Protein Structural Representation

Developing algorithms for protein structural analysis and prediction requires a way to computationally model and represent a protein structure. In order to construct efficient, maintainable software to deal with and manipulate protein structures, a suitable way to store these structures has to be adopted. Depending on the ultimate application, different representations may have advantages and disadvantages from a software perspective. For example, when designing a simple visualization software, the Cartesian (x, y, z) coordinates of each atom are useful and simple to render on the screen. However, if the program is to manipulate bond angles and bond lengths for example, a representation based on the internal degrees of freedom (see below) may be more appropriate. Some applications may even need to store more than one representation at a time; for example a simulation program that needs to compute a protein's Potential Energy, which is a function of both Cartesian and Internal coordinates, would benefit from keeping both representations at the same time. The structure of a protein is the set of atoms it contains, and the bonds that join them, that is, its inherent connectivity. A particular geometric shape of a protein (that is, the spatial arrangement of the atoms in the molecule) is called its conformation. Thus, a given protein structure can have many different conformations. Next, we discuss the two most common ways to model protein structures and conformations for software applications: Cartesian and Dihedral representations.

1 The Protein Databank (PDB)

Now that we have reviewed the basics of protein chemistry, let us turn our attention to the tools. The most important source of information about protein structure is the Protein Databank (PDB). maintained by the Research Collaboratory for Structural Bioinformatics (RCSB). In addition to being an entry point to the structural data itself, the PDB web site, http://www.rcsb.org/pdb, contains links to many tools database you can apply to individual protein structures as you search the database. Information from the database is made available through the Protein Structure Explorer interface. For each protein, you can view the molecular structure using 3D display tools such as JMol and the Java QuickPDB viewer. PDB files and file headers can be viewed as HTML and downloaded in a variety of formats. Links to the protein structure classification databases CATH, FSSP, and SCOP are provided, along with the tools CE (Combinatorial Extension) and VAST (Vector Alignment Search Tool), which search for structures based on structural alignment. Average geometric properties, including dihedral angles, bond angles, and bond lengths can be displayed in tabular format with extremes and deviations noted. Sequences can be viewed and labeled according to secondary structure, and sequence information downloaded in FASTA format. You can go directly to the page for a particular protein of interest by entering that protein's fourletter PDB code in the Explore box on the PDB's main page. The PDB can also be searched using two different search tools, SearchLite and SearchFields. SearchLite is a simple search tool that allows you to enter one or more search terms separated by boolean operators into a single search field. SearchFields is a tool for advanced searches that provides a customizable search form that allows you to use separate keywords to search each PDB header field.

SearchFields supports options for searching a dozen of the most important fields in the PDB header, as well as crystallographic information. SearchFields also allows the database to be searched using FASTA for sequence comparison, as well as secondary structure features or short sequence features. From the individual protein page generated by the Structure Explorer, the PDB provides

a menu of links through which to connect to other tools. These features are still evolving rapidly. Table 9-2 provides a brief overview of the PDB protein page. We also encourage you to explore the PDB site regularly if you are interested in tools for protein structure analysis.

1.1 The PDB file format

Every PDB entry contains a header with information about the molecule. This information includes:

- The name of the protein and what species it came from.
- Details about how the structure was determined X-ray crystallography, NMR etc., the resolution and other experimental and chemical details.
- A literature reference.
- The amino acid sequence and secondary structure information.
- Disulphide bridges.

The header is followed by the structural information itself. The most essential information for modeling a protein structure is the relative position of each atom, given as (x, y, z) Cartesian coordinates. Popular imaging methods such as X-Ray Crystallography, Nuclear Magnetic Resonance (NMR) and Cryogenic Electron Microscopy (Cryo-EM) are used to experimentally obtain relative atom positions from protein crystals and solutions. This is precisely the information provided by Protein Databank (PDB) format coordinate files. PDB format consists of lines of information in a text file. Each line of information in the file is called a record. A PDB file generally contains several different types of records, arranged in a specific order to describe a structure. At the top of the file is an optional header which contains information about the structure: The names of molecules, primary and secondary structure information, sequence database references, where appropriate, and ligand and biological assembly information, details about data collection and structure solution, and bibliographic citations. The header is followed by records that describe, line by line, the atomic coordinates of a protein molecule.

