



Protein Structure Analysis

Thomas Funkhouser
Princeton University
CS597A, Fall 2007

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- Geometry
- Relationship to sequence
- Classifications

Outline



Protein structure databases

- ØPrimary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- Geometry
- Relationship to sequence
- Classifications

Primary Structure Statistics



Main databases:

- UniProtKB/Swiss-Prot ← curated
- UniProtKB/TrEMBL



<http://www.pir.uniprot.org/>

Primary Structure Statistics



Information provided:

- Sequence of amino acid types

Chain IGSA:
Compound Glutathione Synthetase
Type Protein
Molecular Weight 35547
Number of Residues 316

```

1 MIKLGIVNDP IANIKKEDS SFAMLEAQR ROYELHYVMN GELYLNGRA
51 RANHTLAWV QRYKRFSPV GRQLSLADL DVLAKHEDSP FDFEYIATY
101 ILSRAKERT LIVNDQLR DQKELPTAM PSLTPTETLV TRNSAGLAF
151 HESRSLIE ELQNGGLL FWRKSGDM QVATLTER QRYVQAGY
201 LKALDQGR VLVDSEHY VCLARFDQD STKQMLADG RQEPRLTIS
251 DMELAQIQD TLKREHLFV GELDLEMLT RINVTPTCI REIRAEFVS
301 ITOMAGDAE ARLQQD
  
```

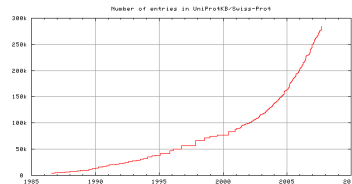
<http://www.uniprot.org/> [Apweiler04]

Primary Structure Statistics



Sequence counts:

- 283,454 sequence entries in UniProtKB/Swiss-Prot
- 4,754,787 sequence entries in UniProtKB/TrEMBL

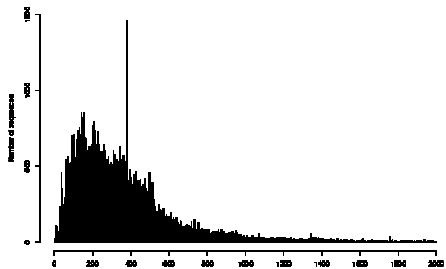


Release 54.2 of 11-Sep-07 of UniProtKB/Swiss-Prot
<http://expasy.org/uniprot/releasenotes/releas11.html>

Primary Structure Statistics



Sequence lengths:



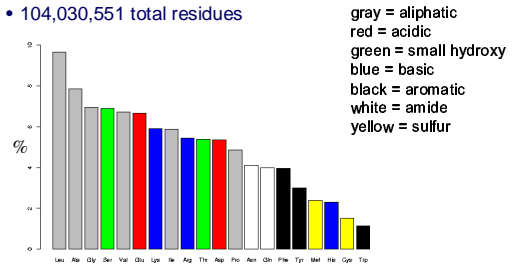
Release 54.2 of 11-Sep-07 of UniProtKB/Swiss-Prot
<http://expasy.org/spot/releasenotes/relestat.html>

Primary Structure Statistics



Amino acid counts:

- 104,030,551 total residues

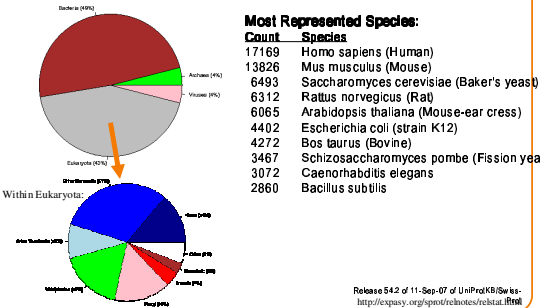


Release 54.2 of 11-Sep-07 of UniProtKB/Swiss-Prot
<http://expasy.org/spot/releasenotes/relestat.html>

Primary Structure Statistics



Taxonomic distribution:



Most Represented Species:

Count	Species
17169	Homo sapiens (Human)
13826	Mus musculus (Mouse)
6493	Saccharomyces cerevisiae (Baker's yeast)
6312	Rattus norvegicus (Rat)
6065	Arabidopsis thaliana (Mouse-ear cress)
4402	Escherichia coli (strain K12)
4272	Bos taurus (Bovine)
3467	Schizosaccharomyces pombe (Fission yeast)
3072	Caenorhabditis elegans
2860	Bacillus subtilis

