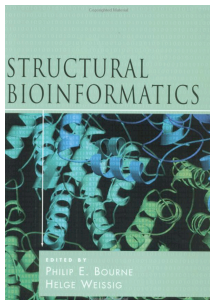


Protein Structure: Data Bases and Classification

Ingo Ruczinski

Department of Biostatistics, Johns Hopkins University

Reference

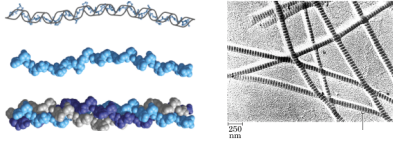


Bourne and Weissig
Structural Bioinformatics
Wiley, 2003

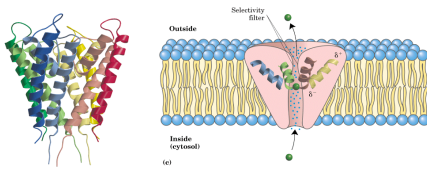
More References



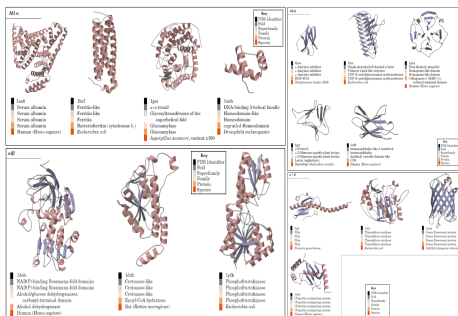
Structural Proteins



Membrane Proteins



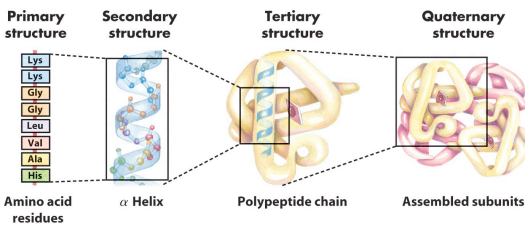
Globular Proteins



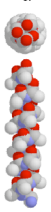
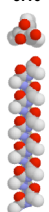
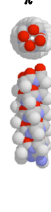
Terminology

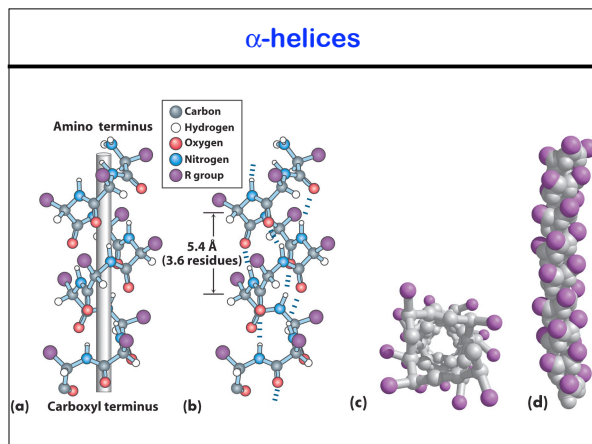
- Primary Structure
- Secondary Structure
- Tertiary Structure
- Quaternary Structure
- Supersecondary Structure
- Domain
- Fold

