



BIOINFORMATICS & BIOMEDICAL IMAGING

Biological Database Modeling

Jake Chen • Amandeep S. Sidhu
editors

Q811.4
B615.2

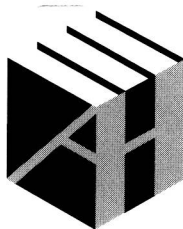
Biological Database Modeling

Jake Chen
Amandeep S. Sidhu

Editors



E2009002943



**ARTECH
HOUSE**

BOSTON | LONDON
artechhouse.com

Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the U.S. Library of Congress.

British Library Cataloguing in Publication Data

A catalogue record for this book is available from the British Library.

ISBN-13: 978-1-59693-258-6

Cover design by Igor Valdman

© 2008 ARTECH HOUSE, INC.

685 Canton Street

Norwood, MA 02062

All rights reserved. Printed and bound in the United States of America. No part of this book may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage and retrieval system, without permission in writing from the publisher.

All terms mentioned in this book that are known to be trademarks or service marks have been appropriately capitalized. Artech House cannot attest to the accuracy of this information. Use of a term in this book should not be regarded as affecting the validity of any trademark or service mark.

10 9 8 7 6 5 4 3 2 1

Biological Database Modeling

Artech House Bioinformatics & Biomedical Imaging
Stephen T. C. Wong and Guang-Zhong Yang, Series Editors

Advanced Methods and Tools for ECG Data Analysis, Gari D. Clifford, Francisco Azuaje, and Patrick E. McSharry, editors

Biomolecular Computation for Bionanotechnology, Jian-Qin Liu and Katsunori Shimohara

Electrotherapeutic Devices: Principles, Design, and Applications, George D. O'Clock

Intelligent Systems Modeling and Decision Support in Bioengineering, Mahdi Mahfouf

Life Science Automation Fundamentals and Applications, Mingjun Zhang, Bradley Nelson, and Robin Felder, editors

Matching Pursuit and Unification in EEG Analysis, Piotr Durka

Microfluidics for Biotechnology, Jean Berthier and Pascal Silberzan

Systems Bioinformatics: An Engineering Case-Based Approach, Gil Alterovitz and Marco F. Ramoni, editors

Text Mining for Biology and Biomedicine, Sophia Ananiadou and John McNaught

List of Contributors

- Chapter 1 Amandeep S. Sidhu,
 Curtin University of Technology, Perth, Australia
 Jake Chen,
 Indiana University School of Informatics, Indianapolis, Indiana, United States
- Chapter 2 Viroj Wiwanitkit,
 Chulalongkorn University, Bangkok, Thailand
- Chapter 3 Ramez Elmasri, Feng Ji, and Jack Fu,
 University of Texas at Arlington, Arlington, Texas, United States
- Chapter 4 Viroj Wiwanitkit,
 Chulalongkorn University, Bangkok, Thailand
- Chapter 5 Amandeep S. Sidhu, Tharam S. Dillon, and Elizabeth Chang,
 Curtin University of Technology, Perth, Australia
- Chapter 6 Cornelia Hedeler and Paolo Missier,
 University of Manchester, Manchester, United Kingdom
- Chapter 7 Greg Butler, Wendy Ding, John Longo, Jack Min, Nick O'Toole, Sindhu Pillai,
 Ronghua Shu, Jian Sun, Yan Yang, Qing Xie, Regis-Olivier Benech, Aleks Spurmanis,
 Peter Ulyczynj, Justin Powlowski, Reg Storms, and Adrian Tsang,
 Concordia University, Canada
- Chapter 8 Willy A. Valdivia-Granda,
 Orion Integrated Biosciences, Inc., New York, United States
 Christopher Dwan,
 BioTeam, Inc., Cambridge, Massachusetts, United States
- Chapter 9 Zhong Yan and Jake Chen,
 Indiana University School of Informatics, Indianapolis, Indiana, United States
 Josh Heyen, Lee W. Ott, Maureen A. Harrington, and Mark G. Goebel,
 Indiana University School of Medicine and the Walther Cancer Institute,
 Indianapolis, Indiana, United States
 Cary Woods,
 Ball State University, Muncie, Indiana, United States
- Chapter 10 Karthik Raman, Yeturu Kalidas, and Nagasuma Chandra,
 Supercomputer Education and Research Centre and Bioinformatics Centre, Indian
 Institute of Science, India
- Chapter 11 Preeti Malik, Tammy Chan, Jody Vandergriff, Jennifer Weisman, Joseph DeRisi, and
 Rahul Singh,
 San Francisco State University, San Francisco, California, United States

