

Structure-Based Prediction of Protein Function

Thomas Funkhouser
Princeton University
CS597A, Fall 2005

Outline

Protein structure databases

- Repositories
- Classifications

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Outline

Protein structure databases

- ~~Ø~~Repositories
- Classifications

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

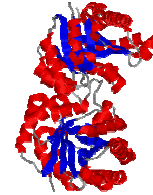
Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Protein Structure Databases

Repositories:

- Primary ← amino acid sequence
- Secondary ← local fold pattern of small subsequence
- Tertiary ← fold of entire protein chain
- Quaternary ← complex of multiple chains



I tim
[Jena]

Protein Structure Databases

Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Protein Structure Databases

Repositories:

- ~~Ø~~Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Chain 1 GSA:_
Compound Glutathione Synthetase
Type Protein
Molecular Weight 35547
Number of Residues 316

```
1 MRLGLVWDV IANINIKKDS SPWALLEAQR RQVELRYMNM GELYLINGDA  
51 RAHRTLANK QVYERFSPV GRDLPLADL DVILMRSDDP FQTERIVAYY  
101 ILSBARERT LVNHPQSLR DNEKLEPTAM FRLTPTSTLV TNSAQLQAF  
151 WKKSDILK PLDQMGASL PRVKEEDNL QVIARTLTH GTRYCHAGNY  
201 LPAEKDQKR VLVVDREVPV YCLARIPQGD EYGNLAAGD RGRPPPLTES  
251 DWKIDARQID VLKELKILPV GDLIGRLRL EHNVTSPCCI REIARFPFIS  
301 ITYDMDAIE ARLQOO
```

<http://www.uniprot.org/> [Apweiler04]

Protein Structure Databases



Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

H = helix
 E = residue in isolated beta bridge
 E = extended beta strand
 G = 310 helix
 T = hydrogen bonded turn
 S = bend

```
Chain 1 GSA:
Compound Glutathione Synthetase
Type Protein
Molecular Weight 35547
Number of Residues 316
Number of Alpha 9 Content of Alpha 27.22
Number of Beta 19 Content of Beta 28.16
1  MKLGLVMDP IANINIKKDS SPFMLEAQR RYVELYVEM GELYLNGEA
   EEEE S  GOOTTTTTT HHRRHHHHH HT EEEE G  QDESEETTEE
51  RAHRTLNVK QVYKRFSPV GRDLKPLADL DVILMRKDDP FOTFIVATY
   EEEEEEE S  SS  EEE  EEEKDDGDS  EEEE  HHHHHHHH
101  ILLEAREKUT LVNEDQGLR DNRKLEPTAM FSULTPFTLV TRNQAQLAF
   HHHHHHHTT  EEE  HHHH  HTTTTGGGG  GTTB  EEE  SS  HHHHHH
151  WKKRSDILK PLDMGGASL FVWRKDDPHL GVIARTLTER QTVYMAQRY
   HHHHEEEE  SS  TTTT  EEE  TTTTTH  HHHHHHHTT  TTS  EEEE
201  LALRDKRER NLVVDGDFV VGLARLDGDS EYHRLMAGD RSDVPLRES
   GGGG  EEE  ESEETTEE S  EEEEE  SS  S  GAT  EEEEE  HH
251  DKARLQSDP TLKRGKLPV GLNIGDRLT EINHTEPTCI REIARFPVE
   HHHHHHHTT  HHHHTT  EE  ESEETTEE  EEE  SS  H  HHHHHSS
301  ITDGMGAIE ARLQQQ
   HHHHHHHHH  HHT
```

[Kabsch83]

Protein Structure Databases



Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS



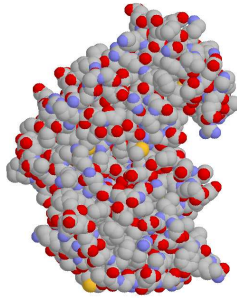
[Kabsch83]

Protein Structure Databases



Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS



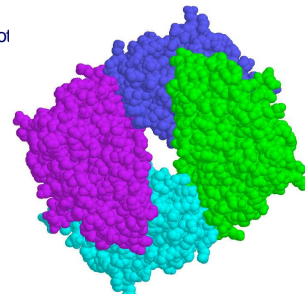
<http://www.rcsb.org/pdb/> [Berman00]

