# Computational Methods for Next Generation Sequencing Data Analysis

EDITED BY
Ion I. Măndoiu • Alexander Zelikovsky

with website

WILEY

# COMPUTATIONAL METHODS FOR NEXT GENERATION SEQUENCING DATA ANALYSIS

Edited by

**ION I. MĂNDOIU**
**ALEXANDER ZELIKOVSKY**

# WILEY

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States 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 may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Cover image courtesy of Gettyimages/Andrew Brookes

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

# CONTRIBUTORS

**Vanessa Aguiar-Pulido,** Bioinformatics Research Group (BioRG), School of Computing and Information Sciences, Florida International University, Miami, FL, USA

**Sahar Al Seesi,** Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA

**Alexander Artyomenko,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Niko Beerenwinkel,** Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

**Adrian Caciula,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**David S. Campo,** Division of Viral Hepatitis, Centers of Disease Control and Prevention, Atlanta, GA, USA

**Michael Campos,** Miller School of Medicine, University of Miami, Miami, FL, USA

**Stefan Canzar,** Center for Computational Biology, McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, and Toyota Technological Institute at Chicago, Chicago, IL, USA

**Jeong-Hyeon Choi,** Cancer Center, Medical College of Georgia, Georgia Regents University, Augusta, GA, USA; Department of Biostatistics and Epidemiology, Medical College of Georgia, Georgia Regents University, Augusta, GA, USA

**Chong Chu,** Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA

**Zoya Dimitrova,** Division of Viral Hepatitis, Centers of Disease Control and Prevention, Atlanta, GA, USA

**Jorge Duitama,** Agrobiodiversity Research Area, International Center for Tropical Agriculture (CIAT), Cali, Colombia

**Eleazar Eskin,** Department of Computer Science, University of California, Los Angeles, CA, USA

**Mitch Fernandez,** Bioinformatics Research Group (BioRG), School of Computing and Information Sciences, Florida International University, Miami, FL, USA

**Liliana Florea,** Center for Computational Biology, McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA

**Olga Glebova,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Xuan Guo,** Department of Computer Science, Department of Biology, Georgia State University, Atlanta, GA, USA

**Steven J. Hallam,** Graduate Program in Bioinformatics and Department of Microbiology and Immunology, University of British Columbia, Vancouver, BC, Canada

**Niels W. Hanson,** Graduate Program in Bioinformatics, University of British Columbia, Vancouver, BC, Canada

**Elena Harris,** Department of Computer Science, California State University, Chico, CA

**Wenrui Huang,** Bioinformatics Research Group (BioRG), School of Computing and Information Sciences, Florida International University, Miami, FL, USA

**Mazhar I. Khan,** Department of Pathobiology and Veterinary Science, University of Connecticut, Storrs, CT, USA

**Yury Khudyakov,** Division of Viral Hepatitis, Centers of Disease Control and Prevention, Atlanta, GA, USA

**Kishori M. Konwar,** Department of Microbiology and Immunology, University of British Columbia, Vancouver, BC, Canada

**Bing Li,** Department of Computer Science, Department of Biology, Georgia State University, Atlanta, GA, USA

**James Lindsay,** Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA

**Rasiah Loganantharaj,** Bioinformatics Research Lab, The Center for Advanced Computer Studies, University of Louisiana, Lafayette, LA, USA

**Stefano Lonardi,** Department of Computer Science and Engineering, University of California, Riverside, CA, USA

**Nicholas Mancuso,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Ion I. Măndoiu,** Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA

**Igor Mandric,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Serghei Mangul,** Department of Computer Science, University of California, Los Angeles, CA, USA

**Tobias Marschall,** Centrum Wiskunde & Informatica, Amsterdam, Netherlands

**Kalai Mathee,** Herbert Wertheim College of Medicine, Florida International University, Miami, FL, USA

**Giri Narasimhan,** Bioinformatics Research Group (BioRG), School of Computing and Information Sciences, Florida International University, Miami, FL, USA

**Ekaterina Nenastyeva,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Rachel O'neill,** Department of Molecular and Cell Biology, University of Connecticut, Storrs, CT, USA

**Yi Pan,** Department of Computer Science, Department of Biology, Georgia State University, Atlanta, GA, USA

**Sumathi Ramachandran,** Division of Viral Hepatitis, Centers of Disease Control and Prevention, Atlanta, GA, USA

**Thomas A. Randall,** Integrative Bioinformatics, National Institute of Environmental Health Sciences, Research Triangle Park, NC, USA

**Juan Riveros,** Bioinformatics Research Group (BioRG), School of Computing and Information Sciences, Florida International University, Miami, FL, USA