Preface

Database management systems (DBMS) are designed to manage large and complex data sets. In the past several decades, advances in computing hardware and software and the need to handle rapidly accumulating data archived in digital media have led to significant progress in DBMS research and development. DBMS have grown from simple software programs that handled flat files on mainframe computers, which were prohibitively expensive to all but a few prestigious institutions, into today's popular form of specialized software platforms underpinning wide ranges of tasks, which include business transactions, Web searches, inventory management, financial forecasts, multimedia development, mobile networks, pervasive computing, and scientific knowledge discovery. Technologies of DBMS have also become increasingly sophisticated, diverging from generic relational DBMS into object-relational DBMS, object-oriented DBMS, in-memory DBMS, semantic Webs data store, and specialized scientific DBMS. Given the sustained exponential data growth rate brought forth by continued adoption of computing in major industries and new inventions of personal digital devices, one can safely predict that DBMS development will continue to thrive in the next millennium.

In this book, we want to share with our readers some fresh research perspectives of post-genome biology data management, a fast-growing area at the intersection of life sciences and scientific DBMS domains. Efficient experimental techniques, primarily DNA sequencing, microarrays, protein mass spectrometers, and nanotechnology instruments, have been riding the wave of the digital revolution in the recent 20 years, leading to an influx of high-throughput biological data. This information overload in biology has created new post-genome biology studies such as genomics, functional genomics, proteomics, and metabolomics—collectively known as “omics” sciences in biology. While most experimental biologists are still making the transition from one-gene-at-a-time type of studies to the high-throughput data analysis mindset, many leaders of the field have already begun exploring new research and industrial application opportunities. For example, managing and interpreting massive omics data prelude ultimate systems biology studies, in which one may analyze disparate forms of biological data and uncover coordinated functions of the underlying biological systems at the molecular and cellular signalling network level. On the practical side, understanding diverse intricate interplays between environmental stimuli and genetic predisposition through omics evidence can help pharmaceutical scientists design drugs that target human proteins with high therapeutic values and low toxicological profiles. With data management tools to handle terabytes of omics data already released in the public domain, the promise of post-genome biology looms large.

Compared with data from general business application domains, omics data has many unique characteristics that make them challenging to manage. Examples of these data management challenges are:

1. Omics data tends to have more complex and more fast-evolving data structures than business data. Biological data representation often depends on scientific application scenarios. For example, biological sequences such as DNA and proteins can be either represented as simple character strings or connected nodes in three-dimensional spatial vectors. Data representation is an essential first step.
2. Omics data is more likely to come from more heterogeneously distributed locations than business data. To study systems biology, a bioinformatics researcher may routinely download genome data from the Genome Database Center at the University of California, Santa Cruz, collect literature abstracts from the PubMed database at the National Library of Medicine in Maryland, collect proteome information from the Swiss-Prot database in Switzerland, and collect pathway data from the KEGG database in Japan. Data integration has to be carefully planned and executed.
3. Omics data tends to reflect the general features of scientific experimental data: high-volume, noisy, formatted inconsistently, incomplete, and often semantically incompatible with one another. In contrast, data collected from business transactions tends to contain far fewer errors, is often more accurate, and shows more consistencies in data formats/coverage. Meticulous data preprocessing before knowledge discovery are required.
4. Omics data also lags behind business data in standard development. For example, Gene Ontology (GO) as a standard to control vocabularies for genes was not around until a decade ago, whereas standards such as industrial product categories have been around for decades. The ontology standards and naming standards for pathway biology are still under development. This makes it difficult to perform mega collaboration, in which cross-validation of results and knowledge sharing are both essential.

