

METHODS IN MOLECULAR BIOLOGY™

Series Editor
John M. Walker
School of Life Sciences
University of Hertfordshire
Hatfield, Hertfordshire, AL10 9AB, UK

For further volumes:
<http://www.springer.com/series/7651>

Evolutionary Genomics

Statistical and Computational Methods, Volume 1

Edited by

Maria Anisimova

*Department of Computer Science, Swiss Federal Institute of Technology (ETHZ),
Zürich, Switzerland*

Swiss Institute of Bioinformatics, Lausanne, Switzerland

 **Humana Press**

Editor

Maria Anisimova, Ph.D.
Department of Computer Science
Swiss Federal Institute of Technology (ETHZ)
Zürich, Switzerland

Swiss Institute of Bioinformatics
Lausanne, Switzerland

The photo used for book cover is made by one of the authors of the book, Wojciech Makalowski.

ISSN 1064-3745 e-ISSN 1940-6029
ISBN 978-1-61779-581-7 e-ISBN 978-1-61779-582-4
DOI 10.1007/978-1-61779-582-4
Springer New York Dordrecht Heidelberg London

Library of Congress Control Number: 2012931926

© Springer Science+Business Media, LLC 2012

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Humana Press, c/o Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper

Humana Press is part of Springer Science+Business Media (www.springer.com)

Preface

Discovery of genetic material propelled the power of classical evolutionary studies across the diversity of living organisms. Together with early theoretical work in population genetics, the debate on sources of genetic makeup initiated by proponents of the neutral theory made a solid contribution to the spectacular growth in statistical methodologies for molecular evolution. The methodology developed focused primarily on inferences from single genes or noncoding DNA segments: mainly on reconstructing the evolutionary relationships between lineages and estimating evolutionary and selective forces. Books offering a comprehensive coverage of such methodologies have already appeared, with Joe Felsenstein's "Inferring Phylogenies" and Ziheng Yang's "Computational Molecular Evolution" among the favorites.

This volume is intended to review more recent developments in the statistical methodology and the challenges that followed as a result of rapidly improving sequencing technologies. While the first sequenced genome (RNA virus Bacteriophage MS2 in 1976) was not even 4,000 nucleotides long, the sequencing progress culminated with the completion of the human genome of about 3.3×10^9 base pairs and advanced to sequence many other species genomes, heading ambitiously towards population sequencing projects such as 1,000 genome projects for humans and *Drosophila melanogaster*. Next-generation sequencing (NGS) technologies sparked the "genomics revolution," which triggered a renewed effort towards the development of statistical and computational methods capable of coping with the flood of genomic data and its inherent complexity.

The challenge of analyzing and understanding the dynamics of large-system data can be met only through an integration of organismal, molecular, and mathematical disciplines. This requires commitment to an interdisciplinary approach to science, where both experimental and theoretical scientists from a variety of fields understand each other's needs and join forces. Evidently, there remains a gap to be breached. This book presents works by top scientists from a variety of disciplines, each of whom embodies the interdisciplinary spirit of evolutionary genomics. The collection includes a wide spectrum of articles—encompassing theoretical works and hands-on tutorials, as well as many reviews with much biological insight.

The evolutionary approach is clearly gaining ground in genomic studies, for it enables inferences about patterns and mechanisms of genetic change. Thus, the theme of *evolution* streams through each chapter of the book, providing statistical models with basic assumptions and illustrated with appealing biological examples. This book is intended for a wide scientific audience interested in a compressed overview of the cutting-edge statistical methodology in evolutionary genomics. Equally, this book may serve as a comprehensive guide for graduate or advanced undergraduate students specializing in the fields of genomics or bioinformatics. The presentation of the material in this volume is aimed to equally suit both a novice in biology with strong statistics and computational skills and a molecular biologist with a good grasp of standard mathematical concepts. To cater for differences in reader backgrounds, *Part I* of *Volume I* is composed of educational primers to help with fundamental concepts in genome biology (Chapters 1 and 2), probability and statistics (Chapter 3), and molecular evolution (Chapter 4). As these concepts reappear repeatedly throughout the books, the first four chapters will help the neophyte to stay "afloat."

The exercises and questions offered at the end of each chapter serve to deepen the understanding of the material. Additional materials and some solutions to exercises can be found online: <http://www.evolutionarygenomics.net>.

