

Computing for Comparative Microbial Genomics

Computational Biology

Editors-in-chief

Andreas Dress
University of Bielefeld (Germany)

Martin Vingron
Max Planck Institute for Molecular Genetics (Germany)

Editorial Board

Gene Myers, Janelia Farm Research Campus, Howard Hughes Medical Institute (USA)
Robert Giegerich, University of Bielefeld (Germany)
Walter Fitch, University of California, Irvine (USA)
Pavel A. Pevzner, University of California, San Diego (USA)

Advisory Board

Gordon Grippen, University of Michigan (USA)
Joe Felsenstein, University of Washington (USA)
Dan Gusfield, University of California, Davis (USA)
Sorin Istrail, Brown University, Providence (USA)
Samuel Karlin, Stanford University (USA)
Thomas Lengauer, Max Planck Institut Informatik (Germany)
Marcella McClure, Montana State University (USA)
Martin Nowak, Harvard University (USA)
David Sankoff, University of Ottawa (Canada)
Ron Shamir, Tel Aviv University (Israel)
Mike Steel, University of Canterbury (New Zealand)
Gary Stormo, Washington University Medical School (USA)
Simon Tavaré, University of Southern California (USA)
Tandy Warnow, University of Texas, Austin (USA)

The *Computational Biology* series publishes the very latest high-quality research devoted to specific issues in computer-assisted analysis of biological data. The main emphasis is on current scientific developments and innovative techniques in computational biology (bioinformatics), bringing to light methods from mathematics, statistics, and computer science that directly address biological problems currently under investigation.

The series offers publications that present the state of the art regarding the problems in question; show computational biology/bioinformatics methods at work, and discuss anticipated demands regarding developments in future methodology. Titles can range from focused monographs, to undergraduate and graduate textbooks, to professional text/reference works.

Author guidelines: springer.com>Authors>Author Guidelines

For other titles published in this series, go to <http://www.springer.com/series/5769>

David W. Ussery • Trudy M. Wassenaar •
Stefano Borini

Computing for Comparative Microbial Genomics

Bioinformatics for Microbiologists

 Springer

David W. Ussery, PhD
Department of Systems Biology
The Technical University of Denmark
Lyngby, Denmark
dave@cbs.dtu.dk

Trudy M. Wassenaar, PhD
Molecular Microbiology
and Genomics Consultants
Zotzenheim, Germany
wassenaar_t@yahoo.co.uk

Stefano Borini, PhD
Laboratory of Physical Chemistry
Swiss Federal Institute of
Technology (ETH)
Zurich, Switzerland
stefano.borini@igc.phys.chem.ethz.ch

Computational Biology Series ISSN 1568–2684
ISBN 978-1-84800-254-8 e-ISBN 978-1-84800-255-5
DOI 10.1007/978-1-84800-255-5

Library of Congress Control Number: 2008940847

© Springer-Verlag London Limited 2009

Apart from any fair dealing for the purposes of research or private study, or criticism or review, as permitted under the Copyright, Designs and Patents Act 1988, this publication may only be reproduced, stored or transmitted, in any form or by any means, with the prior permission in writing of the publishers, or in the case of reprographic reproduction in accordance with the terms of licences issued by the Copyright Licensing Agency. Enquiries concerning reproduction outside those terms should be sent to the publishers.

The use of registered names, trademarks, etc., in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant laws and regulations and therefore free for general use.

The publisher makes no representation, express or implied, with regard to the accuracy of the information contained in this book and cannot accept any legal responsibility or liability for any errors or omissions that may be made.

Printed on acid-free paper

Springer Science + Business Media
springer.com

Preface

Overview and Goals

This book describes how to visualize and compare bacterial genomes. Sequencing technologies are becoming so inexpensive that soon going for a cup of coffee will be more expensive than sequencing a bacterial genome. Thus, there is a very real and pressing need for high-throughput computational methods to compare hundreds and thousands of bacterial genomes.

It is a long road from molecular biology to systems biology, and in a sense this text can be thought of as a path bridging these fields. The goal of this book is to provide a coherent set of tools and a methodological framework for starting with raw DNA sequences and producing fully annotated genome sequences, and then using these to build up and test models about groups of interacting organisms within an environment or ecological niche.