Despite all the challenges, modeling and managing biological data represent significant discovery opportunities in the next several decades. The human genome data bears the ultimate solutions of expanding the several thousand traditional molecular drug targets into tens of thousands genome drug targets; molecular profiling information, based on individuals using either the microarrays or the proteomics platform, promises new types of molecular diagnostics and personalized medicine. As new applications of massive biological data emerge, there will be an increasing need to address data management research issues in biology.

In this compiled volume, we present to our readers a comprehensive view of how to model the structure and semantics of biological data from public literature databases, high-throughput genomics, gene expression profiling, proteomics, and chemical compound screening projects. The idea of compiling this book, which we found to be unique, stems from the editors' past independent work in bioinformatics and biological data management. While topics in this area are diverse and interdisciplinary, we focused on a theme for this book—that is, how to model and manage

omics biological data in databases. By promoting this theme for the past decade among ourselves and the contributing authors of this book, we have contributed to solving complex biological problems and taking biological database management problems to the next level. We hope our readers can extract similar insights by using this book as a reference for future related activities.

There are 11 chapters presented in this book. Individual chapters have been written by selected accomplished research teams active in the research of respective topics. Each chapter covers an important aspect of the fast-growing topic of biological database modeling concepts. Each chapter also addresses its topic with varying degrees of balance between computational data modeling theories and real-world applications.

In Chapters 1 through 5, we introduce basic biological database concepts and general data representation practices essential to post-genome biology. First, biological data management concepts are introduced (Chapter 1) and major public database efforts in omics and systems biology studies are summarized (Chapter 2). Then, biomedical data modeling techniques are introduced (Chapter 3). Next, Gene Ontology as an established basic set of controlled vocabulary in genome database annotations is described (Chapter 4). Finally, the latest research on protein ontology and the use of related semantic webs technologies are presented to enable readers to make the connection between emerging biological data collection and integration trends (Chapter 5).

In Chapters 6 through 9, we examine in detail how to develop data management techniques to process and analyze high-throughput biological data through case studies. First, quality control techniques to reduce variations during experimental data collection steps are described (Chapter 6). Then, biological sequence management experience for a fungi genomics project is discussed (Chapter 7). Next, data management and data integration methods for microarray-based functional genomics studies are investigated (Chapter 8). Finally, data management challenges and opportunities for mass spectrometry based expression proteomics are presented (Chapter 9).

In Chapters 10 and 11, we delve into the practical aspect, demonstrating how to apply biological data management for drug discoveries. First, fundamental drug discovery concepts based on macromolecular structural modeling are introduced (Chapter 10); then, a data management software system that implements high-throughput drug compound screenings is discussed (Chapter 11) to conclude the book.

We hope this book will become a useful resource for bioinformatics graduate students, researchers, and practitioners interested in managing post-genome biological data. By studying the techniques and software applications described in this book, we hope that bioinformatics students will use the book material as a guide to acquire basic concepts and theories of post-genome biological data management, bioinformatics practitioners will find valuable lessons for building future similar biological data management systems, and researchers will find rewarding research data management questions to address in the years to come.

Acknowledgments

We wish to thank all of the authors for sharing their insightful knowledge and making excellent contributions to this book based on their active research in biological data management. This book would not have been completed without the tremendous efforts, held to the highest standards, of all the authors, each of whom spent numerous hours over many drafts in preparing, collating, and revising their writings over the past 2 years. During the publishing process, many colleagues also helped and they deserve our whole-hearted appreciation. They are: David Wong from Indiana University, who provided legal advice for the contract agreement; Dr. Zongmin Ma from Northwestern University of China, who provided assistance in the initial conceptualization and execution of the book publishing process; Susan Lagerstrom-Fife from Springer Science+Business Media, Inc., who guided us from a publisher's perspective while we explored various publication options; Wayne Yuhasz from Artech House Publishers, whose persistence and dedication to timely assistance finally won us over in making our book part of the Artech House Bioinformatics & Biomedical Imaging Series; and Barbara Lovenvirth from Artech House Publishers, who assisted us throughout the final publication process.