Part II of this volume reviews state-of-the-art techniques for genome assembly (Chapter 5), gene finding (Chapter 6), sequence alignment (Chapters 7 and 8), and inference of orthology, paralogy (Chapter 9), and laterally transferred genes (Chapter 10). *Part III* opens with a comparative review of genome evolution in different breeding systems (Chapter 11) and then discusses genome evolution in model organisms based on the studies of transposable elements (Chapters 12 and 13), gene families, synteny (Chapter 14), and gene order (Chapters 15 and 16).

Part I of *Volume 2* is the evidence that, since embracing Darwin's tree-like representation of evolution and pondering over the universal Tree of Life, the field has moved on. Nowadays, the evolutionary biologists are well aware of numerous evolutionary processes that distort the tree, complicating the statistical description of models and increasing computational complexity, often to prohibitive levels. Each taking a different angle, the chapters of *Part I, Volume 2* discuss how to overcome problems with phylogenetic discordance, as the Tree of Life turns out to be more like a "forest" (Chapter 3). The multispecies coalescent model offers one solution to reconciling phylogenetic discord between gene and species trees (Chapter 1); others pursue probabilistic reconciliation for gene families based on a birth–death model along a species phylogeny (Chapter 2). By some perspectives, constraining the understanding of evolution solely with tree-like structures omits many important biological processes that are not tree-like (Chapter 4).

Most fundamental questions in genome biology strive to disentangle the evolutionary forces shaping species genomes, inferring evolutionary history, and understanding how molecular changes affect genomic and phenotypic characteristics. To this goal, *Part II* of the *Volume 2* introduces methods for detecting and reconciling selection (Chapters 5 and 6) and recombination (Chapters 9 and 10), and discusses the mechanisms for the origins of new genes (Chapter 7) and the evolution of protein domain architectures (Chapter 8). The role of *natural selection* in shaping genomes is a pinnacle of the classical neutralist–selectionist debate and sets an important theme of the book; the "neo-selectionist" model of genome evolution is tested on many counts. This theme is also apparent in *Part III* dedicated to *population genomics*, which starts by discussing models for genetic architectures of complex disease and the power of genome-wide association studies (GWAS) for finding susceptibility variants (Chapter 11). With the availability of multiple genomes from closely related species, gleaning the ancestral population history also became possible, as is illustrated in the following chapter (Chapter 12). Most population genetics problems rely on ancestral recombination graphs (ARG), and reducing the redundancy of the ARG structure helps to reduce the computational complexity (Chapter 13).

Entering the era of postgenomics biology, recent years have seen rapid growth of complementary genomic data, such as data on expression and regulation, chemical and metabolic pathways, gene interactions and networks, disease associations, and more. Considering the genome as a uniform collection of coding and noncoding molecular sequences is no longer an option. To address this, great efforts are currently dedicated to embrace the complexity of biological systems through the emerging "-omics" disciplines—the focus of *Part IV* of this volume. Chapter 14 discusses ways to study the evolution of gene expression and regulation based on data from "old-fashioned" microarrays as well as transcriptomics data obtained with NGS such as RNAseq and ChIPseq. Interactomics is the focus of the next chapter. Indeed, better understanding of genes, their diversity and

regulation comes from studies of interaction between their protein products and networks of interacting elements (Chapter 15). Further topics include metabolomics (Chapter 16), metagenomics (Chapter 17), epigenomics (Chapter 18), and the newly reinvented discipline with a mysterious name—genetical genetics (Chapter 19). Despite the effort, complex dependencies and causative effects are difficult to infer. A way forward must be in the integration of complimentary “-omics” information with genomic sequence data to understand the fundamentals of systems biology in living organisms. This cannot be achieved without studying how such information changes over time and across various conditions. Vast amount of multifaceted data promise a big future for machine learning, pattern recognition and discovery, and efficient data mining techniques, as can be seen from many chapters of this book.

Finally, *Part V* of the second volume focuses on challenges and approaches for large and complex data representation and storage (Chapter 20). The rapid pace of computational genomics, as well as research transparency and efficiency, exacerbates the need for sharing of data and programming resources. Fortunately, some solutions already exist (Chapter 21). Handling ever increasing amounts of computation requires efficient computing strategies, which are discussed in the closing chapter of the book (Chapter 22).

