equation (see Eq. [8]) contain self-interaction terms that are very large and highly sensitive to the discretization of the equation. As conformational changes involve movement of charges and dielectric boundaries on the FD or FE mesh, these self-energies are different for each conformation and need to be explicitly removed by reference calculations. Figure 7 depicts the series of calculations needed to accurately calculate the electrostatic contribution to the free energy change going between conformations A and B:

$$\Delta G_{A \rightarrow B} = \Delta G_3 = \Delta G_1 + \Delta G_2 - \Delta G_4 \quad [11]$$



**Figure 7** Thermodynamic cycle illustrating the standard numerical procedure for calculating the electrostatic energy of a conformational change in a molecule. The steps are (1) conformational change in a homogeneous dielectric, (2) transfer of conformation *B* from a homogeneous dielectric into solution with an inhomogeneous dielectric, (3) conformational change in an inhomogeneous dielectric, and (4) transfer of conformation *A* from a homogeneous dielectric into solution with an inhomogeneous dielectric.

The self-energies are explicitly removed in the calculation of the change in solvation energy ( $\Delta G_4 - \Delta G_2$ ) associated with the conformational change. As with the ion-ion PMF example, solvation energies strive to expose polar surfaces and generally oppose the association of any charged molecular surfaces. For the toy example shown in Figure 7, polar solvation forces would be expected to favor conformation *A*. Finally, the remaining term in Eq. [11] ( $\Delta G_1$ ) denotes the usual Coulombic contribution to the electrostatic energy; it is the change in internal electrostatic energy (for a homogeneous dielectric of  $\epsilon_p$ ) caused by the conformational change.

Unlike the ion-ion PMF case, substantial apolar surface area is often buried during biomolecular conformational changes and/or complex formation. Therefore, we must also consider the apolar solvation energies discussed in a previous section. For the current example, a constant apolar coefficient method would provide an additional term for the conformational change:

$$\Delta G^{apolar} = \gamma \Delta A \quad [12]$$

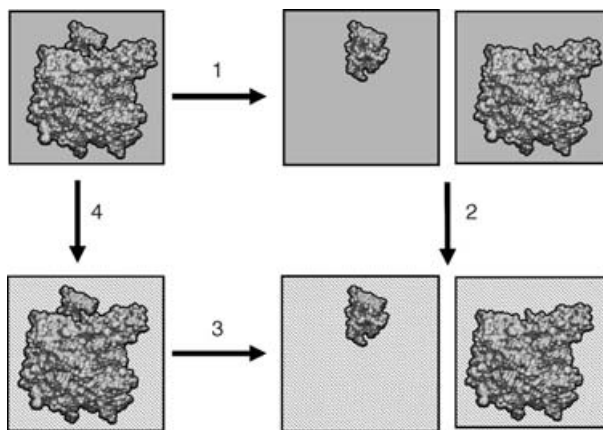
where  $\gamma$  is the apolar coefficient and  $\Delta A$  is the difference in surface area caused by the conformational change. As  $\gamma \geq 0$ , this term favors the burial of molecular surfaces and therefore counteracts the polar solvation forces. In practice, apolar and polar energy terms strike a tenuous balance in stabilizing biomolecular complexes and conformations.

## Binding Free Energies

The calculation of binding free energies is arguably the most common application of PB methods. As described earlier, PB methods provide the electrostatic energies and are often supplemented with additional terms representing apolar and other interactions. These calculations are performed on a wide variety of complexes ranging from small molecule–protein interactions<sup>81</sup> to protein–membrane interactions<sup>82</sup> and the energy of supramolecular assembly.<sup>83</sup>

The simplest type of binding energy calculation treats the associating molecules as rigid entities, i.e., with no conformational change upon binding. This description of the assembly process is clearly not terribly realistic. However, it can often successfully provide relative affinities for families of compounds based on a single starting structure.<sup>81,84</sup> As with the other examples, our discussion of PB methods for binding free energies starts with the thermodynamic cycle shown in Figure 8, which computes the dissociation constant for a complex between mouse acetylcholinesterase (mAChE) and fasciculin-2 (FAS2). The polar (electrostatic) free energy for dissociation of the mAChE-FAS2 complex is given by

$$\Delta G_3^p = \Delta G_1^p + \Delta G_2^p - \Delta G_4^p \quad [13]$$



**Figure 8** Thermodynamic cycle illustrating the numerical procedure that calculates the rigid-body electrostatic contribution to the dissociation energy of a complex between mouse acetylcholinesterase (large molecule) and fasciculin-2 (small molecule) (complex PDB ID: 1MAH<sup>157</sup>). The steps are (1) complex dissociation in a homogeneous dielectric, (2) transfer of isolated components from a homogeneous dielectric into solution with an inhomogeneous dielectric, (3) complex dissociation in an inhomogeneous dielectric, and (4) transfer of complex from a homogeneous dielectric into solution with an inhomogeneous dielectric.

where  $\Delta G_1^p$  is the Coulomb's law dissociation energy in a homogenous dielectric of 12.0,  $\Delta G_2^p$  is the solvation energy of the isolated mAChE and FAS2 molecules,  $\Delta G_4^p$  is the solvation energy of the mAChE-FAS2 complex, and  $\Delta G_3^p$  is the desired dissociation free energy. The calculations illustrated in Figure 8 were carried out with the APBS software with AMBER charges and radii,<sup>85,86</sup> protein dielectric of 12, solvent dielectric of 78.54, molecular surface definition, 150-mM ionic strength,  $129 \times 129 \times 129$  grid points, and two levels of focusing from a coarse  $125 \times 110 \times 130$  Å<sup>3</sup> grid (0.977, 0.859, 1.016 Å spacing) to a fine  $94 \times 84 \times 96$  Å<sup>3</sup> grid (0.734, 0.656, 0.750 Å spacing). The change in solvation energy on dissociation was calculated as  $\Delta G_4^p - \Delta G_2^p = -332.78$  kJ/mol; Coulomb's law gave an energy of dissociation of the complex in a homogeneous dielectric of  $\Delta G_1^p = 364.12$  kJ/mol. Together, these values give a polar dissociation energy of  $\Delta G_3^p = 31.34$  kJ/mol. To obtain the complete dissociation energy, we must also calculate the apolar contribution, which is proportional to the change in surface area when the complex dissociates:

$$\Delta G_3^a = \gamma(A_{\text{FAS2}} + A_{\text{mAChE}} - A_{\text{complex}}) \quad [14]$$

where  $\gamma$  is the surface tension and the various  $A_i$  are the surface areas of the system components. The system gains  $2496.04$  Å<sup>2</sup> on complex dissociation; with a standard surface tension<sup>37,38</sup> of  $58$  cal mol<sup>-1</sup> Å<sup>2</sup>, we obtain an apolar contribution to the dissociation energy of  $\Delta G_3^a = 34.578$  kJ/mol. Together, the polar and apolar dissociation energies give a total dissociation energy of  $65.92$  kJ/mol, which compares very favorably with the experimental dissociation constant of  $6.3$  pM ( $-64$  kJ/mol).<sup>87</sup>

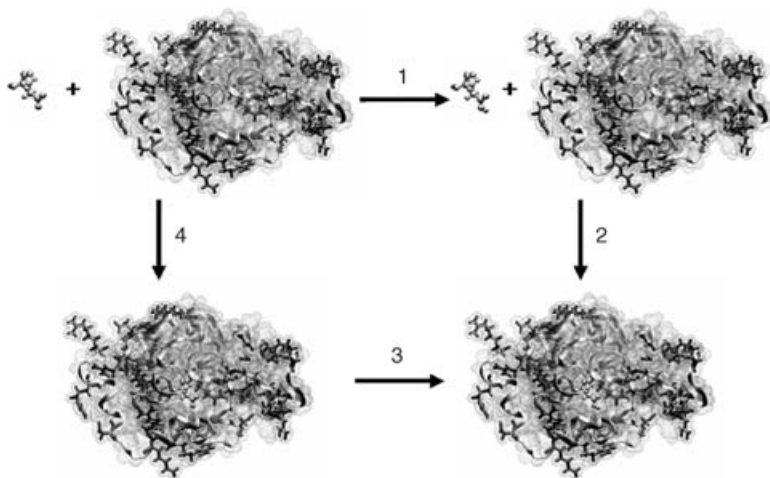
Although rigid-body binding energy calculations are sufficient for some evaluation of affinities at lower levels of accuracy, molecular association necessarily involves conformational change. This change may span a variety of length scales, ranging from simple side chain reorganization to significant protein domain motions to folding/unfolding events. In general, flexibility during molecular association is incorporated into PB calculations indirectly, through the assessment of several different conformations of the participating molecules.<sup>88</sup> These conformations are usually obtained from a molecular dynamics or Monte Carlo simulation but could also be derived from nuclear magnetic resonance (NMR) or X-ray structural data that provides insight into multiple molecular configurations. The molecular mechanics/PBSA (MM/PBSA) method is one of the most popular methods for including conformational degrees of freedom in implicit solvent binding calculations and has been the focus of numerous reviews.<sup>89–91</sup> These reviews not only provide an excellent overview of MM/PBSA, but they also provide an overview of the numerous other methods for combining conformational sampling with PB energy evaluation.

## Titration Calculations

PB and continuum electrostatics theories have been central to protein titration and  $pK_a$  calculations.<sup>92–106</sup> Such calculations inform several aspects of computational biology, including the setup of biomolecular simulations, the investigation of mechanisms for ligand binding and catalysis, and the identification of enzyme active sites. The basic elements of such titration calculations follow in a straightforward manner from the solvation energies described here. Figure 9 illustrates the steps involved in a simple titration calculation of the active site glutamate (GLU35) in hen egg-white lysozyme (PDB ID: 2LZT).<sup>107</sup> This system was chosen because of the available high-quality structural data)<sup>107</sup> as well as both computational<sup>108</sup> and experimental<sup>109</sup>  $pK_a$  measurements. The free energy for protonating the GLU residue in the context of the protein is

$$\Delta G_3 = \Delta G_1 + \Delta G_2 - \Delta G_4 \quad [15]$$

where  $\Delta G_1$  is the acid dissociation energy from the “model”  $pK_a$  of the residue,  $\Delta G_2$  is the energy of charging the unprotonated residue in the protein,  $\Delta G_4$  is the energy of charging the protonated residue, and  $\Delta G_3$  is the desired acid dissociation free energy. The “charging” calculations ( $\Delta G_2$  and  $\Delta G_4$ ) are calculated in a manner analogous to the binding energies described here; i.e., we



**Figure 9** Thermodynamic cycle illustrating the standard numerical procedure for calculating the protonation energy of the active site GLU 35 in lysozyme (CPK representation). Acidic and basic residues are shown in “licorice” representation. The steps are (1) protonation of the residue in isolation, (2) transfer of the protonated residue from isolation into the biomolecule, (3) protonation of the residue in the biomolecule, and (4) transfer of the unprotonated residue from isolation into the biomolecule. See color insert.

calculate the energy of binding the charged residue to the entire biomolecule compared with the uncharged residue. All considerations discussed here for binding calculations are directly applicable to titration calculations. Additionally, hydrogen-bond networks have been shown to play important roles in titration calculations.<sup>92,103,108</sup> Because of the importance of these interactions, hydrogen positions must be carefully optimized to maximize hydrogen-bonding before performing titration calculations. Finally,  $pK_a$  calculations are very sensitive to the choice of protein dielectric constant, with common choices ranging from 2 to 20.<sup>92,97–99,101–103,108</sup> In fact, although values of 8–12 appear to give the best results, no single choice of dielectric is appropriate for all  $pK_a$  calculations. When calculating  $pK_a$  values, it is best to examine the sensitivity of the results to the choice of dielectric and, if possible, calibrate the dielectric against known  $pK_a$  values.

The calculations outlined in Figure 9 were performed with the APBS software package for the active site glutamate (GLU35) in hen egg-white lysozyme (PDB ID: 2LZT)<sup>107</sup> with hydrogens added and optimized with AMBER.<sup>85,86</sup> The PB runs were set up with similar parameters to the work of Nielsen et al.:<sup>108</sup> solvent dielectric of 78.54, solute dielectric of 8.0, molecular surface definition,  $97 \times 97 \times 97$  grid points, no mobile ions, and three levels of focusing from a coarse grid of  $55 \times 70 \times 80 \text{ \AA}^3$  ( $0.573/0.729/0.833 \text{ \AA}$  spacings) centered on the protein to a fine grid of  $(10 \text{ \AA})^3$  ( $0.104 \text{ \AA}$  spacings), and AMBER94 charges and radii.<sup>85,86</sup> The solvation energy changes on charging the protonated and unprotonated glutamate in the protein were calculated as 1.903 and 91.688 kJ/mol, respectively. The Coulombic energy changes on charging the protonated and unprotonated glutamate in the protein were calculated as  $-13.589$  and  $-113.618$  kJ/mol, respectively. Combined, these calculations lead to total energies of  $-11.686$  and  $-21.930$  kJ/mol for charging the protonated and unprotonated GLU35 in the protein, respectively. Therefore, the difference between the acid dissociation constant in the protein and in the solution is

$$\Delta\Delta G_a = \Delta G_1 - \Delta G_3 = \Delta G_2 - \Delta G_4 = 10.244 \text{ kJ/mol} \quad [16]$$

or a  $pK_a$  shift of 1.78 units. When combined with glutamate's model  $pK_a$  of 4.3, this shift implies a  $pK_a$  for GLU35 of 6.1, a value in good agreement with the experimental measurement<sup>109</sup> of 6.2.

## Other Applications

Several applications of PB methods to biomolecular systems have existed, far too many to cover in a single chapter. However, this section highlights some of the PB studies enabled by new solver technology and likely to see increasing use by scientists in the future.

Recently, Luo et al.<sup>110</sup> have investigated many methods for accelerating PB calculations to the point where they are competitive with GB techniques.



Encouragingly, the acceleration methods described by Luo et al. are applicable to many software packages and could improve many PB solvers. Although many applications can benefit from faster PB methods, perhaps the most profound advantage will be in the realm of stochastic dynamics simulations<sup>111</sup> by enabling the efficient implicit solvent simulations of protein dynamics at the PB level of electrostatics theory.

Additionally, some recent efforts have been made to pursue more “informatics”-based approaches to the interpretation of electrostatic properties. Much of this work includes identification of functionally relevant residues in proteins by looking at electrostatic destabilization of conserved residues,<sup>112</sup> highly shifted  $pK_a$  values,<sup>93</sup> clusters of charged residues,<sup>113</sup> and protein-membrane interactions.<sup>82,114–118</sup> Other research has focused on comparisons of electrostatic potential between biomolecules, including analyses of polar and charge group complementarity at biomolecular interfaces<sup>119–124</sup> and similarity of electrostatic potentials both at molecular surfaces<sup>82,114–117,123–132</sup> and in three-dimensional space.<sup>128,133–143</sup> Although the past characterization of electrostatic properties of biomolecules has provided insight into a variety of biomolecular properties, previous applications only focused on a few quantitative measures of electrostatic properties and limited their studies to relatively few biomolecules. However, with the rapid growth of biomolecular structures elucidated by structural genomics efforts<sup>144,145</sup> and the burgeoning interest in understanding biomolecular interactions in a proteomics context, high-throughput tools to facilitate the analysis of electrostatic properties across thousands of biomolecular structures will become increasingly important.

Finally, with the advent of new parallel solution methods for the PB equation<sup>70,76</sup> and the increase in sequential computer capabilities, PB methods are being applied to increasingly larger systems. Recent examples of particularly large-scale PB calculations include the study of microtubule structure and stability,<sup>83</sup> binding calculations of aminoglycoside antibiotics to the small ribosomal subunit,<sup>84</sup> modeling of RNA organization inside viral capsids,<sup>146</sup> and examination of electrostatic interactions in protein-membrane interactions.<sup>82,114–117</sup> Given the rate of increase in computational capabilities, organelle-scale calculations appear feasible in the not-too-distant future. However, it is easy to get carried away with the scale of calculations possible with new methods and technology; we should always consider the fundamental limitations of continuum approaches and make sure the equations we used in the model are appropriate to the systems under consideration.

---

## CONCLUSIONS

An understanding of electrostatic interactions is essential for the study of biomolecular processes. The structures of proteins and other biopolymers are being determined at an increasing rate through structural genomics and other

efforts, whereas specific linkages of these biopolymers in cellular pathways or supramolecular assemblages are being detected by genetic and proteomic studies. To integrate this information in physical models for drug discovery or other applications requires the ability to evaluate the energetic interactions within and between biopolymers. Among the various components of molecular energetics, the electrostatic interactions are of special importance because of the long range of these interactions and the substantial charges of typical biopolymer components.

In this chapter, we have covered some of the basic elements of the Poisson–Boltzmann implicit solvent description of biomolecular electrostatics. Specifically, we have focused on the application of these methods to basic problems in computational biology. The discussion presented here is necessarily incomplete—electrostatics is a very broad field and continually changing. For additional background and more in-depth discussions of some of the principles and limitations of continuum electrostatics, interested readers should see the general continuum electrostatics texts by Jackson<sup>47</sup> and Landau et al.,<sup>147</sup> the electrochemistry text of Bockris et al.,<sup>16</sup> the colloid theory treatise by Verwey and Overbeek,<sup>148</sup> and the fantastic collection of condensed matter electrostatics articles assembled by Holm et al.<sup>13</sup>

---

## ACKNOWLEDGMENTS

Support for this work was provided by the Washington University in St. Louis, National Partnership for Advanced Computational Infrastructure, NIH Grant GM069702, and an Alfred P. Sloan Foundation Research Fellowship.

---

## REFERENCES

1. J. M. Wang, P. Cieplak, and P. A. Kollman, *J. Comput. Chem.*, **21**, 1049 (2000). How Well Does a Restrained Electrostatic Potential (RESP) Model Perform in Calculating Conformational Energies of Organic and Biological Molecules?
2. C. I. Bayly, P. Cieplak, W. D. Cornell, and P. A. Kollman, *J. Phys. Chem.*, **97**, 10269 (1993). A Well-Behaved Electrostatic Potential Based Method Using Charge Restraints for Deriving Atomic Charges—the RESP Model.
3. A. Grossfield, P. Ren, and J. W. Ponder, *J. Am. Chem. Soc.*, **125**, 15671 (2003). Ion Solvation Thermodynamics from Simulation with a Polarizable Force Field.
4. P. Ren and J. W. Ponder, *J. Comput. Chem.*, **23**, 1497 (2002). Consistent Treatment of Inter- and Intramolecular Polarization in Molecular Mechanics Calculations.
5. A. D. MacKerell, Jr. in *Computational Biochemistry and Biophysics*, O. M. Becker, A. D. J. MacKerell, B. Roux, and M. Watanabe, Eds., Marcel-Dekker: New York, 2001, p. 7. Atomistic Models and Force Fields.
6. A. R. Leach, *Molecular Modelling: Principles and Applications*, Prentice Hall, Harlow, England, 2001.
7. T. Schlick, *Molecular Modeling and Simulation: An Interdisciplinary Guide*, Springer-Verlag, New York, 2002.

8. G. Lamm, in *Reviews in Computational Chemistry*, Vol. 19, K. B. Lipkowitz, R. Larter, and T. R. Cundari, Eds., Wiley, Hoboken, New Jersey, 2003, pp. 147–365. The Poisson-Boltzmann Equation.
9. W. Nonner, D. Gillespe, D. Henderson, and D. Eisenberg, *J. Phys. Chem. B*, **105**, 6427 (2001). Ion Accumulation in a Biological Calcium Channel: Effects of Solvent and Confining Pressure.
10. S. M. Bhattacharyya, Z.-G. Wang, and A. H. Zewail, *J. Phys. Chem. B*, **107**, 13218 (2003). Dynamics of Water Near a Protein Surface.
11. F. Figueirido, G. S. Delbuono, and R. M. Levy, *Biophys. Chem.*, **51**, 235 (1994). Molecular Mechanics and Electrostatic Effects.
12. J.-H. Lin, N. A. Baker, and J. A. McCammon, *Biophys. J.*, **83**, 1374 (2002). Bridging the Implicit and Explicit Solvent Approaches for Membrane Electrostatics.
13. C. Holm, P. Kekicheff, and R. Podgornik, Eds., *Electrostatic Effects in Soft Matter and Biophysics*, NATO Science Series. Vol. 46, Kluwer Academic Publishers, Boston, Massachusetts, 2001.
14. T. Simonson, *Curr. Opin. Struct. Biol.*, **11**, 243 (2001). Macromolecular Electrostatics: Continuum Models and Their Growing Pains.
15. B. Roux, in *Computational Biochemistry and Biophysics*, O. M. Becker, A. D. MacKerell, Jr., B. Roux, and M. Watanabe, Eds., Marcel Dekker, New York, 2001, p. 133. Implicit Solvent Models.
16. J. O. Bockris, and A. K. N. Reddy, *Modern Electrochemistry: Ionics*, Plenum Press, New York, 1998.
17. A. D. MacKerell, Jr. and L. Nilsson, in *Computational Biochemistry and Biophysics*, O. M. Becker, A. D. MacKerell, Jr., B. Roux, and M. Watanabe, Eds., Marcel Dekker: New York, 2001, pp. 441–463. Nucleic Acid Simulations.
18. G. Hummer, L. R. Pratt, and A. E. Garcia, *J. Chem. Phys.*, **107**, 9275 (1997). Ion Sizes and Finite-size Corrections for Ionic-solvation Free Energies.
19. I. G. Tironi, R. Sperb, P. E. Smith, and W. F. van Gunsteren, *J. Chem. Phys.*, **102**, 5451 (1995). A Generalized Reaction Field Method for Molecular Dynamics Simulations.
20. N. A. Baker, P. H. Hunenberger, and J. A. McCammon, *J. Chem. Phys.*, **110**, 10679 (1999). Polarization Around an Ion in a Dielectric Continuum with Truncated Electrostatic Interactions.
21. H. S. Ashbaugh and M. E. Paulaitis, *J. Phys. Chem. B*, **102**, 5029 (1998). A Molecular/Continuum Thermodynamic Model of Hydration.
22. M. Schaefer and M. Karplus, *J. Phys. Chem.*, **100**, 1578 (1996). A Comprehensive Analytical Treatment of Continuum Electrostatics.
23. W. C. Still, A. Tempczyk, R. C. Hawley, and T. Hendrickson, *J. Am. Chem. Soc.*, **112**, 6127 (1990). Semianalytical Treatment of Solvation for Molecular Mechanics and Dynamics.
24. B. N. Dominy and C. L. Brooks III, *J. Phys. Chem. B*, **103**, 3765 (1999). Development of a Generalized Born Model Parameterization for Proteins and Nucleic Acids.
25. D. Bashford and D. A. Case, *Annu. Rev. Phys. Chem.*, **51**, 129 (2000). Generalized Born Models of Macromolecular Solvation Effects.
26. K. Osapay, W. S. Young, D. Bashford, C. L. Brooks III, and D. A. Case, *J. Phys. Chem.*, **100**, 2698 (1996). Dielectric Continuum Models for Hydration Effects on Peptide Conformational Transitions.
27. A. Onufriev, D. A. Case, and D. Bashford, *J. Comput. Chem.*, **23**, 1297 (2002). Effective Born Radii in the Generalized Born Approximation: The Importance of Being Perfect.
28. M. Born, *Z. Phys.*, **1**, 45 (1920). Volumen und Hydratationswärme der Ionen.
29. L. David, R. Luo, and M. K. Gilson, *J. Comput. Chem.*, **21**, 295 (2000). Comparison of Generalized Born and Poisson Models: Energetics and Dynamics of HIV Protease.

30. R. Luo, M. S. Head, J. Moulton, and M. K. Gilson, *J. Am. Chem. Soc.*, **120**, 6138 (1998). pK(a) Shifts in Small Molecules and HIV Protease: Electrostatics and Conformation.
31. J. A. Given and M. K. Gilson, *Proteins: Struct. Funct. Genet.*, **33**, 475 (1998). A Hierarchical Method for Generating Low-energy Conformers of a Protein-ligand Complex.
32. M. Feig, A. Onufriev, M. S. Lee, W. Im, D. A. Case, and C. L. Brooks III, *J. Comput. Chem.*, **25**, 265 (2003). Performance Comparison of Generalized Born and Poisson Methods in the Calculation of Electrostatic Solvation Energies for Protein Structures.
33. J.-P. Hansen and I. R. McDonald, *Theory of Simple Liquids*, Academic Press, San Diego, California, 2000.
34. F. Stillinger, *J. Solution Chem.*, **2**, 141 (1973). Structure in Aqueous Solutions of Nonpolar Solutes from the Standpoint of Scaled-particle Theory.
35. R. A. Pierotti, *Chem. Rev.*, **76**, 717 (1976). A Scaled Particle Theory of Aqueous and Nonaqueous Solutions.
36. T. Simonson and A. T. Brunger, *J. Phys. Chem.*, **98**, 4683 (1994). Solvation Free Energies Estimated From Macroscopic Continuum Theory: An Accuracy Assessment.
37. K. A. Sharp, A. Nicholls, R. F. Fine, and B. Honig, *Science*, **252**, 106 (1991). Reconciling the Magnitude of the Microscopic and Macroscopic Hydrophobic Effects.
38. C. Chothia, *Nature*, **248**, 338 (1974). Hydrophobic Bonding and Accessible Surface Area in Proteins.
39. L. Wesson, and D. Eisenberg, *Protein Sci.*, **1**, 227 (1992). Atomic Solvation Parameters Applied to Molecular Dynamics of Proteins in Solution.
40. D. Eisenberg and A. D. McLachlan, *Nature*, **319**, 199 (1986). Solvation Energy in Protein Folding and Binding.
41. D. M. Huang and D. Chandler, *Proc. Nat. Acad. Sci., USA*, **97**, 8324 (2000). Temperature and Length Scale Dependence of Hydrophobic Effects and their Possible Implications for Protein Folding.
42. G. Hummer, *J. Am. Chem. Soc.*, **121**, 6299 (1999). Hydrophobic Force Field as a Molecular Alternative to Surface-area Models.
43. L. R. Pratt, *Annu. Rev. Phys. Chem.*, **53**, 409 (2002). Molecular Theory of Hydrophobic Effects: “She is too Mean to Have Her Name Repeated”.
44. N. A. Baker and J. A. McCammon, in *Structural Bioinformatics*, H. Weissig, Ed., Wiley, New York, 2002, pp. 427–440. Electrostatic Interactions.
45. M. E. Davis and J. A. McCammon, *Chem. Rev.*, **94**, 7684 (1990). Electrostatics in Biomolecular Structure and Dynamics.
46. B. Honig and A. Nicholls, *Science*, **268**, 1144 (1995). Classical Electrostatics in Biology and Chemistry.
47. J. D. Jackson, *Classical Electrodynamics*, Wiley, New York, 1975.
48. B. Lee and F. M. Richards, *J. Mol. Biol.*, **55**, 379 (1971). The Interpretation of Protein Structures: Estimation of Static Accessibility.
49. M. L. Connolly, *J. Mol. Graph.*, **11**, 139 (1993). The Molecular Surface Package.
50. W. Im, D. Beglov, and B. Roux, *Comput. Phys. Commun.*, **111**, 59 (1998). Continuum Solvation Model: Electrostatic Forces From Numerical Solutions to the Poisson–Boltzmann Equation.
51. C. L. Bajaj, V. Pasucci, R. J. Holt, and A. N. Netravali, *Discrete Appl. Math.*, **127**, 23 (2003). Dynamic Maintenance and Visualization of Molecular Surfaces.
52. J. A. Grant, B. T. Pickup, and A. Nicholls, *J. Comput. Chem.*, **22**, 608 (2001). A Smooth Permittivity Function for Poisson–Boltzmann Solvation Methods.
53. K. A. Sharp and B. Honig, *J. Phys. Chem.*, **94**, 7684 (1990). Calculating Total Electrostatic Energies with the Nonlinear Poisson–Boltzmann Equation.

