

Introduction to Proteins

Dr. Fayyaz ul Amir Afsar Minhas

PIEAS Biomedical Informatics Research Lab Department of Computer and Information Sciences Pakistan Institute of Engineering & Applied Sciences PO Nilore, Islamabad, Pakistan http://faculty.pieas.edu.pk/fayyaz/

Importance of proteins

- 50% of the dry weight of the human body is protein
- Up to 92% of the dry weight of the red blood cell is a single protein called Hemoglobin
- Almost all cellular functions and biological functions involve proteins including
 - DNA Transcription
 - RNA Translation
 - Splicing
 - Muscular construction
 - Structural functions...

Binding



The TATA binding protein binds a specific DNA sequence and serves as the platform for a complex that initiates transcription of genetic information. (PDB 1tgh)



Myoglobin binds a molecule of oxygen reversibly to the iron atom in its heme group (shown in grey with the iron in green). It stores oxygen for use in muscle tissues. (PDB 1a6k)

3

Catalysis



HIV protease



Replication of the AIDS virus HIV depends on the action of a protein-cleaving enzyme called HIV protease. This enzyme is the target for protease-inhibitor drugs (shown in grey). (PDB 1a8k)

4

Switching

Signal cell growth



The GDP-bound ("off"; PDB 1pll) state of Ras differs significantly from the GTP-bound ("on"; PDB 121p) state. This difference causes the two states to be recognized by different proteins in signal transduction pathways.

Structural proteins



Its strength comes from the covalent and hydrogen bonds within each sheet; the flexibility from the van der Waals interactions that hold the sheets together. (PDB 1slk) Actin fibers are important for muscle contraction and for the cytoskeleton. They are helical assemblies of actin and actin-associated proteins. (Courtesy of Ken Holmes)

The most versatile cellular molecule

- Proteins are
 - Stable
 - Flexible
 - Versatile in terms of their function

- Protein Function
 - Biochemical function of the molecule in isolation
 - Cellular function as part of larger assembly
 - Resulting Phenotype in the cell or the organism

Uranium Binding Proteins



How are proteins versatile?

- Because of their structure!
- What determines their structure?
 - Composed of amino acids
 - Amino acids are linked to form polymer
 - Significant hydrogen bonding between amino acids
- Primary ideology
 - Sequence determines Structure determines
 Function
 - Proven by Anfinsen 1960

Anfinsen Experiment

- Native structure of the protein is determined by the protein's amino acid
- Principles for protein folding
- Native structure is
 - Unique
 - Stable
 - Kinetically accessible



Protein Composition



Understanding protein structure

- To understand protein structure, we will try to answer the following question
 - Why do proteins take any shape at all?
 - Why do proteins fold?
 - Can proteins take on arbitrary shapes?

The Central Dogma



Prokaryotes

Eukaryotes



Codons

		2nd p	osition	l	
lst position (5' end)	U	C	Α	G	3rd position (3' end)
U	Phe	Ser	Tyr	Cys	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	STOP	STOP	A
	Leu	Ser	STOP	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	GIn	Arg	A
	Leu	Pro	GIn	Arg	G
A	lle Ile Ile Met	Thr Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

Amino acids	Abbrevia	tions	Codons
Alanine	Ala	Α	GCA GCC GCG GCU
Cysteine	Cys	C	UGC UGU
Aspartic acid	Asp	D	GAC GAU
Glutamic acid	l Glu	E	GAA GAG
Phenylalanine	e Phe	F	UUC UUU
Glycine	Gly	G	GGA GGC GGG GGU
Histidine	His	H	CAC CAU
Isoleucine	lle	Ι	AUA AUC AUU
Lysine	Lys	K	AAA AAG
Leucine	Leu	L	UUA UUG CUA CUC CUG CUU
Methionine	Met	М	AUG
Asparagine	Asn	Ν	AAC AAU
Proline	Pro	Р	CCA CCC CCG CCU
Glutamine	Gln	Q	CAA CAG
Arginine	Arg	R	AGA AGG CGA CGC CGG CGU
Serine	Ser	S	AGC AGU UCA UCC UCG UCU
Threonine	Thr	Т	ACA ACC ACG ACU
Valine	Val	V	GUA GUC GUG GUU
Tryptophan	Trp	W	UGG
Tyrosine	Tyr	Y	UAC UAU

Amino Acids



The chemical structure of an amino acid. The backbone is the same for all amino acids and consists of the amino group (NH₂), the alpha carbon and the carboxylic acid group (COOH). Different amino acids are distinguished by their different side chains, R. The neutral form of an amino acid is shown: in solution at pH 7 the amino and carboxylic acid groups ionize, to NH₃⁺ and COO⁻. Except for glycine, where R=H, amino acids are chiral (that is, they have a left–right asymmetry). The form shown is the L-configuration, which is most common.