Organization and Features

The text is divided into four main parts: Introduction, Comparative Genomics, Transcriptomics and Proteomics, and finally Microbial Communities. The first five chapters are introductions of various sorts. Each of these chapters represents an introduction to a specific scientific field, to bring all readers up to the same basic level before proceeding on to the methods of comparing genomes. First, a brief overview of molecular biology and of the concept of sequences as biological information are given. The equivalent in the post-genomics era of the ‘Central Dogma’ of molecular biology (DNA makes RNA makes protein) is that the genome makes the transcriptome, which makes the proteome. Before going on to the details of this, a historical background is provided that pictures the scene of the origins of molecular biology and biological sequences. After this introduction, Chapter 2 describes sequence alignment, the most common procedure used to compare biological sequences. Instead of going into technical details of how exactly these alignments are calculated, the text focuses on their practical use. Chapter 3 introduces bacterial genomes and Chapter 4 deals with the most important databases, whilst Chapter 5 is

an introduction to the computational background of the tools necessary to analyze all of this information.

The second part, on Comparative Genomics (Chapters 6–8), describes some basic methods of comparing genomes. This section introduces various atlases building up to the ‘Genome Atlas,’ which is our standard visualization tool for representing the DNA sequence of a chromosome in a single figure, mapping the most relevant DNA properties along the chromosome. We have found such atlases very useful for mapping newly sequenced genomes and quickly visualizing regions of potential interest. The value of atlas projections is illustrated by the examples provided.

Part three (Chapters 9–11) takes the reader from genome sequences to RNA sequences (transcriptomics) to proteins (proteomics) and regulation of gene expression. An important overview of experimental results can be obtained by mapping back and visualizing the transcriptomic and proteomic data onto physical chromosomal maps. Examples illustrate how important chromosome location is, and which features can be predicted by careful analysis of genes and their surrounding sequences.

The final part (Chapters 12–14) deals with microbial communities. In a sense this can be thought of as ‘population genomics’ (as opposed to the more traditional ‘population biology’ which often focuses on only one or a few genes). First the concept of ‘pan-genome’ and ‘core genome’ is introduced (Chapter 12), followed by metagenomics (Chapter 13), and then evolution of microbial communities (Chapter 14). From a larger perspective, population genomics can provide a framework for modeling ecosystems in terms of interacting biological systems.

Target Audiences and Required Background Knowledge

The reader should have basic knowledge about computers and be able to use web interfaces. For programmers, some general knowledge of microbiology is assumed, but it is our hope that both programmers and more ‘biology-oriented’ readers will find this book helpful. Details on programming were deliberately left out; instead, the text concentrates on the use and interpretation of publicly available web tools. This book has grown out of lectures for the course in Comparative Microbial Genomics,¹ which DWU has taught since 2001 as a full semester length course at the Technical University of Denmark, and as one-week workshops given in Bangkok, Thailand; in Petropolis, Brazil; and in Oslo, Norway.

This book is in a sense merging different scientific languages. The three authors have different scientific and national backgrounds. DWU is from the U.S., studied biochemistry, worked in molecular biology, and for the last 10 years has led a group

¹ <http://www.cbs.dtu.dk/dtucourse/programme27444.php>

in bioinformatics and genomics. SB is from Italy, studied quantum chemistry with focus on scientific programming, data standardization, and software integration; whereas TMW studied biochemistry and worked in molecular biology and later as a consultant in microbiology. These different backgrounds actually helped to develop a common language in science. The subject area of this textbook is extremely interdisciplinary, covering (bio)chemistry, physics, biology, microbiology, mathematics, and computational science, and by the introduction of concepts (and some jargon) from these various disciplines, the different languages used by specialists are bridged.

This book is meant mainly for people studying bacterial genomes, although of course nearly all of the methods described in the text would work for viral, Archaeal, or Eukaryotic genomes as well. There are two main target audiences. The first is the microbiologist who wants to get the most out of a bacterial genome sequence. This could be a university student, or an experienced laboratory microbiologist who enters the field of genomics. This book enables one to get a handle on how to use high-throughput computational methods to compare only a few, or hundreds of sequenced genomes. The second audience comprises the computer programmers who assist these microbiologists in actually carrying out the analyses. From experience we know there can be communication problems between the experimental bacteriologist who is more laboratory-oriented, and the computer scientist who wants to do everything on computers. Both disciplines are essential in present-day research. This book aims to explain to the computational scientist why and how we want to study bacterial genomes, and what questions we hope to answer. At the same time, it explains to the biologist some of the basics behind the bioinformatic tools that are necessary for research in the field. Bringing these two worlds, scientific interests, and languages together is our ultimate goal.