- 
54. F. Fogolari and J. M. Briggs, *Chem. Phys. Lett.*, **281**, 135 (1997). On the Variational Approach to Poisson–Boltzmann Free Energies.
  55. A. M. Micu, B. Bagheri, A. V. Ilin, L. R. Scott, and B. M. Pettitt, *J. Comput. Phys.*, **136**, 263 (1997). Numerical Considerations in the Computation of the Electrostatic Free Energy of Interaction within the Poisson–Boltzmann Theory.
  56. M. K. Gilson, M. E. Davis, B. A. Luty, and J. A. McCammon, *J. Phys. Chem.*, **97**, 3591 (1993). Computation of Electrostatic Forces on Solvated Molecules using the Poisson–Boltzmann Equation.
  57. M. J. Holst and F. Saied, *J. Comput. Chem.*, **16**, 337 (1995). Numerical Solution of the Nonlinear Poisson–Boltzmann Equation: Developing More Robust and Efficient Methods.
  58. M. E. Davis, J. D. Madura, B. A. Luty, and J. A. McCammon, *Comput. Phys. Commun.*, **62**, 187 (1991). Electrostatics and Diffusion of Molecules in Solution—Simulations with the University-of-Houston-Brownian Dynamics Program.
  59. W. Rocchia, S. Sridharan, A. Nicholls, E. Alexov, A. Chiabrera, and B. Honig, *J. Comput. Chem.*, **23**, 128 (2002). Rapid Grid-based Construction of the Molecular Surface and the use of Induced Surface Charge to Calculate Reaction Field Energies: Applications to the Molecular Systems and Geometric Objects.
  60. M. E. Davis and J. A. McCammon, *J. Comput. Chem.*, **12**, 909 (1991). Dielectric Boundary Smoothing in Finite Difference Solutions of the Poisson Equation: An Approach to Improve Accuracy and Convergence.
  61. M. K. Gilson and B. Honig, *Nature*, **330**, 84 (1987). Calculation of Electrostatic Potentials in an Enzyme Active Site.
  62. K. A. Sharp and B. Honig, *Ann. Rev. Biophys.*, **19**, 301 (1990). Electrostatic Interactions in Macromolecules—Theory and Applications.
  63. H. X. Zhou, *Biophys. J.*, **65**, 955 (1993). Boundary Element Solution of Macromolecular Electrostatics: Interaction Energy between two Proteins.
  64. R. J. Zauhar and R. S. Morgan, *J. Mol. Biol.*, **186**, 815 (1985). A New Method for Computing the Macromolecular Electric Potential.
  65. O. Axelsson and V. A. Barker, *Finite Element Solution of Boundary Value Problems. Theory and Computation*, Academic Press, Inc., San Diego, California, 1984.
  66. D. Braess, *Finite Elements. Theory, Fast Solvers, and Applications in Solid Mechanics*, Cambridge University Press, Cambridge, United Kingdom, 1997.
  67. S. C. Brenner and L. R. Scott, *The Mathematical Theory of Finite Element Methods*, Springer, New York, 2002.
  68. M. J. Holst, N. A. Baker, and F. Wang, *J. Comput. Chem.*, **21**, 1319 (2000). Adaptive Multilevel Finite Element Solution of the Poisson–Boltzmann Equation. I. Algorithms and Examples.
  69. N. A. Baker, M. J. Holst, and F. Wang, *J. Comput. Chem.*, **21**, 1343 (2000). Adaptive Multilevel Finite Element Solution of the Poisson–Boltzmann Equation. II. Refinement at Solvent-accessible Surfaces in Biomolecular Systems.
  70. N. A. Baker, D. Sept, M. J. Holst, and J. A. McCammon, *IBM J. Res. Develop.*, **45**, 427 (2001). The Adaptive Multilevel Finite Element Solution of the Poisson–Boltzmann Equation on Massively Parallel Computers.
  71. C. M. Cortis and R. A. Friesner, *J. Comput. Chem.*, **18**, 1591 (1997). Numerical Solution of the Poisson–Boltzmann Equation using Tetrahedral Finite-element Meshes.
  72. C. M. Cortis and R. A. Friesner, *J. Comput. Chem.*, **18**, 1570 (1997). An Automatic Three-dimensional Finite Element Mesh Generation System for the Poisson–Boltzmann Equation.
  73. W. Hackbusch, *Multi-grid Methods and Applications*, Springer-Verlag, Berlin, Germany, 1985.
  74. W. L. Briggs, *A Multigrid Tutorial*, Society for Industrial and Applied Mathematics, Philadelphia, Pennsylvania, 1987.

75. M. J. Holst, *Adv. Comput. Math.*, **15**, 139 (2001). Adaptive Numerical Treatment of Elliptic Systems on Manifolds.
76. N. A. Baker, D. Sept, S. Joseph, M. J. Holst, and J. A. McCammon, *Proc. Natl. Acad. Sci., USA*, **98**, 10037 (2001). Electrostatics of Nanosystems: Application to Microtubules and the Ribosome.
77. R. E. Bank, and M. J. Holst, *SIAM J. Sci. Comput.*, **22**, 1411 (2000). A New Paradigm for Parallel Adaptive Meshing Algorithms.
78. D. K. Ross, *SIAM J. Appl. Math.*, **29**, 699 (1975). The Interaction Energy, Field Strength and Force Acting on a Pair of Dielectric Spheres Embedded in a Dielectric Medium.
79. J. Q. Feng, *Phys. Rev. E*, **62**, 2891 (2000). Electrostatic Interaction between two Charged Dielectric Spheres in Contact.
80. N. V. Sushkin and G. D. J. Phillies, *J. Chem. Phys.*, **103**, 4600 (1995). Charged Dielectric Spheres in Electrolyte Solutions: Induced Dipole and Counterion Exclusion Effects.
81. C. F. Wong, P. H. Hunenberger, P. Akamine, N. Narayana, T. Diller, J. A. McCammon, S. Taylor, and N. H. Xuong, *J. Med. Chem.*, **44**, 1530 (2001). Computational Analysis of PKA-balanol Interactions.
82. D. Murray, N. Ben-Tal, B. Honig, and S. McLaughlin, *Structure*, **5**, 985 (1997). Electrostatic Interaction of Myristoylated Proteins with Membranes: Simple Physics, Complicated Biology.
83. D. Sept, N. A. Baker, and J. A. McCammon, *Protein Sci.*, **12**, 2257 (2003). The Physical Basis of Microtubule Structure and Stability.
84. C. Ma, N. A. Baker, S. Joseph, and J. A. McCammon, *J. Am. Chem. Soc.*, **124**, 1438 (2002). Binding of Aminoglycoside Antibiotics to the Small Ribosomal Subunit: A Continuum Electrostatics Investigation.
85. D. A. Case, D. A. Pearlmann, J. W. Caldwell, J. Wang, W. S. Ross, C. Simmerling, T. Darden, K. M. Merz, R. V. Stanton, A. Cheng, J. J. Vincent, M. Crowley, V. Tsui, H. Gohlke, R. Radmer, Y. Duan, J. Pitera, I. Massova, G. L. Seibel, U. C. Singh, P. Weiner, and P. A. Kollman, *AMBER 7 User's Manual*, University of California. (2002). Available: <<http://amber.scripps.edu/doc7/amber.pdf>>.
86. D. A. Pearlmann, D. A. Case, J. W. Caldwell, W. S. Ross, T. E. Cheatham, 3rd, S. DeBolt, D. Ferguson, G. L. Seibel, and P. A. Kollman, *Comput. Phys. Commun.*, **91**, 1 (1995). AMBER, A Package of Computer Programs for Applying Molecular Mechanics, Normal Mode Analysis, Molecular Dynamics, and Free Energy Calculations to Simulate the Structural and Energetic Properties of Molecules.
87. Z. Radic and P. Taylor, *J. Biol. Chem.*, **276**, 4622 (2001). Interaction Kinetics of Reversible Inhibitors and Substrates with Acetylcholinesterase and its Fasciculin 2 Complex.
88. A. H. Elcock, D. Sept, and J. A. McCammon, *J. Phys. Chem. B*, **105**, 1504 (2001). Computer Simulation of Protein–Protein Interactions.
89. I. Massova and P. A. Kollman, *Perspect. Drug Discov. Design*, **18**, 113 (2000). Combined Molecular Mechanical and Continuum Solvent Approach (MM-PBSA/GBSA) to Predict Ligand Binding.
90. J. Wang, P. Morin, W. Wang, and P. A. Kollman, *J. Am. Chem. Soc.*, **123**, 5221 (2001). Use of MM-PBSA in Reproducing the Binding Free Energies to HIV-1 RT of TIBO Derivatives and Predicting the Binding Mode to HIV-1 RT of Efavirenz by Docking and MM-PBSA.
91. P. A. Kollman, I. Massova, C. Reyes, B. Kuhn, S. Huo, L. Chong, M. Lee, T. Lee, Y. Duan, W. Wang, O. Donini, P. Cieplak, J. Srinivasan, D. A. Case, and T. E. Cheatham, 3rd, *Acc. Chem. Res.*, **33**, 889 (2000). Calculating Structures and Free Energies of Complex Molecules: Combining Molecular Mechanics and Continuum Models.
92. J. E. Nielsen and G. Vriend, *Proteins: Struct. Funct. Genet.*, **43**, 403 (2001). Optimizing the Hydrogen-bond Network in Poisson–Boltzmann Equation-based pK(a) Calculations.
93. M. J. Ondrechen, J. G. Clifton, and D. Ringe, *Proc. Natl. Acad. Sci. (USA)*, **98**, 12473 (2001). THEMATICS: A Simple Computational Predictor of Enzyme Function from Structure.

94. A. Onufriev, D. A. Case, and G. M. Ullmann, *Biochemistry*, **40**, 3413 (2001). A Novel View of pH Titration in Biomolecules.
95. C. Tanford and J. G. Kirkwood, *J. Am. Chem. Soc.*, **79**, 5333 (1957). Theory of Protein Titration Curves. I. General Equations for Impenetrable Spheres.
96. H. W. T. van Vlijmen, M. Schaefer, and M. Karplus, *Proteins: Struct. Funct. Genet.*, **33**, 145 (1998). Improving the Accuracy of Protein pK(a) Calculations: Conformational Averaging versus the Average Structure.
97. A. S. Yang, M. R. Gunner, R. Sampogna, K. A. Sharp, and B. Honig, *Proteins: Struct. Funct. Genet.*, **15**, 252 (1993). On the Calculation of Pk(a)s in Proteins.
98. A. S. Yang and B. Honig, *J. Mol. Biol.*, **231**, 459 (1993). On the pH Dependence of Protein Stability.
99. A. S. Yang and B. Honig, *J. Mol. Biol.*, **237**, 602 (1994). Structural Origins of pH and Ionic Strength Effects on Protein Stability—Acid Denaturation of Sperm Whale Apomyoglobin.
100. J. Antosiewicz, J. M. Briggs, A. H. Elcock, M. K. Gilson, and J. A. McCammon, *J. Comput. Chem.*, **17**, 1633 (1996). Computing Ionization States of Proteins with a Detailed Charge Model.
101. J. Antosiewicz, J. A. McCammon, and M. K. Gilson, *Biochemistry*, **35**, 7819 (1996). The Determinants of pK(a)s in Proteins.
102. D. Bashford and M. Karplus, *Biochemistry*, **29**, 10219 (1990). pK<sub>a</sub>'s of Ionizable Groups in Proteins: Atomic Detail from a Continuum Electrostatic Model.
103. J. E. Nielsen, K. V. Andersen, B. Honig, R. W. W. Hooft, G. Klebe, G. Vriend, and R. C. Wade, *Protein Eng.*, **12**, 657 (1999). Improving Macromolecular Electrostatics Calculations.
104. M. Schaefer, M. Sommer, and M. Karplus, *J. Phys. Chem., B*, **101**, 1663 (1997). pH-Dependence of Protein Stability: Absolute Electrostatic Free Energy Differences between Conformations.
105. H. Oberoi and N. M. Allewell, *Biophys. J.*, **65**, 48 (1993). Multigrid Solution of the Nonlinear Poisson-Boltzmann Equation and Calculation of Titration Curves.
106. D. Morikis, A. H. Elcock, P. A. Jennings, and J. A. McCammon, *Protein Sci.*, **10**, 2379 (2001). Proton Transfer Dynamics of GART: The pH-dependent Catalytic Mechanism Examined by Electrostatic Calculations.
107. M. Ramanadham, L. C. Sieker, and L. H. Jensen, *Acta Crystallogr. B*, **46**, 63 (1990). Refinement of Triclinic Lysozyme. II. The Method of Stereochemically Restrained Least Squares.
108. J. E. Nielsen and J. A. McCammon, *Protein Sci.*, **12**, 313 (2003). On the Evaluation and Optimization of Protein X-ray Structures for pK<sub>a</sub> Calculations.
109. D. J. Vocadlo, G. J. Davies, R. Laine, and S. G. Withers, *Nature*, **412**, 835 (2001). Catalysis by Hen Egg-white Lysozyme Proceeds via a Covalent Intermediate.
110. R. Luo, L. David, and M. K. Gilson, *J. Comput. Chem.*, **23**, 1244 (2002). Accelerated Poisson-Boltzmann Calculations for Static and Dynamic Systems.
111. Q. Lu and R. Luo, *J. Chem. Phys.*, **119**, 11035 (2003). A Poisson-Boltzmann Dynamics Method with Nonperiodic Boundary Condition.
112. A. H. Elcock, *J. Mol. Biol.*, **312**, 885 (2001). Prediction of Functionally Important Residues Based Solely on the Computed Energetics of Protein Structure.
113. Z. Y. Zhu and S. Karlin, *Proc. Natl. Acad. Sci. (USA)*, **93**, 8350 (1996). Clusters of Charged Residues in Protein Three-dimensional Structures.
114. D. Murray, A. Arbuzova, G. Hangyas-Mihalyne, A. Gambhir, N. Ben-Tal, B. Honig, and S. McLaughlin, *Biophys. J.*, **77**, 3176 (1999). Electrostatic Properties of Membranes Containing Acidic Lipids and Adsorbed Basic Peptides: Theory and Experiment.
115. D. Murray and B. Honig, *Mol. Cell.*, **9**, 145 (2002). Electrostatic Control of the Membrane Targeting of C2 Domains.

116. D. Murray, L. H. Matsumoto, C. A. Buser, J. Tsang, C. T. Sigal, N. Ben-Tal, B. Honig, M. D. Resh, and S. McLaughlin, *Biochemistry*, **37**, 2145 (1998). Electrostatics and the Membrane Association of Src: Theory and Experiment.
117. D. Murray, S. McLaughlin, and B. Honig, *J. Biol. Chem.*, **276**, 45153 (2001). The Role of Electrostatic Interactions in the Regulation of the Membrane Association of G Protein Beta Gamma Heterodimers.
118. K. Diraviyam, R. V. Stahelin, W. Cho, and D. Murray, *J. Mol. Biol.*, **328**, 721 (2003). Computer Modeling of the Membrane Interaction of FYVE Domains.
119. L. Lo Conte, C. Chothia, and J. Janin, *J. Mol. Biol.*, **285**, 2177 (1999). The Atomic Structure of Protein-protein Recognition Sites.
120. J. Janin and C. Chothia, *J. Biol. Chem.*, **265**, 16027 (1990). The Structure of Protein–protein Recognition Sites.
121. D. Xu, S. L. Lin, and R. Nussinov, *J. Mol. Biol.*, **265**, 68 (1997). Protein Binding versus Protein Folding: The Role of Hydrophilic Bridges in Protein Associations.
122. V. A. Roberts, H. C. Freeman, A. J. Olson, J. A. Tainer, and E. D. Getzoff, *J. Biol. Chem.*, **266**, 13431 (1991). Electrostatic Orientation of the Electron-transfer Complex between Plastocyanin and Cytochrome c.
123. J. Novotny and K. A. Sharp, *Prog. Biophys. Mol. Biol.*, **58**, 203 (1992). Electrostatic Fields in Antibodies and Antibody/antigen Complexes.
124. A. J. McCoy, V. Chandana Epa, and P. M. Colman, *J. Mol. Biol.*, **268**, 570 (1997). Electrostatic Complementarity at Protein/protein Interfaces.
125. S. A. Botti, C. E. Felder, J. L. Sussman, and I. Silman, *Protein Eng.*, **11**, 415 (1998). Electrotactins: A Class of Adhesion Proteins with Conserved Electrostatic and Structural Motifs.
126. C. E. Felder, S. A. Botti, S. Lifson, I. Silman, and J. L. Sussman, *J. Mol. Graph. Model.*, **15**, 318 (1997). External and Internal Electrostatic Potentials of Cholinesterase Models.
127. E. Demchuk, T. Mueller, H. Oshkinat, W. Sebald, and R. C. Wade, *Protein Sci.*, **3**, 920 (1994). Receptor Binding Properties of Four-helix-bundle Growth Factors Deduced from Electrostatic Analysis.
128. L. T. Chong, S. E. Dempster, Z. S. Hendsch, L. P. Lee, and B. Tidor, *Protein Sci.*, **7**, 206 (1998). Computation of Electrostatic Complements to Proteins: A Case of Charge Stabilized Binding.
129. T. Sulea and E. O. Purisima, *Biophys. J.*, **84**, 2883 (2003). Profiling Charge Complementarity and Selectivity for Binding at the Protein Surface.
130. N. Sinha, S. Mohan, C. A. Lipschultz, and S. J. Smith-Gill, *Biophys. J.*, **83**, 2946 (2002). Differences in Electrostatic Properties at Antibody-antigen Binding Sites: Implications for Specificity and Cross-reactivity.
131. B. Ma, T. Elkayam, H. Wolfson, and R. Nussinov, *Proc. Natl. Acad. Sci. (USA)*, **100**, 5772 (2003). Protein-protein Interactions: Structurally Conserved Residues Distinguish between Binding Sites and Exposed Protein Surfaces.
132. K. Kinoshita, J. Furui, and H. Nakamura, *J. Struct. Funct. Genom.*, **2**, 9 (2002). Identification of Protein's Functions from a Molecular Surface Database, eF-site.
133. N. Blomberg, R. R. Gabdoulline, M. Nilges, and R. C. Wade, *Proteins: Struct. Funct. Genet.*, **37**, 379 (1999). Classification of Protein Sequences by Homology Modeling and Quantitative Analysis of Electrostatic Similarity.
134. R. C. Wade, R. R. Gabdoulline, and B. A. Luty, *Proteins*, **31**, 406 (1998). Species Dependence of Enzyme-substrate Encounter Rates for Triose Phosphate Isomerases.
135. R. C. Wade, R. R. Gabdoulline, S. K. Ludemann, and V. Lounnas, *Proc. Natl. Acad. Sci. (USA)*, **95**, 5942 (1998). Electrostatic Steering and Ionic Tethering in Enzyme-ligand Binding: Insights from Simulations.
136. R. C. Wade, R. R. Gabdoulline, and F. De Rienzo, *Int. J. Quantum Chem.*, **83**, 122 (2001). Protein Interaction Property Similarity Analysis.



137. L. P. Lee and B. Tidor, *Protein Sci.*, **10**, 362 (2001). Optimization of Binding Electrostatics: Charge Complementarity in the Barnase–Barstar Protein Complex.
138. L. P. Lee and B. Tidor, *J. Chem. Phys.*, **106**, 8681 (1997). Optimization of Electrostatic Binding Free Energy.
139. E. Kangas and B. Tidor, *Phys. Rev. E*, **59**, 5958 (1999). Charge Optimization Leads to Favorable Electrostatic Binding Free Energy.
140. A. M. Richard, *J. Comput. Chem.*, **12**, 959 (1991). Quantitative Comparison of Molecular Electrostatic Potentials for Structure-activity Studies.
141. C. Burt, W. G. Richards, and P. Huxley, *J. Comput. Chem.*, **11**, 1139 (1990). The Application of Molecular Similarity Calculations.
142. A. C. Good, E. E. Hodgkin, and W. G. Richards, *J. Comput.-Aided Molec. Design*, **6**, 513 (1992). Similarity Screening of Molecular Data Sets.
143. D. R. Livesay, P. Jambeck, A. Rojnuckarin, and S. Subramaniam, *Biochemistry*, **42**, 3464 (2003). Conservation of Electrostatic Properties within Enzyme Families and Superfamilies.
144. H. M. Berman, T. N. Bhat, P. E. Bourne, Z. Feng, G. Gilliland, H. Weissig, and J. Westbrook, *Nat. Struct. Biol.*, **7 Suppl.**, 957 (2000). The Protein Data Bank and the Challenge of Structural Genomics.
145. S. K. Burley, S. C. Almo, J. B. Bonanno, M. Capel, M. R. Chance, T. Gaasterland, D. Lin, A. Sali, F. W. Studier, and S. Swaminathan, *Nat. Genet.*, **23**, 151 (1999). Structural Genomics: Beyond the Human Genome Project.
146. D. Zhang, R. Konecny, N. A. Baker, and J. A. McCammon, *Biopolymers*, **75**, 325 (2004). Electrostatic Interactions between RNA and Protein Capsid in Cowpea Chlorotic Mottle Virus Simulated by a Coarse-grain RNA Model and Monte Carlo Approach.
147. L. D. Landau, E. M. Lifshitz, and L. P. Pitaevskii, *Electrodynamics of Continuous Media*, Butterworth-Heinemann, Boston, Massachusetts, 1982.
148. E. J. W. Verwey and J. T. G. Overbeek, *Theory of the Stability of Lyophobic Colloids*, Dover Publications, Inc., Mineola, New York, 1999.
149. Adaptive Poisson–Boltzmann Solver (APBS). 1999–2003. Washington University, St. Louis, Missouri. Available: <http://agave.wustl.edu/apbs/>
150. A. Nicholls and B. Honig, *J. Comput. Chem.*, **12**, 435 (1991). A Rapid Finite Difference Algorithm, Utilizing Successive Over-relaxation to Solve the Poisson–Boltzmann Equation.
151. A. Nicholls, K. A. Sharp, and B. Honig, *Proteins*, **11**, 281 (1991). Protein Folding and Association: Insights from the Interfacial and Thermodynamic Properties of Hydrocarbons.
152. D. Bashford, in *Scientific Computing in Object-Oriented Parallel Environments*, Y. Ishikawa, R. R. Oldehoeft, J. V. W. Reynders, and M. Tholburn, Eds., Springer, Berlin, Germany, 1997, pp. 133–140. An Object-oriented Programming Suite for Electrostatic Effects in Biological Molecules.
153. J. D. Madura, J. M. Briggs, R. C. Wade, M. E. Davis, B. A. Luty, A. V. Ilin, J. Antosiewicz, M. K. Gilson, B. Bagheri, L. R. Scott, and J. A. McCammon, *Comput. Phys. Commun.*, **91**, 57 (1995). Electrostatics and Diffusion of Molecules in Solution—Simulations with the University of Houston Brownian Dynamics Program.
154. G. Vacek, J. K. Perry, and J.-M. Langlois, *Chem. Phys. Lett.*, **310**, 189 (1999). Advanced Initial-guess Algorithm for Self-consistent-field Calculations on Organometallic Systems.
155. B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. Swaminathan, and M. Karplus, *J. Comput. Chem.*, **4**, 187 (1983). CHARMM: A Program for Macromolecular Energy, Minimization, and Dynamics Calculations.
156. A. D. MacKerell, Jr., B. R. Brooks, C. L. Brooks III, L. Nilsson, B. Roux, Y. Won, and M. Karplus, in *The Encyclopedia of Computational Chemistry*, P. v. R. Schleyer, Ed. Wiley, Chichester, United Kingdom, 1998, p. 271. CHARMM: The Energy Function and its Parameterization with an Overview of the Program.
157. Y. Bourne, P. Taylor, and P. Marchot, *Cell*, **83**, 503 (1995). Acetylcholinesterase Inhibition by Fasciculin: Crystal Structure of the Complex.



# Data Sources and Computational Approaches for Generating Models of Gene Regulatory Networks

Baltazar D. Aguda,<sup>\*</sup> Georghe Craciun,<sup>†</sup> and  
Rengul Cetin-Atalay<sup>‡</sup>

<sup>\*</sup>*Bioinformatics Institute, 30 Biopolis Street, #07-01 Matrix, Singapore 138671*

<sup>†</sup>*Mathematical Biosciences Institute, The Ohio State University  
231 W. 18th Avenue, Columbus, Ohio 43210*

<sup>‡</sup>*Department of Molecular Biology & Genetics, Bilkent University, Ankara, Turkey and Virginia Bioinformatics Institute, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24061*

---

---

## INTRODUCTION

High-throughput data acquisition technologies in molecular biology, including rapid DNA sequencers, gene expression microarrays, and other microchip-based assays, are providing an increasingly comprehensive parts list of a biological cell. Although this parts list may be far from complete at this time, the so-called “post-genomic era” has now begun in which the goal is to integrate the parts and analyze how they interact to determine the system’s behavior. This integration is being facilitated by the creation of databases, knowledgebases, and other information repositories on the Internet.

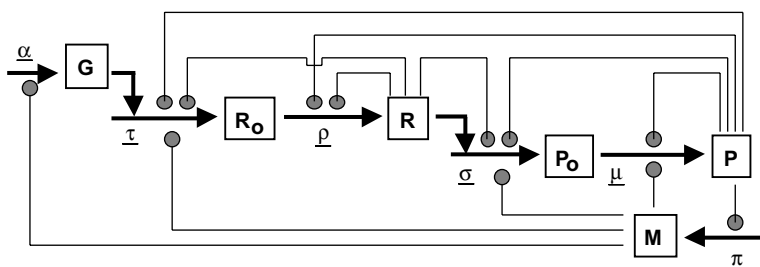
*Reviews in Computational Chemistry, Volume 21*  
edited by Kenny B. Lipkowitz, Raima Larter, and Thomas R. Cundari  
Copyright © 2005 Wiley-VCH, John Wiley & Sons, Inc.

How these huge amounts of information will answer biological questions and predict behavior will keep multidisciplinary teams of scientists busy for many years. A key question is how the expression of genes is regulated in response to various intracellular and external conditions and stimuli. The current paradigm is that the secret to life could be found in the genetic code; however, the expression of genes and the unfolding of the regulatory molecular networks in response to the environment may well be the defining attribute of the living state.

This chapter focuses on gene regulatory networks (GRNs). A “gene regulatory network” refers to a set of molecules and interactions that affect the expression of genes located in the DNA of a cell. Gene expression is the combination of *transcription* of DNA sequences, *processing* of the primary RNA transcripts, and *translation* of the mature messenger RNA (mRNA) to proteins in ribosomes. This picture is often referred to as the “central dogma,” and it has been the canonical model for the flow of information from the genetic code to proteins. These processes are shown schematically as steps labeled  $\tau$ ,  $\rho$ , and  $\sigma$  in Figure 1.

The step labeled  $\mu$  in Figure 1 represents modification of primary proteins to render them functional; examples would be posttranslational covalent modifications (e.g., phosphorylation) and binding with other proteins or other molecules. Represented within the set of steps  $\mu$  are the many regulatory events (other than transcription and translation) affecting gene expression and the overall physiology of the cell.

The complexity of GRNs may develop from the many possible feedback loops shown as gray lines in Figure 1. In step  $\tau$ , proteins could be directly involved in transcription, as in the case of transcription factors binding to upstream regulatory regions of genes. Many RNA and protein molecules cooperate in the translation step  $\sigma$  in Figure 1; examples are tRNA, rRNA, and ribosomal proteins.



**Figure 1** A schematic representation of a gene regulatory network involving modules of molecular classes (shown in boxes); the modules shown are the transcriptional units in the genome (G), primary transcripts ( $R_o$ ), mature transcripts (R), primary proteins ( $P_o$ ), modified proteins (P), and metabolites (M). The labeled steps shown in black lines are transcription ( $\tau$ ), RNA processing ( $\rho$ ), translation ( $\sigma$ ), protein modification ( $\mu$ ), metabolic pathways ( $\pi$ ), and genome replication ( $\alpha$ ). The feedback interactions shown in gray lines are discussed in the text. Filled circles represent either inhibition or activation.

The first goal of this chapter is to survey sources of data and other information that can generate models of GRNs. The focus is on biological databases and knowledgebases that are available on the Internet, especially those that attempt to integrate heterogeneous information including molecular interactions and pathways. The second goal of this review is to summarize current models of GRNs and how they can be extracted from biological databases. Depending on the nature of the data, different granularities of GRN models can be generated, ranging from probabilistic graphical models to detailed kinetic or mechanistic models. A crucial issue in the design of pathways databases is how to represent information having various levels of uncertainty. Because of its central importance in GRN modeling, an extensive discussion of pathway ontology is given. Lastly, the third goal is to discuss theoretical and computational methods for the analysis of detailed models of GRNs. In particular, a summary is given of various tools already developed in the field of reaction network analysis. Particular emphasis of the discussion is on exploiting information on network structure to deduce potential behavior of GRNs without knowing quantitative values of rate parameters.

## FORMAL REPRESENTATION OF GRNs

The GRN of Figure 1 can be formally translated to a set of general dynamical equations. The modules (in boxes) in the GRN represent the following classes of biomolecules:

**G** : Vector of all transcriptional units (TUs) involved in the GRN (in terms, for example, of gene dosage per TU).

**R<sub>o</sub>** : Vector of primary RNA transcripts corresponding to the TUs in **G**.

**R** : Vector of messenger RNA (mRNA), transfer RNA (tRNA), ribosomal RNA (rRNA), and other processed RNAs.

**P<sub>o</sub>** : Vector of newly translated (primary) proteins.

**P** : Vector of modified proteins.

**M** : Vector of metabolites.

Disregarding the replication of the genomic DNA (step  $\alpha$ ) and the changes in the metabolome **M** for now (i.e., assume **G** and **M** to be constant), a mathematical representation of the dynamics of the GRN in Figure 1 would be the following set of vector-matrix equations:

$$\begin{aligned} d\mathbf{R}_o/dt &= \tau\mathbf{G} - \rho\mathbf{R}_o - \delta_1\mathbf{R}_o \\ d\mathbf{R}/dt &= \rho\mathbf{R}_o - \delta_2\mathbf{R} \\ d\mathbf{P}_o/dt &= \sigma\mathbf{R} - \mu\mathbf{P}_o - \delta_3\mathbf{P}_o \\ d\mathbf{P}/dt &= \mu\mathbf{P}_o - \delta_4\mathbf{P} \end{aligned} \quad [1]$$

The “RNA transcription” matrix  $\tau$  is a diagonal matrix (i.e., all off-diagonal entries are 0), with the nonzero entries being, in general, functions of  $\mathbf{R}$ ,  $\mathbf{P}$ , and  $\mathbf{M}$  as depicted by the feedback loops in Figure 1. The “RNA-processing matrix”  $\rho$  is a diagonal matrix with the nonzero entries being, in general, functions of  $\mathbf{R}$  and  $\mathbf{P}$ . The diagonal matrix  $\sigma$  is called the “protein translation” matrix. The diagonal matrix  $\mu$  is called the “protein modification” matrix (which includes all posttranslational modifications and protein–protein interactions). Figure 1 shows the dependence of  $\sigma$  and  $\mu$  on  $\mathbf{R}$ ,  $\mathbf{P}$ , and  $\mathbf{M}$ . The diagonal matrices  $\delta_i$  are “degradation” matrices that account for the degradation of RNA and protein molecules as well as their transport or dilution. Because of the general dependence of the matrices to the variables  $\mathbf{R}$ ,  $\mathbf{P}$ , and  $\mathbf{M}$ , Eq. [1] involves nonlinear equations in these variables.

An example of a GRN is given next to illustrate the formal representation just described. The example also demonstrates the art of modeling and reduction of the network into minimal mathematical models.

