# chapter 1 Proteins

Proteins perform a variety of functions, including enzymatic catalysis, transporting ions and molecules from one organ to another, nutrients, contractile system of muscles, tendons, cartilage, antibodies, and regulating cellular and physiological activities. The functional properties of proteins depend on their three-dimensional structures. The native structure of a protein can be experimentally determined using X-ray crystallography, nuclear magnetic resonance (NMR) spectroscopy, electron microscopy, etc. Over the past 40 years, the structures of more than 53,000 proteins (as of May 12, 2009) have been determined. On the other hand, the amino acid sequences are determined for more than eight million proteins (as of May 5, 2009). The specific sequence of amino acids in a polypeptide chain folds to generate compact domains with a particular three-dimensional structure. Anfinsen (1973) stated that the polypeptide chain itself contains all the information necessary to specify its three-dimensional structure. Deciphering the three-dimensional structure of a protein from its amino acid sequence is a long-standing goal in molecular and computational biology.

# 1.1 Building blocks

Protein sequences consist of 20 different kinds of chemical compounds, known as amino acids, and they serve as building blocks of proteins. Amino acids contain a central carbon atom ( $C_{\alpha}$ ), which is attached to a hydrogen atom, an amino group (NH<sub>2</sub>), and a carboxyl group (COOH) as shown in **Figure 1.1**. The letter R in **Figure 1.1** indicates the presence of a side chain, which distinguishes each amino acid.

# **1.1.1** Amino acids

Amino acids are naturally of 20 different types as specified by the genetic code emerged from DNA sequences. Furthermore, nonnatural amino acids occur, in rare cases, as the products of enzymatic modifications after translocation. The major difference among the 20 amino acids is the *side chain* attached to the  $C_{\alpha}$  through its fourth valance. The variation of side chains in 20 amino acids is shown in **Figure 1.2**. These residues are represented by conventional three- and one-letter codes. Most of the databases use single-letter codes.

The amino acids are broadly divided into two groups, hydrophobic and hydrophilic, based on the tendency of their interactions in the presence of water



FIGURE 1.1 Representation of amino acids. R is the side chain that varies for the 20 amino acids.

molecule. The hydrophobic residues have the tendency of adhering to one another in aqueous environment. Generally, amino acids, Ala (A), Cys (C), Phe (F), Gly (G), Ile (I), Leu (L), Met (M), Val (V), Trp (W), and Tyr (Y), are considered as hydrophobic residues. In this category, Ala, Ile, Leu, and Val contain aliphatic side chains; Phe, Trp, and Tyr contain aromatic side chains; and Cys and Met contain sulfur atom. Gly has no side chain, and it has hydrogen (H) at the fourth position. Two Cys residues in different parts of the polypeptide chain but adjacent to each other in the threedimensional structure of a protein can be oxidized to form a disulfide bridge. The



# Hydrophobic residues

**FIGURE 1.2** The common 20 amino acids in proteins. The three- and one-letter codes for the amino acids are also given. The amino acids are classified into hydrophobic (hydrogen, aliphatic, aromatic, and sulfur containing) and hydrophilic (negatively charged, positively charged, and polar). The side chains are marked with oval boxes.



Hydrophilic residues

FIGURE 1.2 (Continued)

formation of disulfide bridges in protein structures stabilizes the protein, making it less susceptible to degradation.

Amino acids, Asp (D), Glu (E), His (H), Lys (K), Asn (N), Pro (P), Gln (Q), Arg (R), Ser (S), and Thr (T), are classified as hydrophilic residues. In this category, Asp and Glu are negatively charged; His, Lys, and Arg are positively charged; and others are polar and uncharged.

#### **1.1.2** Formation of peptide bonds

The carboxyl group of one amino acid interacts with the amino group of another to form a peptide bond by the elimination of water (**Figure 1.3**). Amino acids are joined end-to-end during protein synthesis by the formation of such peptide bonds. The peptide bond (C–N) has a partial double-bond character due to resonance, and hence there is no rotation about the peptide bond. In **Figure 1.3**, the peptide is represented as a planar unit with the C=O and N–H groups positioning in opposite directions in the plane. This is called *trans*-peptide. There is another form, *cis*-peptide in which the C=O and N–H groups point in the same direction. To avoid steric hindrance, the *trans* form is frequently presented in protein structures for all amino acids except Pro, which has both *trans* and *cis* forms. The *cis* prolines are found in bends of the polypeptide chains.



**FIGURE 1.3** Formation of a peptide bond by the elimination of a water molecule.

A protein chain is formed by several amino acids in which the amino group of the first amino acid and the carboxyl group of the last amino acid remain intact, and the chain is said to extend from the amino (N) to the carboxyl (C) terminus. This chain of amino acids is called a polypeptide chain, main chain, or backbone. Amino acids in a polypeptide chain lack a hydrogen atom at the amino terminal and an OH group at the carboxyl terminal (except at the ends), and hence amino acids are also called **amino acid residues** (simply residues). Nature selects the combination of amino acid residues to form polypeptide chains for their function, similar to the combination of alphabets to form meaningful words and sentences. These polypeptide chains that have specific functions are called **proteins**.

# **1.2** Hierarchical representation of proteins

Depending on their complexity, protein molecules may be described by four levels of structure (Nelson and Cox, 2005): primary, secondary, tertiary, and quaternary (**Figure 1.4**). Because of the advancements in the understanding of protein structures, two additional levels such as supersecondary and domain have been proposed between secondary and tertiary structures. A stable clustering of several elements of secondary structures is referred to as a supersecondary structure. A somewhat higher level of structure is the domain, which refers to a compact region and distinct structural unit within a large polypeptide chain.

# **1.2.1 Primary structure**

Primary structure describes the linear sequence of amino acid residues in a protein. It includes all the covalent bonds between amino acids. The relative spatial arrangement of the linked amino acids is unspecified.

#### Proteins

5



**FIGURE 1.4** Structural organization of proteins.

# **1.2.2 Secondary structure**

Secondary structure refers to regular, recurring arrangements in space of adjacent amino acid residues in a polypeptide chain. It is maintained by hydrogen bonds between amide hydrogens and carbonyl oxygens of the peptide backbone. The major secondary structures are  $\alpha$ -helices and  $\beta$ -structures.

The  $\alpha$ -helical conformation was first proposed by Linus Pauling and co-workers (1951), and a typical  $\alpha$ -helix is shown in **Figure 1.5**. In this structure, the polypeptide backbone is tightly wound around the long axis of the molecule, and R groups of the amino acid residues protrude outward from the helical backbone. The repeating



**FIGURE 1.5** Structure of a typical  $\alpha$ -helix. The hydrogen bonds between the residues *n* and *n* + 4 are shown as dotted lines. Figure was taken as a screenshot from the Web, http://www.food-info. net/uk/protein/structure.htm



**FIGURE 1.6** Structures of (a) antiparallel and (b) parallel  $\beta$ -sheets. The dotted lines show the hydrogen bonds between amino acid residues. The arrows indicate the directions of the polypeptide chain, from N- to C-terminal. Figure was taken as a screenshot from the Web, http://www.food-info.net/uk/protein/structure.htm.

unit is a single turn of a helix, which extends about 0.54 nm along the axis, and the number of amino acid residues required for one complete turn is 3.6. In an  $\alpha$ -helix, each carbonyl oxygen (residue, *n*) of the polypeptide backbone is hydrogen bonded to the backbone amide hydrogen of the fourth residue further toward the C-terminus (residue, *n* + 4). The hydrogen bonds, which stabilize the helix, are nearly parallel to the long axis of the helix.

The other common secondary structure is  $\beta$ -structure that includes  $\beta$ -strands and  $\beta$ -sheets.  $\beta$ -strands are portions of the polypeptide chain that are almost fully extended, and several  $\beta$ -strands constitute  $\beta$ -sheets.  $\beta$ -sheets are stabilized by hydrogen bonds between carbonyl oxygens and amide hydrogens on adjacent  $\beta$ strands (**Figure 1.6**). In  $\beta$ -sheets, the hydrogen bonds are nearly perpendicular to the extended polypeptide chains. The  $\beta$ -strands may be either parallel (running in the same N- to C-terminal) or antiparrallel (running in opposite N- to C-terminal directions).

In a polypeptide chain, the  $\alpha$ -carbon atoms of adjacent amino acids are separated by three covalent bonds arranged as  $C_{\alpha}$ —C—N— $C_{\alpha}$ . In these bonds, rotation is permitted about the N— $C_{\alpha}$  and  $C_{\alpha}$ —C bonds, and the torsional angles are conventionally denoted as  $\phi$  and  $\psi$ , respectively. Every secondary structure is described completely by these two torsional angles that are repeated at each residue. The allowed values for  $\phi$  and  $\psi$  can be shown graphically by simply plotting these values known as Ramachandran plot (Ramachandran et al. 1963). **Figure 1.7** shows the conformations that are permitted for most amino acid residues in Ramachandran plot.



**FIGURE 1.7** Ramachandran plot showing the allowed regions of  $\alpha$ -helical and  $\beta$ -strand conformations. Figure was taken as a screenshot from the Web, http://swissmodel.expasy.org/course/text/chapter1.htm.

# **1.2.3** Tertiary structure

Tertiary structure refers to the spatial relationship among all amino acids in a polypeptide; it is the complete three-dimensional structure of the polypeptide with atomic details. Tertiary structures are stabilized by interactions of side chains of nonneighboring amino acid residues and primarily by noncovalent interactions. The formation of tertiary structure brings the amino acid residues that are far apart in the primary structure close together.