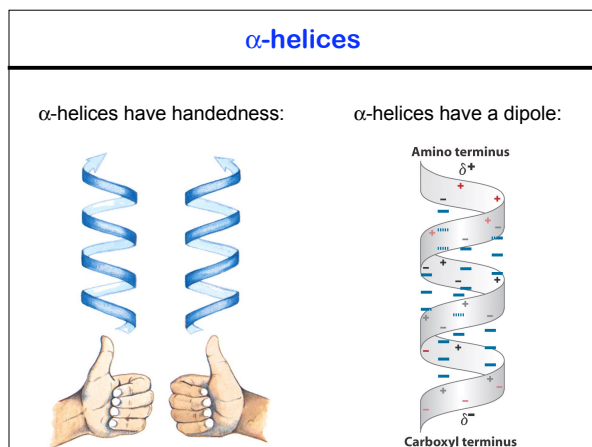
Hierarchy of Protein Structure

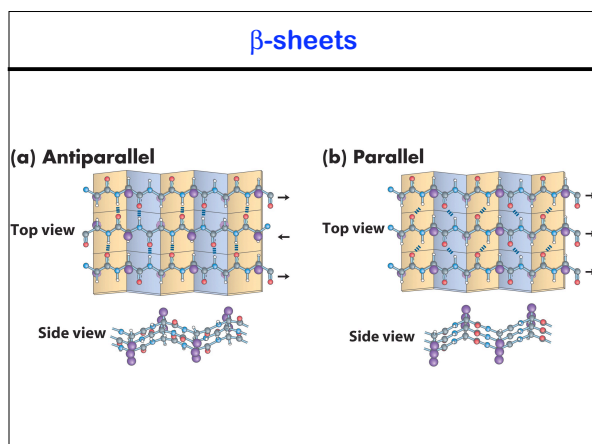


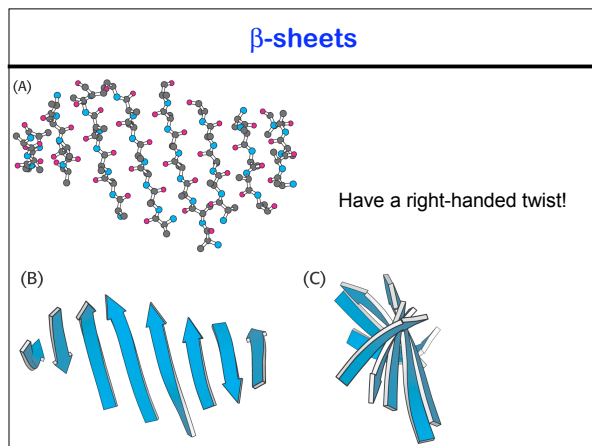
Helices

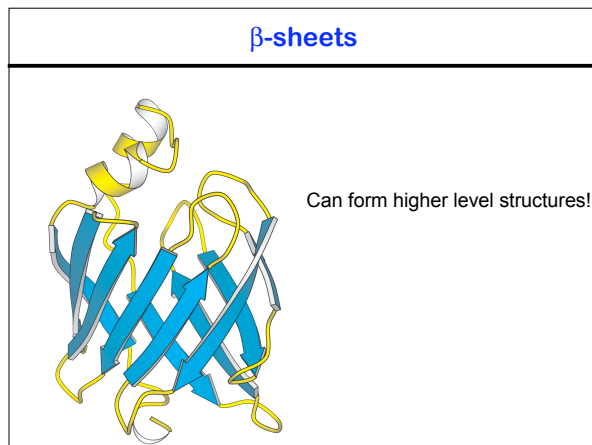
| | α | 3.10 | π |
|-------------------|---|---|---|
| |  |  |  |
| Amino acids/turn: | 3.6 | 3.0 | 4.4 |
| Frequency | ~97% | ~3% | rare |
| H-bonding | $i, i+4$ | $i, i+3$ | $i, i+5$ |

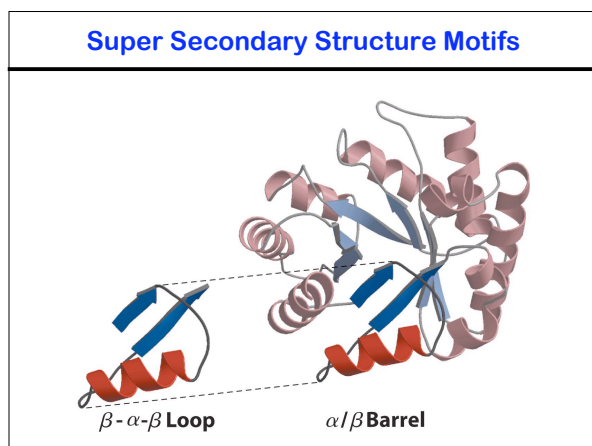




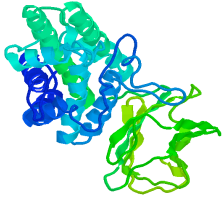








What is a Domain?



