

# The Human Genome in Health and Disease A Story of Four Letters



# The Human Genome in Health and Disease A Story of Four Letters

To my parents

# The Human Genome in Health and Disease A Story of Four Letters

**Tore Samuelsson** 



CRC Press Taylor & Francis Group 6000 Broken Sound Parkway NW, Suite 300 Boca Raton, FL 33487-2742

© 2019 by Taylor & Francis Group, LLC CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper

International Standard Book Number-13: 978-0-367-07633-7 (Hardback) 978-0-8153-4591-6 (Paperback)

This book contains information obtained from authentic and highly regarded sources. While all reasonable efforts have been made to publish reliable data and information, neither the author[s] nor the publisher can accept any legal responsibility or liability for any errors or omissions that may be made. The publishers wish to make clear that any views or opinions expressed in this book by individual editors, authors or contributors are personal to them and do not necessarily reflect the views/opinions of the publishers. The information or guidance contained in this book is intended for use by medical, scientific or health-care professionals and is provided strictly as a supplement to the medical or other professional's own judgement, their knowledge of the patient's medical history, relevant manufacturer's instructions and the appropriate best practice guidelines. Because of the rapid advances in medical science, any information or advice on dosages, procedures or diagnoses should be independently verified. The reader is strongly urged to consult the relevant national drug formulary and the drug companies' and device or material manufacturers' printed instructions, and their websites, before administering or utilizing any of the drugs, devices or materials mentioned in this book. This book does not indicate whether a particular treatment is appropriate or suitable for a particular individual. Ultimately it is the sole responsibility of the medical professional to make his or her own professional judgements, so as to advise and treat patients appropriately. The authors and publishers have also attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (http://www.copyright.com/) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

#### Library of Congress Cataloging-in-Publication Data

Names: Samuelsson, Tore, 1951- author. Title: The human genome in health and disease : a story of four letters / Tore Samuelsson. Description: Boca Raton : Taylor & Francis, 2019. | Includes bibliographical references and index. Identifiers: LCCN 2018040829! ISBN 9780815345916 (pbk. : alk. paper) | ISBN 9780367076337 (hardback : alk. paper) | ISBN 9780429021732 (e-ISBN) Subjects: | MESH: Genome, Human | Disease--genetics Classification: LCC QH447 | NLM QU 460 | DDC 611/.0181663--dc23 LC record available at https://lccn.loc.gov/2018040829

# Visit the Taylor & Francis Web site at http://www.taylorandfrancis.com

## and the CRC Press Web site at http://www.crcpress.com

eResource material is available for this title at https://www.crcpress.com/ 9780367076337

# Contents

#### Preface xi THE GENETIC CODE IS BY AND LARGE **UNIVERSAL** 24 Chapter 1 A SECOND GENETIC CODE 25 Introduction 1 FLOW OF INFORMATION FROM DNA TO **PROTEIN** 25 Chapter 2 SICKLE CELL ANEMIA IS THE RESULT OF A Molecular Disorder 5 A SINGLE NUCLEOTIDE CHANGE IN DNA 26 THE DISCOVERY OF SICKLE CELL ANEMIA 5 27 **Summary** 6 A RECESSIVE INHERITED DISORDER Questions 27 SICKLE CELL ANEMIA AND MALARIA 8 **Further Reading** 28 CHARACTERIZING SICKLE CELL ANEMIA 8 FURTHER: THE ROLE OF HEMOGLOBIN Chapter 4 MAPPING A DELETERIOUS CHANGE The Genome 31 9 IN HEMOGLOBIN DNA IS PRESENT IN THE NUCLEUS A PROTEIN FOLDS INTO A THREE-**OF CELLS** 31 DIMENSIONAL SHAPE BASED ON ITS 9 AMINO ACID SEQUENCE DNA AS SEEN UNDER THE MICROSCOPE 32 SICKLE CELL DISEASE MAY NOW BE DNA IS COMPACTED WITHIN THE UNDERSTOOD IN THE CONTEXT OF NUCLEUS OF A HUMAN CELL 34 10 PROTEIN STRUCTURE TOOLBOX OF THE MOLECULAR A MOLECULAR SEQUENCE: PROTEINS 34 BIOLOGIST CAN BE DEPICTED AS A STRING OF DEVELOPMENT OF DNA SEQUENCING 35 AMINO ACID SYMBOLS 12 HUMAN GENOME SEQUENCING PROJECT 37 14 **Summary** WHOSE GENOME? 40 14 Questions THE MITOCHONDRIAL GENOME 43 **Further Reading** 15 ACCESSING THE GENOME—GENOME Chapter 3 **BROWSERS** 44 A Code of Life 17 THE HUMAN GENOME SPECIFIES PROTEINS AND NONCODING RNAs 17 45 EARLY DAYS OF DNA RESEARCH MOST GENES ARE MOSAICS OF EXONS DNA IS THE CARRIER OF GENETIC AND INTRONS 45 **INFORMATION** 17 **ONLY ABOUT 1.6% OF THE GENOME** THREE-DIMENSIONAL STRUCTURE OF DNA 18 **CODES FOR PROTEIN** 46 POLARITY OF DNA 18 **REPETITIVE SEQUENCES MAKE UP MORE** A MECHANISM OF REPLICATION WAS THAN HALF OF THE HUMAN GENOME **48** SUGGESTED BY THE STRUCTURE OF DNA 20 A COMPLEX PROTEIN UNIVERSE **50 RELATIONSHIP OF DNA AND PROTEIN** 20 Summary 51 **BASIC FEATURES OF THE GENETIC CODE** 21 Questions 52 MOLECULES INVOLVED IN PROTEIN URLS PRODUCTION 21 53 MAKING SENSE OF THE CODONS 23 **Further Reading** 53

Chapter 5 Variants in the Human Genome Sequence and Their Biological	
Significance	55
TYPES OF MUTATIONS	55
GERMLINE AND SOMATIC MUTATIONS	56
MUTATION RATE IN THE HUMAN GENOME	56
HOW DO MUTATIONS COME ABOUT?	57
DNA REPLICATION ERRORS	57
SPONTANEOUS DNA DAMAGE	57
DNA DAMAGE BY REACTIVE OXYGEN SPECIES AND BY IRRADIATION	58
HOW ARE DNA ERRORS REPAIRED?	59
HUMAN GENETIC VARIATION	<mark>62</mark>
INDIVIDUAL DIFFERENCES: SINGLE NUCLEOTIDE POLYMORPHISMS AND SINGLE NUCLEOTIDE VARIANTS	63
INDIVIDUAL DIFFERENCES: STRUCTURAL VARIATION	64
PHENOTYPIC IMPLICATIONS OF INDIVIDUAL VARIATION: SINGLE GENE DISORDERS AND COMPLEX DISEASES	66
PHENOTYPIC IMPLICATIONS OF SNPs MAY BE INFERRED FROM A GENOME- WIDE ASSOCIATION STUDY	67
A RESTRICTED NUMBER OF VARIANTS HAVE A KNOWN DISTINCT FUNCTION	69
CANCER: A GENETIC DISEASE	69
NUMEROUS CANCER GENOMES HAVE BEEN SEQUENCED	72
CANCER: TRANSLOCATIONS AND MORE DRAMATIC GENOMIC ABERRATIONS	73
WHAT ARE THE CAUSES OF CANCER?	74
Summary	75
Questions	76
URLs	76
Further Reading	77
Chapter 6 The Critical Protein Coding	70
Sequences	19
MUTATIONS IN CODING SEQUENCES	79
INHERITED SINGLE GENE DISORDERS ARE OFTEN CAUSED BY NONSYNONYMOUS MUTATIONS	80

MANY DIFFERENT MUTATIONS IN A GENE MAY GIVE RISE TO DISEASE	81
THE GENE <i>HEXA</i> AND TAY-SACHS DISEASE	82
THE GENE <i>HFE</i> AND HEREDITARY HEMOCHROMATOSIS	83
MUTATIONS THAT GENERATE PREMATURE STOP CODONS	85
NONSENSE MUTATIONS RESULT IN mRNA DEGRADATION	85
A NONSENSE MUTATION THAT GIVES RISE TO AGGRESSIVE BEHAVIOR	86
A SPEECH AND LANGUAGE DISORDER	88
A TERMINATION CODON CHANGED TO A SENSE CODON: A LESS COMMON MUTATION	89
INDEL MUTATIONS OF CODING SEQUENCES: TAY-SACHS DISEASE REVISITED	90
FRAMESHIFT ERROR AND RESISTANCE TO HIV INFECTION	90
AN INFRAME DELETION IS A COMMON CAUSE OF CYSTIC FIBROSIS	91
Summary	93
Questions	93
URLS	94
Further Reading	94
Chapter 7 Triplet Repeats and Neurodegenerative Disorders	97

Inplet Repeats and	
Neurodegenerative Disorders	97
THE DEATH OF PHEBE HEDGES	97
GEORGE HUNTINGTON	98
INHERITANCE OF HUNTINGTON'S DISEASE	98
A NEUROPSYCHIATRIC DISORDER	99
THE RESPONSIBLE GENE IS IDENTIFIED	100
HUNTINGTON'S DISEASE IS CAUSED BY AN EXPANSION OF TRIPLET	
REPEATS	102
THE STRUCTURE OF HUNTINGTIN	102
MOLECULAR FUNCTIONS OF	
HUNTINGTIN	104
DNA REPLICATION SLIPPAGE ERRORS	104
OTHER REPEAT DISORDERS	105