# 1.2.4 Quaternary structure

Quaternary structure refers to the spatial relationship of the polypeptides or subunits within the protein. It is the association of two or more polypeptide chains into a multisubunit or oligomeric protein. The polypeptide chains of an oligomeric protein may be identical or different. The quaternary structure also includes the cofactor and other metals, which form the catalytic unit and functional proteins.

# **1.3 Structural classification of proteins**

Proteins are broadly classified into two major groups: fibrous proteins, having polypeptide chains arranged in long strands, and globular proteins, with polypeptide chains folded into a spherical or globular shape.

# **1.3.1** Fibrous proteins

Fibrous proteins are usually static molecules and play important structural roles in the anatomy and physiology of vertebrates, providing external protection, support, shape, and form. They are water insoluble and are typically built upon a single, repetitive structure assembled into cables or threads. Examples of fibrous proteins





**FIGURE 1.8** Ribbon diagram for four typical protein structures in different structural classes (a) all- $\alpha$  (4MBN), (b) all- $\beta$  (3CNA), (c)  $\alpha + \beta$  (4LYZ), and (d)  $\alpha/\beta$  (1TIM). Figure was adapted from Gromiha and Selvaraj (2004).

are  $\alpha$ -keratin, the major component of hair and nails, and collagen, the major protein component of tendons, skin, bones, and teeth.

# **1.3.2** Classification of globular proteins

Globular proteins are categorized into four structural classes: all- $\alpha$ , all- $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  (Levitt and Chothia, 1976). The ribbon diagrams illustrating the structures in each class are shown in **Figure 1.8**.

The all- $\alpha$  and all- $\beta$  classes are dominated by  $\alpha$ -helices ( $\alpha > 40\%$  and  $\beta < 5\%$ ) and by  $\beta$ -strands ( $\beta > 40\%$  and  $\alpha < 5\%$ ), respectively (**Figures 1.8a** and **b**). The  $\alpha + \beta$  class contains both  $\alpha$ -helices (>15%) and antiparallel  $\beta$ -strands (>10%) that do not mix but tend to segregate along the polypeptide chain (**Figure 1.8c**). The  $\alpha/\beta$  class proteins (**Figure 1.8d**) have mixed or approximately alternating segments of  $\alpha$ -helical (>15%) and parallel  $\beta$ -strands (>10%).



**FIGURE 1.9** Representation of (a)  $\alpha$ -helical and (b)  $\beta$ -barrel membrane proteins. The membrane spanning regions are shown within the disc. Protein structures were taken from Protein Data Bank of Transmembrane Proteins (http://pdbtm.enzim.hu/).

#### **1.3.3** Membrane proteins

Membrane proteins, which require embedding into the lipid bilayers, have evolved to have amino acid sequences that will fold with a hydrophobic surface in contact with the alkane chains of the lipids and polar surface in contact with the aqueous phases on both sides of the membrane and the polar head groups of the lipids (**Figure 1.9**). In genomes, 30% of the proteins are suggested to be membrane proteins, and most of the transmembrane helical and strand proteins are identified as targets for drug design. Membrane proteins perform a variety of functions, including cell–cell signaling and mediating the transport of ions and solutes across the membrane. They are of two kinds: (i) transmembrane helical proteins in which they span the cytoplasmic membrane with  $\alpha$ -helices (White and Wimley, 1999) and (ii) transmembrane  $\beta$ -barrel proteins that traverse the outer membranes of gramnegative bacteria with  $\beta$ -strands (Schulz, 2003). **Figure 1.9** shows the structures of membrane proteins with these two different motifs,  $\alpha$ -helices and  $\beta$ -strands.

# **1.4 Databases for protein sequences**

Recombinant DNA techniques have provided tools for the rapid determination of DNA sequences and, by inference, the amino acid sequences of proteins from structural genes. The number of such sequences is increasing exponentially, and these sequences have been deposited in the form of database, generally, known as protein sequence databases. Specifically, Georgetown University, Washington, D.C., USA, developed the database, Protein Information Resource (PIR). The Swiss Institute of Bioinformatics and European Bioinformatics Institute developed SWISS-PROT and TrEMBL databases. Recently, progress has been made to set up a single worldwide database of protein sequence and function, UniProt, by unifying PIR, SWISS-PROT, and TrEMBL database activities.

# **1.4.1** Protein Information Resource

PIR has evolved from the Atlas of Protein Sequence and Structure established in the early 1960s by Margaret O. Dayhoff (Dayhoff et al. 1965). It produces the largest, most comprehensive, annotated protein sequence database in the public

domain, the PIR International Protein Sequence Database, in collaboration with the Munich Information Center for Protein Sequences and the Japan International Protein Sequence Database (Barker et al. 2000). It is freely available at http://pir.georgetown.edu/. PIR offers a wide variety of resources mainly oriented to assist the propagation and standardization of protein annotation on three major aspects: (i) PIRSF, protein family classification system; (ii) iProClass, integrated protein knowledgebase; and (iii) iProLink, literature, information, and knowledge. The iProClass database provides value-added information reports on protein sequences, structures, families, functions, interactions, expressions, and modifications. The sequence information of a specific protein can be searched with simple "Text search" in iProClass (Figure 1.10a). The search yielded 10 proteins, and the correct one has been selected with a click on the left-side box. It is also possible to save the results as a table or in FASTA format. The result obtained for the search with "Human lysozyme" is shown in Figure 1.10b. It includes general information



**FIGURE 1.10** Text search in iProClass of PIR database: (a) the search with "Human lysozyme," along with intermediate steps, and (b) the information provided at the result page are shown.

# Proteins 11

	1 MEMBER		TLALPNRKAVADHLLM LIGCLRNCSAVTAAAKQLAE VTGFSNAKTTAVK				
Protein Information Resc			Tex	t Search:	10		
About PIR	Databases Searcl	1/Analysis Download	Support				
ProClass Summary Report for	r UniProtKB Entry: P616	26	Related Sequences	BioThesaurus	ID Mapping		
GENERAL INFORMATION					<u></u>		
	UpiDrott/8_ID	UniBrotKB Accordion		Brotoin Nama			
	LYSC HUMAN	P61626; P00695; Q13170; Q9UCF8		Lysozyme C precursor			
Protein Name and ID	PIR-PSD: <u>LZHU</u> RefSeq: <u>NP 000230.1</u> GenPept: <u>AA59535.1</u> ; <u>AA</u> IPI: <u>IPI00019038</u>	C63078.1; EAW97222.1; AAA59536.1;	CAA32175.1 ; AAH04147.1 ; EAW	97221.1 ; AAA36188.1	1		
Taxonomy	Source Organism: Homo sap Taxon Group: Euk/mammal NCBI Taxon: <u>9606</u> Lineage: Eukaryota; Metazo Catarrhini; Hominidae; Hom	oiens (Human) oa; Chordata; Craniata; Vertebrata; Eute no.	eleostomi; Mammalia; Eutheria; Ei	uarchontoglires; Primates; F	laplorrhini;		
Gene Name	LYZ; LZM						
Keywords	3d-structure; amyloid; amy bond; glycosidase; hydrola	loidosis; antimicrobial; bacteriolytic enzy se; polymorphism; polysaccharide degra	me; complete proteome; direct pr dation; signal	otein sequencing; disease r	nutation; disulfide		
Function	Lysozymes have primarily a enhance the activity of imm	bacteriolytic function; those in tissues a unoagents.	and body fluids are associated wit	ih the monocyte- macropha(	ge system and		
Subunit	Monomer.						
CROSS-REFERENCES							
Bibliography	► <u>View Bibliography Informa</u> Annotated references: PMID: <u>8105095; 1035048</u> <u>More</u>	tion ►Submit Bibliography 1; 10469827; 10561612; 11887182; 11!	927576; <u>11936950</u> [PDB/GeneRIF	1			
	0 <i>ther references:</i> PMID: 110 <del>8445</del> ; 12675840; 15745733; 8765309; 9659355; 9745729; 18391951; 9359845; 8566845; 17353931; 9883972; 366724; 10534505; 12477932; 10558865; 18591461						
DNA Sequence	GenBank/EMBL/DDBJ: <u>M21</u>	.119; <u>J03801; M19045; X14008; U25677</u>	; <u>BC004147</u>				
Structure	185W         SCOP         CATH         FSSP         MIGD           186W         SCOP         CATH         SSP         MIGD           186W         SCOP         CATH         SSP         MIGD           186Y         SCOP         CATH         SSP         MIGD           186Y         SCOP         CATH         SSP         MIGD           187Y         SCOP         CATH         SSP         MIGD           187Y         SCOP         CATH         SSP         MIGD           187W         SCOP         CATH         SSP         MIGD           187%         SCOP         CATH         SSP         MIGD           187%         SCOP         CATH         SSP         MIGD           1884         SCOP         CATH         SSP         MIGD           1884         SCOP         CATH         SSP         MIGD     <	PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM PDSLIM					
PIR Feature & Post Translational Modifications	rearl; a corve site: sou, App (35,71) [predicted] FEAT2: Indirg: Bise: subtrack doi:10.00 [predicted] FEAT2: Indirg: accuumos (1:10) [predicted] FEAT3; domain: signal accuumos (1:10) [predicted] FEAT5; roduct: spozyme (19-140) [experimental] Phosphoste: <u>PS1526</u>						
FAMILY CLASSIFICATION							
Unikel	UniRef100 P61626; UniRef90 P6	1625; UniRel50_P61625					
Pfam Domain	Pfam: PF00062: C-time lysozyme	/aloha-lactalbumin family (19-146)					
Prosite Motif	Prosite: PS00128: PD0C00119: 4	Ipha-lactalbumin / lysozyme C signature.					
Interfere	Prosite: PS51348: PDOC00119: 4 InterPro: LYSC_HUMAN IPR001916: Glycoside hydrolase,	lipha-lactalbumin / lýsozyme C family profile. family 22					
inter/D	IPR010709: Glycoside hydrolase, IPR000974: Glycoside hydrolase, IPR000974: Glycoside hydrolase,	family 22, conserved site family 22, lysozyme (a+b): <i>Fold:</i> Lysozyme-like: <i>Superfamily</i> : Lysozym	e-like: <i>Femily</i> : C-type lysozyme [1339-4-	134L:A: 185U:A: 185V:A: 185W:4	x 185X:A: 185Y:A: 1857-		
SCOP Fold	1852:8; 187L:A; 187M:A; 187N:A 1CJ8:A; 1CJ9:A; 1CKC:A; 1CKD:A	; 1870-A; 187P:A; 187Q:A; 187R:A; 187S:A; 1885 ; 1CKF:A; 1CKG:A; 1CKG:B; 1CKH:A; 1D6P:A; 1D6	1:43, 1883:8; 1884:4; 1884:4; 1885:4; 1885:4; 18 Q:A; 1DI3:A; 1DI4:A; 1DI5:A; 18Q4:A; 1E Q:A; 1DI3:A; 1DI4:A; 1DI5:A; 1EQ4:A; 1E	B5:B; 1C43:A; 1C45:A; 1C46:A; 10 Q5:A; 1EQE:A; 1GAY:A; 1GAZ:A; 10	C7P:A; 1CJ6:A; 1CJ7:A; GB0:A; 1GB2:A; More]		
Other Classification	PRINTS: PRO132 LYSOZYME PRINTS: PRO135 LYZLACT SMART: SM00263 LYZL HomoloGene: 37278	gnature					
FEATURE & SEQUENCE DISPLAY							
	Length = 140 F01625 FF00062 1 61 121	Click on a bar to show its sequence: to copy and parts it, per 	209 CEA WAN (ES). 145 145 145 145 145 145 145 145				
11		(b)					