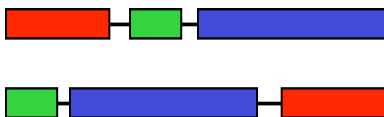
Richardson (1981):

Within a single subunit [polypeptide chain], contiguous portions of the polypeptide chain frequently fold into compact, local semi-independent units called domains.

More About Domains

- Independent folding units.
- Lots of within contacts, few outside.
- Domains create their own hydrophobic core.
- Regions usually conserved during recombination.
- Different domains of the same protein can have different functions.
- Domains of the same protein may or may not interact.

Why Look for Domains?



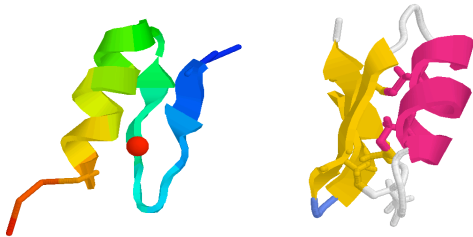
Domains are the currency of protein function!

Domain Size

- Domains can be between 25 and 500 residues long.
- Most are less than 200 residues.
- Domains can be smaller than 50 residues, but these need to be stabilized.

Examples are the zinc finger and a scorpion toxin.

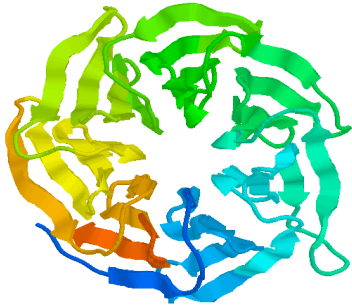
Two Very Small Domains



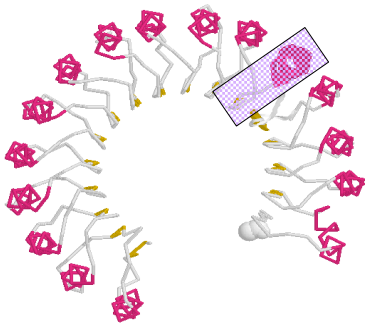
A Humdinger of a Domain



What's the Domain? (Part 1)



What's the Domain? (Part 2)



Homology and Analogy

- Homology: Similarity in characteristics resulting from shared ancestry.
- Analogy: The similarity of structure between two species that are not closely related, attributable to convergent evolution.

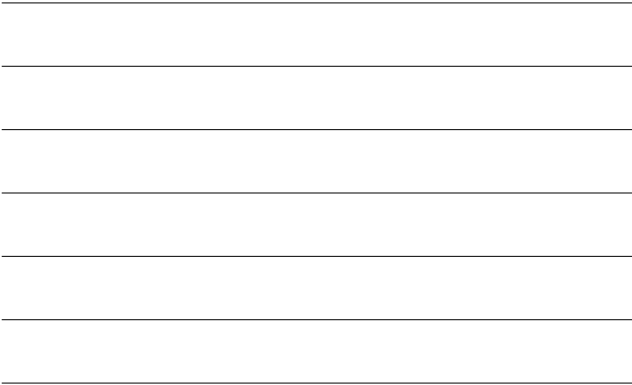
Homologous structures can be divided into orthologues (a result from changes in the same gene between different organisms, such as myoglobin) and paralogues (a result from gene duplication and subsequent changes within an organism and its descendants, such as hemoglobin).

Homology and Analogy

The diagram illustrates the concepts of homology and analogy using protein structures. It shows four protein structures arranged in a 2x2 grid, with labels and distances indicating their relationships.

- Enterotoxin** (top left, grey structure)
- Cholera toxin** (top right, pink structure)
- Remote homologue** (bottom left, pink structure) with a distance of **8.8%/35/2.4Å**
- TSS toxin** (bottom right, grey structure)
- Analogue** (bottom right, pink structure) with a distance of **4.4%/41/2.2Å**
- AA (RNA synthetase)** (bottom right, grey structure)

The diagram highlights the structural similarities between Enterotoxin and Cholera toxin, and the structural differences between Enterotoxin and TSS toxin, and between Enterotoxin and AA (RNA synthetase).