Notes to the Instructor

There are no exercises or questions at the end of the chapters, although at the end of most chapters textboxes present descriptions of essential methods used. From experience we can say that giving small groups of students a project in which they can choose a recently sequenced bacterial genome and compare it to other similar genomes can produce surprisingly successful results. It is very motivating to work with recently published data (new genome sequence papers are being published on an almost daily basis now), and sometimes the students produce important observations that the authors of the scientific papers had missed! In some occasions, such activities have resulted in a real scientific publication by the students, illustrating how ‘easy’ it is to do these kinds of analyses, as long as one asks relevant questions.

Supplemental Resources

A number of web links are mentioned in the book, and since web addresses are not always stable, a dedicated web page is put up on which all web pages presented in the book are summarized, and as necessary, updated. This can be found at <http://comparativemicrobial.com>.

Lyngby, Denmark
Zurich, Switzerland
Zotzenheim, Germany

David Ussery
Stefano Borini
Trudy Wassenaar

Acknowledgements

This book is based on input from many people, including our research team and external collaborators. We are grateful for all the advice, assistance, and help we received throughout this project. We thank all current and former members of the Comparative Microbial Genomics group at CBS: in particular, Peter F. Hallin for his excellent programming skills and help with development of many of the programs mentioned in this book; Flemming Hansen for his work on bacterial replication and his vast knowledge of *E. coli*; Henrik J. Nielsen for his help with *E. coli* genomics; Kristoffer Kiil for his help with phylogeny and work on protein function; and Carsten Friis for his assistance with various analyses and for keeping the group running whilst we were writing.

We thank former group members whose work also contributed to this book, including Tim T. Binnewies for his work with *Vibrio* genomes and secretion systems, and Hanni Willenbrock for her work with developing pan-genome microarrays.

We are grateful to external collaborators, notably Thomas Quinn from Denver University for his work on phylogenetic trees whilst on sabbatical in our group; Karin Lagesen from CMBN, Institute of Medical Microbiology, Rikshospitalet University Hospital in Oslo, who has helped with the rRNA and tRNA searches; and Jon Bohlin from the Norwegian School of Veterinary Science, who has helped with analysis of oligonucleotide usage patterns in bacterial genomes.

We would also like to acknowledge help from the many people at CBS, which is currently one of the largest bioinformatics groups in Europe. In particular, we thank Hans Henrik Stærfeldt from the CBS systems administration group, who wrote the original code for the GeneWiz program that is used to construct the atlas plots, and for his help and support over the past 10 years in updating and maintaining GeneWiz. Jannick D. Bendtsen helped us on the secretome, and Thomas Blicher kindly provided wonderful pictures of protein structures. Finally, Søren Brunak, center director for CBS has established a wonderful place to work (including an excellent coffee machine!) and has been supportive of our group since it was formed in 1998.

David would like to thank his students over the years in his Comparative Microbial Genomics course for their many helpful suggestions and comments. He would also like to thank his wife for helpful editorial comments and for her support during the writing of this book.

Stefano would like to thank his parents, Paola Marani and Walter Padovani, for their constant support and trust, and his dear friends Paolo Soriani and Ruggero Paratelli for their life-long support and understanding.

Trudy would like to thank her son Martijn for inventing the analogy of a road to explain DNA strand direction, both of her sons for their understanding and patience, and her husband for his constant support.

Much of the work described in this textbook has been funded by grants over the past decade from the Danish National Research Foundation (Danmarks Grundforskningsfond), Danish Research Councils, and the EU. Many of the calculations presented in this book have been made on our large computer system at CBS, funded in part by money from the Danish Center for Scientific Computing.

Contents

Preface	v
Acknowledgements	ix

Part I Introductions

1 Sequences as Biological Information: Cells Obey the Laws of Chemistry and Physics	3
Why Study Microbes?.....	3
What is Biological Information and Where Does It Come From	5
How DNA Sequences Code for Information	7
From DNA to Protein: Transcription and Translation.....	9
DNA Sequences: More than Protein-Coding Genes	12
From DNA to DNA: Replication	14
Proteins: Structure and Function.....	14
2 Bioinformatics for Microbiologists: An Introduction	19
Identifying Similarities: Sequence Comparison by Means of Alignments	19
From Alignments to Phylogenic Relationships.....	28
Genome Annotation: the Challenge to Get It Right.....	31
Information Beyond the Single Genome	33
3 Microbial Genome Sequences: A New Era in Microbiology	37
The First Completely Sequenced Microbial Genome.....	37
The Importance of Visualization	38
Genome Atlases to Visualize Chromosomes	42
A Race Against the Clock: The Speed of Sequencing	44
The First Completely Sequenced Bacterial Genome	46
Comparative Bacterial Genomics	47
The Microbial Genome: Not All Bacteria Are Like <i>E. coli</i>	50
4 An Overview of Genome Databases	53
What is a Database?	54