R | (NH₂)-CH-(COOH)



N-terminus

R^L | --(NH)-CH-(COOH)

C-terminus

● Hydrogen ● Carbon ● Oxygen ○ Sulfur ● Nitrogen | bond to functional group (R) | double bond || partial double bond || single bond

C-alpha is the carbon atom on which the functional group is attached

CIS 529: Bioinformatics

Peptide bond formation and hydrolysis



Extended peptide chain



Figure 1-8 Schematic diagram of an extended polypeptide chain The repeating backbone is shown, with schematized representations of the different side chains (R_1 , R_2 and so on). Each peptide bond is shown in a shaded box. Also shown are the individual dipole moments (arrows) associated with each bond. The dashed lines indicate the resonance of the peptide bond.

The carbonyl carbon, oxygen and nitrogen lie in the same plane CIS 529: Bioinformatics PIEAS Biomedical Informatics Research Lab



Types of aminc acids

Amino Acid 🕈	3- Letter ^[116] [≑]	1- Letter ^[116] ^{\$}	Side-chain polarity ^[116] +	Side-chain charge (pH ≑ 7.4) ^[116]	Hydropathy index ^[117]	MW(Weight) ^[119] ÷
Arginine	Arg	R	Basic polar	positive	-4.5	174
Lysine	Lys	К	Basic polar	positive	-3.9	146
Asparagine	Asn	N	polar	neutral	-3.5	132
Aspartic acid	Asp	D	acidic polar	negative	-3.5	133
Glutamic acid	Glu	E	acidic polar	negative	-3.5	147
Glutamine	Gln	Q	polar	neutral	-3.5	146
Histidine	His	н	Basic polar	positive(10%) neutral(90%)	-3.2	155
Proline	Pro	Р	nonpolar	neutral	-1.6	115
Tyrosine	Tyr	Y	polar	neutral	-1.3	181
Tryptophan	Trp	W	nonpolar	neutral	-0.9	204
Serine	Ser	S	polar	neutral	-0.8	105
Threonine	Thr	Т	polar	neutral	-0.7	119
Glycine	Gly	G	nonpolar	neutral	-0.4	75
Alanine	Ala	А	nonpolar	neutral	1.8	89
Methionine	Met	М	nonpolar	neutral	1.9	149
Cysteine	Cys	С	nonpolar	neutral	2.5	121
Phenylalanine	Phe	F	nonpolar	neutral	2.8	165
Leucine	Leu	L	nonpolar	neutral	3.8	131
Valine	Val	V	nonpolar	neutral	4.2	117
Isoleucine	lle	I	nonpolar	neutral	4.5	131

Properties of amino acids

- Different amino acids have different properties
- This allows all types of interactions between these amino acids
 - Charged amino acids would like to be close to oppositely charged ones and away from similarly charged ones
 - Hydrophobic ones don't like water
 - They will like to find other hydrophobic ones so that they can interact with them after displacing the water in between much like how two oil droplets in a cup of water merge to form a bigger oil droplet
 - This effect is called hydrophobic effect
- There are other types of interactions in the protein

Interactions between amino acids

Interactions between backbone atoms

Interactions between side chains

Interactions of amino acids

Chemical Interactions that Stabilize Polypeptides

Interaction	Example	Distance dependence	Typical distance	Free energy (bond dissociation enthalpies for the covalent bonds)
Covalent bond	-C _{\alpha} -C-	-	1.5 Å	356 kJ/mole (610 kJ/mole for a C=C bond)
Disulfide bond	-Cys-S-S-Cys-	-	2.2 Å	167 kJ/mole
Salt bridge	- c (0H-N-H - H+ 0 H	Donor (here N), and acceptor (here O) atoms <3.5 Å	2.8 Å	12.5–17 kJ/mole; may be as high as 30 kJ/mole for fully or partially buried salt bridges (see text), less if the salt bridge is external
Hydrogen bond	N-H0=C	Donor (here N), and acceptor (here O) atoms <3.5 Å	3.0 Å	2–6 kJ/mole in water; 12.5–21 kJ/mole if either donor or acceptor is charged
Long-range electrostatic interaction	- c 0 H-N-H 0 H 0 H	Depends on dielectric constant of medium. Screened by water. 1/r dependence	Variable	Depends on distance and environment. Can be very strong in nonpolar region but very weak in water
Van der Waals interaction	н н -С-н н-С- - I I н н	Short range. Falls off rapidly beyond 4 Å separation. 1/r ⁶ dependence	3.5 Å	4 kJ/mole (4–17 in protein interior) depending on the size of the group (for comparison, the average thermal energy of molecules at room temperature is 2.5 kJ/mole)

