

Rui Jiang  
Xuegong Zhang  
Michael Q. Zhang *Editors*

# Basics of Bioinformatics

Lecture Notes of the Graduate Summer  
School on Bioinformatics of China

# Basics of Bioinformatics



Rui Jiang • Xuegong Zhang • Michael Q. Zhang  
Editors

# Basics of Bioinformatics

Lecture Notes of the Graduate Summer  
School on Bioinformatics of China



*Editors*

Rui Jiang  
Xuegong Zhang  
Department of Automation  
Tsinghua University  
Beijing  
China, People's Republic

Michael Q. Zhang  
Department of Molecular and Cell Biology  
The University of Texas at Dallas  
Richardson, TX, USA

Tsinghua National Laboratory  
for Information Science and Technology  
Tsinghua University  
Beijing, China, People's Republic

ISBN 978-3-642-38950-4

ISBN 978-3-642-38951-1 (eBook)

DOI 10.1007/978-3-642-38951-1

Springer Heidelberg New York Dordrecht London

Jointly published with Tsinghua University Press, Beijing

ISBN: 978-7-302-32359-4 Tsinghua University Press, Beijing

Library of Congress Control Number: 2013950934

© Tsinghua University Press, Beijing and Springer-Verlag Berlin Heidelberg 2013

This work is subject to copyright. All rights are reserved by the Publishers, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publishers' locations, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

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

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publishers can accept any legal responsibility for any errors or omissions that may be made. The publishers make no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

# Foreword

This ambitious volume is the result of the successful 2007 Graduate Summer School on Bioinformatics of China held at Tsinghua University. It is remarkable for its range of topics as well as the depth of coverage. Bioinformatics draws on many subjects for analysis of the data generated by the biological sciences and biotechnology. This foreword will describe briefly each of the 12 chapters and close with additional general comments about the field. Many of the chapters overlap and include useful introductions to concepts such as gene or Bayesian methods. This is a valuable aspect of the volume allowing a student various angles of approach to a new topic.

Chapter 1, “Basics for Bioinformatics,” defines bioinformatics as “the storage, manipulation and interpretation of biological data especially data of nucleic acids and amino acids, and studies molecular rules and systems that govern or affect the structure, function and evolution of various forms of life from computational approaches.” Thus, the first subject they turn to is molecular biology, a subject that has had an enormous development in the last decades and shows no signs of slowing down. Without a basic knowledge of biology, the bioinformatics student is greatly handicapped. From basic biology the authors turn to biotechnology, in particular, methods for DNA sequencing, microarrays, and proteomics. DNA sequencing is undergoing a revolution. The mass of data collected in a decade of the Human Genome Project from 1990 to 2001 can be generated in 1 day in 2010. This is changing the science of biology at the same time. A 1,000 genome project became a 10,000 genome project 2 years later, and one expects another zero any time now. Chromatin Immunoprecipitation or ChIP allows access to DNA bound by proteins and thus to a large number of important biological processes. Another topic under the umbrella of biological sciences is genetics, the study of heredity and inherited characteristics (phenotypes). Heredity is encoded in DNA and thus is closely related to the goals of bioinformatics. This whole area of genetics beginning with Mendel’s laws deserves careful attention, and genetics is a key aspect of the so-called genetic mapping and other techniques where the chromosomal locations of disease genes are sought.

Chapter 2, “Basic Statistics for Bioinformatics,” presents important material for the understanding and analysis of data. Probability and statistics are basic to bioinformatics, and this chapter begins with the fundamentals including many classical distributions (including the binomial, Poisson, and normal). Usually the observation of complete populations such as “all people in China over 35 years old” is not practical to obtain. Instead random samples of the population of interest are obtained and then inferences about parameters of the population are made. Statistics guides us in making those inferences and gaining information about the quality of the estimates. The chapter describes techniques such as method of moments, maximum likelihood, and Bayesian methods. Bayesian methods have become indispensable in the era of powerful computing machines. The chapter treats hypothesis testing which is less used than parameter estimation, but hypothesis testing provides understanding of  $p$ -values which are ubiquitous in bioinformatics and data analysis. Classical testing situations reveal useful statistics such as the  $t$ -statistic. Analysis of variance and regression analysis are crucial for testing and fitting large data sets. All of these methods and many more are included in the free open-source package called *R*.

Chapter 3, “Topics in Computational Genomics,” takes us on a tour of important topics that arise when complete genome information is available. The subject did not begin until nearly 2000 when complete genome sequences became a possibility. The authors present us with a list of questions, some of which are listed next. What are the genes of an organism? How are they turned off and on? How do they interact with each other? How are introns and exons organized and expressed in RNA transcripts? What are the gene products, both structure and function? How has a genome evolved? This last question has to be asked with other genomes and with members of the population comprising the species. Then the authors treat some of the questions in detail. They describe “finding protein coding genes,” “identifying promoters,” “genomic arrays and a CGH/CNP analysis,” “modeling regulatory elements,” “predicting transcription factor binding sites,” and motif enrichment and analysis. Within this last topic, for example, various word counting methods are employed including the Bayesian methods of expectation maximization and Gibbs sampling.

An alert reader will have noticed the prominence of Bayesian methods in the preceding paragraphs. Chapter 4, “Statistical Methods in Bioinformatics,” in this collection focuses on this subject. There is a nice discussion of statistical modeling and then Bayesian inference. Dynamic programming, a recursive method of optimization, is introduced and then employed in the development of Hidden Markov Models (HMMs). Of course the basics of Markov chains must also be covered. The Metropolis-Hastings algorithm, Monte Carlo Markov chains (MCMC), and Gibbs sampling are carefully presented. Then these ideas find application in the analysis of microarray data. Here the challenging aspects of multiple hypothesis testing appear, and false discovery rate analysis is described. Hierarchical clustering and bi-clustering appear naturally in the context of microarray analysis. Then the issues of sequence analysis (especially multiple sequence analysis) are approached using these HMM and Bayesian methods along with pattern discovery in the sequences.

Discovering regulatory sequence patterns is an especially important topic in this section. The topics of this chapter appear in computer science as “machine learning” or under “data mining”; here the subject is called statistical or Bayesian methods. Whatever it is named, this is an essential area for bioinformatics.

The next chapter (Chap. 5), “Algorithms in Computational Biology,” takes up the formal computational approach to our biological problems. It should be pointed out that the previous chapters contained algorithmic content, but there it was less acknowledged. It is my belief that the statistical and algorithmic approaches go hand in hand. Even with the Euclid’s algorithm example of the present chapter, there are statistical issues nearby. For example, the three descriptions of Euclid’s algorithm are analyzed for time complexity. It is easy to ask how efficient the algorithms are on randomly chosen pairs of integers. What is the expected running time of the algorithms? What is the variance? Amazingly these questions have answers which are rather deep. The authors soon turn to dynamic programming (DP), and once again they present clear illustrative examples, in this case Fibonacci numbers. Designing DP algorithms for sequence alignment is covered. Then a more recently developed area of genome rearrangements is described along with some of the impressive (and deep) results from the area. This topic is relevant to whole genome analysis as chromosomes evolve on a larger scale than just alterations of individual letters as covered by sequence alignment.

