

AN INTRODUCTION TO **BIOINFORMATICS ALGORITHMS**

NEIL C. JONES AND PAVEL A. PEVZNER



An
Introduction
to
Bioinformatics
Algorithms

Neil C. Jones

Pavel A. Pevzner

A Bradford Book
The MIT Press
Cambridge, Massachusetts
London, England

© 2004 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

MIT Press books may be purchased at special quantity discounts for business or sales promotional use. For information, please email special_sales@mitpress.mit.edu or write to Special Sales Department, The MIT Press, 5 Cambridge Center, Cambridge, MA 02142.

Typeset in 10/13 Lucida Bright by the authors using L^AT_EX 2_ε.

Printed and bound in the United States of America.

Library of Congress Cataloging-in-Publication Data

Jones, Neil C.

An introduction to bioinformatics algorithms/ by Neil C. Jones and Pavel A.

Pevzner.

p. cm.—(computational molecular biology series)

"A Bradford book."

Includes bibliographical references and index (p.).

ISBN 0-262-10106-8 (hc : alk. paper)

1. Bioinformatics. 2. Algorithms. I. Pevzner, Pavel. II. Title

QH324.2.J66 2004

570'.285—dc22

2004048289

CIP

10 9 8 7 6 5 4 3

An
Introduction
to
Bioinformatics
Algorithms



Computational Molecular Biology

Sorin Istrail, Pavel Pevzner, and Michael Waterman, editors

Computational molecular biology is a new discipline, bringing together computational, statistical, experimental, and technological methods, which is energizing and dramatically accelerating the discovery of new technologies and tools for molecular biology. The MIT Press Series on Computational Molecular Biology is intended to provide a unique and effective venue for the rapid publication of monographs, textbooks, edited collections, reference works, and lecture notes of the highest quality.

Computational Molecular Biology: An Algorithmic Approach
Pavel A. Pevzner, 2000

Computational Methods for Modeling Biochemical Networks
James M. Bower and Hamid Bolouri, editors, 2001

Current Topics in Computational Molecular Biology
Tao Jiang, Ying Xu, and Michael Q. Zhang, editors, 2002

Gene Regulation and Metabolism: Postgenomic Computation Approaches
Julio Collado-Vides, editor, 2002

Microarrays for an Integrative Genomics
Isaac S. Kohane, Alvin Kho, and Atul J. Butte, 2002

Kernel Methods in Computational Biology
Bernhard Schölkopf, Koji Tsuda, and Jean-Philippe Vert, 2004

An Introduction to Bioinformatics Algorithms
Neil C. Jones and Pavel A. Pevzner, 2004

For my Grandfather, who lived to write and vice versa.

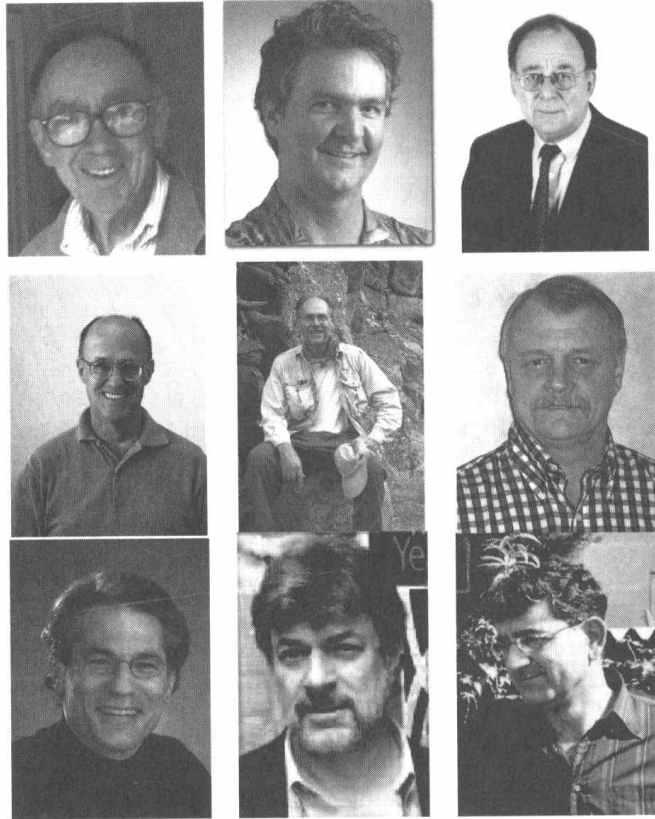
—NCJ

To Chop-up and Manny.