SPINOCEREBELLAR ATAXIAS CAUSED BY CAG	
TRIPLET EXPANSION	106
Summary	108
Questions	108
Further Reading	109
Chapter 8 The Untranslated Parts of a Message	111
ORGANIZATION OF mRNA AND FUNCTIONS WITHIN UTRs	111
IRON-RESPONSIVE ELEMENT REGULATES TRANSLATION AND mRNA STABILITY	112
STRUCTURE OF THE IRON RESPONSIVE ELEMENT	114
AN RNA ELEMENT NECESSARY FOR THE INCORPORATION OF SELENOCYSTEINE IN PROTEINS	116
REGULATION OF TRANSLATION BY 5' UTR OPEN READING FRAMES	117
WHEN UPSTREAM OPEN READING FRAMES ARE MUTATED	117
AU-RICH ELEMENTS ARE INVOLVED IN REGULATION OF mRNA STABILITY	118
A 3' UTR REPEAT EXPANSION AND MYOTONIC DYSTROPHY	119
A POSSIBLE MECHANISM OF MYOTONIC DYSTROPHY TYPE 1	120
A UNIVERSE OF SILENCING RNAs	121
DISEASES ASSOCIATED WITH MUTATIONS IN miRNAS OR TARGET	
3' UTRs	123
Summary	124
Questions	124
URLs	125
Further Reading	125
Chapter 9 Exons, Introns, and a Royal	
Bleeding Disorder	127
QUEEN VICTORIA PASSED HEMOPHILIA ON TO MANY OF HER DESCENDENTS	127
THE ROYAL HEMOPHILIA IS INHERITED THROUGH THE X CHROMOSOME	127
BLOOD CLOTTING DEPENDS ON A CASCADE OF ENZYMATIC REACTIONS	129

DETAILED INFORMATION OF THE ROYAL DISORDER WAS OBTAINED BY DNA ANALYSIS OF THE ROMANOVS	130
EUKARYOTIC GENES ARE MOSAICS OF EXONS AND INTRONS	131
THE ENDS OF INTRONS CONTAIN NUCLEOTIDE SEQUENCES CRUCIAL FOR SPLICING	133
EXONS MAY BE COMBINED IN ALTERNATIVE WAYS	134
DNA ANALYSIS SHOWED THAT THE ROYAL BLOOD DISORDER WAS DUE TO A SPLICING ERROR	135
A HEARING DISORDER PRESENTS ANOTHER CASE OF SPLICING DEFECT	137
A POINT MUTATION CAN HAVE DIFFERENT EFFECTS ON SPLICING	138
PREDICTING SPLICING MUTATIONS	139
Summary	140
Questions	140
URLs	141
Further Reading	141

#### Chapter 10 The Regulation of Transcription 143 DNA SEQUENCES AND GENERAL **TRANSCRIPTION FACTORS** 143 β-THALASSEMIA AND CORE PROMOTER **MUTATIONS** 145 DNA SEQUENCES AND SPECIFIC 146 **TRANSCRIPTION FACTORS** TRANSCRIPTION FACTOR DNA BINDING 146 **SPECIFICITY** THE ENCODE PROJECT MAPPED FUNCTIONAL SITES IN THE HUMAN **GENOME** 148 PROXIMAL PROMOTER MUTATIONS AND 149 DISEASE: GATA-1 AND δ-THALASSEMIA **PROXIMAL PROMOTER MUTATIONS** AND DISEASE: TERT PROMOTER AND CANCER 151 MUTATIONS IN ENHANCERS MAY GIVE **RISE TO DISEASE** 151 SICKLE CELL DISEASE AND ENHANCER **GENETIC VARIANTS** 152

AN ENHANCER GENETIC VARIANT IS ASSOCIATED WITH PARKINSON'S DISEASE	153
TRANSCRIPTIONAL CONTROL IS MEDIATED BY HIGHER-ORDER	154
CREANIZATION OF THE GENOME	154
ENHANCERS THAT GET OUT OF CONTEXT	150
CANCER AND DNA ORGANIZATION	156
START AND TERMINATION SITES	158
EPIGENETIC MECHANISMS OF TRANSCRIPTIONAL CONTROL	158
HISTONE MODIFICATION	159
DNA METHYLATION AND CpG ISLANDS	160
THE EPIGENOME	1 <b>62</b>
FRAGILE X SYNDROME: AN EPIGENETIC DISORDER?	162
Summary	164
Questions	165
URLs	166
Further Reading	166
Chapter 11 The Noncoding RNAs	169
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs	169 169
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD	169 169 170
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs	169 169 170 170
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES	169 169 170 170
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA	169 169 170 170 171
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL nCRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAs	<ul> <li>169</li> <li>169</li> <li>170</li> <li>170</li> <li>171</li> <li>172</li> <li>173</li> </ul>
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAs LONG NONCODING RNAs, A UNIVERSE OF RNAS YET TO BE EXPLORED	169 169 170 170 171 172 173 176
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL nCRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN TRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAs LONG NONCODING RNAS, A UNIVERSE OF RNAS YET TO BE EXPLORED	169 169 170 170 171 172 173 176 178
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAS LONG NONCODING RNAS, A UNIVERSE OF RNAS YET TO BE EXPLORED LncRNAS REGULATE EXPRESSION USING MANY DIFFERENT MECHANISMS	169 169 170 170 171 172 173 176 178 179
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN tRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAS LONG NONCODING RNAS, A UNIVERSE OF RNAS YET TO BE EXPLORED LncRNAS REGULATE EXPRESSION USING MANY DIFFERENT MECHANISMS A REPEAT EXPANSION IN A IncRNA CAUSES A SPINOCEREBELLAR ATAXIA A IncRNA AND CELIAC DISEASE	169 169 170 170 171 172 173 176 178 179 179
Chapter 11 The Noncoding RNAs SMALL NONCODING RNAs THE RNA WORLD BIOSYNTHESIS OF SMALL ncRNAs A NUMBER OF MITOCHONDRIAL DISEASES ARE THE RESULT OF MUTATIONS IN TRNA GENES CARTILAGE-HAIR HYPOPLASIA IS CAUSED BY MALFUNCTIONING MRP RNA PRADER-WILLI SYNDROME AND THE ABSENCE OF SPECIFIC snoRNAs LONG NONCODING RNAS, A UNIVERSE OF RNAS YET TO BE EXPLORED LINCRNAS REGULATE EXPRESSION USING MANY DIFFERENT MECHANISMS A REPEAT EXPANSION IN A INCRNA CAUSES A SPINOCEREBELLAR ATAXIA A INCRNA AND CELIAC DISEASE A INCRNA PLAYS AN IMPORTANT ROLE IN A MUSCLE DISORDER	<ul> <li>169</li> <li>169</li> <li>170</li> <li>170</li> <li>171</li> <li>172</li> <li>173</li> <li>176</li> <li>178</li> <li>179</li> <li>179</li> <li>180</li> </ul>

Summary	183
Questions	<b>184</b>
Further Reading	185

Chapter 12 Computational Methods Are Critical in the Analysis of Molecular	
Sequences	187
BIOINFORMATICS METHODS ARE WIDELY USED	187
THE HUMAN GENOME AND THE REST OF BIOLOGY CAN ONLY BE UNDERSTOOD IN THE LIGHT OF EVOLUTION	188
HISTORY OF SEQUENCE COMPUTING STARTED WITH PROTEINS	191
HOW SEQUENCES ARE COMPARED	191
MULTIPLE ALIGNMENTS	192
MOLECULAR SEQUENCES ARE COLLECTED IN PUBLIC DATABASES	193
HUGE COLLECTIONS OF SEQUENCE DATA MAY EFFECTIVELY BE SEARCHED FOR SEQUENCE SIMILARITY	193
MORE SENSITIVE METHODS TO REVEAL DISTANT EVOLUTIONARY RELATIONSHIPS AMONG PROTEINS	195
BIOLOGICAL FUNCTION MAY BE PREDICTED FROM AMINO ACID SEQUENCES	197
A NETWORK OF GENE AND PROTEIN FUNCTIONS: GENE ONTOLOGY	198
SEQUENCES TO BUILD TREES: PHYLOGENY	199
HUMAN GENOME COMPUTING: PUTTING THE HUMAN GENOME TOGETHER FROM SHORTER SEOUENCES	203
RECONSTRUCTING A GENOME BY ALIGNMENT OF READS TO THE REFERENCE GENOME	204
INFERRING FUNCTIONS OF INDIVIDUAL GENOMIC VARIANTS	206
COMPARISON TO OTHER SPECIES IMPROVES OUR UNDERSTANDING	207
OF THE HUMAN GENOME	207
Summary	208