Interactions between backbone atoms

- Carbonyl oxygen and Amide Hydrogen form nearby amino acids in the chain can form oxygen bonds
 - This is energetically favorable
 - This gives rise to regular conformations
 - Alpha helices
 - Beta sheets

Alpha helices

Figure 1-13 The alpha helix The chain path with average helical parameters is indicated showing (a) the alpha carbons only, (b) the backbone fold with peptide dipoles and (c) the full structure with backbone hydrogen bonds in red. All three chains run from top to bottom (that is, the amino-terminal end is at the top). Note that the individual peptide dipoles align to produce a macrodipole with its positive end at the amino-terminal end of the helix. Note also that the amino-terminal end has unsatisfied hydrogen-bond donors (N-H groups) whereas the carboxy-terminal end has unsatisfied hydrogen-bond acceptors (C=O groups). Usually a polar side chain is found at the end of the helix, making hydrogen bonds to these donors and acceptors; such a residue is called a helix cap.







Alpha helices from top

- 3.6 residues per turn of the helix
- Residues 3-4 amino acids apart in the sequence will project from the same face
- In some proteins the surface face will be polar and the other face will be hydrophobic



Beta sheets

 Sometimes a hydrogen bond will be formed between residues that are faraway in sequence from each other to result in beta sheets. These can also be amphipathic.



Figure 1-17 The structure of the beta sheet The left figure shows a mixed beta sheet, that is one containing both parallel and antiparallel segments. Note that the hydrogen bonds are more linear in the antiparallel sheet. On the right are edge-on views of antiparallel (top) and parallel sheets (bottom). The corrugated appearance gives rise to the name "pleated sheet" for these elements of secondary structure. Consecutive side chains, indicated here as numbered geometric symbols, point from alternate faces of both types of sheet.

CIS 529: Bioinformatics

Beta Hairpin



Shape of protein

- The shape of the protein is determined by competition between self-interactions and interactions with water molecules
- Condensation of multiple secondary structure elements leads to a tertiary structure



Figure 1-23 Highly simplified schematic representation of the folding of a polypeptide chain in water In the unfolded chain, sidechain and main-chain groups interact primarily with water, even if they are hydrophobic and the interaction is unfavorable. Burying the hydrophobic groups in the interior of a compact structure enables them to interact with each other (blue line), which is favorable, and leaves polar side chains on the surface where they can interact with water (red lines). The polar backbone groups that are buried along with the hydrophobic side chains must make hydrogen bonds to each other (not shown), as bulk water is no longer available.

Protein folding

- Although the long chain of amino acids is held together as a chain by strong covalent bonds, this chain folds into a three dimensional structure which is stabilized by weak noncovalent interactions
- Gulliver in Lilliput analogy





Protein folding

CIS 529: Bioinformatics

Protein Folding Video



The energetics of protein folding

- A protein will spontaneously fold only if the folded state is energetically more feasible than the unfolded state
 - Entropy
 - Enthalpy

- Free energy
 - Non-bonding interactions
 - Entropy

Torsion angles

- N, C, O lie in the same plane
 - Rigid Structure
- N–Ca and Ca–C bonds are rotatable and free rotation is allowed
 - Flexible
 - Can rotate as long as no steric clashes arise (i.e., no side chains collide with each other)
- Phi torsion angle
 - The angle of the N–Ca bond to the adjacent peptide bond
- Psi torsion angle
 - the angle of the C−C_a bond to the adjacent peptide bond



Torsion angles

Excellent tool: <u>http://proteopedia.org/wiki/index.php/Psi_and_Phi_Angles</u>



Torsion angles

- A protein has alternating rigid and rotatable flexible covalent bonds
 - This combination greatly reduces the number of possible conformations a protein can have
 - Certain torsion angles will cause clashes
 - Ranges of torsion angles are also dependent upon the type of side chains

Ramachandran plot



Figure 1-11 Ramachandran plot Shown in red are those combinations of the backbone torsion angles phi and psi (see Figure 1-9) that are "allowed" because they do not result in steric interference. The pink regions are allowed if some relaxation of steric hindrance is permitted. Common protein secondary structure elements are marked at the positions of their average phi, psi values. The isolated pink alpha-helical region on the right is actually for a left-handed helix, which is only rarely observed in short segments in proteins. The zero values of phi and psi are defined as the *trans* configuration.