—PAP

Contents in Brief

[illegible]



Featuring historical perspectives from:

Russell Doolittle	Chapter 3
David Haussler	Chapter 11
Richard Karp	Chapter 2
Webb Miller	Chapter 7
Gene Myers	Chapter 9
David Sankoff	Chapter 5
Ron Shamir	Chapter 10
Gary Stormo	Chapter 4
Michael Waterman	Chapter 6

Preface

In the early 1990s when one of us was teaching his first bioinformatics class, he was not sure that there would be enough students to teach. Although the Smith-Waterman and BLAST algorithms had already been developed they had not become the household names among biologists that they are today. Even the term “bioinformatics” had not yet been coined. DNA arrays were viewed by most as intellectual toys with dubious practical application, except for a handful of enthusiasts who saw a vast potential in the technology. A few bioinformaticians were developing new algorithmic ideas for nonexistent data sets: David Sankoff laid the foundations of genome rearrangement studies at a time when there was practically no gene order data, Michael Waterman and Gary Stormo were developing motif finding algorithms when there were very few promoter samples available, Gene Myers was developing sophisticated fragment assembly tools when no bacterial genome has been assembled yet, and Webb Miller was dreaming about comparing billion-nucleotide-long DNA sequences when the 172,282-nucleotide Epstein-Barr virus was the longest GenBank entry. GenBank itself just recently made a transition from a series of bound (paper!) volumes to an electronic database on magnetic tape that could be sent to scientists worldwide.

One has to go back to the mid-1980s and early 1990s to fully appreciate the revolution in biology that has taken place in the last decade. However, bioinformatics has affected more than just biology—it has also had a profound impact on the computational sciences. Biology has rapidly become a large source of new algorithmic and statistical problems, and has arguably been the target for more algorithms than any of the other fundamental sciences. This link between computer science and biology has important educational implications that change the way we teach computational ideas to biologists, as well as how applied algorithmics is taught to computer scientists.

For many years computer science was taught to only computer scientists, and only rarely to students from other disciplines. A biology student in an algorithms class would be a surprising and unlikely (though entirely welcome) guest in the early 1990s. But these things change; many biology students now take some sort of Algorithms 101. At the same time, curious computer science students often take Genetics 101 and Bioinformatics 101. Although these students are still relatively rare, keep in mind that the number of bioinformatics classes in the early 1990s was so small as to be considered nonexistent. But that number is not so small now. We envision that undergraduate bioinformatics classes will become a permanent component at every major university. This is a feature, not a bug.

This is an introductory textbook on bioinformatics algorithms and the computational ideas that have driven them through the last twenty years. There are many important probabilistic and statistical techniques that we do not cover, nor do we cover many important research questions that bioinformaticians are currently trying to answer. We deliberately do not cover all areas of computational biology; for example, important topics like protein folding are not even discussed. The very first bioinformatics textbooks were Waterman, 1995 (108), which contains excellent coverage of DNA statistics and Gusfield, 1997 (44) which includes an encyclopedia of string algorithms. Durbin et al., 1998 (31) and Baldi and Brunak, 1997 (7) emphasize Hidden Markov Models and machine learning techniques; Baxevanis and Ouellette, 1998 (10) is an excellent practical guide to bioinformatics; Mount, 2001 (76) excels in showing the connections between biological problems and bioinformatics techniques; and Bourne and Weissig, 2002 (15) focuses on protein bioinformatics. There are also excellent web-based lecture notes for many bioinformatics courses and we learned a lot about the pedagogy of bioinformatics from materials on the World Wide Web by Serafim Batzoglou, Dick Karp, Ron Shamir, Martin Tompa, and others.

Website

We have created an extensive website to accompany this book at

<http://www.bioalgorithms.info>

This website contains a number of features that complement the book. For example, though this book does not contain a glossary, we provide this service, a searchable index, and a set of community message boards, at the above web address. Technically savvy students can also download practical

bioinformatics exercises, sample implementations of the algorithms in this book, and sample data to test them with. Instructors and students may find the prepackaged lecture notes on the website to be especially helpful. It is our hope that this website be used as a repository of information that will help introduce students to the diverse world of bioinformatics.

Acknowledgements

We are indebted to those who kindly agreed to be featured in the biographical sketches scattered throughout the book. Their insightful and heartfelt responses definitely made these the most interesting part of this book. Their life stories and views of the challenges that lay ahead will undoubtedly inspire students in the exploration of the unknown. There are many more scientists whose bioboxes we would like to have in this book and it is only the page limit (which turned out to be 200 pages too small) that prevented us from commissioning more of them. Special thanks go to Ethan Bier who inspired us to include biographical sketches in this book.