In Chap. 6, “Multivariate Statistical Methods in Bioinformatics Research,” we have a thorough excursion into multivariate statistics. This can be viewed as the third statistical chapter in this volume. Here the multivariate normal distribution is studied in its many rich incarnations. This is justified by the ubiquitous nature of the normal distribution. Just as with the bell-shaped curve which appears in one dimension due to the central limit theorem (add up enough independent random variables and suitably normalized, one gets the normal under quite general conditions), there is also a multivariate central limit theorem. Here detailed properties are described as well as related distributions such as the Wishart distribution (the analog of the chi-square). Estimation is relevant as is a multivariate *t*-test. Principal component analysis, factor analysis, and linear discriminant analysis are all covered with some nice examples to illustrate the power of approaches. Then classification problems and variable selection both give platforms to further illustrate and develop the methods on important bioinformatics application areas.

Chapter 7, “Association Analysis for Human Diseases: Methods and Examples,” gives us the opportunity to look more deeply into aspects of genetics. While this chapter emphasizes statistics, be aware that computational issues also drive much of the research and cannot be ignored. Population genetics is introduced and then the important subjects of genetic linkage analysis and association studies. Genomic information such as single-nucleotide polymorphisms (SNPs) provide voluminous data for many of these studies, where multiple hypothesis testing is a critical issue.

Chapter 8, “Data Mining and Knowledge Discovery Methods with Case Examples,” deals with the area of knowledge discovery and data mining. To quote the authors, this area “has emerged as an important research direction for extracting useful information from vast repositories of data of various types. The basic

concepts, problems and challenges deals with the area of knowledge discovery and data mining that has emerged as an important research direction for extracting useful information from vast repositories of data of various types. The basic concepts, problems and challenges are first briefly discussed. Some of the major data mining tasks like classification, clustering and association rule mining are then described in some detail. This is followed by a description of some tools that are frequently used for data mining. Two case examples of supervised and unsupervised classification for satellite image analysis are presented. Finally an extensive bibliography is provided.”

The valuable chapter on Applied Bioinformatics Tools (Chap. 9) provides a step-by-step description of the application tools used in the course and data sources as well as a list of the problems. It should be strongly emphasized that no one learns this material without actually having hands-on experience with the derivations and the applications. This is not a subject for contemplation only!

Protein structure and function is a vast and critically important topic. In this collection it is covered by Chap. 10, “Foundations for the Study of Structure and Function of Proteins.” There the detailed structure of amino acids is presented with their role in the various levels of protein structure (including amino acid sequence, secondary structure, tertiary structure, and spatial arrangements of the subunits). The geometry of the polypeptide chain is key to these studies as are the forces causing the three-dimensional structures (including electrostatic and van der Waals forces). Secondary structural units are classified into  $\alpha$ -helix,  $\beta$ -sheets, and  $\beta$ -turns. Structural motifs and folds are described. Protein structure prediction is an active field, and various approaches are described including homology modeling and machine learning.

Systems biology is a recently described approach to combining system-wide data of biology in order to gain a global understanding of a biological system, such as a bacterial cell. The science is far from succeeding in this endeavor in general, let alone having powerful techniques to understand the biology of multicellular organisms. It is a grand challenge goal at this time. The fascinating chapter on Computational Systems Biology Approaches for Deciphering Traditional Chinese Medicine (Chap. 11) seeks to apply the computational systems biology (CSB) approach to traditional Chinese medicine (TCM). The chapter sets up parallel concepts between CSB and TCM. In Sect. 11.3.2 the main focus is “on a CSB-based case study for TCM ZHENG—a systems biology approach with the combination of computational analysis and animal experiment to investigate Cold ZHENG and Hot ZHENG in the context of the neuro-endocrine-immune (NEI) system.” With increasing emphasis on the so-called nontraditional medicine, these studies have great potential to unlock new understandings for both CSB and TCM.

Finally I close with a few remarks about this general area. Biology is a major science for our new century; perhaps it will be the major science of the twenty-first century. However, if someone is not excited by biology, then they should find a subject that does excite them. I have almost continuously found the new discoveries such as introns or microRNA absolutely amazing. It is such a young science when such profound wonders keep showing up. Clearly no one analysis subject can

solve all the problems arising in modern computational molecular biology. Statistics alone, computer science alone, experimental molecular biology alone, none of these are sufficient in isolation. Protein structure studies require an entire additional set of tools such as classical mechanics. And as systems biology comes into play, systems of differential equations and scientific computing will surely be important. None of us can learn everything, but everyone working in this area needs a set of well-understood tools. We all learn new techniques as we proceed, learning things required to solve the problems. This requires people who evolve with the subject. This is exciting, but I admit it is hard work too. Bioinformatics will evolve as it confronts new data created by the latest biotechnology and biological sciences.

University of Southern California  
Los Angeles, USA  
March 2, 2013

Michael S. Waterman



# Contents

<b>1</b>	<b>Basics for Bioinformatics .....</b>	<b>1</b>
	Xuegong Zhang, Xueya Zhou, and Xiaowo Wang	
1.1	What Is Bioinformatics .....	1
1.2	Some Basic Biology .....	2
1.2.1	Scale and Time .....	3
1.2.2	Cells .....	3
1.2.3	DNA and Chromosome .....	4
1.2.4	The Central Dogma .....	5
1.2.5	Genes and the Genome .....	7
1.2.6	Measurements Along the Central Dogma .....	10
1.2.7	DNA Sequencing .....	10
1.2.8	Transcriptomics and DNA Microarrays .....	13
1.2.9	Proteomics and Mass Spectrometry .....	16
1.2.10	ChIP-Chip and ChIP-Seq .....	17
1.3	Example Topics of Bioinformatics .....	18
1.3.1	Examples of Algorithmic Topics .....	18
1.3.2	Examples of Statistical Topics .....	19
1.3.3	Machine Learning and Pattern Recognition Examples .....	20
1.3.4	Basic Principles of Genetics .....	21
	References .....	25
<b>2</b>	<b>Basic Statistics for Bioinformatics .....</b>	<b>27</b>
	Yuanlie Lin and Rui Jiang	
2.1	Introduction .....	27
2.2	Foundations of Statistics .....	27
2.2.1	Probabilities .....	27
2.2.2	Random Variables .....	30
2.2.3	Multiple Random Variables .....	32
2.2.4	Distributions .....	34