Release 54.2 of 11-Sep-07 of UniProtKB/Swiss-Prot
<http://expasy.org/spot/releasenotes/relestat.html>

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

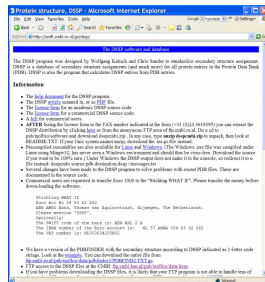
- Geometry
- Relationship to sequence
- Classifications

Secondary Structure Statistics



Main database:

- DSSP



<http://swift.cmbi.ru.nl/gv/dssp/>

Secondary Structure Statistics



Information provided:

- Predicted secondary structure element for every residue

Chain 1GSA:_
 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

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

```

1 MELIIVMD IAINIKKDS SPWLELQD RQVLELVNEM QGLIHLNDA
EEEE S QDITTTTT HHRRHHHS HT EEEK Q QDSEETTE
51 RAHNTLVAK QVYIEPSPV GQDLGLAD DVLMLKSDP PTFEIVATY
EEEEEEK E S EKE EEEKGGGS EEEE HHRRHHH
101 ILERAEKKT LVNYPQLR DNEKLEFAM FSLTPTSLV TNSQLQAF
HHRRHHHT EEE HHHH HTTTGGGG OTTTS EEE EE HHRRHH
151 WEKSDILK DLDQMGASL FVKEKEDNL QVIAITLTH GYVYAGNY
HHRRSEKK SS TTT EEE TTTTH HHRRHHHT TTS EEEE
201 LPAKIDDER VLVDGEPV YCLARIQGD EYKMLAAG RGRPPLTES
GGGG EEE EEEETTS S EEEEE SS S GDTT EEEEE HH
251 DWIARLQD PLKELQPLV GLIILQRLT RINPTDPLI REIASPFD
HHRRHHHT HHHTT EE EEEETTES EEE SS H HHRRHH
301 IYQMGDAE RLQQQ
HHRRHHHH HTT
    
```

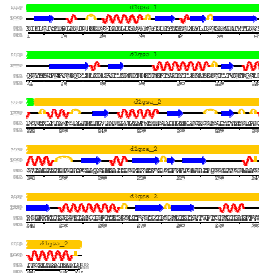
[Kabsch83]

Secondary Structure Statistics



Information provided:

- Predicted secondary structure element for every residue



[Kabsch83]

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- Geometry
- Relationship to sequence
- Classifications

Tertiary Structure Statistics



Main database:

- PDB



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

Tertiary Structure Statistics



Information provided:

- Atomic coordinates for every atom
- Remarks and info about experiment

```

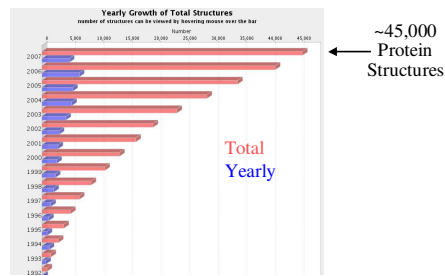
HEADER          1U88          1U88              08-JUN-95
COMPND          1 PROTEIN: POLYMERIZATION PROTEIN:
COMPND          2 PROIN: BILLY
COMPND          3 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          4 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          5 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          6 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          7 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          8 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND          9 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         10 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         11 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         12 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         13 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         14 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         15 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         16 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         17 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         18 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         19 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
COMPND         20 POLYMERIZATION PROTEIN: POLYMERIZATION PROTEIN
    
```

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

Tertiary Structure Statistics



Structure count:



http://www.rcsb.org/pdb/static.do?p=general_information/pdb_statistics/index.html

Tertiary Structure Statistics



Experimental method:

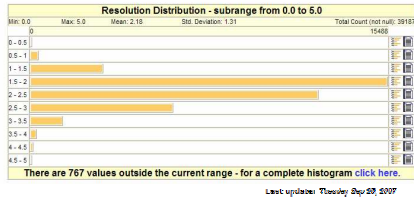
	Molecule Type				Total
	Proteins	Nucleic Acids	Protein/NA Complexes	Other	
X-ray	36466	989	1705	24	39184
NMR	5696	784	134	7	6621
Exp. Electron Microscopy	106	10	38	0	154
Other	82	4	4	2	92
Total	42350	1787	1881	33	46051

http://www.rcsb.org/pdb/static.do?p=general_information/pdb_statistics/index.html