This book would not have been possible without the diligent teaching assistants in bioinformatics courses taught during the winter and fall of 2003 and 2004: Derren Barken, Bryant Forsgren, Eugene Ke, Coleman Mosley, and Degui Zhi all helped find technical errors, refine practical exercises, and design problems in the book. Helen Wu and John Allison spent many hours making technical figures, which is a thankless task like no other. We are also grateful to Vagisha Sharma who was kind enough to read the book from cover to cover and provide insightful comments and, unfortunately, bugs in the pseudocode. Steve Wasserman provided us with invaluable comments from a biologist's point of view that eventually led to new sections in the book. Alkes Price and Haixu Tang pointed out ambiguities and helped clarify the chapters on graphs and clustering. Ben Raphael and Patricia Jones provided feedback on the early chapters and helped avoid some potential misunderstandings. Dan Gilbert, of Dan Gilbert Art Group, Inc. kindly provided us with Triazzles to illustrate the problems of DNA sequence assembly.

Our special thanks go to Randall Christopher, the artist behind the website www.kleemanandmike.com. Randall illustrated the book and designed many unique graphical representations of some bioinformatics algorithms.

It has been a pleasure to work with Robert Prior of The MIT Press. With sufficient patience and prodding, he managed to keep us on track. We also appreciate the meticulous copyediting of G. W. Helfrich.

Finally, we thank the many students in different undergraduate and graduate bioinformatics classes at UCSD who provided comments on earlier versions of this book.

PAP would like to thank several people who taught him different aspects of computational molecular biology. Andrey Mironov taught him that common sense is perhaps the most important ingredient of any applied research. Mike Waterman was a terrific teacher at the time PAP moved from Moscow to Los Angeles, both in science and in life. PAP also thanks Alexander Karzanov, who taught him combinatorial optimization, which, surprisingly, remains the most useful set of skills in his computational biology research. He especially thanks Mark Borodovsky who convinced him to switch into the field of bioinformatics in 1985, when it was an obscure discipline with an uncertain future.

PAP also thanks his former students, postdocs, and lab members who taught him most of what he knows: Vineet Bafna, Guillaume Bourque, Sridhar Hannenhalli, Steffen Heber, Earl Hubbell, Uri Keich, Zufar Mulyukov, Alkes Price, Ben Raphael, Sing-Hoi Sze, Haixu Tang, and Glenn Tesler.

NCJ would like to thank his mentors during undergraduate school—Toshiko Takeuchi, Harry Gray, John Baldeschwieler, and Schubert Soares—for patiently but firmly teaching him that persistence is one of the more important ingredients in research. Also, he thanks the admissions committee at the University of California, San Diego who gambled on a chemist-turned-programmer, hopefully for the best.

Neil Jones and Pavel Pevzner
La Jolla, California, 2004

Contents

Preface	xv
1 Introduction	1
2 Algorithms and Complexity	7
2.1 What Is an Algorithm?	7
2.2 Biological Algorithms versus Computer Algorithms	14
2.3 The Change Problem	17
2.4 Correct versus Incorrect Algorithms	20
2.5 Recursive Algorithms	24
2.6 Iterative versus Recursive Algorithms	28
2.7 Fast versus Slow Algorithms	33
2.8 Big-O Notation	37
2.9 Algorithm Design Techniques	40
2.9.1 Exhaustive Search	41
2.9.2 Branch-and-Bound Algorithms	42
2.9.3 Greedy Algorithms	43
2.9.4 Dynamic Programming	43
2.9.5 Divide-and-Conquer Algorithms	48
2.9.6 Machine Learning	48
2.9.7 Randomized Algorithms	48
2.10 Tractable versus Intractable Problems	49
2.11 Notes	51
Biobox: Richard Karp	52
2.12 Problems	54

3	Molecular Biology Primer	57
3.1	What Is Life Made Of?	57
3.2	What Is the Genetic Material?	59
3.3	What Do Genes Do?	60
3.4	What Molecule Codes for Genes?	61
3.5	What Is the Structure of DNA?	61
3.6	What Carries Information between DNA and Proteins?	63
3.7	How Are Proteins Made?	65
3.8	How Can We Analyze DNA?	67
3.8.1	Copying DNA	67
3.8.2	Cutting and Pasting DNA	71
3.8.3	Measuring DNA Length	72
3.8.4	Probing DNA	72
3.9	How Do Individuals of a Species Differ?	73
3.10	How Do Different Species Differ?	74
3.11	Why Bioinformatics?	75
	Biobox: Russell Doolittle	79
4	Exhaustive Search	83
4.1	Restriction Mapping	83
4.2	Impractical Restriction Mapping Algorithms	87
4.3	A Practical Restriction Mapping Algorithm	89
4.4	Regulatory Motifs in DNA Sequences	91
4.5	Profiles	93
4.6	The Motif Finding Problem	97
4.7	Search Trees	100
4.8	Finding Motifs	108
4.9	Finding a Median String	111
4.10	Notes	114
	Biobox: Gary Stormo	116
4.11	Problems	119
5	Greedy Algorithms	125
5.1	Genome Rearrangements	125
5.2	Sorting by Reversals	127
5.3	Approximation Algorithms	131
5.4	Breakpoints: A Different Face of Greed	132
5.5	A Greedy Approach to Motif Finding	136
5.6	Notes	137

