

***Current Topics in Genome Analysis  
Spring 2005***

***Week 5  
Biological Sequence Analysis II***

*Andy Baxevanis, Ph.D.*



**Overview**

---

- Week 4: Comparative methods and concepts
  - Similarity vs. Homology
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

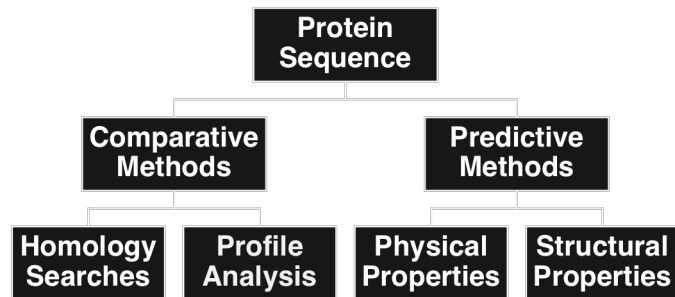


## Protein Conformation

- 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” *ProfileScan*  
*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  
 GVEICIATPG  
 GVDILIGTGG  
 RPHIIVATPG  
 KPHIIVATPG  
 KVQLIATPG  
 RPDIVIVATPG  
 APHIIIVGTGG  
 APHIIIVGTGG  
 GCHVVIVATPG  
 NQDIVVATGG

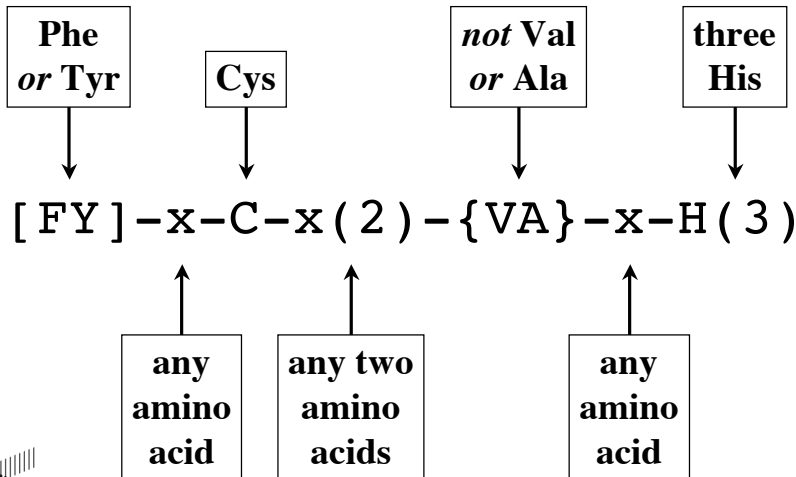
- Which residues are seen at each position?
- What is the frequency of observed residues?
- Which positions are conserved?
- Where can gaps be introduced?

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	0	10	0	0	10	10	0	0	0	0	0	1	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	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	20	30	-10	-10	30	150	20	-60	-30	10
P	21	6	7	6	6	11	10	11	0	6	16	11	11	89	17	17	24	22	9	-50	-48	12
G	10	60	20	70	50	150	-20	-30	-10	-50	-30	40	40	30	20	-30	60	40	20	-100	-70	30



## Patterns



## ProfileScan

- Search sequence against a collection of profiles and patterns
- Databases available
  - PROSITE profiles
  - PROSITE patterns
  - PfamA
  - PfamB
  - InterPro families
  - HAMAP profiles (microbial)
  - TIGRfam protein families
- <http://hits.isb-sib.ch/cgi-bin/PFSCAN>



**Motif Scan Results** user: anonymous  
log in

**Query Protein** temporarily stored here.

**Database of motifs** PROSITE patterns, PROSITE profiles, Pfam HMMs (local models).

Reference: Falquet L, Pagni M, Bucher P, Hulo N, Sigrist CJ, Hofmann K & Bairoch A. (2002) The PROSITE database, its status in 2002. *Nucleic Acids Res.* **30**:235-238

searching PROSITE patterns  
 searching PROSITE profiles  
 searching Pfam HMMs (local models)  
 postprocessing

**Summary**

**Original output** pat, prf, pfam\_fs.

**Matches map** (features from query are above the ruler, matches of the motif scan are below the ruler)

20 40 60 80 100 120 140 160 180 200 220 240 260 280 300 320 340 360 380 400 420 440 460 480 500

pfam\_fs:p450 [1]

Legends: 1, pat:CYTOCHROME\_P450 [1].