FIGURE 1.10 (Continued)

BLAST Search Form		
Retrieve sequences sin	nilar to your query	
1. Select a database:	OuniProtKB (or restricted by organism/taxon group) OuniRef100	
2. Insert the query seque or enter ">" followed	ence below (FASTA format or sequence only) by a UniProtKB identifier:	
>P61626		
	Ontions I	
Submit Reset		

FIGURE **1.11** Utility of similar search option available in PIR. UniprotKB identifier for human lysozyme, ">P61626" is given as input.

(protein name, taxonomy, gene name, keywords, function, and subunit), crossreferences (bibliography, DNA sequence, genome, ontology, function, interaction, structure, and posttranslational modifications), family classification, and feature and sequence display.

It has several features such as similarity search using BLAST and FASTA, peptide match, pattern search, pairwise sequence alignment, and multiple sequence alignment. The similarity search of human lysozyme against UniProtKB (UniProt knowledgebase) using the alignment program BLAST is shown in **Figure 1.11**. It can also be searched using the program FASTA. The partial results obtained with the search option are depicted in **Figure 1.12**. It indicates the sequences and their codes that match the query sequence along with other details, protein name, organism, length, % identity, overlap, e-value, etc (**Figure 1.12a**). Furthermore, it shows the alignment details with other proteins (**Figure 1.12b**). This will be helpful to identify the homologous sequences of any query protein. PIR can also be searched for any specific patterns, for example, alternating hydrophilic and hydrophobic residues as a pattern for  $\beta$ -strands (see **Chapter 2**), and continuous stretches of hydrophobic residues (e.g., AVILLIVWFFGA) in transmembrane helical proteins, etc.

# 1.4.2 SWISS-PROT and TrEMBL

SWISS-PROT (Bairoch and Apweiler, 1996) is an annotated protein sequence database established in 1986 and maintained collaboratively, since 1987, by the Department of Medical Biochemistry of the University of Geneva and the EMBL Data Library. It is a curated protein sequence database, which strives to provide a high level of annotation (such as the description of the function of a protein, its domain structure, posttranslational modifications and variants), a minimal level of redundancy, and a high level of integration with other databases. TrEMBL is a

_					01005 10	SSear	:h	BLAS	ST Search	
	Protein AC/ID	Protein Name	Length 🖯	Organism Name	PIRSF ID 🖯	Overlap 🗦	%iden 🗦	E-Value 🗦	Score 🗦	Alignment
E	P61627/LYSC_PANPA /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Pan paniscus	PIRSF001064	148	100	4e-77	289	
E	P61628/LYSC_PANTR /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Pan troglodytes	PIRSF001064	148	100	4e-77	289	
C	P61626/LYSC_HUMAN /ProClass UniProtKB/Swiss-Prot	Lysozyme C BioThesaurus	148	Homo sapiens	PIRSF001064	148	100	4e-77	289	
E	B2R4C5/B2R4C5_HUMAN	Lysozyme C BioThesaurus	148	Homo sapiens		148	100	4e-77	289	
C	P79179/LYSC_GORGO ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	<u>Gorilla gorilla gorilla</u>	PIRSF001064	148	99	2e-76	287	
C	P79239/LYSC_PONPY /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Pongo pygmaeus	PIRSF001064	148	<u>97</u>	2e-76	287	
٥	P79180/LYSC_HYLLA /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Hylobates lar	PIRSF001064	148	<u>95</u>	5e-75	283	
E	P61634/LYSC_ERYPA /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Erythrocebus patas	PIRSF001064	148	<u>89</u>	3e-72	273	
C	P61633/LYSC_CERAE /ProClass UniProtK8/Swiss-Prot	Lysozyme C BioThesaurus	148	Chlorocebus aethiops	PIRSF001064	148	<u>89</u>	3e-72	273	
٥	P61629/LYSC_PAPAN /ProClass UniProtKB/Swiss-Prot	Lysozyme C BioThesaurus	149	<u>Papio anubis</u>	PIRSF001064	148	88	4e-71	270	
E	P79811/LYSC_NASLA /ProClass UniProtK0/Swiss-Prot	Lysozyme C BioThesaurus	148	<u>Nasalis larvatus</u>	PIRSF001064	148	87	28-71	270	

(a)

Sequence Alignment Generated by Similarity Search



**FIGURE 1.12** Results obtained with the search: (a) details of proteins that have high sequence identity and (b) alignment of residues (see **Chapter 2**) for the two proteins that have high sequence identity.

computer-annotated supplement of SWISS-PROT that contains all the translations of EMBL nucleotide sequence entries, which are not yet integrated in SWISS-PROT. Currently, SWISS-PROT and TrEMBL have 0.5 and 7.6 million sequences, respectively. These databases are freely available at http://www.expasy.org/sprot/ and http://www.ebi.ac.uk/swissprot/.

SWISS-PROT contains the information about the name and origin of the protein, protein attributes, general information, ontologies, sequence annotation, amino acid sequence, bibliographic references, cross-references with sequence, structure and interaction databases, and entry information. An example for human lysozyme is shown in **Figure 1.13**. Furthermore, it has several search options, including sequence retrieval system (SRS), full-text search, advanced search, or by description or identification number.

In bioinformatics, developing a dataset plays a key role for any analysis or prediction. One can easily develop the dataset of amino acid sequences using SWISS-PROT. For example, the procedure for retrieving the data of "transcription factors" is shown in **Figure 1.14a**. Searching in UniProtKB with the keyword "transcription factors" will display all the relevant entries deposited in SWISS-PROT and TrEMBL.