Protein Structure Databases



Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS



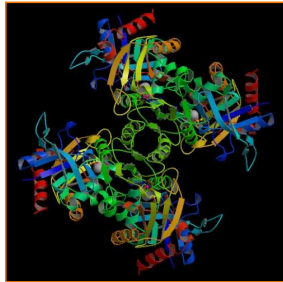
<http://pqs.ebi.ac.uk/> [Hendrick98]

Protein Structure Databases



Repositories:

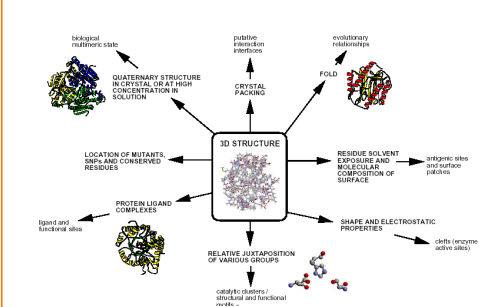
- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS



<http://pqs.ebi.ac.uk/> [Hendrick98]

Protein Structure Databases

Slide courtesy of Philip Bourne



Protein Structure Databases

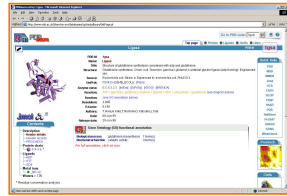


Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Summaries:

- ~~PDBsum~~
- Jena
- MSD



<http://www.ebi.ac.uk/thornton-srv/databases/pdbsum/> [Laskowski05]

Protein Structure Databases

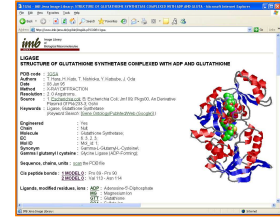


Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Summaries:

- PDBsum
- ~~Jena~~
- MSD



<http://www.imb-jena.de/IMAGE.html>

Protein Structure Databases

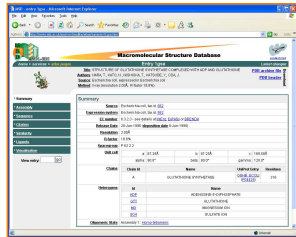


Repositories:

- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Summaries:

- PDBsum
- Jena
- ~~MSD~~



<http://www.ebi.ac.uk/msd/> [Velankar05]

Protein Structure Databases



Repositories:

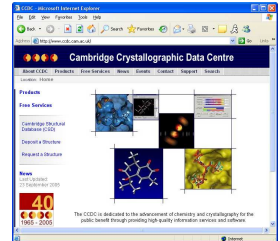
- Primary ← UniProt
- Secondary ← DSSP
- Tertiary ← PDB
- Quaternary ← PQS

Summaries:

- PDBsum
- Jena
- MSD

Small molecules:

- ~~CCDC~~



<http://www.ccdc.cam.ac.uk/>

Outline



Protein structure databases

- Repositories
- ~~Classifications~~

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Protein Structure Classifications

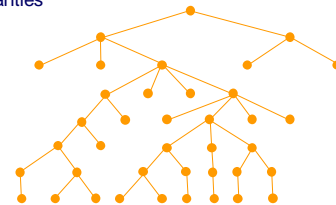


Clustering

- Fold similarities
- Evolutionary relationships
- Sequence similarities

Examples:

- CATH
- SCOP



Protein Structure Classifications



CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100



<http://cathwww.biochem.ucl.ac.uk/> [Orengo97]

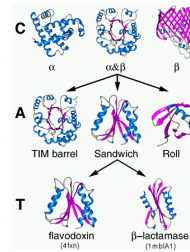
Protein Structure Classifications



CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100

} Structural Layout



<http://cathwww.biochem.ucl.ac.uk/> [Orengo97]

Protein Structure Classifications

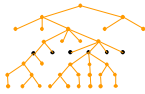


CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100

} Evolution

CATH Domain 1gsa01	
Classification	
Class	3
Alpha Beta	
Architecture	3.40
5-Layered Sandwich	
Topology	3.40.50
Rossmann fold	
Homologous Superfamily	3.40.50.20
LIGASE	
Sequence Family (S35)	3.40.50.20.7
LIGASE	
Non-identical (S95)	3.40.50.20.7.1
LIGASE	
Identical (S100)	3.40.50.20.7.1.1
LIGASE	