PDB File Header

The header contains information about protein and structure, date of the entry, references, crystallographic data, contents and positions of secondary structure elements, etc:

```

HEADER      OXIDOREDUCTASE              03-OCT-02      1M07
TITLE       ATOMIC RESOLUTION STRUCTURE OF CHOLESTEROL OXIDASE
TITLE       2 (STREPTOMYCES SP. SA-COO)
COMPND      MOL: 1; 1-
COMPND      2 MOLECULES: CHOLESTEROL OXIDASE;
COMPND      3 CHAIN: A;
COMPND      4 SYNONYM: CHOD;
COMPND      5 EC: 1.1.3.6;
COMPND      6 ENGINEERED: YES;
COMPND      7 OTHER_DETAILS: FAD COFACTOR NON-COVALENTLY BOUND TO THE
COMPND      8 ENZYME

AUTHOR      A.VRIELINK, P.I.LABIO
REVDAT      1 25-FEB-03 1M07 0
JRNL        AUTH  P.I.LABIO, N.SAMPSON, A.VRIELINK
JRNL        TITL  SUB-ATOMIC RESOLUTION CRYSTAL STRUCTURE OF
JRNL        TITL  2 CHOLESTEROL OXIDASE: WHAT ATOMIC RESOLUTION
JRNL        TITL  3 CRYSTALLOGRAPHY REVEALS ABOUT ENZYME MECHANISM AND
JRNL        1144 4 THE ROLE OF FAD COFACTOR IN REDOX ACTIVITY
JRNL        REF   J.MOL.BIOL. V. 326 1635 2003
JRNL        REFT  ATOM-UNIQUE (W. ISBN 0022-2836
  
```

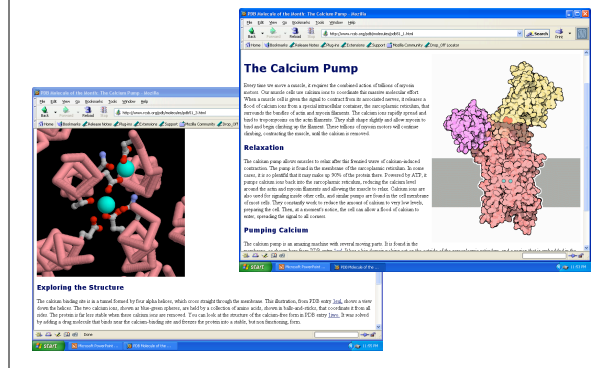
PDB File Body

The body of the PDB file contains information about the atoms in the structure:

```

ATOM       76  N  PRO  A  12      31.129  -4.659  43.245  1.00  9.00  N
ATOM       77  CA  PRO  A  12      32.426  -4.662  42.542  1.00  9.00  C
ATOM       78  C  PRO  A  12      32.423  -4.009  41.182  1.00  8.02  C
ATOM       79  O  PRO  A  12      33.267  -3.177  40.892  1.00  8.31  O
ATOM       80  CB  PRO  A  12      32.791  -6.126  42.592  1.00  10.02  C
ATOM       81  CG  PRO  A  12      32.190  -6.663  43.857  1.00  10.12  C
ATOM       82  CD  PRO  A  12      30.850  -5.927  43.925  1.00  9.87  C
ATOM       90  H  ALA  A  13      31.495  -4.468  40.316  1.00  8.06  H
ATOM       91  CA  ALA  A  13      31.357  -3.854  39.004  1.00  7.28  C
ATOM       92  C  ALA  A  13      29.947  -3.309  38.814  1.00  7.21  C
ATOM       93  O  ALA  A  13      28.969  -3.932  39.200  1.00  7.56  O
ATOM       94  CB  ALA  A  13      31.636  -4.879  37.897  1.00  8.54  C
  