#### CONTENTS ix

Questions	209	GENE THERAPY WITH ADENO-	
URLs	209	ASSOCIATED VIRUS	230
Further Reading	209	CANCER GENE THERAPY WITH p53	230
		CANCER GENE THERAPY BY BOOSTING	
Chapter 13		THERAPY	231
Diagnosing the Genome	211	GENOME EDITING: CRISPR/Cas9	231
A DIFFICULT INFLAMMATORY BOWEL DISEASE RESISTS DIAGNOSIS	211	CRISPR/Cas9 AND SICKLE CELL ANEMIA	232
THE FAULTY GENE IS IDENTIFIED BY GENOME SEQUENCING	212	THE NUMBER OF CRISPR APPLICATIONS IS GROWING RAPIDLY	234
A CURE FROM TRANSPLANTATION	212	MUTATIONS IN DUCHENNE MUSCULAR	
DIAGNOSIS OF GENETIC ERRORS	213	ELIMINATION	235
HIGH-THROUGHPUT DNA SEQUENCING METHODS	213	RESTORING DYSTROPHIN EXPRESSION IN DMD BY ELIMINATION	
WHOLE GENOME AND EXOME		OF EXONS	237
SEQUENCING	216	CRISPR/Cas9 APPLIED TO PATIENTS	237
DIAGNOSIS OF A NEUROMUSCULAR DISORDER: ANOTHER SUCCESS STORY OF GENOME SEQUENCING	217	GENE THERAPY TO CHANGE GENE EXPRESSION WITH OLIGONUCLEOTIDES	238
GENOME SEQUENCING TO IDENTIFY GENETIC ERRORS	219	ARTIFICIAL MICRORNAS EXPLOITED AS A THERAPY FOR HUNTINGTON'S	
MANY PEOPLE ARE AFFECTED BY A RARE DISORDER	220	DISEASE	238
RARE GENETIC DEFECTS ARE BEING ANALYZED AT A LARGE SCALE	220	ERRONEOUS SPLICE SITE: USHER SYNDROME	239
Summary	222	SPLICING PLAYS A CRITICAL ROLE IN	240
Questions	223	SPINAL MUSCULAR AIROPHY	240
URLs	223	SPINAL MUSCULAR ATROPHY	241
Further Reading	223	OLIGONUCLEOTIDES TO PRODUCE EXON SKIPPING IN DUCHENNE MUSCLE	2/12
Chapter 14	~~~		243
Correcting Genome Errors	225	PROTEIN	244
THERAPY	225	Summary	244
ADMINISTERING A FUNCTIONAL GENE		Questions	245
TO A PATIENT BY TRANSPLANTATION OF HEMATOPOIETIC STEM CELLS	226	Further Reading	246
CORRECTING A GENETIC DISORDER WITH THE HELP OF A "SAVIOR SIBLING"	227	Chapter 15	240
VIRAL VECTORS FOR GENE THERAPY	228	Ephogue	249
FIRST CASE OF GENE THERAPY: CURING ADENOSINE DEAMINASE		UNA SEQUENCES HAVE A VARIETY OF FUNCTIONAL ROLES	249
DEFICIENCY	228	YOUR GENOME IS NOT YOUR FATE	252
WITH LENTIVIRAL VECTORS LEUKEMIA IS AVOIDED	229	WHO HAS ACCESS TO A PERSON'S GENOME?	254

х

HOW MUCH DO YOU WANT TO KNOW ABOUT YOUR OWN GENOME?	256
SELECTING AGAINST DISEASE MUTATIONS	257
ETHICS OF GERMLINE GENE THERAPY	258
CONCLUDING REMARKS	259
Further Reading	260

Appendixes	261
Glossary	263
Recommended Textbooks for	
Further Reading	271
Index	273

# Preface

The DNA molecule contains genetic information in the form of sequences of nucleotides abbreviated A, T, C, and G. Therefore, genetic information is essentially a very long sequence of these four letters. In a single human individual, the genetic content amounts to about 3 billion letters. Furthermore, information in DNA is used to direct the production of proteins and RNAs. Sequences of nucleotides in DNA and RNA, as well as sequences of amino acids in proteins, are all examples of molecular sequences. This book attempts to present an accessible account of the molecular sequence information contained in the human genome and how it is being used to direct the production of RNA and protein. It addresses the question of what important biological signals are found in the linear sequence of nucleotides in DNA and shows how specific DNA sequences have distinct functions. Such functional sequence elements, typically in a small size range like 3-20 nucleotides, may be classified into a number of different functional categories. Examples are the three-nucleotide sequences that specify amino acids, or short DNA sequences that are targets for proteins that regulate transcription.

To provide biological motivation, the importance of the DNA sequence for biological function is illustrated with a variety of inherited disorders or with other genetic disorders such as cancer. With such examples, different functional elements of a gene, as well as different aspects of genetic information transfer within the cell, are introduced. For instance, a point mutation in a coding sequence of a globin gene gives rise to sickle cell anemia, a splicing mutation results in a form of hemophilia, and a mutation in an untranslated region of an mRNA leads to an iron metabolism disorder.

When discussing the functional consequences of DNA sequence, it is also of interest to consider the molecular impact of mutations and sequence variation. For instance, changes in DNA sequence may affect the recognition by a nucleic acid, such as in the case of a codon being read by a tRNA anticodon. Alternatively, a DNA variant will affect the recognition by a protein, such as a protein regulating transcription. There are indeed throughout the book several examples showing in structural detail the interaction of a nucleic acid with another nucleic acid or with a protein, illustrating the consequences of mutations at the level of molecular interactions.

In essence, therefore, this book has a *molecular sequence perspective*, and recurring themes are *functional DNA sequence elements*, illustration of *functional impact* with genetic disorders, and *molecular interactions* affected by sequence variation.

Why is it important to know about the different DNA sequence elements of the human genome and their functional impact? There are many medical areas where such knowledge is critical, as shown in this book. Examples include pharmacology, gene therapy, and the diagnosis of rare inherited disorders and cancer.

The book is organized such that a major part (Chapters 6 through 12) discusses specific functional sequence elements of the human genome. The preceding chapters offer an introduction to basic concepts with regard to the human genome. Chapters 2 and 3 introduce human genes and their relationship to human disease using sickle cell anemia as an example. Proteins and their structure are introduced (Chapter 2), and the focus is then on DNA structure, the flow of genetic information from DNA to protein, as well as the basic principles of translating an mRNA using the genetic code (Chapter 3). Human genome sequencing and the organization and overall

content of the human genome are topics of Chapter 4. Individual genetic variation, the phenotypic effects of mutations, as well as cancer are then discussed in the following chapter.

Chapters 6 through 12 provide details about functional sequence elements of the genome and their role in the flow of genetic information. The importance of these elements is illustrated with inherited disorders or cancer. First, the role of the coding regions of mRNA molecules is discussed in Chapter 6. Chapter 7 covers repeat expansions in coding regions. Then the role of mRNA untranslated regions (Chapter 8) and the process of splicing to form mature mRNA (Chapter 9) is considered. The complex but important area of transcriptional control is then discussed in Chapter 10. The noncoding RNAs are the subject of Chapter 11. The analysis of DNA and protein sequences requires a substantial amount of computing. Therefore, in Chapter 12, computational methods used with biological sequences are briefly reviewed.

Chapters 13 through 15 provide additional biological motivation as they further illustrate why careful studies of the relationship between sequence and function are essential. These chapters deal with important medical applications of human genome sequencing. Thus, experimental methods to diagnose errors in DNA sequences are described, including successful efforts to reveal causative mutations in rare inherited disorders (Chapter 13). Furthermore, there is recent progress in the area of gene therapy. Some of the most important methods are covered, as well as a few cases of successful therapy for inherited disease (Chapter 14). Finally, some of the applications of human genome sequencing raise a number of ethical concerns, and these are discussed in the final chapter.

When it comes to level of difficulty, I have attempted to present the material at a basic level to make it understandable for a reader without previous studies of genetics and molecular biology. However, the reader will benefit from a basic knowledge about chemistry and biochemistry. Following the principle that a picture says more than a thousand words, the book is richly illustrated. Furthermore, a web supplement to the book that includes scientific updates and answers to selected chapter questions is available at http://toresamuelsson.se/hg.

Why did this book come about? We currently see a dramatic development in terms of human genome sequencing. Research laboratories around the world generate a wealth of genomic sequence data. Sequencing is becoming widely used in the clinic to analyze a variety of genetic disorders, including cancer. In addition, you can order your own genome sequence from a company ("direct-to-consumer" sequencing). With this wealth of genetic information, it becomes increasingly important to know what the human genome sequence is all about, how the sequence should be understood in terms of biological function, and how particular variants in the genome should be interpreted. I wanted to write a book providing an introduction to these topics. In the early days of my scientific career, I worked as a molecular biologist. Eventually, I turned to bioinformatics with a focus on molecular sequences. For a long time I have been intrigued by the digital nature of genetic information, that the genome sequence can be handled with computers as a long string of letters, and that computing with molecular sequences may be used to address a variety of biological problems. There are already very good books dealing with the human genome, but this book has a focus on molecular sequences and it examines in a systematic manner the functional role of DNA sequence elements as illustrated with human genetic disorders.

#### ACKNOWLEDGMENTS

Illustrations of the book were typically created with CorelDRAW version X4. All images of protein or nucleic acid structures were made with the UCSF Chimera package and structures available in the Protein Data Bank (www. rcsb.org; Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. 2000. The Protein Data Bank. *Nucleic Acids Res* 28:235–242). Chimera is developed by the Resource for Biocomputing, Visualization, and Informatics at the University of California, San Francisco (supported by NIGMS P41-GM103311). Chimera is described in Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE. 2004. UCSF Chimera—A visualization system for exploratory research and analysis. *J Comput Chem* 25(13):1605–1612. I am indebted to the KEGG/GenomeNet Support Team for permission to use a metabolic pathway image.