2.2.5	Random Sampling .....	37
2.2.6	Sufficient Statistics .....	39
2.3	Point Estimation .....	40
2.3.1	Method of Moments.....	41
2.3.2	Maximum Likelihood Estimators .....	41
2.3.3	Bayes Estimators .....	42
2.3.4	Mean Squared Error.....	44
2.4	Hypothesis Testing.....	44
2.4.1	Likelihood Ratio Tests .....	45
2.4.2	Error Probabilities and the Power Function .....	46
2.4.3	<i>p</i> -Values.....	48
2.4.4	Some Widely Used Tests .....	50
2.5	Interval Estimation .....	52
2.6	Analysis of Variance .....	54
2.6.1	One-Way Analysis of Variance .....	55
2.6.2	Two-Way Analysis of Variance .....	59
2.7	Regression Models.....	61
2.7.1	Simple Linear Regression.....	62
2.7.2	Logistic Regression .....	65
2.8	Statistical Computing Environments .....	66
2.8.1	Downloading and Installation .....	66
2.8.2	Storage, Input, and Output of Data .....	67
2.8.3	Distributions .....	67
2.8.4	Hypothesis Testing .....	68
2.8.5	ANOVA and Linear Model .....	68
	References .....	68
3	<b>Topics in Computational Genomics .....</b>	69
	Michael Q. Zhang and Andrew D. Smith	
3.1	Overview: Genome Informatics .....	69
3.2	Finding Protein-Coding Genes .....	71
3.2.1	How to Identify a Coding Exon? .....	72
3.2.2	How to Identify a Gene with Multiple Exons? .....	72
3.3	Identifying Promoters.....	73
3.4	Genomic Arrays and aCGH/CNP Analysis.....	75
3.5	Introduction on Computational Analysis of Transcriptional Genomics Data .....	76
3.6	Modeling Regulatory Elements .....	77
3.6.1	Word-Based Representations .....	77
3.6.2	The Matrix-Based Representation .....	78
3.6.3	Other Representations.....	79
3.7	Predicting Transcription Factor Binding Sites.....	79
3.7.1	The Multinomial Model for Describing Sequences .....	80
3.7.2	Scoring Matrices and Searching Sequences .....	81

3.7.3	Algorithmic Techniques for Identifying High-Scoring Sites .....	82
3.7.4	Measuring Statistical Significance of Matches .....	83
3.8	Modeling Motif Enrichment in Sequences .....	84
3.8.1	Motif Enrichment Based on Likelihood Models.....	84
3.8.2	Relative Enrichment Between Two Sequence Sets .....	86
3.9	Phylogenetic Conservation of Regulatory Elements .....	88
3.9.1	Three Strategies for Identifying Conserved Binding Sites .....	88
3.9.2	Considerations When Using Phylogenetic Footprinting .....	90
3.10	Motif Discovery .....	91
3.10.1	Word-Based and Enumerative Methods.....	92
3.10.2	General Statistical Algorithms Applied to Motif Discovery .....	93
3.10.3	Expectation Maximization .....	94
3.10.4	Gibbs Sampling .....	95
	References .....	96
<b>4</b>	<b>Statistical Methods in Bioinformatics .....</b>	<b>101</b>
	Jun S. Liu and Bo Jiang	
4.1	Introduction .....	101
4.2	Basics of Statistical Modeling and Bayesian Inference .....	102
4.2.1	Bayesian Method with Examples .....	102
4.2.2	Dynamic Programming and Hidden Markov Model ..	104
4.2.3	Metropolis–Hastings Algorithm and Gibbs Sampling ..	107
4.3	Gene Expression and Microarray Analysis .....	109
4.3.1	Low-Level Processing and Differential Expression Identification.....	110
4.3.2	Unsupervised Learning .....	113
4.3.3	Dimension Reduction Techniques .....	117
4.3.4	Supervised Learning .....	119
4.4	Sequence Alignment .....	126
4.4.1	Pair-Wise Sequence Analysis .....	126
4.4.2	Multiple Sequence Alignment .....	129
4.5	Sequence Pattern Discovery .....	133
4.5.1	Basic Models and Approaches .....	133
4.5.2	Gibbs Motif Sampler .....	136
4.5.3	Phylogenetic Footprinting Method and the Identification of <i>Cis</i> -Regulatory Modules .....	138
4.6	Combining Sequence and Expression Information for Analyzing Transcription Regulation .....	140
4.6.1	Motif Discovery in ChIP-Array Experiment .....	140
4.6.2	Regression Analysis of Transcription Regulation .....	141
4.6.3	Regulatory Role of Histone Modification .....	143

4.7	Protein Structure and Proteomics .....	144
4.7.1	Protein Structure Prediction .....	145
4.7.2	Protein Chip Data Analysis .....	146
	References .....	147
<b>5</b>	<b>Algorithms in Computational Biology .....</b>	<b>151</b>
	Tao Jiang and Jianxing Feng	
5.1	Introduction .....	151
5.2	Dynamic Programming and Sequence Alignment .....	153
5.2.1	The Paradigm of Dynamic Programming .....	153
5.2.2	Sequence Alignment .....	155
5.3	Greedy Algorithms for Genome Rearrangement .....	157
5.3.1	Genome Rearrangements .....	157
5.3.2	Breakpoint Graph, Greedy Algorithm and Approximation Algorithm .....	159
	References .....	161
<b>6</b>	<b>Multivariate Statistical Methods in Bioinformatics Research .....</b>	<b>163</b>
	Lingsong Zhang and Xihong Lin	
6.1	Introduction .....	163
6.2	Multivariate Normal Distribution .....	163
6.2.1	Definition and Notation .....	163
6.2.2	Properties of the Multivariate Normal Distribution .....	164
6.2.3	Bivariate Normal Distribution .....	165
6.2.4	Wishart Distribution .....	167
6.2.5	Sample Mean and Covariance .....	167
6.3	One-Sample and Two-Sample Multivariate Hypothesis Tests .....	168
6.3.1	One-Sample $t$ Test for a Univariate Outcome .....	168
6.3.2	Hotelling's $T^2$ Test for the Multivariate Outcome .....	169
6.3.3	Properties of Hotelling's $T^2$ Test .....	170
6.3.4	Paired Multivariate Hotelling's $T^2$ Test .....	171
6.3.5	Examples .....	172
6.3.6	Two-Sample Hotelling's $T^2$ Test .....	174
6.4	Principal Component Analysis .....	178
6.4.1	Definition of Principal Components .....	178
6.4.2	Computing Principal Components .....	179
6.4.3	Variance Decomposition .....	179
6.4.4	PCA with a Correlation Matrix .....	180
6.4.5	Geometric Interpretation .....	181
6.4.6	Choosing the Number of Principal Components .....	183
6.4.7	Diabetes Microarray Data .....	184
6.5	Factor Analysis .....	187
6.5.1	Orthogonal Factor Model .....	187
6.5.2	Estimating the Parameters .....	188
6.5.3	An Example .....	190