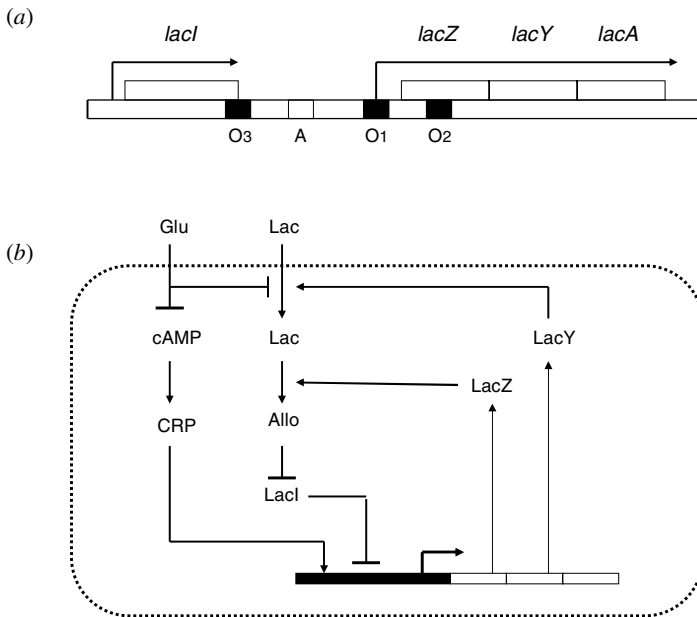
---

## AN EXAMPLE OF A GRN: THE LAC OPERON

The *lac operon* in the bacterium *Escherichia coli* is a well-studied GRN. This prokaryotic gene network has been the subject of numerous reviews;<sup>1–4</sup> it is discussed here primarily to illustrate the various aspects of GRN modeling, starting with the information on genome organization (operon structure) to knowledge on protein–DNA interactions, protein–protein interactions, and the influence of metabolites.

Understanding the *lac operon* begins by looking at the genome organization of *E. coli*. The complete genome sequence of various strains of this bacterium can be accessed through the webpage of the National Center for Biotechnology Information (<http://www.ncbi.nlm.nih.gov>). From the homepage menu, clicking *Entrez* followed by *Genome* gives the link to complete bacterial genomes including *E. coli*. Genes in the circular chromosome of *E. coli* are organized into “operons.” An operon is a cluster of genes whose expression is controlled by a common set of operator sequences and regulatory proteins.<sup>5</sup> The genes in the cluster are usually involved in the synthesis of enzymes needed for the metabolism of a molecule. Several reviews on the influence of operon structure on the dynamical behavior of GRNs are available.<sup>6,7</sup>

The *lac operon* is shown in Figure 2(a). The GRN involves the gene set  $\{lacZ, lacY, lacA, lacI\}$  and the regulatory sequences  $\{O1, O2, O3, A\}$  as shown in Figure 2(a). The gene *lacI* encodes a repressor protein that binds the operator sequences  $O1$ ,  $O2$ , and  $O3$ , thereby repressing the synthesis of the *lacZ*-*lacY*-*lacA* transcript. Gene *lacZ* encodes the  $\beta$ -galactosidase enzyme, gene *lacY* encodes a permease, and gene *lacA* encodes a transacetylase. The CRP/cAMP complex binds the sequence  $A$  and enhances transcription.



**Figure 2** The lac operon. (a) The expression of the genes *lacZ*, *lacY*, and *lacA* as one transcriptional unit is controlled by the upstream regulatory sequences including the operator regions O1, O2, and O3 where the repressor protein (the product of the *lacI* gene) binds. The CRP/cAMP protein complex binds the sequence A, which results in increased transcription. (b) A schematic representation of the key pathways regulating the lac operon. (Figure is modified from Ozbudak et al.<sup>3</sup>)

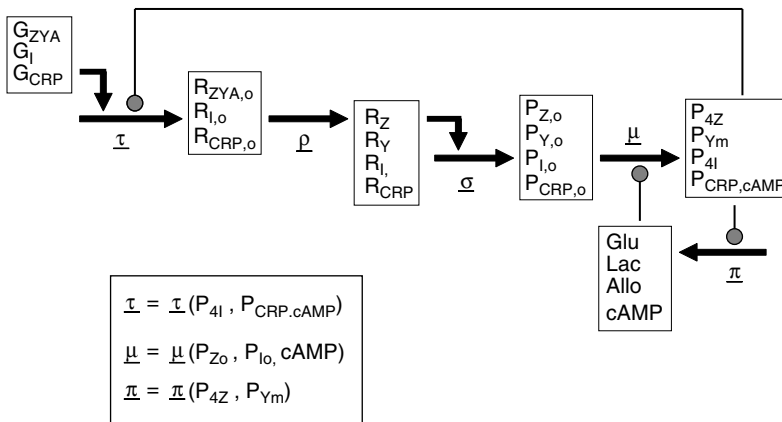
The key pathways that generate the switching behavior of the GRN are shown in Figure 2(b). This switching behavior of the lac operon explains the diauxic growth (shift from glucose to lactose utilization) of *E. coli*. If glucose is in the growth medium, the operon is always OFF because glucose inhibits cyclic adenosine monophosphate (cAMP) and lactose transport into the cell. If glucose is absent, the operon would remain OFF unless some lactose is present inside the cell (which is true when glucose is depleted and lactose from the outside can now enter the cell); an initially small amount of internal lactose increases rapidly because of at least two positive feedback loops as shown in Figure 2(b). It is the positive feedback loop involving lactose transport that ultimately controls the influx of lactose.

For the formal representation of the lac operon according to Figure 1, the vectors of variables corresponding to the model shown in Figure 2(b) are the following:  $\mathbf{G} = [\mathbf{G}_{ZYA} \mathbf{G}_I \mathbf{G}_{CRP}]^T$ , where  $\mathbf{G}_{ZYA}$  is the base sequence on DNA that includes genes *lacZ*, *lacY*, and *lacA* and transcribed as one transcriptional unit ( $[\ ]^T$  means ‘transpose’);  $\mathbf{G}_I$  is the DNA sequence containing gene *lacI*, and

$G_{CRP}$  is the transcribed DNA sequence containing gene *CRP*.  $\mathbf{R}_o = [R_{ZYA,o} \ R_{I,o} \ R_{CRP,o}]^T$  is the vector of primary transcripts;  $\mathbf{R} = [R_Z \ R_Y \ R_I \ R_{CRP}]^T$  is the vector of mature transcripts. Note that the transcript  $R_A$  (corresponding to gene A) is not included because it is not considered further in the dynamics of the GRN.  $\mathbf{P}_o = [P_{Z,o} \ P_{Y,o} \ P_{I,o} \ P_{CRP,o}]^T$  is the vector of primary protein translates;  $\mathbf{P} = [P_{4Z} \ P_{Ym} \ P_{4I} \ P_{CRP,cAMP}]^T$  is the vector of mature, modified, and active proteins; the protein  $P_Z$  ( $\beta$ -galactosidase) is tetrameric in its functional form, the permease  $P_Y$  acts at the plasma membrane (hence the subscript “m” in  $P_{Ym}$ ), the repressor protein  $P_I$  is tetrameric, and CRP’s binding with cAMP is necessary for its DNA-binding activity.  $\mathbf{M} = [\text{Glu} \ \text{Lac} \ \text{Allo} \ \text{cAMP}]^T$  is the vector of metabolites (Glu = glucose, Lac = lactose, Allo = allolactose, cAMP = cyclic adenosine monophosphate). The GRN for the lac operon model with the representation of Figure 1 is shown in Figure 3. The first equation in [1] would look like this:

$$\begin{pmatrix} dR_{ZYA,o}/dt \\ dR_{I,o}/dt \\ dR_{CRP,o}/dt \end{pmatrix} = \begin{pmatrix} \tau_{11} & 0 & 0 \\ 0 & \tau_{22} & 0 \\ 0 & 0 & \tau_{33} \end{pmatrix} \begin{pmatrix} G_{ZYA} \\ G_I \\ G_{CRP} \end{pmatrix} - \begin{pmatrix} \rho_{11} & 0 & 0 \\ 0 & \rho_{22} & 0 \\ 0 & 0 & \rho_{33} \end{pmatrix} \begin{pmatrix} R_{ZYA,o} \\ R_{I,o} \\ R_{CRP,o} \end{pmatrix} - \begin{pmatrix} \delta_{11} & 0 & 0 \\ 0 & \delta_{22} & 0 \\ 0 & 0 & \delta_{33} \end{pmatrix} \begin{pmatrix} R_{ZYA,o} \\ R_{I,o} \\ R_{CRP,o} \end{pmatrix} \quad [2]$$

where  $\tau_{11}$  would be a function of  $P_I$  and  $P_{CRP,cAMP}$ . For example, we could choose the function  $\tau_{11} = (c_1 + c_2 P_{CRP,cAMP}) / (c_3 + c_4 P_I^n)$  to represent the



**Figure 3** The lac operon in accordance with the scheme shown in Figure 1. See text for details.



activation of transcription by the protein complex  $P_{CRP.cAMP}$  and inhibition by the tetrameric repressor  $P_1$  (the  $n$  and  $c_i$ 's are constant parameters;  $n$  should be greater than 1 because of the tetrameric complex of  $P_1$ ).

New mathematical models and reviews on the lac operon have appeared recently.<sup>2-4</sup> Yildirim and Mackey<sup>2</sup> used delay differential equations to account for the transcriptional and translational steps that are missing in their model. An earlier detailed kinetic model was proposed and analyzed by Wong, Gladney and Keasling.<sup>1</sup> Recently, Vilar, Guet, and Leibler<sup>4</sup> used a four-variable model that captures many of the essential dynamics of the lac operon. Note that the Vilar–Guet–Leibler model is essentially a three-variable model. The bistability exhibited by the model was the explanation for the ON-OFF behavior of the lac operon.

At the single-cell level, the operon is either ON or OFF (all-or-nothing induction) as shown in the recent experimental report of Ozbudak et al.<sup>3</sup> These authors exploited the positive feedback loop between the permease ( $y$ ) and the inducer ( $x$ ), and they used the following mathematical model to represent the positive feedback loop:

$$\begin{aligned}\tau_y dy/dt &= \alpha(1/[1 + R/R_o]) - y \\ \tau_x dx/dt &= \beta y - x\end{aligned}\tag{3}$$

where  $R/R_T = 1/[1 + (x/x_0)^n]$  and  $R$  = concentration of active LacI,  $R_o$  = initial concentration of active LacI,  $R_T$  = total concentration of LacI tetramers,  $x_o$  = initial concentration of LacY (permease), and the rest of the symbols are parameters. The parameter  $n$  allows consideration that the repressor is a tetramer. This simple model generates bistability in which the all-or-nothing transition is associated with a saddle-node bifurcation. The simple set of Eq. [3] was useful in guiding the authors' experiments in showing ON-OFF behavior as well exploring the phase diagram (coordinates of which are the variables  $x$  and  $y$ , for example) for bistable and monostable regions.

The lac operon illustrates several important points in modeling GRNs. Although the operon structure is not a general property of all genomes, we can expect that genomic DNA sequence organization affects the dynamics of the GRN; this is primarily because of coexpression of genes found in the same transcriptional units or coregulation of genes by transcription factors that recognize promoter regions having similar regulatory sequences. Another lesson from the lac operon is that abstraction of the complex GRN may be sufficient to understand the behavior of the system. This abstraction was facilitated by prior knowledge of the influence of network topology on dynamical behavior, e.g., bistability originating from positive feedback loops.<sup>8</sup> A discussion on how network structure alone influences system behavior is provided in the penultimate section of this chapter.

---

## HIERARCHIES OF GRN MODELS: FROM PROBABILISTIC GRAPHS TO DETERMINISTIC MODELS

The general representation of GRNs in Figure 1 considers groups of molecules according to their chemical classes (DNA, RNA, proteins, metabolites) whose “interactions” merely encode the broad concepts of transcription, posttranscriptional processing, translation, and posttranslational modifications. Depending on the nature of available experimental information, specific models of gene regulatory networks can be constructed at various levels of detail.

Networks, in general, are described by their graphical structures. A graph is basically a set of “nodes” and a set of “edges,” the latter being the representations of the interactions or associations among nodes. Progressively more detailed mechanistic information can be added to a graph as they become available. At one extreme of the spectrum of models, the nodes in the graph could be just a set of genes (and no other kinds of objects), with certain pairs of genes linked by undirected edges if these pairs are known to “interact” or are “associated” in some way. Sometimes the nodes could be proteins, and the edges represent physical interactions. Because of the correspondence between proteins and genes (albeit not generally one-to-one), protein–protein interaction networks may imply some underlying GRN structure. In general, nodes in a graph can be defined according to the level of detail that is sufficient to describe a particular feature, function, or behavior of the system. For example, the nodes in Figure 1 represent various classes of molecules. A node could also represent a subnetwork or module with specific cellular function.

An edge of a graph is assigned a direction if information on causality exists, i.e. that one node affects the state of the other. A directed edge can be further characterized as either “activating” or “inhibiting.” As more quantitative data are available, it may be possible to identify the “strength” of an edge. For dynamic models, the strength of an edge would, for example, require identification of rate expressions as functions of the states of the nodes. At this point, a dynamic model encoded in deterministic differential equations is possible. Finally, at the other extreme in the spectrum of GRN models, microscopic details of the interactions between individual molecular species are known and molecular dynamics simulations are possible.

As the example of the lac operon illustrates, abstract models involving differential equations that do not necessarily reflect the detailed mechanism are sometimes used by scientists when the goal is primarily to explore possible system dynamics originating from the structure of the network. Associated with the process of “abstraction” is the problem of reducing the network into a smaller set of “modules” and their interactions. Modules can range from individual molecules or genes, to a set of genes or proteins, or to functional

subnetworks with definable cellular functions. Similar ideas have been discussed recently by Vilar et al.<sup>4</sup> in their work on the lac operon. The lac operon is an example of a well-defined small model system in which a considerable amount of biological knowledge and mechanistic understanding have already accumulated so that refined mathematical modeling can be carried out. Many other focused models and corresponding mathematical model formalisms have been reviewed recently by de Jong.<sup>9</sup> In contrast, constructing the network graph of gene interactions from large-scale gene expression measurements is just beginning and is, at times, controversial. As this field has been reviewed<sup>10–13</sup> recently, only a brief account is given below.

High-throughput gene expression measurements with DNA microarrays provide global snapshots of the dynamics of gene networks at the RNA level. Expression data are intrinsically noisy, and conclusions derived from them are probabilistic in nature. Furthermore, the mRNA levels are averages from cell populations. Gene network reconstruction from microarray data also suffers from the so-called “dimensionality problem”<sup>11</sup> because the number of genes is much greater than the number of microarray experiments. Statistical analysis of gene expression data usually involve clustering methods to find genes with similar expression patterns across time series or across different experimental conditions (e.g., see D’haeseleer et al.<sup>14</sup> and Eisen et al.<sup>15</sup>). The assumption is that clustered or coexpressed genes are somehow coregulated or perhaps share similar functions. The results of clustering for GRN modeling could therefore be a coarse-grain network composed of modules (nodes), each module representing a set of genes with similar functions.

Graphical models that combine probability theory and graph theory are suitable frameworks for inferring GRNs from gene expression data.<sup>10,16</sup> In general, these graphical models are probability models for multivariate random variables whose independence structure can be represented by a conditional independence graph. Recently, Friedman<sup>10</sup> reviewed the field of probabilistic graphical models for gene networks, including Bayesian networks. In a Bayesian network, the nodes represent random variables (e.g., genes and their expression levels), whereas the edges show conditional dependence relations. Husmeier<sup>17,18</sup> has also reviewed the applications of Bayesian networks to microarray data. Bayesian networks were first applied to the problem of reverse engineering of GRNs from microarray expression data by Friedman et al.,<sup>19</sup> Pe’er et al.,<sup>20</sup> and Hartemink et al.<sup>21</sup> Other examples of graphical models involving various statistical methods are discussed by Wang, Myklebost and Hovig.<sup>16</sup>

Zak et al.<sup>22</sup> have argued that inferring the GRN structure from expression data alone is impossible. However, promising results come from more recent work showing that properly designed perturbation experiments do permit network reconstruction (see Stark et al.,<sup>12,13</sup> Husmeier,<sup>18</sup> de La Fuente et al.,<sup>23</sup> Kholodenko et al.,<sup>24</sup> and Gardner et al.<sup>25</sup>). Two papers<sup>23,24</sup> extended ideas from metabolic control analysis to suggest perturbation experiments

designed to determine the direction and strengths of interactions between genes. Also, Gardner et al.<sup>25</sup> used systematic perturbations combined with least-squares regression to infer the gene network topology and weights of interactions.

In general, the issues encountered during the creation of a GRN graph are similar to those faced when designing a pathway or interaction database. These issues will be discussed in more detail in the section on pathway ontology. An extensive discussion on this ontology is provided because it is a crucial stepping stone for future projects concerned with the extraction of GRN models from pathways databases. Pathways databases are relatively recent developments in bioinformatics. These databases are built from more elementary databases, and it is important to be aware of the many heterogeneous bioinformatics resources available, most of them on the Internet. Thus, a brief guide is given next.

---

## A GUIDE TO DATABASES AND KNOWLEDGBASES ON THE INTERNET

The field of bioinformatics has naturally originated to cope with the deluge of data generated by high-throughput technologies in genomics, transcriptomics, proteomics, and other omics. These data are organized into databases (DBs) and knowledgebases (KBs), many of which are publicly available on the Internet. Comprehensive and realistic modeling of GRNs should tap into the information contained in these DBs and KBs. Thus, it is expected that the next generation of modelers will have to be sufficiently aware of bioinformatics resources. For this reason, an overview of the major bioinformatics DBs and KBs is provided here, although their utility for modeling GRNs may not be direct and obvious at this time. It was alluded to in the discussion of the lac operon that understanding the operon structure of the genomic DNA was necessary to understand the dynamics of the network. In general, relating genome organization to GRN dynamics is a difficult and still a much open problem. This section begins with genomic sequence databases in anticipation of their future use in helping scientists predict GRN structures; a specific example would be that of finding regulatory sequences where transcription factors bind, thereby linking one gene product to the transcription of another gene.

To date, the genomes of more than 150 organisms have been sequenced, and many more sequencing projects are currently going on or planned. Publicly available DNA sequence data as well as functional and structural data on proteins are accumulating at an exponential rate, virtually doubling every year. The major sequence and structure repositories that are regularly updated are listed in Table 1.

The partners of the *International Nucleotide Sequence Databases (INSD)*, namely *GenBank*, *EMBL*, and *DDBJ*, share their nucleic acid

**Table 1** Major Sequence and Structure Repositories

Database	Description	URL
GenBank	Repository of all publicly available annotated nucleotide and protein sequences	<a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a>
EMBL Database	Repository of all publicly available annotated nucleotide and protein sequences	<a href="http://www.ebi.ac.uk/embl.html">http://www.ebi.ac.uk/embl.html</a>
DDBJ (DNA Data Bank of Japan)	Repository of all publicly available annotated nucleotide and protein sequences	<a href="http://www.ddbj.nig.ac.jp">http://www.ddbj.nig.ac.jp</a>
PIR	Protein information resource: Protein sequence database	<a href="http://pir.georgetown.edu/">http://pir.georgetown.edu/</a>
Swiss-Prot	Highly annotated curated protein sequence database	<a href="http://www.expasy.org/sprot">http://www.expasy.org/sprot</a>
PDB	Protein structure databank: Collection of publicly available 3-D structures of proteins and nucleic acids	<a href="http://www.rcsb.org/pdb">http://www.rcsb.org/pdb</a>

sequence data for a comprehensive coverage of all available genome information. *Swiss-Prot* is a manually curated protein sequence database with a high-level annotation of protein function and protein modifications, including links to property, structure, and pathways databases. *PIR* is similar to *Swiss-Prot*, with the former providing some options for sequence analysis. Recently, *UniProt Knowledgebase* (<http://www.uniprot.org>) was established with the aim of unifying and linking protein databases with cross-references and query options.

Some of the major protein sequence and structure property databases are listed in Table 2. Although many more general or specialized property databases are available,<sup>26</sup> the list given in Table 2 is a good start for exploring protein property databases. Table 3 gives a list of gene expression repositories.

It is difficult for one person to keep up with the rapidly increasing number of genomics, proteomics, and interactomics and metabolomics databases, let alone their intended usage.<sup>26</sup> To alleviate this problem, an increasing number of integrated database retrieval and analysis systems tools are being developed for the purpose of data management, acquisition, integration, visualization, sharing, and analysis. Table 4 lists promising examples of these tools, which are regularly maintained and updated. *GeneCards* is an integrated database of human genes, genomic maps, proteins, and diseases, with software that retrieves, combines, searches, and displays human genome information. *GenomeNet* is of particular interest because its analytical tools are tightly linked with the *KEGG* pathways database (discussed in the next section).

**Table 2** Protein Sequence and Structure Property Databases

Database	Description	URL
eMOTIF	Protein sequence motif database	<a href="http://motif.stanford.edu/emotif">http://motif.stanford.edu/emotif</a>
InterPro	Integrated resource of protein families, domains	<a href="http://www.ebi.ac.uk/interpro">http://www.ebi.ac.uk/interpro</a>
iProClass	Integrated protein classification database	<a href="http://pir.georgetown.edu/iproclass/">http://pir.georgetown.edu/iproclass/</a>
ProDom	Protein domain families	<a href="http://www.toulouse.inra.fr/prodom.html">http://www.toulouse.inra.fr/prodom.html</a>
CDD	Conserved domain database: Covers protein domain information from Pfam, SMART, and COG databases	<a href="http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml">http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml</a>
CATH	Protein structure classification database	<a href="http://www.biochem.ucl.ac.uk/bsm/cath/">http://www.biochem.ucl.ac.uk/bsm/cath/</a>
CE	Repository of 3-D protein structure alignments	<a href="http://cl.sdsc.edu/ce.html">http://cl.sdsc.edu/ce.html</a>
SCOP	Structural classification of proteins	<a href="http://scop.mrc-lmb.cam.ac.uk/scop">http://scop.mrc-lmb.cam.ac.uk/scop</a>

*ToolBus* comprises several data analysis software platforms such as multiple sequence alignment, phylogenetic trees, generic XML viewer, pathways, and microarray analysis, which are linked to each other as well as to major databases. *SRS* and *NCBI* serve as general data retrieval portals as well as provide links to specific analysis tools.

**Table 3** Gene Expression Databases

Database	Description	URL
ArrayExpress	Microarray gene expression data collection database	<a href="http://www.ebi.ac.uk/arrayexpress">http://www.ebi.ac.uk/arrayexpress</a>
CIBEX	Center for Information Biology gene: A public repository for high-throughput experimental data in gene expression	<a href="http://cibex.nig.ac.jp">http://cibex.nig.ac.jp</a>
GeneNote	Database of human genes expression profiles in healthy tissues	<a href="http://genecards.weizmann.ac.il/gene-note/">http://genecards.weizmann.ac.il/gene-note/</a>
GEO	Gene Expression Omnibus: A high-throughput gene expression data repository	<a href="http://ncbi.nlm.nih.gov/geo">http://ncbi.nlm.nih.gov/geo</a>
SMD	Stanford Microarray Database: Raw and normalized data from microarray experiments	<a href="http://genome-www.stanford.edu/microarray">http://genome-www.stanford.edu/microarray</a>

**Table 4** Integrated Database Retrieval and Analysis Systems

Database	Description	URL
GeneCards	Database of human genes, proteins and their involvement in diseases	<a href="http://bioinfo.weizmann.ac.il/cards">http://bioinfo.weizmann.ac.il/cards</a>
GenomeNet	Network of database and computational services for genome research	<a href="http://www.genome.ad.jp/">http://www.genome.ad.jp/</a>
NCBI	Retrieval system for searching several linked databases	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
PathPort/Tool-Bus	Collection of web-services for gene prediction and multiple sequence alignment, along with visualization tools	<a href="https://www.vbi.vt.edu/pathport">https://www.vbi.vt.edu/pathport</a>
SRS-EBI	Integration system for both data retrieval and applications for data analysis	<a href="http://srs.ebi.ac.uk">http://srs.ebi.ac.uk</a>

## PATHWAYS DATABASES AND PLATFORMS

Along with recent advances in genomics and proteomics, requirements for analysis, expansion, and visualization of cell signaling, GRNs and protein–protein interaction maps are leading to the development of data representation and integration tools. Pathways databases can be classified into four groups according to their interactome data content and representation as listed in Table 5. Only those websites that are regularly maintained are included in the list. The first group of databases represents *binary interaction databases*. *BIND*, *DIP*, and *MINT* document experimentally determined protein–protein interactions from peer-reviewed literature or from other curated databases. *BIND* and *MINT* store experimental conditions used to observe the interaction, chemical action, kinetics, and other information linked to the original research articles.

*Static image databases* are very good sources of pathway diagrams that provide a broad introductory view of cell regulatory pathways along with good reviews and links. *ACSF*, *STKE*, and *Biocarta* are comprehensive knowledgebases on signal transduction pathways and other regulatory networks.

*Metabolic signaling databases* contain detailed information on metabolic pathways. These DBs have well-established data structures but have nonuniform ontologies. *BioCyc* is a collection of pathway/genome databases for many bacteria and up to 14 species of other organisms. Enzyme catalyzed

**Table 5** Pathways Databases and Platforms

	Database	Description	URL
Binary interactions	BIND	Biomolecular interaction network database	<a href="http://www.bind.ca">http://www.bind.ca</a>
	BindingDB	Collection on experimental data on the noncovalent association of molecules in solution	<a href="http://www.bindingdb.org">http://www.bindingdb.org</a>
	BRENDA	Enzyme Information System: Sequence, structure, specificity, stability, reaction parameters, isolation data, and molecular functions ontology	<a href="http://www.brenda.uni-koeln.de">http://www.brenda.uni-koeln.de</a>
	DIP	Database of interacting proteins	<a href="http://dip.doe-mbi.ucla.edu">http://dip.doe-mbi.ucla.edu</a>
	IntAct project	Public repository for annotated protein–protein interaction data	<a href="http://www.ebi.ac.uk/intact">http://www.ebi.ac.uk/intact</a>
	InterDom	Putative interacting protein domain database derived from multiple sources	<a href="http://interdom.lit.org.sg">http://interdom.lit.org.sg</a>
Static images	MINT	A molecular interaction database	<a href="http://mint.bio.uniroma2.it/mint/">http://mint.bio.uniroma2.it/mint/</a>
	ACSF	Signaling resource for signal transduction elements	<a href="http://www.signaling-gateway.org/">http://www.signaling-gateway.org/</a>
	BioCarta	Molecular relationship map pages from areas of active research	<a href="http://www.biocarta.com">http://www.biocarta.com</a>
Metabolic signaling	STKE	Signal transduction knowledge environment	<a href="http://stke.org/">http://stke.org/</a>
	BRITE	Biomolecular relations in information transmission and expression	<a href="http://www.genome.ad.jp/brite">http://www.genome.ad.jp/brite</a>
	KEGG	Kyoto encyclopedia of genes and genomes: Molecular interaction networks of metabolic and regulatory pathways	<a href="http://www.genome.ad.jp/kegg">http://www.genome.ad.jp/kegg</a>
	BioCyc	A collection of databases that describes the genome and metabolic pathways of a single organism	<a href="http://biocyc.org/">http://biocyc.org/</a>
Regulatory signaling	PathDB	A data repository and a system for building and visualizing cellular networks	<a href="http://www.ncgr.org/pathdb">http://www.ncgr.org/pathdb</a>
	aMAZE	A system for the representation, annotation, management, and analysis of biochemical and gene regulatory networks	<a href="http://www.amaze.ulb.ac.be/">http://www.amaze.ulb.ac.be/</a>
	Cytoscape	Software platform for visualizing molecular interaction networks	<a href="http://www.cytoscape.org/">http://www.cytoscape.org/</a>
	GeneNet	Database on gene network components and a program for the data visualization	<a href="http://www.mgs.bionet.n-sc.ru/mgs/gnw/genenet">http://www.mgs.bionet.n-sc.ru/mgs/gnw/genenet</a>
	PATIKA	Software platform for pathway analysis tool for integration and knowledge acquisition	<a href="http://www.patika.org/">http://www.patika.org/</a>
	Pathway Assist	Tool for analysis, expansion and visualization of biological pathways, gene regulation networks, and protein interaction maps	<a href="http://www.ariadnegenomics.com/products/pathway.html">http://www.ariadnegenomics.com/products/pathway.html</a>
	TRANS-PATH	Gene regulatory network and microarray analysis system	<a href="http://www.biobase.de/pages/products/databases.html">http://www.biobase.de/pages/products/databases.html</a>



reactions, or the gene that encodes that enzyme or the structures of chemical compounds in pathways and reactions, can be displayed by *BioCyc* ontology-based software for a given biochemical pathway. In addition, *BioCyc* supports computational tools for simulation of metabolic pathways.

*KEGG* is a frequently (daily) updated group of databases for the computerized knowledge representation of molecular interaction networks in metabolism, genetic information processing, environmental information processing, cellular processes, and human diseases. The data objects in the *KEGG* databases are all represented as graphs and various computational methods for analyzing and manipulating these graphs are available.

The fourth category of the DBs and software platforms listed in Table 5 is concerned with *regulatory signaling* networks. *GeneNet*, *aMAZE*, and *PATIKA* have very similar ontologies for representing and analyzing molecular interactions and cellular processes. *PATIKA* and *GeneNet* provide graphical user interfaces for illustrating signaling networks. The *aMAZE* tool called *LightBench*<sup>27</sup> allows users to browse information stored in the database, which covers chemical reactions, genes and enzymes involved in metabolic pathways, and transcriptional regulation. Another *aMAZE* tool called *SigTrans* is a database of models and information of signal transduction pathways.

Both *GeneNet* and *PATIKA* are composed of a server-side with a database and client-side. In addition to its database components, a *PATIKA* client-side editor software provides an integrated, multiuser environment for visualizing, entering, and manipulating networks of cellular events independent of an additional web-browser.

