



HANDBOOK OF COMPARATIVE GENOMICS

PRINCIPLES AND METHODOLOGY

Cecilia Saccone
Graziano Pesole

HANDBOOK OF COMPARATIVE GENOMICS

Principles and Methodology

CECILIA SACCONE

Department of Biochemistry and Molecular Biology
University of Bari
Italy

GRAZIANO PESOLE

Department of Physiology and General Biochemistry
University of Milan
Italy

 **WILEY-LISS**

A JOHN WILEY & SONS PUBLICATION

Copyright © 2003 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 978-750-4470, or on the web at 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, e-mail: permreq@wiley.com.

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 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.

Library of Congress Cataloging-in-Publication Data:

Saccone, Cecilia.

Handbook of comparative genomics: Principles and methodology /
Cecilia Saccone, Graziano Pesole,
p. cm.

Includes bibliographical references and index.

ISBN 0-471-39128-X (cloth : alk. paper)

1. Genomics—Handbooks, manuals, etc. 2. Evolutionary genetics—Handbooks, manuals, etc. I. Pesole, Graziano.

II. Title.

[DNLM: 1. Genomics. 2. Computational Biology. 3. Evolution, Molecular. 4. Sequence Analysis, DNA. QU 58.5 S119h
2003]

QH447.S23 2003

572.8—dc21

2002011158

Printed in the United States of America

10 9 8 7 6 5 4 3 2

HANDBOOK OF COMPARATIVE GENOMICS

To
Ernesto Quagliariello

PREFACE

No sooner had the genomics era begun when a post-genomics era was declared. Since the beginning of large-scale sequencing, there has been a need for new approaches and tools for the study of living matter at the molecular level, and particularly of new hypotheses for progress in biology. There are tremendous expectations about future benefits of sequencing complete genomes; realistically, however, much still remains to be done. Indeed, the exploitation of genome information is still in its infancy and new methodologies may be required to make the most of them.

Some fundamental questions are: To what extent has large-scale sequencing increased our knowledge of the properties and functions of species and, in general, of the structure and function of the genome? What value-adding notions have emerged since sequencing? How has our knowledge of the molecular basis of biological processes improved, if at all? What is our capability to perform comparative genomic studies? What general rules can we derive?

It is well known that advances in comparative biology are based largely on the concepts of analogy and homology. For this reason, comparative analyses of complete genomes rather than single genes, gene families, or specific genome regions must be performed. The sequencing of complete genomes from prokaryotes, eukaryotes, and organelles has aided research on the structure and evolution of the genome as a unit, as opposed to the previous focus on its components. Hence, the new discipline of evolutionary genomics is emerging and will make a revolutionary impact in terms of helping to disclose the linear structure of nucleic acids and proteins, their three-dimensional folding, cytogenetics, gene expression, and regulatory pathways.

We are persuaded that comparative evolutionary genomics is the key to unraveling the hidden messages of living matter. Our use of *genomics* includes its expression and the regulatory mechanisms that are at the basis of all biological processes. In other words, genomics also includes transcriptomics, proteomics, and other “omics,” concepts that interconnect and overlap. “Omics” research needs broad

databases and dynamic technologies capable of facilitating collaborative efforts across many disciplines, such as biology, chemistry, informatics, mathematics, and physics.

The ambitious goal of this book is to offer a tool for trained biologists who want to tackle the new genomic dimension of modern biology. It will also be useful for technology developers, managers, industries, and funding agencies—in essence anyone interested in the exciting new applications in this new dimension. We are aware that the book is a personal vision and because of its extent and complexity, cannot cover all the literature. Nonetheless, we have tried to highlight milestones, emerging principles, key methods used, and the most urgent needs of this new field. Beginning with a description of complete genomes sequenced from living organisms, this handbook pinpoints new concepts emerging from available data. At the same time we describe the leading methods used to study complete genomes and their evolution.

In summary, we have written this book as a guide for students and researchers not necessarily specialized in genomics. It is written for anyone who approaches, with theoretical and practical aims, this intriguing new chapter of biology—perhaps already a new discipline. We have organized the book in three parts.

