

## Overview

---

- Week 2: Comparative methods and concepts
  - Similarity vs. Homology
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 3: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## Protein Conformation

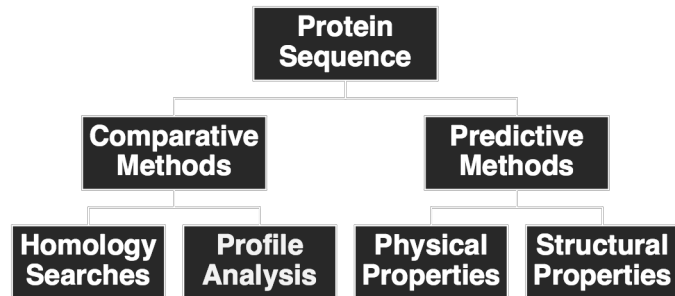
---

- Christian Anfinsen  
Studies on reversible denaturation →  
“Sequence specifies conformation”
- Chaperones and disulfide  
interchange enzymes:  
involved but not controlling final state
- “Starting with a newly-determined sequence,  
what can be determined computationally about  
its possible function and structure?”



## Protein Sequence Analysis

---



- *Common structure?*
- *Common function?*
- *Evolutionary relationship?*
- *Global or local similarity?*



## Sequence Comparisons

---

- Homology searches
  - Usually “one-against-one” *BLAST, FASTA*
  - Allows for comparison of individual sequences against databases comprised of individual sequences
- Profile searches
  - Uses collective characteristics of a family of proteins
  - Search can be “one-against-many” *Pfam, InterPro, CDD*
  - or “many-against-one” *PSI-BLAST*



## Profiles

- Numerical representations of multiple sequence alignments
- Depend upon *patterns* or *motifs* containing conserved residues
- Represent the common characteristics of a protein family
- Can find similarities between sequences with little or no sequence identity
- Allow for the analysis of distantly-related proteins



## Profile Construction

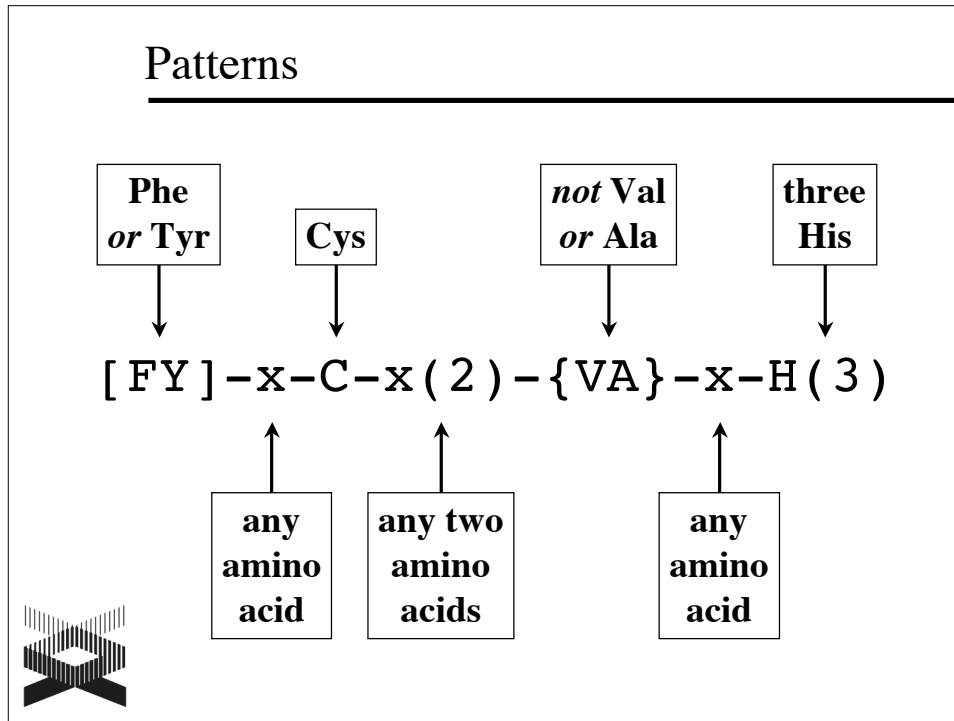
APHIIVATPG  
 GCEIIVATPG  
 GVEICIVATPG  
 GVDILIGTTG  
 RPHIIVATPG  
 KPHIIVATPG  
 KVQLIIVATPG  
 RPDIVIVATPG  
 APHIVGTTPG  
 APHIVGTTPG  
 GCHVVIVATPG  
 NQDIVVATTG


- Which residues are seen at each position?
- What is the frequency of observed residues?
- Which positions are conserved?
- Where can gaps be introduced?

Position-Specific Scoring Table

Cons	A	B	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y	Z
G	17	18	0	19	14	-22	31	0	-9	12	-15	-5	15	10	9	6	18	14	1	-15	-22	11
P	10	9	23	2	2	2	2	2	2	2	2	2	2	23	2	-2	12	11	17	-31	-8	1
H	5	24	-12	29	25	-20	8	32	-9	9	-10	-9	22	7	30	10	0	4	-8	-20	-7	27
I	-1	-12	6	-13	-11	33	-12	-13	63	-11	40	29	-15	-9	-14	-15	-6	7	50	-17	8	-11
V	3	-11	1	-11	-9	22	-3	-11	46	-9	37	30	-13	-3	-9	-13	-6	6	50	-19	2	-8
V	5	-9	9	-9	-9	19	-1	-13	57	-9	35	26	-13	-2	-11	-13	-4	9	58	-29	0	-9
A	54	15	12	20	17	-24	44	-6	-4	-1	-11	-5	12	19	9	-13	21	19	9	-39	-20	10
T	40	20	20	20	-30	40	-10	20	20	-10	0	20	30	30	-10	-10	30	150	20	-60	-30	10
F	21	6	27	6	6	13	10	13	0	6	16	13	13	89	17	17	24	22	9	-50	-48	12
G	10	60	20	10	30	150	-20	-30	-10	-50	-30	40	30	30	20	-30	60	40	20	-100	-70	30





- ### Pfam
- 
- Collection of multiple alignments of protein domains and conserved protein regions (regions which probably have structural or functional importance)
  - Each Pfam entry contains:
    - Multiple sequence alignment of family members
    - Protein domain architectures
    - Species distribution of family members
    - Information on known protein structures
    - Links to other protein family databases
- 

# Pfam

- Pfam A
  - Based on curated multiple alignments
  - Given the method used to construct the alignments, hits are highly likely to be true positives
  - > 74% of all known protein sequences have at least one match to Pfam
- Pfam B
  - Large number of small families taken from the PRODOM database; these families do not overlap with PfamA
  - Deemed “lower quality”, but can be useful when no Pfam A family is identified



wellcome trust  
sanger institute

Pfam

By UniProt Identifier

Enter a UniProt name or accession number

Submit Query Reset Example

Pfam has pre-calculated the domain structure of the proteins in UniProt. If you know the name or accession number (e.g. YAV\_HUMAN or Q91437) then you can see the Pfam domains on the sequence instantaneously.

By Protein sequence

Single sequence searches

If you don't know the UniProt identifier for your sequence, you can perform a slower, HMM search by giving your sequence below. Cut and Paste your sequence here (This search will take 1-5 minutes)

MAFSQYISLAPPELLLATAIFCLVFWLRCRTTQVPKGLKSPF YGDVLIQIRIGSTPVVVLISGLMTTKQALVKQGDDFKGRPDLYS DALKSFSIASDPTSPVSSCVLEEHVSKEMHHSKFKQLMAYV KSEEMNLVSKSKDFENVNTSGNAVDFPVLRYLPLNPKRF DITGALFKHSENYKDNGLIIPQEKIVNIVNDIFGAGPETVTI RDROPRLSDRPQLPYLEAF ILEIYRYTSFVPFTIPHSTTRDI DPFVFRPERFLTNDNTAIDKTLSEKVMFLGLGKRRICIGIPIF RSYGLTWKFRTCHEHVQAWRFK

Pfam Search Options

Search type: Both Global & Fragment Pfam search

Output format: Graphical output

\* Searching against SMART and TIGR hmms has been disabled. It should return shortly. \*

E-value cutoff level: 1.0

For help on the scores in Pfam, and the difference between standard and fragment searches, click here

Search Pfam Reset Example

Other regions to search for:

low-complexity (seq)

Large batch searches

To do large scale searching against Pfam, you can upload a TEXT file (Not Word) in FASTA format. This resource is primarily for people who do not have access to large computing facilities or personnel to install HMMER locally.

Browse... email address: Search type: Both Global & Fragment Pfam search

Search file against Pfam Reset

\* Searches larger than 1000 proteins, please split into separate files and upload each one separately.  
\* Please do not search proteins that are already in Pfam.

Done

NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II

Pfam: Results for Userseq

http://www.sanger.ac.uk/cgi-bin/Pfam/getblast?id=208L8381F76958976M916#UNKNOWN-QUERY:41:506:ls

welcome trust **sanger** institute

**Pfam** Search Pfam

Pfam Home Search by Browse by FTP IPfam Help About

Trusted matches - domains scoring higher than the gathering threshold (A) [513 residues]

Domain	Start	End	Bits	Evalue	Alignment	Mode
p450	41	506	368.50	9.2e-108	Align	ls

Alignments of Pfam-A domains to HMMs

Format for fetching alignments to seed | Jalview Java alignment viewer

Alignment of p450 vs UNKNOWN-QUERY/41-506

```

*->PpgptpPlfGnllqlgrgrlkdnlhsvftklakkyGpiftlylGpk
Ppdp +lP++G++l lg +n+h +tkl++ YG+++ ++G++
UNKNOWN-QU 41 PFGFWGLPFIHMLTLG-----KNPHLSLTKLSQQYGDVLQIRIGST 82

pvVlsgpeavkevLikkgeefsgrgdeawfytllvpflgkivfang.G
pvVlsg+ +k +L+k+g++f+g+d +y+++ ++gk + E+ ++G
UNKNOWN-QU 83 PVVVLSGLNTIKQALVKVQDDPKGRPD---LYSFTLITNGKSMFTNPDsG 129

erWrqlRrfltptrfrsfmgkklk.....sleprigeardLveklrkt
+W Rr+ ++ sf + +++ ++ le+ +ea+ L+ k++k
UNKNOWN-QU 130 PFWAARRRLAQDALKSFSI-ASDptsvsscYLEEHVSKEANHLISFKQL 178

agepgsGlviDitflkskaalnvIcsilFgkrfdeledpkflelvkavge
+e g +++++ + + nvI+ +Pgk+f + + l lvk ++
UNKNOWN-QU 179 MAEVBGH---FEPVNVQVSVANVIGAMCFGNFP-RKSEMLNLVKSSKD 224

lfsllspspgllldlfpilkkyfpgphlrlklrarkkldldklier
+++ +s++ +d+fp +l+y+p+p l++k+ + + lk +e++
UNKNOWN-QU 225 FVENV--TSGNAVDFFP-VLRYLPLPALKRFPNFDNVLSQLQTVQEHY 271

etldsagleeekkkksprDfdallLaknemekekdggeeskldeelr
+++++ s D + al ++ kd+g + e ++
UNKNOWN-QU 272 QDFNKN-----SIQDITGALFKH---SENKDNQ--GLIPEKIV 306

atvldlffAGteTTSsTLswaLyeLakhPevQeklreEidqvigdhrkei
v d+I+AG+eT ++++++ ++L + P+vQ+k++eS+d+vig++r
UNKNOWN-QU 307 NIVNDIFGAGFETVTTAIFWSILLVTEPKVQRKIHEELDTVIGRDR--- 353
    
```

Done

Pfam: p450

http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00067

welcome trust **sanger** institute

**Pfam** Search Pfam

Pfam Home Search by Browse by FTP IPfam Help About

Pfam entry p450

Accession number: PF00067 [Add Annotation](#)

**Cytochrome P450**

Cytochrome P450s are haem-thiolate proteins [6] involved in the oxidative degradation of various compounds. They are particularly well known for their role in the degradation of environmental toxins and mutagens. They can be divided into 4 classes, according to the method by which electrons from NAD(P)H are delivered to the catalytic site. Sequence conservation is relatively low within the family - there are only 3 absolutely conserved residues - but their general topography and structural fold are highly conserved. The conserved core is composed of a coil termed the 'meander', a four-helix bundle, helices 1 and 4, and two sets of beta-sheets. These constitute the haem-binding loop (with an absolutely conserved cysteine that serves as the 5th ligand for the haem iron), the proton-transfer groove and the absolutely conserved EXXR motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with microsomal membranes. their general enzymatic function is to catalyse regio-specific and stereospecific oxidation of non-activated hydrocarbons at physiological temperatures [6].

This family forms **interactions** with other Pfam families, to view them click [here](#)

**INTERPRO description (entry IPR001128)**

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes. P450 proteins of prokaryotic 3-component systems and fungal P450s (CYP56) belong to the B-class, all other known P450 proteins from distinct systems are of the E-class [PUBMED:7678494](#).

**QuickGO**

Function	heme binding (GO:0020037) iron ion binding (GO:0005506) monooxygenase activity (GO:0004497)
Process	electron transport (GO:0006118)

For additional annotation, see the [PROSITE](#) document PDOC00081 [ [Expasy](#) | [SRS-UK](#) ]

Alignment	Domain organisation
<input checked="" type="radio"/> Seed (51) <input type="radio"/> Full (6037) Format Coloured alignment <a href="#">Get alignment</a> <a href="#">View HMM logo</a> Further alignment options <a href="#">here</a> Help relating to Pfam alignments <a href="#">here</a>	<input checked="" type="radio"/> View 48 representative architectures <input type="radio"/> View architectures for 6037 proteins Zoom 0.5 pixels/aa <a href="#">View Graphic</a>

Done

NHGRI Current Topics in Genome Analysis 2006  
Biological Sequence Analysis II

Firefox browser window showing the Pfam database entry for Cytochrome P450 (PF00067). The page title is "Pfam: p450" and the URL is "http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00067".