*Cytoscape* and *PathwayAssist* are similar software tools for automated analysis, integration, and visualization of protein interaction maps. In these tools, automated methods for mining *PubMed* and other public literature databases are incorporated to facilitate the discovery of possible interactions or associations between genes or proteins.

---

## ONTOLOGIES FOR GRN MODELING

Bioinformatics is now moving toward the direction of creating tools, languages, and software for the integration of heterogeneous biological data and their analysis at the level of cellular systems and beyond. This direction requires establishing appropriate “ontologies” to annotate the various parts and events occurring in the system. An ontology is a set of controlled and unambiguous vocabulary for describing objects and concepts.<sup>28</sup>

### Current Gene, Interaction, and Pathway Ontologies

At the genome level, the Gene Ontology<sup>TM</sup> (GO) Consortium (<http://www.geneontology.org>) introduced a comprehensive bio-ontology that is

aimed to cover genes in all organisms. GO provides unique identifiers for each concept related to “molecular function,” “biological process,” and “cellular component” searchable through the *AmiGO* tool (<http://www.godatabase.org>). Note that these three concepts (especially the concept of “biological process”) can be interpreted in terms of memberships of genes in cellular pathways; hence, GO can be considered as part of a pathway ontology.

A conventional approach for representing cellular pathways is the application of static diagrams such as those found in the websites of *ACSF*, *BioCarta*, and *STKE* (see Table 5). These diagrams are often not reusable, and the pathway representations are far from being uniform and consistent among different websites, because the various representations carry implicit conventions rather than explicit rules as required by formal ontologies. Because pathways are basically composed of components and steps or processes, the development of *interaction databases* is a logical first step (see sample databases in Table 5). These databases provide a diverse amount of binary interaction data, which could then be used by scientists for building networks.

Among the cellular pathways, metabolic pathways are generally more detailed and structured because of more advanced knowledge about metabolism in cells (see Table 5). In all of these databases, the proteins are classified according to the Enzyme Commission list of enzymes (EC numbers). These metabolic DBs have strict ontologies that are focused on protein activities relevant to metabolic pathways. Because of a detailed knowledgebase and ontology, metabolic pathways are amenable to kinetic modeling and computer simulations.<sup>29</sup>

## Whole-Cell Modeling Platforms

Several whole-cell modeling and simulation software environments exist (e.g., *Virtual Cell*, *E-Cell* and *CellWare*) with their specific ontologies. *Virtual Cell*<sup>30</sup> provides a subcellular localization-based visual environment for modeling cellular events. The ontology is mainly based on a single mechanistic physiological model that encodes the general structure and function of a cellular event such as release of calcium and its effects on the cell. In *Virtual Cell*, a cell is considered as distinct geometrical subdomains containing specific cellular components with known concentration. This model allows users to proceed through *Virtual Cell* simulation tools. Even though *Virtual Cell* has some applications relevant to GRNs, the platform may have difficulties in modeling events that occur only in one cellular compartment with unknown molecular concentrations.

*E-Cell*<sup>31,32</sup> is a generic software platform for visualization, modeling, and simulation of whole cell events. *E-Cell* provides several graphical interfaces for user-definable models of certain cellular states. A cell model can be constructed with three classes of objects (entities): substances, genes, and

reaction rules. The *E-Cell* ontology shares several similarities with the *PATIKA* ontology, which is discussed in the next section.

*CellWare*<sup>33</sup> is a multi-algorithmic software platform for modeling and simulation of cellular events. It has several toolboxes including tools for user-dependent model description, definition, and construction with a graph editor. A simulation toolbox contains various simulation algorithms and interfaces from which a user can choose.

## Ontology for Modeling Multiscale and Incomplete Networks

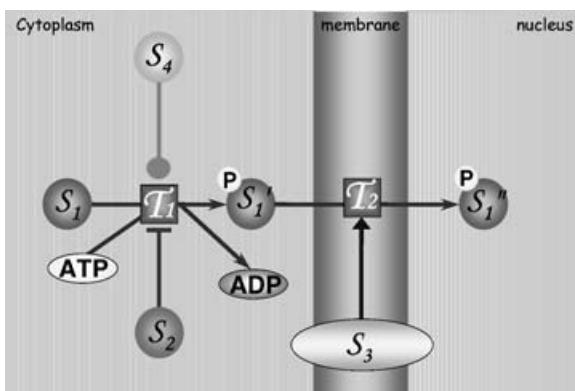
The current state of our knowledge on cellular regulatory pathways is still fragmented, incomplete, and uncertain in many respects despite accumulating data. A pathway ontology should be able to represent available information even when it is incomplete, thus allowing incremental construction of pathways. In addition, the ontology must have the flexibility for continuous modification of data without compromising the integrity of the network being built. Therefore, the ontology must describe integrity rules of the pathway data, enabling the construction of a robust model of the system. A data integrity rule should state that for every instance of a bioentity (see below), a primary key with an accession number (such as *SwissProt* ID) must exist and be unique. The seamless integration of various hierarchies of detail or scale is a key problem in modeling and in the representation of complex systems like a cell.

Pathway visualization with diagrams or graphs facilitates the creation of a mathematical model of a GRN. An efficient visualization scheme is generated when an ontology uses intuitive images. The ontology should offer ways to reduce the complexity of the information at some stage of the modeling process.

The discussion in the next subsection focuses on an ontology that is suitable for modeling incomplete information and abstractions of varying levels of complexity. This ontology has been recently implemented in a pathway database tool named *PATIKA* (Pathway Analysis Tool for Integration and Knowledge Acquisition).<sup>34,35</sup> The *Pathway Database System (PDS)* developed by Krishnamurthy et al.<sup>36</sup> shares several basic similarities with *PATIKA* for database organization and visualization. As in *PATIKA*, *PDS* provides tools for modeling, storing, analyzing, visualizing, and querying biological pathways. However, *PDS* does not define a formal ontology for GRNs but instead follows the rules of *KEGG* metabolic pathway ontology and includes *KEGG* data.

## An Ontology for Cellular Processes

*States and bioentities.* Components of a GRN are macromolecules (e.g., DNAs, RNAs, or proteins), small molecules (e.g., ions, GTP, or ATP), or

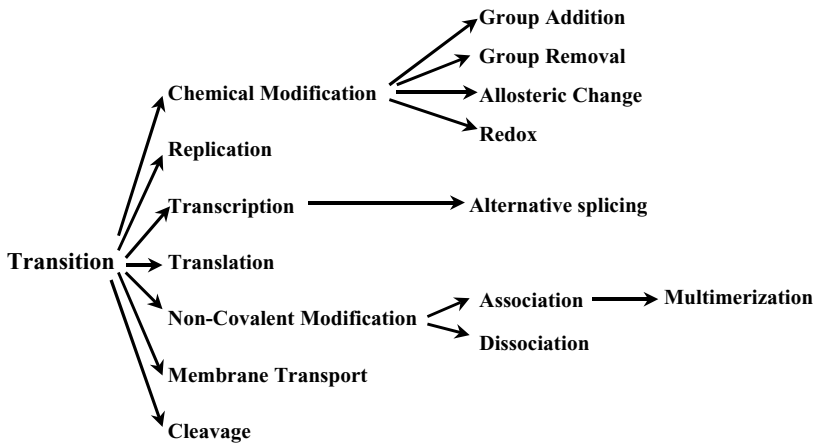


**Figure 4** An illustration of the basic features of the *PATIKA* ontology. States, transitions, and interactions are represented by circles, rectangles, and lines, respectively. The bioentity “ $S_1$ ” has three states (namely,  $S_1$ ,  $S'_1$ , and  $S''_1$ ) located in two distinct subcellular compartments (cytoplasm and nucleus), which are separated by a third compartment, the nuclear membrane.  $S_1$  and  $S'_1$  are both in the cytoplasm.  $S_1$  is phosphorylated through transition  $T_1$  giving rise to a new state, the phosphorylated  $S'_1$ .  $S'_1$  is translocated to the nucleus through transition  $T_2$  and becomes  $S''_1$ .  $T_1$  has two effector states,  $S_2$  (inhibitor) and  $S_4$  (unspecified effect).  $T_2$  has an activator type of effector ( $S_3$ ) representing, for example, the nuclear pore complex.

physical events (e.g., heat, radiation, or mechanical stress). Often, these players share a common synthesis pathway and/or are chemically similar. For example, the p53 protein has many *states*, including its native, phosphorylated, nuclear, or MDM2-bound forms. These states are represented as nodes in the network graph, while maintaining their biological or chemical groupings under a common *bioentity*.

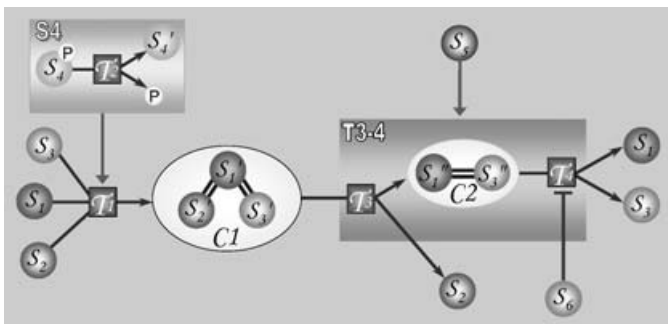
**Transitions.** A transition represents a cellular event, and each is represented as a separate node in the graph (see Figures 4 and 5). A state may go through a certain transition, may be produced by a transition, or may affect a transition as being an activator or inhibitor. When a transition occurs, all of its products are generated.

**Compartments.** Transitions also include transport of molecules between cell compartments. The set of transitions that a state can be involved in is strictly related to its compartment; accordingly, a change in the compartment means a change in the state’s information context. The state’s compartment is a part of the ontology. As the compartments and their vicinity are cell-type dependent, compartmental structure can be modeled as part of the ontology. Cell membranes create an additional complexity because not only can a molecule be located completely inside the membrane, but it may also communicate with both sides of the membrane as part of the events involved in adjacent compartments. So membranes are considered as separate compartments in the ontology.



**Figure 5** Proposed tree structure that classifies transitions in the *PATIKA* ontology. If the nature of a transition cannot be defined in the existing ontology, it can be considered as generic transition to be defined and added in the ontology.

*Molecular complexes.* In biological systems, molecules often form complexes to perform certain tasks (Figure 6). Each member of a molecular complex can be considered as a new state of its associated bioentity. The intrinsic specific binding relations affect the function of a molecular complex. Therefore, these binding relations must be represented in the model ontology. Moreover, members of a molecular complex may independently participate in different transitions; thus, we should be able to address each member



**Figure 6** A pathway containing two abstractions and a molecular complex *C1* (composed of three states  $S_1$ ,  $S_2$ , and  $S_3$ ). Superstate  $S_4$  is an example of an abstraction in which the state  $S_4$ -P or  $S_4'$  may act as an activator of transition  $T_2$ .  $S_5$  leads to the dissociation of complex *C1* acting on either before or after the dissociation of  $S_2$ . Therefore,  $S_5$  may be an activator of either  $T_3$  or  $T_4$ ; thus,  $S_5$  is illustrated as the activator of supertransition  $T_{3-4}$ .

individually (Figure 6). In addition, a molecular complex may contain members from neighboring compartments (e.g., receptor-ligand complexes).

*Abstractions.* Various levels of abstractions are involved in the analysis of complex cellular events. A set of transitions can be described as a single “process” (e.g., the MAPK pathway), and a set of related processes may be classified under one “cellular mechanism” (e.g., apoptosis). Some explicit examples of abstractions are shown in Figure 6. In cases in which it is not identified which state among a set of states constitutes the substrate or effector of a transition, or in which target transition of an effector is unclear, we may need to abstract these states (or transitions) as a single state (or transition) to represent the available information despite its incomplete nature.

## The PATIKA Pathway Ontology

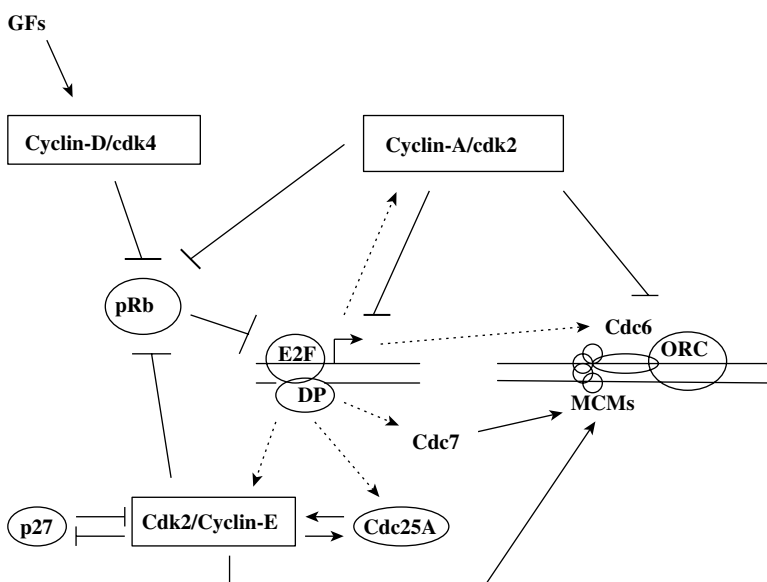
A pathway is an abstraction of a certain biological event and is the primary abstraction in the *PATIKA* ontology.<sup>35</sup> The context of this abstraction can change from a single molecule–molecule interaction to a complete network of all interactions in a cell. In *PATIKA*, a pathway is represented by a *pathway graph*, which is a compound graph.<sup>37</sup> A *pathway graph* is defined by an interaction graph  $G = (V, E)$  along with several rules on the topology;  $V$  is the union of a finite set of *states*  $V_s$  and a finite set of *transitions*  $V_t$ .  $E$  is the union of *interactions* of five sets: *substrate* edges  $E_s$ , *product* edges  $E_p$ , *activator effector* edges  $E_a$ , *inhibitor effector* edges  $E_i$ , *effector of unknown type* edges  $E_u$ , and each directed edge belonging to either  $V_t \times V_s$  (for product edges) or to  $V_s \times V_t$  (for remaining interaction edge types). Every state has a defined *type*: DNA, RNA, protein, small molecule, or physical factor. States are also associated with a specific *compartment*. Identical states in different compartments are considered as separate states. States of the same biological origin and/or similar chemical structure are grouped under a *biological entity* or simply *bioentity* that acts as state and transition connectivity data holders in *PATIKA*.

Every transition must be affiliated with at least one substrate and one product edge. It may have an arbitrary number of effectors, a combination of which defines the exact behavior for the transition. Transitions are classified according to the tree shown in Figure 5. A transition is not associated with a specific compartment; instead, its compartment is determined by its interacting states. Different types of molecules (e.g., protein, DNA, and RNA) have distinct user interfaces for easier visual discrimination in *PATIKA*. Compartmental information is also modeled. *PATIKA* also implements collaborative construction and modification to existing regulatory signaling data on the database. Therefore, *PATIKA* maintains version numbers as part of the ID of each graph object. Thus, although a user is working on a *PATIKA* graph locally, others might change the topology and/or properties of states and transitions in the *PATIKA* database.

## EXTRACTING MODELS FROM PATHWAYS DATABASES

A clear pathway ontology, as discussed in the previous section, will allow systematic methods for extracting GRN models from the interactions stored in a pathways database. The specific model would, of course, depend on the particular biological question being asked. Here, a brief example is given of how a model is extracted from a network of interactions taken from some of the databases listed in Table 5. The work of Aguda and Tang<sup>38</sup> on the G1 checkpoint of the cell cycle is an example. A cell cycle checkpoint is a surveillance mechanism that arrests or slows down cell cycle progression if something goes wrong, e.g., DNA damage. The significance of elucidating the control mechanism of the G1 checkpoint lies in the observation that many human cancers are associated with nonfunctional G1 checkpoints.

A qualitative network of the G1-S transition is shown in Figure 7. The network was generated by integrating information from the published literature, including sequence analysis of upstream regulatory regions of genes that are targeted by the E2F transcription factor family. Aguda and Tang<sup>38</sup> were interested in finding a minimal subnetwork that is sufficient to explain the



**Figure 7** A qualitative network involving key interactions in the G1-S transition of the mammalian cell cycle. Solid lines are posttranslational modifications or protein–protein interactions. Dashed arrows are transcriptional steps. Arrows mean “activation,” and hammerheads mean “inhibition.” GFs = growth factors, cdk = cyclin-dependent kinase, pRb = retinoblastoma protein, ORC = origin recognition complex.

switching behavior of the G1 checkpoint. The key step toward finding this subnetwork was the hypothesis that a core set of interactions with an intrinsic instability ultimately generates a switching behavior (see Aguda and Tang<sup>38</sup> and Aguda<sup>39</sup> for details; network stability analysis is discussed in the next section). Experimentally, the activity of cyclin E/CDK2 is a marker for the entry into the S phase of the cell cycle. Hence, this minimal set of interactions must include cyclin E/CDK2.

In the network graph shown in Figure 7, the arrows are interpreted as “activation” and the hammerheads as “inhibition.” From this qualitative network, a network stability analysis pointed to a core mechanism involving cyclin E/CDK2, Cdc25A, p27Kip1, and their interactions. These interactions involve two coupled positive feedback loops, namely, between the pair (Cdk2/Cyclin E, Cdc25A) and the pair (Cdk2/Cyclin E, p27). This core mechanism was then the basis for a more detailed mechanistic model. The dynamics of the model was coded into differential equations and solved in a computer. The computer simulations reproduced the experimentally observed qualitative behavior of the G1 checkpoint.<sup>38</sup>

A discussion of the mathematical and computational tools already available for the analysis of GRN models is given in the next section. Most of the models extracted from pathways databases are expected to be qualitative and incomplete in nature; hence, the discussion focuses on qualitative network structures and how these structures influence the capacity of the system to exhibit certain dynamical behavior.

---

## PATHWAY AND DYNAMIC ANALYSIS TOOLS FOR GRNs

Selection of the appropriate network analysis tool depends on the questions being asked and the scale or size of the network being considered. Questions of robustness of the entire system against perturbations require more consideration of global network properties and less of the attributes of individual processes or reactions. Questions focusing on particular phenomena, such as the switching behavior of a particular set of genes, may require more attention to the local network details involving these genes. How the global and local network properties interplay to produce local or system-level behavior is an important problem that requires multiscale analysis, both in time and space. In this section, a brief account is given on global network properties, how large networks can be analyzed or reduced by identifying recurring network motifs and extreme pathways, and how topology or network structure alone may already determine a network’s stability and its capacity to exhibit certain dynamical behavior. The goal of this section is not to provide a comprehensive review of the aforementioned topics (as they are broad and recent



reviews will be cited) but, instead, to point out particular directions of analysis of a GRN model once it has been constructed.

## Global Network Properties

Considering the very large number of interacting genes, proteins, and other molecules in a living cell, we would first like to ask questions about global features and properties of the entire network. How connected are the nodes in the network, and what is the mean path length between any two nodes? Are there clusters of interactions so that we may subdivide the network into modules? How robust is the system to perturbations, i.e., could redundant pathways take over if a pathway is cut off, so that the system's function is still intact? In general, the aim is to identify global network topological features that affect system function or behavior independent of the details of the individual nodes or interactions. Various attempts at searching for quantifiable structural features of metabolic networks, signaling networks, and GRNs had been carried out (see Barabasi and Oltvai<sup>40</sup> for a review). Some basic network descriptors are the *degree distribution*, the *path length distribution*, and the *clustering coefficient*.

The degree distribution  $P(k)$  is the probability that a node is linked to  $k$  other nodes. The  $P(k)$  of *random networks* exhibits a Poisson distribution, whereas that of *scale-free networks* approximates a power law of the form  $P(k) \sim k^{-\gamma}$ . An interesting suggestion is that most cellular networks approximate a scale-free topology<sup>41,42</sup> with an exponent  $\gamma$  between 2 and 3.<sup>43,44</sup> The interpretation of this suggestion is not clear.

The path length distribution of a network tells us how far nodes are from each other. Scale-free networks are “ultra-small” because they have an average path length of the order  $\log(\log N)$ , where  $N$  is the number of nodes. Random networks are “small” because their mean path length is of the order  $\log N$ .<sup>43,44</sup>

The clustering coefficient of a particular node  $A$  of a network is defined by  $C(A) = 2n(A)/(k(A)(k(A) - 1))$ , where  $k(A)$  is the number of neighbors of  $A$ , and  $n(A)$  is the number of connections between the neighbors of  $A$ .<sup>40</sup> The average clustering coefficient characterizes the tendency of a network to form node clusters and is a measure of the network's modularity. The average clustering coefficient of most real networks is larger than that of same-size random networks.<sup>45</sup> Cellular networks have a high average clustering coefficient, which indicates a highly modular structure.<sup>46,47</sup>

## Recurring Network Motifs

One approach that could simplify the analysis of a large network is to look for recurring *motifs*, which are subgraphs that are overrepresented in the network.<sup>48–51</sup> The motivation is that each motif can be analyzed separately

for its intrinsic properties, and the original network may be reduced to a set of motif interactions. Recent analysis<sup>48</sup> show that three-node feed-forward motifs are abundant in transcriptional regulatory networks and neural networks, whereas four-node feedback loops are characteristic of electric circuits, but not of biological networks. Remarkable evolutionary conservation of motifs<sup>52</sup> and convergent evolution toward the same motif types in transcriptional regulatory networks of diverse species<sup>53,54</sup> show that motifs are indeed significant biologically.

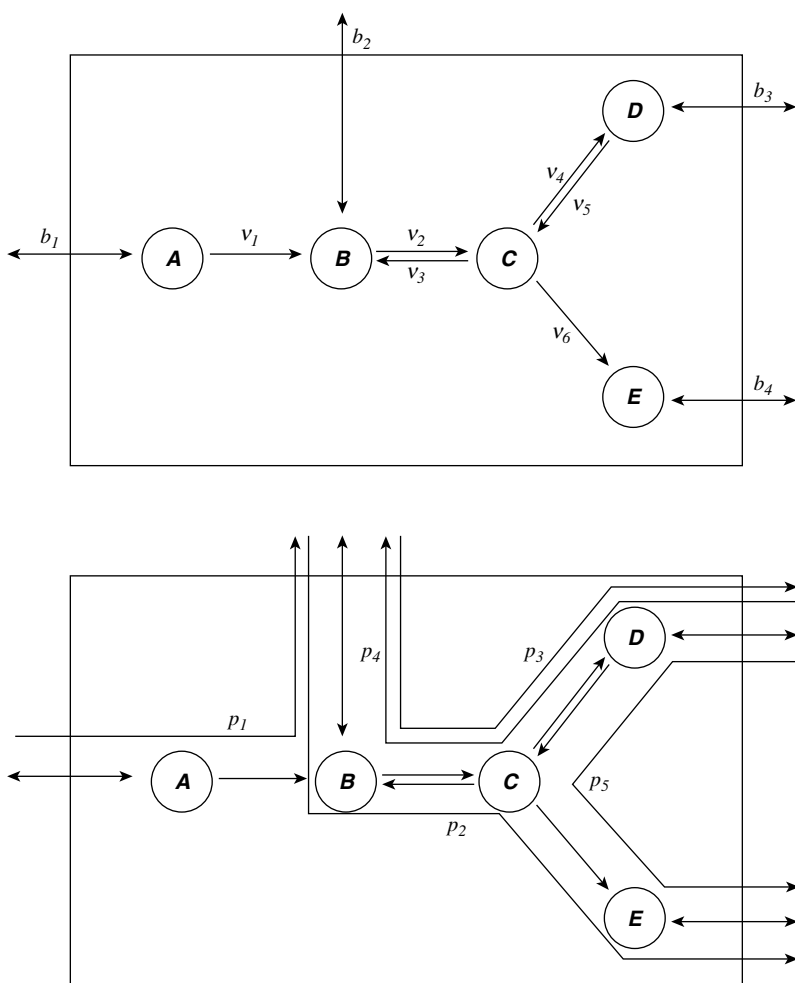
### Identifying Pathway Channels in Networks: Extreme Pathway Analysis

Another way of coping with large networks involves breaking down the network into channels through which distinct processes are carried out. Clarke<sup>55</sup> developed a formalism called “stoichiometric network analysis” and was the first to show that all steady-state fluxes are found in a convex set called the “current cone”; furthermore, he showed that each cone has a certain number of edge vectors that can be uniquely determined from the stoichiometric matrix. Clarke referred to the pathways corresponding to the edge vectors as “extreme currents”; alternatively, these currents are called *extreme pathways* in this chapter. Recent algorithms for computing extreme pathways can be found in Papin et al.<sup>56</sup> and Schilling et al.<sup>57</sup> The network shown in Figure 8 serves as an illustration of the basic ideas of extreme pathway analysis. In the network of Figure 8, six internal fluxes (labeled  $v_1$ – $v_6$ ) and four exchange fluxes (the reversible arrows  $b_1$ – $b_4$  showing exchange across the rectangular boundary) are present. Except for the two cycling pathways corresponding to the two reversible reactions ( $v_2$ – $v_3$  and  $v_4$ – $v_5$ ), the five extreme pathways are shown in the lower panel of Figure 8. Extreme pathway analysis has been extensively applied to metabolic networks.<sup>56,57</sup>

### Network Stability Analysis

One of the usual purposes of GRN modeling is to determine the origins of switching or threshold behaviors. These behaviors are often associated with the stability properties of the system against perturbations. Would initial perturbations of a species or a reaction in the network die out, or would it reverberate throughout the network? At least near steady states, the stability of the network is influenced by the network structure to a large extent.

More often than not, kinetic or other rate parameters are unknown in GRNs. Only the qualitative interactions between species are usually known, e.g., “X activates Y” or “V inhibits W.” As mentioned earlier, we can interpret the meaning of these qualitative interactions as follows:  $\partial(dY/dt)/\partial X > 0$  and  $\partial(dW/dt)/\partial V < 0$ , respectively. A “qualitative network” can be defined as a set of nodes (species) and a set of qualitative interactions (“activation” and



**Figure 8** A reaction network and its extreme pathways labeled  $p_1$ – $p_5$  (adapted from Schilling et al.<sup>57</sup>)

“inhibition”). Note that these qualitative interactions are none other than the elements of the Jacobian matrix of a linearized system of differential equations. The state  $\mathbf{x}$  of a linear dynamical system varies according to the differential equation

$$d\mathbf{x}/dt = \underline{\mathbf{A}}\mathbf{x} \quad [4]$$

where  $\underline{\mathbf{A}} = [a_{ij}]$  is an  $n \times n$  matrix and  $n$  is the number of species. For the case of a biochemical network,  $\mathbf{x}$  is the vector of perturbations from a steady state.

This dynamical system is *stable* if each solution  $\mathbf{x}(t)$  approaches zero for  $t$  large enough. A weaker condition is that the dynamical system is *semistable*, which means that, as  $t$  becomes larger and larger, the solution  $\mathbf{x}(t)$  could increase, but not at an exponential rate. It is well known that the dynamical system in Eq. [4] is stable if and only if all eigenvalues of  $\underline{\mathbf{A}}$  have negative real parts, and it is semistable if and only if all eigenvalues of  $\underline{\mathbf{A}}$  have nonpositive real part. The eigenvalues  $\lambda$  of the matrix  $\underline{\mathbf{A}}$  is given by the roots of the characteristic polynomial:

$$\det(\lambda \underline{\mathbf{I}} - \underline{\mathbf{A}}) = \lambda^n + \alpha_1 \lambda^{n-1} + \alpha_2 \lambda^{n-2} + \cdots + \alpha_{n-1} \lambda + \alpha_n = 0 \quad [5]$$

The coefficients  $\alpha_i$  are functions of the elements of  $\underline{\mathbf{A}}$ , and more importantly, the  $\alpha_i$ 's are functions of cycles in the qualitative network graph.<sup>58</sup> An example of cycles would be the three-cycle ( $a_{12}a_{23}a_{31}$ ) and the one-cycle ( $a_{ii}$ ). The eigenvalues, and therefore the linear stability of a network, are determined only by cycles in the qualitative network graph.

Suppose that only the *sign pattern* of the matrix  $\underline{\mathbf{A}}$  is known; i.e., the magnitudes of  $a_{ij}$  matrix entries are unknown, but their algebraic signs are known. If all matrices that have the same sign pattern as  $\underline{\mathbf{A}}$  are stable, then  $\underline{\mathbf{A}}$  is referred to as *sign-stable*. If all matrices that have the same sign pattern as  $\underline{\mathbf{A}}$  are semistable, then  $\underline{\mathbf{A}}$  is *sign-semistable*.

The notion of sign-semistability has a simple characterization for signs of the entries of the matrix  $\underline{\mathbf{A}}$  (the notion of sign-stability can also be characterized for signs of the entries of the matrix  $\underline{\mathbf{A}}$ , but in a more complicated way<sup>59</sup>). A useful theorem<sup>60</sup> states that  $\underline{\mathbf{A}}$  is sign-semistable if and only if three conditions are met: (1) there is no excitatory one-cycle, (2) any two-cycle must be negative (i.e., one edge must be inhibitory and the other excitatory), and (3) no cycles of length three exist. Note that because lack of sign-semistability implies lack of sign-stability, the theorem<sup>60</sup> on sign-semistability also gives a set of necessary conditions for sign-stability.

## Predicting Dynamics and Bistability from Network Structure Alone

We should identify or classify classes of network structures that, from their structures alone, it is possible to tell whether they have the capacity to exhibit certain behavior. Given a biochemical reaction network, we can ask the following question: Under which circumstances would this network exhibit phenomena like periodic oscillations and/or bistability? For example, we would want to know the answer to this question when modeling the cell division cycle and circadian rhythm where periodic oscillations are required. For mass-action kinetics models, an extensive theoretical work already exists that answers this type of question for large classes of reaction networks. One such set of results is the *deficiency theory*.<sup>61–64</sup> The deficiency  $\delta$  of a reaction

network is a function of the number of objects and linkages in the network and can be computed easily even if the rate expressions and kinetic parameters are unknown. For reaction networks with  $\delta = 0$ , they do not have the capacity to exhibit cyclic variation or bistability.<sup>63</sup> Feinberg also showed that some networks with deficiency  $\delta > 0$  do not have the capacity for bistability, if they have some additional properties.<sup>64</sup> These methods are implemented in the software package called *Chemical Reaction Network Toolbox*.<sup>64,65</sup> Recently, other methods of deciding on the capacity for bistability of biochemical reaction networks were developed.<sup>66,67</sup> The *SR graph method* of Craciun and Feinberg<sup>67</sup> allows us to draw conclusions on the capacity of a network to exhibit bistability based on the properties of cycles in the graph.