We also want to thank the general support from our home institutions: Indiana University School of Informatics at Indianapolis, Indiana; Purdue University School of Science, Department of Computer and Information Sciences at Indianapolis, Indiana; and Curtin University of Technology, Perth, Australia.

Last, but not least, Jake Chen wishes to thank his family, including his wife, Mabel Liu, for her unbounded love and support in assuming additional responsibilities so that the book project could be completed.

Jake Chen
Indianapolis, Indiana
Amandeep S. Sidhu
Perth, Australia
Editors
October 2007

Contents

Preface	<i>xiii</i>
Acknowledgments	<i>xvii</i>
CHAPTER 1	
Introduction to Data Modeling	1
1.1 Generic Modern Markup Languages	1
1.2 Modeling Complex Data Structures	3
1.3 Data Modeling with General Markup Languages	3
1.4 Ontologies: Enriching Data with Text	4
1.5 Hyperlinks for Semantic Modeling	5
1.6 Evolving Subject Indexes	6
1.7 Languages	6
1.8 Views	7
1.9 Modeling Biological Data	7
References	8
CHAPTER 2	
Public Biological Databases for -Omics Studies in Medicine	9
2.1 Introduction	9
2.2 Public Databases in Medicine	10
2.3 Application of Public Bioinformatics Database in Medicine	11
2.3.1 Application of Genomic Database	11
2.3.2 Application of Proteomic Database	16
2.3.3 Application of the Metabolomics Database	18
2.3.4 Application of Pharmacogenomics Database	19
2.3.5 Application of Systemics Database	21
References	21
CHAPTER 3	
Modeling Biomedical Data	25
3.1 Introduction	25
3.2 Biological Concepts and EER Modeling	27
3.2.1 Sequence Ordering Concept	27
3.2.2 Input/Output Concept	29
3.2.3 Molecular Spatial Relationship Concept	30
3.3 Formal Definitions for EER Extensions	31
3.3.1 Ordered Relationships	31

3.3.2	Process Relationships	33
3.3.3	Molecular Spatial Relationships	34
3.4	Summary of New EER Notation	35
3.5	Semantic Data Models of the Molecular Biological System	35
3.5.1	The DNA/Gene Model	36
3.5.2	The Protein 3D Structure Model	36
3.5.3	The Molecular Interaction and Pathway Model	40
3.6	EER-to-Relational Mapping	41
3.6.1	Ordered Relationship Mapping	41
3.6.2	Process Relationship Mapping	42
3.6.3	Molecular Spatial Relationship Mapping	43
3.7	Introduction to Multilevel Modeling and Data Source Integration	45
3.8	Multilevel Concepts and EER Modeling	46
3.9	Conclusion	48
	References	49

CHAPTER 4

	Fundamentals of Gene Ontology	51
4.1	Introduction to Gene Ontology	51
4.2	Construction of an Ontology	52
4.3	General Evolution of GO Structures and General Annotation Strategy of Assigning GO Terms to Genes	56
4.3.1	General Evolution of GO Structures	56
4.3.2	General Annotation Strategy of Assigning GO Terms to Genes	57
4.4	Applications of Gene Ontology in Biological and Medical Science	57
4.4.1	Application of Gene Ontology in Biological Science	57
4.4.2	Application of Gene Ontology in Medical Science	58
	References	60

CHAPTER 5

	Protein Ontology	63
5.1	Introduction	63
5.2	What Is Protein Annotation?	64
5.3	Underlying Issues with Protein Annotation	64
5.3.1	Other Biomedical Ontologies	65
5.3.2	Protein Data Frameworks	66
5.3.3	Critical Analysis of Protein Data Frameworks	68
5.4	Developing Protein Ontology	68
5.5	Protein Ontology Framework	69
5.5.1	The ProteinOntology Concept	70
5.5.2	Generic Concepts in Protein Ontology	70
5.5.3	The ProteinComplex Concept	71
5.5.4	Entry Concept	71
5.5.5	Structure Concept	72
5.5.6	StructuralDomains Concept	72
5.5.7	FunctionalDomains Concept	73
5.5.8	ChemicalBonds Concept	74