**Cytochrome P450**

Cytochrome P450s are haem-thiolate proteins [6] involved in the oxidative degradation of various compounds. They are particularly well known for their role in the degradation of environmental toxins and mutagens. They can be divided into 4 classes, according to the method by which electrons from NAD(P)H are delivered to the catalytic site. Sequence conservation is relatively low within the family - there are only 3 absolutely conserved residues - but their general topography and structural fold are highly conserved. The conserved core is composed of a coil termed the 'meander', a four-helix bundle, helices J and K, and two sets of beta-sheets. These constitute the haem-binding loop (with an absolutely conserved cysteine that serves as the 5th ligand for the haem iron), the proton-transfer groove and the absolutely conserved EXXR motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with microsomal membranes. Their general enzymatic function is to catalyse regioselective and stereospecific oxidation of non-activated hydrocarbons at physiological temperatures [6].

This family forms **interactions** with other Pfam families, to view them click [here](#)

**Figure 1: 1f4t Oxidoreductase**  
Thermophilic p450: cyp119 from *Sulfolobus solfataricus* with 4-phenylimidazole bound  
1akd  Display pdb

**INTERPRO description (entry IPR001128)**

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class, all other known P450 proteins from distinct systems are of the E-class [PubMed/572645](#).

**QuickGO**

Function: heme binding (GO:0020037), iron ion binding (GO:0005506), monooxygenase activity (GO:0004492)  
Process: electron transport (GO:0006118)

For additional annotation, see the [PROSITE](#) document PDOC00081 [[Expasy](#) | [SRS](#) | [UK](#)]

**Alignment** (Seed (51) / Full (6037))  
Format: Coloured alignment  
Get alignment | View HMM logo  
Further alignment options [help](#)  
Help relating to Pfam alignments [here](#)

**Domain organisation**  
View 48 representative architectures  
View architectures for 6037 proteins  
Zoom 0.5 pixels/aa  
View Graphic

**Species Distribution**  
NEW! View alignments & domain organisation by species  
Tree depth: Show all levels  
View Species Tree

**Phylogenetic tree**  
Seed (51) / Full (6037)  
Download tree | ATV Applet  
The trees were generated using [QuickTree](#)  
To find out more about ATV phylogenetic tree-viewer [click here](#)


Mozilla Firefox browser window showing a sequence alignment of Cytochrome P450 proteins. The URL is "http://www.sanger.ac.uk/Software/Pfam/data/tm1/seed/PF00067.shtm.gz".

Sequence alignment showing conserved residues across various species. The alignment is color-coded to highlight conserved regions. The species listed on the left include:

- TCMO\_HBLR/34-499
- C75A2\_SOLME/37-498
- C76A2\_SOLME/36-500
- C77A2\_SOLME/43-509
- C71A\_ONMY/34-500
- CP17A\_CHICK/33-496
- CP17A\_MOUSE/28-497
- CP17A\_HUMAN/28-493
- CP17A\_BOVIN/28-493
- CP21A\_CHICK/33-488
- CP21A\_HUMAN/33-489
- CP27O\_RAY/30-486
- CP2C1\_RABIT/30-486
- CP2CC\_RAY/30-487
- CP2C7\_RAY/30-487
- CP2C1\_RABIT/30-487
- CP2CN\_RAY/34-491
- CP2AL\_RAY/33-489
- CP2P1\_HUMAN/31-488
- CP2DE\_BOVIN/24-484
- CP2DA\_RAY/37-497
- CP2D1\_RAY/37-497
- CP4F1\_RAY/52-515
- CP491\_HUMAN/47-501
- CP4A2\_RAY/52-499
- CP4AA\_RAY/52-504
- CP4C1\_BLAZI/37-502
- CP4D1\_DROME/34-507
- CP6B1\_PAPPO/31-495
- CP6A1\_MUSBO/35-500
- CP6A2\_DROME/32-502
- CP3A3\_HUMAN/38-493
- CP3A6\_RABIT/37-491
- CP3A1\_RAY/38-496
- CP19A\_MOUSE/48-488
- CP19A\_CHICK/47-487
- CP18A\_MOUSE/42-496
- C11B1\_BOVIN/42-499
- CP11A\_ONMY/50-507
- CP11A\_HUMAN/52-511
- CP11A\_BOVIN/52-510
- CP27A\_HUMAN/41-526
- P1D6\_FUSO/51-503
- CP53\_AGNPC/36-509
- CP17A\_HUMAN/32-497
- CP51\_YEAS7/57-521
- CPX1\_BACME/12-406

The alignment shows conserved residues across these species, with some residues highlighted in red and others in green. The alignment is shown in a standard format with the protein name and residue number on the left, and the amino acid sequence on the right.

Pfam: p450  
 http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00067



**Cytochrome P450** Add Annotation

Cytochrome P450s are haem-thiolate proteins [6] involved in the oxidative degradation of various compounds. They are particularly well known for their role in the degradation of environmental toxins and mutagens. They can be divided into 4 classes, according to the method by which electrons from NAD(P)H are delivered to the catalytic site. Sequence conservation is relatively low within the family - there are only 3 absolutely conserved residues - but their general topography and structural fold are highly conserved. The conserved core is composed of a coil termed the 'meander', a four-helix bundle, helices J and K, and two sets of beta-sheets. These constitute the haem-binding loop (with an absolutely conserved cysteine that serves as the 5th ligand for the haem iron), the proton-transfer groove and the absolutely conserved EXXR motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with microsomal membranes. Their general enzymatic function is to catalyse regioselective and stereospecific oxidation of non-activated hydrocarbons at physiological temperatures [6].  
 This family forms **interactions** with other Pfam families, to view them click [here](#)

**Figure 1: 1f4t Oxidoreductase**  
 Thermophilic p450: cyp119 from *Sulfolobus solfataricus* with 4-phenylimidazole bound  
 1akd

**INTERPRO description (entry IPR001128)**

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class. [PubMed:7678494](#)

**QuickGO**

Function	heme binding (GO:0020037) iron ion binding (GO:0005506) monooxygenase activity (GO:0004497)
Process	electron transport (GO:0006118)

For additional annotation, see the PROSITE document PDOC0081 [[Expasy](#) | [SRS-UK](#)]

**Alignment**

Seed (51)  Full (6037)

Format: Coloured alignment

Further alignment options [help](#)  
[Help relating to Pfam alignments here](#)

**Species Distribution**

NEW! View alignments & domain organisation by species  
 Tree depth: Show all levels

**Domain organisation**

View 48 representative architectures  
 View architectures for 6037 proteins

Zoom: 0.5 pixels/aa

**Phylogenetic tree**

Seed (51)  Full (6037)

The trees were generated using [Quicktree](#)  
 To find out more about ATV phylogenetic tree-viewer [click here](#)

Done

Pfam: Distinct domain architectures for p450  
 http://www.sanger.ac.uk/cgi-bin/Pfam/getallproteins.pl?name=p450&acc=PF00067&verbose=true&type=full&dc

**Distinct domain architectures for p450**

Domain Image Key (in order of priority): ss\_p, pfam, contab, mbnt, cold\_coil, transmembrane, low\_complexity, pfam

**5262 proteins with p450 architecture**

Q4STK8\_TETNG[ tetraodon nigrovindis (green puffer) chromosome undetermined scaf14158, whole genome shotgun sequence.(fragment)  
 [612 residues]

**328 proteins with p450, p450 architecture**

Q4S3E8\_TETNG[ tetraodon nigrovindis (green puffer) chromosome 1 scaf14751, whole genome shotgun sequence.(fragment)  
 [2030 residues]

**29 proteins with p450, Flavodoxin\_1, FAD\_binding\_1, NAD\_binding\_1 architecture**

Q89R90\_BRAJA[ bradyrhizobium japonicum] b12882 protein  
 [1078 residues]

**9 proteins with p450, FAD\_binding\_6, NAD\_binding\_1, Fer2 architecture**

Q4IRN1\_GIBZE[ gibberella zeae (fusarium graminearum)] hypothetical protein  
 [756 residues]

**5 proteins with p450, adh\_short architecture**

Q629N7\_BURMA[ burkholderia mallei (pseudomonas mallei)] cytochrome p450-related protein  
 [1373 residues]

**4 proteins with p450, p450, p450 architecture**

Q4TC47\_TETNG[ tetraodon nigrovindis (green puffer) chromosome undetermined scaf7053, whole genome shotgun sequence  
 [622 residues]

Done



Pfam: 29 proteins with p450~Flavodoxin\_1~FAD\_binding\_1~NAD\_binding\_1 architecture

http://www.sanger.ac.uk/cgi-bin/Pfam/getallproteins.pl?zoom\_factor=&domain\_view=arch&name=p450&type=ft

welcome trust sanger institute Pfam Search Pfam

Pfam Home Search by Browse by FTP IPfam Help About

**29 proteins with p450~Flavodoxin\_1~FAD\_binding\_1~NAD\_binding\_1 architecture**

Domain Image Key (in order of priority)

sg\_p pfam context smart coiled\_coil transmembrane low\_complexity pfam

Q415B2\_GIBZE[ gibberella zeae (fusarium graminearum)] hypothetical protein  
 p450 [1066 residues]

Q41LD6\_GIBZE[ gibberella zeae (fusarium graminearum)] c505\_fusox bifunctional p-450:nadph-p450 reductase (fatty acid omega-hydroxylase) (p450foxy)  
 p450 [1069 residues]

C505\_FUSOX[ fusarium oxysporum] bifunctional p-450:nadph-p450 reductase (fatty acid omega-hydroxylase)(p450foxy) [includes: cytochrome p450 505 (ec 1.14.14.1); nadph-cytochrome p450 reductase (ec 1.6.2.4)]  
 p450 [1066 residues]

Q8KUIJ\_ACTPA[ actinosynnema pretiosum subsp. auranticum] cytochrome p450  
 p450 [1005 residues]

Q7S8E0\_NEUCR[ neurospora crassa] hypothetical protein (probable bifunctional p-450:nadph-p450reductase)  
 p450 [1108 residues]

Q9HGE0\_GIBMO[ gibberella moniliformis (fusarium verticillioides)] fun6p  
 p450 [1115 residues]

Q3RZG6\_9BURK[ raletonia metallidurans ch34] cytochrome p450:oxidoreductase fad/nad(p)-binding:fad-binding:flavodoxin/nitric oxide synthase  
 p450 [1064 residues]

Q4WXE3\_ASPEU[ aspergillus fumigatus (sartorya fumigata)] fatty acid hydroxylase, putative  
 p450 [1120 residues]

Q81BF4\_BACCR[ bacillus cereus (strain atcc 14579 / dsm 31)] nadph-cytochrome p450 reductase (ec 1.6.2.4)  
 p450 [1065 residues]

Q82QD5\_STRAW[ streptomyces avermitilis] putative cytochrome p450  
 p450 [1073 residues]

Done

Pfam: p450

http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00067

welcome trust sanger institute Pfam Search Pfam

Pfam Home Search by Browse by FTP IPfam Help About

**Pfam entry p450**

Accession number: PF00067

**Cytochrome P450** [Add Annotation](#)

Cytochrome P450s are haem-thiolate proteins [6] involved in the oxidative degradation of various compounds. They are particularly well known for their role in the degradation of environmental toxins and mutagens. They can be divided into 4 classes, according to the method by which electrons from NAD(P)H are delivered to the catalytic site. Sequence conservation is relatively low within the family - there are only 3 absolutely conserved residues - but their general topography and structural fold are highly conserved. The conserved core is composed of a coil termed the 'meander', a four-helix bundle, helices 1 and K, and two sets of beta-sheets. These constitute the haem-binding loop (with an absolutely conserved cysteine that serves as the 5th ligand for the haem iron), the proton-transfer groove and the absolutely conserved EXXR motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with microsomal membranes. their general enzymatic function is to catalyse regio-specific and stereospecific oxidation of non-activated hydrocarbons at physiological temperatures [6].  
 This family forms **interactions** with other Pfam families, to view them click [here](#)

**INTERPRO description (entry IPR001128)**

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes. P450 proteins of prokaryotic 3-component systems and fungal P450s (CYP50) belong to the B-class, all other known P450 proteins from distinct systems are of the E-class [PUBMED:7678494](#).

**QuickGO**

Function	heme binding (GO:0020037) iron ion binding (GO:0005506) monooxygenase activity (GO:0004497)
Process	electron transport (GO:0006118)

For additional annotation, see the [PROSITE](#) document PDOC00081 [ [Expasy](#) | [SRS-UK](#) ]

Alignment	Domain organisation
<input checked="" type="radio"/> Seed (51) <input type="radio"/> Full (6037) Format Coloured alignment <a href="#">Get alignment</a> <a href="#">View HMM logo</a> Further alignment options <a href="#">here</a> Help relating to Pfam alignments <a href="#">here</a>	<input checked="" type="radio"/> View 48 representative architectures <input type="radio"/> View architectures for 6037 proteins Zoom 0.5 pixels/aa <a href="#">View Graphic</a>

Done

InterPro: IPR001128 Cytochrome P450

http://www.ebi.ac.uk/interpro/

EMBL-EBI European Bioinformatics Institute

InterPro home | Text Search | Sequence Search | Databases | Documentation | FTP site | Protein of the month

Search: [ ] Search Entries Search InterPro

InterPro IPR001128 Cytochrome P450

Overview: sorted by AC, sorted by name, of known structure, proteins with splice variants  
 Detailed: sorted by AC, sorted by name, of known structure, proteins with splice variants  
 Table: For all matching proteins, of known structure

Accession: IPR001128 Cytochrome\_P450 Matches: 7139 proteins

Type: Family

Database	ID	Name	Proteins
Gene3D	G3DSA:1.10.630.10	Cytochrome_P450	6884
Pfam	PF00067	p450	6584
PRINTS	PR00385	P450	5099
PROSITE pattern	PS00086	CYTOCHROME_P450	5227
PANTHER	PTHR19383	Cytochrome_P450	6582
SuperFamily	SF48264	Cytochrome_P450	6804

Signatures: [FW]-[SGNH]-x-[GD]-{F}-[RKHPT]-{P}-C-[LIVMFAP]-[GAD]

Children: IPR002397 B-class P450, IPR002399 Mitochondrial P450, IPR02401 E-class P450, group I, IPR02402 E-class P450, group II, IPR02403 E-class P450, group IV

Process: GO:0006118 electron transport

Function: GO:0004497 monooxygenase activity, GO:0005506 iron ion binding, GO:0020037 heme binding

Abstract: The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [1].