I am grateful to a number of people for their help in the context of this book. Jan Korbel, European Molecular Biology Laboratory, Heidelberg, and Evan Eichler, University of Washington, Seattle, helped out with comments about human structural variation. Johanna Rommens at the Hospital for Sick Children, Toronto, Canada, provided information about the Cystic Fibrosis Mutation Database, and Catherine Porcher, University of Oxford, United Kingdom, about the ELK/GATA sequence logo. Oxford Nanopore Technology allowed the use of a photograph of the MinION apparatus. Aravinda Chakravarti and Sumantra Chatterjee, Johns Hopkins University School of Medicine, Baltimore, Marvland, informed me on the role of transcription factor binding sites with regard to Hirschsprung disease. Eric Ottesen, Iowa State University, Ames, provided information regarding SMN2 exon 8 and nusinersen. I am grateful to Stefan Mundlos, Max Planck Institute for Molecular Genetics, Berlin, Germany, for providing clinical photographs in the context of developmental disorders resulting from genomic rearrangements in the EPHA4 locus. Sally Heywood, Cardiff University, United Kingdom, provided statistics of the collection of mutations in the Human Gene Mutation Database (HGMD). Mark Johnson, reporter at the Milwaukee Journal Sentinel, helped out with questions regarding the Nicholas Volker case. The National Center for Biotechnology Information (NCBI) User Services provided help regarding human singlenucleotide variants that are part of the NCBI dbSNP. I am grateful to Retta Beery for permission to use a photograph of her family. Niclas Juth, Karolinska Institute, Sweden, commented on the ethical topics of Chapter 15.

Colleagues at my department were also very helpful. Gunnar Hansson helped out with the parts of the book dealing with cystic fibrosis, and the noncoding RNA chapter was examined by Chandrasekhar Kanduri. I am also grateful to Per Elias for discussions on a variety of topics. Erik Larsson carefully examined specific parts of the manuscript and had very useful suggestions regarding the contents of Chapter 10.

A number of anonymous reviewers recruited by the publisher provided important comments to the manuscript. Furthermore, I am indebted to staff at Garland Science/CRC Press. Elizabeth Owen, Senior Editor at Garland Science, carefully read all of the text and offered very helpful comments in an early phase of the book writing. I am also grateful for the work and support of Developmental Editor Jordan Wearing and Senior Editor Chuck Crumly at CRC Press.

Finally, I'm very grateful for a stipend from The Royal Society of Arts and Sciences in Gothenburg, Sweden, that allowed me to work on this book during one month in 2017 at Hotel Chevillon, Grez-sur-Loing, France.

#### **Tore Samuelsson**



# Introduction

NA, short for deoxyribonucleic acid, is a universal carrier of hereditary information. In all life forms—viruses, bacteria, fungi, plants, and animals—it carries important instructions for the design of the organism. And not only does it carry information—it is also a molecule designed so that it may be accurately copied to the next generation. DNA is built from simple units, referred to as nucleotides, that are joined to form very long molecules. Each nucleotide contains any of four different nitrogenous bases: adenine, thymine, cytosine, or guanine, abbreviated A, T, C, and G, respectively. It is the sequence of these bases that forms the actual genetic message. Thus, the information in DNA may be expressed as a long sequence of the letters A, T, C, and G—for an example, see **Figure 1.1**.

We refer to the complete genetic material of an organism as its **genome**. The human genome is an astounding three billion letters. An important milestone was reached in biomedical research in 2001 when, for the first time, a draft of the human genome was presented and the complete sequence of letters could be read. A small fraction of the human genome is shown in Figure 1.1. Consider the whole genome printed as a physical book. A total of 6,400 bases are in Figure 1.1. You would need in the order of 500,000 pages like this to cover the full human genome. That would correspond to more than 1,600 books, each with 300 pages. For more on printing the human genome on paper, see Figure 1.2.

The issues addressed by this textbook are related to the three billion letter sequence of the human genome. How are we to make sense of and understand this vast information? What different biological signals are contained in the DNA? How important are different regions of the sequence? Are some regions more important than others? What are the effects in the event the sequence of letters in DNA is changed? In molecular biology laboratories, scientists have carried out experiments to address these questions. In addition, as changes or mutations in DNA are natural components of evolution, nature has by itself carried out experiments during billions of years that may guide us in understanding the relationship between genetic information and biological function. For instance, mutations in DNA can



**Figure 1.1 Portion of the human genome.** Letters A, T, C, and G represent the DNA bases adenine, thymine, cytosine, and guanine, respectively. This page has 6,400 bases. It would take 500,000 pages like this to cover the full human genome. This would correspond to more than 1,600 books assuming one book contains 300 pages. The magnifying glass is to indicate that the object of research concerning the human genome is to associate sequence elements with biological functions. Scientists use experimental methods, as well as computational methods, to do this—a work in progress.

3

give rise to specific inherited diseases as well as cancer. What are the changes in the DNA sequence that cause such deleterious effects?

To answer these questions, we need to understand the organization of the human genome, as well as the different functional sequence elements in that genome. The flow of genetic information is crucial. Hence, DNA specifies what RNA molecules are to be made. One subclass of these RNAs is subject to processing to form messenger RNA (mRNA) molecules. These mRNA molecules in turn act as templates for the production of proteins. Another abundant class of RNA molecules has functions other than to specify proteins. Throughout the elaborate flow of genetic information that includes copying of DNA sequences to RNA, RNA processing, as well as the synthesis of protein using mRNA, specific nucleotide sequences have distinct functions.

In essence, this book explores the information in the human genome and all of the important biological signals that are present. It illustrates various functions of DNA sequences. Examples include protein coding sequences and sequences that regulate the flow of genetic information. For all of the different sequence elements, the relationship between sequence and function is illustrated with disorders of a genetic background.



Figure 1.2 Human genome printed on paper. Scientists at the University of Leicester printed the whole human genome on paper. It resulted in 130 book volumes that would take 95 years to read. (Published under CC BY-SA 2.0.)



# A Molecular Disorder

The theme of this book is the information contained within the human genome as outlined in Chapter 1. As a first element of information, we consider regions in the genome that specify the proteins to be made. Proteins are molecules built from amino acids, and the sequence of amino acids is determined by the sequence of nucleotides in the genome. Regions specifying proteins make up only a minute portion of the entire genome but are nevertheless significant.

We first consider how amino acid sequences are related to inherited disorders. As an example, we discuss the disorder sickle cell anemia. It is caused by a mutation that gives rise to a replacement of the amino acid glutamic acid to valine in the protein hemoglobin. There are multiple reasons why we discuss this particular disorder in some detail in the first book chapters. Early studies of sickle cell anemia were based on nonmolecular clinical observations, and it seemed likely that the disease is inherited. But as research progressed, we eventually obtained a detailed molecular understanding of sickle cell anemia from the changes in DNA to detailed structural information about the hemoglobin protein. Sickle cell anemia is historically significant as it was the first disease to be characterized where a genetic change is associated with a well-defined change in a protein molecule. This finding gave an early clue as to the power of molecular medicine. In addition, very few inherited disorders have been so thoroughly examined as sickle cell anemia, and information about the disease is still being collected today. There is a significant medical impact from this knowledge, which is why we return to this disorder also in other contexts such as gene therapy (Chapter 14).

## THE DISCOVERY OF SICKLE CELL ANEMIA

Abotutuo. Chwechweechwe. Nwiiwii. Nuiduidui. These are all names of a disease common in Western Africa—a disease we now know also as sickle cell anemia. Its history in Africa may be tracked as far back as the seven-teenth century. Classification of the disease was difficult on this continent, because the symptoms were closely related to those of other diseases in tropical areas. It was not until the early twentieth century that sickle cell



Figure 2.1 Sickle or crescent shape of red blood cells characteristic of sickle cell anemia. A sickle cell is shown (to the left) along with normal red blood cells. (Published under CC BY 3.0.)

anemia was first described in a medical publication. The affected individual was Walter Clement Noel.

Noel was born in 1884 on a large estate on Grenada. At this time, this island was a British colony. Noel was from a wealthy black family. He suffered from sickle cell anemia but was still able to attend school, and he completed his undergraduate studies in 1904. The same year, he sailed to New York. During this week-long journey, he developed a leg ulcer, a common complication of sickle cell anemia. When he arrived in New York, he immediately sought medical attention. His ulcer was treated with iodinethis chemical had the effect of killing the bacteria of the ulcer. Noel was thus cured and then traveled to Chicago, Illinois, where he was to study dentistry at the Chicago College of Dental Surgery. Noel was an unusual student at this time in the sense that it was uncommon for students of African descent to reach higher studies. In November 1904, his disease unfortunately got more severe as he developed respiratory problems. These problems persisted for more than 1 month. He finally sought medical attention at the Presbyterian Hospital in Chicago and was examined by an intern, Ernest Irons. Among other tests, Irons examined a blood sample from Noel. Under the microscope he could see that the sample contained, as he phrased it, "many pear-shaped and elongated forms-some small." Iron discussed these findings with his supervisor James Brian Herrick. Despite further investigations, Herrick and Irons could not reveal the cause of these unusual cells.

Noel eventually recovered from his respiratory problems and returned to continue his studies at the dentistry school. However, he experienced additional illnesses during a period of more than two years of studies. Thus, he once was hospitalized for bronchitis and once for painful muscular crises and gallstones. He was then under the care of Irons who kept dutiful notes. When Irons was done with his training, he gave these notes to Herrick. Herrick then took care of Noel for two years.

Herrick was at a national meeting in 1910 and there presented the case of Noel—incidentally without giving credit to Irons. He published a report the same year. To describe the blood cells of the disease, Herrick used the term "sickle shaped cells" (**Figure 2.1**).