6.6	Linear Discriminant Analysis .....	193
6.6.1	Two-Group Linear Discriminant Analysis .....	194
6.6.2	An Example .....	198
6.7	Classification Methods .....	200
6.7.1	Introduction of Classification Methods .....	200
6.7.2	$k$ -Nearest Neighbor Method .....	202
6.7.3	Density-Based Classification Decision Rule .....	205
6.7.4	Quadratic Discriminant Analysis .....	208
6.7.5	Logistic Regression .....	212
6.7.6	Support Vector Machine .....	214
6.8	Variable Selection .....	219
6.8.1	Linear Regression Model .....	220
6.8.2	Motivation for Variable Selection .....	221
6.8.3	Traditional Variable Selection Methods .....	222
6.8.4	Regularization and Variable Selection .....	223
6.8.5	Summary .....	231
	References .....	231
7	<b>Association Analysis for Human Diseases: Methods and Examples</b> .....	233
	Jurg Ott and Qingrun Zhang	
7.1	Why Do We Need Statistics? .....	233
7.2	Basic Concepts in Population and Quantitative Genetics .....	234
7.3	Genetic Linkage Analysis .....	236
7.4	Genetic Case-Control Association Analysis .....	237
7.4.1	Basic Steps in an Association Study .....	238
7.4.2	Multiple Testing Corrections .....	239
7.4.3	Multi-locus Approaches .....	241
7.5	Discussion .....	241
	References .....	241
8	<b>Data Mining and Knowledge Discovery Methods with Case Examples</b> .....	243
	S. Bandyopadhyay and U. Maulik	
8.1	Introduction .....	243
8.2	Different Tasks in Data Mining .....	245
8.2.1	Classification .....	245
8.2.2	Clustering .....	248
8.2.3	Discovering Associations .....	252
8.2.4	Issues and Challenges in Data Mining .....	254
8.3	Some Common Tools and Techniques .....	256
8.3.1	Artificial Neural Networks .....	256
8.3.2	Fuzzy Sets and Fuzzy Logic .....	258
8.3.3	Genetic Algorithms .....	258

8.4	Case Examples .....	259
8.4.1	Pixel Classification .....	260
8.4.2	Clustering of Satellite Images .....	262
8.5	Discussion and Conclusions .....	267
	References .....	267
<b>9</b>	<b>Applied Bioinformatics Tools .....</b>	<b>271</b>
	Jingchu Luo	
9.1	Introduction .....	271
9.1.1	Welcome .....	271
9.1.2	About This Web Site .....	273
9.1.3	Outline .....	274
9.1.4	Lectures .....	275
9.1.5	Exercises .....	276
9.2	Entrez .....	277
9.2.1	PubMed Query .....	277
9.2.2	Entrez Query .....	278
9.2.3	My NCBI .....	278
9.3	ExPASy .....	278
9.3.1	Swiss-Prot Query .....	278
9.3.2	Explore the Swiss-Prot Entry HBA_HUMAN .....	279
9.3.3	Database Query with the EBI SRS .....	279
9.4	Sequence Alignment .....	280
9.4.1	Pairwise Sequence Alignment .....	280
9.4.2	Multiple Sequence Alignment .....	281
9.4.3	BLAST .....	281
9.5	DNA Sequence Analysis .....	282
9.5.1	Gene Structure Analysis and Prediction .....	282
9.5.2	Sequence Composition .....	283
9.5.3	Secondary Structure .....	283
9.6	Protein Sequence Analysis .....	283
9.6.1	Primary Structure .....	283
9.6.2	Secondary Structure .....	283
9.6.3	Transmembrane Helices .....	284
9.6.4	Helical Wheel .....	284
9.7	Motif Search .....	284
9.7.1	SMART Search .....	284
9.7.2	MEME Search .....	284
9.7.3	HMM Search .....	285
9.7.4	Sequence Logo .....	285
9.8	Phylogeny .....	285
9.8.1	Protein .....	285
9.8.2	DNA .....	286

9.9	Projects .....	286
9.9.1	Sequence, Structure, and Function Analysis of the Bar-Headed Goose Hemoglobin.....	286
9.9.2	Exercises .....	287
9.10	Literature .....	287
9.10.1	Courses and Tutorials .....	287
9.10.2	Scientific Stories .....	288
9.10.3	Free Journals and Books .....	288
9.11	Bioinformatics Databases .....	289
9.11.1	List of Databases .....	289
9.11.2	Database Query Systems.....	289
9.11.3	Genome Databases .....	290
9.11.4	Sequence Databases .....	291
9.11.5	Protein Domain, Family, and Function Databases .....	292
9.11.6	Structure Databases .....	293
9.12	Bioinformatics Tools.....	294
9.12.1	List of Bioinformatics Tools at International Bioinformatics Centers.....	295
9.12.2	Web-Based Bioinformatics Platforms .....	295
9.12.3	Bioinformatics Packages to be Downloaded and Installed Locally .....	295
9.13	Sequence Analysis .....	296
9.13.1	Dotplot.....	296
9.13.2	Pairwise Sequence Alignment .....	296
9.13.3	Multiple Sequence Alignment .....	296
9.13.4	Motif Finding .....	297
9.13.5	Gene Identification .....	297
9.13.6	Sequence Logo .....	297
9.13.7	RNA Secondary Structure Prediction .....	297
9.14	Database Search.....	298
9.14.1	BLAST Search .....	298
9.14.2	Other Database Search .....	298
9.15	Molecular Modeling .....	299
9.15.1	Visualization and Modeling Tools .....	299
9.15.2	Protein Modeling Web Servers .....	300
9.16	Phylogenetic Analysis and Tree Construction .....	300
9.16.1	List of Phylogeny Programs .....	300
9.16.2	Online Phylogeny Servers .....	300
9.16.3	Phylogeny Programs .....	301
9.16.4	Display of Phylogenetic Trees .....	301
	References .....	301

<b>10 Foundations for the Study of Structure and Function of Proteins</b>	....	303
Zhirong Sun		
10.1	Introduction .....	303
10.1.1	Importance of Protein .....	303
10.1.2	Amino Acids, Peptides, and Proteins.....	304
10.1.3	Some Noticeable Problems .....	306
10.2	Basic Concept of Protein Structure .....	306
10.2.1	Different Levels of Protein Structures.....	306
10.2.2	Acting Force to Sustain and Stabilize the High-Dimensional Structure of Protein .....	308
10.3	Fundamental of Macromolecules Structures and Functions .....	310
10.3.1	Different Levels of Protein Structure.....	310
10.3.2	Primary Structure.....	311
10.3.3	Secondary Structure.....	312
10.3.4	Supersecondary Structure .....	314
10.3.5	Folds .....	319
10.3.6	Summary .....	321
10.4	Basis of Protein Structure and Function Prediction .....	322
10.4.1	Overview .....	322
10.4.2	The Significance of Protein Structure Prediction .....	322
10.4.3	The Field of Machine Learning.....	323
10.4.4	Homological Protein Structure Prediction Method .....	331
10.4.5	Ab Initio Prediction Method .....	334
Reference .....	.....	336
<b>11 Computational Systems Biology Approaches for Deciphering Traditional Chinese Medicine</b>	.....	337
Shao Li and Le Lu		
11.1	Introduction .....	337
11.2	Disease-Related Network .....	338
11.2.1	From a Gene List to Pathway and Network .....	338
11.2.2	Construction of Disease-Related Network .....	340
11.2.3	Biological Network Modularity and Phenotype Network .....	346
11.3	TCM ZHENG-Related Network .....	349
11.3.1	“ZHENG” in TCM .....	350
11.3.2	A CSB-Based Case Study for TCM ZHENG .....	352
11.4	Network-Based Study for TCM “Fu Fang” .....	358
11.4.1	Systems Biology in Drug Discovery .....	358
11.4.2	Network-Based Drug Design .....	359
11.4.3	Progresses in Herbal Medicine .....	360
11.4.4	TCM <i>Fu Fang</i> (Herbal Formula) .....	361
11.4.5	A Network-Based Case Study for TCM <i>Fu Fang</i> .....	361
References .....	.....	364