Structural links: CATH: 1.10.630.10, SCOP: 2.104.1.1

InterPro: IPR001128 Cytochrome P450

http://www.ebi.ac.uk/interpro/DisplayproEntry?ac=IPR001128

EMBL-EBI European Bioinformatics Institute

InterPro home | Text Search | Sequence Search | Databases | Documentation | FTP site | Protein of the month

Search: [ ] Search Entries Search InterPro

InterPro IPR001128 Cytochrome P450

Overview: sorted by AC, sorted by name, of known structure, proteins with splice variants  
 Detailed: sorted by AC, sorted by name, of known structure, proteins with splice variants  
 Table: For all matching proteins, of known structure

Accession: IPR001128 Cytochrome\_P450 Matches: 7139 proteins

Type: Family

Database	ID	Name	Proteins
Gene3D	G3DSA:1.10.630.10	Cytochrome_P450	6884
Pfam	PF00067	p450	6584
PRINTS	PR00385	P450	5099
PROSITE pattern	PS00086	CYTOCHROME_P450	5227
PANTHER	PTHR19383	Cytochrome_P450	6582
SuperFamily	SF48264	Cytochrome_P450	6804

Signatures: [FW]-[SGNH]-x-[GD]-{F}-[RKHPT]-{P}-C-[LIVMFAP]-[GAD]

Children: IPR002397 B-class P450, IPR002399 Mitochondrial P450, IPR02401 E-class P450, group I, IPR02402 E-class P450, group II, IPR02403 E-class P450, group IV

Process: GO:0006118 electron transport

Function: GO:0004497 monooxygenase activity, GO:0005506 iron ion binding, GO:0020037 heme binding

Abstract: The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [1].

Structural links: CATH: 1.10.630.10, SCOP: 2.104.1.1

Parent-Child Relationships (Subfamilies)

Child entries are more specific than the parent  
 A match to the child entry implies a match to the parent  
 Signatures for the parent and child entries must overlap

NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II

InterPro: IPR001128 Cytochrome P450

Structural links: CATH: 1.10.630.10  
 SCOP: a.104.1.1  
 PDB - click here

Database links: COME: PRX000238  
 PANDIT: PF00067  
 PROSITE doc: PDOC00081  
 Enzyme: 1.14  
 MSDsite: PS00088

**Taxonomic coverage**

**Overlapping InterPro entries**

InterPro ID	Numbers of overlapping proteins	Average numbers of overlap
IPR001128	6103	1036
IPR002397	7108	31
IPR002399	2797	4342
IPR002401	7069	70
IPR002402	6278	861
IPR002403	7124	15
IPR002439	7043	96
IPR002974	6950	189
IPR008066	7065	74
IPR008068	7084	55
IPR008069	7008	131
IPR008070	7102	37
IPR008071	7102	37
IPR008072	7015	124
IPR011347	7013	176

Center: Tree root  
 Inner circles: Tree nodes  
 Outer circles: Representative model organisms

*There is no significance to the placement of individual nodes on the circles*

InterPro: IPR001128 Cytochrome P450

**Example proteins**

Q22203 Cytochrome P450 98A3 (EC 1.14.-.-)

Q31449 Cytochrome P450 152A1 (EC 1.14.-.-) (P450BsBeta) (Fatty acid beta-hydroxylase)

Q48051 Probable cytochrome P450 4d14 (EC 1.14.-.-) (CYP1VD14)

P08684 Cytochrome P450 3A4 (EC 1.14.13.67) (Quinine 3-monooxygenase) (CYP11A4) (Nifedipine oxidase) (Taurochenodeoxycholate 6-alpha-hydroxylase) (EC 1.14.13.97) (NF-25) (P450-PCN1)

P23295 Cytochrome P450 55A1 (EC 1.14.-.-) (CYPLVA1) (P450 DNIR) (Nitric-oxide reductase) (P450nor)

**More proteins**

- IPR001128 Cytochrome P450
- IPR002397 B-class P450
- IPR002401 E-class P450, group I
- IPR008072 E-class P450, CYP3A
- ModBase
- CATH Domain
- SCOP Domain
- PDB Chain

**Publications**

1. Nelson D.R., Kamataki T., Waxman D.J., Guengerich F.P., Estabrook R.W., Feyereisen R., Gonzalez F.J., Coon M.J., Gunsalus I.C., Gotoh O. The P450 superfamily: update on new sequences, gene mapping, accession numbers, early trivial names of enzymes, and nomenclature. DNA Cell Biol. 12: 1-51 (1993) [PubMed: 7678494]

## Conserved Domain Database (CDD)

- Identify conserved domains in a protein sequence
- “Secondary database”
  - Pfam A and B
  - Simple Modular Architecture Research Tool (SMART)
  - Clusters of Orthologous Groups
- Search performed using RPS-BLAST
  - Query sequence is used to search a database of precalculated position-specific scoring tables
  - *Not* the same method used by ProfileScan



NCBI Conserved Domains

Query sequence: [(local sequence)|NP\_005206]  
 deleted in colorectal carcinoma [Homo sapiens]

Concise Result Full Result Show Search Information

Descriptions

Title	Pssmid	Multi-Dom	E-value
Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	2e-15
Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	5e-13
Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	9e-13
Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	6e-12
Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	2e-11
Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	2e-10
Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	2e-9
Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	5e-7
Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	6e-7
Hpfam06583, Neogenin_C, Neogenin C-terminus. This family represents the C-terminus of e...	46472	No	3e-105

CD Search Reference:  
 Marchler-Bauer A, Bryant SH (2004), "CD-Search: protein domain annotations on the fly.", *Nucleic Acids Res.*32(W)327-331.

NCBI Conserved Domains

Query sequence: [(local sequence)|NP\_005206]  
 deleted in colorectal carcinoma [Homo sapiens]

Concise Result Full Result Show Search Information

Descriptions

Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ... 28983 6e-15

cd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members are components of neural cell adhesion molecules (N-CAM L1), Fasciclin II and the insect immune protein Hemolin. The subfamily also includes receptor domains such as the extracellular ligand binding domain of Fibroblast Growth Factor Receptor 2. Members are phylogenetically diverse, occurring throughout metazoa, and are not components of the adaptive immune system molecules found in jawed vertebrates. A predominant feature of most Ig domains is a disulfide bridge connecting 2 beta-sheets with a Trp packing against the disulfide bond.

CD Length: 89, Pct. Aligned: 100, Bit Score: 79.775909, E-value: 2e-15

```

      10      20      30      40      50      60      70      80
query  330  PPWFLNHPNSLYAYESNDIEFECTVSGKVPVTVNMKMGDVPVPSD--YFQIVGGSNLRILGVVKSDEGFYQCVARENAG 407
consensus 1  PTFYQKPPFDIVVAGGEDVLECRASGNPPTTWLKNRGRFLSLDQGYTVLDNNGTLISNVTKEDAGTYTCVATNSAG 80
    
```

Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ... 28945 No 5e-13

Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ... 28945 No 9e-13

Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ... 28983 No 6e-12

Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ... 28945 No 2e-11

Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ... 28983 No 2e-10

Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ... 28945 No 2e-9

Hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ... 28945 No 5e-7

Hcd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ... 28983 No 6e-7

Hpfam06583, Neogenin\_C, Neogenin C-terminus. This family represents the C-terminus of e... 46472 No 3e-105

NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II

NCBI CDD cd00931

**Conserved Domains**

**cd00931.3** **IGcam**

**Links:**  
 Source: Smart  
 Taxonomy: Bilateria  
 PubMed: 3 links  
 Protein: cd00931 related architectures representatives

**Related CD:** 6 links

**Statistics:**  
 PSSM-ID: 28983  
 Aligned: 69 rows  
 Status: curated CD  
 Created: 1-Nov-2000  
 Updated: 18-Jul-2003

**Structure:**  
 Show Structure  
 Program: Cn3D  
 Drawing: All Atoms  
 Aligned Rows: up to 10 (download Cn3D)

cd00931 is part of a hierarchy of related CD models.  
 Use the graphical representation to navigate this hierarchy.  
 [mouse over icons to display CD accession and name]

Show Sequence Tree | cd00931 Branch | Whole Hierarchy

**Feature 1:** FGF/FGF-Receptor Interaction  
**Evidence:** Structure: 1E2: Receptor domain (chain F) contacts FGF (chain B) - View structure with Cn3D 4.1

Show Alignment | Format: Compact Hypertext | Row Display: up to 10 | Color Bits: 2.0 bits

NCBI CDD cd00931

Show Alignment | Format: Hypertext | Row Display: up to 20 | Color Bits: 2.0 bits

Type Selection: top listed sequences | Feature Display: FGF/FGF-Recep

```

    ....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
    Feature 1
    1BIH A      309 FKYeqk---pekVIVVkgq---gDVTIPCKVTGlp-aPNVVVShnakp-----lsggRATVTD-----SG 361
    1E2 F       8  PYWntekmekLHAVPaa---nTVKFRCPAGGnp-mPTMRWLKngkefk-----gehriGKVRNq-----hWS 69
    3NCM A      2  VMFkna---ptPQEFKeg---eDAVICDVVSal-pPTIIWKHgrdvi-----lkkdvRFIVLSn-----NY 57
    gi 462073   428 PSFadt---pgkEVWGWeg---hAVNLTCLAHslp-nATISWHFngadlfrgr--egqelqgtGYTLFGsg---prST 493
    gi 8134650   412 PSWlkk---pgdSQLEEG---kPGYLDCLTQAtp-kPTVVVYRngmlt-----sedsRFEVFNk-----GT 465
    gi 12644418 1098 PAFkqk---lqgVHVAEG---kKLLQCVSSdp-pATIWTLnkgt-----lktKFILSgeg---sLCS 1154
    gi 1169233   240 LYFlqr---psnVVALEg---kDAVLCCVSGyp-pPSTWLRgeevi-----glrskKYSLLGg---sSN 294
    gi 3024083   226 FVItv---prPRLGQal---qyDMDLECHIEAyp-pPAIVWTKddiqlannq-hysishfatADEYTD-----ST 287
    gi 2497323   219 FTites---ksnEATTGr---QASLKEASAvp-aDFEWTRddtrins-----anglelKSTEGQ-----SS 274
    gi 12644418 721 PWFis---kprSVTAsl---ggSVLISCALAGlp-fPTVWLRdgkal-----ckdtgHFEVLQne---dvFT 778
    gi 1705571    35 FVItqgs---prLWVFPt---dDISLKEASGkp-eQERWTRgdvghfk-----pkeelGVTVYQaph---sgFT 96
    gi 462740     41 PTItqg---spkDYIVDpr---enIVIQCEAKGkp-pPESWTRngthf-----didkDAQVTMkpn---sgTLV 100
    gi 3334268   401 PKFis---ngTIYYSwe---gpnPINISCDVKSnp-pASIHWRrdklvlpak---nttnlktYSTGRK-----MI 459
    gi 10720048 126 TVLvpv---skpTISVpssvtignRAVLTCEHDqspSEYSWFKdgismItadakktrafmsSPTIDPK-----sGD 196
    gi 1709202   139 RVLgg---lpDvvtIqg---gkALNLTcnVWGdp-pPEVSWLKnekalaq-----tdhcnlKFEAGRt-----AY 136
    gi 20178015 834 VRFapqvdhptpLTKVAaagdstsSATLHCRARGvp-nIDPTWTkngvpl-----dlqdpRYTEHkYhgqvvhSL 903
    gi 10720348 134 FAVpsc---evpSSALSg---tvVELRCQDKGpnPEVTFWKdgirllenp-rlgsqstnsYSTMNTk-----tGT 198
    gi 9789817   457 KIFsk---peKALIKhg---kKLRIVCEVSGgp-pPKVTFWKdkds-----inrKRNLYQEkhh-krrSE 513
    gi 12643619 313 PLKAt---isPRKVKsvs---gQVSLSCSVTGte-dQELSWYRngeiln-----pgknVRITGINh-----EN 369
    gi 8134650   320 AEIedmp-lfepRVFTAg---sEERVTCPLPKglpePSVWVWEGHagvr-----lptHGRVYQk-----gHE 375
    
```

90 100 110

```

    ....*.....|.....*.....|.....*.....|.....*.....
    Feature 1
    1BIH A      362 LVIKGVKn---gDKGYGCRATne---hgDKYF-ETLQV 393
    1E2 F       70 LIMESVvp---sDKGNVTCVVEne---ygSINH-TYHLD 101
    3NCM A      58 LQIRGIKk---tDEGTYRCEGRlArgeINFk-DIQVI 91
    gi 462073   494 LQVIFPNr---kMYGNYKCTATn---kHGTA-VHEIM 523
    gi 8134650   466 LRINSVEV---yDGTWYRCMSTp---agSIEaQAVLQ 497
    gi 12644418 1155 VSIEKALp---eDRGLYKCVAKnd---agQAEc-SCQVT 1186
    gi 1169233   295 LLISNVTk---dDSGMVTCVVYtykn-enISASaELTLV 328
    gi 3024083   288 LRVITVEk---rQYGDVVCkATnr---fGEAEaRVNLF 319
    gi 2497323   275 LVTNVTte---eHYGNVTCVAAnkl-gvTNAS-LVLF 307
    gi 12644418 779 LVLKKVOn---wHAGOVETILLKnr---vQECSC-OVSLM 810
    
```

Conserved Domains

Query sequence: [(local sequence)c|NP\_005206]  
 deleted in colorectal carcinoma [Homo sapiens]

CD Search Reference:  
 Marchler-Bauer A, Bryant SH (2004), "CD-Search: protein domain annotations on the fly.", *Nucleic Acids Res.*32(W)327-331.