---

## CONCLUDING REMARKS

With a well-defined pathway ontology, we could envisage a computer program that automates the analysis of complex gene regulatory networks and the extraction or building of GRN models; these models can then be analyzed by computer simulation and other mathematical methods. An investigator would most likely start with a short list of genes or even a short list of specific cellular processes (from which a gene list could be derived with existing gene annotations such as GeneOntology). By scouring databases, the computer program would then try to establish pathways among the initial set of genes; this step will increase the number of genes in the network and include proteins and other molecules regulating the pathways. At this point, the GRN is a static graph, perhaps a qualitative network containing some information about how the nodes affect each other. The computer program can now apply network analysis tools to study the topology of the GRN and to identify stabilizing or destabilizing cycles, extreme pathways, or even try to reduce the size of the network without removing the capacity for certain behaviors of interest. Databases (including the published literature) containing experimental information will have to be consulted to validate the significance or strength of contribution of the pathways and cycles present in the reduced model. The rate expressions and associated kinetic parameters of the individual steps in the model are then supplied to a solver of the dynamical equations to simulate the temporal evolution of the model system. Predictions of the model will have to be compared with experimental data, and the process of model refinement and experimental validation could be iterated.

As the work of Ozbudak et al.<sup>3</sup> and Vilar, Guet and Leibler<sup>4</sup> on the lac operon demonstrated, abstract kinetic models with a few variables are sometimes sufficient to capture the essential behavior of the system, e.g., the bistable switch in the lac operon. Some arbitrariness may seem to exist in how these simple lac operon models<sup>3,4</sup> were generated, because they seem to look very different and the dynamical variables are not the same. However,

both models preserve the common property of having a positive feedback loop. The presence of such a loop has long been known, in dynamical systems theory, to give a system the ability to generate bistability given the right parameters. As the work of Ozbudak et al.<sup>3</sup> showed, a low-dimensional abstract model can indeed be predictive. In the future, the process of extracting an abstract model from a complex GRN may well be carried out systematically. The key will be the application of the mathematical fields of nonlinear dynamics and reaction network analysis. As mentioned, possible behavior of networks may already be predicted from their qualitative network structures regardless of the values of rate parameters. The development of a pathway ontology that can interface with network structure analysis tools will be crucial for the integration and application of the huge amounts of data stored in databases.

Significant advances toward understanding gene networks are coming from recent work on synthetic gene networks (see Weiss et al.<sup>68</sup> for a recent review); the goal here is the construction and engineering control of genetic circuits built from well-understood building blocks of small gene modules. What is being learned from these man-made gene networks will be very useful to scientists in future analysis of the very complex GRN repertoire of a living cell.

---

## REFERENCES

1. P. Wong, S. Gladney, and J. D. Keasling, *Biotechnol. Prog.*, **13**, 132 (1997). Mathematical Model of the Lac Operon: Inducer Exclusion, Catabolite Repression, and Diauxic Growth on Glucose and Lactose.
2. N. Yildirim and M. C. Mackey, *Biophys. J.*, **84**, 2841 (2003). Feedback Regulation in the Lactose Operon: A Mathematical Modeling Study and Comparison with Experimental Data.
3. E. M. Ozbudak, M. Thattai, H. N. Lim, B. I. Shraiman, and A. van Oudenaarden, *Nature*, **427**, 737 (2004). Multistability in the Lactose Utilization Network of *Escherichia coli*.
4. J. M. G. Vilar, C. C. Guet, and S. Leibler, *J. Cell Biol.*, **161**, 471 (2003). Modeling Network Dynamics: The Lac Operon, A Case Study.
5. F. Jacob, D. Perrin, C. Sanchez, and J. Monod, *Compt. Rendu. Acad. Sci.*, **245**, 1727 (1960). L'operon: Groupe de Genes a l'Expression Coordonne par un Operateur.
6. J. J. Tyson and H. G. Othmer, in *Progress in Biophysics*, R. Rosen, Ed., Academic Press, New York, **5**, 1 (1978). The Dynamics of Feedback Control Circuits in Biochemical Pathways.
7. D. M. Wolf and F. H. Eeckman, *J. Theor. Biol.*, **195**, 167 (1998). On the Relationship between Genomic Regulatory Element Organization and Gene Regulatory Dynamics.
8. J. J. Tyson, K. C. Chen, and B. Novak, *Curr. Opin. Cell Biol.*, **15**, 221 (2003). Sniffers, Buzzers, Toggles and Blinkers: Dynamics of Regulatory and Signaling Pathways in the Cell.
9. H. de Jong, *J. Comput. Biol.*, **9**, 67 (2002). Modeling and Simulation of Genetic Regulatory Systems: A Literature Review.
10. N. Friedman, *Science*, **303**, 799 (2004). Inferring Cellular Networks Using Probabilistic Graphical Models.
11. E. P. van Someren, L. F. A. Wessels, E. Backer, and M. J. T. Reinders, *Pharmacogenomics*, **3**, 507 (2002). Genetic Network Modeling.

12. J. Stark, D. Brewer, M. Barenco, D. Tomescu, R. Callard, and M. Hubank, *Biochem. Soc. Trans.*, **31**, 1519 (2003). Reconstructing Gene Networks: What are the Limits?
13. J. Stark, R. Callard, and M. Hubank, *Trends Biotech.*, **21**, 290 (2003). From the Top Down: Towards a Predictive Biology of Signaling Networks.
14. P. D'haeseleer, S. Liang, and R. Somogyi, *Bioinformatics*, **16**, 707 (2000). Genetic Network Inference: From Co-expression Clustering to Reverse Engineering.
15. M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein, *Proc. Natl. Acad. Sci. USA*, **95**, 14863 (1998). Cluster Analysis and Display of Genome-Wide Expression Patterns.
16. J. Wang, O. Myklebost, and E. Hovig, *Bioinformatics*, **19**, 2210 (2003). MGraph: Graphical Models for Microarray Data Analysis.
17. D. Husmeier, *Biochem. Soc. Trans.*, **31**, 1516 (2003). Reverse Engineering of Genetic Networks with Bayesian Networks.
18. D. Husmeier, *Bioinformatics*, **19**, 2271 (2003). Sensitivity and Specificity of Inferring Genetic Regulatory Interactions from Microarray Experiments with Dynamic Bayesian Networks.
19. N. Friedman, M. Linial, I. Nachman, and D. Pe'er, *J. Comput. Biol.*, **7**, 601 (2000). Using Bayesian Networks to Analyze Expression Data.
20. D. Pe'er, A. Regev, G. Elidan, and N. Friedman, *Bioinformatics*, **17**, S215 (2001). Inferring Subnetworks from Perturbed Expression Profiles.
21. A. J. Hartemink, D. K. Gifford, T. S. Jaakkola, and R. A. Young, *Pac. Symp. Biocomput.*, **6**, 422 (2001). Using Graphical Models and Genomic Expression Data to Statistically Validate Models of Genetic Regulatory Networks.
22. D. E. Zak, F. J. Doyle, and J. S. Schwaber, *Proceedings of the Third International Conference on Systems Biology*, Karolinska Institute, Sweden, 2002, pp. 236–237. Local Identifiability: When Can Genetic Networks be Identified from Microarray Data?
23. A. de la Fuente, P. Brazhnik, and P. Mendes, *Trends Genet.*, **18**, 395 (2002). Linking the Genes: Inferring Quantitative Gene Networks from Microarray Data.
24. B. N. Kholodenko, A. Kiyatkin, F. J. Bruggeman, E. Sontag, H. V. Westerhoff, and J. B. Hoek, *Proc. Natl. Acad. Sci. USA*, **99**, 12841 (2002). Untangling Wires: A Strategy to Trace Functional Interactions in Signaling and Gene Networks.
25. T. S. Gardner, D. di Bernardo, D. Lorenz, and J. J. Collins, *Science*, **301**, 102 (2003). Inferring Genetic Networks and Identifying Compound Mode of Action via Expression Profiling.
26. M. Y. Galperin, *Nucleic Acids Res.*, **32**, D3 (2004). The Molecular Biology Database Collection: 2004 Update.
27. C. Lemer, E. Antezana, F. Couche, F. Fays, X. Santolaria, R. Janky, Y. Deville, J. Richelle, and S. J. Wodak, *Nucleic Acids Res.*, **32**, D443 (2004). The aMAZE LightBench: A Web Interface to a Relational Database of Cellular Processes.
28. T. R. Gruber, *Knowledge Acquisition*, **5**, 199 (1993). A Translation Approach to Portable Ontologies.
29. P. Mendes, *Trends Biochem. Sci.*, **22**, 361 (1997). Biochemistry by Numbers: Simulation of Biochemical Pathways with Gepasi 3.
30. B. M. Slepchenko, J. C. Schaff, I. Macara, and L. M. Loew, *Trends Cell Biol.*, **13**, 570 (2003). Quantitative Cell Biology with the Virtual Cell.
31. M. Tomita, K. Hashimoto, K. Takahashi, T. Shimizu, Y. Matsuzaki, F. Miyoshi, K. Saito, S. Tanida, K. Yugi, J. C. Venter, and C. A. Hutchison, *Bioinformatics*, **15**, 72 (1999). E-CELL: Software Environment for Whole-Cell Simulation.
32. K. Takahashi, N. Ishikawa, Y. Sadamoto, H. Sasamoto, S. Ohta, A. Shiozawa, F. Miyoshi, Y. Naito, Y. Nakayama, and M. Tomita, *Bioinformatics*, **19**, 1727 (2003). E-Cell 2: Multi-Platform E-Cell Simulation System.
33. P. Dhar, T. C. Meng, S. Somani, L. Ye, A. Sairam, M. Chitre, Z. Hao, and K. Sakthar, *Bioinformatics*, **20**, 1319 (2004). Cellware: A Multi-Algorithmic Software for Computational Systems Biology.

34. E. Demir, O. Babur, U. Dogrusoz, A. Gursoy, G. Nisanci, R. Cetin-Atalay, and M. Ozturk, *Bioinformatics*, **18**, 996 (2002). PATIKA: An Integrated Visual Environment for Collaborative Construction and Analysis of Cellular Pathways.
35. E. Demir, O. Babur, U. Dogrusoz, A. Gursoy, A. Ayaz, G. Gulesir, G. Nisanci, and R. Cetin-Atalay, *Bioinformatics*, **20**, 349 (2004). An Ontology for Collaborative Construction and Analysis of Cellular Pathways.
36. L. Krishnamurthy, J. Nadeau, G. Ozsoyoglu, M. Ozsoyoglu, G. Schaeffer, M. Tasan, and W. Xu, *Bioinformatics*, **22**, 930 (2003). Pathways Database System: An Integrated System for Biological Pathways.
37. K. Fukuda and T. Takagi, *Bioinformatics*, **17**, 829 (2001). Knowledge Representation of Signal Transduction Pathways.
38. B. D. Aguda and Y. Tang, *Cell Proliferation*, **32**, 321 (1999). The Kinetic Origins of the Restriction Point in the Mammalian Cell Cycle.
39. B. D. Aguda, *Oncogene*, **18**, 2846 (1999). Instabilities in Phosphorylation Dephosphorylation Cascades and Cell Cycle Checkpoints.
40. A.-L. Barabasi and Z. N. Oltvai, *Nat. Rev. Genet.*, **5**, 101 (2004). Network Biology: Understanding the Cell's Functional Organization.
41. H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A.-L. Barabasi, *Nature*, **407**, 651 (2000). The Large-Scale Organization of Metabolic Networks.
42. A. Wagner and D. A. Fell, *Proc. R. Soc. Lond. B*, **268**, 1803 (2001). The Small World Inside Metabolic Networks.
43. F. Chung and L. Lu, *Proc. Natl. Acad. Sci. USA*, **99**, 15879 (2002). The Average Distances in Random Graphs with Given Expected Degrees.
44. R. Cohen and S. Havlin, *Phys. Rev. Lett.*, **90**, 058701 (2003). Scale-Free Networks are Ultra Small.
45. D. J. Watts and S. H. Strogatz, *Nature*, **393**, 440 (1998). Collective Dynamics of 'Small-World' Networks.
46. A. Wagner, *Mol. Biol. Evol.*, **18**, 1283 (2001). The Yeast Protein Interaction Network Evolves Rapidly and Contains Few Redundant Duplicate Genes.
47. S. Wuchty, *Mol. Biol. Evol.*, **18**, 1694 (2001). Scale-Free Behaviour in Protein Domain Networks.
48. R. Milo, S. S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon, *Science*, **298**, 824 (2002). Network Motifs: Simple Building Blocks of Complex Networks.
49. S. Itzkovitz, R. Milo, N. Kashtan, G. Ziv, and U. Alon, *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.*, **68**, 026127 (2003). Subgraphs in Random Networks.
50. E. Yeger-Lotem, S. Sattath, N. Kashtan, S. Itzkovitz, R. Milo, R. Y. Pinter, and U. Alon, *Proc. Natl. Acad. Sci. USA*, **101**, 5934 (2004). Network Motifs in Integrated Cellular Networks of Transcription-Regulation and Protein-Protein Interaction.
51. J. Ihmels, G. Friedlander, S. Bergmann, O. Sarig, Y. Ziv, and N. Barkai, *Nature Genet.*, **31**, 370 (2002). Revealing Modular Organization in the Yeast Transcriptional Network.
52. S. Wuchty, Z. N. Oltvai, and A.-L. Barabasi, *Nature Genet.*, **35**, 176 (2003). Evolutionary Conservation of Motif Constituents within the Yeast Protein Interaction Network.
53. G. C. Conant and A. Wagner, *Nature Genet.*, **34**, 264 (2003). A Convergent Evolution of Gene Circuits.
54. V. F. Hinman, A. T. Nguyen, R. A. Cameron, and E. H. Davidson, *Proc. Natl. Acad. Sci. USA*, **100**, 13356 (2003). Developmental Gene Regulatory Network Architecture Across 500 Million Years of Echinoderm Evolution.
55. B. L. Clarke, *Adv. Chem. Phys.*, **43**, 1 (1980). Stability of Complex Reaction Networks.
56. J. A. Papin, N. D. Price, and B. Ø. Palsson, *Genome Res.*, **12**, 1889 (2000). Extreme Pathway Lengths and Reaction Participation in Genome-Scale Metabolic Networks.



- 
57. C. Schilling, D. Letscher, and B. Ø. Palsson, *J. Theor. Biol.*, **203**, 229 (2000). Theory for the Systemic Definition of Metabolic Pathways and their use in Interpreting Metabolic Function from a Pathway-Oriented Perspective.
  58. C. J. Puccia and R. Levins, *Qualitative Modeling of Complex Systems: An Introduction to Loop Analysis and Time Averaging*, Harvard University Press, Cambridge, Massachusetts, 1985.
  59. V. Klee, in *Applications of Combinatorics and Graph Theory to the Biological and Social Sciences*, F. Roberts, Ed., IMA Volumes in Mathematics and Its Applications, Vol. 17, Springer, New York, 1989, pp. 203-219. Sign-Patterns and Stability.
  60. J. Quirk and R. Ruppert, *Rev. Economic Studies*, **32**, 311 (1965). Qualitative Economics and the Stability of Equilibrium.
  61. M. Feinberg, *Arch. Rational Mech. Anal.*, **49**, 187 (1972). Complex Balancing in General Kinetic Systems.
  62. F. Horn and R. Jackson, *Arch. Rational Mech. Anal.*, **47**, 81 (1972). General Mass Action Kinetics.
  63. M. Feinberg, *Lectures on Chemical Reaction Networks*, written version of lectures given at the Mathematical Research Center, University of Wisconsin, Madison, Wisconsin, 1979.
  64. M. Feinberg, *Arch. Rational Mech. Anal.*, **132**, 311 (1995). The Existence and Uniqueness of Steady States for a Class of Chemical Reaction Networks.
  65. M. Feinberg. *Chemical Reaction Network Toolbox*, Version 1.02 (1995), and Version 1.1 (1999, with P. Ellison.) Available: <http://www.che.eng.ohio-state.edu/~feinberg/crnt>.
  66. G. Craciun and M. Feinberg, *SIAM J. Appl. Math.*, (2005). In press. Multiple Equilibria in Complex Chemical Reaction Networks: The Injectivity Property.
  67. G. Craciun and M. Feinberg, *Mathematical Biosciences Institute Technical Report*, 22 (2004). Multiple Equilibria in Complex Chemical Reaction Networks: The SR Graph.
  68. R. Weiss, S. Basu, S. Hooshangi, A. Kalmbach, D. Karig, R. Mehreja, and I. Netravali, *Natural Computing*, **2**, 47 (2003). Genetic Circuit Building Blocks for Cellular Computation, Communications, and Signal Processing.



---

# Author Index

---

- Abe, H., 281  
Abraham, M. M., 124  
Achenie, L. E. K., 286  
Affi, A., 344, 345  
Agrafiotis, D. K., 285, 348  
Aguda, B. D., 410  
Aguero, R., 278  
Ahlrichs, R., 113  
Ahrens, H. K., 119  
Ahuja, R., 119  
Ajay, 285  
Akamine, P., 376  
Akutsu, T., 277  
Albert, H., 344  
Albert, R., 410  
Albrecht, S., 116  
Alexov, E., 375  
Alfredsson, M., 120  
Allan, D. C., 119, 120  
Allan, N. L., 118, 123, 124, 202  
Allen, F. H., 203  
Allewell, N. M., 377  
Almo, S. C., 379  
Alon, U., 410  
Alouani, M., 116  
Al-Sharif, A. I., 116  
Alvarez, J., 278  
Amat, L.I., 197, 199, 201, 203, 204, 205, 206, 207  
Ancian, R., 280  
Andersen, K. V., 377  
Andersen, O. K., 114  
Anderson, J. B., 115  
Anderson, R. E., 344  
André, J. M., 114  
Andreoni, W., 125  
Andrés, J., 118, 121, 122  
Andrew, A. W., 343  
Ángyán, J. G., 118  
Antezana, E., 409  
Antosiewicz, J., 377, 379  
Aprá, E., 117, 118  
Aradi, B., 123  
Arbuzova, A., 377  
Arias, G., 278  
Armstrong, D. E., 343  
Arnau, J., 199  
Arnaud, B., 116  
Artacho, E., 125  
Ashbaugh, H. S., 373  
Ashcroft, N. W., 117  
Auton, T., 285  
Axelsson, O., 375  
Ayala, Q. Y., 115  
Ayaz, A., 410  
Ayers, P. W., 200  
Azavant, P., 124  
  
Babur, O., 410  
Bachelet, G. B., 113  
Bachrach, S. M., 198  
Backer, E., 408  
Bader, R. F. W., 204  
Badertscher, M., 283  
Baerends, E. J., 125, 200  
Bagheri, B., 375, 379  
Bagus, P. S., 123  
Bajaj, C. L., 374  
Bajorath, J., 198, 348  
Baker, N. A., 373, 374, 375, 376, 379

*Reviews in Computational Chemistry, Volume 21*  
edited by Kenny B. Lipkowitz, Raima Larter, and Thomas R. Cundari  
Copyright © 2005 Wiley-VCH, John Wiley & Sons, Inc.

- Balaban, A. T., 198, 276, 277, 278, 279, 282  
Balasubramanian, K., 278  
Baldwin, K. F., 345  
Balkenhohl, F., 284  
Ballinger, R. A., 113  
Bandura, A. V., 120  
Bangov, I. P., 280  
Bank, R. E., 376  
Banner, D., 285  
Barabasi, A.-L., 410  
Baraldi, I., 276  
Baranek, P., 124  
Barberis, F., 283  
Barden, C. J., 200  
Barenco, M., 409  
Barkai, N., 410  
Barker, V. A., 375  
Barnard, J. M., 203, 285  
Barone, R., 283  
Baroni, S., 119, 125  
Bartlett, R. J., 115  
Bartolotti, L. J., 204  
Basak, S. C., 344  
Bashford, D., 373, 377, 379  
Baskin, I. I., 281  
Bassani, F., 113  
Basu, S., 411  
Bauman, N., 281  
Baumann, K., 344, 348  
Bayly, C. I., 372  
Beavers, M. P., 286  
Becke, A. D., 116  
Becker, O. M., 372, 373  
Becker, U., 121, 122  
Beglov, D., 374  
Bell, P. M., 119  
Bellaiche, L., 114  
Belsley, D. A., 344  
Beltrán, A., 118, 121, 122  
Bendel, R. B., 345  
Bender, C. F., 113  
Benecke, C., 280  
Benfenati, E., 347  
Bennett, B. I., 114  
Bennewitz, R., 120  
Benson, M. T., 201  
Ben-Tal, N., 376, 377  
Bergmann, S., 410  
Berman, H. M., 379  
Bernardi, F., 117  
Bernasconi, L., 116  
Bertz, S. H., 197  
Besalú, E., 197, 199, 200, 203, 204, 205, 206, 207  
Beuken, J.-M., 119  
Bhat, T. N., 379  
Bhattacharyya, S. M., 373  
Bianchi, R., 118  
Bickelhaupt, F. M., 125  
Biczó, G., 114  
Bielawski, J., 285  
Bihlmayer, G., 125  
Black, W. C., 344  
Blaha, P., 125  
Blair, C., 276  
Blaney, J. M., 284  
Blinov, K. A., 283, 284  
Bloch, F., 116  
Blöchl, P. E., 116, 125  
Block, R., 119  
Blomberg, N., 378  
Blügel, S., 125  
Blum, H., 124  
Bockris, J. O., 373  
Bockstedte, M., 125  
Boggia, R., 346  
Bohanec, S., 280, 283  
Bohm, H.-J., 285, 343  
Bolton-Smith, C., 345  
Bonanno, J. B., 379  
Boon, G., 202  
Bordiga, S., 124  
Born, M., 119, 199, 373  
Börnstein, R., 116  
Borstel, G., 121  
Botstein, D., 409  
Botti, S. A., 378  
Bourne, P. E., 379  
Bourne, Y., 379  
Boyd, D. B., 117, 198, 201, 203, 276, 281, 284, 346  
Boys, S. F., 117  
Bradshaw, J., 285  
Braess, D., 375  
Bransden, B. H., 202  
Brause, M., 120  
Brazhnik, P., 409  
Bredow, T., 122, 199  
Breimann, L., 345  
Brenner, S. C., 375  
Brewer, D., 409  
Briggs, J. M., 374, 377, 379  
Briggs, W. L., 375  
Brinkmann, G., 279

- 
- Brodholt, J. P., 115  
 Brooks, B. R., 379  
 Brooks, III, C. L., 373, 374, 379  
 Brooks, R. J., 345  
 Brown, D. C., 286  
 Brown, P. O., 409  
 Brown, R., 285  
 Brown, R. D., 197  
 Bruccoleri, R. E., 379  
 Bruggeman, F. J., 409  
 Bruggemann, R., 286  
 Bruneau, P. P., 343, 348  
 Brunger, A. T., 374  
 Brunnvoll, J., 276, 277, 278, 279, 282, 283  
 Buchanan, B. G., 277, 278  
 Bultinck, P., 197, 198, 199, 200, 201, 202, 203, 204  
 Bunce, J. D., 206  
 Burden, F. R., 347  
 Burggraf, L. W., 123  
 Burkhard, L. P., 343  
 Burley, S. K., 379  
 Burt, C., 379  
 Buser, C. A., 377  
 Bush, I. J., 115, 125  
 Bussche-Hunnefeld, C. v. d., 284  
 Bussemaker, F. C., 279  
 Busso, M., 115, 122, 123  
 Bussolin, G., 121  
 Bytautas, L., 277  
 Byun, K. S., 123
- Calabuig, B., 197, 199, 200, 204  
 Calatayud, M., 121, 122  
 Caldwell, J. W., 376  
 Callard, R., 409  
 Camblor, M. A., 118  
 Cameron, R. A., 410  
 Camus, S., 118  
 Capecchi, G., 115  
 Capel, M., 379  
 Caporossi, G., 279  
 Car, R., 115  
 Carabedian, M., 280, 283  
 Caracas, R., 119  
 Carbó, R., 197, 199, 200, 201, 202, 203, 204, 205, 206, 207  
 Carbó-Dorca, R., 197, 199, 200, 201, 202, 203, 204, 205, 206, 207  
 Cárdenas-Jirón, G. L., 118  
 Carhart, R. E., 280, 281, 285
- Carlsen, L., 286  
 Cartier, A., 344  
 Caruthers, J. M., 282  
 Casassa, S., 115, 118, 120, 121, 122, 123  
 Case, D. A., 373, 374, 376  
 Casida, M. E., 116  
 Catlow, C. R. A., 114, 118, 120  
 Catti, M., 118, 119  
 Causà, M., 114, 117, 118, 120, 121, 122, 123  
 Cayley, A., 275, 276  
 Cedeno, W., 348  
 Cetin-Atalay, R., 410  
 Chadwick, A. V., 122  
 Chan, K., 282  
 Chance, M. R., 379  
 Chandana Epa, V., 378  
 Chandler, D., 374  
 Chang, C., 113  
 Chanon, M., 283  
 Charifson, P. S., 284  
 Chartier, A., 119  
 Charton, M., 343  
 Chatfield, C., 344  
 Chatterjee, N., 118  
 Chaudry, U. A., 204  
 Cheatham, III, T. E., 376  
 Chen, K. C., 408  
 Chen, L., 198  
 Chen, X., 286  
 Chen, Y., 124  
 Cheng, A., 376  
 Chiabrera, A., 375  
 Chitre, M., 409  
 Chklovskii, D., 410  
 Cho, S. J., 347  
 Cho, W., 378  
 Chong, L., 376  
 Chong, L. T., 378  
 Chothia, C., 374, 378  
 Christie, B., 285  
 Christie, B. D., 280, 284  
 Chung, F., 410  
 Churchwell, C. J., 281, 286  
 Cieplak, P., 372, 376  
 Cioslowski, J., 203, 204  
 Civalieri, B., 115, 117, 118, 119, 120, 122  
 Clark, D., 285  
 Clark, D. E., 284  
 Clark, R., 284  
 Clark, S. J., 124  
 Clark, V., 344  
 Clarke, B. L., 410

- Clementi, E., 113  
Clifton, J. G., 376  
Cobeljic, S., 279  
Cohen, R., 410  
Colle, R., 112  
Collins, A. J., 344  
Collins, J. J., 409  
Colman, P. M., 378  
Combariza, J. E., 117  
Comley, J. C. W., 346  
Conant, G. C., 410  
Connolly, J. W., 200  
Connolly, M. L., 374  
Consonni, V., 343  
Constans, P., 201, 203  
Contreras, M. L., 278  
Cooper, D. L., 202  
Copas, J. B., 346  
Corà, F., 120, 124  
Cornell, W. D., 372  
Cortis, C. M., 375  
Couche, F., 409  
Cracium, G., 411  
Cramer, III, R. D., 206, 284  
Craven, B. N., 125  
Crippen, G. M., 204  
Crowley, M., 376  
Cruz Hernandez, N., 122  
Culbertson, J. C., 348  
Cundari, T. R., 201, 373  
Curioni, A., 125  
Curtiss, L. A., 117, 118, 124  
Custodio, R., 121  
Cvetkovic, D. M., 279  
Cyvin, B. N., 276, 277, 279, 282, 283  
Cyvin, S. J., 276, 277, 278, 279, 282, 283
- D'Arco, P., 115, 119  
D'Ercole, A., 120  
D'haeseleer, P., 409  
Dagane, L., 283  
Dailey, R. S., 281  
Dal Corso, A., 115, 119, 125  
Dallal, G. E., 346  
Damin, A., 122, 124  
Dannenberg, J. J., 117, 118  
Dapprich, S., 123  
Darden, T., 376  
David, L., 373, 377  
Davidson, E. H., 410  
Davidson, E. R., 117, 201  
Davidson, S., 283
- Davies, G. J., 377  
Davis, M. E., 118, 374, 375, 379  
de Gironcoli, S., 119, 125  
de Graaf, C., 119  
de Jong, H., 408  
De Jong, L. J., 119  
de la Fuente, A., 409  
De Leeuw, S. W., 114  
De Proft, F., 198, 202  
De Rienzo, F., 378  
de Vries, J., 279  
De Winter, H., 197, 198, 199, 200, 202  
Deák, P., 123  
Dean, P. M., 197, 199  
Dearden, J. C., 343, 347, 348  
DeBolt, S., 376  
Defranceschi, M., 206  
Del Sole, R., 116  
Delbuono, G. S., 373  
Delley, B., 125  
Del Re, G., 114  
Demchuk, E., 378  
Deming, S. N., 203  
Demir, E., 410  
Demmell, J., 125  
Dempster, S. E., 378  
Demuth, T., 118  
Derringer, G. C., 281  
Detraux, F., 119  
Deville, Y., 409  
Devillers, J., 344, 346  
Dhar, P., 409  
di Bernardo, D., 409  
di Bona, A., 122  
Diallo, M. S., 284  
Dias, J. R., 277, 282, 283  
Diercksen, G. H. F., 113  
Diller, T., 376  
Dillon, W. R., 344  
Dinte, B. P., 117  
Dirac, P. A. M., 116, 206  
Diraviyam, K., 378  
Ditchfield, R., 113  
Dixon, J. D., 281  
Dixon, W. J., 345  
Djerassi, C., 277, 280, 281  
Dobson, J. F., 117  
Dogrusoz, U., 410  
Dolg, M., 115  
Dolhaine, H., 283  
Doll, K., 115, 120, 121, 124  
Dominy, B. N., 285, 373