Despite his illnesses, Noel graduated from dental school. He then went back to Grenada where he set up a general dentistry practice in St. George's, the capital of Grenada. Not much is known of him from then on. However, in 1916, he overexerted himself. He attended a horse race on Grenada a long way from home and traveled home, all on the same day. As a result of this, he developed a serious respiratory infection. Only a few weeks after the horse race, he died from pneumonia at 32 years old. Noel is buried in a churchyard with a view of the Caribbean Sea. He is next to his sister Jane, who also died young of respiratory problems.

Only three months after Herrick published the case of Noel in 1910, a second case of the same disorder was described. Blood samples from a 25-yearold woman, Ellen Anthony, a resident of Virginia, showed the same strange shape of red blood cells as was observed in Noel. As more cases were identified in the 1920s, it was noted that all individuals with the disease were of African origin. The disease was eventually to be named *sickle cell anemia* (SCA).

#### A RECESSIVE INHERITED DISORDER

Significant advances in understanding sickle cell anemia were made by the end of the 1940s—both with respect to genetic inheritance and as to the molecular basis of the disease. Already in 1923, John Huck studied families with sickle cell anemia and noted that the disease was probably inherited, although his studies did not provide any firm evidence of this theory. Studies were complicated

by the fact that some individuals had symptoms related to sickle cell anemia but had a much milder form of the disease that did not shorten their lives. These individuals were said to have the **sickle cell trait**. It was often difficult to distinguish between these two categories of patients. But in 1949, James Neel carried out a careful examination of families affected by sickle cell anemia and was able to conclude that the disorder is indeed hereditary.

How is sickle cell anemia inherited? We turn to basic principles of genetics and inheritance, first elucidated by Gregor Mendel in the nineteenth century. The genetic makeup of an individual is referred to as the genotype, whereas the **phenotype** is the collection of observable characteristics. The phenotype is determined by the genotype and/or environmental factors. Human individuals-like all other animals-have two copies of each gene, one of paternal and one of maternal origin. A gene may have two or more variants-in the language of genetics, such variants are referred to as alleles. If two individuals have the same allele, they are said to be homozygous for that allele, and if they have two different alleles, they are heterozygous. From the perspective of an allele, one basic principle of inheritance is outlined in Figure 2.2. In this example, two different alleles "A" and "a" are considered, and the two parents are heterozygous, since they both have the allele configuration (genotype) Aa. During the formation of sperm and egg cells that occurs during meiosis, each cell ends up with any one of the two alleles. During fertilization, alleles are combined, and a child of the two parents may in this case have any of the genotypes AA, Aa, and aa.

For the discussion of sickle cell anemia, the two different alleles we consider are the sickle cell variant and the normal form not associated with disease. Individuals with two copies of the sickle gene develop sickle cell anemia, whereas patients with one copy of the sickle gene and one normal copy of the hemoglobin gene have milder symptoms and express the sickle cell trait. This is illustrated by the pedigree (family tree) with members affected by the disorder in **Figure 2.3**.

The inheritance of sickle cell anemia follows the rules of a **recessive** disorder. In such a disorder, two copies of the disease allele are required to develop the disease. The rules of inheritance also inform us on probabilities on inheriting disorders as explained in the diagram in **Figure 2.4**,



Figure 2.2 Inheritance of alleles. Two different alleles (gene variants) "A" and "a" are considered. In this example, both parents have the same setup of these alleles. The sperm and egg cells formed during meiosis have only one copy of each allele. The probability that a certain sperm or egg cell has a specific allele is about 50%. The fertilized egg will have two copies of each autosomal chromosome, and in the example shown here, "a" from the father is combined with the same allele from the mother. However, other outcomes are possible given the parent genotypes. Thus, offspring may be of three different genotypes: AA, Aa, or aa.



trait)

anemia)

Figure 2.3 Sickle cell anemia is a recessive inherited disorder. The pedigree illustrates inheritance of sickle cell anemia. Males are represented by squares, females by circles. The disease gene is indicated in red. In recessive disorders, two copies of the affected allele are required to develop the disease.



Figure 2.4 Probabilities related to sickle cell anemia. Squares have been arranged to discover all possible genotypes that occur in children, given the genotypes of their parents. From this diagram, we may also infer the probability of each offspring genotype. The sickle cell allele is represented by "S" and the normal allele as "N." The top row has the genotype of one parent, and the leftmost column the genotype of the other parent. The other boxes are obtained by copying the parent letters across or down. In this way, we get the predicted frequencies of all the potential genotypes. (a) This is a case where one parent has sickle cell anemia (genotype SS) and the other parent is heterozygous. As the SS genotype occurs in two of the four cells, we estimate that the probability that a child has sickle cell anemia is 50%. In (b) both parents are heterozygous. Here, the SS genotype occurs in one of the four cells, and we estimate that the probability is 25% that a child is affected by the disease. The example shows sickle cell anemia, a recessive disorder, but the method of inferring probabilities may be used to examine just about any case of recessive or dominant inheritance.



**Figure 2.5 Chemical structure of an amino acid.** All amino acids contain a carbon atom to which is connected an amino group, a carboxyl group, a hydrogen atom, and a side chain (R) specific to the amino acid.

Figure 2.6 Chemical structures of selected amino acids. Glycine is the smallest amino acid with a hydrogen as its side chain. Lysine and glutamic acid are amino acids that are charged because of their amino and carboxyl groups, respectively. Valine belongs to a group of nonpolar amino acids. known as a Punnett square. For instance, if one parent is homozygous for the sickle variant and the other parent is heterozygous, the probability that a child of these parents develops sickle cell anemia is 50%. If both parents are heterozygous, the probability is instead 25%.

As opposed to recessive disorders, in **dominant** disorders, only one disease gene is sufficient to develop the disease. Recessive and dominant diseases may be distinguished because they have different inheritance patterns (for an example of a dominant disorder and its inheritance, see Chapter 7).

## SICKLE CELL ANEMIA AND MALARIA

Sickle cell anemia affects many hundreds of thousands around the world. In particular, it is common among individuals of Sub-Saharan African descent. In the United States, it occurs among about 1 out of every 365 black American births, and about 1 in 13 black American babies is born with sickle cell trait. In Africa, sickle cell disease is even more common—in some parts of Sub-Saharan Africa, up to 1 in 30 of all newborns are affected by the disease.

Why is sickle cell anemia such a common inherited disorder in parts of Africa? It was demonstrated in the 1950s that people with the sickle cell trait are more resistant against malaria caused by the protozoan parasite *Plasmodium falciparum*. Therefore, in areas of malaria, the sickle cell gene has a selective advantage. It is still not clear why the sickle cell gene protects against malaria.

## CHARACTERIZING SICKLE CELL ANEMIA FURTHER: THE ROLE OF HEMOGLOBIN

What is actually the cause of sickle cell anemia? Research eventually zoomed in on **hemoglobin**, a protein.

Proteins are important molecules in biology, as they carry out many critical functions. For instance, they act as enzymes, transporters, and receptors mediating hormonal response. By the turn of the twentieth century, it was shown that proteins are composed of **amino acids**. All amino acids have a carboxyl group and an amino group in addition to a side chain specific to the amino acid (**Figure 2.5**). The chemical structures of selected amino acids are shown in **Figure 2.6** (see Appendix Figure A.1 for the structure of all 20 amino acids). In proteins, the amino acids are joined through **peptide bonds** (**Figure 2.7**).

Hemoglobin belongs to the family of transporter proteins. It is an important protein responsible for transporting oxygen from the lungs to the tissues; in addition, it carries carbon dioxide back to the lungs. The oxygenbinding properties of hemoglobin had already been discovered in the nine-teenth century. Could the amino acid composition of hemoglobin somehow



be changed in sickle cell anemia? In one classic experiment in 1949, Linus Pauling compared hemoglobin from healthy individuals to that of sickle cell patients. It turned out that sickle cell globin had a different electrophoretic mobility—that is, when the proteins were subjected to an electric field, the sickle cell hemoglobin moved as a positively charged protein, in contrast to the normal hemoglobin (**Figure 2.8**). Pauling suggested that the difference was due to a change in the number of charged amino acids in the sickle cell version of the protein. The amino acids glutamic acid and aspartic acid are negatively charged because they carry a carboxyl group, and lysine and arginine are positively charged as they have an amino group (see Appendix Figure A.1). It seemed likely that sickle cell hemoglobin was affected in one or more of these amino acids. In the context of this finding, Pauling also used the term "molecular disease." This expression referred to the fact that for the first time the molecular background to a disease was identified.

### MAPPING A DELETERIOUS CHANGE IN HEMOGLOBIN

Further details as to the molecular basis of the disease were elucidated in the 1950s. Vernon Ingram did one crucial experiment in 1956. A protein may be cut up into smaller pieces-peptides-using an enzyme that is able to break up peptide bonds. One such enzyme is trypsin that cuts on the carboxyl-terminal side of lysine and arginine residues. Ingram used this enzyme to fragment hemoglobin and was able to separate the resulting peptides with a newly developed technique of two-dimensional separation by electrophoresis and chromatography. When comparing normal hemoglobin to the sickle cell variant, it turned out that one single peptide was different (Figure 2.9a). The peptide in sickle cell hemoglobin was more positively charged than that of normal hemoglobin. Analysis of its amino acid composition showed that it had less glutamic acid and more of valine, suggesting that in sickle cell anemia, glutamic acid had been replaced by valine. The exact sequence of amino acids in the peptide was soon after determined. In the sickle cell protein, the peptide had valine instead of glutamic acid in one position (Figure 2.9b).