Image: Sequence status       Complete         Protein attributes       Eukaryota > Metazoa > Condata > Craniata > Carliata > Euteriata > Euteria > Eute	★ Reviewed, UniProtK Last modified July 28, 2009. Vers	B/Swiss-Prot P61626 (LYSC_HUMAN)
Names and origin       Protein attributes       General annotation (Comments)       Ontologies       Einsry interactions       Sequence annotation (Features)       Sequences       References       Cross-references       Ertry         Names and origin       Protein names       Recommended name: Lysozyme C EC=3.21.17       Sequence       EC=3.21.17         Afternative name(s): 1.4-beta-N-acetylmuramidase C       Afternative name(s): 1.4-beta-N-acetylmuramidase C       Sequence	55 Clusters with 100%, 90%, 50	6 identity   🗅 Documents (7)   🍘 Third-party data   🗟 Customize display
Names and origin         Hide   Top           Protein names         Recommended name: Lysozyme C EC=3.21.17         Lysozyme C EC=3.21.17         EC=3.21.17           Alternative name(s): 1.4-beta-N-acetylmuramidase C         Synonyms: LZM         Corganism         Homo sapiens (Human) [Complete proteome]           Taxonomic identifier         9606 [NCBI]         Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Homo         Protein attributes         Hide   Top           Sequence length         148 AA.         Sequence status         Complete.         Sequence status         Complete.           Sequence status         Complete.         Evidence at protein level.         Hide   Top           Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Hide   Top	Names and origin · Protein attribute information · Relevant documents	s General annotation (Comments) Ontologies - Binary interactions - Sequence annotation (Features) - Sequences - References - Web resources - Cross-references - Entry
Protein names       Recommended name: Lysozyme C EC=3.2.1.17 Alternative name(s): 1,4-beta-N-acetylmuramidase C         Gene names       Name:       LYZ Synonyms: LZM         Organism       Homo sapiens (Human) [Complete proteome]         Taxonomic identifier       9606 [NCBI]         Taxonomic lineage       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Eukaryota > Metazoa > Complete.         Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hilde [Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acetylmuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Names and origin	Hide   Top
Gene names       Name:       LYZ         Synonyms:       LZM         Organism       Homo sapiens (Human) [Complete proteome]         Taxonomic identifier       9606 [NCBI]         Taxonomic lineage       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Hide [ Top         Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hide [ Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acetylmuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Protein names	Recommended name: Lysozyme C EC=3.2.1.17 Alternative name(s): 1,4-beta-N-acetylmuramidase C
Organism       Homo sapiens (Human) [Complete proteome]         Taxonomic identifier       9606 [NCBI]         Taxonomic lineage       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Hide   Top         Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hide   Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acety/Imuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Gene names	Name: LYZ Synonyms: LZM
Taxonomic identifier       9606 [NCBI]         Taxonomic lineage       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Hide   Top         Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hide   Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acety/muramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Organism	Homo sapiens (Human) [Complete proteome]
Taxonomic lineage       Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Homo         Protein attributes       Hide   Top         Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hide   Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acety/muramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Taxonomic identifier	9606 [NCBI]
Protein attributes         Hide   Top           Sequence length         148 AA.           Sequence status         Complete.           Sequence status         Complete.           Sequence processing         The displayed sequence is further processed into a mature form.           Protein existence         Evidence at protein level.           General annotation (Comments)         Hide   Top           Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.           Catalytic activity         Hydrolysis of (1->4)-beta-linkages between N-acety/Inuramic acid and N-acetyI-D-glucosamine residues in a peptidoglycan and between	Taxonomic lineage	Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Haplorrhini > Catarrhini > Hominidae > Horno
Sequence length       148 AA.         Sequence status       Complete.         Sequence processing       The displayed sequence is further processed into a mature form.         Protein existence       Evidence at protein level.         General annotation (Comments)       Hide [ Top         Function       Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Catalytic activity       Hydrolysis of (1->4)-beta-linkages between N-acety/Inuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Protein attributes	Hide   Top
Sequence status         Complete.           Sequence processing         The displayed sequence is further processed into a mature form.           Protein existence         Evidence at protein level.           General annotation (Comments)         Hide [ Top           Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.           Catalytic activity         Hydrolysis of (1->4)-beta-linkages between N-acety/Inuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Sequence length	148 AA.
Sequence processing         The displayed sequence is further processed into a mature form.           Protein existence         Evidence at protein level.           General annotation (Comments)         Hide [ Top           Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Hide [ Top           Catalytic activity         Hydrolysis of (1->4)-beta-linkages between N-acety/Inuramic acid and N-acetyI-D-glucosamine residues in a peptidoglycan and between	Sequence status	Complete.
Protein existence     Evidence at protein level.       General annotation (Comments)     Hide   Top       Function     Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.       Catalytic activity     Hydrolysis of (1->4)-beta-linkages between N-acety/Inuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Sequence processing	The displayed sequence is further processed into a mature form.
General annotation (Comments)         Hide   Top           Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.         Image: Catalytic activity           Catalytic activity         Hydrolysis of (1->4)-beta-linkages between N-acety/Image: acid and N-acetyI-D-glucosamine residues in a peptidoglycan and between	Protein existence	Evidence at protein level.
Function         Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.           Catalytic activity         Hydrolysis of (1->4)-beta-linkages between N-acetylmuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	General annotation (Comr	nents) Hide   Top
Catalytic activity Hydrolysis of (1->4)-beta-linkages between N-acetylmuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between	Function	Lysozymes have primarily a bacteriolytic function; those in tissues and body fluids are associated with the monocyte-macrophage system and enhance the activity of immunoagents.
N-acetvi-D-glucosamine residues in chitodextrins.	Catalytic activity	Hydrolysis of (1->4)-beta-linkages between N-acetylmuramic acid and N-acetyl-D-glucosamine residues in a peptidoglycan and between N-acetyl-D-glucosamine residues in chitodextrins.

FIGURE 1.13 Sample entry for human lysozyme in UniProtKB/SWISS-PROT.

The retrieval system has the options to restrict the data deposited only in SWISS-PROT and the data at different sequence identities, 90% or 50%. Furthermore, the final data can be downloaded in different file formats, such as Tab delimited, Excel, FASTA, GFF, Flat text, XML, RDF/XML, and list. The result obtained for the sequences of transcriptions factors in FASTA format is shown in **Figure 1.14b**. This set of sequences can be used for the analysis of DNA-binding proteins. In a similar way, sequences for any kind of proteins can be easily obtained with SWISS-PROT.

#### **1.4.3 UniProt: The Universal Protein Resource**

Recently, the SWISS-PROT, TrEMBL, and PIR protein database activities have united to form the Universal Protein Knowledgebase (UniProt) consortium. It provides the scientific community with a single, centralized, authoritative resource for protein sequences and functional information (Bairoch et al. 2005). The UniProt produces three layers of protein sequence databases: UniProt Archive, Knowledgebase, and Reference database. The UniProt Knowledgebase is a comprehensive, fully classified, richly and accurately annotated protein sequence knowledgebase with extensive cross-references. It is freely available at: http://www.uniprot.org/

# **1.4.4** Other protein sequence databases

There are several other protein sequence databases, which aim for specific classes or functions. EXProt (http://www.cmbi.kun.nl/EXProt/) is a nonredundant protein database containing a selection of entries from genome annotation projects and

uniProt → UniProtKB		Downloads · Contact · Documentation/Help
Search in Protein Knowledgebase (UniProtKB)	Cuery transcription factors	Search Clear Fields »
	Search	Blast Align Retrieve ID Mapping *

#### 1 - 25 of 11,439 results for transcription B AND factors in UniProtKB sorted by score descending

🚓 Browse by taxonomy, keyword, gene ontology, enzyme class or pathway | 👬 Reduce sequence redundancy to 100%, 90% or 50% | 📾 Customize Dow

- > Show only reviewed (UniProtKB/Swiss-Prot) or unreviewed (UniProtKB/TrEMBL) entries
- > Quote terms: "transcription factors"
- Restrict term "transcription" to protein family, gene ontology, keyword, protein name
- > Restrict term "factors" to protein family, gene ontology, protein name, web resource

Gene names None Entry name Protein names <sup>‡</sup> Organism Length <sup>‡</sup> Accession Status 8 P18850 ATF6A\_HUMAN Cyclic AMP-dependent ATF6 Homo sapiens 670 transcription factor ATF-6 (Human) alpha (cAMP-dependent transcription factor ATF-6 alpha) (Activating transcription factor 6 alpha) (ATF6-alpha) [Cleaved into: Processed cyclic AMP-dependent transcription factor ATF-6 alpha]  $\checkmark$ P0A4H2 BVGA\_BORPE 1 Virulence factors putative bvgA Bordetella 209 positive transcription (BP1878) pertussis regulator bygA  $\checkmark$  $\hat{\mathbf{n}}$ Cyclic AMP-dependent Homo sapiens Q9Y2D1 ATE5 HUMAN ATF5 (ATFX) 282 transcription factor ATF-5 (Human) (cAMP-dependent transcription factor ATF-5)

(a)

>sp|P18850|ATF6A HUMAN Cyclic AMP-dependent transcription factor ATF-6 alpha OS=Homo s MGEPAGVAGTMESPFSPGLFHRLDEDWDSALFAELGYFTDTDELOLEAANETYENNFDNL DFDLDLMPWESDIWDINNQICTVKDIKAEPQPLSPASSSYSVSSPRSVDSYSSTQHVPEE LDLSSSSQMSPLSLYGENSNSLSSAEPLKEDKPVTGPRNKTENGLTPKKKIQVNSKPSIQ PKPLLLPAAPKTQTNSSVPAKTIIIQTVPTLMPLAKQQPIISLQPAPTKGQTVLLSQPTV VQLQAPGVLPSAQPVLAVAGGVTQLPNHVVNVVPAPSANSPVNGKLSVTKPVLQSTMRNV GSDIAVLRRQQRMIKNRESACQSRKKKKEYMLGLEARLKAALSENEQLKKENGTLKRQLD  ${\tt EVVSEN} QRLKVPSPKRRVVCVMIVLAFIILNYGPMSMLEQDSRRMNPSVSPANQRRHLLG$ FSAKEAQDTSDGIIQKNSYRYDHSVSNDKALMVLTEEPLLYIPPPPCQPLINTTESLRLN HELRGWVHRHEVERTKSRRMTNNQQKTRILQGALEQGSNSQLMAVQYTETTSSISRNSGS ELQVYYASPRSYQDFFEAIRRRGDTFYVVSFRRDHLLLPATTHNKTTRPKMSIVLPAINI NENVINGQDYEVMMQIDCQVMDTRILHIKSSSVPPYLRDQQRNQTNTFFGSPPAATEATH VVSTIPESLO >sp|PDA4H2|BVGA BORPE Virulence factors putative positive transcription regulator bvgA MYNKVLIIDDHPVLRFAVRVLMEKEGFEVIGETDNGIDGLKIAREKIPNLVVLDIGIPKL DGLEVIARLQSLGLPLRVLVLTGQPPSLFARRCLNSGAAGFVCKHENLHEVINAAKAVMA GYTYFPSTTLSEMRMGDNAKSDSTLISVLSNRELTVLQLLAQGMSNKDIADSMFLSNKTV STYKTRLLQKLNATSLVELIDLAKRNNLA >sp/Q9Y2D1/ATF5 HUMAN Cyclic AMP-dependent transcription factor ATF-5 OS=Homo sapiens MSLLATLGLELDRALLPASGLGWLVDYGKLPPAPAPLAPYEVLGGALEGGLPVGGEPLAG DGFSDWMTERVDFTALLPLEPPLPPGTLPQPSPTPPDLEAMASLLKKELEQMEDFFLDAP  ${\tt LLPPPSPPPLPPPPLPPAPSLPLSLPSFDLPQPPVLDTLDLLAIYCRNEAGQEEVGMPPL$ PPPQQPPPPSPPQPSRLAPYPHPATTRGDRKQKKRDQNKSAALRYRQRKRAEGEALEGEC QGLEARNRELKERAESVEREIQYVKDLLIEVYKARSQRTRSC