Fold-it

http://fold.it/portal/info/about

- Christopher B Eiben, Justin B Siegel, Jacob B Bale, Seth Cooper, Firas Khatib, Betty W Shen, Foldit Players, Barry L Stoddard, Zoran Popović, David Baker (2012). Increased Diels-Alderase activity through backbone remodeling guided by Foldit players. Nature Biotechnology, 2012.
 [PDF] [PROJECT]
- Firas Khatib, Seth Cooper, Michael D. Tyka, Kefan Xu, Ilya Makedon, David Baker, Foldit players (2011). Algorithm discovery by protein folding game players. Proceedings of the National Academy of Sciences of the United States of America vol. 108 no. 47 18949-18953, 2011. [PDF] [PROJECT]
- Firas Khatib, Frank DiMaio, Foldit Contenders Group, Foldit Void Crushers Group, Seth Cooper, Maciej Kazmierczyk, Miroslaw Gilski, Szymon Krzywda, Helena Zabranska, Iva Pichova, James Thompson, Mariusz Jaskolski, David Baker (2011). Crystal structure of a monomeric retroviral protease solved by protein folding game players. Nature Structural and Molecular Biology 18, 1175-1177, 2011. [PDF] [PROJECT]
- Seth Cooper, Firas Khatib, Ilya Makedon, Hao Lu, Janos Barbero, David Baker, James Fogarty, Zoran Popović, Foldit players (2011). Analysis of Social Gameplay Macros in the Foldit Cookbook. Foundations of Digital Games, 2011. [PDF] [PROJECT]
- Seth Cooper, Firas Khatib, Adrien Treuille, Janos Barbero, Jeehyung Lee, Michael Beenen, Andrew Leaver-Fay, David Baker, Zoran Popović, Foldit players (2010). Predicting protein structures with a multiplayer online game. Nature 446 p. 756-760, 05 August 2010. [PDF] [PROJECT]
- Seth Cooper, Adrien Treuille, Janos Barbero, Andrew Leaver-Fay, Kathleen Tuite, Firas Khatib, Alex Cho Snyder, Michael Beenen, David Salesin, David Baker, Zoran Popović, >57,000 Foldit players (2010). The challenge of designing scientific discovery games. International Conference on the Foundations of Digital Games, 2010. [PDF] [PROJECT]
- Other

Miroslaw Gilski, Maciej Kazmierczyk, Szymon Krzywda, Helena Zábranská, Seth Cooper, Zoran Popović, Firas Khatib, Frank DiMaio, James Thompson, David Baker, Iva Pichová, Mariusz Jaskolskia (2011). High-resolution structure of a retroviral protease folded as a monomer. Acta Crystallographica D67, 907-914, 2011. [PDF] [PROJECT]

How big is the protein conformation space?

- Assume that:
 - each of the torsion angles can take 2 or 3 values
 - In practice these angles are continuous
 - Side chains are fixed
 - We have a protein of length L
- Thus the number of possible "conformations" is:
 - 2L or 3L
 - If L = 100, this means 10^{30} or 10^{47} conformations

Levinthal's paradox

- The number of possible folding conformations of a protein is very large
- But the protein folds in microseconds or less

• How?

Formation of local interactions

Supersecondary structures

- Motif
 - A supersecondary structure that appears in a variety of molecules
 - Can appear in proteins with dissimilar functions
 - May not be associated with a sequence motifs
- Examples
 - Helix-Turn-Helix
 - Zinc Finger
 - $-\beta \alpha \beta$
 - β-hairpin
- Databases
 - PROSITE



Super-secondary structures

• Fold

- A protein fold is defined by the way the secondary structure elements of the structure are arranged relative to each other in space.
- Proteins can fold in only so many ways
- Examples
 - TIM Barrel Fold (PDB: 8TIM)
 - Beta-Barrel
 - Beta Propellor (PDB: 1ERJ)



CIS 529: Bioinformatics

Growth in PDB



Growth of unique folds in PDB 1,250 Growth Of Unique Folds Per Year As Defined By SCOP (v1.75) number of folds can be viewed by hovering mouse over the bar 1,000 Number Year