- In Part I we describe the state-of-the-art in molecular knowledge of the main biological processes achieved with modern biotechnological tools. We introduce recent insights brought about by genome sequencing. In particular, we summarize the major features of the genomes completely sequenced in prokaryotes, eukaryotes, and organelles.
- In Part II we illustrate the most recent methodologies used in genomics. We describe the available experimental and bioinformatics tools with particular emphasis on molecular biology techniques, biological databases, and computational methods for the analysis of sequence data.
- Part III contains data derived from comparative studies. We discuss fundamental, cutting-edge topics such as the evolution of genome size, base compositional constraints, and the structure and origin of organisms at the molecular level. We conclude by addressing recent advances in molecular phylogenetics.

ACKNOWLEDGMENTS

We thank our dearly missed friend and colleague Giuliano Preparata, a distinguished theoretical physicist who introduced us to the multidisciplinary approach in the study of molecular evolution. Warm thanks also to several colleagues and friends who offered advice and suggestions on several topics covered in the book; they include Marcella Attimonelli, Giorgio Bernardi, Rita Casadio, Nicola Cataldo, Victor De Lorenzo, Annamaria D'Erchia, Ilenia D'Errico, Carmela Gissi, Alessandro Minelli, Aurelio Reyes, Teresa M. R. Regina, Elisabetta Sbisà, and Apollonia Tullo. Finally, special thanks to Alessandra Larizza, for recovering the bibliography and organizing the material, and to Marilina Lonigro, for assistance in English.

CONTENTS

PREFACE	xiii
I GENOME FEATURES	1
1 PROKARYOTES	3
1.1 Introduction / 3	
1.2 Morphology and Classification / 8	
1.3 Genome Shape and Size / 10	
1.4 Gene Content and Organization / 16	
1.5 Base Composition / 20	
1.6 Codon Use / 27	
1.7 Replication and Expression / 29	
2 EUKARYOTES	36
2.1 Introduction / 36	
2.2 Classification and Time Scale / 36	
2.3 Genome Shape and Size / 40	
2.4 Base Composition / 46	
2.5 Replication, Repair, and Recombination / 48	
2.6 Gene Expression / 53	
2.6.1 Transcription and Posttranscriptional Regulation / 53	
2.6.2 Genetic Code and Codon Use / 59	
2.6.3 Translation and Posttranslation Modifications / 61	

- 2.7 Completely Sequenced Eukaryotic Genomes / 64
 - 2.7.1 *Saccharomyces cerevisiae* Genome / 66
 - 2.7.2 *Schizosaccharomyces pombe* Genome / 69
 - 2.7.3 *Caenorhabditis elegans* Genome / 70
 - 2.7.4 *Drosophila melanogaster* Genome / 72
 - 2.7.5 *Arabidopsis thaliana* Genome / 74
 - 2.7.6 *Oryza sativa* Genome / 77
 - 2.7.7 *Homo sapiens* Genome / 78

3 ORGANELLES 85

- 3.1 Mitochondria / 85
 - 3.1.1 General Structure and Function / 85
 - 3.1.2 DNA and Genetic System / 88
 - 3.1.3 Genome Features / 97
- 3.2 Chloroplasts and Other Plastids / 125

II METHODOLOGIES 131

4 MOLECULAR BIOLOGY TECHNIQUES FOR GENOMICS 133

- 4.1 Genome DNA Sequencing / 133
 - 4.1.1 DNA-Sequencing Techniques / 133
 - 4.1.2 The Human Genome Project / 136
- 4.2 Analysis of the Transcriptome / 137
 - 4.2.1 Analysis of Gene Expression / 137
 - 4.2.2 Expressed Sequence Tags / 138
 - 4.2.3 Serial Analysis of Gene Expression / 139
 - 4.2.4 Differential Display / 141
 - 4.2.5 Representational Difference Analysis / 142
 - 4.2.6 DNA Microarrays / 143
- 4.3 Analysis of the Proteome / 149
 - 4.3.1 Two-Dimensional Gel Electrophoresis / 151
 - 4.3.2 Protein Identification / 151
 - 4.3.3 Study of Protein–DNA and Protein–Protein Interactions / 155
 - 4.3.4 Proteome Analysis Using Biochips / 157