```

Molecule of the Month

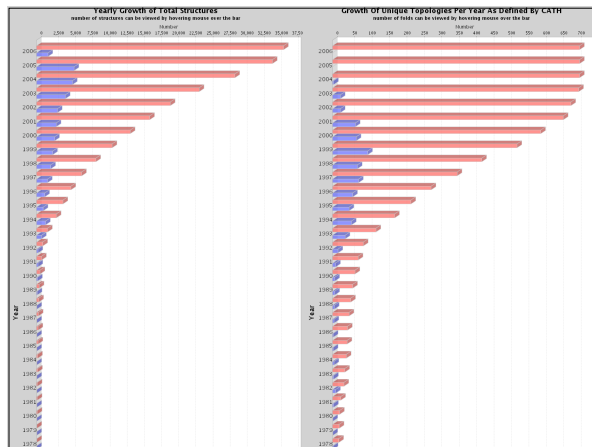


The Calcium Pump

Every time we move a muscle, it requires the constant action of billions of calcium pumps. Our muscle cells are calcium-rich to maintain the intense muscular effort. When a muscle cell gives the signal to contract, it releases calcium, which is a key factor in the contraction. This is a very important process, and the calcium pump is the key to maintaining the calcium level in the cell. The calcium pump is a transmembrane protein that is found in the cell membrane. It is a very important protein, and it is the key to maintaining the calcium level in the cell. The calcium pump is a transmembrane protein that is found in the cell membrane. It is a very important protein, and it is the key to maintaining the calcium level in the cell.

Exploring the Structure

The calcium binding site is a cleft formed by four glutamates, which cross straight through the membrane. The structure from PDB entry 1j2d shows a view down the helix. The two calcium ions, shown as blue spheres, are held by a network of water molecules, shown as red and white spheres. The protein structure is shown as a ribbon diagram. The calcium pump is a transmembrane protein that is found in the cell membrane. It is a very important protein, and it is the key to maintaining the calcium level in the cell.