<http://cathwww.biochem.ucl.ac.uk/> [Orengo97]

Protein Structure Classifications



CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100

} Sequence Identity

CATH Domain 1gsa01	
Classification	
Class	3
Alpha Beta	
Architecture	3.40
5-Layered Sandwich	
Topology	3.40.50
Rossmann fold	
Homologous Superfamily	3.40.50.20
LIGASE	
Sequence Family (S35)	3.40.50.20.7
LIGASE	
Non-identical (S95)	3.40.50.20.7.1
LIGASE	
Identical (S100)	3.40.50.20.7.1.1
LIGASE	



<http://cathwww.biochem.ucl.ac.uk/> [Orengo97]

Protein Structure Classifications



SCOP hierarchy:

- Class
- Fold
- Superfamily
- Family
- Protein Domain
- Species
- PDB

SCOP: 1gsa

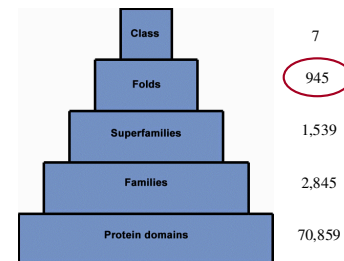
1. Root: [scop](#)
2. Class: [Alpha and beta proteins \(ab\)](#) [51349]
3. Fold: [PreATP-grasp domain](#) [52439]
4. Superfamily: [PreATP-grasp domain](#) [52440]
5. Family: [Prokaryotic glutathione synthetase, N-terminal domain](#) [52457]
6. Protein: [Prokaryotic glutathione synthetase, N-terminal domain](#) [52458]
7. Species: [Escherichia coli](#) [52459]

<http://scop.mrc-lmb.cam.ac.uk/scop/> [Murzin95]

Protein Structure Classifications



SCOP hierarchy:



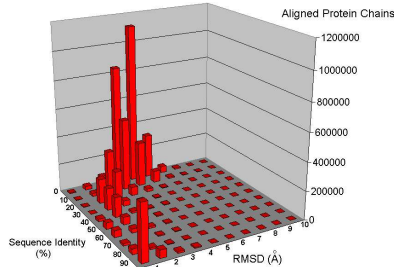
SCOP: Structural Classification of Proteins (1.69 release)

Protein Structure Classifications

Slide courtesy of Philip Bourne



Protein folds are highly redundant



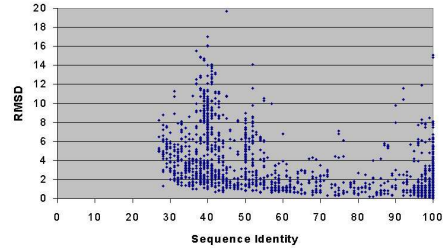
Structure Alignments using CE with $z > 4.0$

Sequence → Structure → Function

Slide courtesy of Philip Bourne



Sequence determine structure, but ...



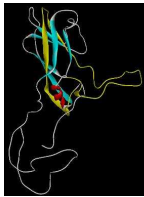
Structure Comparison of 30% of PDBSelect Set

Sequence → Structure → Function

Slide courtesy of Philip Bourne



Similar sequence, different structure & function



IPIV:1 (Viral Capsid Protein)



IHMP:A (Glycosyltransferase)

80 Residue Stretch (Yellow) with Over 40% Sequence Identity

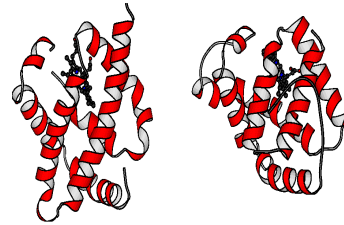


Sequence → Structure → Function

Slide courtesy of Philip Bourne



Different sequence, similar structure & function



The globin fold is resilient to amino acid changes. *V. stercoraria* (bacterial) hemoglobin (left) and *P. marinus* (eukaryotic) hemoglobin (right) share just 8% sequence identity, but their overall fold and function is identical.