- Dominy, B. W., 198  
 Donini, O., 376  
 Dorfman, S., 123  
 Doroslovacki, R., 278  
 Dovesi, R., 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124  
 Dow, E. R., 346  
 Downs, G. M., 197, 203, 285  
 Doyle, F. J., 409  
 Draper, N. R., 345  
 Dress, A. W., 279  
 Duan, Y., 376  
 Dubois, J. E., 280, 283  
 Dubrovinskaia, N. A., 119  
 Dubrovinsky, L. S., 119  
 Duffield, A. M., 277  
 Dunlap, B. I., 200  
 Duran, M., 201, 205, 206, 207  
 Dury, L., 202  
  
 Eeckman, F. H., 408  
 Efroymson, M. A., 345  
 Eglitis, R. I., 121  
 Ehrenreich, H., 112  
 Eichler, U., 123  
 Eisen, M. B., 409  
 Eisenberg, B., 373, 374  
 Elcock, A. H., 376, 377  
 Elidan, G., 409  
 Elkayam, T., 378  
 Ellinger, Y., 206  
 Ellis, D. E., 200  
 Ellman, J. A., 284, 286  
 Elyashberg, M. E., 281, 283, 284  
 Enting, I. G., 279  
 Erdemir, A., 118  
 Evans, D. A., 344  
 Evarestov, R. A., 120, 121, 124  
 Ewald, P. P., 114  
 Ewema, R. N., 114  
 Ewing, T., 284, 286  
 Ewing, T. J. A., 284  
 Eyring, H., 206  
  
 Fachinger, W., 280  
 Fan, B. T., 347  
 Faradzev, I. A., 277  
 Faulon, J.-L., 276, 280, 281, 284, 286  
 Fays, F., 409  
 Fedyaev, K. S., 286  
 Feenay, P. J., 198  
 Feeney, P., 285  
  
 Fei, Y., 118  
 Feig, M., 373  
 Feigenbaum, E. A., 277, 278  
 Feinberg, M., 411  
 Felder, C. E., 378  
 Fell, D. A., 410  
 Feller, D., 117  
 Felton, R. H., 200  
 Feng, J., 348  
 Feng, J. Q., 376  
 Feng, Z., 379  
 Ferguson, D., 376  
 Ferrari, A. M., 119, 120, 123  
 Feuston, B. P., 348  
 Figueirido, F., 373  
 Fine, R. F., 374  
 Firth, M., 285  
 Fisk, J. B., 281  
 Flannery, B. P., 203  
 Fleming, P., 285, 347  
 Flurey, B., 344  
 Fogolari, F., 374  
 Fonseca Guerra, C., 125  
 Fontana, P., 283  
 Ford, M. G., 344, 345, 346, 347, 348  
 Forés, M., 207  
 Först, C. J., 125  
 Foster, R., 344  
 Fowler, P. W., 275, 276, 279, 282  
 Fradera, X., 201, 203, 204, 205  
 Fraga, S., 200  
 Frank, I. E., 345  
 Frauenheim, T., 123  
 Frazer, J. W., 280, 281  
 Frechard, F., 121  
 Freeman, A. J., 114  
 Freeman, H. C., 378  
 Freund, H. J., 122  
 Freyria-Fava, C., 115, 118, 124  
 Friedlander, G., 410  
 Friedman, J. H., 345  
 Friedman, N., 408, 409  
 Friesner, R. A., 375  
 Fripertinger, H., 276  
 Frisch, M. J., 123  
 Fu, L., 121  
 Fuchs, M., 119  
 Fuji, Z., 282  
 Fujita, S., 276  
 Fujita, T., 204, 343  
 Fuks, D., 123  
 Fukuda, K., 410

- Fulde, P., 115  
Funatsu, K., 280, 347  
Furui, J., 378
- Gaasterland, D., 379  
Gabdouline, R. R., 378  
Gale, J. D., 120, 121, 125  
Gallegos, A., 199  
Gallop, M. A., 284  
Galperin, M. Y., 409  
Gambhir, A., 377  
Gans, D. J., 347  
Garcia, A., 125  
Garcia, A. E., 373  
Gardner, T. S., 409  
Garey, M. R., 276  
Garg, S., 286  
Garrone, E., 117, 122  
Gassman, P., 284  
Gasteiger, J., 206, 286  
Geerlings, P., 198, 202  
Gennard, S., 120  
Gerstmann, U., 123  
Gervais, F., 120  
Getzoff, E. D., 378  
Ghose, A. K., 204  
Ghosez, P., 115, 119  
Giannozzi, P., 119, 125  
Gifford, D. K., 409  
Gilat, G., 117  
Gill, P. M. W., 202  
Gillan, M. J., 123  
Gillespe, D., 373  
Gillet, V. J., 285, 286, 347  
Gilliland, G., 379  
Gilson, M. K., 373, 374, 375, 377, 379  
Giordano, L., 122  
Giovanardi, C., 122  
Giraldo, J., 347  
Gironés, X., 197, 199, 200, 201, 203, 205, 206  
Given, J. A., 374  
Gladney, S., 408  
Glen, R. C., 343  
Globus, A., 282  
Glover, F., 346  
Gluzman, I., 286  
Goddard, III, W. A., 284  
Godden, J. W., 348  
Gohlke, H., 376  
Golbraikh, A., 346, 348  
Goldberg, D., 286  
Goldberg, L. A., 277, 286  
Goldstein, M., 344  
Gombar, V., 348  
Gomes, J. R. B., 122  
González, R. C., 203  
Gonze, X., 115, 116, 119, 120  
Good, A. C., 285, 379  
Gordon, A., 123  
Gordon, E. M., 284  
Gordon, M. S., 123  
Görling, A., 116  
Gould, T., 117  
Gouverneur, L., 114  
Gramaccioli, C. M., 115  
Grant, J. A., 374  
Graovac, A., 279  
Gray, N. A. B., 278, 280  
Green, D., 285  
Green, D. V. S., 347  
Grier, D., 285  
Gross, E. K. U., 115, 116  
Grossfield, A., 372  
Gruber, T. R., 409  
Gruen, D. M., 124  
Grund, R., 280, 283  
Grüner, T., 279  
Guet, C. C., 408  
Gugisch, R., 285  
Gulesir, G., 410  
Gunner, M. R., 377  
Gunster, J., 120  
Gursoy, A., 410  
Gute, B. D., 344  
Gutman, I., 198, 277, 282  
Guttmann, A. J., 279
- Habas, M. P., 119  
Habashita, H., 286  
Hackbusch, W., 375  
Hafner, J., 118, 125  
Hahn, T., 116  
Hair, J. F., 344  
Hall, G. G., 200  
Hall, L. H., 198, 280, 281, 286, 343  
Hamann, D. R., 113  
Handy, N. C., 199  
Hangyas-Mihalyne, G., 377  
Hann, M., 285  
Hansch, C., 204, 343  
Hansen, J.-P., 374  
Hansen, P., 275, 279, 282  
Hao, J., 283



- 
- Hao, Z., 409  
 Haque, T., 286  
 Harary, F., 276, 277  
 Harding, J. H., 122  
 Harris, F. E., 114, 115  
 Harris, K. D. M., 118  
 Harrison, N. M., 115, 118, 119, 120, 121, 122, 123, 124  
 Härtel, U., 122  
 Hartemink, A. J., 409  
 Hartshough, D., 347  
 Hasegawa, K., 347  
 Hashimoto, K., 409  
 Hasnip, P. J., 124  
 Hassan, M., 285  
 Hatcher, P. G., 284  
 Havlin, S., 410  
 Hawley, R. C., 373  
 Hay, P. J., 113, 201  
 Hayasaka, H., 281  
 Haynes, W., 123  
 Hays, T. R., 280  
 He, W. C., 278  
 He, W. J., 278  
 Head, M. S., 373  
 Heaney, P. J., 118  
 Heberger, K., 348  
 Hedin, L., 112  
 Hehre, W. J., 113  
 Heifets, E., 121, 122  
 Helgaker, T., 199  
 Helland, I. S., 345  
 Hempel, J., 285  
 Henderson, B., 123  
 Henderson, D., 373  
 Hendrickson, J. B., 283, 373  
 Hendrickson, T.  
 Hendsch, Z. S., 378  
 Henze, H. R., 276  
 Herman, F., 114, 276  
 Hermansson, K., 118, 120  
 Hermsmeier, M. A., 347  
 Herndon, W. C., 197  
 Herring, C., 113  
 Hess, W. P., 120  
 Hetzer, G., 115  
 Hewat, A. W., 116  
 Hill, T. L., 117  
 Hinman, V. F., 410  
 Hirano, H., 204  
 Hirono, S., 204  
 Hirshfeld, F. L., 200  
 Hochella, M. F., 121, 122  
 Hodgkin, E. E., 203, 379  
 Hoek, J. B., 409  
 Hoerl, A. E., 345  
 Hoffman, B. T., 347  
 Hoffmeyer, R. E., 117  
 Hohberger, R., 280  
 Hohenberg, P., 112, 199  
 Holland, J. H., 346  
 Holm, B., 119  
 Holm, C., 373  
 Holst, M. J., 375, 376  
 Holt, R. J., 374  
 Holthausen, M. C., 200  
 Holton, W. C., 124  
 Holzwarth, N. A. W., 125  
 Honig, B., 374, 375, 376, 377, 378, 379  
 Honig, H., 283  
 Hooft, R. W. W., 377  
 Hooshanagi, S., 411  
 Hopcroft, J. E., 278  
 Hopfinger, A. J., 346  
 Horn, F., 411  
 Horner, D. A., 124  
 Hosoya, H., 198  
 Hovig, E., 409  
 Hu, C., 283  
 Hu, Z. D., 347  
 Huang, D. M., 374  
 Hubank, M., 409  
 Hudson, A. T., 346  
 Hui, Y., 284  
 Huisinga, M., 120  
 Hummer, G., 373, 374  
 Hunenberger, P. H., 373, 376  
 Huo, S., 376  
 Husmeier, D., 409  
 Hutchison, C. A., 409  
 Hutter, J., 116, 125  
 Huxley, P., 379  
 Hybertson, M. S., 116  
 Hyde, R. M., 343, 344  
 Ihmels, J., 410  
 Ilin, A. V., 375, 379  
 Illas, F., 119, 122  
 Im, W., 374  
 Iroff, L. D., 204  
 Ishikawa, N., 409  
 Ishikawa, Y., 379  
 Itzkovitz, S., 410

- Iwasa, J., 204  
Izrailev, S., 348
- Jaakkola, T. S., 409  
Jackson, J. E., 344  
Jackson, J. D., 374  
Jackson, P., 346  
Jackson, R., 411  
Jacob, F., 408  
Jacobs, P. W. M., 122  
Jambeck, P., 379  
James, C. A., 285  
Jamois, E., 285  
Janak, J. F., 113  
Jandu, K. S., 346  
Janin, J., 378  
Janky, R., 409  
Janzen, W. P., 284  
Jeanloz, R., 119  
Jeanvoine, Y., 118  
Jennings, P. A., 377  
Jensen, F., 198  
Jensen, I., 279  
Jensen, L. H., 377  
Jeong, H., 410  
Jericevi, Z., 278  
Jerrum, M., 286  
Joachain, C. J., 202  
Johansson, B., 119  
Johnson, B. G., 202  
Johnson, D. S., 276  
Johnson, J. J. H., 284  
Johnson, M. A., 197, 198, 199  
Johnson, R. L., 118  
Jollet, F., 119  
Jones, A. D., 205  
Jorgensen, P., 199  
Joseph, S., 376  
Judson, R., 203, 281  
Junghans, M., 283  
Junker, J., 283  
Junquera, J., 125
- Kahl, S. D., 284  
Kahn, O., 119  
Kaijser, P., 202  
Kalmbach, A., 411  
Kangas, E., 379  
Karasev, Y. Z., 284  
Karelson, M., 197  
Karig, D., 411  
Karlin, S., 377
- Karplus, M., 347, 373, 377, 379  
Kashtan, N., 410  
Katritzky, A. R., 198  
Katz, J. L., 347  
Kaufman, J. J., 278  
Kearsley, S. K., 282  
Keasling, J. D., 408  
Kekicheff, P., 373  
Kempter, V., 120  
Kennard, R. W., 345  
Kerber, A., 280, 281, 283, 285  
Kertesz, M., 125  
Kestner, N. R., 117  
Kholodenko, B. N., 409  
Kholodovych, V. V., 347  
Kick, E., 284  
Kier, L. B., 198, 280, 281, 286, 343  
Kikuchi, O., 343  
Kimball, G. E., 206  
Kimura, T., 347  
King, P. J., 203  
Kinoshita, K., 378  
Kirby, E. C., 279  
Kirkwood, J. G., 377  
Kiyatkin, A., 409  
Klebe, G., 377  
Klee, V., 411  
Klein, D. J., 276, 279, 282  
Kley, A., 125  
Klopman, G., 204  
Klopper, W., 199  
Knop, J. V., 278, 278  
Koch, W., 200  
Kock, M., 283  
Koga, T., 117  
Kohn, H., 348  
Kohn, W., 112, 113, 199  
Kokalj, A., 121  
Kollman, P. A., 372, 376  
Kölmel, C. M., 123  
Komáromi, I., 123  
Konecny, R., 379  
Koopmans, T., 116  
Kopajtíc, T., 347  
Korringa, J., 113  
Kort, H. M., 204  
Korytko, A., 280, 283  
Koski, W. S., 278  
Kotomin, E., 118, 120, 121, 123  
Kotomin, E. A., 121, 122, 123, 124  
Kotu, A., 286  
Kovacevic, G., 198

- 
- Kovalishyn, V. V., 347  
 Kresse, G., 125  
 Krishnaiah, P. R., 345  
 Krishnamurthy, L., 410  
 Krishnan, R., 113  
 Kroemer, R. T., 203  
 Krüger, P., 116  
 Kubinyi, H., 285, 343, 346  
 Kudin, K. N., 115, 125  
 Kudo, Y., 276, 279  
 Kuh, E., 344  
 Kuhara, S., 277  
 Kühlenbeck, H., 122  
 Kuhn, B., 376  
 Kuntz, I. D., 284, 286  
 Kuppens, T., 203  
 Kurti, J., 125  
 Kutzelnigg, W., 113  
 Kvasnicka, V., 280  
  
 Laberty, C., 118  
 Ladd, M. F. C., 116  
 Ladik, J., 114  
 Laine, R., 377  
 Lam, C.-W., 282  
 Lamberti, C., 124  
 Lamm, G., 373  
 Landau, L. D., 379  
 Landolt, H., 116  
 Langenaeker, W., 197, 198, 199, 200, 202  
 Langlois, J.-M., 379  
 Langone, J. J., 344  
 Lansky, A., 284  
 Lapeña, F., 202  
 Larrieu, C., 124  
 Larson, R. S., 286  
 Larter, R., 373  
 Lathan, W. A., 113  
 Laue, R., 279, 280, 281, 283, 285  
 Lawless, J. F., 345  
 Lawless, M., 284  
 Lawton, J., 282  
 Lax, M., 117  
 Le Bret, C., 283  
 Le, H. M., 117  
 Leach, A. R., 276, 285, 372  
 Leardi, R., 346  
 Lebégue, S., 116  
 Lederberg, J., 277, 278, 279  
 Lee, B., 374  
 Lee, C., 116, 286  
 Lee, K.-H., 348  
 Lee, L. P., 378, 379  
 Lee, M., 376  
 Lee, M. S., 374  
 Lee, T., 376  
 Leherter, L., 202, 203  
 Leibler, S., 408  
 Leigh Lutz, M., 201  
 Leite, J. R., 114  
 Leland, B., 285  
 Lemer, C., 409  
 Lemmen, C., 203  
 Lengauer, T., 203, 284  
 Leo, A. J., 204, 343  
 Leroy, G., 114  
 Leslie, M., 123, 124  
 LeTiran, A., 348  
 Letscher, D., 411  
 Levine, I. N., 198  
 Levins, R., 411  
 Levy, M., 116  
 Levy, R. M., 373  
 Lewis, D. F. V., 343  
 Lewis, R. A., 284, 284  
 Leyda, L., 199  
 Li, H.-T., 348  
 Li, J., 285  
 Li, R., 276  
 Liang, S., 409  
 Liaw, A., 348  
 Lichanot, A., 119, 121, 124  
 Lide, D. R., 117  
 Lifshitz, E. M., 379  
 Lifson, S., 378  
 Lim, H. N., 408  
 Lin, D., 379  
 Lin, G., 276  
 Lin, J.-H., 373  
 Lin, S. L., 378  
 Lin, T.-H., 348  
 Lindan, P. J. D., 124  
 Lindel, T., 283  
 Lindsay, R. K., 278  
 Linial, M., 409  
 Lipinski, C. A., 198, 285  
 Lipkowitz, K. B., 117, 198, 201, 203, 276,  
     281, 284, 346, 373  
 Lippert, G., 125  
 Lipschultz, C. A., 378  
 Lischka, H., 113  
 Liu, G., 284  
 Liu, H. X., 347  
 Liu, H.-L., 347

- Liu, M. C., 347  
Liu, S.-S., 347  
Liu, X., 279  
Livesay, D. R., 379  
Livingstone, D. J., 343, 344, 345, 346, 347, 348  
Llunell, M., 115, 118, 119  
Lo Conte, L., 378  
Lobanov, V. S., 198, 285  
Lobato, M., 197, 207  
Loew, L. M., 409  
Lombardo, F., 198, 285  
Longo, E., 121  
Lopez Gejo, F., 120  
Lorenz, D., 409  
Louie, S. G., 115, 116  
Lounnas, V., 378  
Lovell, M. C., 345  
Löwdin, P.-O., 206  
Lu, L., 410  
Lu, Q., 377  
Lüchow, A., 115  
Luck, E. M., 276  
Ludemann, S. K., 378  
Luik, A. I., 346, 347  
Luinge, H. J., 283  
Lukovits, I., 278  
Lundqvist, S., 112  
Luo, R., 373, 377  
Luque, F. J., 207  
Luty, B. A., 375, 378, 379  
  
Ma, B., 378  
Ma, C., 376  
Ma, K., 284  
Macara, I., 409  
MacKerell, Jr., A. D., 372, 373, 379  
Mackey, M. C., 408  
Mackrodt, W. C., 114, 118, 120, 123, 124  
Madison, M., 280, 283  
Madura, J. D., 375, 379  
Maffett, T., 285  
Maggiora, G. M., 197, 198, 199, 203  
Maier, J., 121, 123  
Makov, G., 123  
Malinowski, E., 344  
Mallia, G., 118, 120, 123, 124  
Mallows, C. L., 346  
Maloney, P. P., 343  
Malpass, J. A., 345  
Mammone, J. F., 119  
Manallack, D. T., 346  
Manaut, F., 347  
Manly, B. F. J., 346  
Mannhold, R., 198  
Manolopoulos, D. E., 279  
Mao, H. K., 119  
Marchot, P., 379  
Marcus, P., 113  
Markham, R. L., 281  
Martin, D., 200  
Martin, E. M., 284  
Martin, G., 284  
Martin, M., 205  
Martin, S., 286  
Martin, Y. C., 284, 285  
Martínez, A., 199  
Martins, J. B. L., 121  
Martirosian, E. R., 283, 284  
Marzari, N., 115  
Maschio, L., 115  
Maseras, F., 123  
Masinter, L. M., 279  
Mason, C. H., 345  
Mason, J. S., 197  
Massova, I., 376  
Masulovic, D., 279  
Masunov, A., 117, 118  
Matsumoto, L. H., 377  
Matsuzaki, Y., 409  
Matthews, G. E., 125  
Mauri, F., 115  
Mays, P. K., 343  
Mazzatorta, P., 347  
McCammon, J. A., 373, 374, 375, 376, 377, 379  
McCarthy, M. I., 120, 123  
McCoy, A. J., 378  
McDonald, I. R., 374  
McFarland, J. W., 347  
McIntyre, S. H., 346  
McKay, B. D., 276, 277, 279  
McLachlan, A. D., 374  
McLaughlin, S., 376, 377, 378  
McLennan, K., 117  
McMahon, A. J., 203  
McMullan, R. K., 125  
McWeeny, R., 206  
Mehreja, R., 411  
Meiler, J., 282  
Melle-Franco, M., 122  
Mendes, P., 409  
Meng, T. C., 409  
Merawa, M., 119, 124

- 
- Meringer, M., 279, 285  
 Mermin, N. D., 117  
 Merz, Jr., K. M., 376  
 Mestres, J., 197, 201, 203, 205, 206, 207  
 Methfessel, M., 117  
 Metropolis, N., 281  
 Meurice, N., 202  
 Meyer, W., 113  
 Mezey, P. G., 197, 201, 202, 204, 205, 207  
 Michel, F., 120  
 Micu, A. M., 375  
 Mikami, M., 119  
 Milanesio, M., 118  
 Milicevic, A., 198  
 Miller, A. J., 345  
 Miller, R., 278  
 Millini, R., 124  
 Milo, R., 410  
 Mishra, R. K., 286  
 Mitchell, A. S., 117  
 Miyabayashi, N., 280  
 Miyano, S., 277  
 Miyashita, Y., 281  
 Miyoshi, F., 409  
 Mohan, S., 378  
 Moia, T. S., 122  
 Molchanova, M. S., 280  
 Molodtsov, S. G., 280, 283  
 Monkhurst, H. J., 114, 115, 117  
 Monod, J., 408  
 Montgomery, D. C., 345  
 Montgomery, D. B., 346  
 Moreira, I. d. P. R., 119  
 Morgan, R. S., 375  
 Morgan, S. L., 203  
 Moriguchi, I., 204  
 Morikis, D., 377  
 Morin, P., 376  
 Morokuma, K., 123  
 Morris, J. J., 343  
 Moruzzi, V. L., 113  
 Moser, D. A., 280  
 Moss, R. E., 206  
 Moulton, J., 373  
 Mueller, T., 378  
 Muir, R. M., 343  
 Müller, W. R., 278  
 Munk, M. E., 280, 283, 284  
 Murcko, M., 285  
 Murray, D., 285, 376, 377, 378  
 Murray, W. J., 343  
 Muscat, J., 119, 121, 122  
 Myklebost, O., 409  
 Nachman, I., 409  
 Nadeau, J., 410  
 Naito, Y., 409  
 Nakagome, I., 204  
 Nakamura, H., 378  
 Nakayama, Y., 409  
 Namboodiri, K., 204  
 Nanayakkara, A., 204  
 Narayana, N., 376  
 Navrotsky, A., 118  
 Nembra, R. M., 276  
 Netravali, I., 374, 411  
 Neugebauer, J., 125  
 Newman, A. H., 347  
 Newton, M. D., 113  
 Nguyen, A. T., 410  
 Nicholls, A., 374, 375, 379  
 Nicolotti, O., 347  
 Nielsen, J. E., 376, 377  
 Nijenhuis, A., 281  
 Nikolic, S., 198, 278  
 Nilakantan, R., 281  
 Nilges, M., 378  
 Nilsson, L., 373, 379  
 Nisanci, G., 410  
 Noguera, C., 123  
 Nonner, W., 373  
 Norman, Z., 276  
 Norvig, P., 275  
 Nourse, J. G., 280, 281, 285  
 Novak, B., 408  
 Novak, I., 282  
 Novotny, J., 378  
 Nusair, M., 112  
 Nussinov, R., 378  
 Nyberg, M., 120  
 Nye, J. F., 119  
 Nygren, M. A., 120  
 O'Brien, L. A., 204, 347  
 O'Dowd, A. B., 346  
 Oberoi, H., 377  
 Oganov, A. R., 115  
 Ohta, S., 409  
 Ojamäe, L., 118, 122  
 Olafson, B. D., 379  
 Oldehoeft, R. R., 379  
 Oldenburg, K. R., 284  
 Oliva, J. M., 205  
 Olsen, J., 199  
 Olson, A. J., 378

- Oltvai, Z. N., 410  
Ondrechen, M. J., 376  
Onida, G., 116  
Onufriev, A., 373, 374, 376  
Oppenheimer, J. R., 119  
Oprea, T. I., 198  
Ordejón, P., 125  
Orlando, R., 115, 117, 118, 119, 120, 122, 123, 124  
Orozco, M., 207  
Osapay, K., 373  
Osawa, E., 282  
Oschkinat, H., 378  
Ostlund, N. S., 198  
Othmer, H. G., 408  
Otter, R., 276  
Ouazzani, T., 121  
Overbeek, J. T. G., 379  
Overhof, H., 123  
Ozbudak, E. M., 408  
Ozsoyoglu, G., 410  
Ozsoyoglu, M., 410  
Ozturk, M., 410
- Pacchioni, G., 122, 123  
Pack, J. D., 117  
Palmer, E. M., 277  
Palsson, B. Ø., 410, 411  
Palyulin, V. A., 281, 286  
Papin, J. A., 410  
Parker, Jr., L. R., 203  
Parks, C. A., 283  
Parmigiani, F., 123  
Parr, R. G., 116, 199, 204  
Parretti, M. F., 203  
Parrinello, M., 115, 125  
Pascale, F., 120  
Pascual, J.-L., 120  
Pasquarello, A., 115  
Pastori Parravicini, P., 113  
Pasucci, V., 374  
Patel, D. V., 284  
Patterson, D., 284  
Patterson, D. E., 206  
Paulaitis, M. E., 373  
Pavone, P., 119  
Paxton, A. T., 117  
Payne, M. C., 123, 124  
Pe'er, D., 409  
Pearlman, R. S., 286  
Pearlmann, D. A., 376  
Peck, E. A., 345
- Peeters, A., 200  
Peng, C., 284  
Perdew, J. P., 112  
Perego, G., 124  
Perram, J. W., 114  
Perreault, W. D., 345  
Perrey, S. W., 279  
Perrin, D., 408  
Perry, J. K., 379  
Persson, P., 122  
Petersilka, M., 116  
Peterson, K. L., 346  
Petke, J. D., 203  
Petrovic, I., 118  
Pettersson, L. G. M., 120  
Pettitt, B. M., 375  
Pfrommer, B. G., 115, 125  
Phillies, G. D. J., 376  
Piccione, P. M., 118  
Pickard, C. J., 124  
Pickett, S. D., 197, 284  
Pickup, B. T., 374  
Pierotti, R. A., 374  
Pilar, F. L., 198  
Pinter, R. Y., 410  
Piriou, B., 120  
Pisani, C., 114, 115, 116, 118, 120, 121, 122, 123  
Pisanski, T., 279  
Pitaevskii, L. P., 379  
Pitera, J., 376  
Podgornik, R., 373  
Pollack, P., 279  
Pollmann, J., 116  
Pólya, G., 275, 276  
Ponder, J. W., 372  
Ponec, R., 202, 204, 205  
Pons, V., 205  
Popelier, P. L. A., 204  
Pophale, R. S., 281, 286  
Pople, J. A., 113, 117, 202  
Portmann, P., 283  
Pospichal, J., 280  
Prabhakar, Y. S., 347  
Pratt, L. R., 373, 374  
Prencipe, M., 117  
Press, W. H., 203  
Pretsch, E., 283  
Preuss, H., 114  
Price, D., 115  
Price, N. D., 410  
Probert, M. J., 124

- 
- Puccia, C. J., 411  
 Puchin, V. E., 120, 124  
 Puchina, A. V., 120, 124  
 Pudenz, S., 286  
 Pulay, P., 113  
 Pun, F. C., 346  
 Purisima, E. O., 378  
  
 Quirk, J., 411  
  
 Radmer, R., 376  
 Raghavachari, K., 117  
 Rahr, E., 344  
 Rajczy, P., 125  
 Rajko, R., 348  
 Ralston, A., 345  
 Ramanadham, M., 377  
 Ramer, N. J., 114  
 Ramnarayan, K., 285  
 Randic, M., 198, 343, 376  
 Rappé, A. M., 114  
 Rarey, M., 284  
 Raty, J.-Y., 119  
 Raubenheimer, P., 117  
 Read, R. C., 276, 277  
 Reddy, A. K. N., 373  
 Redfern, P. C., 124  
 Redmond, D. B., 276  
 Regev, A., 409  
 Rehbein, C., 120  
 Reichling, M., 120, 124  
 Reinders, M. J. T., 408  
 Reinhardt, P., 119  
 Reining, L., 116  
 Rekker, R. F., 204  
 Ren, P., 372  
 Rencher, A. C., 346  
 Resca, L., 121  
 Resh, M. D., 377  
 Resta, R., 116, 121  
 Reyes, C., 376  
 Reynders, J. V. W., 379  
 Ricca, F., 120  
 Richard, A. M., 379  
 Richards, F. M., 374  
 Richards, W. G., 203, 379  
 Richelle, J., 409  
 Richert, J., 283  
 Richert, J. R., 280  
 Richet, P., 119  
 Riedwyl, H., 344  
 Riekel, C., 116  
  
 Rignanese, G.-M., 119  
 Ringe, D., 376  
 Rintoul, M. D., 286  
 Rioux, F., 202  
 Rivail, J.-L., 344  
 Riveros, M., 278  
 Robert, D., 197, 199, 203, 204, 205, 206  
 Roberts, F., 411  
 Roberts, V. A., 378  
 Robertson, A. V., 277  
 Rocchia, W., 375  
 Rodriguez-Rodrigo, M., 122  
 Roe, D. C., 284  
 Roetti, C., 114, 115, 116, 117, 118, 119, 120, 121, 123, 124  
 Rogers, D., 346  
 Rohlfing, M., 116  
 Rojnuckarin, A., 379  
 Roman, R. V., 280  
 Rongsi, C., 282  
 Roos, B. O., 199  
 Ros, P., 200  
 Rose, V. S., 343, 346  
 Rosen, R., 408  
 Rosenbluth, A. W., 281  
 Ross, D. K., 376  
 Ross, W. S., 376  
 Rosso, K. M., 121, 122  
 Rostoker, N., 113  
 Rothman, J. H., 203  
 Rouvray, D. H., 197, 275  
 Roux, B., 372, 373, 374, 379  
 Roy, A., 119  
 Royle, G. F., 279  
 Rozas, R., 278  
 Rubio, A., 117  
 Rudolph, C., 286  
 Ruecker, C., 198  
 Ruecker, G., 198  
 Runge, E., 115  
 Ruppert, R., 411  
 Rurali, R., 123  
 Russell, S. J., 275  
  