5.5.9	Constraints Concept	74
5.5.10	Comparison with Protein Annotation Frameworks	75
5.6	Protein Ontology Instance Store	76
5.7	Strengths and Limitations of Protein Ontology	77
5.8	Summary	78
	References	78

CHAPTER 6

	Information Quality Management Challenges for High-Throughput Data	81
6.1	Motivation	81
6.2	The Experimental Context	84
6.2.1	Transcriptomics	86
6.2.2	Qualitative Proteomics	88
6.3	A Survey of Quality Issues	89
6.3.1	Variability and Experimental Design	89
6.3.2	Analysis of Quality Issues and Techniques	91
6.3.3	Specificity of Techniques and Generality of Dimensions	93
6.3.4	Beyond Data Generation: Annotation and Presentation	94
6.4	Current Approaches to Quality	96
6.4.1	Modeling, Collection, and Use of Provenance Metadata	96
6.4.2	Creating Controlled Vocabularies and Ontologies	97
6.5	Conclusions	98
	Acknowledgments	98
	References	98

CHAPTER 7

	Data Management for Fungal Genomics: An Experience Report	103
7.1	Introduction	103
7.2	Materials Tracking Database	109
7.3	Annotation Database	110
7.4	Microarray Database	111
7.5	Target Curation Database	111
7.6	Discussion	112
7.6.1	Issue of Data and Metadata Capture	113
7.7	Conclusion	116
	Acknowledgments	116
	References	116

CHAPTER 8

	Microarray Data Management: An Enterprise Information Approach	119
8.1	Introduction	119
8.2	Microarray Data Standardization	122
8.2.1	Gene Ontologies	123
8.2.2	Microarray Ontologies	125
8.2.3	Minimum Information About a Microarray Experiment	125
8.3	Database Management Systems	126
8.3.1	Relational Data Model	127

8.3.2	Object-Oriented Data Model	128
8.3.3	Object-Relational Data Model	131
8.4	Microarray Data Storage and Exchange	131
8.4.1	Microarray Repository	133
8.4.2	Microarray Data Warehouses and Datamarts	133
8.4.3	Microarray Data Federations	134
8.4.4	Enterprise Microarray Databases and M-KM	135
8.5	Challenges and Considerations	136
8.6	Conclusions	138
	Acknowledgments	138
	References	139

CHAPTER 9

	Data Management in Expression-Based Proteomics	143
9.1	Background	143
9.2	Proteomics Data Management Approaches	147
9.3	Data Standards in Mass Spectrometry Based Proteomics Studies	149
9.4	Public Repositories for Mass Spectrometry Data	152
9.5	Proteomics Data Management Tools	154
9.6	Expression Proteomics in the Context of Systems Biology Studies	155
9.7	Protein Annotation Databases	159
9.8	Conclusions	159
	References	160

CHAPTER 10

	Model-Driven Drug Discovery: Principles and Practices	163
10.1	Introduction	163
10.2	Model Abstraction	165
10.2.1	Evolution of Models	166
10.3	Target Identification	168
10.3.1	Sequence-to-Function Models	170
10.3.2	Sequence Alignments and Phylogenetic Trees	170
10.3.3	Structure-to-Function Models	172
10.3.4	Systems-Based Approaches	173
10.3.5	Target Validation	176
10.4	Lead Identification	177
10.4.1	Target Structure-Based Design	177
10.4.2	Ligand-Based Models	179
10.5	Lead to Drug Phase	182
10.5.1	Predicting Drug-Likeness	182
10.5.2	ADMET Properties	182
10.6	Future Perspectives	183
	Acknowledgments	184
	References	184

CHAPTER 11

Information Management and Interaction in High-Throughput Screening for Drug Discovery	189
11.1 Introduction	189
11.2 Prior Research	191
11.3 Overview of Antimalarial Drug Discovery	192
11.4 Overview of the Proposed Solution and System Architecture	193
11.5 HTS Data Processing	194
11.5.1 Introduction to HTS	194
11.5.2 Example of HTS for Antimalarial Drug Screening	195
11.6 Data Modeling	199
11.6.1 The Database Design	202
11.7 User Interface	204
11.8 Conclusions	206
Acknowledgments	207
References	207
Selected Bibliography	208
About the Authors	209
Index	217