**FIGURE 1.14** Sequence retrieval in UniProtKB/SWISS-PROT: (a) The query and retrieved results and (b) sequences of "transcription factors" in FASTA format (the first line starts with ">" followed by amino acid sequences in single-letter code. Each line has 60 amino acid residues; see **Chapter 2**).

oad...

of 458 | Next »

Page 1

public databases, aiming at including only proteins with an experimentally verified function. The NCBI Entrez Protein database (http://www.ncbi.nlm.nih.gov/ entrez) comprises sequences taken from a variety of sources, including SWISS-PROT, PIR, the Protein Research Foundation, the Protein Data Bank, and translations from annotated coding regions in the GenBank and RefSeq databases. Protein sequence records in Entrez have links to precomputed protein BLAST alignments, protein structures, conserved protein domains, nucleotide sequences, genomes, and genes. The Transport Classification Database (TCDB) is a curated, relational database containing sequence, classification, structural, functional, and evolutionary information about transport systems from a variety of living organisms (Busch and Saier, 2002; http://www.tcdb.org/).

Pongor's group (Pongor et al. 1993; Vlahovicek et al. 2005) developed a database of annotated protein domain sequences, SBASE. It facilitates the detection of domain homologies based on direct sequence database search using BLAST or a highspeed Smith Waterman algorithm, and returns a predicted domain architecture of the query sequence. Unlike traditional consensus representations, such as HMM (hidden Markov models), profiles and regular expressions, the SBASE domain library approach gives equal weights to all representatives, and a search against this library will detect both the typical and atypical (rare) representatives. It is to be noted that the domain library approach does not require either multiple alignment or learning algorithms to achieve accuracy. SBASE is freely available at http://www.icgeb.trieste.it/sbase.

# **1.5** Protein structure databases

Kendrew et al. (1958) solved the first three-dimensional structure of the protein myoglobin using X-ray crystallography. Subsequently, several structures have been determined with this technique. After two decades, NMR spectroscopy has been used to determine the protein structures in solution. Recently, neutron diffraction and electron microscopy have been used to determine protein structures. The number of protein structures has been increased every year, and they are deposited in a database, Protein Data Bank (PDB). On the basis of the structures available in PDB, several other databases have been established for the structural classification of proteins (SCOP), topology and architecture, and amino acid properties, etc.

#### **1.5.1** Protein Data Bank

The PDB was established at Brookhaven National Laboratories, USA, in 1971 as an archive for biological macromolecular crystal structures (Bernstein et al. 1977). Recently, the management of PDB became the responsibility of the Research Collaboratory for Structural Bioinformatics (RCSB), and it has six mirror sites at San Diego Supercomputer Center and Rutgers University in the USA; Cambridge Crystallographic Data Center, UK; National University of Singapore; Osaka University, Japan; and Max Delbruck Center for Molecular Medicine, Germany (Berman et al. 2000). PDB is available at http://www.rcsb.org/.

PDB stores the data in a uniform format atomic coordinates and partial bond connectivities, as derived from crystallographic studies. Text included in each data entry gives pertinent information for the structure at hand (e.g., species from which



**FIGURE 1.15** Snapshot showing the details of Protein Data Bank. The data for human lysozyme are shown.

the molecule has been obtained, resolution of diffraction data, literature citations and specifications of secondary structure). In addition to atomic coordinates and connectivities, the PDB stores the temperature factor for each atom. The PDB has been widely used in structural analysis on various aspects, such as atomic and residue contacts, amino acid clusters, developing potentials, and amino acid properties.

Currently, PDB has more than 57,000 structures, and 53,000 of them are proteins and their complexes (May 12, 2009). PDB has the search option with its code, authors, or full-text search. It contains the summary information about the name and source of the protein, experimental method, authors, and references. In addition, the details about crystallographic conditions are also provided. An example to human lysozyme is shown in **Figure 1.15**. Furthermore, it has the options to view the structure, display and download the files with/without three-dimensional coordinates, and it provides information about structural neighbors, geometry, sequence details, and other sources.

The sequence provides the detailed information about the number of amino acid residues in the protein, chain information, and amino acid sequence for each chain along with their secondary structures derived using DSSP (Kabsch and Sander, 1983), and number and percentage of each secondary structure through external links. For example, human lysozyme has four  $\alpha$ -helices and three  $\beta$ -strands, and the content of  $\alpha$ -helices is about 31%. Furthermore, it has the facility to save the sequence in FASTA format, which can be used in other programs.

PDB has the option to see the structural coordinates and view the structure on the screen or to save the structural information in a disk. It also has the possibility to download the data by ftp (ftp://ftp.wwpdb.org/pub/pdb/) or request the structures in CD. The structures can also be visualized using other programs, SWISS-PDB viewer (http://spdbv.vital-it.ch/), RASMOL (Sayle and Milner-White, 1995; http://www.umass.edu/microbio/rasmol/), Jmol (http://jmol.sourceforge.net/), KiNG (http://kinemage.biochem.duke.edu/software/king.php), PyMOL (De-Lano, 2002; http://www.pymol.org), etc. Structural data have the information about the protein, references, refinement details, sequence and secondary structure information, translation and rotation matrices for oligomers, disulfide bonds, and the atomic positions. The atomic position has the X, Y, and Z coordinates and temperature factors (**Figure 1.16**). The atomic coordinates have been extensively used to analyze the principles governing the folding and stability of protein structures, mechanism of protein folding and its interactions with other molecules, etc.

In addition, PDB has plenty of external links to other software for structural analysis and verification, modeling and simulation, molecular graphics, and so on.

HEADER TITLE COMPND COMPND COMPND COMPND SOURCE SOURCE SOURCE SOURCE KEYWDS EXPDTA AUTHOR	HYDROLASE (O-GLYCOSYL) 12-OCT-84 1L21 REFINEMENT OF HUMAN LYSOZYME AT 1.5 ANGSTROMS RESOLUTION. 2 ANALYSIS OF NON-BONDED AND HYDROGEN-BOND INTERACTIONS MOL ID: 1; 2 MOLECULE: HUMAN LYSOZYME; 3 CHAIN: A; 4 EC: 3.2.1.17; 5 ENGINEERED: YES MOL ID: 1; 2 ORGANISM SCIENTIFIC: HOMO SAPIENS; 3 ORGANISM SCIENTIFIC: HOMO SAPIENS; 3 ORGANISM TAXID: 9606 HYDROLASE (O-GLYCOSYL); X-RAY DIFFRACTION P.J.ARTYMIUK,C.C.F.BLAKE 2 24.0EP.00 LIZI VERN	
REVDAT	2 01-APR-03 1LZ1 1 JRNL	X Y 7
REVDAT JRNL JRNL JRNL JRNL JRNL JRNL JRNL JRNL		NTCM         1         N LVS A         1         1.933         12.0611         21.051         1.00         13.61         N           NTCM         2         CA LVS A         1         1.933         12.0611         21.051         1.00         13.61         N           NTCM         3         C LVS A         1         3.971         19.943         20.6611         21.051         1.00         13.31         CC           NTCM         4         C LVS A         1         3.971         19.943         20.432         1.00         14.70         CC           NTCM         5         CE LVS A         1         2.255         21.919         21.722         1.00         18.86         CC           NTCM         6         CG LVS A         1         3.661         22.552         21.419         1.00         12.65         CC         CC         NTCM         10.18.86         CC         CC         CC         NTCM         13.611         23.957         22.071         1.00         12.032.40         CC         NTCM         10.022.40         CC         NTCM         10.07.15         NTCM         13.00         1.01         15.15         N         N         NTCM         11.00
SHEET	3 B 3 ILE A 59 SER A 61 -1 N SER A 61 0 THR A 52 A	ATOM 31 C GLUA 4 10.892 24.090 17.810 1.00 14.24 C ATOM 32 CB GLUA 4 9.244 24.965 15.198 1.00 35.30 C
SSBOND	1 CYS A 6 CYS A 128 1555 1555 2.07 A 2 CYS A 30 CYS A 116 1555 1555 2.06 A	ATOM 33 CG GLUA 4 7.835 25.412 14.765 1.00 71.50 C
SSBOND	3 CYS A 65 CYS A 81 1555 1555 2.08 A	ATOM 34 CD GLUA 4 7.531 24.964 13.264 1.00 94.85 C ATOM 35 OE1 GLUA 4 8.204 25.884 12.678 1.00 76.94 C
SSBOND	4 CYS A 77 CYS A 95 1555 2.04 A	ATOM 36 OE2 GLU A 4 7.053 23.956 12.978 1.00107.70 C
	(a)	(b) Temperature factor



#### **1.5.2** Database for nonredundant PDB structures

A dataset with several similar structures influences the characteristic features of such data, which cause a bias on the result. Hence, a dataset of nonredundant proteins is necessary for analyzing the features of protein structures. The level of redundancy may vary depending upon the nature of the problem and the number of proteins in the database. The ASTRAL database provides the PDB codes and amino acid sequences for two cutoff values, i.e., less than 40% and 95% sequence identities. Furthermore, it has the option of selecting the cutoff values from 10% to 100%. The description about the selection of nonredundant dataset of proteins with 25% sequence identity is shown in **Figure 1.17**. The PDB identifiers for this cutoff have been selected to display the results, and it is also possible to get the sequences for these protein structures. Furthermore, ASTRAL has the options to get the sequences with e-values, fold, superfamily, etc.