Tertiary Structure Statistics



Resolution:



http://www.rcsb.org/pdb/static.do?p=general_information/pdb_statistics/index.html

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

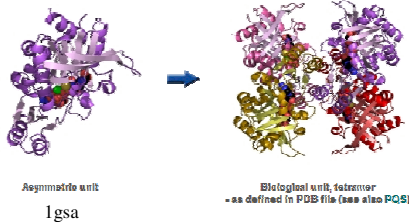
- Geometry
- Relationship to sequence
- Classifications

Quaternary Structure Statistics



Main databases:

- PQS
- PISA

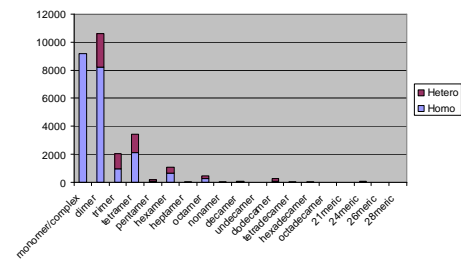


<http://pqs.ebi.ac.uk/> [Hendrick98]

Quaternary Structure Statistics



Counts of monomeric / oligomeric proteins in PQS:



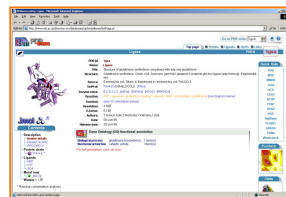
PQS

Protein Structure Databases



Useful resources:

- [PDBsum](#)
- Jena
- MSD



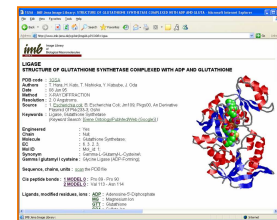
<http://www.ebi.ac.uk/thornton-srv/databases/pdbsum/> [Laskowski05]

Protein Structure Databases



Useful resources:

- PDBsum
- [Jena](#)
- MSD



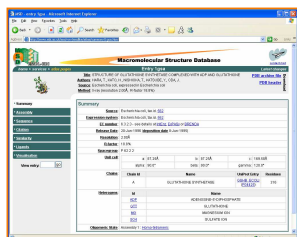
<http://www.imb-jena.de/IMB.html>

Protein Structure Databases



Useful resources:

- PDBsum
- Jena
- [MSD](http://www.ebi.ac.uk/msd/)



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

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

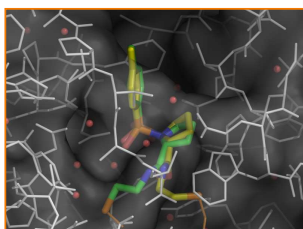
- Geometry
- Relationship to sequence
- Classifications

Protein Structure Visualization



Some tools:

- PyMOL
- OpenRasMOL
- Protein Explorer
- Grasp
- etc.



[pymol.sourceforge.net]

Protein Structure Visualization



Demo!

by Meghan Bellows

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- [Geometry](#)
- Relationship to sequence
- Classifications

Protein Structure Analysis



Bond types (for 1gsa):

• The color code is based on FDS (fold deviation score), defined as a multiple of the standard deviation for a specific reference value.

Bond Length							
Bond Type	Chain Id	Tot Num	Cal Ave	Cal StdDev	Std Val	Std StdDev	Minimum Maximum
C-N	A	296	1.33	0.009	1.329	0.014	1.30 1.36
C-NP	A	17	1.36	0.010	1.341	0.016	1.32 1.38
C-O	A	314	1.23	0.007	1.231	0.02	1.20 1.25
CA-C	A	289	1.52	0.011	1.525	0.021	1.48 1.56
CA-C(O)	A	25	1.58	0.013	1.516	0.018	1.50 1.54
CA-CB	A	207	1.52	0.017	1.53	0.02	1.47 1.56
CA-CB(A)	A	22	1.51	0.015	1.521	0.033	1.48 1.53
CA-CB(TV)	A	60	1.56	0.020	1.54	0.027	1.50 1.58
N-CA	A	272	1.45	0.009	1.458	0.019	1.44 1.47
N-CA(O)	A	25	1.45	0.009	1.451	0.016	1.44 1.47
N-CAP	A	17	1.46	0.009	1.466	0.015	1.44 1.48