The atomic coordinates in a PDB file are shown in Figure 2. The meaning of each record is given in Figure 3

Size of the PDB over the years: The PDB was established in 1971. The number of deposited structures began to increase dramatically due to improved experimental techniques. In October 1998, the management of the PDB became the responsibility of the Research Collaboratory for Structural Bioinformatics (RCSB). As of 2017, there are over 135,000 structures in the PDB, with thousands of new structures being added every year.

The breakdown by experimental method shows that the majority of structures were obtained by X-ray crystallography. Out of the rest, most of the structures were obtained using NMR. Other methods are still not as popular.

PDB-101: In recent years the RCSB has provided extensive research and educational tools through the PDB-101 portal (https://pdb101.rcsb.org/). Most prominently, the "Molecule of the Month" series provides a curated introduction to the structures available in the PDB. It presents a short description of selected molecules from the PDB. Each article includes an introduction to

```
HEADER
          CHROMOSOMAL PROTEIN
                                                   02-JAN-87
                                                               1UB0
          STRUCTURE OF UBIQUITIN REFINED AT 1.8 ANGSTROMS RESOLUTION
TITLE
COMPND
          MOL_ID: 1;
         2 MOLECULE: UBIQUITIN;
COMPND
COMPND
         3 CHAIN: A;
COMPND
         4 ENGINEERED: YES
SOURCE
          MOL_ID: 1;
         2 ORGANISM_SCIENTIFIC: HOMO SAPIENS;
SOURCE
SOURCE
         3 ORGANISM_COMMON: HUMAN;
         4 ORGANISM_TAXID: 9606
SOURCE
KEYWDS
          CHROMOSOMAL PROTEIN
EXPDTA
          X-RAY DIFFRACTION
AUTHOR
          S.VIJAY-KUMAR, C.E.BUGG, W.J.COOK
REVDAT
         5
             09-MAR-11 1UBQ
                               1
                                        REMARK
REVDAT
         4
             24-FEB-09 1UBQ
                                1
                                        VERSN
REVDAT
         3
             01-APR-03 1UBQ
                                1
                                        JRNL
REVDAT
         2
             16-JUL-87 1UBQ
                                1
                                        JRNL
                                               REMARK
REVDAT
         1
             16-APR-87 1UBQ
                                Θ
JRNL
            AUTH
                   S.VIJAY-KUMAR, C.E.BUGG, W.J.COOK
JRNL
            TITL
                   STRUCTURE OF UBIQUITIN REFINED AT 1.8 A RESOLUTION.
JRNL
            REF
                   J.MOL.BIOL.
                                                  V. 194 531 1987
JRNL
            REFN
                                    ISSN 0022-2836
JRNL
            PMID
                   3041007
JRNL
            DOI
                   10.1016/0022-2836(87)90679-6
REMARK
         1
REMARK
         1 REFERENCE 1
REMARK
            AUTH
                   S.VIJAY-KUMAR, C.E.BUGG, K.D.WILKINSON, R.D.VIERSTRA,
         1
REMARK
         1
            AUTH 2 P.M.HATFIELD, W.J.COOK
REMARK
         1
            TITL
                   COMPARISON OF THE THREE-DIMENSIONAL STRUCTURES OF HUMAN,
REMARK
         1
            TITL 2 YEAST, AND OAT UBIQUITIN
REMARK
            REF
                   J.BIOL.CHEM.
         1
                                                  V. 262 6396 1987
REMARK
            REFN
                                    ISSN 0021-9258
         1
REMARK
         1
           REFERENCE 2
            AUTH
REMARK
                   S.VIJAY-KUMAR, C.E.BUGG, K.D.WILKINSON, W.J.COOK
         1
REMARK
         1
            TITL
                   THREE-DIMENSIONAL STRUCTURE OF UBIQUITIN AT 2.8 ANGSTROMS
REMARK
            TITL 2 RESOLUTION
         1
REMARK
         1
            REF
                   PROC.NATL.ACAD.SCI.USA
                                                  V. 82 3582 1985
REMARK
         1
            REFN
                                    ISSN 0027-8424
REMARK
           REFERENCE 3
         1
REMARK
                   W.J.COOK, F.L.SUDDATH, C.E.BUGG, G.GOLDSTEIN
         1
           AUTH
REMARK
                   CRYSTALLIZATION AND PRELIMINARY X-RAY INVESTIGATION OF
         1
            TITL
                   UBIQUITIN, A NON-HISTONE CHROMOSOMAL PROTEIN
REMARK
            TITL 2
         1
REMARK
         1
            REF
                   J.MOL.BIOL.
                                                  V. 130
                                                           353 1979
REMARK
            REFN
                                    ISSN 0022-2836
         1
REMARK
           REFERENCE 4
         1
REMARK
         1
            AUTH
                   D.H.SCHLESINGER, G.GOLDSTEIN
REMARK
            TITL
                   MOLECULAR CONSERVATION OF 74 AMINO ACID SEQUENCE OF
         1
REMARK
                   UBIQUITIN BETWEEN CATTLE AND MAN
         1
            TITL 2
                                                  V. 255
REMARK
         1
            REF
                   NATURE
                                                           423 1975
REMARK
         1
                                    ISSN 0028-0836
            REFN
REMARK
REMARK
         2 RESOLUTION.
                          1.80 ANGSTROMS.
```