Three Databases Storing Sequences and a Lot More.....	57
Data Files and Formats	61
RNA Databases	62
Protein Databases.....	64
5 The Challenges of Programming: a Brief Introduction	69
Part 1: A Brief Overview of Computer Science Concepts.....	69
A Look at the Most Common Bioinformatic Procedures.....	73
Achieving Better Automation	81
Part 2: Some Technical Details and Future Directions	83
Programming Languages	83
Markup Languages.....	86
Service Oriented Architecture.....	88
Specific Tools for Bioinformatic Use.....	89

Part II Comparative Genomics

6 Methods to Compare Genomes: the First Examples	95
Genomic Comparisons: The Size of a Genome	95
Pairwise Alignment of Genomes	99
Comparing Gene Content and Annotation Quality	100
RNA Comparisons: A Look at rRNAs	102
Proteome Comparisons: What Makes a Family?.....	103
7 Genomic Properties: Length, Base Composition and DNA Structures.....	111
Length of Genomes: the ‘C-Value Paradox’	112
Genome Average Base Composition: The Percentage of AT.....	114
GC Skew—Bias Towards The Replication Leading Strand	118
Global Chromosomal Bias of AT Content	122
DNA Structures.....	125
The Structure Atlas.....	128
Bias In Purines—A-DNA Atlases.....	129
More on Structure Atlases.....	131
8 Word Frequencies and Repeats	137
Analyzing Word Frequencies in a Genome.....	137
DNA Repeats Within a Chromosome	139
Introduction to the DNA Repeat Atlas	143
Local DNA Repeats are Related to Chromosomal AT Content	146
DNA Structures Related to Repeats in Sequences.....	147
The Genome Atlas: Our Standard Method for Visualization	147

Part III Transcriptomics and Proteomics

9 Transcriptomics: Translated and Untranslated RNA..... 153
 Counting rRNA and tRNA Genes 154
 A Closer Look at Ribosomal RNA..... 155
 Genes Encoding Transfer RNA..... 160
 Genes Coding mRNA: Comparing Codon Usage Between Bacteria 161
 Other Non-coding RNA: tmRNA 164

10 Expression of Genes and Proteins 167
 Comparing Gene Expression and Protein Expression 168
 Part 1: Regulation of Transcription..... 169
 Part 2: Regulation of Translation 179
 Part 3: Protein Modification and Cellular Localization 180
 Antigen and Epitope Prediction 185

11 Of Proteins, Genomes, and Proteomes 189
 Part 1: Analysis of Individual Protein-Coding Genes..... 190
 Part 2: How to Annotate a Complete Genome 197
 Part 3: Proteome Comparisons..... 203

PART IV MICROBIAL COMMUNITIES

12 Microbial Communities: Core and Pan-Genomics..... 213
 Defining Pan-Genomes and Core Genomes 214
 Current Data Available for Pan- and Core Genome Analysis..... 218
 The Pan- and Core Genome of *Streptococcus* 219
 The Current *Bacillus* Pan- and Core Genome..... 221
 An Overview of Some Proteobacterial Pan- and Core Genomes 222
 The *Burkholderia* Pan- and Core Genome..... 223

13 Metagenomics of Microbial Communities..... 229
 Metagenomics Based on 16S rRNA Analysis..... 230
 Metagenomics Based on Complete DNA Sequencing..... 232
 Environmental Influences on Base Composition..... 234
 Visualization of Environmental Metagenomic Data..... 235
 Marine Metagenomics 240
 Other Metagenomics Applications..... 241

14 Evolution of Microbial Communities; or, On the Origins of Bacterial Species 243
 Where Does Diversity Come From?..... 244

Evolution Takes Time 245
Evidence of Evolution in a Single Genome..... 247
Genome Islands..... 249
Evolution on a Chip 252
Species and Speciation: *Vibrio cholerae*..... 253
Can We Predict Evolution? *Escherichia coli* Genome Reduction..... 253

Abbreviations 257

Index..... 263