<http://www.rcsb.org/pdb/explore/geometryDisplay.do?structureId=1GSA>

Protein Structure Analysis



Bond angles (for 1gsa):

Bond Angle	Chain Id	Tot Num	Cal Ave	Cal StdDev	Std Val	Std StdDev	Minimum	Maximum
C-N-CA	A	271	122.56	1.732	121.7	1.8	117.99	125.14
C-N-CA(S)	A	25	122.47	1.621	120.6	1.7	119.92	125.14
C-N-CA(F)	A	17	123.85	1.408	122.6	5.0	120.92	126.15
C-N-C	A	271	118.87	2.938	116.2	2.0	109.37	128.03
C-C-N(S)	A	25	118.14	2.494	116.4	2.1	112.92	121.37
C-C-N(F)	A	17	119.21	3.021	116.9	1.5	110.91	123.90
C-C-O	A	289	119.99	1.853	120.8	1.7	115.95	125.11
C-C-O(S)	A	25	120.19	1.390	120.8	2.1	117.54	123.07
CS-C-A	A	207	109.98	2.421	110.1	1.9	102.95	115.92
CS-C-A(A)	A	22	108.36	1.736	110.5	1.5	103.92	118.96
CS-C-A(T,V)	A	60	110.87	2.088	109.1	2.2	103.92	117.92
N-C-C	A	272	109.81	3.766	111.2	2.8	101.96	119.87
N-C-O(S)	A	25	112.84	3.143	112.5	2.9	106.37	120.06
N-C-O(F)	A	17	113.36	3.726	113.8	2.5	107.47	122.22
N-C-O	A	160	110.66	2.092	110.5	1.7	106.17	116.96
N-C-O(A)	A	22	110.21	1.444	110.4	1.5	107.56	119.30
N-C-O(B,T,V)	A	60	109.19	2.097	111.5	1.7	103.72	114.34
N-C-O(F)	A	17	103.76	1.082	103.0	1.1	100.94	106.80
O-C-N	A	296	121.16	1.772	123.0	1.6	114.96	126.08
O-C-N(F)	A	17	121.01	1.447	122.0	1.4	118.53	124.55

<http://www.rcsb.org/pdb/explore/geometryDisplay.do?structureId=1GSA>

Protein Structure Analysis



Dihedral angles (for 1gsa):

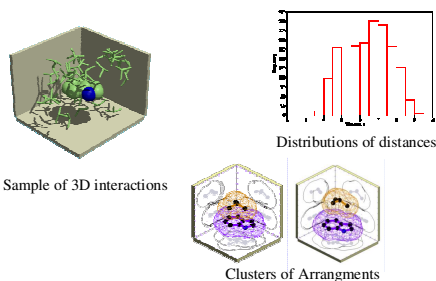
Dihedral Angle	Chain Id	Tot Num	Cal Ave	Cal StdDev	Std Val	Std StdDev	Minimum	Maximum
Chi1 (phi)	A	148	63.25	14.074	-66.7	15.0	114.90	7.50
Chi1 (psi)	A	47	54.73	14.908	64.1	15.7	21.20	86.19
Chi1 (tau)	A	21	168.68	11.202	183.6	16.8	165.92	169.90
Omega	A	313	177.51	14.847	180	5.8	8.90	193.60
Psi	A	180	30.54	62.456	-65.3	11.9	166.70	126.50
Psi (helix)	A	100	67.27	13.311	-65.3	11.9	164.88	126.50
Psi (loop)	A	14	58.48	9.022	-65.4	11.2	82.90	-58.20
Psi	A	183	36.11	75.600	-39.4	11.3	179.49	177.80
Psi (helix)	A	100	-31.19	16.235	-39.4	11.3	28.70	82.80
Psi (loop)	A	21	28.98	91.944	-39.4	11.3	177.49	161.90

<http://www.rcsb.org/pdb/explore/geometryDisplay.do?structureId=1GSA>

Protein Structure Analysis



Side-chain arrangements: TRP→LEU



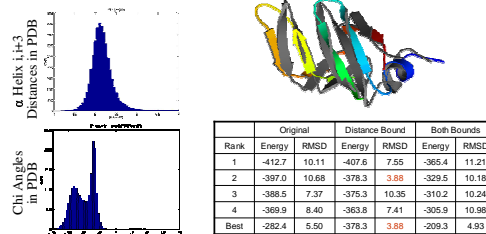
<http://www.biochem.ucl.ac.uk/bsm/sidechains/Trp/Leu/tp.html>

Protein Structure Analysis



Example application:

- Scott McAllister: "Generating Likely Distance Bounds for Efficient Protein Structure Prediction" (COS 597A)



Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- Geometry
- Relationship to sequence
- Classifications

Sequence → Structure?