Figure 1: The beginning of the header of Ubiquitin (PDB:1UBQ)

ATOM	1	N	PRO	А	1	-3.190	7.728	33.820	1.00 21.66	N
ATOM	2	CA	PRO	А	1	-2.220	6.922	34.499	1.00 18.48	С
ATOM	3	С	PRO	А	1	-0.802	7.080	34.031	1.00 17.67	С
ATOM	4	0	PRO	А	1	-0.530	7.806	33.045	1.00 18.49	0
ATOM	5	СВ	PRO	А	1	-2.727	5.495	34.165	1.00 20.72	С
ATOM	6	CG	PRO	А	1	-3.834	5.651	33.165	1.00 20.84	С
ATOM	7	CD	PRO	А	1	-4.438	7.016	33.499	1.00 19.67	C
ATOM	8	N	GLN	А	2	0.091	6.450	34.755	1.00 14.65	N
ATOM	9	CA	GLN	А	2	1.526	6.384	34.480	1.00 17.51	С
ATOM	10	С	GLN	А	2	1.753	4.880	34.129	1.00 18.83	С
ATOM	11	0	GLN	А	2	1.442	3.982	34.963	1.00 19.75	0
ATOM	12	СВ	GLN	А	2	2.519	6.960	35.431	1.00 17.46	С
ATOM	13	CG	GLN	А	2	3.943	6.608	35.023	1.00 20.07	C
ATOM	14	CD	GLN	А	2	4.890	7.376	35.931	1.00 26.75	С
ATOM	15	OE1	GLN	А	2	5.366	6.856	36.946	1.00 31.80	0
ATOM	16	NE2	GLN	А	2	5.172	8.611	35.545	1.00 29.41	N

Figure 2: An example of atomic coordinates in the PDB

						Chain :	name			
	Amino	Aci	id			/ Seq	uence Nu	umber		
			N			1 1				
	Eleme	ent	<u>۱</u>			/ /	Co	ordinate	s	
		\ \	<u>۱</u>	<u>۱</u>	1	1	X	Y	Z	(etc.)
ATOM		1	N	ASP	L	1	4.060	7.307	5.186	•••
ATOM		2	CA	ASP	L	1	4.042	7.776	6.553	•••
ATOM		3	С	ASP	L	1	2.668	8.426	6.644	•••
ATOM		4	0	ASP	L	1	1.987	8.438	5.606	•••
ATOM		5	СВ	ASP	L	1	5.090	8.827	6.797	•••
ATOM		6	CG	ASP	L	1	6.338	8.761	5.929	•••
ATOM		7	OD1	ASP	L	1	6.576	9.758	5.241	
ATOM		8	OD2	ASP	L	1	7.065	7.759	5.948	•••
			- N	Ν						
			1	Elem	ent	position	within	amino ac	id	

Atomic Coordinates: PDB Format

Figure 3: An example of atomic coordinates in the PDB

Table 1:	Breakdown of	f structures o	on the PDB	as of 2017	(from the RCSB website)
					1	

Experimental Method	Proteins	Nucleic Acids	Protein/NA Complexes	Other	Total
X-RAY	113840	1906	5818	4	121568
NMR	10571	1229	246	8	12054
ELECTRON MICROSCOPY	1328	30	473	0	1831
HYBRID	105	3	2	1	111
other	200	4	6	13	223
Total	126044	3172	6545	26	135787



Figure 4: Number of structures in the PDB over the years, as of 2017

the structure and function of the molecule, a discussion of the relevance of the molecule to human health and welfare, and suggestions for how visitors might view these structures and access further details.