The analysis by Pauling and Ingram of sickle cell hemoglobin in the 1940s and 1950s was an important milestone in molecular biology. For the first time, a genetic inherited disorder was explained in terms of a specific amino acid substitution in a protein. We now know many more examples of inherited disorders associated with amino acid replacements as will be apparent in forthcoming chapters.

### A PROTEIN FOLDS INTO A THREE-DIMENSIONAL SHAPE BASED ON ITS AMINO ACID SEQUENCE

To understand properly the molecular basis of sickle cell anemia, we also need to know the structural consequences of the amino acid substitution. Proteins are built from one or more **polypeptide** chains—each such chain is a string of amino acids joined by peptide bonds. The sequence of amino acids in each of the polypeptide chains will determine the three-dimensional shape of the protein. We commonly refer to four different levels of protein structure: **primary**, **secondary**, **tertiary**, and **quaternary** (**Figure 2.10**). The primary structure refers to the amino acid sequence of the polypeptide chain (**Figure 2.10a**). Noncovalent interactions between amino acids in the same polypeptide chain give rise to structures known



**Figure 2.7** Amino acids are connected through peptide bonds. A unit with CONH represents the peptide bond. R1, R2, and R3 represent amino acid side chains.



Figure 2.8 Electrophoretic mobility of normal hemoglobin and its sickle cell variant as studied by Linus Pauling. Pauling used a method of electrophoresis where the protein samples to be separated are in solution. By applying a current, the proteins will move in the solution dependent on their charge. The position with respect to the electric field is on the x-axis, and the peak height (y-axis) is proportional to protein concentration. The vertical line represents the position of an electrically neutral molecule. The results of Pauling's experiment demonstrated that the sickle cell hemoglobin moved as a positively charged molecule, whereas the normal protein migrated as a negatively charged molecule. The four different samples are (1) normal hemoglobin, (2) sickle cell hemoglobin, (3) hemoglobin prepared from individuals with the sickle cell trait (such individuals have approximately equal proportions of sickle and normal protein), and (4) a synthetic mixture with equal amounts of normal and sickle cell protein.

Figure 2.9 Two-dimensional separation of hemoglobin peptides. (a) Normal hemoglobin and sickle cell globin were treated with trypsin to generate a number of peptides. These peptides were separated in two dimensions using paper electrophoresis and paper chromatography. The patterns obtained are identical except for one peptide (N and S, for normal and sickle cell hemoglobin, respectively). (b) Determination of the peptide sequences showed that glutamic acid in position 6 of the normal protein had been replaced by valine in the sickle cell protein.



as  $\alpha$ -helices and  $\beta$ -sheets. These are secondary structure elements (Figure 2.10b). The tertiary structure of a protein refers to the global folding of the entire polypeptide chain of a protein. This structure may contain  $\alpha$ -helices and/or  $\beta$ -sheets, as well as loop regions in between (Figure 2.10c). In the event the protein has at least two polypeptide chains, also known as **subunits**, we also need to consider the number and arrangement of these polypeptide chains—the quaternary structure (Figure 2.10d).

There are three major experimental methods whereby the threedimensional structure of a protein may be inferred. In **x-ray crystallography**, the protein is in the form of a crystal. The crystal is irradiated with x-rays, and the diffraction pattern obtained is used to calculate the electron density of the crystal. The electron density is finally used to elucidate the coordinates of the different atoms in the protein. Most available protein structures have been determined with x-ray crystallography. Second, **nuclear magnetic resonance** (**NMR**) is a method for protein structure determination where the protein is in solution. Third, **cryo-electron microscopy** is a more recent method used for structural studies of large proteins or complexes involving proteins and other large molecules. Structure may also be predicted from the sequence of amino acids using computational tools, although this is much less reliable than the experimental methods.

## SICKLE CELL DISEASE MAY NOW BE UNDERSTOOD IN THE CONTEXT OF PROTEIN STRUCTURE

In 1959, the structure of hemoglobin was elucidated by Max Perutz using the technique of x-ray crystallography (**Figure 2.11**). The hemoglobin molecule is built from two different polypeptides,  $\alpha$ -globin and  $\beta$ -globin. In one molecule of hemoglobin, there are two  $\alpha$ -chains and two  $\beta$ -chains. It is the  $\beta$ -chain that is affected in sickle cell anemia.

Hemoglobin and a related protein myoglobin were in fact the first protein structures ever to be presented. The structure of hemoglobin, as well as that of the sickle cell variant, have been further refined in later work.



Figure 2.10 The four levels of protein structure. (a) Primary structure is the linear sequence of amino acids along the polypeptide chain. (b) Protein secondary structure elements. An  $\alpha$ -helix is a common helical structure in proteins in which every backbone N-H group donates a hydrogen bond to the backbone C = O group of the amino acid located three or four residues earlier along the protein sequence. Another common secondary structure element in proteins is the  $\beta$ -sheet . It consists of strands connected laterally by backbone hydrogen bonds. (c) The tertiary structure of a protein is its global fold, including  $\alpha$ -helices (blue),  $\beta$ -sheets (orange), and loop regions without ordered structure. The protein shown is a subunit of the enzyme Akt2 with both  $\alpha$ -helical and  $\beta$ -sheet structures. (d) The quaternary structure of Akt2 is shown— a dimer of two identical polypeptide chains.

Figure 2.11 Three-dimensional structure of hemoglobin. Hemoglobin is built from two  $\alpha$ -chains (orange) and two  $\beta$ -chains (cyan). Each of the four subunits contains heme groups (shown in ball-and-stick representation). (Adapted from Paoli M et al. 1996. *J Mol Biol* 256:775. PDB ID: 1GZX.)



The replacement in sickle cell anemia of glutamic acid—a negatively charged amino acid (see Figure 2.6 and Figure 2.12a)—with valine has a significant effect on the properties of hemoglobin. Valine is a hydrophobic amino acid and interacts with a pair of hydrophobic amino acids—phenylalanine in position 85 and leucine in position 88 of the protein chain—that are located in the  $\beta$ -chain of a neighboring globin molecule (Figure 2.12b). As a result, an aggregation process is initiated where globin molecules form fibers. These fibers extend through the red blood cells such that the cells are distorted. Because of this alteration in shape and because sickle cells tend to stick to the wall of blood vessels, normal blood flow is prevented. Pain may arise when blood flow is restricted. Sickled cells have a shorter life span, giving rise to a lowered amount of red blood cells and poor oxygen transport (anemia).

### A MOLECULAR SEQUENCE: PROTEINS CAN BE DEPICTED AS A STRING OF AMINO ACID SYMBOLS

Whereas Ingram focused on the amino acid sequence of selected peptides, the complete amino acid sequences of the human  $\alpha$ - and  $\beta$ -globins were elucidated in the early 1960s. The sequence of the  $\beta$ -chain is shown in **Figure 2.13**, where the amino acids are shown with a background color based on their physical and chemical properties. A polypeptide chain built from amino acids has a distinct polarity, where one end of the polypeptide has a free amino-terminal group (the N-terminus) and the other end has a free carboxyl terminal group (C-terminus). Amino acid sequences are, as a rule, displayed with the N-terminus to the left and the C-terminus to the right. The amino acids are typically abbreviated with one letter code (see Figure 2.13) or three letter codes. The reader is referred to Figure A.1 for details of these codes.

Amino acid sequences of proteins may be compared using a computational procedure of **alignment**. Such an alignment is shown in **Figure 2.14**, where  $\alpha$ - and  $\beta$ -chains of five different vertebrates are compared. From this alignment, we observe that all the globin sequences are similar, reflecting the fact that they have been well conserved during evolution. We also identify features that distinguish the  $\alpha$ - and  $\beta$ -chains,



Figure 2.12 Sickle cell hemoglobin molecules interact with each other in an abnormal manner. (a) Space-filling model of normal  $\beta$ -chain of hemoglobin showing surface location of glutamic acid (colored with side-chain carboxyl oxygen atoms in red) in position 6 of the polypeptide chain. For comparison is shown the structure of valine, the amino acid replacing glutamic acid in sickle cell anemia. (b) Two molecules of sickle cell hemoglobin showing interaction between hydrophobic amino acids in neighboring subunits. The  $\beta$ -chains are in green and  $\alpha$ -chains in red (wireframe representation). Highlighted in space-fill representation is an interaction between valine 6 in one polypeptide with a pair of two other hydrophobic amino acids-phenylalanine 85 and leucine 88-in a neighboring polypeptide. This interaction initiates aggregation of globin molecules into harmful fibers. (Adapted from Tame J, Vallone B. 2000. Acta Crystallogr 56:805-811. PBD ID: 1A3N. Harrington DJ et al. 1997. J Mol Biol 272:398-407. PDB ID: 2HBS.)

 N
 1
 V
 H
 L
 T
 P
 E
 K
 S
 A
 V
 T
 A
 L
 W
 G
 K
 V
 N
 V
 D
 E
 V
 G
 G
 E
 A
 L
 G
 R
 30

 31
 L
 L
 V
 V
 Y
 P
 W
 T
 Q
 R
 F
 E
 S
 F
 Q
 D
 L
 S
 T
 P
 D
 A
 V
 M
 G
 N
 P
 K
 V
 60

 61
 K
 A
 H
 G
 A
 K
 V
 L
 G
 A
 H
 L
 D
 N
 L
 K
 G
 T
 F
 A
 T
 L
 S
 E
 90

 91
 L
 H
 C
 D
 K
 L
 H
 F
 G
 K
 120
 L
 L
 A
 N
 A
 L
 A
 H
 H
 F
 G
 K
 120

such as an insertion of five amino acids in the  $\beta$ -chains, as compared to the  $\alpha$ -chains. However, there is a lot of additional information that we may extract from amino acid sequences and alignments as further discussed in Chapter 12.