<b>12 Advanced Topics in Bioinformatics and Computational Biology .....</b>	<b>369</b>
Bailin Hao, Chunting Zhang, Yixue Li, Hao Li, Liping Wei, Minoru Kanehisa, Luhua Lai, Runsheng Chen, Nikolaus Rajewsky, Michael Q. Zhang, Jingdong Han, Rui Jiang, Xuegong Zhang, and Yanda Li	
12.1 Prokaryote Phylogeny Meets Taxonomy .....	369
12.2 Z-Curve Method and Its Applications in Analyzing Eukaryotic and Prokaryotic Genomes .....	372
12.3 Insights into the Coupling of Duplication Events and Macroevolution from an Age Profile of Transmembrane Gene Families .....	374
12.4 Evolution of Combinatorial Transcriptional Circuits in the Fungal Lineage .....	375
12.5 Can a Non-synonymous Single-Nucleotide Polymorphism (nsSNP) Affect Protein Function? Analysis from Sequence, Structure, and Enzymatic Assay .....	377
12.6 Bioinformatics Methods to Integrate Genomic and Chemical Information .....	379
12.7 From Structure-Based to System-Based Drug Design .....	381
12.8 Progress in the Study of Noncoding RNAs in <i>C. elegans</i> .....	383
12.9 Identifying MicroRNAs and Their Targets .....	385
12.10 Topics in Computational Epigenomics .....	387
12.11 Understanding Biological Functions Through Molecular Networks .....	389
12.12 Identification of Network Motifs in Random Networks .....	390
12.13 Examples of Pattern Recognition Applications in Bioinformatics .....	392
12.14 Considerations in Bioinformatics .....	394
<b>Erratum .....</b>	<b>E1</b>

# Chapter 1

## Basics for Bioinformatics

Xuegong Zhang, Xueya Zhou, and Xiaowo Wang

### 1.1 What Is Bioinformatics

Bioinformatics has become a hot research topic in recent years, a hot topic in several disciplines that were not so closely linked with biology previously. A side evidence of this is the fact that the 2007 Graduate Summer School on Bioinformatics of China had received more than 800 applications from graduate students from all over the nation and from a wide collection of disciplines in biological sciences, mathematics and statistics, automation and electrical engineering, computer science and engineering, medical sciences, environmental sciences, and even social sciences. So what is bioinformatics?

It is always challenging to define a new term, especially a term like bioinformatics that has many meanings. As an emerging discipline, it covers a lot of topics from the storage of DNA data and the mathematical modeling of biological sequences, to the analysis of possible mechanisms behind complex human diseases, to the understanding and modeling of the evolutionary history of life, etc.

Another term that often goes together or close with bioinformatics is computational molecular biology, and also computational systems biology in recent years, or computational biology as a more general term. People sometimes use these terms to mean different things, but sometimes use them in exchangeable manners. In our personal understanding, computational biology is a broad term, which covers all efforts of scientific investigations on or related with biology that involve mathematics and computation. Computational molecular biology, on the other hand, concentrates on the molecular aspects of biology in computational biology, which therefore has more or less the same meaning with bioinformatics.

---

X. Zhang (✉) • X. Zhou • X. Wang

MOE Key Laboratory of Bioinformatics and Bioinformatics Division, TNLIST/Department of Automation, Tsinghua University, Beijing 100084, China

e-mail: [zhangxg@tsinghua.edu](mailto:zhangxg@tsinghua.edu)

Bioinformatics studies the storage, manipulation, and interpretation of biological data, especially data of nucleic acids and amino acids, and studies molecular rules and systems that govern or affect the structure, function, and evolution of various forms of life from computational approaches. The word “computational” does not only mean “with computers,” but it refers to data analysis with mathematical, statistical, and algorithmic methods, most of which need to be implemented with computer programs. As computational biology or bioinformatics studies biology with quantitative data, people also call it as quantitative biology.

Most molecules do not work independently in living cells, and most biological functions are accomplished by the harmonic interaction of multiple molecules. In recent years, the new term systems biology came into being. Systems biology studies cells and organisms as systems of multiple molecules and their interactions with the environment. Bioinformatics plays key roles in analyzing such systems. People have invented the term computational systems biology, which, from a general viewpoint, can be seen as a branch of bioinformatics that focuses more on systems rather than individual elements.

For a certain period, people regarded bioinformatics as the development of software tools that help to store, manipulate, and analyze biological data. While this is still an important role of bioinformatics, more and more scientists realize that bioinformatics can and should do more. As the advancement of modern biochemistry, biophysics, and biotechnologies is enabling people to accumulate massive data of multiple aspects of biology in an exponential manner, scientists begin to believe that bioinformatics and computational biology must play a key role for understanding biology.

People are studying bioinformatics in different ways. Some people are devoted to developing new computational tools, both from software and hardware viewpoints, for the better handling and processing of biological data. They develop new models and new algorithms for existing questions and propose and tackle new questions when new experimental techniques bring in new data. Other people take the study of bioinformatics as the study of biology with the viewpoint of informatics and systems. These people also develop tools when needed, but they are more interested in understanding biological procedures and mechanisms. They do not restrict their research to computational study, but try to integrate computational and experimental investigations.

## 1.2 Some Basic Biology

No matter what type of bioinformatics one is interested in, basic understanding of existing knowledge of biology especially molecular biology is a must. This chapter was designed as the first course in the summer school to provide students with non-biology backgrounds very basic and abstractive understanding of molecular biology. It can also give biology students a clue how biology is understood by researchers from other disciplines, which may help them to better communicate with bioinformaticians.