1.2 The wwPDB

In recent years, the major database for macromolecular structures is the worldwide PDB (wwPDB) at http://www.wwpdb.org/. It is a joint effort of the RCSB, the Protein Data Bank Europe (at the European Bioinformatics Institute, EBI), the Protein Databank Japan (based at Osaka University), and the Biological Magnetic Resonance Data Bank (BMRB).

2 Protein Representation

Computer graphics, scientific visualization and geometry to create a 3-D visual model of molecular structures. This facilitates structure, dynamic and function analysis. The simplest way to represent a protein chain is to store the Cartesian (x, y, z) coordinates of each atom, as they appear in the PDB file. These coordinates are relative to some arbitrary global coordinate frame. The specific frame is unimportant, for example that in which the atomic positions were obtained by X-Ray crystallography and which are typically read from the PDB files, as shown in Figure 2. These coordinates can be subjected to geometric transformations, if so desired. Common changes are to remove the center of mass (thus centering the protein at the global origin), subtract the position of the anchor atom (to center the protein at this atom), etc. Cartesian representation is good for visualization and analysis of the protein structure. Numerous ways are available for visualizing the structures stored in the PDB and other repositories, and many tools utilize computer graphics and scientific visualization to create an image of protein structures, which facilitates their analysis. The software reads the structural information provided by the PDB file. The most basic information contained in the file is the geometric location of each atom stored in the file. To visualize it, the



Figure 5: Different ways to visualize a molecule: a) Every dot represents the location of an atom. b) A wireframe representing atomic bonds. c) A molecular surface. d) A cartoon representation of the protein.

software has to simply draw a dot in the position of each atom, as seen in Figure 5 (a). Each dot is colored according to a coloring convention: Carbon atoms are colored in cyan, Nitrogen atoms in blue, Oxygen in red, Sulphur in yellow and Hydrogen (not shown) in white. The arbitrary axis system is also shown. While this is a possible way to depict a structure, it is usually not very informative. It does not show us the bonds between the atoms, the secondary structure elements, and it only gives us a general idea about the 3D structure of the protein.

Fortunately, most such tools also allow a detailed examination of the molecule in a variety of rendering modes. For example, we sometimes want a ball-and-stick representation to visualize each and every atom and bond in the structure (Figure 5 b). Alternatively, sometimes it may be useful to have a detailed image of the surface of the molecule as experienced by a molecule of water (Figure 5 c). For other purposes, a simple, cartoonish representation of the major structural features, such as secondary structures, may be sufficient (Figure 5 d). Many other options exist. We can combine several representations into the same image, we can also change the color to emphasize certain parts of the structure. The possibilities are (almost) endless!

Sometimes we want to use other ways to represent a molecular structure, such as Internal coordinates, using bond lengths, angles and dihedrals (more about them later...). It is possible to switch back and forth between cartesian and internal coordinates.

3 Structural Classification of Proteins

Protein structure classification is important because it gives you an entry point into the world of protein structure that is independent of sequence similarity. Proteins are grouped not by functional families, but according to what kind of secondary structure (alpha helix, beta sheet, or both) they have. Within those larger classes, subclasses are defined based on how the secondary structures in the protein are arranged. The focus in protein classification is on finding proteins that have similar chemical architectures; it doesn't matter if their sequences are related. Over the years, we've learned from classification that there are far fewer unique protein folds than there are protein sequence families. Protein chemists often are interested in the information that can be extracted from broader structural classes of proteins, since analyzing that information can help them better understand how proteins fold.

here isn't really a consensus as to how to classify protein structures quantitatively. Instead, structures end up in qualitatively named classes such as "greek key," "helix bundle," and "alphabeta barrel." These fold classes are useful in that they draw attention to prominent structural features and create a frame of reference for classifying structure. However, qualitative classifications don't lend themselves to automated analysis, and such protein classification databases still require the involvement of expert curators. If you're simply concerned with finding the close structural relatives of a published protein structure, there are a number of online classification databases in which existing structures have been annotated by a combination of automated analysis and input from protein structure experts. There are also automated tools for finding structural neighbors by structure alignment, though like any alignment method, these tools require you to understand the significance of comparison scores when analyzing results. If you are interested in doing your own analysis of a protein structure, there are several structure classification processes and tools that might help.