Amino acids in proteins are our first instance of a **molecular sequence**. There will be more of these as we enter into the world of nucleic acids—DNA and RNA—in the next chapter. We then address the important problem of how, at a molecular level, glutamic acid was changed into valine in sickle cell anemia.

Figure 2.13 Amino acid sequence of the  $\beta$ -chain of hemoglobin. Amino acids are colored according to their physical and chemical properties. For instance, the basic (positively charged) amino acids arginine (R) and lysine (K) have a red background. Cleavage sites for trypsin as used in the experiment of Figure 2.9 are indicated with arrows.

α	HUMAN BOVINE MOUSE	1 - VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF - DLSH GSAQVKGHGKKVADALT67 1 - VLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHF - DLSH GSAQVKGHGAKVAAALT67 1 - VLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHF - DVSH GSAQVKGHGKKVADALA67
	FROG	
	ZEBRAFISH	1 - SLSDTDKAVVKAIWAKISPKADEIGAEALARMLTVYPOTKTYFSHWADLSPGSGPVKKHGKTIMGAVG68
	HUMAN	1 VHLTPEEKSAVTALWGKVN - VDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLGAFS72
	BOVINE	1 LTAEE <mark>K</mark> AAVTAFWG <mark>K</mark> VK VDEVGGEALG <mark>R</mark> LLVVYPWTQRFFESFGDLSTADAVMNNPKV <mark>K</mark> AHGKKVLDSFS70
ß	MOUSE	1 VHLTDAEKAAVSCLWGKVN - SDEVGGEALG <mark>R</mark> LLVVYPWTQRYFDSFGDLSSASAIMGNAKVKAHGKKVITAFN.72
Ρ	CHICKEN	1 VHWTAEEKQLITGLWGKVNVAECGAEALARLIVYPWTQRFFASFGNLSSPTAILGNPMVRAHGKKVLTSFG72
	FROG	1 VNLTAKERQLITGTWSKICAKTLGKQALGSMLYTYPWTQRYFSSFGNLSSIEAIFHNAAVATHGEKVLTSIG72
	ZEBRAFISH	1 VEWTDAERTAILGLWGKLNIDEIGPQALSECLIVYPWTQRYFATFGNLSSPAAIMGNPKVAAHGRTYMGGLE72
	ΗΠΜΑΝ	68 NAVAHVDDMPNAL SALSDI HAHKI RVDPVNEKI I SHCI I VTI AAHI - PAFETPAVHASI DKELASVSTVI TSKY 14
	HUMAN	68 NAVAHVDDMPNAL SAL SDLHAHKLRVDPVNFKLLSHCLLVTLAAHL - PAEFTPAVHASLDK FLASVSTVLTSKY 14 68 KAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHL - PSDFTPAVHASLDKFLANVSTVLTSKY 14
Q	HUMAN BOVINE MOUSE	68 NAVAHVDDMPNAL SAL SDLHAHKLRVDPVNFKLLSHCLLVTLAAHL - PAEFTPAVHASLDK FLASVSTVLTSKY 14 68 KAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHL - PSDFTPAVHASLDK FLANVSTVLTSKY 14 68 SAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHH - PADFTPAVHASLDKFLASVSTVLTSKY 14
α	HUMAN BOVINE MOUSE CHICKEN	68 NAVAHVDDMPNAL SAL SDLHAHKLRVDPVNFKLLSHCLLVTLAAHL - PAEFTPAVHASLDKFLASVSTVLTSKY14 68 KAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHL - PSDFTPAVHASLDKFLANVSTVLTSKY14 68 SAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHH - PADFTPAVHASLDKFLASVSTVLTSKY14 68 EAANHIDDIAGTLSKLSDLHAHKLRVDPVNFKLLSHCLVVAIHH - PAALTPEVHASLDKFLCAVGTVLTAKY14
α	HUMAN BOVINE MOUSE CHICKEN FROG	68 NAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHL - PAEFTPAVHASLDKFLASVSTVLTSKY14 68 KAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHL - PSDFTPAVHASLDKFLANVSTVLTSKY14 68 SAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHH - PADFTPAVHASLDKFLASVSTVLTSKY14 68 EAANHIDDIAGTLSKLSDLHAHKLRVDPVNFKLLSHCLLVTLASHH - PAALTPEVHASLDKFLCAVGTVLTAKY14 68 EACNHLDNIAGCLSKLSDLHAHKLRVDPVNFKLLGQCFLVVVAIHH - PAALTPEVHASLDKFLCAVGTVLTAKY14
α	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH	68 NAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHL - PAEFTPAVHASLDKFLASVSTVLTSKY 14 68 KAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHL - PSDFTPAVHASLDKFLANVSTVLTSKY 14 68 SAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHH - PADFTPAVHASLDKFLASVSTVLTSKY 14 68 EACHHIDDIAGTLSKLSDLHAHKLRVDPVNFKLLGCCFLVVVAIHH - PAALTPEVHASLDKFLCAVGTVLTAKY 14 68 EACHHLDNIAGCLSKLSDLHAHKLRVDPGNFFLLAHQILVVVAIHH - PAQFTPAVHASLDKFLCAVGTVLTSKY 14 69 EAISKIDDLVGGLAALSELHAFKLRVDPANFKILSHVIVVIAMLF - PADFTPEVHVSVDKFFNNLALASEKY 14
α	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH HUMAN	68 NAVAHVDDMPNAL SAL SDLHAHKL RVDPVNFKLL SHCLLVTLAAHL - PAEFTPAVHASLDK FLASVSTVL TSKY 14 68 KAVEHLDDL PGAL SEL SDLHAHKL RVDPVNFKLL SHSLLVTLASHL - PSDFTPAVHASLDK FLANVSTVL TSKY 14 68 SAAGHLDDL PGAL SAL SDLHAHKL RVDPVNFKLL SHCLLVTLASHH - PADFTPAVHASLDK FLASVSTVL TSKY 14 68 EAANHIDD I AGTLSKL SDLHAHKL RVDPVNFKLL GOCFLVVVAIHH - PAALTPEVHASLDK FLASVSTVL TSKY 14 68 EAANHIDD I AGTLSKL SDLHAHKL RVDPONFKLL GOCFLVVVAIHH - PAALTPEVHASLDK FLASVSTVL TSKY 14 69 EAISKIDDLVGGLAAL SELHAFKL RVDPANFK I SHNVIVVIAMLF - PADFTPEVHVSVDK FFNNLALAL SEKY 14 73 DGLAHLDNLKGTFATL SELHCDKLHVDPENFRLL GNVLVCVLAHHF - GKEFTPPVQAAYQK VVAG VANALAHKY 14
α	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH HUMAN BOVINE	68 NAVAHVDDMPNAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLAAHL - PAEFTPAVHASL DKFLASVSTVL TSKY 14 68 KAVEHL DDL PGAL SEL SDL HAHKL RVDPVNFKLL SHSLLVTLASHL - PSDFTPAVHASL DKFLASVSTVL TSKY 14 68 SAAGHL DDL PGAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLASHH - PADFTPAVHASL DKFLASVSTVL TSKY 14 68 EAANHI DD I AGTLSKL SDL HAHKL RVDPVNFKLL GLCFLVVVA I HH - PAALTPEVHASL DKFLASVSTVL TSKY 14 68 EACNHL DN I AGCL SKL SDL HAHKL RVDPGNFPLL AHQI LVVVA I HH - PAALTPEVHASL DKFL ASVSTVL TSKY 14 69 EACNHL DN I AGCL SKL SDL HAHKL RVDPGNFPLL AHQI LVVVA I HH - PAALTPEVHASL DKFL ASVSTVL TSKY 14 73 DGLAHL DN I AGCL SKL SDL HAY DL RVDPGNFKLL GNVL VVI AMLF - PADFTPEVHVSVDK FFNNL ALAL SEKY 14 74 NGMKHL DDL KGTFAAL SEL HCDKL HVDPENFKLL GNVL VVI ARNF - GKEFTPPVQAAYQK VVAG VANAL AHKY 14
αß	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH HUMAN BOVINE MOUSE	68 NAVAHVDDMPNAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLAAHL - PAEFTPAVHASL DKFLASVSTVL TSKY 14 68 KAVEHLDDL PGAL SEL SDL HAHKL RVDPVNFKLL SHSLLVTLASHL - PSDFTPAVHASL DKFLASVSTVL TSKY 14 68 SAAGHLDDL PGAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLASHH - PADFTPAVHASL DKFLASVSTVL TSKY 14 68 EAANHIDD IAGTLSKL SDL HAHKL RVDPVNFKLL GQCFLVVVAI HH - PAALTPEVHASL DKFLASVSTVL TSKY 14 68 EACNHLDNI AGCLSKL SDL HAHKL RVDPVNFKLL GQCFLVVVAI HH - PAALTPEVHASL DKFLASVSTVL TSKY 14 69 EACNHLDNI AGCLSKL SDL HAYDL RVDPGNFPLL AHQILVVVAI HF - PKQFDPATHKAL DKFLVSVSNVL TSKY 14 69 EAISKI DDL VGGLAAL SELHAFKL RVDPANFKIL SHNVI VVI AMLF - PADFTPEVHVSVDKFFNNL ALAL SEKY 14 73 DGLAHLDNL KGTFATL SELHCDKL HVDPENFRLLGNVL VVVL ANNF - GKEFTPPVQAAYQK VVAG VANALAHKY 14 73 DGLNHLDSL KGTFASL SELHCDKL HVDPENFRLLGNVL VVVL ANNF - GKEFTPVLQADFQK VVAG VANALAHKY 14
α	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH HUMAN BOVINE MOUSE CHICKEN	68 NAVAHVDDMPNAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLAAHL - PAEFTPAVHASL DK FLASVSTVL TSKY 14 68 KAVEHL DDL PGAL SEL SDL HAHKL RVDPVNFKLL SHSLLVTLASHL - PSDFTPAVHASL DK FLASVSTVL TSKY 14 68 SAAGHL DDL PGAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLASHH - PADFTPAVHASL DK FLASVSTVL TSKY 14 68 EAANHI DD I AGTL SKL SDL HAHKL RVDPVNFKLL GQCFLVVVA I HH - PAALTPEVHASL DK FLASVSTVL TSKY 14 68 EACNHL DNI AGCL SKL SDL HAHKL RVDPVNFKLL GQCFLVVVA I HH - PAALTPEVHASL DK FLASVSTVL TSKY 14 69 EA I SKI DDL VGGL AAL SEL HAFKL RVDPGNFPLL AHUL VVVA I HF - PKQF DPATHKAL DK FL YVS SNVL TSKY 14 73 DGL AHL DNL KGTFAAL SEL HAFKL RVDPANFK I L SHNVI VVI AMLF - PADFTPEVHVS VDK FFNNL AL AL SEKY 14 73 NGKHL DDL KGTFAAL SEL HCDKL HVDPENFRLL GNVL VVVL ARNF - GKEFTPVQAAYQK VVAG VANAL AHKY 14 73 DGL NHL DSL KGTFASL SEL HCDKL HVDPENFRLL GNMI VI VLGHHL - GKDFTPAAQAAFQK VVAG VAAL AHKY 14 73 DGL NHL DSL KGTFASL SEL HCDKL HVDPENFRLL GNI VI VLGHHL - GKDFTPAAQAAFQK VVAG VAAL AHKY 14 73 DGL NHL DSL KGTFASL SEL HCDKL HVDPENFRLL GNI VI VLGHHL - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 73 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 74 NGKKL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 75 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 76 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 77 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 78 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 79 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG VAAL AHKY 14 79 DAVKNL DNI KNTFSQL SEL HCDKL HVDPENFRLL GDI LI I VLAAHF - SKDFTPECQAAWQKL VRVVAG AHAL AKKY 14
α	HUMAN BOVINE MOUSE CHICKEN FROG ZEBRAFISH HUMAN BOVINE MOUSE CHICKEN FROG	68 NAVAHVDDMPNAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLAAHL - PAEFTPAVHASL DK FLASVSTVL TSKY 14 68 KAVEHL DDL PGAL SEL SDL HAHKL RVDPVNFKLL SHSLLVTLASHL - PSDFTPAVHASL DK FLASVSTVL TSKY 14 68 SAAGHL DDL PGAL SAL SDL HAHKL RVDPVNFKLL SHCLLVTLASHH - PADFTPAVHASL DK FLASVSTVL TSKY 14 68 EAANHI DDI AGTL SKL SDL HAHKL RVDPVNFKLL GQCFLVVVAI HH - PAALTPEVHASL DK FLASVSTVL TSKY 14 68 EAANHI DDI AGTL SKL SDL HAHKL RVDPVNFKLL GQCFLVVVAI HH - PAALTPEVHASL DK FLASVSTVL TSKY 14 69 EAANHI DDI AGCL SKL SDL HAY DL RVDPGNFPLL AHQI LVVVAI HH - PAALTPEVHASL DK FL CAVGTVL TAKY 14 69 EAISKI DDL VGGLAAL SEL HAFKL RVDPANFKIL SHNVI VVI AML F - PADFTPEVHVSVDK FFNNL AL AL SE KY 14 73 DGL AHL DNL KGTFATL SEL HCDKL HV DPENFRLL GNVL VCVL AHHF - GKEFTPVQAAYQK VVAG VANAL AHKY 14 74 NGMKHL DDL KGTFAAL SEL HCDKL HV DPENFRLL GNNI VIVL GHHL - GKDFTPAAQAAFQK VVAG VANAL AHKY 14 73 DGL NHL DSL KGTFASL SEL HCDKL HV DPENFRLL GNUI VIVL GHHL - GKDFTPAAQAAFQK VVAG VANAL AHKY 14 73 DAVKNL DNI KNTFSQL SEL HCDKL HV DPENFRLL GDI LI IVLAAHF - SKDFTPECQAAWQKL VRVVAHAL ARKY 14 73 AVKNL DNI KNTFSQL SEL HCDKL HV DPENFRLL GDI LI IVLAAHF - SKDFTPECQAAWQKL VRVVAHAL ARKY 14 74 NGMKHL DDI KGYYAQL SKYHSETL HV DPYNFKRL GDI LI IVLAAHF - SKDFTPECQAAWQKL VRVVAHAL ARKY 14 75 ALKHMDDI KGYYAQL SKYHSETL HV DPYNFKRFCSCTI I SMAQTL - QEDFTPELQAAFEKLFAAIADALGKGY 14