📕 Total 📕 Yearly

Supersecondary structures

- Chain
 - A protein can have multiple independent chains of amino acids
- Domains
 - A protein domain is a conserved part of a given protein sequence and <u>(tertiary)</u> <u>structure</u> that can <u>evolve</u>, function, and exist independently of the rest of the protein chain.
 - Two proteins having the same domain are said to belong to the same family
- SCOP database

Pyruvate kinase, a protein with three domains (<u>PDB 1PKN</u>

Computational Problems in Protein Biology

- <u>Visualization</u>
- Databases
- Alignments
 - Sequence based: Local, global, homology, phylogenetics, co-evolution, motif-finding
 - <u>Structure based alignments</u>
- <u>Structure prediction</u>
 - Secondary and tertiary, contact maps, fold, domain localization, family identification
 - Prediction of structures of larger protein assemblies
 - Molecular dynamics, normal mode analysis
 - Protein folding, folding-pathways
- Properties: <u>Surface area</u>, trans-membrane portions, signal peptides, functional sites, post-translational modifications
- Cellular Function prediction, phenotype prediction
- Cellular localization
- Biological process prediction
- <u>Mutation analysis</u>
- Protein interactions: Binding, binding sites, protein, nucleic acid, water, ions, ligands
- Motif Finding
- Drug discovery and design
- Protein design and engineering

PyMOL Visualization

- Load 1CLL
- Visualize 1CLL
 - Show the structure as a cartoon
 - Hide lines
 - Get rid of het-atoms
 - Show the Calcium ions as spheres
 - Show lines
 - Find polar contacts of the calcium ions
 - Select all atoms in an alpha helix and show the hydrogen bonds between them
 - Do the same for a beta sheet

46

Visualization types

- Types:
 - Lines
 - Sticks
 - Cartoon
 - Space-fill
 - Surface
- How are secodnary structures assigned?
 - DSSP
 - STRIDE (http://en.wikipedia.org/wiki/STRIDE_%28protein%29)

PyMOL Visualization

- Load 1CFD & 1DMO
 - Is it showing hydrogens?
- Align the two structures

http://pymol.sourceforge.net/newman/user/toc.html

Protein databases

- Genbank: All DNA sequences
- Uniprot
- SCOP: Structural classification of proteins
 - manual classification of protein structural domains based on similarities of their structures and amino acid sequences.
- Pfam
 - Pfam is a database of protein families that includes their annotations and multiple sequence alignments generated using hidden Markov models.

http://en.wikipedia.org/wiki/UniProt

http://en.wikipedia.org/wiki/Structu ral Classification of Proteins datab ase

UniProt

http://en.wikipedia.org/wiki/Pfam

Uniprot Statistics



Usage of amino acids



gray = aliphatic, red = acidic, green = small hydroxy, blue = basic, black = aromatic, white = amide, yellow = sulfur

CIS 529: Bioinformatics

Uniprot structure

- "The primary mission of Universal Protein Resource (UniProt) is to support biological research by maintaining a stable, comprehensive, fully classified, richly and accurately annotated protein sequence knowledgebase, with extensive crossreferences and querying interfaces freely accessible to the scientific community"
- Each protein has a uniprot accession number

<u>"Ongoing and future developments at the Universal Protein Resource," *Nucleic Acids* <u>Res.</u>, vol. 39, no. Database issue, pp. D214–D219, Jan. 2011.</u>

Protein databases: Sequence

- Sequence
 - UniProt
 - UniProt Knowledgebase (UniProtKB): Expertly curated annotations
 - The UniProt Archive (UniParc): Reflects history of all proteins
 - UniProt Reference Clusters (UniRef): Closely related proteins are merged
 - UniProt Metagenomic and Environmental Sequences Database (UniMES)



"Ongoing and future developments at the Universal Protein Resource," Nucleic Acids Res., vol. 39, no. Database issue, pp. D214–D219, Jan. 2011.

CIS 529: Bioinformatics

Uniprot structure

- UniProt Knowledgebase (UniProtKB)
 - Expertly curated database,
 - Central access point for integrated protein information
 - Cross-references to multiple sources.
- The UniProt Archive (UniParc)
 - is a comprehensive sequence repository
 - Reflects the history of all protein sequences
- UniProt Reference Clusters (UniRef)
 - Merge closely related sequences based on sequence identity to facilitate sequence similarity searches
 - Using CD-HIT
- UniProt Metagenomic and Environmental Sequences Database (UniMES)
 - Caters for the expanding area of metagenomic data.