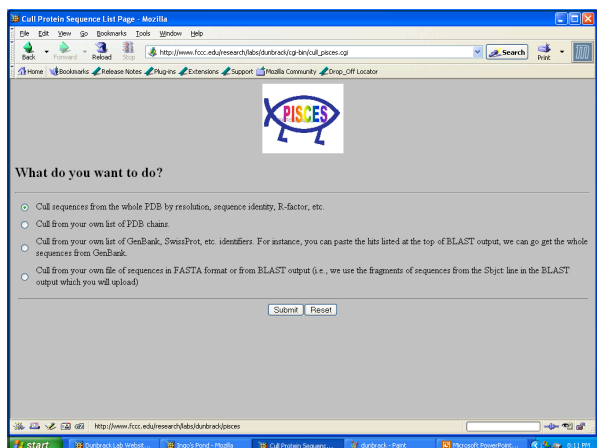


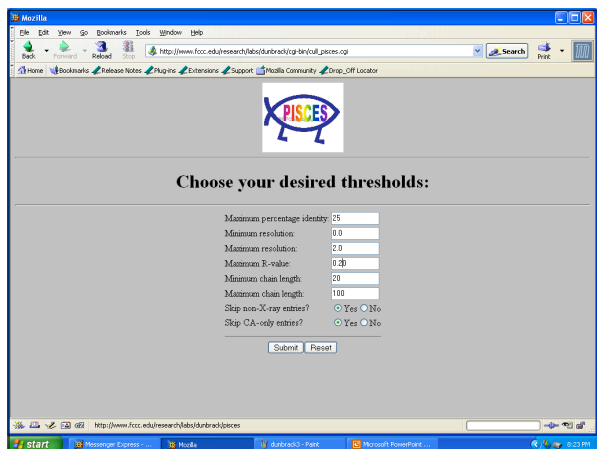
Dunbrack Lab

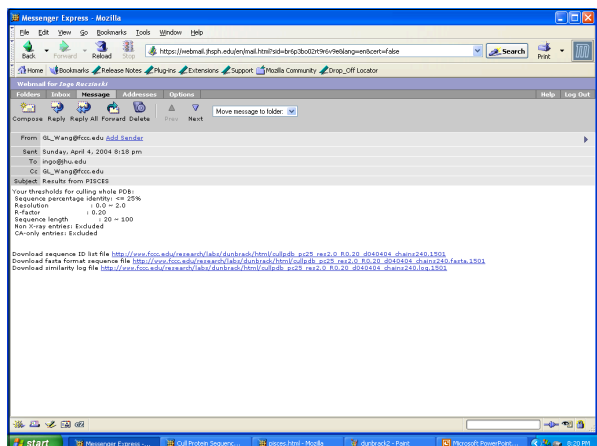
Home People Projects Publications Software Links

SCWRL3 PISCES BOLD

start Dunbrack Lab website Sign Up! Home







| ID# | length | Exp. 1 | resolution | R-factor | FreeRvalue |
|-------|--------|--------|------------|----------|------------|
| 17JSL | 52 | 2XAT | 1.900 | 0.20 | 0.28 |
| 18JQA | 91 | 2XAT | 0.970 | 0.14 | 0.15 |
| 1L9LA | 74 | 2XAT | 0.500 | 0.14 | 0.19 |
| 1Q9BA | 23 | 2XAT | 1.540 | 0.19 | 0.23 |
| 1G2YA | 32 | 2XAT | 1.000 | 0.20 | 0.20 |
| 1HT5A | 81 | 2XAT | 1.700 | 0.20 | 0.21 |
| 1Q7NB | 95 | 2XAT | 1.200 | 0.17 | 0.19 |
| 1Q0GA | 20 | 2XAT | 1.450 | 0.19 | 0.22 |
| 1K5YA | 46 | 2XAT | 0.540 | 0.09 | 0.09 |
| 1H24A | 72 | 2XAT | 1.700 | 0.19 | 0.22 |
| 1Q0BA | 99 | 2XAT | 1.900 | 0.18 | 0.21 |
| 1Q0EZ | 79 | 2XAT | 1.700 | 0.20 | 0.25 |
| 1Q0FY | 93 | 2XAT | 1.700 | 0.20 | 0.25 |
| 1HFOA | 95 | 2XAT | 1.250 | 0.13 | 0.17 |
| 1Q2BA | 62 | 2XAT | 1.120 | 0.15 | 0.20 |
| 1P7LA | 81 | 2XAT | 2.000 | 0.20 | 1.00 |
| 1R99D | 53 | 2XAT | 0.920 | 0.07 | 1.00 |
| 1RFPD | 80 | 2XAT | 1.800 | 0.19 | 1.00 |
| 1LATA | 82 | 2XAT | 1.900 | 0.20 | 0.28 |
| 11GQA | 62 | 2XAT | 1.700 | 0.20 | 0.23 |
| 1Q7LA | 69 | 2XAT | 1.800 | 0.20 | 0.22 |
| 1P7EA | 63 | 2XAT | 1.800 | 0.17 | 1.00 |
| 1E2GA | 84 | 2XAT | 1.400 | 0.16 | 0.20 |
| 1J8EA | 44 | 2XAT | 1.850 | 0.19 | 0.22 |
| 1JVEB | 77 | 2XAT | 1.800 | 0.17 | 1.00 |
| 1Q0GA | 74 | 2XAT | 0.930 | 0.10 | 0.13 |
| 1Q2BA | 100 | 2XAT | 1.350 | 0.16 | 0.18 |
| 1LARA | 77 | 2XAT | 2.000 | 0.19 | 0.22 |
| 1C0DA | 30 | 2XAT | 1.850 | 0.17 | 1.00 |
| 1PACD | 99 | 2XAT | 1.330 | 0.15 | 1.00 |
| 1H0QA | 93 | 2XAT | 1.240 | 0.17 | 0.19 |
| 1C75A | 71 | 2XAT | 0.970 | 0.12 | 1.00 |
| 112TA | 61 | 2XAT | 1.040 | 0.15 | 0.17 |
| 2R2DQ | 62 | 2XAT | 1.400 | 0.18 | 1.00 |
| 1RFPD | 55 | 2XAT | 1.700 | 0.17 | 0.23 |

SCOP

Structural Classification of Proteins

- Proteins are classified (manually!) taking both the structural and evolutionary relationship into account.
- There are 7 classes of proteins, the main ones being all alpha, all beta, alpha/beta, and alpha+beta.
- The principle levels in the hierarchy of SCOP are fold, superfamily, and family.