For a novice in the field, this book is certainly a treasure chest of state-of-the-art methods to study genomic and omics data. I hope that this collection will motivate both young and experienced readers to join the interdisciplinary field of evolutionary genomics. But even the experienced bioinformatician reader is certain to find a few surprises. On behalf of all authors, I hope that this book will become a source of inspiration and new ideas for our readers. Wishing you a pleasant reading!

Zürich, Switzerland

Maria Anisimova, Ph.D.

Acknowledgments

The foremost gratitude goes to the authors of this book who came together to make this resource possible and who were enthusiastic and encouraging about the whole project. Over 100 reviewers have contributed to improving the quality and the clarity of the presentation with their constructive and detailed comments. Some reviewers have accepted to be acknowledged by their name. With great pleasure, I list them here:

Tyler Alioto, Peter Andolfatto, Miguel Andrade, Irena Artamonova, Richard M. Badge, David Balding, Mark Beaumont, Chris Beecher, Robert Beiko, Adam Boyko, Katarzyna Bryc, Kevin Bullaughey, Margarida Cardoso-Moreira, Julian Catchen, Annie Chateau, Karen Cranston, Karen Crow, Tal Dagan, Dirk-Jan de Koning, Christophe Dessimoz, Mario dos Reis, Katherine Dunn, Julien Y. Dutheil, Toni Gabaldon, Nicolas Galtier, Mikhail Gelfand, Josefa Gonzalez, Maja Greminger, Stephane Guindon, Michael Hackenberg, Carolin Kosiol, Mary Kuhner, Anne Kupczok, Nicolas Lartillot, Adam Leache, Gerton Lunter, Thomas Mailund, William H. Majoros, James McInerney, Gabriel Musso, Pjotr Prins, David A. Ray, Igor Rogozin, Mikkel H. Schierup, Adrian Schneider, Daniel Schoen, Cathal Seoighe, Erik Sonnhammer, Andrea Splendiani, Tanja Stadler, Manuel Stark, Krister Swenson, Adam M. Szalkowski, Gergely J. Szöllösi, Jijun Tang, Todd Treangen, Oswaldo R. Trelles Salazar, Albert Vilella, Rutger Vos, Tom Williams, Carsten Wiuf, Yuri Wolf, Xuhua Xia, S. Stanley Young, Olga Zhaxybayeva, and Stefan Zoller.

My colleagues from the Computational Biochemistry Research Group at ETH Zurich deserve much credit for being a constant source of inspiration and for providing such an enjoyable working environment. Finally, but no less importantly, I would like to thank my family for their love and for tolerating the overtime that this project required.

Contents

<i>Preface</i>	v
<i>Contributors</i>	xiii
PART I INTRODUCTION: BIOINFORMATICIAN'S PRIMERS	
1 Introduction to Genome Biology: Features, Processes, and Structures	3
<i>Aidan Budd</i>	
2 Diversity of Genome Organisation	51
<i>Aidan Budd</i>	
3 Probability, Statistics, and Computational Science	77
<i>Niko Beerenwinkel and Juliane Siebourg</i>	
4 The Essentials of Computational Molecular Evolution	111
<i>Stéphane Aris-Brosou and Nicolas Rodrigue</i>	
PART II GENOMIC DATA ASSEMBLY, ALIGNMENT, AND HOMOLOGY INFERENCE	
5 Next-Generation Sequencing Technologies and Fragment Assembly Algorithms	155
<i>Heewook Lee and Haixu Tang</i>	
6 Gene Prediction	175
<i>Tyler Alioto</i>	
7 Alignment Methods: Strategies, Challenges, Benchmarking, and Comparative Overview	203
<i>Ari Löytynoja</i>	
8 Whole-Genome Alignment	237
<i>Colin N. Dewey</i>	
9 Inferring Orthology and Paralogy	259
<i>Adrian M. Altenhoff and Christophe Dessimoz</i>	
10 Detecting Laterally Transferred Genes	281
<i>Rajeev K. Azad and Jeffrey G. Lawrence</i>	
PART III GENOME EVOLUTION: INSIGHTS FROM STATISTICAL ANALYSES	
11 Genome Evolution in Outcrossing Versus Selfing Versus Asexual Species	311
<i>Sylvain Glémin and Nicolas Galtier</i>	
12 Transposable Elements and Their Identification	337
<i>Wojciech Makalowski, Amit Pande, Valer Gotea, and Izabela Makalowska</i>	