For a general analysis, one can get the information about the nonredundant structures from ASTRAL. Furthermore, the nonredundant dataset of any

# **ASTRAL SCOP** Genetic Domain Sequences 1.73

- Notes on changes in this release.
- All ASTRAL SCOP genetic domain sequences, based on PDB SEQRES records: <u>astral-scopdom-seqres-gd-all-1.73.fa</u> (27 MB)
- All ASTRAL SCOP genetic domain sequences, based on PDB ATOM records: <u>astral-scopdom-atom-gd-all-1.73.fa</u> (26 MB) These are not recommended unless you have a special need.
- Percentage identity filtered ASTRAL SCOP cenetic domain sequence subsets, based on PDB SEQRES records Get identifiers only v with less than 25 v percent identity (using "in both" criterion).
- E-value filtered ASTRAL SCOP genetic domain sequence subsets, based on PDB SEQRES records.
   Get identifiers only v with E-values >= 0.01 v (in a database size of 100,000,000 residues).
- SCOP filtered ASTRAL SCOP genetic domain sequence subsets, based on PDB SEQRES records
   Get identifiers only v which represent each Superfamily v in SCOP.

dlejga	
dlucsa	
d2dsxa1	Adapter a 12 1 1 (2) Graphin (Shumainian ashbaga (Grapha shumainian) (Mauti, 2721)
d1r6ja	<pre>&gt;dlejgagill.ll.i (A:) Crampin (Abyssinian Cabbage (Crampe abyssinica) [Taxid: 5/21]; ttccnsiversnfwuchadtabelcatutaciiingatchaduan</pre>
dlus0a	>dli71a g.14.1.1 (A:) Apoliopprotein A (Human (Homo sapiens), IV-7 variant [TaxId: 9606])
d2b97a1	dcyhgdgqsyrgsfsttvtgrtcqswssmtphwhqrtteyypnggltrnycrnpdaeirp
dlacia	wcytmdpsvrweycnltqcpvme
dladna	>dlh8pal g.14.1.2 (A:22-67) PDC-109, collagen-binding type II domain (Cow (Bos taurus) [TaxId: 9913])
dliuss	eecvipivymmrkniadcivngsingadyvgrwkydaddaya Nall6ia a 14 1 2 (B-275-331) Galatinaea B (MMD-0) tupe II modules (Human (Homo samiens) [Tavid: 06061)
div6zal	rlytrdgnadgkpcgfpfifgggsvacttdgrsdgvrwcattanydrdklfgfcpt
d2bEgo1	>dlrOri_ g.68.1.1 (I:) Ovomucoid domains (Turkey (Meleagris gallopavo) [TaxId: 9103]}
dilucar	vdcseypkpactleyrplcgsdnktygnkcnfcnavvesngtltlshfgkc
diwuna_	>dltgsi g.68.1.1 (I:) Secretory trypsin inhibitor {Pig (Sus scrofa) [TaxId: 9823]}
din55a_	tspdreatctsevsgcpklynpvcgtdgltysnecvlcsenkkrdtpvlldksgpc Ndliwie – c 69 1 1 (A:) Ascidian tymsin inbistor (See squirt (Helogumthie roretzi) [Tevtd: 7729])
dlnwza_	ahudte finlicavidakeehrmicalicebaafevanace
d1mc2a_	>dlhdla g.68.1.2 (A:) Serine proteinase inhibitor lekti (Human (Homo sapiens) [TaxId: 9606]}
d1p9ga_	knedgem Chefqafmkngklfcpqdkkffqsldgim finkcatckmilekeaksq
dlpjxa_	>d4sgbi_ g.69.1.1 (I:) Plant chymotrypsin inhibitor {Potato tuber (Solanum tuberosum) [TaxId: 4113]}
d1x8qa	pictnccagykgcnyysangaiicegqsdpkkpkacpincdphiayskopr
d1dy5a	Ausphar (A.1953) ranceaute spasmolytic putypeptide (Fig (Sus Sciola) [lakid. 9023])
d1g6xa	enhance of distance when a decode on and the theory when a decode of the theory of the second s
d1gwea	
d1muwa	

**FIGURE 1.17** Retrieval of nonredundant protein sequences from PDB using ASTRAL. It uses SCOP classification for the domains, and the data obtained with the sequence identity of less than 25% are shown.



**FIGURE 1.18** Retrieval of nonredundant protein structures using PISCES. The users have the feasibility of selecting the identity, resolution, etc. and provide own dataset of proteins. (a) The search options and (b) results.

specific class of proteins (e.g., all proteins, membrane proteins, DNA-binding proteins, protein–protein complexes, etc.) can be derived by matching the dataset of interest with ASTRAL nonredundant dataset. ASTRAL is available at http://astral.berkeley.edu/.

Wang and Dunbrack (2005) developed a sequence-culling database server, PISCES, for producing lists of nonredundant proteins from the PDB using entry- and chain-specific criteria and mutual sequence identity. It uses a combination of PSI-BLAST and structure-based alignments to determine sequence identities. An example is shown in **Figure 1.18**. It takes the PDB codes with chain information along with other conditions, such as % sequence identity, resolution, R-factor, inclusion of non–X-ray structures, number of amino acid residues, etc (**Figure 1.18a**), and sends the results via e-mail. The output has four files: (i) original id, (ii) culled id, (iii) FASTA format sequence, and (iv) similarity log. The result obtained for the culled sequences from the list of 34  $\beta$ -barrel membrane proteins is shown in **Figure 1.18b**. PISCES is available at http://dunbrack.fccc.edu/pisces/.

Noguchi and Akiyama (2003) developed a database of representative chains from PDB. The search options to retrieve nonredundant set of proteins are shown in **Figure 1.19a**. It has several features, including the limits with resolution, R-factor, number of residues, fragments, complex structures, and membrane proteins. Furthermore, the users have the options to use any cutoff for sequence identity and root-mean-square deviation between the sequences and structures, respectively as shown in **Figure 1.19b**. The results obtained with the search are displayed in

(based on PDB Rel.#2007_11_14 . s (Apr. 15, 2009 updated.) Eliminate and Sort Chains	99469 chains )		2		
					PDB-REPRDB
factors	apply constraints	5	threshold	priority	Database of representative protein chains from PDB
Resolution	ONo ⊙Yes	X > 3.0	will be eliminated.	1	by Tanotsu NOGUCHI and Yutaka AKIYAMA
R-factor	ONo ⊙Yes	X > 0.3	will be eliminated.	2	(Computational Biology Research Center)
number of chain break	⊙No ⊖Yes	x>0	will be eliminated.	3	Reference -> <u>About PDB-REPRDB</u>
ratio of non-standard residues	⊙No ⊖Yes	x > 0	%will be eliminated.	4	ID : EDR aster ID & abaia ID
ratio of residues with only CA coordinates	⊙No ⊙Yes	X>0	%will be eliminated.	5	(click to show the Protein 3D viewer)
e ratio of residues with only backbone coordinates	ONo ⊙Yes	X>0	%will be eliminated.	6	naa : the number of amino acids ( from SEQRES line of PDB )
number of residues	ONo ⊙Yes	X < 40	will be eliminated.		Rfao: R-fator
e include MUTANT	⊙No OYes			8	Methd : experimental method n sid : the number of residues with side chain coordinates
© COMPLEX	Oonly COMPLEX ⊙exclude COMP OAll	( LEX		9	n.bok the number of residues with backbone coordinates n.ca : the number of residues with DA coordinates n.naa : the number of non-standard amino acid residues
FRAGMENT	Oonly FRAGMEN ⊙exclude FRAGM OAII	1T MENT		10	brk: the number of chain breaks mutant: mutant or valid complex : complex or not
include NMR	⊙No ⊖Yes				scop : SOOF(ver.1.73) socs(concise classification strings)
include membrane proteins	⊙No OYes				
Make List Service status Reset this form	(a)				Full Leng Real Devices Your request ID is 090825-20262 esc: Threshold : ID% >= 20 %, RMSD <= 10 A esc
Submit Submit	e similarity : IC (%) e similarity : Ri 10 (A) e status Re	0% MSD or eset this fo	Dmax		$ \begin{array}{c c c c c c c c c c c c c c c c c c c $
	(b)				(c)

**FIGURE 1.19** Retrieval of nonredundant proteins with PDB-REPRDB. This server will have several options to retrieve protein sequences, including the sequence identity and RMSD between two structures. (a) The initial parameters, (b) similarities, and (c) results are shown in the figure.