 Sabin, J. R., 200  
 Sadamoto, Y., 409  
 Sadowski, J., 206, 285, 286  
 Saffern, M. M., 114  
 Sah, C.-H., 279  
 Saied, F., 375  
 Sairam, A., 409  
 Saito, K., 409

- Saito, M., 117  
Sakharkar, K., 409  
Salasco, L., 115  
Sali, A., 379  
Salt, D. W., 345  
Salvetti, O., 112  
Sambe, H., 200  
Sambrano, J. R., 121  
Sampogna, R., 377  
Sanchez, C., 408  
Sánchez-Portal, D., 125  
Sandall, J. P. B., 276  
Sandrone, G., 119  
SanFeliciano, S. G., 282  
Santolaria, X., 409  
Sanz, F., 347  
Sanz, J. F., 122  
Sarig, O., 410  
Sasaki, S.-I., 276, 279, 280, 281  
Sasamoto, H., 409  
Sattath, S., 410  
Sauer, J., 118, 123  
Saunders, M. R., 343, 344  
Saunders, V. R., 113, 114, 115, 117, 118, 119,  
120, 123, 124  
Sautet, P., 121  
Savin, A., 116  
Savrasov, S. Y., 119  
Saxena, A. K., 118, 343  
Schaefer, III, H. F., 200  
Schaefer, M., 373, 377  
Schaeffer, G., 410  
Schaff, J. C., 409  
Scheffler, M., 125  
Scheraga, H. A., 276  
Schilling, C., 411  
Schimpl, J., 125  
Schirmer, O. F., 124  
Schlegel, H. B., 113  
Schleyer, P. v. R., 379  
Schlick, T., 372  
Schlüter, M., 113  
Schneider, G., 343  
Schochet, M., 204  
Schrijver, A., 280  
Schröder, K. P., 118  
Schulz, K. P., 280, 283  
Schütz, M., 115  
Schoor, J., 286  
Schwaber, J. S., 409  
Schwarz, K., 125  
Scott, L. R., 375, 379  
Scuseria, G. E., 115, 125  
Searle, B., 120  
Sebald, W., 378  
Segall, M. D., 124  
Seibel, G. L., 376  
Seidel, H., 123  
Seidel, J. J., 279  
Seitz, F., 112  
Selwood, D. L., 346  
Selzer, P., 286  
Seminario, I., 116  
Sen, K. D., 199  
Sensato, F. R., 121  
Sept, D., 375, 376  
Sgroi, M., 122  
Sham, L. J., 112  
Shao, X., 283  
Sharp, K. A., 374, 375, 377, 378, 379  
Shavitt, I., 113  
Shcherbukhin, V. V., 280  
Shelley, C. A., 280  
Shen, G., 118  
Shen, M., 348  
Shen-Orr, S. S., 410  
Sheridan, R. P., 282, 348  
Shimizu, T., 409  
Shiozawa, A., 409  
Shluger, A. L., 120  
Shoemaker, J. R., 123  
Shraiman, B. I., 408.  
Sieker, L. C., 377  
Sierka, M., 118  
Sigal, C. T., 377  
Sillerud, L. O., 286  
Silman, I., 378  
Silvi, B., 118  
Simmerling, C., 376  
Simon, H., 125  
Simonson, T., 373, 374  
Simpson, A., 284  
Simpson, W. J., 201  
Sindic, L., 119  
Singh, U. C., 376  
Sinha, N., 378  
Sittampalam, G. S., 284  
Skillman, A. G., 284, 286  
Skvortsova, M. I., 281, 286  
Slater, J. C., 113, 114  
Slepchenko, B. M., 409  
Slovokhotova, O. L., 281  
Smirnov, V. P., 124  
Smith, Jr., V. H., 202



- Smith, C. M., 200  
 Smith, D., 285  
 Smith, D. H., 279, 280, 281  
 Smith, E. R., 114  
 Smith, G., 345  
 Smith, P. E., 373  
 Smith, P. J., 204  
 Smith, W. C. S., 345  
 Smith-Gill, S. J., 378  
 Snijders, J. G., 125  
 So, S.-S., 347  
 Solà, M., 197, 201, 205, 206, 207  
 Soler, J. M., 125  
 Soltanshahi, F., 284  
 Somani, S., 409  
 Somborg, H., 284  
 Sommer, M., 377  
 Sommerer, S. O., 201  
 Somogyi, R., 409  
 Sontag, E., 409  
 Sorensen, P. B., 286  
 Spackman, M. A., 117  
 Spellman, P. T., 409  
 Spellmeyer, D. C., 284  
 Sperb, R., 373  
 Sprik, M., 116  
 Sridharan, S., 279, 375  
 Srinivasan, J., 376  
 Srinivasan, V., 346  
 Stables, J. N., 346  
 Stahelin, R. V., 378  
 Stämmeler, V., 113  
 Stanton, D. T., 344  
 Stanton, R. V., 376  
 Stark, J., 409  
 States, D. J., 379  
 Stefanov, B. B., 203  
 Steinbeck, C., 282  
 Steinhauer, L., 286  
 Steinhauer, V., 286  
 Stevanovic, D., 282  
 Stewart, J. J. P., 198  
 Stewart, R. F., 201  
 Stiefl, N., 348  
 Still, W. C., 373  
 Stillinger, F., 374  
 Stojmenovi, I., 278, 279  
 Stoll, H., 114, 115  
 Stone, M., 345  
 Stoneham, A. M., 122, 123  
 Stove, J., 279  
 Strachan, A., 284  
 Strnad, M., 202  
 Strogatz, S. H., 410  
 Studier, F. W., 379  
 Subramaniam, S., 379  
 Suhai, S., 115  
 Sulea, T., 378  
 Sullivan, J. J., 205  
 Sun, G. Y., 125  
 Sun, J. Q., 115  
 Sun, Y., 284  
 Suñé, E., 202  
 Surratt, G. T., 114  
 Sushkin, N. V., 376  
 Sussman, J. L., 378  
 Sutcliffe, B. T., 113, 206  
 Sutherland, G. L., 277  
 Sutherland, J. J., 347  
 Svetnik, V., 348  
 Swaminathan, S., 125, 379  
 Swamy, V., 119  
 Sykes, R., 285  
 Szabo, A., 198  
 Szymanski, K., 278  
 Tackett, A. R., 125  
 Tainer, J. A., 378  
 Takagi, T., 410  
 Takahashi, K., 409  
 Tanford, C., 377  
 Tang, Y., 410  
 Tanida, S., 409  
 Tanji, K. K., 205  
 Tarjan, R. E., 278, 344  
 Tasan, M., 410  
 Tasker, P. W., 120  
 Tathan, R. L., 344  
 Taylor, D. P., 120  
 Taylor, P., 376, 379  
 Taylor, S., 202, 376  
 te Velde, G., 125  
 Teig, S., 285  
 Tempczyk, A., 373  
 Terrile, M., 346  
 Teter, M. P., 120  
 Tetko, I. V., 346, 347  
 Teukolsky, S. A., 203  
 Thakkar, A. J., 117  
 Thatcher, E., 278  
 Thattai, M., 408  
 Thiele, H., 284  
 Tholburn, M., 379  
 Thompson, L. A., 284

- Thompson, M. L., 345  
Thomsen, M., 286  
Thornton, G., 121  
Tidor, B., 378, 379  
Tironi, I. G., 373  
Todeschini, R., 343, 347  
Toi, R., 278  
Tollenaere, J. P., 197, 198, 199, 200, 202  
Tombor, B., 410  
Tomescu, D., 409  
Tomita, M., 409  
Tong, C., 348  
Torrent, M., 119  
Tosic, R., 279  
Tou, J. T., 203  
Towler, M. D., 118, 123, 124  
Trevisiol, G. M., 278  
Trinajstić, N., 198, 277, 278  
Tropsha, A., 346, 347, 348  
Trucks, G. W., 117  
Tsai, K.-C., 348  
Tsang, J., 377  
Tsui, V., 376  
Turnbull, D., 112,  
Turnstall-Pedoe, H., 345  
Tversky, A., 197  
Tyson, J. J., 408  
  
Ugliengo, P., 117, 118, 119, 120, 121, 122, 123  
Ullmann, G. M., 376  
Umari, P., 115  
Umrigar, C. J., 116  
Unruh, W. P., 124  
  
Vacek, G., 379  
Valdivia, R., 278  
Valeri, S., 122  
Valerio, G., 118, 119  
van Almsick, M., 283  
Van Alsenoy, C., 200, 201, 202  
Van Damme, S., 201  
van de Waterbeemd, H., 198, 345, 347, 348  
van Gisbergen, S. J. A., 125  
van Gunsteren, W. F., 373  
Van Lenthe, J. H., 113  
Van Nostrand, R. C., 345  
van Oudenaarden, A., 408  
van Someren, E. P., 408  
van Vlijmen, H. W. T., 377  
Vanderbilt, D., 114, 115  
  
Vanossi, D., 276  
Veillard, A., 113  
Veithen, M., 115  
Venkataraghavan, R., 281  
Venkatasubramanian, V., 282  
Venter, J. C., 409  
Vera, L., 200  
Vercauteren, D. P., 202, 203  
Verstaete, M., 119  
Verwey, E. J. W., 379  
Vetterling, W. T., 203  
Vilar, J. M. G., 408  
Villa, A. E. P., 346, 347  
Vincent, J. J., 376  
Vinogradov, I. M., 199  
Visco, Jr., D. P., 281, 286  
Viterbo, D., 118  
Vocadlo, D. J., 377  
Voge, M., 279  
von Korff, M., 344, 348  
von Neumann, J., 199  
Vosko, S. H., 112  
Vracko, M., 347  
Vreven, T., 123  
Vriend, G., 376, 377  
  
Wade, R. C., 377, 378, 379  
Wadt, W. R., 201  
Wagener, M., 206  
Wagner, A., 410  
Waldman, M., 285  
Walter, H. J., 206  
Walters, F. H., 203  
Walters, W., 285  
Wander, A., 120, 121, 122  
Wang, F., 375  
Wang, J., 117, 276, 285, 376, 408  
Wang, J. M., 372  
Wang, L.-S., 347  
Wang, P., 345  
Wang, Q. X., 278  
Wang, S., 276  
Wang, W., 376  
Wang, Z.-G., 373  
Waszkowycz, B., 285  
Watanabe, M., 372, 373  
Watts, D. J., 410  
Weaver, D. F., 347  
Weber, L., 285  
Wegner, J. K., 347  
Weidinger, J., 285  
Weiner, P., 376  
Weiss, R., 411

- 
- Weissig, H., 374, 379  
 Weisstein, E. W., 199  
 Weitz, B. A., 346  
 Welford, S., 295  
 Welsch, R. E., 344  
 Wen-sheng, C., 283  
 Wepfer, G. G., 114  
 Werner, H. J., 115  
 Wessels, L. F. A., 408  
 Wesson, L., 374  
 Westbrook, J., 379  
 Westerhoff, H. V., 409  
 Westhead, D., 285  
 Whitley, D. C., 344, 347  
 Wichtendahl, R., 122  
 Wieland, T., 280, 281  
 Wiener, H., 198  
 Wigner, E., 112  
 Wikel, J. H., 346  
 Wilf, H. S., 281, 345  
 Wilhite, D. L., 114  
 Wilk, L., 112  
 Wilkins, C. L., 198  
 Wilkinson, L., 345, 346  
 Will, M., 280, 282, 283  
 Willett, P., 203, 285, 344  
 Williams, A., 284  
 Williams, A. J., 283  
 Williams, A. R., 113  
 Wilson, S., 113, 123  
 Winkler, D. A., 347  
 Wipke, T., 282  
 Withers, S. G., 377  
 Wodak, S. J., 409  
 Wold, H., 345  
 Wolf, D. M., 408  
 Wolf, H. C., 123  
 Wolfson, H., 378  
 Won, Y., 379  
 Wong, C. F., 376  
 Wong, P., 408  
 Woodward, M., 345  
 Wormald, N. C., 281  
 Wu, H., 284  
 Wuchty, S., 410  
 Wynn, E. W., 345  
  
 Xiao, Y.-D., 348  
 Xiao, Z., 348  
 Xiaofeng, G., 277, 282  
 Xu, D., 378  
 Xu, L., 283  
 Xu, W., 410  
  
 Xuong, N. H., 376  
 Yang, A. S., 377  
 Yang, S. Y., 118  
 Yang, W., 116, 199, 200  
 Yao, X. J., 347  
 Yaschenko, E., 121  
 Yasri, A., 347  
 Ye, L., 409  
 Yeger-Lotem, E., 410  
 Yeh, C. Y., 276, 277  
 Yildirim, N., 408  
 Yin, C.-S., 347  
 Yoshida, M., 282  
 Yosida, K., 119  
 Young, R. A., 409  
 Young, S., 285  
 Young, W. S., 373  
 Yuan, S., 284  
 Yugi, K., 409  
  
 Zak, D. E., 409  
 Zaltina, L. A., 281  
 Zapol, P., 118, 124  
 Zauhar, R. J., 375  
 Zecchina, A., 122, 124  
 Zechel, C., 284  
 Zeffrov, N. S., 280, 281, 286  
 Zell, A., 347  
 Zerah, G., 119  
 Zerovnik, J., 279  
 Zewail, A. H., 373  
 Zhang, D., 379  
 Zhang, F., 276, 277  
 Zhang, J. P., 283  
 Zhang, M., 283  
 Zhang, R. S., 347  
 Zheng, C., 284  
 Zheng, M., 279  
 Zheng, W., 347  
 Zhou, H. X., 375  
 Zhu, S. Y., 283  
 Zhu, Z. Y., 377  
 Zhukovskii, Y. F., 122, 123  
 Zicovich-Wilson, C. M., 115, 117, 118, 120,  
     122, 124  
 Ziegler, T., 125  
 Ziv, G., 410  
 Ziv, Y., 410  
 Zivkovic, T. P., 282  
 Zones, S. I., 118  
 Zunger, A., 112, 113, 114  
 Zupan, A., 117  
 Zupan, J., 280, 283



---

# Subject Index

---

Computer programs are denoted in boldface; databases and journals are in italics.

- Ab initio ASA, 150
- Ab initio calculations, 151
- Ab initio MQSM, 144
- Ab initio quantum simulation, 1, 40
- Absolute configuration, 177
- Abstract models, 388
- Abstractions, 400
- ACD/Structure Elucidator**, 270
- Achiral compounds, 228
- Acoustic modes, 64
- $\alpha$ -Cristabolite, 52
- ACSF*, 394
- Active sites, 358
- Acyclic compounds, 224, 227
- Acyclic molecular graph enumeration, 238
- Additive group additions, 171
- ADEPT**, 272
- Adhesion energy, 77
- Adiabatic approximation, 58
- Adjacency matrix, 212, 215, 241
- ADME, 270
- Adsorption energy, 44, 74
- Aggregated trees, 341
- Aldehydes, 224
- Algebraic multigrid methods, 359
- Algorithm for the Exhaustive Generation of Irredundant Structures (AEGIS), 267
- Alignment, 154, 157, 163, 191
- Alignment algorithm, 194
- Alignment free methods, 163
- Aliphatic systems, 287
- Alkaline earth oxides, 90
- Alkane isomers, 226, 266
- Alkanes, 225, 261, 274
- n-Alkanes, 289
- Alkenes, 261
- Alkyl group, 222
- Alkynes, 224, 261
- All subset regression, 326, 332
- All-electron basis sets, 45
- aMAZE*, 394
- AMBER**, 368, 370
- AmiGO**, 396
- Ammonia, 7, 10
- Amorphous solid, 81
- Analysis of variance, 299
- Analytical gradients, 2
- Anisotropic coupling constants, 96
- Anisotropy, 4, 5
- Annealing schedule, 258
- Anthracene, 221, 266
- Antibonding crystal orbital, 28
- Antiferromagnetic (AFM) phase, 54
- APBS**, 360, 361, 368, 370
- Apolar interactions, 353
- Apolar solvation energies, 366
- Apolar surface area, 366
- Approximate electron density, 144
- Aqueous environment, 350
- ArrayExpress*, 392
- Artificial ant colonies, 341
- Artificial intelligence, 210
- Artificial neural networks (ANNs), 295, 338
- ASSEMBLE**, 248, 249, 267
- Asymmetric unit, 10
- Atom equivalent classes, 248

*Reviews in Computational Chemistry, Volume 21*  
edited by Kenny B. Lipkowitz, Raima Larter, and Thomas R. Cundari  
Copyright © 2005 Wiley-VCH, John Wiley & Sons, Inc.

- Atom valence, 212  
Atomic charges, 131  
Atomic orbitals, 16  
Atomic shell approximation (ASA), 144, 148, 159, 177  
Atomic signatures, 250  
Atom-in-molecule densities, 154  
Atom-in-molecule similarity, 167, 168  
 $\alpha$ -tridymite, 52  
Attractor, 169  
Attractor basins, 169  
Augmented plane waves (APW), 2  
Automatic relevance determination (ARD), 340  
Automorphism, 214  
Automorphism group, 215, 218, 219  
  
B3LYP, 3, 35, 47, 50, 59, 65, 89, 94, 148  
Backpropagation, 338  
Backward elimination, 323  
Bacterial genomes, 384  
Band gap, 36  
Band structure, 2, 21, 33, 89, 95  
Band theory, 6, 3, 17  
Bandwidth, 33  
Basic vectors, 8  
Basis set, 2, 62, 108, 152  
Basis set incompleteness, 45  
Basis sets  
  2-1G, 99  
  3-21G, 108, 175  
  6-21G, 99, 108  
  6-311G(d, p), 52, 108  
  6-31G, 108  
  6-31G(d, p), 50, 52, 108  
  6-31G\*, 148, 191  
Basis set superposition error (BSSE), 50, 51, 75, 77  
Bayesian networks, 389  
Bayesian neural networks, 340  
Benzene, 216, 218, 221, 266  
  1,2-Dichlorobenzene, 212  
  1,4-Dichlorobenzene, 212  
Benzenoid hydrocarbons, 221, 228, 241, 242, 263, 264  
Benzodiazepine, 271  
Beryllium, 34, 38, 47, 59  
Best subset regression, 332  
Best subset variable selection, 318  
BFGS, 160  
Bias neurons, 338  
Biased regression, 312, 319  
  
Biased regression coefficient, 319  
Binary interaction database, 393  
Binary particle swarms, 341  
*BIND*, 394  
Binding energy, 51  
Binding free energies, 367  
Binding sites, 358  
*BindingDB*, 394  
*BioCarta*, 394  
*BioCyc*, 394  
Bioentity, 398  
Bioinformatics, 390  
Biological processes, 396  
Biological response, 291  
Biomolecular interiors, 352  
Bio-ontology, 395  
Bistability, 387, 406  
Bjerrum length, 354  
Bloch function, 12, 10  
Bloch theorem, 12, 14, 19, 63  
BLYP, 89, 94  
Boiling point, 260  
Bond connectivity, 288  
Bond critical points, 169  
Bond order, 212  
Bond order switch, 259  
Bonded valence graph, 215  
Bonding crystal orbital, 28  
Bootstrapped trees, 341  
Born ion, 352, 363  
Born-Haber cycle, 48  
Born-Oppenheimer approximation, 16, 63  
Boron nitride, 3, 33  
Boundary edges code (BEC), 242, 245  
Boundary element (BE) methods, 358  
 $\beta$ -quartz, 52  
**BRABO**, 144, 194  
Bravais lattices, 8  
*BRENDA*, 394  
Brillouin zone, 11, 5, 17, 24, 37, 38  
*BRITE*, 394  
Bulk basis set, 46  
Bulk modulus, 58, 61, 62  
Bulk total energy, 46  
  
 $C_{60}$ , 221  
Candidate compounds, 270  
Canonical labeling, 214  
Canonical n-tuple, 238  
Canonization procedure, 214, 236, 238  
Canonized molecular graphs, 214  
CaO, 60

- Carbó index, 136, 141, 161, 165, 194  
Carbon monoxide (CO), 75  
Car-Parrinello (CP) methodology, 5  
Catacondensed benzenoids, 228  
Catafusenes, 230  
Categorization, 314  
*CATH*, 392  
Cayley's counting formula, 218, 228  
CCSD, 47  
CCSD(T), 47, 59  
*CDD (Conservative Domain Database)*, 392  
*CE*, 392  
Cell compartments, 398  
Cell cycle checkpoint, 401  
Cell parameters, 8  
Cellular component, 396  
Cellular functions, 389  
Cellular process ontology, 397  
Cellular regulatory pathways, 397  
**CellWare**, 396  
Chabazite, 52, 57  
Chain connectivity, 289  
Charge cloud, 355  
Charge distributions, 349  
Charged defect, 84  
Charge-transfer complexes, 292, 329  
**CHARMM**, 360, 361  
Chemical diversity, 270  
Chemical information, 261  
Chemical information content, 162  
Chemical intuition, 163  
**Chemical Reaction Network Toolbox**, 407  
Chemical relevance, 162  
Chemical structures, 211  
Chemical topology, 152, 161  
**CHEMICS**, 247, 248, 249, 267  
Chemisorption, 68, 74  
Chi indices, 273  
Chiral compounds, 228  
Chirality, 169, 170, 177, 263  
Chirality metric, 177  
Chlorobenzene isomers, 220  
Cholesky factorization, 319  
Chromosomes, 160  
*CIBEX*, 392  
Circular chromosome, 384  
**CISOC-SES**, 269  
City-block distance, 135  
Classification, 128  
Classification trees, 341  
Cluster analysis, 295  
Cluster approach, 81, 82, 102  
Cluster connectivity, 289  
Cluster significance analysis, 340  
Cluster-in-cluster scheme, 82  
Cluster-in-crystal embedding, 77  
Clustering coefficient, 403  
Clustering methods, 389  
Coarse-grain network, 389  
**COCON**, 268  
Cohesive energy, 44, 47, 51  
Colinearity, 289  
**CombiBUILD**, 271  
**CombiDOCK**, 271  
**CombiLibMaker**, 272  
Combinatorial chemistry, 210  
Combinatorial libraries, 213, 217, 222, 235, 261, 272  
Combinatorics, 209  
Compartments, 398  
Competition bias, 320  
Compressed adjacency matrix (CAM), 241  
Computer-aided drug design (CADD), 287, 290, 339  
**Computer-Assisted Molecular Generation and Counting (CAMGEC)**, 266  
Condensed phase, 7  
Conduction bands, 33  
Conduction-valence gap, 33, 34, 38  
Conductor, 34  
Configuration generating function, 219  
Conformational analysis, 174  
Conformational changes, 365, 367  
Conformational degrees of freedom, 163  
Conformational free energies, 365  
Connected labeled graphs, 217  
Connectivity degree, 288  
Connectivity stack, 213, 214, 215, 258  
**Constrained Combination of Atom-centered fragments (COCOA)**, 268  
**CONstrained GENerator (CONGEN)**, 247, 253, 257, 267  
Constraints, 249, 273  
Construction set, 328  
Continuum modeling, 350  
Continuum regression (CR), 310, 314, 315  
Convergence, 92  
Convergent structure generation, 252  
**CORCHOP**, 308  
Core bands, 33  
Core electron density, 150, 151, 158  
Core electrons, 33  
Coronoids, 228  
Correlation effects, 47

- Correlation matrix, 290, 304  
Correspondence principle, 132  
Coulomb's law, 352, 364  
Counterion correlation, 356  
Counterion fluctuation, 356  
Counterions, 356  
Counterpoise (CP) method, 50, 51, 75  
Counting, 211  
Counting formula, 230  
Counting series, 219, 227, 231, 232  
Covalent crystals, 44  
Covariance, 316  
Critical points, 154  
Crossover, 160, 260  
Cross-validation (CV), 328, 329  
**CRYSTAL**, 3, 62, 103, 106, 110  
Crystal formation energy, 47  
Crystal morphology, 71  
Crystalline orbitals, 16, 11, 45  
Crystalline urea, 49, 108  
Crystallography, 4  
Cubic graphs, 245  
Cusp condition, 134  
Cycle decomposition, 218  
Cycle index, 218, 219  
Cyclic adenosine monophosphate (cAMP), 385  
Cyclic graphs, 256  
Cyclic substructures, 247  
1,2,3,4-tetrachlorocyclobutane, 254  
*Cytoscape*, 394
- Dangling bonds, 68, 71, 82, 102  
Darwin's evolutionary rules, 260  
Data management, 391  
Data matrix rank, 298  
Databases, 381, 383, 390  
Dataset dimensionality, 291  
Dataset redundancy, 291  
*DDBJ (DNA Databank of Japan)*, 391  
Debye length, 354  
Debye-Hückel law, 351  
Decane, 215, 227, 233, 255  
Defect formation energy, 85, 90, 97  
Defect zone, 81  
Defect-defect interaction, 92  
Defective ionic surfaces, 74  
Defects, 80, 84  
Deficiency theory, 406  
Deformed MQSM, 158  
Degree distribution, 403  
Degree of chirality, 177  
Degree of multicollinearity, 305, 336  
Degree of similarity, 129  
Degree sequence, 212  
**DelPhi**, 360, 361  
**DENDRITic ALgorithm (DENDRAL)**, 210, 238, 247, 248, 250, 267  
Dendrograms, 141, 143, 194, 295  
Density fitting, 146  
Density functional perturbation theory, 36  
Density functional theory (DFT), 1, 35  
Derivatives, 43, 57  
Descriptor inter-relationships, 294  
Descriptor tables, 288  
Descriptors, 272, 294  
Desolvation energy, 364  
DFT, 133  
Diamond, 3  
Dielectric boundary pressure, 356  
Dielectric constant, 351  
Dielectric saturation, 363  
Diffuse AOs, 45  
Dimension reduction, 290, 291  
Diophantine integration schemes, 2  
*DIP (Database of Interacting Proteins)*, 394  
Dipolar surfaces, 68  
Direct lattice, 7  
Direct space, 11, 39  
Discontinuities, 153  
Discretization methods, 357  
Dislocations, 80  
Disordered solid, 81  
Dispersion in k-space, 33  
Dissimilarity measures, 136, 141  
Dissociation free energy, 368  
Distance-dependent dielectric functions, 352  
Divide-and-conquer, 272  
DNA, 356, 382  
DNA damage, 401  
DNA microarrays, 389  
Donor-acceptor complexes, 292  
Drug database, 257  
Drug design, 134  
Drug discovery, 270  
Drug-like molecule, 129, 130, 171  
Drug-likeness, 270  
Dummy variables, 296
- E-Cell**, 396  
Edge multiplicity, 211, 212  
Edges, 80, 211  
Edingtonite, 52  
Effective core potentials (ECPs), 151  
Effective mass, 2



- Elastic constants, 61, 62
- Elastic tensor, 58, 62
- Electron correlation, 47, 153
- Electron density, 129, 132, 133, 177, 194
- Electrostatic focusing, 358
- Electrostatic free energies, 355
- Electrostatic interactions, 349
- Electrostatic potential, 4, 153, 165, 352, 354
- Electrostatic properties, 371
- Elementary cycles, 247
- Elitism, 160
- Elucidation by Progressive Intersection of Ordered Substructures (EPIOS)**, 250, 267
- Embedded cluster, 83
- EMBL database*, 391
- eMOTIF*, 392
- Enantiomeric form, 227
- Enantiomers, 170, 177
- ENDOR, 86, 87, 90, 96
- Energy differences, 43
- Enumerating molecules, 209, 237, 261
- Enumeration, 211, 244, 272
- Environmental studies, 288
- Enzyme Commission list of enzymes (EC numbers), 396
- EPR, 86, 87, 90, 96
- Equilibrium geometry, 43, 58, 84
- Equilibrium unit cell, 58
- Equivalent classes, 215, 249
- Escherichia coli*, 384
- Esters, 225, 263
- Ethane, 175
- Euclidean distance, 134, 163, 167, 295, 312
- Euler angles, 157, 159
- Ewald's method, 21
- Excitation spectra, 2, 33
- Excited states, 5
- Expert system, 238, 266
- Expert System for the Elucidation of the Structures of Organic Compounds (ESESOC)**, 268
- Explicit enumeration, 233
- Explicit particles, 350
- Explicit solvent methods, 351
- Extended defects, 80
- Extreme pathways, 404, 405
- Factor analysis (FA), 291, 292
- Factor loadings plot, 294
- False-positive tests, 271
- Fast PB methods, 371
- Faujasite, 40, 52
- F-center, 85, 86, 92
- Feature counts, 130
- Feedback interactions, 382
- Feedback loops, 384, 387, 402
- Fermi energy, 19, 33, 39
- Fermi surface, 2, 39
- Ferromagnetic (FM) phase, 54
- Field-based descriptors, 131
- Figure generating function, 219, 222
- Filters, 271
- Finite basis set, 51
- Finite cluster, 81
- Finite difference (FD) methods, 358
- Finite element (FE) discretization, 359
- First-order connectivity index, 289
- First-order electron density, 132, 143, 153
- Fisher's discriminant ratio, 341
- Fitness, 160
- Fitness function, 260, 340
- Fixed charge distribution, 354
- Flexibility, 368
- Fock matrix, 41
- Focused libraries, 270, 272
- Folding events, 368
- Force constants, 58
- Forward inclusion, 323
- Fractional coordinates, 10, 16
- Fragment marking, 271
- Fragment self-similarity, 173
- Fragments, 257
- Fukui function, 151, 163
- Fullerenes, 233, 244, 245, 263, 264, 266, 274
- Fullgen**, 263
- GalvaStructures**, 267
- Garnets, 4, 42
- Gas phase, 7
- Gaussian basis sets, 18
- Gaussian function exponents, 108
- GAUSSIAN70**, 2
- GAUSSIAN98**, 194
- Gaussian-type primitives, 144
- GEN**, 267
- GenBank*, 391
- Gene expression, 382
- Gene expression data, 389
- Gene interactions, 389, 390
- Gene network dynamics, 389
- Gene networks, 213, 217, 235
- Gene ontology (GO) consortium, 395
- Gene regulatory networks (GRNs), 381, 382, 383