13	Evolution of Genome Content: Population Dynamics of Transposable Elements in Flies and Humans	361
	<i>Josefa González and Dmitri A. Petrov</i>	
14	Detection and Phylogenetic Assessment of Conserved Synteny Derived from Whole Genome Duplications	385
	<i>Shigehiro Kuraku and Axel Meyer</i>	
15	Analysis of Gene Order Evolution Beyond Single-Copy Genes	397
	<i>Nadia El-Mabrouk and David Sankoff</i>	
16	Discovering Patterns in Gene Order	431
	<i>Laxmi Parida and Niina Haiminen</i>	
	<i>Index</i>	457

Contributors

- TYLER ALIOTO • *Centro Nacional de Análisis Genómico, Barcelona, Spain*
- ADRIAN M. ALTENHOFF • *Department of Computer Science, ETH Zurich, Zurich, Switzerland; Swiss Institute of Bioinformatics, Switzerland*
- STÉPHANE ARIS-BROUSO • *Departments of Biology and Mathematics & Statistics and Center for Advanced Research in Environmental Genomics, University of Ottawa, Ottawa, ON, Canada*
- RAJEEV K. AZAD • *Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA, USA; Departments of Biological Sciences and Mathematics, University of North Texas, Denton, TX, USA*
- NIKO BEERENWINKEL • *Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland*
- AIDAN BUDD • *European Molecular Biology Laboratory (EMBL), Heidelberg, Germany*
- CHRISTOPHE DESSIMOZ • *Department of Computer Science, ETH Zurich, Zurich, Switzerland; Swiss Institute of Bioinformatics, Switzerland*
- COLIN N. DEWEY • *Biostatistics and Medical Informatics and Computer Sciences, Genome Center of Wisconsin, University of Wisconsin-Madison, Madison, WI, USA*
- NADIA EL-MABROUK • *Département d'Informatique et de Recherche Opérationnelle, Université de Montréal, Montréal, QC, Canada*
- NICOLAS GALTIER • *Institut des Sciences de l'Evolution, UMR5554, Université Montpellier II, Montpellier, France*
- SYLVAIN GLÉMIN • *Institut des Sciences de l'Evolution, UMR5554, Université Montpellier II, Montpellier, France*
- JOSEFA GONZÁLEZ • *Department of Biology, Stanford University, Stanford, CA, USA; Institute of Evolutionary Biology (CSIC-UPF), Barcelona, Spain*
- VALER GOTEA • *National Human Genome Research Institute, National Institutes of Health, Rockville, MD, USA*
- NIINA HAIMINEN • *IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA*
- SHIGEHIRO KURAKU • *Genome Resource and Analysis Unit, RIKEN Center for Developmental Biology, Chuo-ku, Kobe, Japan*
- JEFFREY G. LAWRENCE • *Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA, USA*
- HEEWOOK LEE • *School of Informatics and Computing, Indiana University, Bloomington, IN, USA*
- ARI LÖYTYNOJA • *European Bioinformatics Institute (EMBL), Hinxton, UK; Institute of Biotechnology, University of Helsinki, Helsinki, Finland*
- IZABELA MAKALOWSKA • *Laboratory of Bioinformatics, Adam Mickiewicz University, Poznań, Poland*

WOJCIECH MAKALOWSKI • *Institute of Bioinformatics, University of Muenster, Muenster, Germany*

AXEL MEYER • *Department of Biology, University of Konstanz, Constance, Germany*

AMIT PANDE • *Institute of Bioinformatics, University of Muenster, Muenster, Germany*

LAXMI PARIDA • *IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA*

DMITRI A. PETROV • *Department of Biology, Stanford University, Stanford, CA, USA*

NICOLAS RODRIGUE • *Department of Biology and Center for Advanced Research in Environmental Genomics, University of Ottawa, Ottawa, ON, Canada;*

Quebec Center for Biodiversity Science, McGill University, Montreal, QC, Canada;

Agriculture and Agri-Food Canada, Eastern Cereal and Oilseeds Research Center,

Central Experimental Farm, Ottawa, ON, Canada

DAVID SANKOFF • *Department of Mathematics and Statistics,*

University of Ottawa, Ottawa, ON, Canada

JULIANE SIEBOURG • *Department of Biosystems Science and Engineering,*

ETH Zurich, Basel, Switzerland

HAIXU TANG • *School of Informatics and Computing, Indiana University,*

Bloomington, IN, USA