Sequence → Structure → Function

Slide courtesy of Philip Bourne



Some folds have many functions



NAD binding domain

This is a double beta-alpha-beta-alpha-beta motif, and is a common structural motif of enzymes binding NAD, NADP and other related cofactors, for example, NAD is found in dehydrogenases as the hydrogen acceptor. The domain is found as a common core unit in many structures, with other structural units at the periphery.

Number of EC numbers associated with this fold (to the third level): 5



F-loop NTP hydrolase fold

This fold consists of alpha-beta-alpha, parallel or mixed beta sheets of variable size. The fold binds the phosphate of ATP or GTP and is found in ATP and GTP binding proteins such as adenylate kinase. The F-loop is a phosphate binding loop which binds the phosphate groups of ATP and GTP and is a glycine-rich sequence with the consensus sequence (A/G)xxx(G/K)S. The F-loop residues are shown in detail (left) in guanylate kinase.

Number of EC numbers associated with this fold (to the third level): 5



Ferredoxin-like fold

This fold consists of an alpha-beta sandwich with an antiparallel beta sheet. The ferredoxin-like fold is associated with predominantly with non-enzymatic ferredoxins, like the example shown (Ferredoxin II from *D. gigas*, left). Ferredoxins are iron-sulfur clusters involved in electron transport, and often form part of multi-subunit assemblies. An example of an enzyme with this fold is muconolactone isomerase (EC 5.3.3.4).

Number of EC numbers associated with this fold (to the third level): 5

Outline



Protein structure databases

- Repositories
- Classifications

Protein function databases ←

- Gene Ontology (GO)
- Enzyme Commission (EC)

Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Protein Function



Possible levels of functional characterization

Molecular Function = elemental activity/task

- the tasks performed by individual gene products; examples are *carbohydrate binding* and *ATPase activity*

Biological Process = biological goal or objective

- broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

Cellular Component = location or complex

- subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *RNA polymerase II holoenzyme*

[Karen Christie et al., GO Annotation Camp 2005]

Protein Function



Possible levels of functional characterization

Molecular Function = elemental activity/task

- the tasks performed by individual gene products; examples are *carbohydrate binding* and *ATPase activity*

Biological Process = biological goal or objective

- broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

Cellular Component = location or complex

- subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *RNA polymerase II holoenzyme*

Structure is mostly useful for predicting molecular function

[Karen Christie et al., GO Annotation Camp 2005]

Protein Function Databases



Gene Ontology

- Directed acyclic graph of gene product attributes – provides functional terms at multiple levels for genes

Enzyme Commission (EC)

- Hierarchical classification of enzymes based on the chemical reactions they catalyze

Protein Function Databases



Gene Ontology

- Directed acyclic graph of gene product attributes – provides functional terms at multiple levels for genes

Enzyme Commission (EC)

- Hierarchical classification of enzymes based on the chemical reactions they catalyze

Protein Function Databases



Gene Ontology (top-level terms)

```
GO:0008159 | biological_process (137592)
- GO:0007610 | behavior (4716)
- GO:0000004 | biological_process_unknown (31206)
- GO:0000987 | cellular_process (87883)
- GO:0007255 | development (18172)
- GO:0040067 | growth (3591)
- GO:0044445 | interaction_betweenorganisms (871)
- GO:0007882 | physiological_process (92179)
- GO:0044573 | pigmentation (174)
- GO:0009789 | regulation_of_biological_process (18806)
- GO:0009091 | reproduction (512)
- GO:0046822 | viral_life_cycle (305)
GO:0008579 | cellular_component (124814)
- GO:0005623 | cell (94481)
- GO:0008712 | cellular_component_unknown (25532)
- GO:0010162 | extracellular_matrix (868)
- GO:0005756 | extracellular_region (9981)
- GO:0043226 | organelle (65695)
- GO:0043214 | protein_complex (13646)
- GO:0009632 | vesicle (11)
GO:0003674 | molecular_function (138336)
- GO:0004209 | oxidoreductase_activity (142)
- GO:0004814 | binding (14687)
- GO:0003824 | catalytic_activity (44201)
- GO:0003018 | chaperone_regulator_activity (50)
- GO:0002014 | cysteine_regulator_activity (2515)
- GO:0005514 | molecular_function_unknown (14234)
- GO:0003774 | motor_activity (594)
- GO:0045735 | nutrient_receptor_activity (43)
- GO:0001186 | protein_tag (18)
- GO:0004476 | signal_transducer_activity (10883)
- GO:0005208 | structural_molecular_activity (3968)
- GO:0008259 | transcription_regulator_activity (9526)
- GO:0045182 | translation_regulator_activity (136)
- GO:0006215 | transporter_activity (11567)
- GO:0000533 | triplet_codon-amino_acid_adaptor_activity (1217)
```