Help | Disclaimer | Write to the Help Desk  
 NCBI | NLM | NIH

Title	PssmId	Multi-Dom	E-value
cd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	2e-15
cd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	5e-13
cd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	9e-13
cd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	6e-12
cd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	2e-11
cd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	2e-10
cd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	2e-9
cd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	5e-7
cd00931, Igcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	6e-7
cd006583, Neogenin_C, Neogenin C-terminus. This family represents the C-terminus of e...	46472	No	3e-105

CDART: Conserved Domain Architecture Retrieval Tool

Query: Neogenin\_C, FN3

Similar domain architectures:

- 3 Sequences: cellular program factor of CD ppa
- 33 Sequences: Neogenin
- 3 Sequences: Bilirubin PREDICTED similar
- 3 Sequences: Eukaryotic PREDICTED similar
- 74 Sequences: Dufanin Neogenin-binding pro
- 25 Sequences: Eukaryotic neogenin binding pro
- 2 Sequences: Caenorhabditis elegans Tumor-inhib Reson
- 4 Sequences: Caenorhabditis elegans SIKK4\_2 protein
- 5 Sequences: Caenorhabditis elegans Uncoordinated proSH3

Result page: Previous 1 2 3 4 5 6 7 8 9 10 11 Next

Subset by Taxonomy

Subset by selected domains:

- cd00063 Fibronectin type 3 domain; One of three types of ...
- includes: pfam00041 smart00060
- cd00931 Immunoglobulin domain cell adhesion molecule (cam...
- includes: cd00096 cd00099 pfam00047 smart00406 smart00407 smart00408 smart00409

## PSI-BLAST

- Position-Specific Iterated BLAST search
- Easy-to-use version of a profile-based search
  - Perform BLAST search against protein database
  - Use results to calculate a position-specific scoring matrix
  - PSSM replaces query for next round of searches
  - May be iterated until no new significant alignments are found
    - Convergence – all related sequences deemed found
    - Divergence – query is too broad, make cutoffs more stringent



NCBI BLAST

http://www.ncbi.nlm.nih.gov/BLAST

Latest news: 7 May 2004

The **Basic Local Alignment Search Tool (BLAST)** finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

<p><b>Nucleotide</b></p> <ul style="list-style-type: none"> <li>• Quickly search for highly similar sequences (megablast)</li> <li>• Quickly search for divergent sequences (discontiguous megablast)</li> <li>• Nucleotide-nucleotide BLAST (blastn)</li> <li>• Search for short, nearly exact matches</li> <li>• Search trace archives with megablast or discontiguous megablast</li> </ul>	<p><b>Protein</b></p> <ul style="list-style-type: none"> <li>• Protein-protein BLAST (blastp)</li> <li>• Position-specific iterated and pattern-hit initiated BLAST (PSI- and PHI-BLAST)</li> <li>• Search for short, nearly exact matches</li> <li>• Search the conserved domain database (rpsblast)</li> <li>• Protein homology by domain architecture (cdart)</li> </ul>
<p><b>Translated</b></p> <ul style="list-style-type: none"> <li>• Translated query vs. protein database (blastx)</li> <li>• Protein query vs. translated database (tblastn)</li> <li>• Translated query vs. translated database (tblastx)</li> </ul>	<p><b>Genomes</b></p> <ul style="list-style-type: none"> <li>• Human, mouse, rat, chimp, cow, pig, dog, sheep, cat</li> <li>• Chicken, puffer fish, zebrafish</li> <li>• Fly, honey bee, other insects</li> <li>• Microbes, environmental samples</li> <li>• Plants, nematodes</li> <li>• Fungi, protozoa, other eukaryotes</li> </ul>
<p><b>Special</b></p> <ul style="list-style-type: none"> <li>• Search for gene expression data (GEO BLAST)</li> <li>• Align two sequences (bl2seq)</li> <li>• Screen for vector contamination (VecScreen)</li> <li>• Immunoglobulin BLAST (IgBlast)</li> <li>• SNP BLAST</li> </ul>	<p><b>Meta</b></p> <ul style="list-style-type: none"> <li>• Retrieve results</li> </ul>

Disclaimer  
 Privacy statement  
 Accessibility  
 This page is valid XHTML 1.0.

Done



NHGRI Current Topics in Genome Analysis 2006  
Biological Sequence Analysis II

NCBI BLAST - Netscape

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi?CMD=Web&LAYOUT=TwoWindows&AUTO\_FORMAT=Semiauto&ALIGNMENTS=250&ALIGNMENT\_VIEW=Pairwise&CLIENT=

NCBI BLAST

protein-protein BLAST

Nucleotide Protein Translations Retrieve results for an RID

Search

>NP\_002119.1 high-mobility group box 1 [Homo sapiens]  
MGKGDPKKPRGKMSSYAFFVOTCREEHKKHPDASVNFSEFSKCKSERNKTMSAKEKGFED  
YEREMKTYIPPKGEBTKRKFDPNAPKRPSPAFFLFCSEYRPKIKGEHPGLSIGDVAKKLGEM  
KQPYEKRAAKLREKYEKDIAYRARGKPDAAKRGVVRAEKSKKKKEDEDEDEDEDEDEDE  
DDDD

Set subsequence From: To:

Choose database swissprot

Do CD-Search

Now: BLAST or FAST QUERY FAST3D

**Options** for advanced blasting

Limit by entrez query or select from: All organisms

Composition-based statistics

Choose filter  Low complexity  Mask for lookup table only  Mask lower case

Expect 10

Word Size 3

Matrix BLOSUM62 Gap Costs Existence: 11 Extension: 1

NCBI BLAST - Netscape

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi?CMD=Web&LAYOUT=TwoWindows&AUTO\_FORMAT=Semiauto&ALIGNMENTS=250&ALIGNMENT\_VIEW=Pairwise&CLIENT=

NCBI BLAST

Other advanced

PHI pattern

**Format**

Show  Graphical Overview  Linkout  Sequence Retrieval  NCBI-g Alignment in HTML format

Use new formatter  Masking Character Default (X for protein, n for nucleotide) Masking Color Black

Number of: Descriptions 500 Alignments 250

Alignment view Pairwise

Format for PSI-BLAST  with inclusion threshold: 0.001

Limit results by entrez query or select from: All organisms

Expect value range:

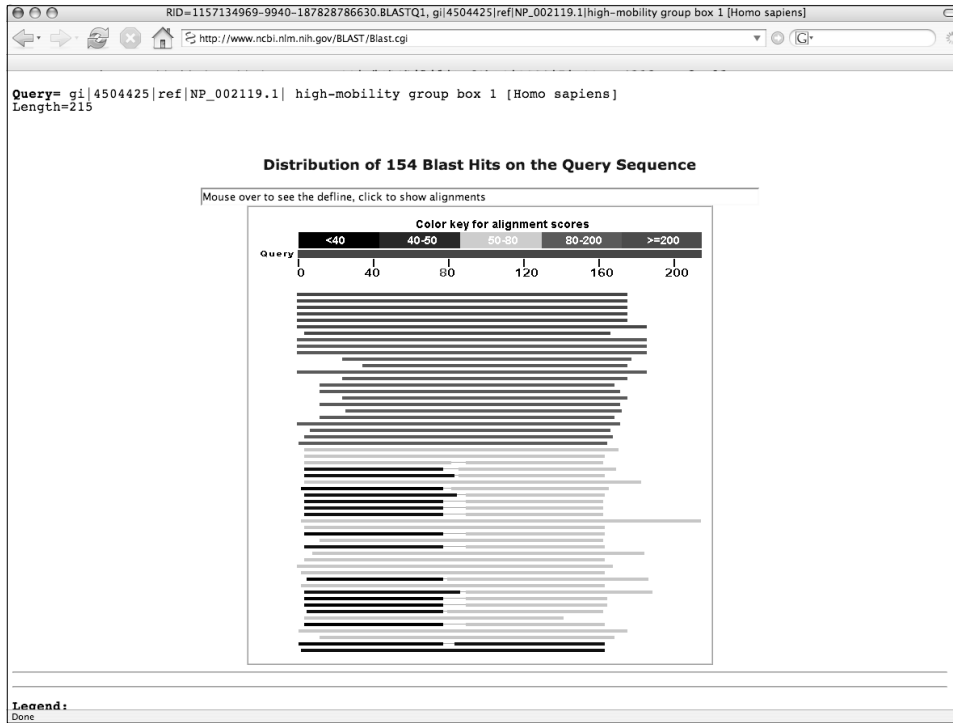
Layout: One Window Formatting options on page with results: None

Autoformat Semi-auto

BLAST FAST3D

Generate URL with preset values? FAST3D

NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II



RID=1157134969-9940-187828786630.BLASTQ1, gi|4504425|ref|NP\_002119.1|high-mobility group box 1 [Homo sapiens]

Legend:  
 NEW - means that the alignment score was below the threshold on the previous iteration  
 ● - means that the alignment was checked on the previous iteration

Run PSI-Blast iteration 2

Hit list size 500  
 Distance tree of results NEW

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

	Score	E	
	(Bits)	Value	
NEW <input type="checkbox"/> gi 123371 sp P12682 HMGB1_PIG High mobility group protein B1 ...	239	4e-63	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123367 sp P10103 HMGB1_BOVIN High mobility group protein B...	239	5e-63	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 75076928 sp Q4R844 HMGB1_MACFA High mobility group protein...	239	5e-63	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 52783618 sp P63159 HMGB1_RAT High mobility group protein B...	239	5e-63	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 20138433 sp Q9UGV6 HMGI1X_HUMAN High mobility group protein 1-	230	3e-60	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123373 sp P26584 HMGB2_CHICK High mobility group protein B...	203	4e-52	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123382 sp P07746 HMGT_ONCMY High mobility group-T protein (HM	201	2e-51	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 1708250 sp P52925 HMGB2_RAT High mobility group protein B2...	194	2e-49	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 1708259 sp P30681 HMGB2_MOUSE High mobility group protein ...	194	2e-49	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123374 sp P26583 HMGB2_HUMAN High mobility group protein B...	194	2e-49	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 13878931 sp P23497 SP100_HUMAN Nuclear autoantigen Sp-100 ...	193	4e-49	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123368 sp P07156 HMGB1_CRIGR High mobility group protein B...	189	5e-48	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 123375 sp P17741 HMGB2_PIG High mobility group protein B2 ...	187	2e-47	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 23396868 sp Q9N1Q6 SP100_GORGO Nuclear autoantigen Sp-100 ...	181	2e-45	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 729728 sp P40618 HMGB3_CHICK High mobility group protein B...	174	1e-43	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 20138160 sp O54879 HMGB3_MOUSE High mobility group protein...	174	2e-43	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 23396869 sp Q9N1Q7 SP100_PANTR Nuclear autoantigen Sp-100 ...	174	2e-43	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 85701353 sp O15347 HMGB3_HUMAN High mobility group protein...	174	2e-43	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 23396867 sp Q9N1Q5 SP100_HYLLA Nuclear autoantigen Sp-100 ...	170	2e-42	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 547652 sp P36194 HMGB1_CHICK High mobility group protein B...	170	4e-42	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 20138434 sp Q9UJ13 HMGL4_HUMAN High mobility group protein 4-	161	2e-39	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 17366497 sp Q24537 HMGL2_DROME High mobility group protein DSP	159	4e-39	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 729735 sp P40644 HMGH_STRPU High mobility group protein 1 hom	128	1e-29	<input type="checkbox"/>
NEW <input type="checkbox"/> gi 21903502 sp O09380 HMG12_CAEEL High mobility group protein 1	108	9e-24	<input type="checkbox"/>

Legend:  
Done

NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II

RID=1157134969-9940-187828786630.BLASTQ1, gj[4504425][ref][NP\_002119.1][high-mobility group box 1 [Homo sapiens]

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi

NEW	gi	sp	acc	description	bits	e-value	bits
NEW	gi	75263808	sp Q9LEF5 SSRP1_MAIZE	FACT complex subunit SSRP1 ...	43.1	6e-04	6
NEW	gi	729737	sp P40620 HMGL_VICFA	HMGI/2-like protein	43.1	6e-04	6
NEW	gi	47117886	sp Q04887 SOX9_MOUSE	Transcription factor SOX-9	43.1	7e-04	6
NEW	gi	11135387	sp Q9W757 SOX10_CHICK	Transcription factor SOX-10 (c	43.1	7e-04	6
NEW	gi	1351090	sp P48430 SOX2_CHICK	Transcription factor SOX-2	43.1	7e-04	6
NEW	gi	61216727	sp Q9BG91 SOX9_CALJA	Transcription factor SOX-9	42.7	7e-04	6
NEW	gi	1351096	sp P48436 SOX9_HUMAN	Transcription factor SOX-9 >g...	42.7	8e-04	6
NEW	gi	61216612	sp Q7YR77 SOX9_CANFA	Transcription factor SOX-9	42.7	8e-04	6
NEW	gi	10720294	sp P57073 SOX8_HUMAN	Transcription factor SOX-8	42.7	8e-04	6
NEW	gi	2506519	sp P35693 FPR1_PODAN	MAT+ sexual cell fertilization-p	42.7	8e-04	6
NEW	gi	12644232	sp P35713 SOX18_HUMAN	Transcription factor SOX-18	42.7	8e-04	6
NEW	gi	6175039	sp Q42569 SOX2_XENLA	Transcription factor SOX-2 (XLSO	42.7	9e-04	6
NEW	gi	38503365	sp Q9BG89 SOX9_PANTR	Transcription factor SOX-9	42.7	9e-04	6
NEW	gi	82186099	sp Q6P0E1 SOX2_BRARE	Transcription factor Sox-2	42.7	0.001	6