Murzin AG, Brenner SE, Hubbard T, and Chothia C (1995)

SCOP Levels

- Family:** Clear evolutionary relationship. In general >30% pairwise residue identities between the proteins.
- Superfamily:** Probable common evolutionary origin. Proteins have low sequence identities, but structural and functional features suggest that a common evolutionary origin is probable.
- Fold:** Major structural similarity. Proteins have the same major secondary structures in same arrangement and with the same topological connections.





SCOP: Structural Classification of Proteins

Scop Classification Statistics

SCOP: Structural Classification of Proteins, **1.69** release
25973 PDB Entries (1 Oct 2004), 70859 Domains, 1 Literature Reference
(excluding nucleic acids and theoretical models)

| Class | Number of folds | Number of superfamilies | Number of families |
|------------------------------------|-----------------|-------------------------|--------------------|
| All alpha proteins | 218 | 376 | 608 |
| All beta proteins | 144 | 290 | 560 |
| Alpha and beta proteins (a/b) | 136 | 222 | 629 |
| Alpha and beta proteins (a+b) | 279 | 409 | 717 |
| Multi-domain proteins | 46 | 46 | 61 |
| Membrane and cell surface proteins | 47 | 88 | 99 |
| Small proteins | 75 | 108 | 171 |
| Total | 945 | 1539 | 2845 |

Some Maybe Surprising Results

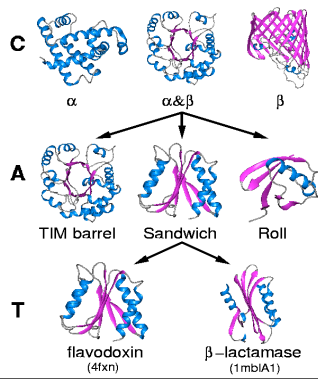
| 5NLL | 1AMO | 1CHN | 1FNB |
|---|---|---|---|
|  |  |  |  |
| Flavodoxin | Cytochrome reductase | Protein CHEY | Ferredoxin reductase |

CATH
Protein Structure Classification

- The CATH database is a hierarchical domain classification of protein structures in the Brookhaven protein databank. Only NMR structures and crystal structures solved to resolution better than 3.0 angstroms are considered.
- There are four major levels in this hierarchy: Class, Architecture, Topology (fold family) and Homologous superfamily.
- Multidomain proteins are subdivided into their domains using a consensus procedure. All the classification is performed on individual protein domains.

Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, and Thornton JM (1997)

The CATH Hierarchy



SCOP versus CATH

| SCOP | CATH |
|-------------|------------------------|
| Class | Class |
| | Architecture |
| Fold | Topology |
| | Homologous superfamily |
| Superfamily | |
| Family | Sequence family |
| Domain | Domain |

The screenshot shows the CATH database website. The main content area displays the 'CATH v2.6.0' release information, including the version number (2.6.0) and the release date (11-04-2005). Below this, there is a table showing the distribution of protein structures across different classes and architectures.

| Class | Architecture | Topology | Homologous superfamily | Sequence family | Domain | | |
|---------------------------------------|--------------|----------|------------------------|-----------------|--------|------|-------|
| Mainly Alpha | 5 | 251 | 465 | 1402 | 2189 | 3705 | 14105 |
| Mainly Beta | 19 | 160 | 311 | 1443 | 2061 | 4329 | 18771 |
| Alpha Beta | 14 | 414 | 706 | 3014 | 4781 | 7660 | 33080 |
| Few Secondary Structures | 1 | 82 | 90 | 144 | 232 | 285 | 1098 |
| Preliminary single domain assignments | 10 | 808 | 809 | 906 | 967 | 1090 | 3012 |
| Multi-domain domains | 1 | 12 | 12 | 16 | 25 | 36 | 109 |
| CATH-35 Sequence families | 1 | 4707 | 4707 | 4719 | 4768 | 4862 | 6168 |
| | 1 | 22 | 22 | 27 | 33 | 38 | 198 |