Uniprot structure

- UniProtKB/Swiss-Prot
 - 542,782 sequence entries in 2014_03 release
 - Reviewed, manually annotated entries
 - Annotations include
 - Protein and gene names
 - Function
 - <u>Enzyme</u>-specific information such as <u>catalytic activity</u>, <u>cofactors</u> and <u>catalytic residues</u>
 - <u>Subcellular location</u>
 - <u>Protein-protein interactions</u>
 - Pattern of expression
 - Locations and roles of significant domains and sites
 - <u>lon</u>-, <u>substrate</u>- and cofactor-binding sites
 - Protein variant forms produced by natural genetic variation, <u>RNA editing</u>, alternative splicing, <u>proteolytic</u> processing, and post-translational modification
- UniProtKB/TrEMBL

Automatic annotation

Example

- <u>http://www.uniprot.org/uniprot/P62158</u>
- Breuza, Lionel, Sylvain Poux, Anne Estreicher, Maria Livia Famiglietti, Michele Magrane, Michael Tognolli, Alan Bridge, Delphine Baratin, Nicole Redaschi, and The UniProt Consortium, 2016. "The UniProtKB Guide to the Human Proteome." Database 2016 (January): bav120. doi:10.1093/database/bav120.

Protein Databank

- <u>http://proteopedia.org/wiki/index.php/Believe_It_or_Not</u>
- <u>http://en.wikipedia.org/wiki/Protein_Data_Bank</u>

Experimental Method	Proteins	Nucleic Acids	Protein/Nucleic Acid complexes	Other	Total
X-ray diffraction	85406	1558	4547	5	91516
NMR	9281	1092	218	7	10598
Electron microscopy	579	67	190	0	836
Hybrid	63	3	2	1	69
Other	157	4	6	13	180
Total:	95486	2724	4963	26	103199

Tools for PDB

Database name	Database description
DSSP	Secondary structure of proteins
	http://swift.cmbi.ru.nl/gv/dssp/
HSSP	Multiple sequence alignments of UniProtKB against PDB files
	http://swift.cmbi.ru.nl/gv/hssp/
PDBREPORT	Reports about errors and anomalies in macromolecules
	http://swift.cmbi.ru.nl/gv/pdbreport/
PDB_REDO	Re-refined PDB files solved by X-ray crystallography
_	http://www.cmbi.ru.nl/pdb_redo/
PDBFINDER	Searchable summaries of PDB file information
	ftp://ftp.cmbi.ru.nl/pub/molbio/data/pdbfinder/
PDBFINDER2	As PDBFINDER, but with much extra information added
	ftp://ftp.cmbi.ru.nl/pub/molbio/data/pdbfinder2/
PDB_SELECT	Quality-sorted culled lists of protein chains in the PDB
_	http://swift.cmbi.ru.nl/gv/select/
WHY_NOT	Explanation why entries in other databases cannot exist
_	http://www.cmbi.ru.nl/WHY_NOT/

R. P. Joosten, T. A. H. te Beek, E. Krieger, M. L. Hekkelman, R. W. W. Hooft, R. Schneider, C. Sander, and G. Vriend, "A series of PDB related databases for everyday needs," *Nucleic Acids Res.*, vol. 39, no. suppl 1, pp. D411–D419, Jan. 2011.

Other tools

- OCA
- FirstGlance
- PDBSum
- ccPDB

Gene Ontology (GO)

- Gene ontology, or GO, is a major bioinformatics initiative to unify the representation of gene and gene product attributes across all species. More specifically, the project aims to:
 - Maintain and develop its controlled vocabulary of gene and gene product attributes;
 - Annotate genes and gene products, and assimilate and disseminate annotation data;
 - Provide tools for easy access to all aspects of the data provided by the project, and to enable functional interpretation of experimental data using the GO, for example via enrichment analysis.

http://geneontology.org/page/documentation

GO





http://geneontology.org/page/current-go-statistics

CIS 529: Bioinformatics

Protein databases: Function annotations

- Gene Ontology Database
 Hierarchically organized
 - Ontologies
 - Subcellular localization
 - Molecular Function
 - Biological Process
 - Example: "Calmdoulin Binding"
 - http://amigo.geneontology.org/a migo/term/GO:0005516







goslim_yeast

QuickGO - http://www.ebi.ac.uk/QuickGO



http://cbl-gorilla.cs.technion.ac.il/exa@[18:529: Bioinformatics

End of Lecture