Run PSI-Blast iteration 2

Sequences with E-value WORSE than threshold

gi	6175075	sp P56693 SOX10_HUMAN	Transcription factor SOX-10	42.4	0.001	6
gi	2495255	sp Q03435 NHP10_YEAST	Non-histone protein 10 (High mo	42.4	0.001	6
gi	6175054	sp P36389 SRY_HORSE	Sex-determining region Y protein	42.4	0.001	6
gi	22654148	sp Q912W1 TFAM_RAT	Transcription factor A, mitochond	42.4	0.001	6
gi	6175076	sp Q04888 SOX10_MOUSE	Transcription factor SOX-10 ...	42.4	0.001	6
gi	82582249	sp Q6I248 SOX8_TETNG	Transcription factor SOX-8	42.4	0.001	6
gi	82183737	sp Q6EJ77 SOX3_BRARE	Transcription factor Sox-3	42.4	0.001	6
gi	6094380	sp Q55170 SOX10_RAT	Transcription factor SOX-10	42.4	0.001	6
gi	729738	sp P40621 HMGL_WHEAT	HMGI/2-like protein	42.4	0.001	6
gi	6831689	sp Q95416 SOX14_HUMAN	Transcription factor SOX-14 ...	42.0	0.001	6
gi	2506521	sp P48434 SOX9_CHICK	Transcription factor SOX-9	42.0	0.001	6
gi	24638225	sp Q9W7R6 SOX14_CHICK	Transcription factor SOX-14 (S	42.0	0.002	6
gi	19862533	sp Q04892 SOX14_MOUSE	Transcription factor SOX-14	42.0	0.002	6
gi	1351091	sp P48431 SOX2_HUMAN	Transcription factor SOX-2	42.0	0.002	6
gi	1711465	sp P54231 SOX2_SHEEP	Transcription factor SOX-2	42.0	0.002	6
gi	3913481	sp Q24533 DICH_DROME	SOX-domain protein dicaete (Pro	42.0	0.002	6
gi	12644266	sp P43267 SOX15_MOUSE	SOX-15 protein	42.0	0.002	6
gi	1723428	sp Q10241 CMB1_SCHPO	Mismatch-binding protein cmb1	41.6	0.002	6
gi	6094324	sp P48432 SOX2_MOUSE	Transcription factor SOX-2	41.6	0.002	6
gi	136654	sp P25977 UBF1_RAT	Nucleolar transcription factor 1...	41.6	0.002	6
gi	74684398	sp Q5KEP6 NHP6_CRYNE	Nonhistone chromosomal protein	41.6	0.002	6
gi	136652	sp P17480 UBF1_HUMAN	Nucleolar transcription factor...	41.6	0.002	6

Done

RID=1157135131-21179-60038174750.BLASTQ4, high-mobility group box 1 [Homo sapiens]

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi

**No new sequences were found above the 0.001 threshold!**

Taxonomy reports

Query= high-mobility group box 1 [Homo sapiens]  
 Length=215

**Distribution of 184 Blast Hits on the Query Sequence**

Mouse-over to show define and scores, click to show alignments

Color key for alignment scores

Color key for alignment scores
<40
40-50
60-80
80-200
>=200

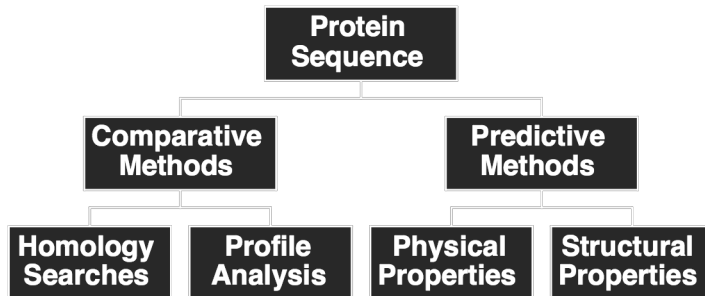
Query

0 40 80 120 160 200

① 154  
↓  
④ 184

Done

## Protein Sequence Analysis

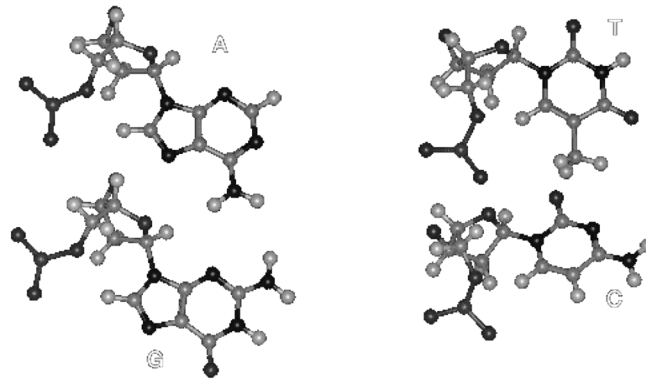


- *Common structure?*
- *Common function?*
- *Evolutionary relationship?*
- *Global or local similarity?*

- *Composition*
- *Hydrophobicity*
- *Secondary structure*
- *Specialized structures*
- *Tertiary structure*



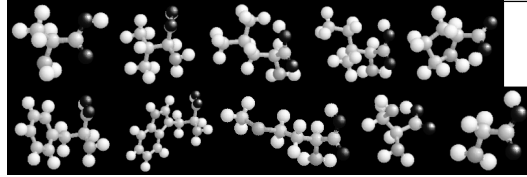
## Information Landscape



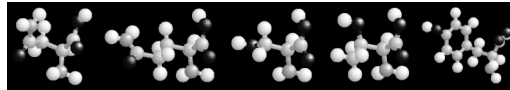
## Information Landscape

---

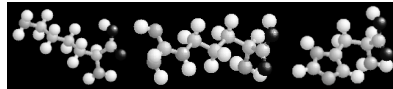
*Nonpolar*



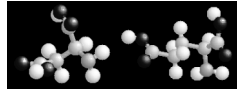
*Polar Neutral*



*Polar Basic*



*Polar Acidic*



## ProtParam

---

- Computes physicochemical parameters
  - Molecular weight
  - Theoretical pI
  - Amino acid composition
  - Extinction coefficient
- Simple query
  - SWISS-PROT accession number
  - User-entered sequence, in single-letter format
- <http://www.expasy.ch/tools/protparam.html>



## ProtParam Query

```
MNGEADCPTDLEMAAPKQDRWSQEDMLTLECMKNNLPSNDSSKFKTTESHMDWEKVAFKDFSGDMCKL  
KWVEISNEVRKFRTLTELEILDAQEHVKNPYKGGKLLKKHPDFPKKPLTPYFRFFMEKRAKYAKLHPEM...
```

↓ Compute parameters

Number of amino acids: 727  
Molecular weight: 84936.8  
Theoretical pI: 5.44

Amino acid composition:

Ala (A)	35	4.8%	Leu (L)	57	7.8%
Arg (R)	39	5.4%	Lys (K)	97	13.3%
Asn (N)	28	3.9%	Met (M)	25	3.4%
Asp (D)	58	8.0%	Phe (F)	18	2.5%
Cys (C)	6	0.8%	Pro (P)	39	5.4%
Gln (Q)	36	5.0%	Ser (S)	67	9.2%
Glu (E)	98	13.5%	Thr (T)	22	3.0%
Gly (G)	26	3.6%	Trp (W)	11	1.5%
His (H)	11	1.5%	Tyr (Y)	20	2.8%
Ile (I)	18	2.5%	Val (V)	16	2.2%
Asx (B)	0	0.0%			
Glx (Z)	0	0.0%			
Xaa (X)	0	0.0%			

Total number of negatively charged residues (Asp + Glu): 156  
Total number of positively charged residues (Arg + Lys): 136



## Expert Protein Analysis System (ExpASy)

- All tools available through a single Web front-end, at <http://us.expasy.org/tools>
- Primary sequence analysis tools include:

ProtParam

Compute pI/Mw

Titration Curve

ProtScale

*Plot any measurable characteristic*

*(e.g., hydrophobicity) by sequence position*

HelixWheel/HelixDraw

*Display protein sequence as a helical wheel*



## Secondary Structure Prediction

---

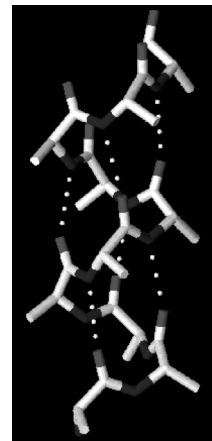
- Deduce the most likely position of alpha-helices and beta-strands
- Confirm structural or functional relationships when sequence similarity is weak
- Determine guidelines for rational selection of specific mutants for further laboratory study
- Basis for further structure-based studies



## Alpha-helix

---

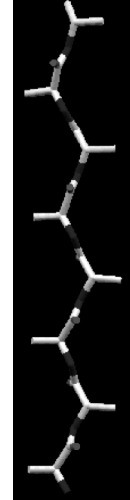
- Corkscrew
- Main chain forms backbone, side chains project out
- Hydrogen bonds between CO group at  $n$  and NH group at  $n+4$
- Helix-formers: Ala, Glu, Leu, Met
- Helix-breaker: Pro



## Beta-strand

---

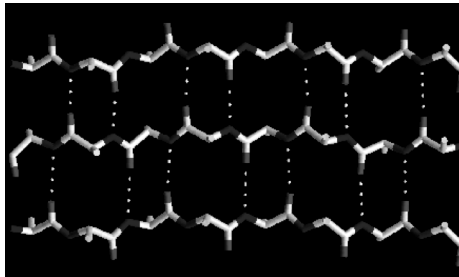
- Extended structure (“pleated”)
- Peptide bonds point in opposite directions
- Side chains point in opposite directions
- No hydrogen bonding *within* strand



## Beta-sheet

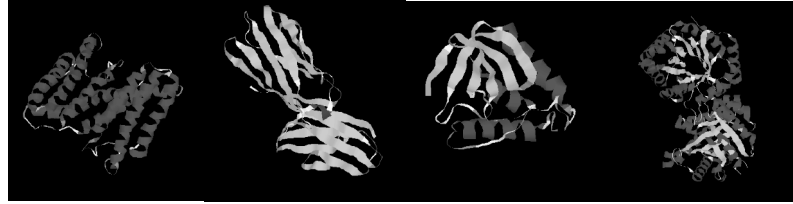
---

- Stabilization through hydrogen bonding
- Parallel or antiparallel
- Variant: beta-turn





## Folding Classes



$\alpha$

$\beta$

$\alpha+\beta$

$\alpha/\beta$

*Cyt c*

*CD4*

*Staph  
nuclease*

*Triose  
phosphate  
isomerase*

Globins  
Orthogonal  
EF-hand  
Up-Down  
Cytochrome

Orthogonal  
Super-barrel  
Greek key  
Sandwich  
Jelly roll

Split sandwich  
Meander  
Metal-rich  
Open roll  
OB/UB roll

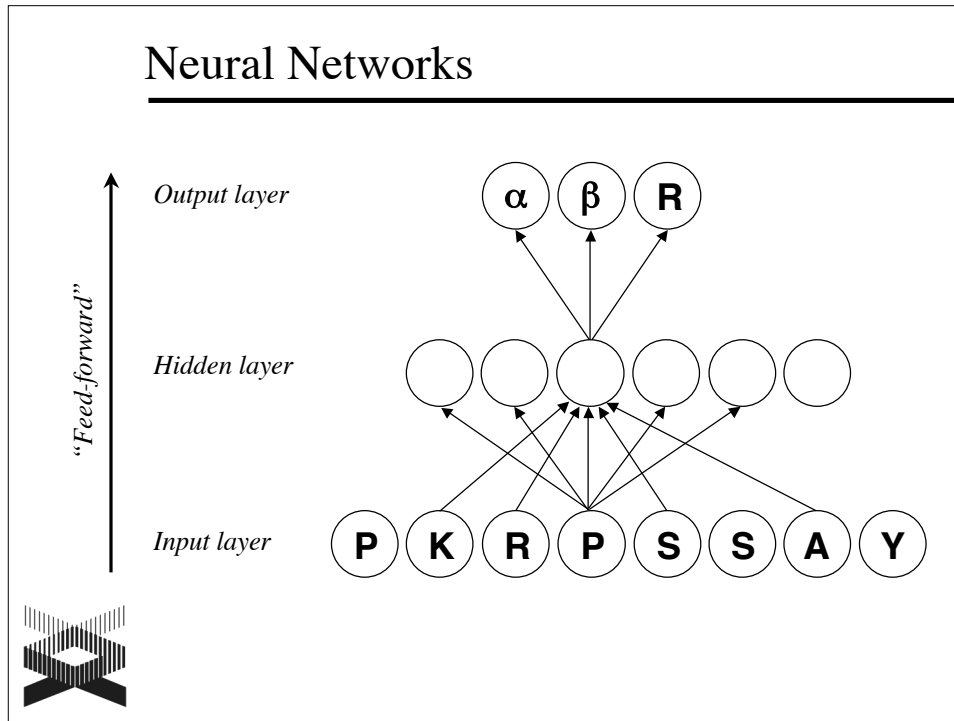
TIM barrel  
Doubly-wound



## Neural Networks


- Used when direct cause-and-effect rules between the beginning and end states are not known
  - Beginning and end states must be related
  - Neural networks attempt to deduce the relationship between the beginning and end states
- Supervised learning approach
  - Involves use of “training sets” where relationship is known
  - Based on data in training sets, network attempts to “learn” the relationship between input and output layers