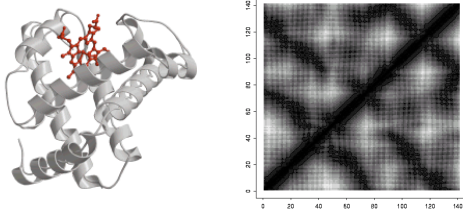
DALI

Distance Matrix Alignment

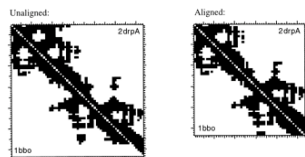
- DALI generates alignments of structural fragments, and is able to find alignments involving chain reversals and different topologies.
- The algorithm uses distance matrices to represent each structure to be compared.
- Application of DALI to the entire PDB produces two classifications of structures: FSSP and DDD (3D).

Holm L, and Sander C (1993)

DALI



DALI



Unaligned:

```
1bbo 1 KYICRQDIDKVFESHLGQILDTQVRFIDTVCHFSEKTDGLTFERSKAKGK 57
2drpa 103 PTKRGRHQIQRVCRVYVTHLQNECEKQYVZKGRVVYVYQPFCKRERKQNDKQOLLIK 165
```

Aligned:

```
1bbo 1 .....KYICRQDIDKVRKESHLGQILDTQVRFIDTVCHFSEKTDGLTFERSKAKK 57
2drpa 103 ftkegehZQQRVCRVYVTHLQNECEKQYVZKGRVVYVYQPFCKRERKQNDKQOLLIK... 165
```

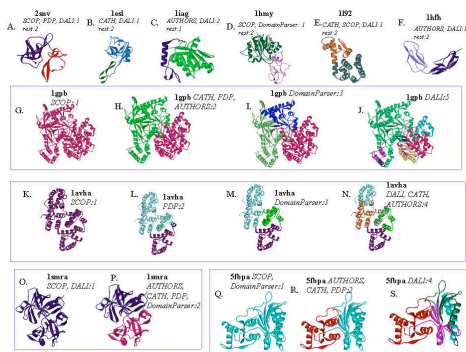
FSSP and DDD

- The families of structurally similar proteins (FSSP) is a database of structural alignments of proteins in the protein data bank (PDB). It presents the results of applying DALI to (almost) all chains of proteins in the PDB.
- The DALI domain dictionary (DDD) is a corresponding classification of recurrent domains automatically extracted from known proteins.

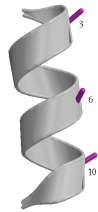
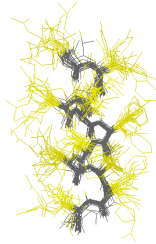
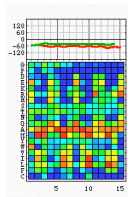
Other Algorithms for Domain Decomposition

- The Protein Domain Parser (PDP) uses compactness as a chief principle.
<http://123d.ncicfcrf.gov/pdp.html>
- DomainParser is graph theory based. The underlying principle used is that residue-residue contacts are denser within a domain than between domains.
<http://compbio.ornl.gov/structure/domainparser/>

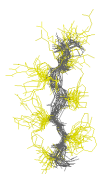
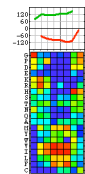
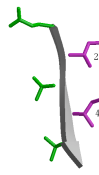
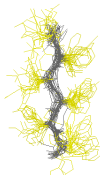
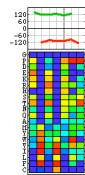
Oh Dear...



I-Sites



I-Sites



I-Sites