If proteins have similar sequences they probably have similar structures

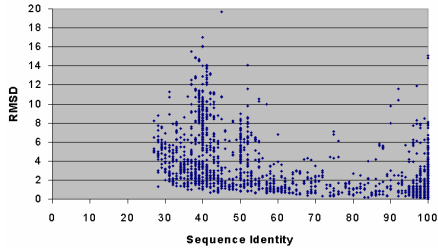
- >30% sequence identity
 - § Usually same structure & function
- 20-30% sequence identity
 - § Maybe related structure & function
 - § "Twilight zone"
- <20% sequence identity
 - § Unlikely to be related
 - § "Midnight zone"

Sequence → Structure?

Slide courtesy of Philip Bourne



Relationships between sequence and structure



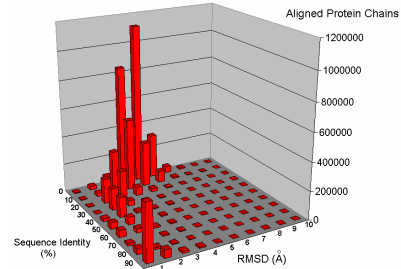
Structure Comparison of 30% of PDBSelect Set

Sequence → Structure?

Slide courtesy of Philip Bourne



Relationships between sequence and structure



Structure Alignments using CE with $\alpha > 4.0$

Sequence → Structure?

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

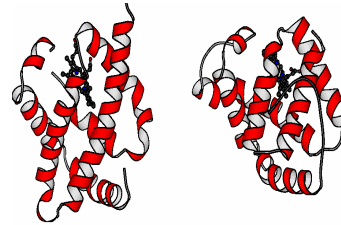


Sequence → Structure?

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?



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

Outline



Protein structure databases

- Primary
- Secondary
- Tertiary
- Quaternary

Protein structure visualization

- Demo

Protein structure analysis

- Geometry
- Relationship to sequence
- Ø Classifications

Protein Structure Classifications



Main databases:

- CATH
- SCOP

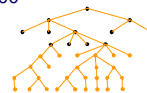
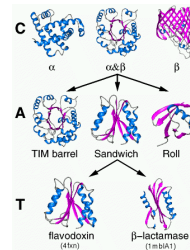
Protein Structure Classifications



CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S60
- S95
- S100

Structural Layout



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

Protein Structure Classifications

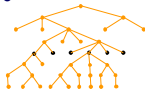


CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S60
- S95
- S100

Evolution

CATH Domain 1gsa01	
Classification	
Class	3
Alpha Beta	
Architecture	3.40
2-Layer(β) 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)
- S60
- S95
- S100

Sequence Identity

CATH Domain 1gsa01	
Classification	
Class	3
Alpha Beta	
Architecture	3.40
2-Layer(β) 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 (4)
- Architecture (40)
- Topology (1084)
- Homology (2091)
- S35 (7794)
- S60 (10363)
- S95 (13781)
- S100 (25491)

CATH Domain 1gsa01	
Classification	
Class	3
Alpha Beta	
Architecture	3.40
2-Layer(β) 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

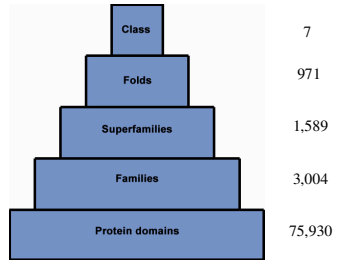
1. Root:	scop
2. Class:	Alpha and beta proteins (αβ) [51349]
3. Fold:	PreATP-grasp domain [52439]
4. Superfamily:	PreATP-grasp domain [52440]
5. Family:	Prokaryotic glutathione synthetase, N-terminal domain [52457]
6. Protein:	Prokaryotic glutathione synthetase, N-terminal domain [52458]
7. Species: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

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

Protein Structure Classifications



SCOP hierarchy:



SCOP: Structural Classification of Proteins (1.71 release)

Summary



Protein structure databases

- Lots of structural data is available

Protein structure analysis

- Structure provides information about function
- Need good tools for analysis!