### 1.2.1 Scale and Time

Biology is the science about things that live in nature. There are many forms of life on the earth. Some forms are visible to human naked eyes, like animals and plants. Some can only be observed under light microscope or electron microscope, like many types of cells in the scale of  $1.100\text{ }\mu\text{m}$  and some virus in the scale of  $100\text{ nm}$ . The basic components of those life forms are molecules of various types, which scale around  $1.10\text{ nm}$ . Because of the difficulty of direct observation at those tiny scales, scientists have to invent various types of techniques that can measure some aspects of the molecules and cells. These techniques produce a large amount of data, from which biologists and bioinformaticians infer the complex mechanisms underlying various life procedures.

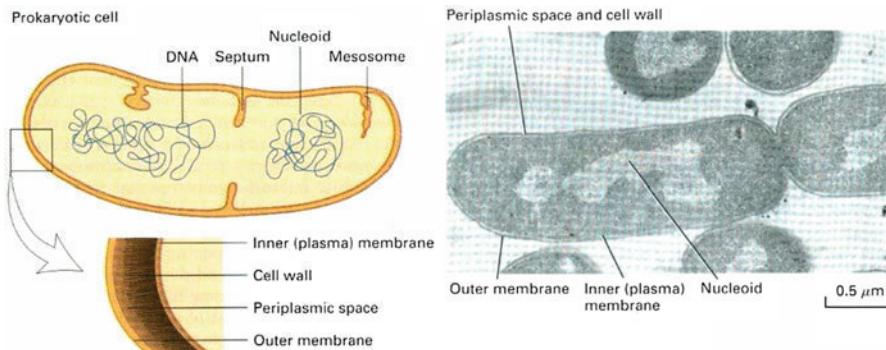
Life has a long history. The earliest form of life appeared on the earth about 4 billion years ago, not long after the forming of the earth. Since then, life has experienced a long way of evolution to reach today's variety and complexity. If the entire history of the earth is scaled to a 30-day month, the origin of life happened during days 3–4, but there has been abundant life only since day 27. A lot of higher organisms appeared in the last few days: first land plants and first land animals all appeared on day 28, mammals began to exist on day 29, and birds and flowering plants came into being on the last day. Modern humans, which are named homo sapiens in biology, appeared in the last 10 min of the last day. If we consider the recorded human history, it takes up only the last 30 s of the last day. The process that life gradually changes into different and often more complex or higher forms is called evolution. When studying the biology of a particular organism, it is important to realize that it is one leaf or branch on the huge tree of evolution. Comparison between related species is one major approach when investigating the unknown.

### 1.2.2 Cells

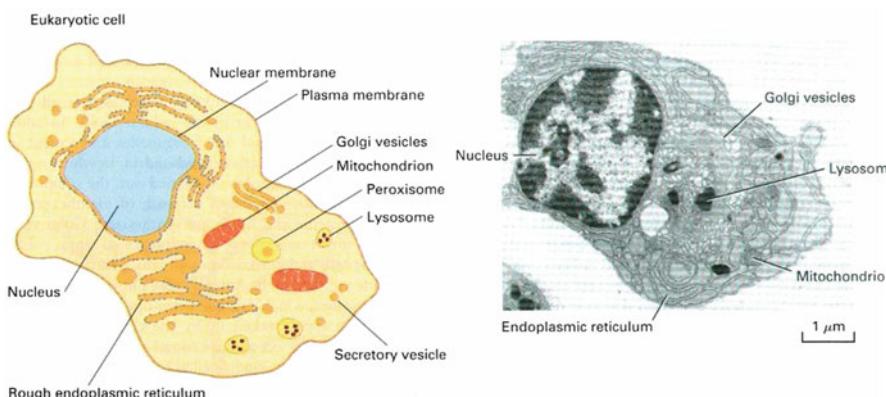
The basic component of all organisms is the cell. Many organisms are unicellular, which means one cell itself is an organism. However, for higher species like animals and plants, an organism can contain thousands of billions of cells.

Cells are of two major types: prokaryotic cells and eukaryotic cells. Eukaryotic cells are cells with real nucleus, while prokaryotic cells do not have nucleus. Living organisms are also categorized as two major groups: prokaryotes and eukaryotes according to whether their cells have nucleus. Prokaryotes are the earlier forms of life on the earth, which includes bacteria and archaea. All higher organisms are eukaryote, including unicellular organisms like yeasts and higher organisms like plants and animals. The bacteria *E. coli* is a widely studied prokaryote. Figure 1.1 shows the structure of an *E. coli* cell, as a representative of prokaryotic cells.

Eukaryotic cells have more complex structures, as shown in the example of a human plasma cell in Fig. 1.2. In eukaryotic cells, the key genetic materials, DNA, live in nucleus, in the form of chromatin or chromosomes. When a cell is



**Fig. 1.1** A prokaryotic cell



**Fig. 1.2** An eukaryotic cell

not dividing, the nuclear DNA and proteins are aggregated as chromatin, which is dispersed throughout the nucleus. The chromatin in a dividing cell is packed into dense bodies called chromosomes. Chromosomes are of two parts, called the P-arm and Q-arm, or the shorter arm and longer arm, separated by the centromere.

### 1.2.3 DNA and Chromosome

DNA is the short name for deoxyribonucleic acid, which is the molecule that stores the major genetic information in cells. A nucleotide consists of three parts: a phosphate group, a pentose sugar (ribose sugar), and a base. The bases are of four types: adenine (A), guanine (G), cytosine (C), and thymine (T). A and G are purines with two fused rings. C and T are pyrimidines with one single ring. Besides DNA,

5'-A T T A C G G T A C C G T -3'  
3'-T A A T G C C A T G G C A -5'

**Fig. 1.3** An example segment of a double-strand DNA sequence

there is another type of nucleotide called RNA or ribonucleic acid. For RNA, the bases are also of these four types except that the T is replaced by the uracil (U) in RNA.

DNA usually consists of two strands running in opposite directions. The backbone of each strand is a series of pentose and phosphate groups. Hydrogen bonds between purines and pyrimidines hold the two strands of DNA together, forming the famous double helix. In the hydrogen bonds, a base A always pairs with a base T on the other stand and a G always with a C. This mechanism is called base pairing. RNA is usually a single strand. When an RNA strand pairs with a DNA strand, the base-pairing rule becomes A-U, T-A, G-C, and C-G.

The ribose sugar is called pentose sugar because it contains five carbons, numbered as 1'-5', respectively. The definition of the direction of a DNA or RNA strand is also based on this numbering, so that the two ends of a DNA or RNA strand are called the 5' end and the 3' end. The series of bases along the strand is called the DNA or RNA sequence and can be viewed as character strings composed with the alphabet of "A," "C," "G," and "T" ("U" for RNA). We always read a sequence from the 5' end to the 3' end. On a DNA double helix, the two strands run oppositely. Figure 1.3 is an example of a segment of double-strand DNA sequence. Because of the DNA base-pairing rule, we only need to save one strand of the sequence.