<http://www.godatabase.org/>

Protein Function Databases



Gene Ontology (molecular function terms)

Term name	Total Gene Products	Percent
binding	46487	33.5
catalytic activity	44201	31.9
molecular function unknown	34234	24.7
transporter activity	11587	8.37
signal transducer activity	10883	7.86
transcription regulator activity	9526	6.88
structural molecule activity	3986	2.88
enzyme regulator activity	2515	1.81
triplet codon-amino acid adaptor activity	1217	0.87
translation regulator activity	836	0.60
motor activity	594	0.42
antioxidant activity	542	0.39
chaperone regulator activity	50	0.03
nutrient reservoir activity	43	0.03
protein tag	18	0.01
molecular_function	1	0.00

<http://www.godatabase.org/>

Protein Function Databases



Gene Ontology (molecular function terms)

- Most proteins have several GO terms associated with them

Gene Ontology (GO) functional annotation for PDB entry 1gsa

GO Term	Chain(s)
0006750 glutathione biosynthesis	
0003624 catalytic activity	
0004363 glutathione synthase activity	
0005524 ATP binding	

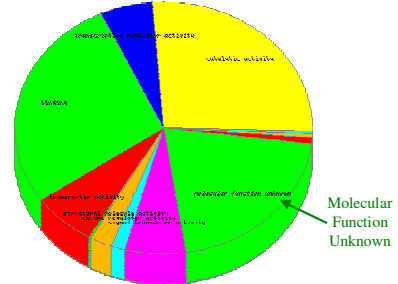
PDBsum

GO Terms

Protein Function Databases



Gene Ontology (molecular function terms)



Protein Function Databases



Gene Ontology

- Directed acyclic graph of gene product attributes – provides functional terms at multiple levels for genes

Enzyme Commission (EC)

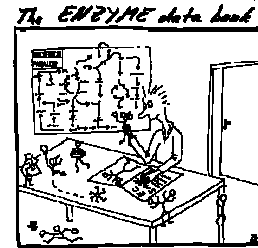
- Hierarchical classification of enzymes based on the chemical reactions they catalyze

Protein Function Databases



Enzyme Commission (EC) numbers

- Hierarchical classification of enzymes based on the chemical reactions they catalyze



Protein Function Databases



Enzyme Commission (EC) numbers

EC 1 Oxidoreductases

EC 1.1 Acting on the CH-OH group of donors

EC 1.1.1 With NAD⁺ or NADP⁺ as acceptor

EC 1.1.2 With a cytochrome as acceptor

EC 1.1.3 With oxygen as acceptor

EC 1.1.4 With a disulfide as acceptor

EC 1.1.5 With a quinone or similar compound as acceptor

EC 1.1.99 With other acceptors

EC 1.2 Acting on the aldehyde or oxo group of donors

EC 1.2.1 With NAD⁺ or NADP⁺ as acceptor

EC 1.2.2 With a cytochrome as acceptor

EC 1.2.3 With oxygen as acceptor

EC 1.2.4 With a disulfide as acceptor

EC 1.2.7 With an iron-sulfur protein acceptor

EC 1.2.99 With other acceptors

EC 1.3 Acting on the CH-OH group of donors

EC 1.3.1 With NAD⁺ or NADP⁺ as acceptor

EC 1.3.2 With a cytochrome as acceptor

EC 1.3.3 With oxygen as acceptor

EC 1.3.5 With a quinone or related compound as acceptor

EC 1.3.9 With an iron-sulfur protein as acceptor

EC 1.3.99 With other acceptors

EC 1.4 Acting on the CH-NH2 group of donors

EC 1.4.1 With NAD⁺ or NADP⁺ as acceptor

EC 1.4.2 With a cytochrome as acceptor

EC 1.4.3 With oxygen as acceptor

EC 1.4.4 With a disulfide as acceptor

EC 1.4.7 With an iron-sulfur protein as acceptor

EC 1.4.99 With other acceptors

etc.