**Figure 1.19c**. It provides several details, including the links to SCOP database (see **Section 1.5.4**). It is available at http://www.cbrc.jp/pdbreprdb/.

# 1.5.3 PDB-related databases

On the basis of PDB, several other databases have been developed for different sets of proteins, such as membrane proteins, protein-protein complexes, protein-nucleic acid complexes, and ligand binding proteins. Tusnady et al. (2005) developed a database of transmembrane proteins, PDBTM (http:// pdbtm.enzim.hu/), which includes the sequences and structures of redundant and nonredundant  $\alpha$ -helical and  $\beta$ -barrel membrane proteins along with their membrane-spanning segments. Jayasinghe et al. (2001) compiled the structures of known membrane proteins, MPtopo (http://blanco.biomol.uci.edu/mptopo), and classified them into several groups, such as monotopic, GPCRs, rhodopsins, and  $\beta$ -barrel membrane proteins. They have included the PDB codes, structures, and their respective references. Ikeda et al. (2003) developed a database of transmembrane protein topologies, TMPDB (http://bioinfo.si.hirosaki-u.ac.jp/ ~TMPDB/), which is based on the experimental evidences from X-ray crystallography, NMR spectroscopy, etc. Sarai and colleagues (An et al. 1998) developed a database for protein-nucleic acid complex structures, and these structures have been classified into different groups based on the recognition motif of proteins and DNA involved in the complex. Puvanendrampillai and Mitchell (2003) developed the Protein Ligand Database (PLD), which has the PDB codes for protein-ligand complexes along with binding information.



**FIGURE 1.20** Sample entry for human lysozyme in SCOP database. It is classified as  $\alpha + \beta$  protein. Furthermore, links are available for the different folds and families within this class.

#### **1.5.4 Databases for structural classification of proteins**

Murzin et al. (1995) constructed the SCOP database, which provides a detailed and comprehensive description of the structural and evolutionary relationships of the proteins of known structures. For each protein, the classification has the hierarchical levels, family, superfamily, fold, and structural class. An example for human lysozyme is shown in **Figure 1.20**. It belongs to the family of C-type lysozyme and the fold of lysozyme-like under  $\alpha + \beta$  structural class. The structure can be identified with a six-letter code (1lz1\_). The first four letters of the code (1lz1) represent the PDB code followed by the chain name (\_; it indicates there are no multiple chains in this protein) and domain name (\_). SCOP database has been linked from the PDB to obtain the structural class information of each protein directly. On the other hand, one can search the SCOP database for obtaining the structural class, fold, and domain information.

Orengo et al. (1997) developed a semiautomatic procedure for deriving a novel hierarchical classification of protein domain structures (CATH) and created a database in providing the class information for all the structures in PDB. The four main levels of CATH classification are protein class (C), architecture (A), topology (T), and homologous superfamily (H). Class is the simplest level, and it essentially describes the secondary structure composition of each domain. Architecture summarizes the shape revealed by the orientations of the secondary structure units,

Home Search Domain 1b:1A00       Carbo: V3_2.0 (chan         CATH Domain: 1lz1A00 mms       Carbo: V3_2.0 (chan         DB 1lz1, Chain A, Domain 0       Carbo: V3_2.0 (chan <ul> <li>CATH Code</li> <li>Many, Abbs</li> <li>Many</li></ul>				CATHV
ATH Domain: 1/z1A00 xxxx         DB 1/z1, Chain A, Domain 0         Image: CATH Code         Image: C	lome Search Domain 1lz1A00			CathDB: V3_2_0 (chang
DB 1121, Chain A, Domain 0       Image: CATH Code       Level Description       Links         1       Marthy Alpha       Image: Cath Code       Image: Cath Code         1.10       Orthogonal Bundle       Image: Cath Code       Image: Cath Code         1.10.530       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         1.10.530.0       Image: Cath Code       Image: Cath Code       Image: Cath Code         Structure       Segment boundaries for domain	ATH Domain: 1Iz1A00	XML		
CATH Code         Level Description         Links           1         Mainy Alpha         Initial State           1.10         Orthogonal Bundle         Initial State           1.10.530         Lysopinte         Initial State           1.10.530.10         Icensol         Icensol           1.10.530.10.12         Icensol         Icensol           1.10.530.10.12         Icensol         Icensol           1.10.530.10.12.1         Icensol         Icensol           1.10.530.10.12.1         Icensol         Icensol           1.10.530.10.12.1         Icensol         Icensol           1.10.530.10.12.1.1         Icensol         Icensol           1.10.530.10.12.1.1         Icensol         Icensol           1.10.530.10.12.1.1         Icensol         Icensol           Structure         Segment boundaries for domain 1t/1A00         Icensol	0B 1Iz1. Chain A. Domain 0			
CALL         Level Description         Links           1         Mainty Abria         Initial           21.10         Orthogonal Bundle         Initial           1.10.550         Orthogonal Bundle         Initial           1.10.550         LivisS20100         (Cenes0)           1.10.550.10.1         Icens20)         Icens20)           1.10.550.10.1.2.1         Icens20)         Icens20)           1.10.550.10.1.2.1         Icens20)         Icens20)           1.10.550.10.1.2.1.1         Icens20)         Icens20)           Structure         Segment boundaries for domain 1t/1400         Icens20)	CATUCAN	Laural Descaringtion	Linka	
110         Orthogonal Bundle           @ 1.10.530         Lrsozhme           @ 1.10.530.10         [Gene3D]           @ 1.10.530.10.1.2		Mainly Alpha	LINKS	<b>B</b>
Into 530       Lisson max         Into 10,530       Lisson max         Into 11,0530,10       Identify         Into 10,530,10,12       Identify         Into 10,530,10,12,11       Identify         Into 10,530,10,12,1,11       Identify         Introductor       Identify         Segment boundaries for domain 12/1A00       Identify	1.10	Orthogonal Bundle		
P. 1.10.530.10.1         [Cene3D]           P. 1.10.530.10.1.2	Q 1.10.530	Lysozyme		
1.10.530.10.1                  0.1.10.530.10.1.2                  0.1.10.530.10.1.2.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                  0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 0.1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.530.10.1.2.1.1                 1.10.10                 1.10.10 <td< td=""><td>. 1.10.530.10</td><td></td><td>[Gene3D]</td><td></td></td<>	. 1.10.530.10		[Gene3D]	
@ 1.10.530.10.1.2.1	1.10.530.10.1			
•••••••••••••••••••••••••	@ 1.10.530.10.1.2			
Interview     Iteration       Iteration     Iteration       Iteration     Iteration	U 1.10.530.10.1.2.1			0
Segment boundaries for domain 1/2/A00	D 1 10 530 10 1 2 1 1 1	1	[Gene3D]	1lz1A00
Structure Sequence History Segment boundaries for domain 1/z1A00			10000001	A
Structure   Sequence   History Segment boundaries for domain 1z1A00				
Segment boundaries for domain 1/21A00	tructure Sequence History			
	Segment boundaries for domain 1	:1A00		
			1 81 101	
The manual II a Manual Land Manual Land Manual Land Manual Land Manual	Demain ID Start Des Sten Des blam	a Lawath		
	11ZTAUU  1  130	130		

FIGURE 1.21 Structural classification of human lysozyme with CATH.

such as barrels and sandwiches. At the topology level, sequential connectivity is considered, such that members of the same architecture might have quite different topologies. The homologous superfamilies contain proteins with highly similar structures and functions. The CATH classification for human lysozyme is shown in **Figure 1.21**. It has the architecture of an orthogonal bundle, and it is designated an all- $\alpha$  class proteins in CATH. It may be possible to have different assignments in SCOP and CATH for some proteins, especially if the percentage of helical content is high and strand content is less or vice versa. In human lysozyme, the content of  $\alpha$ -helical structures is 31% and that of  $\beta$ -strands is 8%. On the basis of the high content of  $\alpha$ -helix it was classified as all- $\alpha$  proteins in CATH, and due to the presence of  $\alpha$ -helices and  $\beta$ -strands SCOP classified it as  $\alpha + \beta$  protein.

# **1.6 Literature databases**

Literature databases play a vital role to understand the current status of research on various fields of interest. At present, several databases provide the information with/without restrictions. These databases include PubMed, Google Scholar, Scopus, Science Citation Index, Chemical abstracts, and so on. The PubMed literature database covers most of the protein-related works published in international journals, and it has different features to search the database and display the results. It can be searched with author's name, year of publication, name of the journal, and key words. As an example, a search for the articles published in *Nature* about protein structure is shown in **Figure 1.22**. The journal name is specified with *Nature*. The search retrieved 3240 records, including 45 reviews. **Figure 1.22** shows the summary of results, and it is possible to display abstract, abstract plus with



**FIGURE 1.22** Retrieval of results using PubMed literature database. The search with keywords and journal name, display options, and number of retrieved records are indicated. The link to related articles is shown with an arrow.

key words, etc. It is also possible to see the related articles for all the references listed in the result page. Each hit is linked with its respective journal site, and one can get the complete article if the journal has open access options or the users have the license to access the full text of the journal. PubMed can be freely accessed at http://www.ncbi.nlm.nih.gov/pubmed/.