5 BIOLOGICAL DATABASES IN THE GENOMIC ERA 159

- 5.1 Introduction / 159
- 5.2 Primary and Specialized Databases / 159
- 5.3 Database Structures / 162
- 5.4 Linked Databases and Database Interoperability / 163

- 5.5 Database Annotation / 166
- 5.6 Retrieval Systems / 169
 - 5.6.1 SRS / 169
 - 5.6.2 Entrez / 169
 - 5.6.3 Other Retrieval Systems / 170
- 5.7 Nucleotide Databases / 170
- 5.8 Protein Databases / 171
- 5.9 Other Protein Databases / 171
- 5.10 Genomic Databases and Resources / 174
- 5.11 Gene Databases and Resources / 179
- 5.12 Transcriptome Databases / 180
- 5.13 Metabolism Databases / 181
- 5.14 Mutation Databases / 182
- 5.15 Mitochondrial Databases and Resources / 184

6 COMPUTATIONAL METHODS FOR THE ANALYSIS OF GENOME SEQUENCE DATA

187

- 6.1 Introduction / 187
- 6.2 Dot-Plot Matrix / 188
- 6.3 Sequence Pairwise Alignment / 189
 - 6.3.1 Needleman–Wunsch Global Alignment Algorithm / 191
 - 6.3.2 Smith–Waterman Algorithm for the Identification of
Common Molecular Subsequences / 193
 - 6.3.3 Alignment of cDNA and Genomic DNA Sequences / 196
 - 6.3.4 Genome Alignment / 198
 - 6.3.5 Cleanup of Sequence Databases from Redundancy / 201
 - 6.3.6 Measure of the Similarity Degree between Homologous
Sequences / 203
- 6.4 Database Searching / 208
 - 6.4.1 FASTA / 209
 - 6.4.2 BLAST / 211
 - 6.4.3 BLAST and FASTA Family of Programs / 215
 - 6.4.4 Filtering Matches to Unwanted Sequences / 216
 - 6.4.5 Filtering Matches to Repetitive Sequences / 221
 - 6.4.6 Statistical Significance of Alignment Scores / 223
- 6.5 Multiple Alignment / 226
- 6.6 Alignment Profiles to Recognize Distantly Related Protein or
Protein Modules / 230
- 6.7 Methods for Sequence Assembly / 234
 - 6.7.1 Sequence Cleaning / 236
 - 6.7.2 Sequence Clustering / 237
 - 6.7.3 Construction of Alignment Consensus / 238

6.7.4	Sequence Mapping by Electronic PCR /	240
6.7.5	Sequence Assembly in Genome and EST Projects /	240
6.7.6	Sequence Assembly for Gene Index Construction /	243
6.8	Linguistic Analysis of Biosequences /	245
6.8.1	Biosequences as Markov Chains /	248
6.8.2	Linguistic Complexity of Biosequences /	249
6.8.3	Identification of Repeats in Genomic Sequences /	253
6.8.4	Pattern Searching in Biosequences /	254
6.8.5	Identification of Promoter Regions in Chromosomal Sequences /	261
6.8.6	Pattern Discovery for the Identification of Gene Regulatory Elements and of Protein Motif Models /	262
6.8.7	Gene Prediction /	265
6.8.8	Identification of CpG Islands in Genomic Sequences /	269
6.8.9	Analysis of Codon Use Strategy /	271
6.9	Prediction of RNA Secondary Structures /	273
6.10	Protein Sequence Analysis /	280
6.10.1	Analysis of Protein Primary Sequences /	281
6.10.2	Prediction of Transmembrane Protein Helices /	285
6.10.3	Identification of Protein Signal Peptides and Prediction of Their Subcellular Location /	290
6.10.4	Prediction of Protein Secondary Structure /	293
6.10.5	Prediction of Coiled-Coil and Helix-Turn-Helix Structures /	298
6.10.6	Prediction of Protein Tertiary Structure /	300
6.10.7	Protein Fold Recognition and Classification /	302
6.10.8	Comparative Evolutionary Genomic Tools for Predicting Protein Function /	303
6.11	Evolutionary and Phylogenetic Analysis /	306
6.11.1	Estimating Genetic Distances between Homologous Sequences /	306
6.11.2	Molecular Phylogeny /	309