**Alexander Schönhuth,** Centrum Wiskunde & Informatica, Amsterdam, Netherlands

**Jonathan Segal,** Herbert Wertheim College of Medicine, Florida International University, Miami, FL, USA

**Huidong Shi,** Cancer Center, Medical College of Georgia, Georgia Regents University, Augusta, GA, USA Department of Biochemistry, Medical College of Georgia, Georgia Regents University, Augusta, GA, USA

**Pavel Skums,** Division of Viral Hepatitis, Centers of Disease Control and Prevention, Atlanta, GA, USA

**Ren Sun,** Department of Molecular and Medical Pharmacology, University of California, Los Angeles, CA, USA

**Sing-hoi Sze,** Department of Computer Science and Engineering and Department of Biochemistry and Biophysics, Texas A&M University, College Station, TX, USA

**Yvette Temate-tiagueu,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Armin Töpfer,** Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

**Bassam Tork,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Nicholas C. Wu,** Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA, USA

**Shang-ju Wu,** Department of Computer Science, University of British Columbia, Vancouver, BC, Canada

**Yufeng Wu,** Department of Computer Science and Engineering, University of Connecticut, Storrs, CT, USA

**Ning Yu,** Department of Computer Science, Department of Biology, Georgia State University, Atlanta, GA, USA

**Alexander Zelikovsky,** Department of Computer Science, Georgia State University, Atlanta, GA, USA

**Erliang Zeng,** Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN, USA

**Jin Zhang,** McDonnell Genome Institute, Washington University in St. Luis, MO, USA

# PREFACE

Massively parallel DNA sequencing and RNA sequencing have become widely available, reducing the cost by several orders of magnitude and placing the capacity to generate gigabases to terabases of sequence data into the hands of individual investigators. These so-called *next-generation sequencing (NGS)* technologies have dramatically accelerated biological and biomedical research by enabling the comprehensive analysis of genomes and transcriptomes to become inexpensive, routine, and widespread. The ensuing explosion in the volume of data has spurred numerous advances in computational methods for NGS data analysis.

This book aims to provide an in-depth survey of some of the most important recent developments in this area. It is neither intended as an introductory text nor as a comprehensive review of existing bioinformatics tools and active research areas in NGS data analysis. Rather, our intention is to make a carefully selected set of advanced computational techniques accessible to a broad readership, including graduate students in bioinformatics and related areas and biomedical professionals who want to expand their repertoire of computational techniques for NGS data analysis. We hope that our emphasis on in-depth presentation of both algorithms and software for computational data analysis of current high-throughput sequencing technologies will best prepare the readers for developing their own algorithmic techniques and for successfully implementing them in existing and novel NGS applications.

The book features 18 chapters authored by bioinformatics experts who are active contributors to the respective subjects. The chapters are intended to be largely independent, so that readers do not have to read every chapter nor have to read them in a particular order. The chapters are grouped into the following four parts:

- Part I focuses on computing and experimental infrastructure for NGS data analysis, including chapters on cloud computing, a modular pipeline for metabolic pathway reconstruction, pooling strategies for massive viral sequencing, and high-fidelity sequencing protocols.

- Part II concentrates on analyses of DNA sequencing data and includes chapters on the classic scaffolding problem, detection of genomic variants, two chapters on finding insertions and deletions, and two chapters on the analysis of DNA methylation sequencing data.
- Part III is devoted to analyses of RNA-seq data. Two chapters describe algorithms and compare software tools for transcriptome assembly: one chapter focuses on methods for alternative splicing analysis and the other chapter focuses on tools for transcriptome quantification and differential expression analysis.
- Part IV explores computational tools for NGS applications in microbiomics. The first chapter concentrates on error correction of NGS reads from viral populations, then two chapters describe methods for viral quasispecies reconstruction, and the last chapter surveys the state of the art and future trends in microbiome analysis.

We are grateful to all the authors for their excellent contributions, without which this book would not have been possible. We hope that their deep insights and fresh enthusiasm will help in attracting new generations of researchers to this dynamic field. We would also like to thank Yi Pan and Albert Y. Zomaya for nurturing this project since its inception, and the editorial staff at Wiley Interscience for their patience and assistance throughout the project. Finally, we wish to thank our friends and families for their continuous support.

ION I. MĂNDOIU
Storrs, Connecticut
ALEXANDER ZELIKOVSKY
Atlanta, Georgia

# ABOUT THE COMPANION WEBSITE

This book is accompanied by a companion website:

**www.wiley.com/go/Mandoiu/NextGenerationSequencing**

The book companion website contains the color version of a few selected figures

# CONTENTS IN BRIEF

# CONTENTS