# 1.7 Exercises

**1.** Retrieve the sequence of E. coli Triose phosphate isomerase in FASTA format. *Hint:* Use the text search in iProClass, and search with Triose phosphate isomerase and E. coli; select the correct one among the displayed results; "show" the selected one; and click on "save the result" as FASTA.

- 2. Obtain the amino acid sequences of membrane proteins located in outer membrane using SWISS-PROT/UniProtKB database. *Hint:* use the keyword, "membrane protein" and location, "outer membrane." Click on Download and select the format. Number of hits is about 500.
- **3.** Find the three-dimensional coordinates of porin from Rhodobactor capsulatus and store the data. *Hint:* Search at the top of the page of PDB using the key words. Select the best

match (e.g., 3POR or 2POR). Click on Display files and select PDB text. Click on Download files to save the coordinates.

- **4.** Get the list of PDB codes for RNA-binding proteins. *Hint:* Search with RNA-binding proteins. In the result page, click on "Results ID List" to get the list of codes.
- **5.** Check the interacting partners (protein and DNA) in 1RXW. *Hint:* Use Jmol to view the protein–DNA complex, and find the chain information.
- **6.** Obtain the PDB codes that have ligands and the sequence identity is less than 30%.

*Hint:* Go to Advanced Search and search with "ligand" in Structure title. Select the sequence identity as 30%.

- **7.** How many helices are there in 2LZM? *Hint:* Check the PDB file of 2LZM and the presence of helices and strands.
- **8.** What is the average fluctuation for the residue Phe at position 4 in 2LZM? *Hint:* Identify the temperature factors of all atoms in Phe4, and get the average. The information is shown in **Figure 1.16**.
- **9.** List the major transport systems in Transport classification database? *Hint:* Check the TCDB database and click on TC system
- **10.** What is the major class of transporters? *Hint:* Find the number of proteins in each transporter, and find the one with the maximum number of proteins.
- **11.** Get the dataset of protein structures with less than 40% sequence identity. *Hint:* Click on ASTRAL 1.73 with less than 40% sequence identity.
- **12.** Get the dataset of protein structures with less than 20% sequence identity. *Hint:* Select the sequences and follow **Figure 1.17**.
- 13. Get the PDBs with less than 20% sequence identity using the IDs obtained for RNA-binding proteins (question 4).*Hint:* Follow Figure 1.18 using the list of PDB codes obtained as the answer to question 4.
- **14.** Obtain the list of nonmembrane proteins without mutation and not complexed with other molecules and the sequence identity of less than 25%. *Hint:* Use PDB-REPRDB and provide the necessary conditions.
- **15.** Get the nonredundant structures of β-barrel membrane proteins. *Hint:* Go to PDBTM and click on download. The codes are listed with the link at the end of the page.
- **16.** Retrieve the proteins with seven α-helical segments in PDBTM. *Hint:* Search with type of transmembrane protein, and select the number of segments.

**17.** Compare the folding types and number of proteins belonging to different structural classes of proteins.

*Hint:* SCOP has the information.

- **18.** What are the major folding types in four structural classes? *Hint:* Compare the folding types with number of entries.
- **19.** Compare the heme-binding proteins (family: globin) and contents in SCOP and CATH databases.

Hint: SCOP: All-Globin like-Globin; CATH: 1.10.490.10.

- **20.** Find the domain information for outer membrane protein TolC. *Hint:* search CATH database with FASTA sequence, which can be found from UniProtKB/SWISS-PROT (P02930) or PDB.
- **21.** Find the articles published about "protein interactions" and "shape complimentarity."

*Hint:* Search PubMed with relevant keywords.

**22.** Identify the papers published in the journal *Cell* about mitochondrial  $\beta$ -barrel membrane proteins.

*Hint:* Search with mitochondrial  $\beta$ -barrel membrane proteins AND *Cell* [Journal].

# References

- An J, Nakama T, Kubota Y, Sarai A. 3DinSight: an integrated relational database and search tool for the structure, function and properties of biomolecules. Bioinformatics. 1998;14(2):188–195.
- Anfinsen CB. Principles that govern the folding of protein chains. Science. 1973;181(96): 223–230.
- Bairoch A, Apweiler R, Wu CH, Barker WC, Boeckmann B, Ferro S, Gasteiger E, Huang H, Lopez R, Magrane M, Martin MJ, Natale DA, O'Donovan C, Redaschi N, Yeh LS. The Universal Protein Resource (UniProt). Nucleic Acids Res. 2005;33(Database issue):D154–159.
- Bairoch A, Apweiler R. (1996) The SWISS-PROT protein sequence data bank and its new supplement TREMBL. Nucleic Acids Res. 1996;24(1):21–25.
- Barker WC, Garavelli JS, Huang H, McGarvey PB, Orcutt BC, Srinivasarao GY, Xiao C, Yeh LS, Ledley RS, Janda JF, Pfeiffer F, Mewes HW, Tsugita A, Wu C. The protein information resource (PIR). Nucleic Acids Res. 2000;28(1):41–44.
- Berman HM, Westbrook JZ, Feng G, Gilliland TN, Bhat H, Weissig IN, Shindyalov, Bourne PE. The Protein Data Bank. Nucleic Acids Res. 2000;28:235–242.
- Bernstein FC, Koetzle TF, Williams GJ, Meyer EF Jr, Brice MD, Rodgers JR, Kennard O, Shimanouchi T, Tasumi M. The Protein Data Bank: a computer-based archival file for macromolecular structures. J Mol Biol. 1977;112(3):535–542.
- Busch W, Saier MH Jr. The transporter classification (TC) system, 2002. Crit Rev Biochem Mol Biol. 2002;37:287–337.
- DeLano WL. The PyMOL Molecular Graphics System. San Carlos, CA: DeLano Scientific; 2002. Available online at http://www.pymol.org.
- Dayhoff MO, Eck RV, Chang MA, Sochard MR. Atlas of Protein Sequence and Structure, Vol. 1. Silver Spring, MD: National Biomedical Research Foundation; 1965.
- Gromiha MM, Selvaraj S. Inter-residue interactions in protein folding and stability. Prog Biophys Mol Biol. 2004;86(2):235–277.
- Ikeda M, Arai M, Okuno T, Shimizu T. TMPDB: a database of experimentallycharacterized transmembrane topologies. Nucleic Acids Res. 2003;31(1):406–409.

- Jayasinghe S, Hristova K, White SH. MPtopo: A database of membrane protein topology. Protein Sci. 2001;10(2):455–458.
- Kabsch W, Sander C, Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. Biopolymers. 1983;22:2577–2637.
- Kendrew JC, Bodo G, Dintzis HM, Parrish RG, Wyckoff H, Phillips DC. A threedimensional model of the myoglobin molecule obtained by x-ray analysis. Nature. 1958;181(4610):662–666.
- Levitt M, Chothia C. Structural patterns in globular proteins. Nature. 1976;261(5561): 552–558.
- Murzin AG, Brenner SE, Hubbard T, Chothia C. SCOP: a structural classification of proteins database for the investigation of sequences and structures. J Mol Biol. 1995;247(4):536–540.
- Nelson DL, Cox MM. Lehninger Principles of Biochemistry. New York: W.H. Freeman and Company; 2005.
- Noguchi T, Akiyama Y. PDB-REPRDB: a database of representative protein chains from the Protein Data Bank (PDB) in 2003. Nucleic Acids Res. 2003;31(1):492–493.
- Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, Thornton JM. CATH: a hierarchic classification of protein domain structures. Structure. 1997;5(8):1093–1108.
- Pauling L, Corey RB, Branson HR. The structure of proteins: two hydrogen-bonded helical configurations of the polypeptide chain. Proc Natl Acad Sci US A. 1951;37(4):205– 211.
- Pongor S, Skerl V, Cserzö M, Hátsági Z, Simon G, Bevilacqua V. The SBASE domain library: a collection of annotated protein segments. Protein Eng. 1993;6(4):391–395.
- Puvanendrampillai D, Mitchell JB. Protein Ligand Database (PLD): additional understanding of the nature and specificity of protein–ligand complexes. Bioinformatics. 2003;19(14):1856–1857.
- Ramachandran GN, Ramakrishnan C, Sasisekharan V. Stereochemistry of polypeptide chain configurations. J Mol Biol. 1963;7:95–99.
- Sayle RA, Milner-White EJ. RASMOL: biomolecular graphics for all. Trends Biochem Sci. 1995;20(9):374.
- Schulz GE. Transmembrane beta-barrel proteins. Adv Protein Chem. 2003;63:47–70.
- Tusnady GE, Dosztanyi Z, Simon I. PDB\_TM: selection and membrane localization of transmembrane proteins in the protein data bank. Nucleic Acids Res. 2005;33(Database issue):D275–278.
- Vlahovicek K, Kaján L, Agoston V, Pongor S. The SBASE domain sequence resource, release 12: prediction of protein domain-architecture using support vector machines. Nucleic Acids Res. 2005;33(Database issue):D223–D225.
- Wang G, Dunbrack RL Jr. PISCES: a protein sequence culling server. Bioinformatics. 2003;19:1589–1591.
- White SH, Wimley WC. Membrane protein folding and stability: physical principles. Annu Rev Biophys Biomol Struct. 1999;28:319–365.