DNA molecules have very complicated structures. A DNA molecule binds with histones to form a vast number of nucleosomes, which look like "beads" on DNA "string." Nucleosomes pack into a coil that twists into another larger coil and so forth, producing condensed supercoiled chromatin fibers. The coils fold to form loops, which coil further to form a chromosome. The length of all the DNA in a single human cell is about 2 m, but with the complicated packing, they fit into the nucleus with diameter around 5  $\mu\text{m}$ .

#### ***1.2.4 The Central Dogma***

The central dogma in genetics describes the typical mechanism by which the information saved in DNA sequences fulfills its job: information coded in DNA sequence is passed on to a type of RNA called messenger RNA (mRNA). Information in mRNA is then passed on to proteins. The former step is called transcription, and latter step is called translation.

Transcription is governed by the rule of complementary base pairing between the DNA base and the transcribed RNA base. That is, an A in the DNA is transcribed to a U in the RNA, a T to an A, a G to a C, and vice versa.

		Second letter								
		U	C	A	G					
First letter	U	UUU UUC UUA UUG	Phenylalanine Serine	UAU UAC UAA UAG	Tyrosine Stop codon Stop codon	UGU UGC UGA UGG	Cysteine Stop codon Tryptophan	U C A G		
	C	CUU CUC CUA CUG	Leucine	CCU CCC CCA CCG	Proline	CAU CAC CAA CAG	Histidine Glutamine	CGU CGC CGA CGG	Arginine	U C A G
	A	AUU AUC AUA AUG	Isoleucine Methionine; start codon	ACU ACC ACA ACG	Threonine	AAU AAC AAA AAG	Asparagine Lysine	AGU AGC AGA AGG	Serine Arginine	U C A G
	G	GUU GUC GUA GUG	Valine	GCU GCC GCA GGC	Alanine	GAU GAC GAA GAG	Aspartic acid Glutamic acid	GGU GGC GGA GGG	Glycine	U C A G

Fig. 1.4 The genetic codes

Proteins are chains of amino acids. There are 20 types of standard amino acids used in lives. The procedure of translation converts the information from the language of nucleotides to the language of amino acids. The translation is done by a special dictionary: the genetic codes or codon. Figure 1.4 shows the codon table. Every three nucleotides code for one particular amino acid. The three nucleotides are called a triplet. Because three nucleotides can encode 64 unique items, there are redundancies in this coding scheme, as shown in Fig. 1.4. Many amino acids are coded by more than one codon. For the redundant codons, usually their first and second nucleotides are consistent, but some variation in the third nucleotide is tolerated. AUG is the start codon that starts a protein sequence, and there are three stop codons CAA, CAG, and UGA that stop the sequence.

Figure 1.5a illustrates the central dogma in prokaryotes. First, DNA double helix is opened and one strand of the double helix is used as a template to transcribe the mRNA. The mRNA is then translated to protein in ribosome with the help of tRNAs.

Figure 1.5b illustrates the central dogma in eukaryotes. There are several differences with the prokaryote case. In eukaryotic cells, DNAs live in the nucleus, where they are transcribed to mRNA similar to the prokaryote case. However, this mRNA is only the preliminary form of message RNA or pre-mRNA. Pre-mRNA is processed in several steps: parts are removed (called splicing), and ends of 150–200 As (called poly-A tail) are added. The processed mRNA is exported outside the nucleus to the cytoplasm, where it is translated to protein.

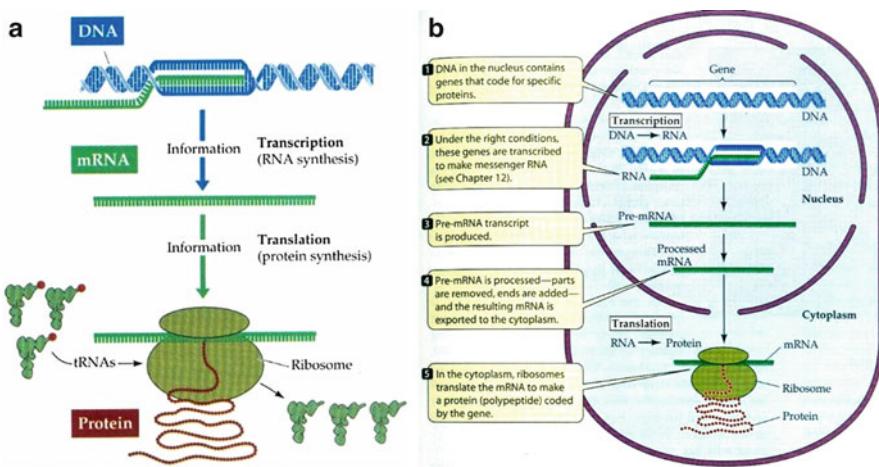


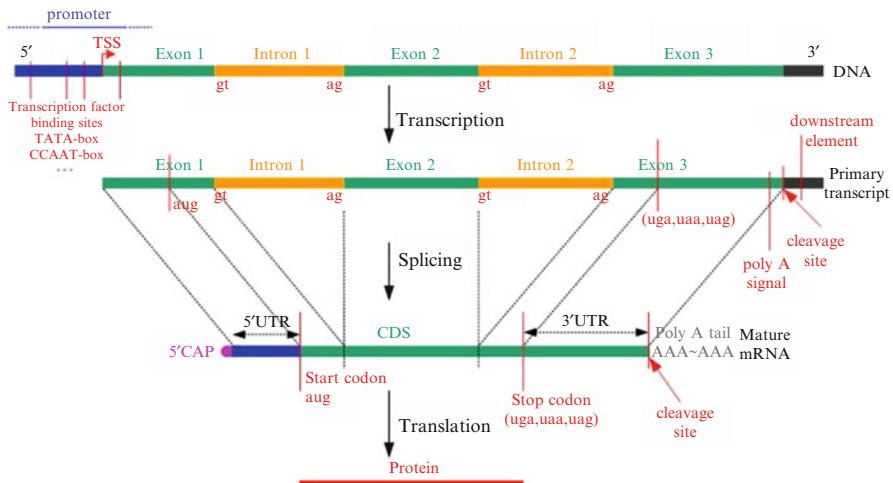
Fig. 1.5 The central dogma

The procedure that genes are transcribed to mRNAs which are then translated to proteins is called the expression of genes. And the abundance of the mRNA molecules of a gene is usually called the expression value (level) of that gene, or simply the expression of the gene.

### 1.2.5 Genes and the Genome

We believe that the Chinese translation “基因” of the term “gene” is one of the best scientific term ever translated. Besides that the pronunciation is very close to the English version, the literal meaning of the two characters is also very close to the definition of the term: basic elements. Genes are the basic genetic elements that, together with interaction with environment, are decisive for the phenotypes.