III COMPARATIVE GENOMICS 325

7 MOLECULAR EVOLUTION 327

7.1	Introduction /	327
7.2	Evolution of Genome Size /	328
7.3	Role of Base Composition in Evolution /	332
7.4	Evolution of the Prokaryotic Genome /	337
7.5	From Prokaryotes to Eukaryotes /	338
7.5.1	Origin of the Eukaryotic Cell /	338

7.5.2	Evolution of the Mitochondrial Genome /	340
7.5.3	Origin and Evolution of Plastids /	346
7.6	From Unicellular to Multicellular State /	349
7.7	Evolution of the Nuclear Genome /	351
7.7.1	Introns /	351
7.7.2	Gene and Protein Number /	352
7.7.3	Noncoding Elements /	353
7.7.4	Expansion of Gene Families /	353
7.7.5	Genome Duplication /	360
7.7.6	Conclusion /	361
8	MOLECULAR PHYLOGENY	362
8.1	Introduction /	362
8.2	Molecular Clock /	363
8.3	Similarity Measure: Orthology Versus Paralogy /	364
8.4	Molecular Phylogeny in the Genomics Era /	367
8.5	Interrelationships between Distant Taxa: The Tree of Life /	369
8.6	Phylogeny of Metazoans /	370
8.6.1	Organellar versus Nuclear Taxonomy /	370
8.6.2	Phylogeny of Mammals /	371
	APPENDIX: URLs Cited in the Text	375
	REFERENCES	377
	INDEX	423

PART I

GENOME FEATURES

CHAPTER 1

PROKARYOTES

1.1 INTRODUCTION

While this book is being written, complete sequences of bacterial genomes are being produced at a rate of about two genomes per month, and the National Center for Biotechnology Information (NCBI) Web site (see the URL in Table 5.1) reports about 60 completely sequenced prokaryotic genomes. Data reported in this chapter refer to the status of completely sequenced genomes, summarized in Table 1.1. Obviously, by the time you read this book, many more will have been sequenced and perhaps some of the aspects dealt with could be viewed differently, although we do not expect dramatic changes in our knowledge unless technology speeds its pace considerably.

Table 1.1 reports the prokaryotic genomes completely sequenced up to now and includes such features as species name, EMBL data library accession number, size, shape, presence of extrachromosomal elements, and bibliographic references. From a look at this list, one can gain an appreciation of the diverse reasons for promoting the sequencing of one species rather than another. Bacterial species are sequenced according to their research interest in basic or applied science: their importance for phylogenetic investigations, to shed light into the metabolic machinery (mainly Archaea) as well as for their importance as human and/or animal pathogens, and for their role as a source of industrial enzymes. In other words, priority has been given to species already well known or species presenting attractive opportunities in applied fields; thus from a phylogenetic point of view, the choice turns out to be very random.

We know we are at the infancy of the genomic era; despite the fact that completely sequenced organisms are still tiny in number, they have already turned out to be full of surprises. In this chapter we summarize the principal sequencing achievements that have improved our knowledge of the prokaryotic genomes and have contributed to outlining methods and approaches to be used in such studies.