<http://www.expsy.org/enzyme/>

Protein Function Databases



Enzyme Commission (EC) numbers

- Specify exact reaction catalyzed by enzyme

PDBsum

EC

PDBsum

Outline

Protein structure databases

- Repositories
- Classifications

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

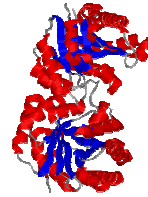
Sequence → Structure → Function ←

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Sequence → Structure → Function

Goal:

- Given a protein sequence/structure, predict its function



Protein Structure



? ? ?
? ? ?
? ? ?
? ? ?

Protein Function

Sequence → Structure → Function

General strategy:

1. Given a protein with unknown function
2. Match it to proteins/templates with known functions
3. Transfer function from statistically significant matches



Protein Structure

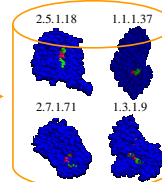
Sequence → Structure → Function

General strategy:

1. Given a protein with unknown function
2. Match it to proteins/templates with known functions
3. Transfer function from statistically significant matches



Protein Structure



Database

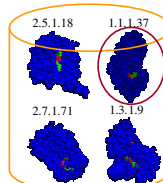
Sequence → Structure → Function

General strategy:

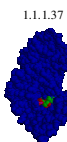
1. Given a protein with unknown function
2. Match it to proteins/templates with known functions
3. Transfer function from statistically significant matches



Protein Structure



Database



Statistically Significant Match

Sequence → Structure → Function

Evolution:

- Divergent evolution
 - § Homology: proteins share a common ancestor
 - Orthology: separated by a speciation event
 - Paralogy: separated by a gene duplication event
- Convergent evolution
 - § Analogy: similar structure evolves independently in two species due to similar selective pressures



a. Subtilisin EC 3.4.21.62



b. Chymotrypsin EC 3.4.21.1

Sequence Database Search



Compute sequence alignment for query with all in database, and report statistically significant matches

Smith-Waterman score	k-clip	a.a.	z	E	PDB code	Protein name
identity	overlap	length	score	value		
			318			
1914	90.1%	313	315	2355.5	3.4e-11	1lp5A Theoretical model of glutathione synthetase
1906	94.3%	314	299	1729.1	7.1e-9	1glv Glutathione synthase loopless mutant with residues
1887	95.6%	314	296	1737.5	3.2e-8	1gll Structure of escherichia coli glutathione synthetase
						1gsh Structure of escherichia coli glutathione synthetase
						1glt Glutathione synthase

FASTA

[Pearson88]

Outline



Protein structure databases

- Repositories
- Classifications

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

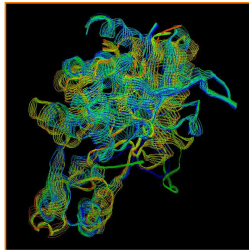
Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Structure Alignment



Align protein structures
Report score and significance



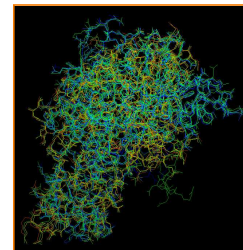
SSM

[Krissne104]

Structure Alignment



Align protein structures
Report score and significance



SSM

[Krissne104]