Figure 2.14 Alignment of vertebrate hemoglobin amino acid sequences. Coloring of amino acids is as in Figure 2.13.

#### **SUMMARY**

- Sickle cell anemia is a recessive inherited disorder. For the disease to be expressed, individuals must carry two copies of the sickle cell allele. Individuals with only one copy have the sickle cell trait with less severe symptoms.
- Protein structure may be described at the four levels: primary, secondary, tertiary, and quaternary.
- Hemoglobin is built from two chains of  $\alpha$ -globin and two chains of  $\beta$ -globin.
- In sickle cell anemia, the amino acid glutamic acid of β-globin is replaced by the hydrophobic amino acid valine, resulting in an aggregation process.
- Sickle cell anemia was the first disease to be characterized where a genetic change is associated with a well-defined change in a protein molecule.
- The primary structure of a protein—the amino acid sequence—may be depicted as a string of amino acid symbols. Amino acid sequences may be compared using a computational procedure of alignment.

## QUESTIONS

- 1. What is characteristic of recessive and dominant inheritance, respectively?
- 2. What is the difference between sickle cell trait and sickle cell anemia?
- **3.** Explain the concepts *genotype* and *phenotype*.
- **4.** Explain what is meant when we say that an individual is *homozygous* for a specific genetic variant (allele).
- 5. Use a Punnett square to show the offspring when one parent is homozygous for the sickle cell allele and one parent has two normal globin alleles.
- 6. Why is the sickle cell gene variant (allele) fairly common in specific populations?
- **7.** Draw the chemical structure of an *amino acid*. Also draw a dipeptide (two amino acids joined with a *peptide bond*).
- 8. What are the acidic and basic amino acids, respectively?
- 9. Give examples of hydrophobic amino acids.

- **10.** What were the experiments carried out in the 1940s and 1950s to elucidate the molecular basis of sickle cell anemia?
- **11.** Explain what is meant by *primary*, *secondary*, *tertiary*, and *quaternary* structures of proteins.
- **12**. What are the methods used to infer the three-dimensional structure of proteins and other large molecules?
- 13. In sickle cell anemia, one amino acid is replaced with another in the  $\beta$ -globin subunit of hemoglobin. What are the consequences in terms of amino acid interactions and protein structure? What are the physiological consequences?
- 14. What is meant by a molecular sequence?

## **FURTHER READING**

#### General on sickle cell anemia and its history

- Bjorklund R. 2011. Sickle cell anemia. Marshall Cavendish Benchmark, New York.
- Orkin SH, Higgs DR. 2010. Medicine. Sickle cell disease at 100 years. *Science* 329(5989):291–292.
- Serjeant GR. 2001. The emerging understanding of sickle cell disease. Br J Haematol 112(1):3–18.
- Serjeant GR. 2010. One hundred years of sickle cell disease. Br J Haematol 151(5):425–429.

#### Early work on human hemoglobin and myoglobin

- Baglioni C. 1961. An improved method for the fingerprinting of human hemoglobin. *Biochim Biophys Acta* 48:392–396.
- Braunitzer G, Hilschmann N, Rudloff V et al. 1961. The hæmoglobin particles. *Nature* 190:480–482.
- Braunitzer G, Gehring-Mueller R, Hilschmann N et al. 1961. The structure of normal adult human hemoglobins. *Hoppe Seylers Z Physiol Chem* 325:283–286.

- Ingram VM. 1959. Abnormal human haemoglobins. III. The chemical difference between normal and sickle cell haemoglobins. *Biochim Biophys Acta* 36:402–411.
- Kendrew JC, Bodo G, Dintzis HM et al. 1958. A three-dimensional model of the myoglobin molecule obtained by x-ray analysis. *Nature* 181(4610):662–666.
- Perutz MF, Rossmann MG, Cullis AF et al. 1960. Structure of haemoglobin: A three-dimensional Fourier synthesis at 5.5-A. resolution, obtained by x-ray analysis. *Nature* 185(4711):416–422.

#### Sickle cell anemia: Early genetic studies

Neel JV. 1949. The inheritance of sickle cell anemia. *Science* 110(2846):64–66.

#### Sickle cell anemia: Early paper by Pauling

Pauling L, Itano HA, Singer SJ, Wells IC. 1949. Sickle cell anemia: A molecular disease. *Science* 110(2865):543–548.