### PredictProtein

- Multi-step predictive algorithm (*Rost et al., 1994*)
  - Protein sequence queried against SWISS-PROT
  - MaxHom used to generate iterative, profile-based multiple sequence alignment (*Sander and Schneider, 1991*)
  - Multiple alignment fed into neural network (PROFsec)
- Accuracy
  - Average > 70%
  - Best-case > 90%



PredictProtein - [Submission]

http://www.predictprotein.org

# PredictProtein

Home Submission Docs Downloads Register MetaPP

**If you would prefer using the old version of the site please click here**

Description of field (click on description to get help) Type the required information into the fields

Your email address [ watch typos -! ] andy@nhgri.nih.gov

Password (only for commercial users)

Results on the Predict Protein site, NOT in email (our current default)

**Input options (click to expand)**

**Prediction Type (click to expand)**

Which type of prediction do you require? secondary structure only (PHDsec)

Specify the format for the returned multiple-sequence alignment no alignment returned

Run iterated PSI-BLAST  PSI-BLAST

**Output options (click to expand)**

Return result in HTML  HTML formatted results  HTML for printouts  
 HTML with PHD graphs  HTML with PHD graphs for printouts

Switch off default methods (e.g. to reduce output, or to save time)  NO prosite  NO prodom  NO seg  NO coils  NO disulfind  
 NO PredictNLS  NO PHD  NO PROF  NO ASP

Return BLAST output from SWISS-PROT search  BLAST

Return additional PROF output  PROF rdb  PROF col  PROF casp  PROF msf  PROF only casp

Return additional PHD output  PHD msf  PHD rdb  PHD col  PHD casp

**Additional Services (click to expand)**

**Alignment options (click to expand)**

**Expert options (click to expand)**

PredictProtein - [Submission]

http://www.predictprotein.org/newwebsite/submit.html#%22xp\_opt\_click

Return result in HTML  HTML formatted results  HTML for printouts  
 HTML with PHD graphs  HTML with PHD graphs for printouts

Switch off default methods (e.g. to reduce output, or to save time)  NO prosite  NO prodom  NO seg  NO coils  NO disulfind  
 NO PredictNLS  NO PHD  NO PROF  NO ASP

Return BLAST output from SWISS-PROT search  BLAST

Return additional PROF output  PROF rdb  PROF col  PROF casp  PROF msf  PROF only casp

Return additional PHD output  PHD msf  PHD rdb  PHD col  PHD casp

**Additional Services (click to expand)**

**Alignment options (click to expand)**

**Expert options (click to expand)**

Paste, or type your sequence

- amino acids in one-letter code
- any number of white spaces allowed
- non-standard amino acids to 'X'

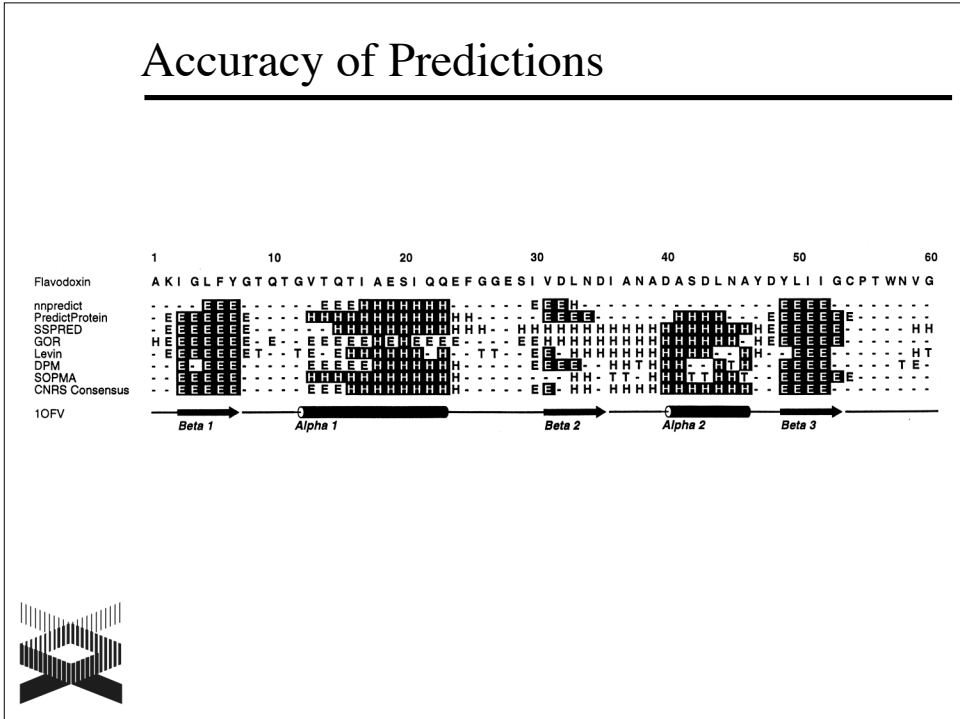
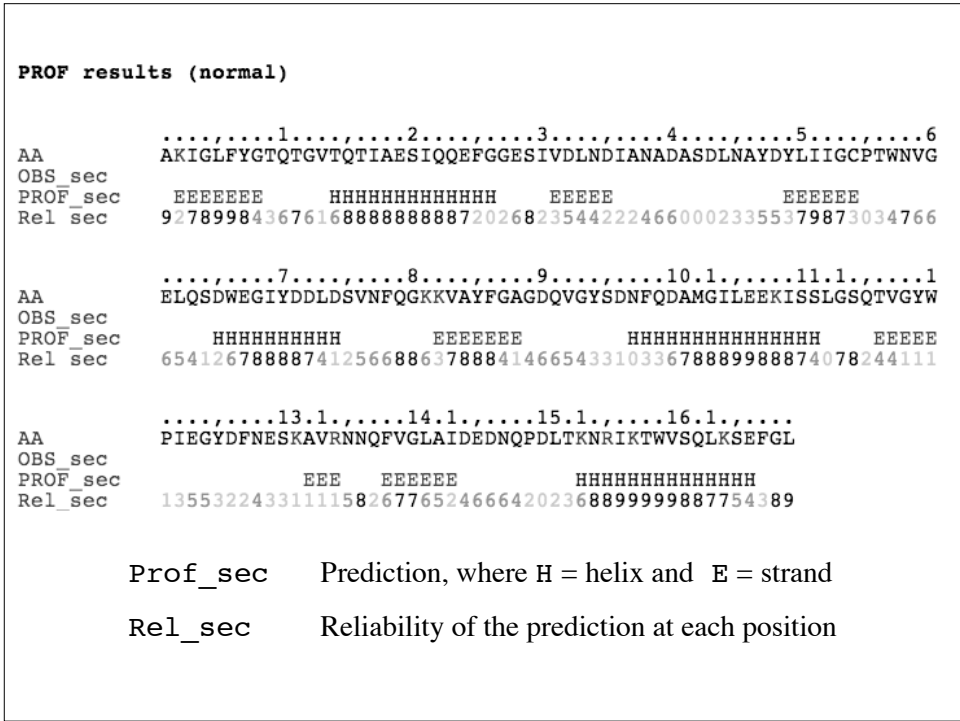
Use SRS6 to get your sequence from a public database

>Flavodoxin - Anacystis nidulans  
 AKIGLPVGTQTGVVQTIAESIQBFGGESIVDLNDIANADASDLNAYDYLIIGCPFTWNVGEI  
 DDLSVNFQKKVAYFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYWPIEGYDFNESKFLAIDEDNQDPLTKNRITWVSQLKSEFGL

Batch or interactive? interactive

Final action

Search



## Secondary Structure Prediction Methods

---

- PredictProtein  
*<http://www.predictprotein.org>*
- PSIPRED  
*<http://bioinf.cs.ucl.ac.uk/psipred/>*
- SAM-T99  
*<http://www.cse.ucsc.edu/research/compbio/HMM-apps/T99-query.html>*
- Jpred  
*<http://www.compbio.dundee.ac.uk/~www-jpred/submit.html>*



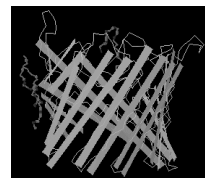
## Transmembrane Classes

---



- Helix bundles
  - Long stretches of apolar amino acids
  - Fold into transmembrane alpha-helices
  - “Positive-inside rule”

*Cell surface receptors*  
*Ion channels*  
*Active and passive transporters*



- Beta-barrel
  - Anti-parallel sheets rolled into cylinder

*Outer membrane of Gram-negative bacteria*  
*Porins (passive, selective diffusion)*



## TopPred

- Combines hydrophobicity analysis with the analysis of electrical charges
  - Calculates hydrophobicity profile
  - Hydrophobic-rich regions marked as “transmembrane”
  - Hydrophobic regions that fail to exceed a predefined cutoff are considered “putative transmembrane”
  - Topology prediction with and without putative helices



TopPred

<http://bioweb.pasteur.fr/seqan/interfaces/toppred.html>

**TopPred : Topology prediction of membrane proteins (Heijne, Wallin, Claros, Deveaud, Schuerer )**

Reset   your e-mail

( \* = required, \* = conditionally required )

• Sequence : please enter either :

1. the name of a file:

```
>gi|21431740|sp|Q18007
MPNYTVPPDPADTSWDSYPYIPVQIVVWIIIVLSLETIIGNAMVVMAYR
DLIIGIEGFPFTVYVNLGDRPLGWVACQWLEFLDYTCLVSIITVLLI
TKTQLLIVMSWLLPAIFGIMYIYQWAMTQOSTMSGAECSAPFLSNPYV
KGIHQAAKNLEKKAKAKERRHIALILSORLGTQVGVSLMLQSKAEKRAE
```

2. or the **actual data** here:

(sequence format)

Produce hydrophobicity graph image (-g)

Produce image of each topology (-t)

Control options

Output options

---

**Control options**

GES-scale (Goldman Engelman Steitz)  Hydrophobicity scale (-H)

Done

**TopPred : Topology prediction of membrane proteins (Heijne, Wallin, Claros, Deveaud, Schuerer )**

**Results:**

- gi\_21431740\_sp\_Q18007-1.png (4.33 Ko) ← **Models**
- gi\_21431740\_sp\_Q18007-2.png (4.12 Ko) ← **Models**
- gi\_21431740\_sp\_Q18007.png (6.29 Ko)
- gi\_21431740\_sp\_Q18007.hydro (7.72 Ko) ← **Hydrophobicity plot**
- toppred.out (3.36 Ko) ← **Text output**
- standard error file

From now, this files will remain accessible for 10 days at: <http://bioweb.pasteur.fr/seqanal/tmp/toppred/A25282311300796/>  
 You can save them individually by the **Save file** function if needed.

[Job summary](#) | [default format](#)

**Unix exact command:**  
 toppred -H GES-scale -g png query.data

**Your input data:**  
 query.data

[Help](#)

**References:**  
 von Heijne, G. (1992) Membrane Protein Structure Prediction: Hydrophobicity Analysis and the 'Positive Inside' Rule. J.Mol.Biol. 225, 487-494.  
 Claros, M.G., and von Heijne, G. (1994) TopPred II: An Improved Software For Membrane Protein Structure Predictions. CABIOS 10, 685-686.  
 Deveaud and Schuerer (Pasteur Institute) new implementation of the original toppred program, based on G. von Heijne algorithm.

*Pise CGI.generator version 5.a (04 Dec 2004 13:20)*  
 Done

# TopPred

---

Algorithm specific parameters:

```

Full window size : 21
Core window size : 11
Wedge window size: 5
Using hydrophobicity file: GES-scale

Cutoff for certain transmembrane segments: 1.00
Cutoff for putative transmembrane segments: 0.60
Critical distance between 2 transmembrane segments: 2

Critical loop length: 60

Kingdom: procaryote

Using cyt/ext file: CYTEXT-scale
    
```

Sequence : gi\_21431740\_sp\_Q18007 (713 res)

```

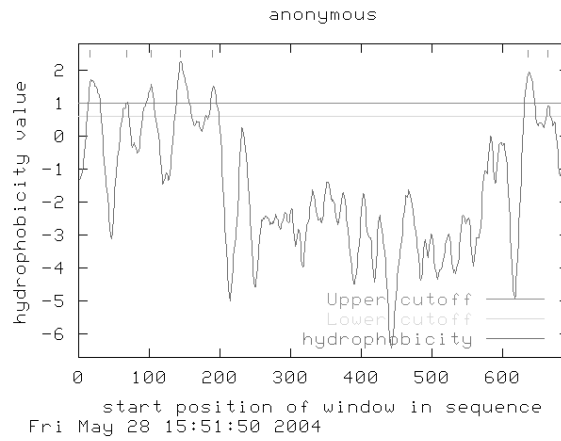
MPNVTVPDPADTSDSPYSIPVQIVVWIIIVLSLETHIGNAMVVMAYRIERNISKQVS
NRYIVSLAISDLIIGIEGPFPTVVLNGDRWFLGWACQWFLDYTLCLVSLITVLLI
TADRYLSVCHTAKYLKQSPKTKQLLVMSWLLPAIIFGIMYQWQAMTQSTSMGAC
SAPFLSNPYVNMGMVYAYWTTLVAMLILYKGIHQAAKLEKAKAKERRHIALILSQR
LGTQVGVSLMLOSAAEKAEAEAAQKDSGYTSNQAGDANLRRFGFSEPEPSTQFRVDPNSNN
NLNVEGSLWTEMDQNLGVIEERSGLSRRESNESYYPGPHFAANSRRCSEMEKVSLS
ESDGVSPTRPAKSYGRSLRSRYASSESITTTHEHDEKEVEKADSLQKLFADDELGSVLN
FKEEKLNKTDSDNDSDTTSVILQRSRKYKKNKPRSSKRSEHSTPRQIAKVQAEQTAQ
LIEESVDDDDQTEIEVKRTRDVRVSMKKRIARALIRRRSTRPERGSSNSDDSSSEVE
GEEKPEVRNNGLKIPQLTNNENRGETSSQPRDRLAPPNKTDTFLSASGVSRKISTIST
VITREKVISSIFAPIAVFNRRKQTKAEKRAHAFERTITFVGFALLWSPYIMATVYG
FCGEGCIPFLTLYSMCYLNSGPNFPAYALANRQFRSFAFMRMRGNFKVA
    
```