Introduction to Data Modeling

Amandeep S. Sidhu and Jake Chen

Scientific data is often scattered among heterogeneous data repositories. Exploring data across multiple data repositories requires the ability to understand and correlate their structures (schemas). Such correlations need to address the diversity of views of the scientific domain represented by different data repositories as well as the diversity of data modeling languages used for expressing these views. In this chapter, we introduce the concepts of *data modeling* and discuss its application to *biological databases*.

1.1 Generic Modern Markup Languages

Modern markup languages, such as Standard Generalized Markup Language (SGML) [1] and eXtensible Markup Language (XML) [2], which were initially conceived for modeling texts, are now receiving increasing attention as formalisms for data and knowledge modeling. XML is currently establishing itself as a successor of HyperText Markup Language (HTML) for a better modeling of texts as well as of other kinds of data. There are several reasons for this evolution. Even though multiple databases may cover the same data, their focus might be different. Modern markup languages such as SGML and XML are generic in that:

- They serve to specify the semantic structure, not the layout, of documents or data items.
- They make it possible to freely specify application-dependent document or data structures.

In the following, the term “data” refers also, but not exclusively, to text data. Thus, a data item may consist of: (1) text only (such data items are also known as human-readable documents); (2) nontext only (such data items are also known as data-oriented documents); or (3) both (such data items are also known as mixed-model documents). In the terminology of generic markup languages, data items are called documents. In the following, the term “data item” is used in lieu of “document” for stressing that not only (structured) texts are meant, but more generally (structured) data of any kind.

Widespread specific markup languages such as PostScript or Rich Text Format (RTF), whose conceptual roots go back to the 1970s, serve to specify the layout of data items. Here, layout is not exclusively meant as the appearance of a data item when printed on paper, but more generally as any kind of presentation of a data item to human perception. Examples of such an extended notion of layout include the formats of data items as they are displayed on a terminal screen, rendered in the script on an output device, or presented by any other means on any device.

The family of generic markup languages started in the late 1980s with the conception of its first specimen, SGML. The purpose of a generic markup language is to specify the semantic—or logical—structure of data items, not their layout. In the following, the term “presentation” is reserved to refer to the layout of a data item in the extended sense above, while the term “representation” refers to how semantics is conveyed through structural elements of the underlying data modeling formalism.

The distinction between layout and structure is important, for a layout format is device or system dependent, whereas a semantic structure should not be. It is desirable that the semantic structure of data items be specified independently of any layout. This ensures both:

- Independence of data modeling from data usage;
- Independence of data modeling from presentation devices.

The first property, data independence from usage, is important because data is rarely used in a single manner only. The second property, data independence from presentation devices, is important for several reasons. To begin with, different kinds of presentation devices require different layouts. For example, a structurally complex data item is likely not to be displayed using identical layouts on standard size screens and on small screens like those of cellular phones. Also, such devices are likely to become technically obsolete sooner than data. Moreover, a presentation format does not necessarily fully convey data semantics. For instance, it is common practice to rely on printed text layout for conveying semantic structure when using text processing systems or the markup language HTML. This practice often leads to semantic losses, especially when files are transferred from one text processing system to another, because the layout of the one system cannot always be faithfully mapped into that of the other system.

In order to specify layouts for classes of documents specified in a generic markup language, so-called style-sheet languages are used in addition. These languages basically allow the definition of layouts for those structural elements specified with the markup language. Such definitions do not have to be unique, thus ensuring the desired independence of the data from their presentations in various contexts.

Generic markup languages (like the XML family of languages) do not impose any predefined structure, nor any predefined names for the structural elements occurring in data items. Structure and names can be freely chosen, hence the denomination of generic markup language. Thus, using generic markup languages it is possible to faithfully model the structure of data items needed in applications and to name the structural elements of a chosen structure in a way that is natural in the application context.