- GeneCards*, 391, 393  
*GeneNet*, 394  
*GeneNote*, 392  
General position, 10  
Generalized Born (GB) models, 352  
Generalized gradient approximation (GGA), 3, 78  
**Generation of Isomers (GI)**, 267  
**GENeration with Overlapping Atoms (GENOA)**, 267  
Genetic algorithm (GA), 159, 250, 258, 260, 272, 339  
Genetic codes, 260  
Genetic function approximation, 340  
**GENM**, 268  
**GENMAS**, 268  
Genome organization, 384  
*GenomeNet*, 391, 393  
Genomes, 390  
**GENSTR**, 268  
*GEO*, 392  
Geometry optimization, 2  
Ghost atoms, 51  
Gilat's grid, 39  
Girona index, 166  
Global extremum, 160  
Global maximum, 159  
Global network properties, 402  
Gradient-corrected functionals, 52  
Graph, 210, 211, 388  
Graph theory, 209  
Graphite, 23  
**GRASP**, 360, 361  
Green's function, 359  
GRN analysis tools, 402  
GRN modeling ontologies, 395  
GRN representation, 383  
  
Hadamard product, 181  
Hamiltonian, 3, 35, 41, 109, 134  
Hammett  $\sigma$  descriptor, 130, 173  
Hartree-Fock (HF), 2, 3, 36, 47, 50, 59, 65  
Hasse diagrams, 273  
Hausmannite, 54  
Heat conduction, 63  
Heaviside step function, 20  
Helicenes, 228  
Helmholtz free energy, 65  
Hen egg-white lysozyme, 369  
2,2,3-trimethylhexane, 238, 239  
Hidden neurons, 338  
High throughput screening, 270  
  
Highest occupied molecular orbital (HOMO), 33  
Highly charged biomolecules, 356  
High-pressure phases, 59  
High-throughput data acquisition, 381  
High-throughput gene expression measurements, 389  
Hirshfeld promolecular electron density, 145, 168  
HIV-1 protease inhibitors, 274  
Hodgkins-Richards index, 165  
Hohenberg-Kohn theorems, 134  
Holographic electron density theorem, 177, 178, 194  
Homogeneous dielectric, 363, 368  
Homogeneous dielectric medium, 352  
Homogeneous polarizable medium, 351  
Homomorphisms, 248  
Hosoya index, 273  
**HOUDINI**, 269  
Hückel Hamiltonian, 21  
Hybrid functionals, 52  
Hydrogen-bond networks, 370  
Hydrogen-suppressed molecular graph, 213, 288  
Hyperfine couplings, 90  
Hyperfine isotropic coupling constants, 96  
Hyperstructure, 251  
  
**IBMOL**, 2  
Ice XI, 69  
Ideal surface, 66, 77  
Implicit solvent binding calculation, 368  
Implicit solvent methods, 350, 351, 371  
Impurity defect, 82  
Infinite crystal, 13  
Information, 289, 292, 295, 296  
Information theory, 341  
Informative design, 270  
Infrared spectra, 267  
Inhomogeneous dielectric, 363  
Inhomogeneous medium, 362  
Input neurons, 338  
Insulator, 34, 36, 38  
*IntAct project*, 394  
Integrated database analysis, 391, 393  
Integrated database retrieval, 391, 393  
Interaction databases, 396  
Interaction energy, 44  
Interactions, 398  
*InterDOM*, 394  
Interface energy, 44

- Interfaces, 66, 77  
 Internal pressure, 58  
*INSD (International Nucleotide Sequence Databases)*, 390  
 International Tables of Crystallography, 10, 14  
*InterPro*, 392  
 Intractable problems, 339  
 Inverse imaging, 249  
 Inverse-QSAR (I-QSAR), 272, 173  
 Inward matrix product (IMP), 181  
 Ion channels, 351  
 Ion size, 363  
 Ion solvation, 363  
 Ionic crystals, 44, 71  
 Ionic strength, 351  
 Ionization energy, 36  
*iProClass*, 392  
 Irreducible representations, 41  
 Ising model, 56  
**ISOGEN**, 267  
 Isomer enumeration, 247, 261  
 Isomer generation, 263  
 Isomer lists, 261  
 Isomers, 215, 221  
 Isomorphic graphs, 235  
 Isomorphism, 214  
 Isotopically labeled alkanes, 225  
  
 Jacobi rotation, 148  
**Jaguar**, 360, 361  
  
 Kappa-shape descriptors, 273  
 KBr, 47, 59  
*KEGG (Kyoto Encyclopedia of Genes and Genomes)*, 391, 394  
 Kekulé structures, 228  
 Ketones, 224, 263  
 Kinetic MQSM, 151  
 Kinks, 74, 80  
 k-Nearest Neighbor method, 340  
 Knowledgebases, 381, 383, 390  
 Koopmans theorem, 36  
 Korrington-Kohn-Rostoker (KKR), 2  
 Kruskal trees, 141  
 Kurtosis, 308  
  
 Labeled enumeration, 236, 262  
 Labeled graph, 212, 213, 215, 233, 256, 257  
*Lac operon*, 384  
 Lagrange multiplier, 146  
 Lamarckian genetic algorithm, 159, 160  
 Large-scale PB calculations, 371  
  
 Latent variable regression analysis, 314, 316  
 Latent variables, 314, 315  
 Lattice, 7  
 Lattice deformation, 63  
 Lattice dynamics, 63  
 Lattice energy, 44  
 Lattice enumeration, 241, 242  
 Lattice origin, 24  
 Lattice parameters, 7, 12  
 Lattice vector, 8  
 Lattice vibration frequencies, 61, 64  
 Layer groups, 11  
 Lead compound, 270  
 Least squares regression, 314  
 Leave-multiple-out (LMO) CV, 329  
 Leave-one-out (LOO) CV, 329  
 Lexicographic coloring, 234  
 Library design, 270  
 Library sizes, 222  
 LiF, 47, 59, 86, 88, 92  
**LightBench**, 395  
 Linear combination of atomic orbitals (LCAO), 17, 15  
 Linear combinations of descriptor variables, 314  
 Linear defects, 80  
 Linear genetic codes, 261  
 Linearized augmented plane waves (LAPW), 2  
 Linearized PB equation, 355, 359  
 Lipinski rule-of-five, 130, 171  
 Local basis set, 39, 42, 67  
 Local density approximation (LDA), 1, 35, 47, 50, 52, 59, 65, 78  
 Local extremum, 160  
 Local minima, 324  
 Local minimization, 148  
 Local softness, 151, 152  
 Locally restricted graphs, 233, 237  
 logP, 130, 171, 290  
 Lone pairs, 150  
 Long-range forces, 349  
 Low charge-density ions, 363  
  
**MacroDox**, 360, 361  
 Magnesium oxide (MgO), 34, 38, 47, 59, 74, 75  
 Magnetic coupling constants, 96  
 Magnetic properties, 36, 54  
 Mallows Cp, 327  
 Manhattan distance, 135  
 Maximal connectivity stack, 214  
**MEAD**, 360, 361  
 Mean centered data matrix, 298

- Mean centering of data, 298
- Mean square PRESS, 328
- Measurement error, 292
- Membranes, 351
- Messenger RNA (mRNA), 382, 383
- Metabolic signaling database, 393
- Metabolites, 383, 384
- Metabolome, 383
- Metallic crystals, 44
- Metals, 70
- Miller indices, 67
- Minerals, 59
- Minkowski distance, 135
- MINT (Molecular Interactions)*, 394
- Mobile charge distribution, 355
- Mobile counterion distribution, 356
- Modified proteins, 383
- Molecular
  - alignment, 144, 155, 161, 164
  - binding energies, 44
  - biology, 381
  - cages, 233, 244
  - complexes, 399
  - crystals, 49
  - descriptors, 29, 132, 152, 249, 252, 290
  - design, 210, 249, 272
  - dynamics, 368
  - electrostatic potential, 161
  - enumeration, 210
  - flexibility, 157
  - fragments, 215, 249
  - function, 396
  - graph, 212, 237, 247, 249
  - graph enumeration, 249
  - mechanics/PBSA, 368
  - quantum chemistry, 1
  - quantum self-similarity measure (MQSSM), 135, 139, 154, 173
  - quantum similarity, 127, 134, 143, 167
  - quantum similarity measure (MQSM), 135, 148
  - shape, 289
  - similarity, 128
  - superposition, 144, 150, 178
  - surface, 355, 366
  - tree, 215
- Molecule counting, 220
- Molecular connectivity indices, 288
- MOLGEN**, 245, 248, 253, 255, 257, 268
- MOLGEN-COMB**, 271
- MOLGRAPH**, 268
- Møller-Plesset (MP), 5
- MolSurfer**, 360
- Momentum space, 152
- Momentum space density, 152
- Momentum transformation, 152
- Monte Carlo (MC) techniques, 258, 368
- Monte Carlo sampling, 258
- Most fit member, 160
- Mouse acetylcholinesterase (mAChE), 367
- MP2, 77
- MPP-CRYSTAL**, 110
- Muffin tin potentials, 2
- Mulliken charge, 89, 95, 100
- Multicollinearity, 289, 296, 299, 304, 311, 334, 336
- Multigraph, 211
- Multigrid methods, 359
- Multiple edges, 211
- Multiple linear regression (MLR), 310
- Multiscale modeling ontology, 397
- Multivariate dataset, 290
- Multivariate statistics, 291
- Multivariate techniques, 290
- Mutation, 160, 260, 261
- MUSEUM (Mutation and SElection Uncover Models)*, 340
- Nanotubes, 233, 244
- Naphthalene, 221
- Nauty**, 214, 237
- NCBI**, 393
- Neighboring defects, 85
- Network graph, 389, 402
- Network stability analysis, 404
- Network topology, 387
- Neural net, 272
- Neutral defect, 84, 90, 103
- Neutron elastic scattering, 64
- Newton-Raphson, 148
- Nickel oxide, 36
- NIPALS**, 318
- NMR, 368
- NMR spectra, 215, 250, 260, 261
- Noise, 292, 294
- Nondrug-like molecule, 130
- Nonlinear modeling, 338
- Nonnegative garrote (NNG), 310
- Nonoverlapping fragments, 250
- Nonpolar interactions, 352
- Nonprimitive lattices, 8, 39
- Normal modes, 64
- n-tuple code, 258
- Numerical instabilities, 39, 45

- Object coloring, 219, 222  
Object set, 180  
Octanol/water partition coefficient (logP), 171  
Offspring, 160, 260  
Omission bias, 318  
Omission of control variables, 320  
Omission of observation variables, 320  
One-electron integrals, 108  
ONIOM, 82  
Ontology, 383, 395  
Operators, 137  
Operon structure, 384  
Optical modes, 64  
Optimization, 148, 159  
Optimization methods, 258  
Orbital localization, 6  
Orderly generation, 235, 236, 237, 241, 245  
Ordinary least squares (OLS), 310, 311  
Organic molecules, 215  
Origin cell, 8, 10  
Orthogonalized plane waves (OPW), 2  
Osmotic pressure, 356  
Outer electron density, 150  
Overfitting, 338  
Overinflated significance values, 324  
Overlapping fragments, 250  
  
Pack-Monkhorst grid, 39  
Pair correlations, 341  
Paragon, 334  
Parent, 260  
Partial least squares (PLS), 310, 314, 315  
Path connectivity, 289  
Path length distribution, 403  
*PathDB*, 394  
*PathPort/Toolbus*, 393  
Pathway, 400  
**Pathway Assist**, 394  
Pathway channels, 404  
Pathway database models, 401  
Pathway Database System (PDS), 397  
Pathway graph, 400  
Pathway ontology, 397, 400  
Pathways databases, 390, 393, 394  
**PATIKA**, 394, 397  
**PCRYSTAL**, 110  
*PDB*, 391  
Peptidomimetic, 271  
Perception of similarity, 129  
Perfect crystal, 6, 81, 82, 85  
Perfect similarity, 136  
Pericondensed benzenoids, 229  
Perifusenes, 230  
Periodic boundary conditions, 12, 63  
Periodic oscillations, 406  
Periodic replica, 83  
Periodic systems, 6  
Periodicity, 82  
Permutation, 214  
Permutation group, 227, 231  
Peroskivites, 54  
Perturbation technique, 259  
Phase diagram, 59  
Phase stability, 59  
Phase transition, 59, 60  
Phenylene, 229  
Phenanthrene, 221, 228  
Phonon density of states, 65  
Phonon dispersion, 64  
Phonon frequencies, 64  
Phonon spectra, 65  
Physicochemical descriptors, 130, 287, 291  
Physisorption, 74  
Piezoelectricity, 5  
*PIR (Protein Information Resource)*, 391  
 $pK_a$  calculations, 369, 370  
Planar defects, 80  
Plane wave (PW) basis sets, 2, 67  
Plane waves, 16  
Point defects, 80, 85  
Point symmetry, 4, 41, 84  
Poisson-Boltzmann equation (PBE), 89, 94, 350, 352, 354, 357  
Poisson's equation, 354  
Polar interactions, 351  
Polar solvation forces, 366  
Polar surfaces, 71  
Polarization, 149  
Polarization functions, 46, 63  
Polarization of energy, 356  
Pólya counting theory, 218  
Pólya's theorem, 219, 231, 271  
Polychlorinated biphenyls, 289  
Polycyclic aromatic hydrocarbons (PAHs), 229  
Polyhex enumeration, 241  
Polyhex hydrocarbons, 228, 241, 261  
Polymer design, 261  
Polymorphism, 52  
Population, 160  
Population analysis, 144  
Position isomorphous phases, 56  
Position space electron density, 152

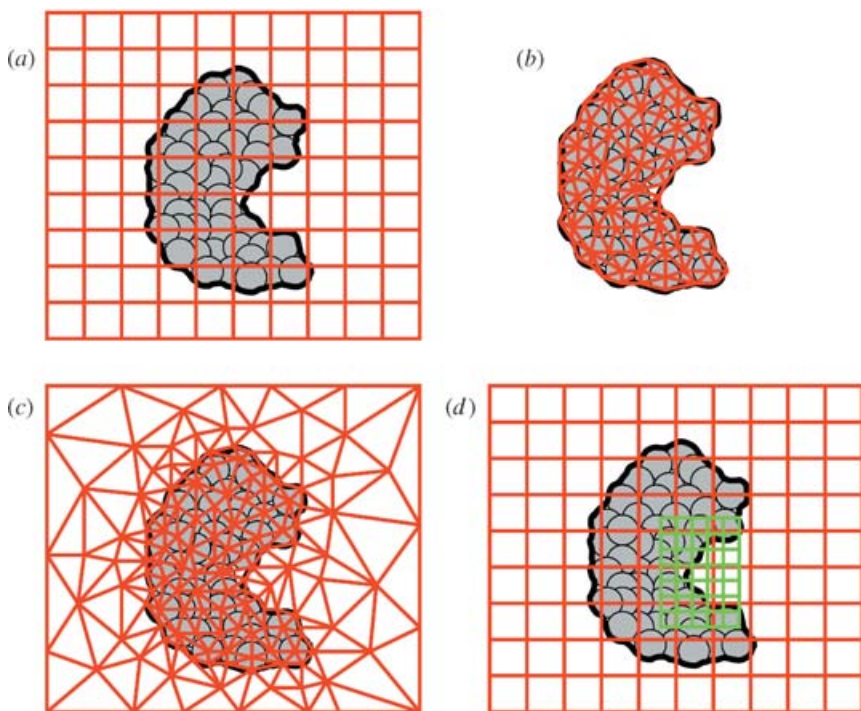
- Posttranslational covalent modifications, 382  
Postulates of quantum mechanics, 132  
Potential of mean force (PMF), 363, 364  
Predicted residual error of sum of squares (PRESS), 327  
Pressure, 59  
Primary alcohols, 224, 264  
Primitive lattices, 8, 40  
Principal component analysis (PCA), 289, 291  
Principal component regression (PCR), 310, 314  
Principal components (PC), 292  
**PRO\_SELECT**, 272  
Probabilistic graphical models, 389  
*ProDom (Protein Domain)*, 392  
Promolecular atomic shell approximation (PASA), 145, 148, 159, 194  
Promolecular electron density, 144, 145, 148  
Promoter regions, 387  
Protein domain motion, 368  
Protein systems, 356  
Protein-DNA interactions, 384  
Protein-protein interaction networks, 388  
Protein-protein interactions, 384  
Proteins, 382  
Proteomics, 371  
Pruning, 338  
Pseudopotential (PP), 2, 151  
*PubMed*, 395  
PW91, 47, 50, 59  
Pyrene, 229  
Pyrope, 42  
  
Quantitative structure-activity relationships (QSAR) 129, 171, 173, 191, 252, 257, 272, 287, 289  
Quantitative structure-property relationships (QSPR) 129, 171, 260  
Quantum chemistry, 128  
Quantum mechanical simulation, 1  
Quantum Monte Carlo methods, 6  
Quantum object, 141, 142, 180  
Quantum QSAR, 141, 180, 191, 194  
Quantum similarity indices, 164  
Quantum similarity matrix, 190  
Quantum similarity maximization (MaxiSim), 157, 163, 178  
Quantum similarity superposition algorithm (QSSA), 157, 159, 163, 178, 191  
Quantum topological molecular similarity, 169  
 $\alpha$ -quartz, 52, 65  
**QuaSAR-CombiGen**, 272  
  
Quasi-Newton-Raphson, 160  
Quasi-periodic structures, 77, 82  
Quaternions, 157, 159  
  
Random error, 310  
Random networks, 403  
Random sampling, 257  
Random search, 157  
Random structures, 257  
Rate parameters, 383  
Reaction field treatments, 352  
Reaction network analysis, 383  
Reaction networks, 213, 217, 405  
Reaction transform, 272  
Real crystals, 34  
Real surfaces, 74  
Reciprocal lattice, 11, 24  
Reciprocal space, 11, 39, 40  
Redundancy free isomers, 268  
Redundant variables, 296, 299  
Regression coefficient, 299  
Regression modeling, 309  
Regular graph, 244, 248  
3-Regular graphs, 245  
3-Regular spherical map, 245  
Regulatory events, 382  
Regulatory molecular networks, 382  
Regulatory sequences, 387  
Relative stability, 44  
Repressor protein, 384  
Ribosomal RNA (rRNA), 383  
Ribosomes, 382  
Ridge regression (RR), 310, 311  
Ridge traces, 313  
Rigid bodies, 157  
Rigid entities, 367  
Rigid-body binding energy, 368  
RNA, 356, 382, 389  
Rod groups, 11  
Rooted catafusene, 231  
Rooted labeled graphs, 217  
Rooted tree, 215, 222, 225, 256  
Roothaan equations, 18  
Rutile, 67  
  
Saddle-node bifurcation, 387  
Sampling molecules, 257  
Sampling structures, 255  
Scaffold, 213, 222, 235, 270  
Scaffold of electron density, 134  
Scaled-particle theories (SPTs), 353  
Scale-free networks, 403

- Schönflies notation, 10, 19  
 Schrödinger equation, 12, 16, 39, 134  
*SCOP (Structural Classification of Proteins)*, 392  
 Scores plot, 294  
 Scoring function, 162  
 Secondary alcohols, 224, 264  
 Second-order electron density, 133  
 Second-order saddle point, 50  
 Selection, 260  
 Selection bias, 320, 322, 324, 336  
 Self-energy terms, 356  
 Self-interaction terms, 362  
 Self-similarity, 171, 174  
 Self-similarity measures, 135, 139, 165  
 Selwood data set, 339  
 Semiconductor, 33, 38, 70  
 Semigroup, 183  
**SENECA**, 260, 268  
 Sequential agglomerative hierarchical nonoverlapping (SAHN) dendrogram, 141  
**SESAMI**, 269  
 Shape functions, 167  
 Shell model, 82  
 Short-range forces, 349  
 Shrinkage regression technique, 312  
 Side chain reorganization, 368  
**SIGNATURE**, 268, 269, 274  
 Signature equation, 250  
 Significance level, 322  
**SigTrans**, 395  
 Silica, 52  
 Silicon, 34, 38, 47, 59, 71  
 Similarity, 127, 129, 295  
 Similarity integral, 135  
 Similarity matrix, 140  
 Similarity measure, 273, 295  
 Simple graph, 211, 258  
 Simple linear regression, 310  
 Simplex method, 160  
 Simulated annealing (SA), 258  
 Single response regression method, 310  
 Singular-value decomposition (SVD), 305, 342  
 Size-extensivity, 85  
**SKEL\_GEN**, 266  
 Skewness, 308  
 Slab model, 66  
 Slow-moving electrons, 152  
**SMARTS**, 272  
*SMD (Stanford Microarray Database)*, 392  
**SMOG**, 248  
 Sodalite, 52  
 Sodium Chloride (NaCl), 45, 47, 59  
 Software, 104  
 Solid state chemistry, 1  
 Solid state physics, 2, 20  
 Solid state reaction energy, 44  
 Solids, 1, 45  
 Solute dielectric, 363  
 Solvation free energies, 362  
 Solvent dielectric, 363  
 Solvent electrorestriction, 363  
 Solvent polarization, 351  
 Solvent-accessible surface area (SASA), 353  
 Solvent-mediated interactions, 352  
 Space group, 10, 41  
 Special position, 10  
**SpecSolv**, 250, 268  
 Spherical cellular schemes, 2  
 Spin, 133  
 Spin contamination, 56  
 Spin density, 87, 89  
 Spin moments, 95  
 Spin polarization, 55  
 Spiral canonical code, 246  
 SPZ, 89, 94  
**SRS-EBI**, 393  
 States, 398  
 Static image database, 393  
 Static point charges, 350  
 Statistical analysis, 389  
 Steps, 74, 80  
 Stepwise regression, 324, 330  
 Stereoalkanes, 261  
 Stereoisomer enumeration, 253  
 Stewart atoms, 149  
 Stewart charges, 149  
*STKE (Signal Transduction Knowledge Environment)*, 394  
 Stochastic dynamics simulations, 371  
 Stochastic matrix, 140, 190  
 Stochastic sampling, 272  
 Stockholder charges, 168  
 Stopping rules for variable inclusion, 327  
 Strain tensor, 62  
 Stress, 5  
**StrucEluc**, 268  
 Structural alignment, 161  
 Structural genomics, 371  
 Structural isomers, 142, 221, 227, 247  
 Structural n-tuple, 238  
 Structure counting, 215  
 Structure elucidation, 210, 257, 266, 269  
 Structure elucidation successes, 269

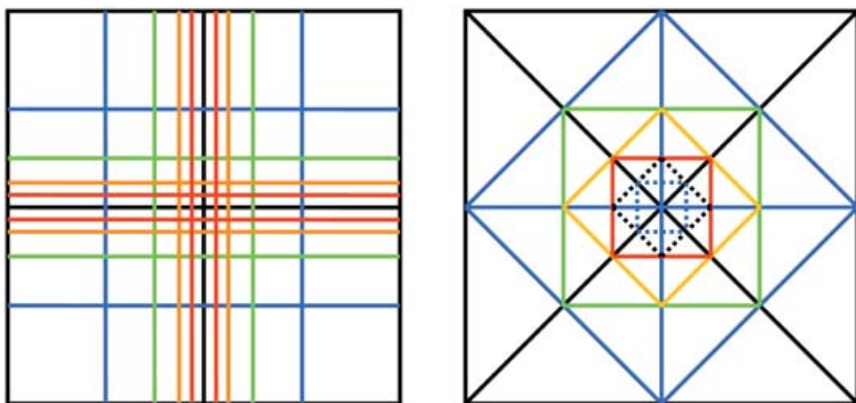
- Structure-activity relationship, 129  
Sublimation energies, 49  
Subnetwork, 388  
Substituent constants, 290, 294  
Substituent effects, 173  
Substitution energy, 44  
Substitutional defects, 80, 84  
Substructure search systems, 130, 162  
Substructures, 238, 247, 252  
Subvalence electrons, 151  
Subvalence regions, 148  
Supercells, 54, 83  
Super-exchange coupling constants, 56  
Super-exchange energy, 44  
Super-exchange interactions, 55  
Superposed atoms, 162  
Superposition, 150  
Supervised elimination, 309  
Supervised forward selection (SFS), 335  
Supervised variable selection, 307, 334  
Surface defects, 81  
Surface formation energy, 44, 70  
Surface reactivity, 80  
Surface reconstruction, 68, 71  
Surface relaxation, 71  
Surface stability energy, 44  
Surface termination, 68  
Surfaces, 66, 74  
Survival of the fittest, 260  
*Swiss-Prot*, 391  
Symmetry, 40  
Symmetry adapted basis sets, 41  
Symmetry operations, 218  
Systematic search, 157  
  
Tabu search (TS), 324  
Tagged sets, 136, 180  
Tanimoto index, 165  
Targeted libraries, 261  
Tartaric acid, 253  
Taylor expansion, 358  
Tensors, 5  
Tertiary alcohols, 224, 264  
Test set, 273, 328  
Tetracene, 221  
Thermodynamic cycle, 367  
Third-order electron density, 153  
Time-dependent density functional theory (TDDFT), 6, 36  
Time-reversal symmetry, 41  
Titratable residues, 358  
Titration calculations, 368  
  
Tolerance, 304  
Topo-geometrical superposition approach (TGSA), 162, 191  
Topographical indices, 130  
Topological descriptors, 252, 288  
Topological indices, 130, 274  
Toroidal polyhex, 247  
Total energy, 37, 43  
Training set, 273, 296, 328  
Transcription, 382  
Transcriptional regulatory networks, 403  
Transfer RNA (tRNA), 383  
Transitions, 398  
Translation, 382  
Translation invariance, 4, 38  
*TRANS-PATH*, 394  
Trapped hole centers, 86, 90  
Tree, 210, 215, 218, 250  
Tree graphs, 141  
 $\alpha$ -tridymite, 52  
Triphenylene, 221  
True model, 309, 318  
Two-electron integrals, 108  
  
**UHBD**, 360, 361  
Ultrathin oxide films, 78  
Uncertainty, 383  
Unfolding events, 368  
*UniProt knowledgebase*, 391  
Unit cell, 8, 39  
Unlabeled graphs, 214, 215, 217, 233, 256  
Unlabeled trees, 225  
Unnecessary variables, 320  
Unrestricted Hartree-Fock (UHF), 55, 88, 89, 94  
Unrooted trees, 256  
Unsupervised elimination, 307  
Unsupervised forward selection (UFS), 308, 335  
Unsupervised learning, 291  
Unsupervised variable selection, 307, 334  
Urea, 49  
  
Vacancy, 80  
Valence bands, 33  
Valence density, 150, 161  
Validation set, 273, 296, 328  
Variable elimination, 291, 293, 295, 339  
Variable selection, 287, 291, 309, 318, 339  
Variable standardization, 298  
Variable subset selection (VSS), 318



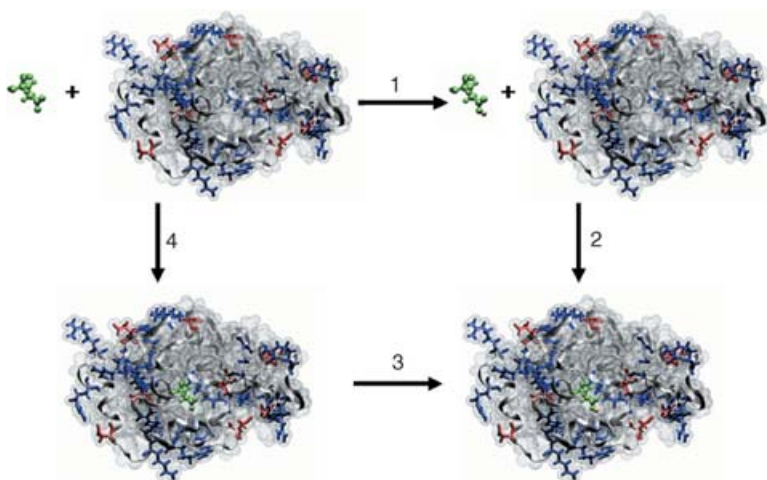
- Variable weighting, 341
- Variables, 290
- Variance, 316
- Variance inflation factor, 304
- Vector semispaces, 136, 167, 180, 186
- Vertex degree, 211, 212
- Vertices, 80, 211
- Vibration spectra, 63
- Vibration frequencies, 64
- Virtual chemistry, 270
- Virtual combinatorial library, 271
- VirtualCell**, 396
- Wave vector, 12
- Wavefunction, 2, 132, 152
- Wavelet transforms, 154
- Web-based services, 360
- Whole-cell modeling, 396
- Wiener index, 260
- X- $\alpha$ , 2
- X-ray structural data, 368
- Zeolites, 4, 40
- Zero point energy, 48
- Zero-gap semiconductor, 31



**Figure 2** Popular discretization schemes for numerical solution of the Poisson-Boltzmann equation. The solid black line and circles denote a model protein; other lines denote the mesh on which the system is discretized. (a) Finite difference. (b) Boundary element. (c) Finite Element. (d) Focusing on finite difference grids.



**Figure 3** Adaptive refinement for finite different (left) and finite element (right) methods. Shading denotes successive levels of refinement.



**Figure 9** Thermodynamic cycle illustrating the standard numerical procedure for calculating the protonation energy of the active site GLU 35 in lysozyme (CPK representation). Acidic and basic residues are shown in “licorice” representation. The steps are (1) protonation of the residue in isolation, (2) transfer of the protonated residue from isolation into the biomolecule, (3) protonation of the residue in the biomolecule, and (4) transfer of the unprotonated residue from isolation into the biomolecule.