Found: 7 segments

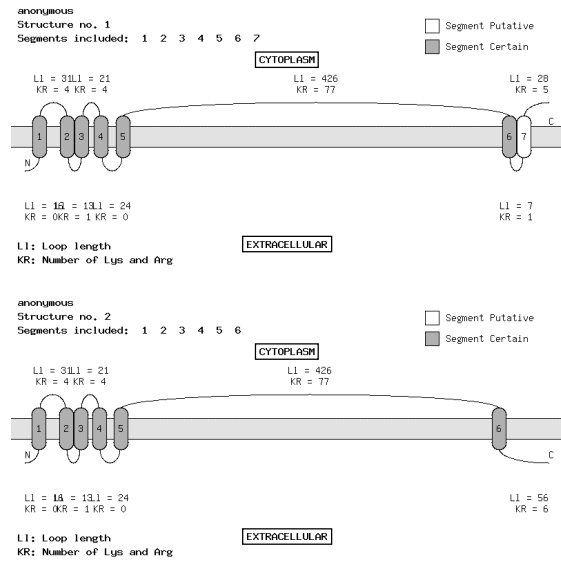
Candidate membrane-spanning segments:

Helix	Begin	End	Score	Certainty
1	17	37	1.717	Certain
2	69	89	1.024	Certain
3	103	123	1.555	Certain
4	145	165	2.264	Certain
5	190	210	1.531	Certain
6	637	657	1.931	Certain
7	665	685	0.920	Putative

# TopPred



# TopPred

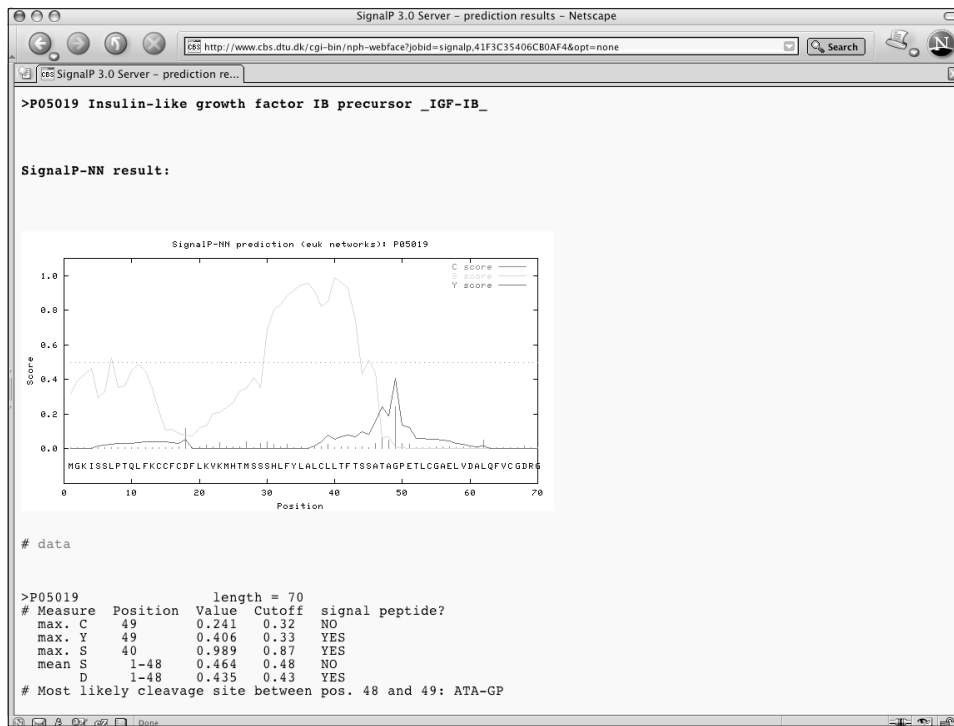




## SignalP

- Neural network trained based on phylogeny
  - Gram-negative prokaryotic
  - Gram-positive prokaryotic
  - Eukaryotic
- Predicts secretory signal peptides  
(*not* those involved in intracellular signal transduction)





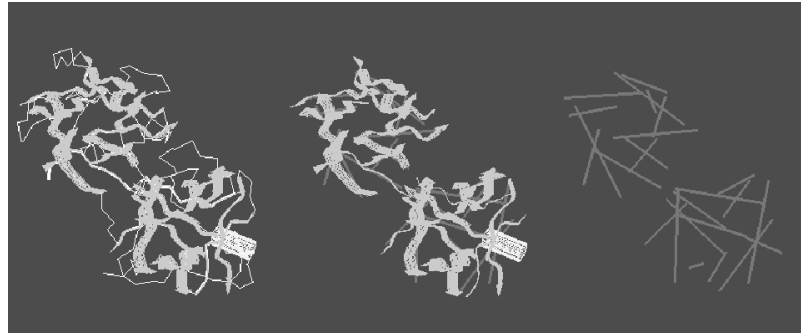
## Predicting Tertiary Structure

- Sequence specifies conformation, *but* conformation does *not* specify sequence
- Structure is conserved to a much greater extent than sequence
- Similarities between proteins may not necessarily be detected through “traditional” methods



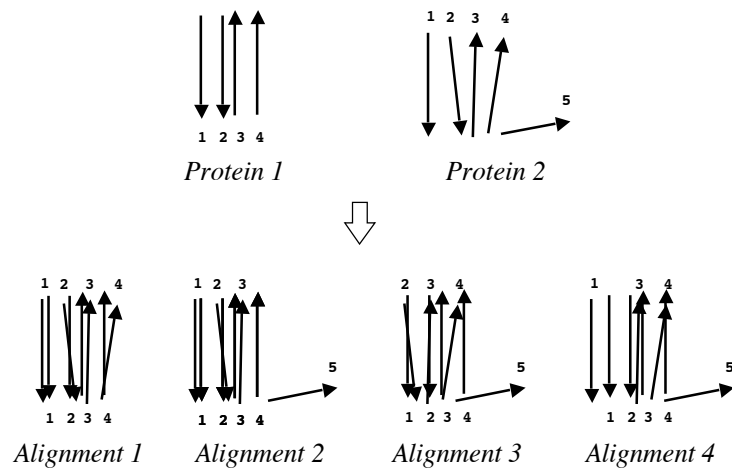
## VAST Structure Comparison

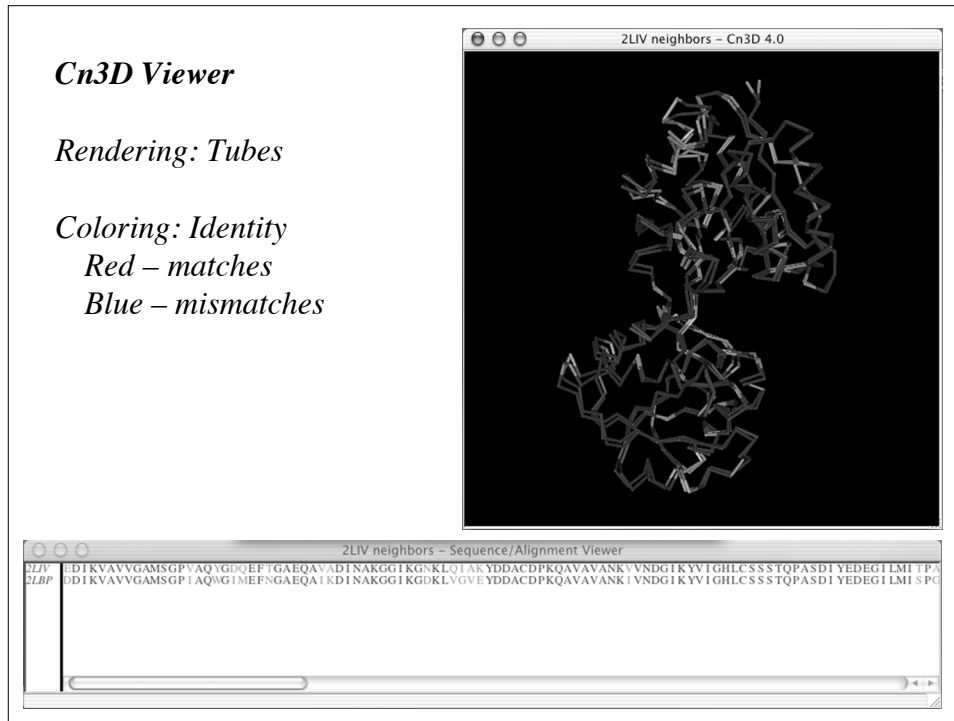
Step 1: Construct vectors for secondary structure elements



## VAST Structure Comparison

Step 2: Optimally align structure element vectors





## VAST Shortcomings

- Not the best method for determining structural similarities
- Reducing a structure to a series of vectors necessarily results in a loss of information (less confidence in prediction)
- Regardless of the “simplicity” of the method, provides a simple and fast first answer to the question of structural similarity



NHGRI Current Topics in Genome Analysis 2006  
 Biological Sequence Analysis II

The screenshot shows the NCBI Home Page with the following elements:

- Navigation Menu:** PubMed, All Databases, BLAST, OMIM, Books, TaxBrowser, Structure.
- Search Bar:** Search Structure for 2LIV Go
- Left Sidebar:**
  - SITE MAP: Alphabetical List, Resource Guide
  - About NCBI: An Introduction to NCBI
  - GenBank: Sequence submission support and software
  - Literature databases: PubMed, OMIM, Books, and PubMed Central
  - Molecular databases: Sequences, structures, and taxonomy
  - Genomic biology: The human genome, whole genomes, and related resources
  - Tools: Data mining
  - Research at NCBI: People, projects, and seminars
  - Software
- Main Content:**
  - What does NCBI do?** Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. More...
  - Hot Spots:**
    - Assembly Archive
    - Clusters of orthologous groups
    - Coffee Break, Genes & Disease, NCBI Handbook
    - Electronic PCR
    - Entrez Home
    - Entrez Tools
    - Gene expression omnibus (GEO)
    - Human genome resources
    - Influenza Virus Resource
    - Map Viewer
    - dbMHC
    - Mouse genome resources
    - My NCBI
    - ORF finder
    - Rat genome resources
    - Reference
  - Whole Genome Association:** The NCBI Whole Genome Association (WGA) resource provides researchers with access to genotype and associated phenotype information that will help elucidate the link between genes and disease. For more information, click here to see the the WGA resource page and click here to read the press release.
  - 100 Gigabases:** GenBank and its collaborating databases, the European Molecular Biology Laboratory and the DNA Data Bank of Japan, have reached a milestone of 100 billion bases from over 165,000 organisms. See the press release or find more information on GenBank.
  - PubMed Central:** An archive of life sciences journals
    - Free fulltext
    - Over 500,000 articles from over 200 journals
    - Linked to PubMed and fully searchable
 Use of PubMed Central requires no registration or fee. Access it from any computer with an Internet connection.

The screenshot shows the Entrez Structure search results page for the query 2LIV. The search results are as follows:

- Search Bar:** Search Structure for 2LIV Go Clear Save Search
- Display Options:** Display Summary Show 20 Send to
- Filters:** All: 1 Bacterial: 1 Eukaryotic: 0 Ligand: 0 NMR: 0 X-ray: 1
- Search Results:**
  - 1: 2LIV ← Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP) [mmbid:2778] VAST, Links
- Disclaimer:** Write to the Help Desk NCBI | NLM | NIH
- Footer:** Nov 3 2006 08:10:01

Structure Summary, 2LIV, 2778

http://www.ncbi.nlm.nih.gov/Structure/mmdb/mmdbsrv.cgi?form=6&db=t&Dopt=s&uid=2778

**NCBI** **MMDB**  
**Structure Summary**

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

**Reference:** Sack JS, Saper MA, Quiócho FA Periplasmic binding protein structure and function. Refined X-ray structures of the leucine/isoleucine/valine-binding protein and its complex with leucine J. *Mol. Biol.* v206, p.171-191  
 All References

**Description:** Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP).  
**Deposition:** 1989/4/10  
**Taxonomy:** Escherichia coli  
**MMDB:** 2778 **PDB:** 2LIV **Structure Neighbors:** VAST

View 3D Structure of All Atom Model Cn3D Display Download Cn3D!

Molecular components in the MMDB structure are listed below. The icons indicate macromolecular chains, 3D domains, protein classifications and ligands. Please hold the mouse over each icon for more information on the component.

Protein Chain 1 2 1 2  
 3d Domains 1 2  
 Domain Family LivK

Back to Home Page

Citing MMDB: Chen J, Anderson JB, DeWeese-Scott C, Fedorova ND, Geer LY, He S, Hurwitz DJ, Jackson JD, Jacobs AR, Lanczycki CJ, Liebert CA, Liu C, Madaj T, Marchler-Bauer A, Marchler GH, Mazumder R, Nikolskaya AN, Rao BS, Panchenko AR, Shoemaker

Vast Neighbor Summary

http://www.ncbi.nlm.nih.gov/Structure/vast/vastsrv.cgi?sdid=6728

**NCBI** **VAST**  
**Structure Neighbors**

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

VAST neighbors for: **MMDB 2778, 2LIV**

**Overview:** There are two main sections to this page. The first section consists of the alignment view controls, the list controls, and the advanced neighbor search controls. The second section is the VAST neighbor list itself.

View 3D Alignment of All Atoms with Cn3D Display Download Cn3D!