	Biobox: David Sankoff	139
5.7	Problems	143
6	Dynamic Programming Algorithms	147
6.1	The Power of DNA Sequence Comparison	147
6.2	The Change Problem Revisited	148
6.3	The Manhattan Tourist Problem	153
6.4	Edit Distance and Alignments	167
6.5	Longest Common Subsequences	172
6.6	Global Sequence Alignment	177
6.7	Scoring Alignments	178
6.8	Local Sequence Alignment	180
6.9	Alignment with Gap Penalties	184
6.10	Multiple Alignment	185
6.11	Gene Prediction	193
6.12	Statistical Approaches to Gene Prediction	197
6.13	Similarity-Based Approaches to Gene Prediction	200
6.14	Spliced Alignment	203
6.15	Notes	207
	Biobox: Michael Waterman	209
6.16	Problems	211
7	Divide-and-Conquer Algorithms	227
7.1	Divide-and-Conquer Approach to Sorting	227
7.2	Space-Efficient Sequence Alignment	230
7.3	Block Alignment and the Four-Russians Speedup	234
7.4	Constructing Alignments in Subquadratic Time	238
7.5	Notes	240
	Biobox: Webb Miller	241
7.6	Problems	244
8	Graph Algorithms	247
8.1	Graphs	247
8.2	Graphs and Genetics	260
8.3	DNA Sequencing	262
8.4	Shortest Superstring Problem	264
8.5	DNA Arrays as an Alternative Sequencing Technique	265
8.6	Sequencing by Hybridization	268
8.7	SBH as a Hamiltonian Path Problem	271

8.8	SBH as an Eulerian Path Problem	272
8.9	Fragment Assembly in DNA Sequencing	275
8.10	Protein Sequencing and Identification	280
8.11	The Peptide Sequencing Problem	284
8.12	Spectrum Graphs	287
8.13	Protein Identification via Database Search	290
8.14	Spectral Convolution	292
8.15	Spectral Alignment	293
8.16	Notes	299
8.17	Problems	302
9	Combinatorial Pattern Matching	311
9.1	Repeat Finding	311
9.2	Hash Tables	313
9.3	Exact Pattern Matching	316
9.4	Keyword Trees	318
9.5	Suffix Trees	320
9.6	Heuristic Similarity Search Algorithms	324
9.7	Approximate Pattern Matching	326
9.8	BLAST: Comparing a Sequence against a Database	330
9.9	Notes	331
	Biobox: Gene Myers	333
9.10	Problems	337
10	Clustering and Trees	339
10.1	Gene Expression Analysis	339
10.2	Hierarchical Clustering	343
10.3	k -Means Clustering	346
10.4	Clustering and Corrupted Cliques	348
10.5	Evolutionary Trees	354
10.6	Distance-Based Tree Reconstruction	358
10.7	Reconstructing Trees from Additive Matrices	361
10.8	Evolutionary Trees and Hierarchical Clustering	366
10.9	Character-Based Tree Reconstruction	368
10.10	Small Parsimony Problem	370
10.11	Large Parsimony Problem	374
10.12	Notes	379
	Biobox: Ron Shamir	380
10.13	Problems	384

11 Hidden Markov Models	387
11.1 <i>CG</i> -Islands and the “Fair Bet Casino”	387
11.2 The Fair Bet Casino and Hidden Markov Models	390
11.3 Decoding Algorithm	393
11.4 HMM Parameter Estimation	397
11.5 Profile HMM Alignment	398
11.6 Notes	400
Biobox: David Haussler	403
11.7 Problems	407
12 Randomized Algorithms	409
12.1 The Sorting Problem Revisited	409
12.2 Gibbs Sampling	412
12.3 Random Projections	414
12.4 Notes	416
12.5 Problems	417
Using Bioinformatics Tools	419
Bibliography	421
Index	428