Structure Alignment



Align protein structures
Report score and significance

H#	Select	Z-score	RMSD	SE	SI	entry	Name
1	<input type="checkbox"/>	22.2	28	0.00	100.0%	1glt	Structure of glutathione synthetase complexed with adp and glutathione
2	<input type="checkbox"/>	19.9	27	0.09	99.2%	1gsh	Structure of escherichia coli glutathione synthetase at pH 7.0
3	<input type="checkbox"/>	19.4	27	0.42	99.2%	2glv	Structure of escherichia coli glutathione synthetase at pH 5.0
4	<input type="checkbox"/>	19.3	27	0.66	99.0%	1glt	Glutathione synthase loopless mutant with residues 226-244 replaced
5	<input type="checkbox"/>	19.0	26	0.86	99.0%	1lp5A	Theoretical model of glutathione synthetase
6	<input type="checkbox"/>	7.8	18	1.49	32.0%	1glt	Crystal structure of a lysine biosynthesis enzyme, lysx, from thermotoga
7	<input type="checkbox"/>	7.6	17	2.43	26.4%	1glt	Crystal structure of a lysine biosynthesis enzyme, lysx, from thermotoga
8	<input type="checkbox"/>	6.7	18	2.07	18.7%	1glt	Crystal structure of a lysine biosynthesis enzyme, lysx, from thermotoga
9	<input type="checkbox"/>	6.4	18	2.71	21.0%	1glt	Crystal structure analysis of the complex of the C domain of lysine biosynthesis
10	<input type="checkbox"/>	6.4	18	2.71	18.2%	1glt	Crystal structure of the C domain of lysine biosynthesis from bacillus

SSM

[Krissne104]

Outline



Protein structure databases

- Repositories
- Classifications

Protein function databases

- Gene Ontology (GO)
- Enzyme Commission (EC)

Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- Structure motifs

Sequence Motifs

Recognize local patterns (motifs) in protein sequences associated with specific functions

Example: Prokaryotic glutathione synthetase ATP-binding domain

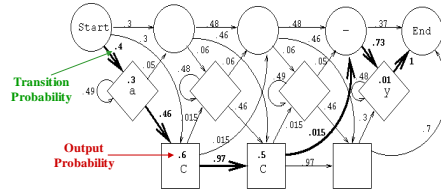
ProDom

<http://prodes.toulouse.inra.fr/prodom/>

Sequence Motifs

Recognize local patterns (motifs) in protein sequences associated with specific functions

- Represent as Hidden Markov Model (HMM)



Sequence Motifs

Match query sequence against all motifs

Examples:

- InterPro
- PROSITE
- PRINTS
- Pfam-A
- TIGRFAM
- PROFILES
- PRODOM

ProFunc: InterProScan results for 1gsa

DB Scan	code	Residue range	Residues	Motif name
1	InterPro	124-180	57	glu-amy glutathione synthase
2	SMART	124-180	57	Prokaryotic glutathione synthetase, domain
3	SMART	124-300	176	Prokaryotic glutathione synthetase, ATP-grp
4	SMART	124-259	136	ATP-grp
5	SMART	124-259	136	ATP-grp
6	SMART	124-259	136	ATP-grp
7	SMART	124-259	136	ATP-grp
8	SMART	124-259	136	ATP-grp
9	SMART	124-259	136	ATP-grp
10	SMART	124-259	136	ATP-grp

InterPro [Zdobnov01]

Outline

Protein structure databases

- Repositories
- Classifications

Protein function databases

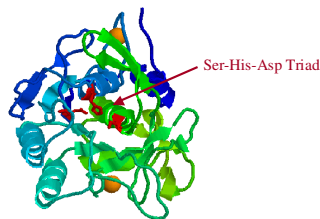
- Gene Ontology (GO)
- Enzyme Commission (EC)

Sequence → Structure → Function

- Sequence alignment
- Structure alignment
- Sequence motifs
- ◊ Structure motifs

Structure Motifs

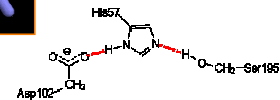
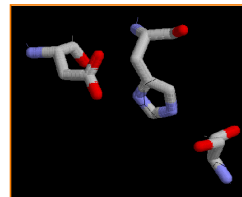
Recognize local patterns (motifs) in protein structures associated with specific functions



<http://chemistry.umeche.maine.edu/>

Structure Motifs

Example: serine proteases



<http://chemistry.umeche.maine.edu/>

Going Forward



Study algorithms

- Next meeting: structural alignment

Think about new methods

- e.g., new structural motifs

Think about how methods interact

Acknowledgements



Slides:

- James Watson
- Philip Bourne