View Sequence Alignment using Hypertext for Selected VAST neighbors

List All sequences subset, sorted by Vast E-value in Table

**Advanced neighbor search**

Move the mouse over the red alignment footprints in the graphics below and click, you will obtain a structure-based sequence alignment.

Total neighbors: 4663; 1 - 60 of 667 representatives from the Medium redundancy subset displayed. Page: 1

Click to: Check All Uncheck All

2LIV 3d Dom. Protein Family Chain LivK Rli\_len

1215 B 344  
 1EHT B 332  
 10P4 C 306  
 1008 B 256  
 1600 B 253  
 1215 B 1 252  
 12HH B 245  
 1215 B 2 228

Vast Neighbor Summary

http://www.ncbi.nlm.nih.gov/Structure/vast/vastsvr.cgi?reqid=&sdid=6728&llbfid=14573001%2C4883801%2C42:

**VAST**  
 Structure Neighbors

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

VAST neighbors for: **MMDB 2778, 2LIV**

Overview: There are two main sections to this page. The first section consists of the alignment view controls, the list controls, and the advanced neighbor search controls. The second section is the VAST neighbor list itself.

View 3D Alignment of All Atoms with Cn3D Display Download Cn3D

View Sequence Alignment using Hypertext for Selected VAST neighbors

List All sequences subset, sorted by Vast E-value in Table

Advanced neighbor search

1 - 60 of 4663 neighbors displayed. Page: 1

Click to: Check All Uncheck All

PDB C D	Ali. Len	Score	E_Val	Rmsd	%Id	MMDB Date	LHM	GSP	Description
<input type="checkbox"/> 1Z15 A	344	42.1	10e-48.8	1.3	99.7	10/2005	0.0	0.4	Crystal Structure Analysis Of Periplasmic LeuLEVAL-Binding Protein In Superopen Form
<input checked="" type="checkbox"/> 2LBP	344	39.8	10e-44.6	0.9	79.1	03/2001	0.2	0.3	Leucine-Binding Protein (LBP)
<input type="checkbox"/> 1USGA	343	40.1	10e-42.4	2.0	79.0	01/2004	0.2	0.6	L-Leucine-Binding Protein, Apo Form
<input type="checkbox"/> 1YK0 A	323	29.9	10e-22.6	4.6	14.6	05/2006	6.2	1.5	Structure Of Natriuretic Peptide Receptor-C Complexed With Atrial Natriuretic Peptide
<input type="checkbox"/> 1JDP B	310	29.9	10e-22.6	4.3	14.8	10/2001	6.2	1.5	Crystal Structure Of Hormone RECEPTOR

Done

P-value  $\leq 0.001$   
 and  
 % Identity > 25  
 over at least 20 residues

Read the descriptions!

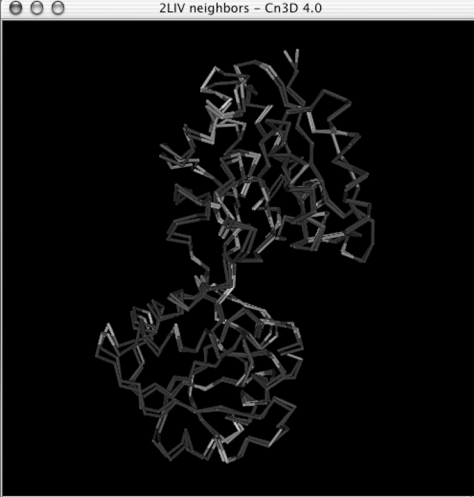
**Cn3D Viewer**

Rendering: Tubes

Coloring: Identity

Red – matches

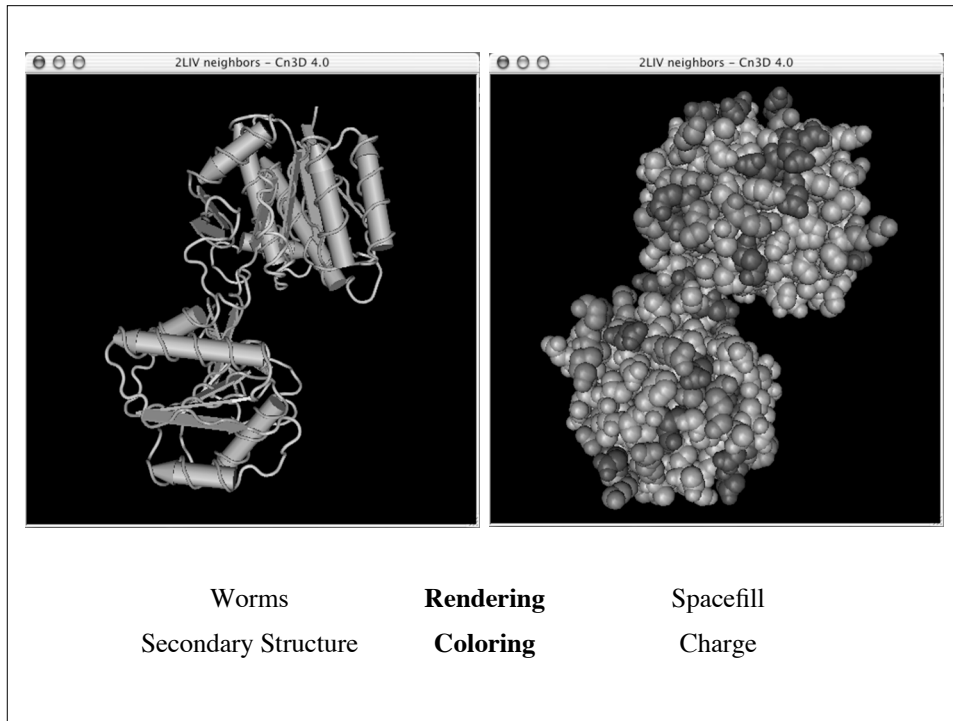
Blue – mismatches



2LIV neighbors - Sequence/Alignment Viewer

```

2LIV | D I K V A V V G A M S G P V A Q Y G D Q E F T G A E Q A V A D I N A K G G I R G N K L O I A K Y D D A C D P K Q A V A V A N K I V N D G I K Y V I G H L C S S S T Q P A S D I Y E D E G I L M I T P A
2LBP | D D I K V A V V G A M S G P I A Q W G I M E F N G A E Q A I K D I N A K G G I R G D K L V G V E Y D D A C D P K Q A V A V A N K I V N D G I K Y V I G H L C S S S T Q P A S D I Y E D E G I L M I S P G
    
```



## SWISS-MODEL

- Automated comparative protein modelling server
- Web front-end at <http://www.expasy.org/swissmod>  
Results returned by E-mail

BLAST search to find similarities in PDB *by sequence*

Select templates with sequence identity > 25% and  
projected model size > 20 amino acids

Generate models

Do energy minimization

Generate PDB file for new protein model






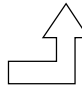
```

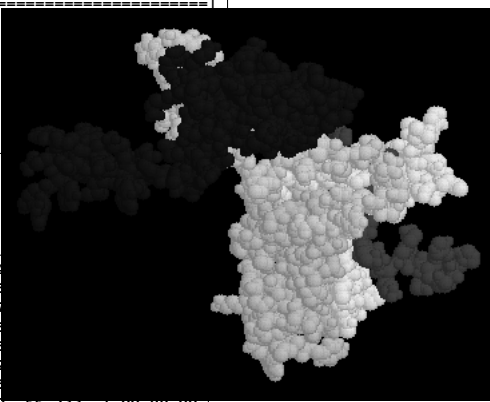
21DJH.pdb: 42.77 % identity
21DJG.pdb: 42.77 % identity
11DJG.pdb: 42.22 % identity
11QAS.pdb: 44.17 % identity
11QAT.pdb: 43.52 % identity
21QAT.pdb: 43.52 % identity
21QAS.pdb: 43.52 % identity

Target:
-----
21DJH.pdb
21DJG.pdb
11DJG.pdb
11QAS.pdb
11QAT.pdb
21QAT.pdb
21QAS.pdb
    
```



ATOM	1	H1	SER	1	24.219	22.954			
ATOM	2	H2	SER	1	24.770	21.435			
ATOM	3	N	SER	1	24.355	22.187			
ATOM	4	H3	SER	1	23.466	21.925			
ATOM	5	CA	SER	1	25.266	22.675			
ATOM	6	CB	SER	1	24.826	24.072			
ATOM	7	OG	SER	1	24.857	25.006			
ATOM	8	HG	SER	1	24.717	25.929	-55.233	1.00	99.00
ATOM	9	C	SER	1	25.471	21.750	-53.751	1.00	25.00
ATOM	10	O	SER	1	25.923	22.169	-52.684	1.00	25.00
ATOM	11	N	LYS	2	25.227	20.460	-53.972	1.00	25.00
ATOM	12	H	LYS	2	24.961	20.142	-54.878	1.00	99.00
ATOM	13	CA	LYS	2	25.366	19.408	-52.943	1.00	25.00
ATOM	14	CB	LYS	2	24.003	18.772	-52.622	1.00	25.00



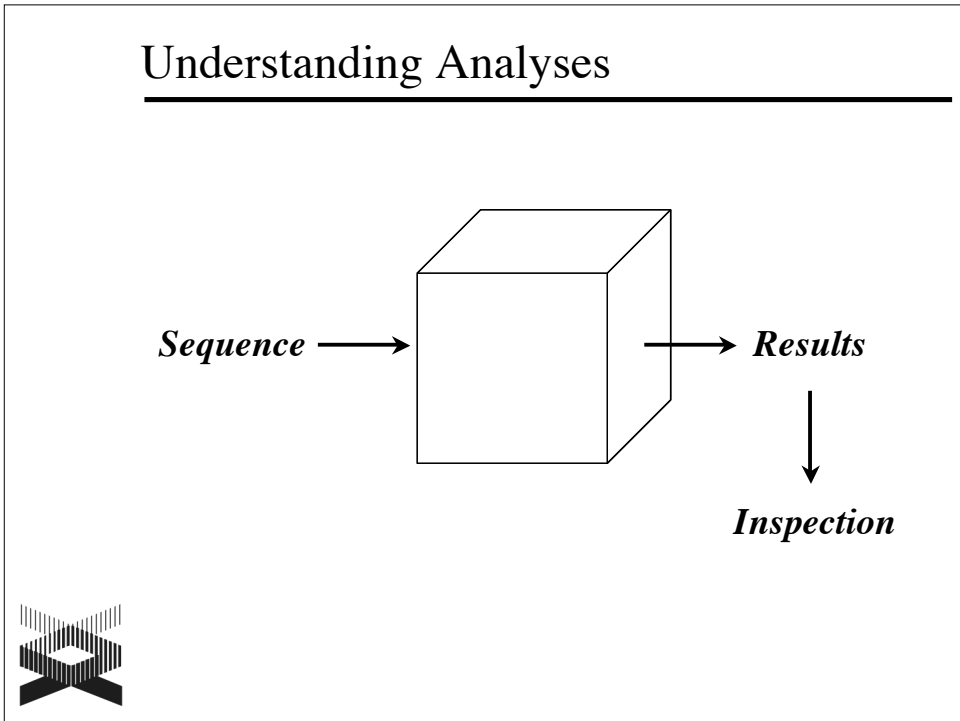
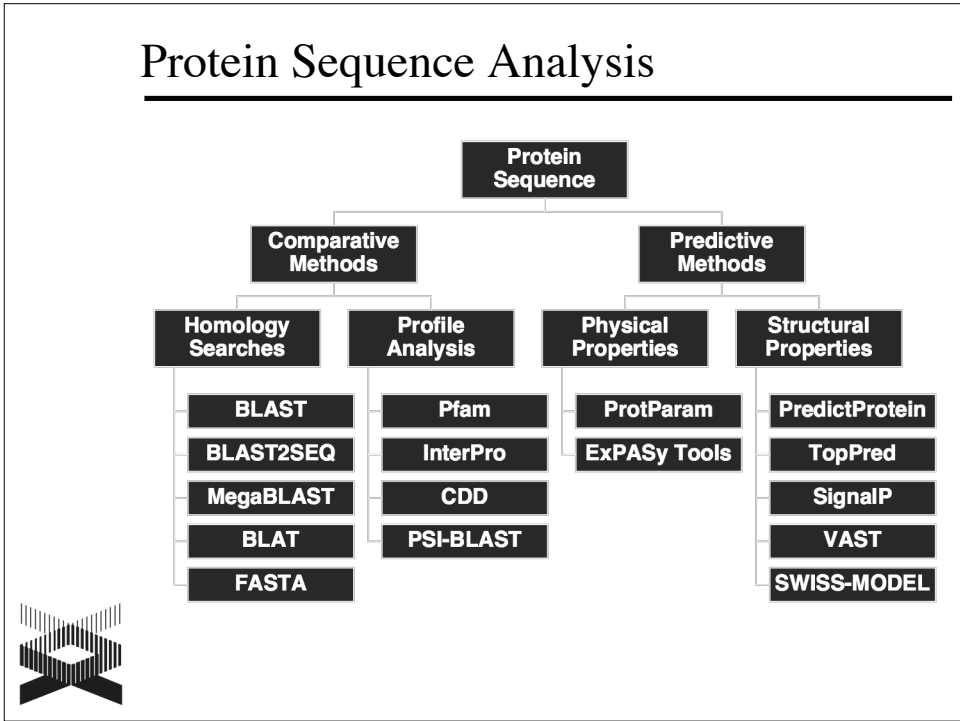


## Structural Modeling Software

---

- Modeller  
[http://www.salilab.org/modeller/about\\_modeller.html](http://www.salilab.org/modeller/about_modeller.html)
- DeepView  
<http://us.expasy.org/spdbv/>
- WHAT IF  
<http://swift.cmbi.kun.nl/whatif/>





## A User's Guide to the Human Genome II

---

*[http://www.nature.com/  
ng/supplements/](http://www.nature.com/ng/supplements/)*

***Commentary:  
Keeping Biology  
in Mind***