3.1 The SCOP Database

The Structural Classification of Proteins (SCOP) database is a largely manual classification of protein structural domains based on similarities of their structures and amino acid sequences [SCOP]. A motivation for this classification is to determine the evolutionary relationship between proteins. SCOP was created in 1994 in the Centre for Protein Engineering and the Laboratory of Molecular Biology.

The source of protein structures is the Protein Data Bank. The unit of classification of structure in SCOP is the protein *domain*. What the SCOP authors mean by "domain" is suggested by their statement that small proteins and most medium sized ones have just one domain, and by the observation that human hemoglobin, which has an $\alpha 2\beta 2$ structure, is assigned two SCOP domains, one for the α and one for the β subunit.

The shapes of domains are called "folds" in SCOP. Domains belonging to the same fold have the same major secondary structures in the same arrangement with the same topological connections. 1195 folds are given in SCOP version 1.75. Short descriptions of each fold are given. For example, the "globin-like" fold is described as core: 6 helices; folded leaf, partly opened. The fold to which a domain belongs is determined by inspection, rather than by software.

The hierarchical levels of SCOP are as follows.

- 1. Class: Types of folds, e.g., beta sheets.
- 2. Fold: The different shapes of domains within a class.
- 3. Superfamily: The domains in a fold are grouped into superfamilies, which have at least a distant common ancestor.
- 4. Family: The domains in a superfamily are grouped into families, which have a more recent common ancestor.
- 5. Protein domain: The domains in families are grouped into protein domains, which are essentially the same protein.



Figure 6: A SCOP2 graph structure example.

6. Species: The domains in "protein domains" are grouped according to species.

7. Domain: part of a protein. For simple proteins, it can be the entire protein.

Work on SCOP concluded in June 2009 with the release of SCOP 1.75. SCOP is still available online but is no longer maintained or updated. The prototype of a new Structural Classification of Proteins 2 (SCOP2) database has been made publicly available. SCOP2 defines a new approach to the classification of proteins that is essentially different from SCOP, but retains its best features [SCOP2]. Rather than using a simple tree hierarchy, the classification of proteins is described in terms of a directed acyclic graph in which each node defines a relationship of particular type and is exemplified by a region of protein structure and sequence. Importantly, there can be more than one parental node for a child node that allows multiple routes to a particular relationship. Figure 6 shows an example of a graph structure.

By 2009, the original SCOP database manually classified 38,000 PDB entries into a strictly hierarchical structure. With the accelerating pace of protein structure publications, the limited automation of classification could not keep up, leading to a non-comprehensive dataset. The Structural Classification of Proteins extended (SCOPe) database was released in 2012 with far greater automation of the same hierarchical system and is full backwards compatible with SCOP [SCOPe]. In 2014, manual curation was reintroduced into SCOPe to maintain accurate structure assignment. As of February 2015, SCOPe 2.05 classified 71,000 of the 110,000 total PDB entries. Figure 7 shows a screenshot of the SCOPe database.



Figure 7: The SCOPe database.

	10010 2:	The four main levels of efficit etassification
#	Level	Description
1	Class	Overall secondary-structure content of the domain.
		(Equivalent to SCOP class)
2	Architecture	High structural similarity but no evidence of homology.
		(Equivalent to SCOP fold)
3	Topology	A large-scale grouping of topologies which share
		particular structural features
4	Homologous superfamily	Indicative of a demonstrable evolutionary relationship.
		(Equivalent to SCOP superfamily)

Table 2: The four main levels of CATH classification

3.2 The CATH Database

CATH is another database which classifies protein structures downloaded from the Protein Data Bank. It is a semi-automatic, hierarchical classification of protein domains initially published in 1997. The name CATH is an acronym of the four main levels in the classification. The four main levels of the CATH hierarchy are as in Table 2