Armed with knowledge of central dogma and genetic code, people had long taken the concept of a gene as the fragments of the DNA sequence that finally produce some protein products. This is still true in many contexts today. More strictly, these DNA segments should be called protein-coding genes, as scientists have found that there are some or many other parts on the genome that do not involve in protein products but also play important genetic roles. Some people call them as nonprotein-coding genes or noncoding genes for short. One important type of noncoding genes is the so-called microRNAs or miRNAs. There are several other types of known noncoding genes and may be more unknown. In most current literature, people still use gene to refer to protein-coding genes and add attributes like “noncoding” and “miRNA” when referring to other types of genes. We also follow this convention in this chapter.



**Fig. 1.6** The structure of a gene

The length of a DNA segment is often counted by the number of nucleotides (nt) in the segment. Because DNAs usually stay as double helix, we can also use the number of base pairs (bp) as the measurement of the length. For convenience, people usually use “k” to represent “1,000.” For example, 10 kb means that the sequence is of 10,000 bp. A protein-coding gene stretches from several hundreds of bp to several k bp in the DNA sequence. Figure 1.6 shows an example structure of a gene in high eukaryotes.

The site on the DNA sequence where a gene is started to be transcribed is called the transcription start site or TSS. The sequences around (especially the upstream) the TSS contain several elements that play important roles in the regulation of the transcription. These elements are called cis-elements. Transcription factors bind to such factors to start, enhance, or repress the transcription procedure. Therefore, sequences upstream the TSS are called promoters. Promoter is a loosely defined concept, and it can be divided into three parts: (1) a core promoter which is about 100 bp long around the TSS containing binding sites for RNA polymerase II (Pol II) and general transcription factors, (2) a proximal promoter of several hundred base pairs long containing primary specific regulatory elements located at the immediately upstream of the core promoter, and (3) a distal promoter up to thousands of base pairs long providing additional regulatory information. In eukaryotes, the preliminary transcript of a gene undergoes a processing step called splicing, during which some parts are cut off and remaining parts are joined. The remaining part is called exon, and the cut part is called intron. There can be multiple exons and introns in a gene. After introns are removed, the exons are connected to form the processed mRNA. Only the processed mRNAs are exported to the cytoplasm, and only parts of the mRNAs are translated to proteins. There may be

untranslated regions (UTRs) at both ends of the mRNA: one at the TSS end is called 5'-UTR, and the other at the tail end is called 3'-UTR. The parts of exons that are translated are called CDS or coding DNA sequences. Usually exons constitute only a small part in the sequence of a gene.

In higher eukaryotes, a single gene can have more than one exon-intron settings. Such genes will have multiple forms of protein products (called isoforms). One isoform may contain only parts of the exons, and the stretch of some exons may also differ among isoforms. This phenomenon is called alternative splicing. It is an important mechanism to increase the diversity of protein products without increasing the number of genes.

The term “genome” literally means the set of all genes of an organism. For prokaryotes and some low eukaryotes, majority of their genome is composed of protein-coding genes. However, as more and more knowledge about genes and DNA sequences in human and other high eukaryotes became available, people learned that protein-coding genes only take a small proportion of all the DNA sequences in the eukaryotic genome. Now people tend to use “genome” to refer all the DNA sequences of an organism or a cell. (The genomes of most cell types in an organism are the same.)

The human genome is arranged in 24 chromosomes, with the total length of about 3 billion base pairs ( $3 \times 10^9$  bp). There are 22 autosomes (Chr.1-22) and 2 sex chromosomes (X and Y). The 22 autosomes are ordered by their lengths (with the exception that Chr.21 is slightly shorter than Chr.22): Chr.1 is the longest chromosome and Chr.21 is the shortest autosome. A normal human somatic cell contains 23 pairs of chromosomes: two copies of chromosomes 1-22 and two copies of X chromosome in females or one copy of X and one copy of Y in males. The largest human chromosome (Chr.1) has about 250 million bp, and the smallest human chromosome (Chr.Y) has about 50 million bp.

There are about 20,000–25,000 protein-coding genes in the human genome, spanning about 1/3 of the genome. The average human gene consists of some 3,000 base pairs, but sizes vary greatly, from several hundred to several million base pairs. The protein-coding part only takes about 1.5–2 % of the whole genome. Besides these regions, there are regulatory sequences like promoters, intronic sequences, and intergenic (between-gene) regions. Recent high-throughput transcriptomic (the study of all RNA transcripts) study revealed that more than half of the human genomes are transcribed, although only a very small part of them are processed to mature mRNAs. Among the transcripts are the well-known microRNAs and some other types of noncoding RNAs. The functional roles played by majority of the transcripts are still largely unknown. There are many repetitive sequences in the genome, and they have not been observed to have direct functions.

Human is regarded as the most advanced form of life on the earth, but the human genome is not the largest. Bacteria like *E. coli* has genomes of several million bp, yeast has about 15 million bp, *Drosophila* (fruit fly) has about 3 million bp, and some plants can have genomes as large as 100 billion bp. The number of genes in a genome is also not directly correlated with the complexity of the organism’s

complexity. The unicellular organism yeast has about 6,000 genes, fruit fly has about 15,000 genes, and the rice that we eat everyday has about 40,000 genes. In lower species, protein-coding genes are more densely distributed on the genome. The human genome also has a much greater portion (50 %) of repeat sequences than the lower organisms like the worm (7 %) and the fly (3 %).

### ***1.2.6 Measurements Along the Central Dogma***

For many years, molecular biology can only study one or a small number of objects (genes, mRNAs, or proteins) at a time. This picture was changed since the development of a series of high-throughput technologies. They are called high throughput because they can obtain measurement of thousands of objects in one experiment in a short time. The emergence of massive genomic and proteomic data generated with these high-throughput technologies was actually a major motivation that promotes the birth and development of bioinformatics as a scientific discipline. In some sense, what bioinformatics does is manipulating and analyzing massive biological data and aiding scientific reasoning based on such data. It is therefore crucial to have the basic understanding of how the data are generated and what the data are for.

### ***1.2.7 DNA Sequencing***

The sequencing reaction is a key technique that enables the completion of sequencing the human genome. Figure 1.7 illustrates the principle of the widely used Sanger sequencing technique.

The technique is based on the complementary base-pairing property of DNA. When a single-strand DNA fragment is isolated and places with primers, DNA polymerase, and the four types of deoxyribonucleoside triphosphate (dNTP), a new DNA strand complementary to the existing one will be synthesized. In the DNA sequencing reaction, dideoxyribonucleoside triphosphate (ddNTP) is added besides the above components, and the four types of ddNTPs are bound to four different fluorescent dyes. The synthesis of a new strand will stop when a ddNTP instead of a dNTP is added. Therefore, with abundant template single-strand DNA fragments, we'll be able to get a set of complementary DNA segments of all different lengths, each one stopped by a colored ddNTP. Under electrophoresis, these segments of different lengths will run at different speeds, with the shortest segments running the fastest and the longest segments running the slowest. By scanning the color of all segments ordered by their length, we'll be able to read the nucleotide at each position of the complementary sequence and therefore read the original template sequence. This technique is implemented in the first generation of sequencing machines.