TABLE 1.1. Prokaryotic Genomes Completely Sequenced

Species	Main Chromosome		Extrachromosomal Elements		References
	Accession Number	Size (bp)	Accession Number	Size (bp)	
Archaea					
<i>Aeropyrum pernix</i>	BA000002	1,669,695			Kawarabayasi, Hino et al. (1999)
<i>Archaeoglobus fulgidus</i>	AE000782	2,178,400			Klenk, Clayton et al. (1997)
<i>Halobacterium</i> sp. NRC-1 (3 chromosomes)	AE004437	2,014,239	AF016485	191,346	Ng, Kennedy et al. (2000)
	AE004438	365,425	AE004438	365,425	
	AF016485	191,346			
<i>Methanobacterium thermoautotrophicum</i>	AE000666	1,751,377			Smith, Doucette-Stamm et al. (1997)
	L77117	1,664,970	L77118	58,407	
<i>Methanococcus jannaschii</i>			L77119	16,550	Bult, White et al. (1996)
<i>Methanococcus kandleri</i> AV19	AE009439	1,694,969			Slesarev, Mezhevaya et al. (2002)
<i>Methanosarcina acetivorans</i> str. C2A	AE010299	5,751,492			Galagan, Nusbaum et al. (2002)
<i>Methanosarcina mazei</i> Goel	AE008384	4,096,345			Deppenmeier, Johann et al. (2002)
<i>Pyrobaculum aerophilum</i>	AE009441	2,222,430			Fitz-Gibbon, Ladner et al. (2002)
<i>Pyrococcus abyssi</i>	AL096836	1,765,118			Lecompte, Ripp et al. (2001)
<i>Pyrococcus furiosus</i> DSM 3638	AE009950	1,908,256			Robb, Maeder et al. (2001)
<i>Pyrococcus horikoshii</i>	AP000001–AP000007	1,738,505			Kawarabayasi, Sawada et al. (1998)
<i>Sulfolobus solfataricus</i>	AE006641	2,992,245			She, Singh et al. (2001)
<i>Sulfolobus tokodaii</i>	BA000023	2,694,765	AJ010405	41,229	Kawarabayasi, Hino et al. (2001)
<i>Thermoplasma acidophilum</i>	AL445063–AL445067	1,564,905			Ruepp, Graml et al. (2000)
	AP000991–AP000996	1,584,799			
<i>Thermoplasma volcanium</i>					Kawashima, Amano et al. (2000)
Bacteria					
<i>Agrobacterium tumefaciens</i> str. C58(Cereon)	AE007869	2,841,581			Goodner, Hinkle et al. (2001)
<i>Agrobacterium tumefaciens</i> str. C58(U.Washington)	AE008688	2,841,490	AE008687 AE008690	542,780 214,234	Wood, Setubal et al. (2001)

<i>Aquifex aeolicus</i>	AE000657	1,551,335	AE000667	39,456	Deckert, Warren et al. (1998)
<i>Bacillus halodurans</i>	BA000004	4,202,353			Takami, Nakasone et al. (2000)
<i>Bacillus subtilis</i>	AL009126	4,214,814			Kunst, Ogasawara et al. (1997)
<i>Borrelia burgdorferi</i> ⁿ	AE000783	910,725	AE000791	9,386	Fraser, Casjens et al. (1997);
			AE000792	26,498	Casjens, Palmer et al. (2000)
			AE001575	30,750	
			AE001576	30,223	
			AE001577	30,299	
			AE001578	29,838	
			AE001579	30,800	
			AE001580	30,885	
			AE001581	30,651	
			AE001583 ^a	5,228	
			AE000793 ^a	16,823	
			AE001582 ^a	18,753	
			AE000785 ^a	24,177	
			AE000794 ^a	26,921	
			AE000786 ^a	29,766	
			AE000784 ^a	28,601	
			AE000789 ^a	27,323	
			AE000788 ^a	36,849	
			AE000787 ^a	38,829	
			AE000790 ^a	53,561	
			AE001584 ^a	52,971	
<i>Brucella melitensis</i>	AE008917	2,117,144			DeVecchio, Kapatral et al. (2002)
<i>Buchnera</i> sp. APS	AP000398	640,681	AP001070	7,258	Shigenobu, Watanabe et al. (2000)
			AP001071	7,786	
<i>Campylobacter jejuni</i>	AL111168	1,641,481			Parkhill, Wren et al. (2000)
<i>Caulobacter crescentus</i>	AE005673	4,016,947			Nierman, Feldblyum et al. (2001)
<i>Chlamydia pneumoniae</i> AR39	AE002161	1,229,853		4,524	Read, Brunham et al. (2000)
<i>Chlamydia pneumoniae</i> CWL029	AE001363	1,230,230			Kalman, Mitchell et al. (1999)
<i>Chlamydia trachomatis</i> MoPn	AE002160	1,069,412	AE002162	7,501	Read, Brunham et al. (2000)
<i>Chlamydia trachomatis</i> serovar D	AE001273	1,042,519			Stephens, Kalman et al. (1998)
<i>Chlamydophila pneumoniae</i> J138	BA000008	1,226,565			Shirai, Hirakawa et al. (2000)
<i>Clostridium acetobutylicum</i>	AE001437	3,940,880	NC_001988	192,000	Nolling, Breton et al. (2001)