Much of the work is done by automatic methods, however there are important manual elements to the classification. The very first step is to separate the proteins into domains. It is difficult to produce an unequivocal definition of a domain and this is one area in which CATH and SCOP differ. The domains are automatically sorted into classes and clustered on the basis of sequence similarities. These groups form the **H** levels of the classification. The topology level is formed by structural comparisons of the homologous groups. Finally, the **A**rchitecture level is assigned manually. Class Level classification is done on the basis of 4 criteria:

1. Secondary structure content;

- 2. Secondary structure contacts;
- 3. Secondary structure alternation score; and
- 4. Percentage of parallel strands.

CATH defines four classes: mostly- α , mostly- β , α and β , few secondary structures.

4 Molecular Visualization Tools

Visualizing Protein Structures Numerous tools are available for visualizing the structures stored in the PDB and other repositories. Most such tools allow a detailed examination of the molecule in a variety of rendering modes. For example, sometimes it may be useful to have a detailed image of the surface of the molecule as experienced by a molecule of water. For other purposes, a simple, cartoonish representation of the major structural features may be sufficient.

4.1 A Few Molecular Visualization Programs

Visual Molecular Dynamics (VMD) [VMD] was originally developed for viewing molecular simulation trajectories. It is a very powerful, full-featured, and customizable molecular viewing package. Customization is available using Tcl/Tk scripting. Information on Tcl/Tk scripting can be found at this Tcl/Tk website: https://www.tcl.tk/. PyMol [Pymol] is an open-source molecular viewer that can be used to generate professional-looking images. PyMol is highly customizable through the Python scripting language. Protein Explorer [ProtExplorer] is an easy-to-use, web browser-based visualization tool. Protein explorer is built using the MDL Chime18 browser plugin, which in turn is based on the RasMol [RasMol] viewer. Because Chime only works under Windows and Macintosh OS, the use of Protein Explorer is restricted to those platforms. JMol [JMol] is a Java-based molecular viewer. In applet form, it can be downloaded on-the-fly to view structures from the web. A stand-alone version also exists, which can be used independently of a web browser. Chimera [Chimera] is a powerful visualizer and analysis tool that can be comfortably used with very large molecular complexes. It can also produce very high-quality images for use in presentations and publications. It provides a python interface.

References

- [SCOP] Loredana Lo Conte, Bart Ailey, Tim J. P. Hubbard, Steven E. Brenner, Alexey G. Murzin, and Cyrus Chothia. Scop: a structural classification of proteins database. *Nucleic Acids Research*, 28(1):257–259, 2000.
- [SCOP2] Antonina Andreeva, Dave Howorth, Cyrus Chothia, Eugene Kulesha, and Alexey G. Murzin. Scop2 prototype: a new approach to protein structure mining. Nucleic Acids Research, 2013.
- [SCOPe] Naomi K. Fox, Steven E. Brenner, and John-Marc Chandonia. Scope: Structural classification of proteins – extended, integrating scop and astral data and classification of new structures. Nucleic Acids Research, 42(D1):D304–D309, 2014. URL http: //nar.oxfordjournals.org/content/42/D1/D304.abstract.

- [JMol] Jmol: an open-source java viewer for chemical structures in 3d. URL http://www.jmol. org/.
- [Pymol] Schrödinger, LLC. The PyMOL molecular graphics system, version 1.3r1. August 2010.
- [RasMol] R. A. Sayle and Milner E. J. White. RASMOL: biomolecular graphics for all. Trends Biochem Sci, 20(9), 1995.
- [ProtExplorer] E. Martz. Protein explorer: easy yet powerful macromolecular visualization. Trends Biochem Sci, 27(2):107–109, 2002.
- [Chimera] E.F. Pettersen, T.D. Goddard, C.C. Huang, G.S. Couch, D.M. Greenblatt, E.C. Meng, and T.E. Ferrin. Ucsf chimera-a visualization system for exploratory research and analysis. J. Comput Chem., 25(13):1605–1612, 2004.
- [VMD] W. Humphrey, A. Dalke, and K. Schulten. VMD Visual Molecular Dynamics. J. Molec. Graphics, 14:33-38, 1996. URL http://www.ks.uiuc.edu/Research/vmd/.