List of matches		FT	MYHIT	449	458	pat:CYTOCHROME_P450 [1]
		FT	MYHIT	41	506	pfam_fs:p450 [1]

**Match details**

match detail	match score	motif information
heme_iron		pat:CYTOCHROME_P450 Cytochrome P450 cysteine heme-iron

Status: !

**Status: !**

pos.: 41-506  
 raw-score = 450.5  
 N-score = 147.344  
 E-value = 9.6e-141

pfam\_fs:p450  
 Cytochrome P450  
 [ entry ]

Question or comment about this page.

NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis II

pfam\_fs.p450 - Netscape

http://myhits.isb-sib.ch/cgi-bin/view\_mot\_entry?name=pfam\_fs.p450

myhits Entry pfam\_fs:p450 user: anonymous log in

HMMER2.0 [2.3.1] [ documentation at Sanger Institute ]

NAME p450  
 ACC PF00067.10  
 DESC Cytochrome P450  
 LENG 499  
 ALPH Amino  
 RF no  
 CS no  
 NAP yes  
 COM hmmbuild -f -F HMM fs.ann SEED.ann  
 COM hmmscalibrate --seed 0 HMM fs.ann  
 NSEQ 52  
 DATE Sat Oct 2 17:06:11 2004  
 CKSUM 4876  
 GA 13.00 13.00  
 TC 13.20 13.00  
 NC 4.48 4.48  
 XT -8455 -4 -1000 -1000 -8455 -4 -8455 -4  
 NULT -4 -8455  
 NULE 595 -1558 85 338 -294 453 -1158 197 249 902 -1085 -142 -21 -313 45  
 EVD -10.980991 0.710359  
 HMM

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R
m->m	m->i	m->d	i->m	i->i	d->m	d->d	b->m	m->>e							
	-40	*	-5208												
1	-5254	-4976	-7829	-7450	-3505	-6831	-6481	-2766	-7126	-434	251	-7039	4109	-6156	-6707
-	-149	-500	233	43	-381	399	106	-626	210	-466	-720	275	394	45	96
2	-3	-10334	-11376	-894	-1115	-701	-1378	-1040	-9960						
-	-341	-4053	-2487	-1927	-4384	1004	941	-4116	-92	-1591	-3147	-2225	3021	416	809
-	-149	-500	233	43	-381	399	106	-626	210	-466	-720	275	394	45	96
3	-3	-10334	-11376	-894	-1115	-701	-1378	-10000	-9959						
-	-2696	127	-1233	-1340	-350	2514	-2929	-2249	-3303	-167	-1918	-3424	2110	-3085	-3415
-	-149	-500	233	43	-381	399	106	-626	210	-466	-720	275	394	45	96
4	-3	-10334	-11376	-894	-1115	-701	-1378	-10000	-9957						
-	-1240	-3459	-2690	-105	398	736	-2324	-3227	-1118	-294	-2601	-2392	3242	-185	-2484
-	-149	-500	233	43	-381	399	106	-626	210	-466	-720	275	394	45	96
-	-283	-10334	-2506	-894	-1115	-701	-1378	-10000	-9956						

Pfam: p450 - Netscape

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

Protein families database of alignments and HMMs

Wellcome Trust Sanger Institute

p450 Home Search by Browse by ftp IPfam Help

Accession number: PF00067

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

Cytochrome P450s are involved in the oxidative degradation of various compounds. Particularly well known for their role in the degradation of environmental toxins and mutagens. Structure is mostly alpha, and binds a heme cofactor.

**NEW!** 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 P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class MEDLINE:7678494.

**QuickGO**

**PROCESS :** electron transport (GO:0006118)

For additional annotation, see the PROSITE document PDOC00081 (ExPasy|SRS-UK|SRS-USA)

**Alignment** **Domain organisation**

Seed (52) Full (3878)

Format Coloured alignment

Get alignment View HMM logo

View 13 representative architectures  
 View architectures for 3878 proteins

Zoom 0.5 pixels/aa.

View Graphic

Further alignment options here

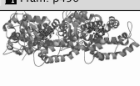




Pfam: p450 - Netscape

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

Pfam: p450



Cytochrome P450s are involved in the oxidative degradation of various compounds. Particularly well known for their role in the degradation of environmental toxins and mutagens. Structure is mostly alpha, and binds a heme cofactor.

**NEW!** 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 P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class MEDLINE:7678494.

**QuickGO**

**PROCESS :** electron transport (GO:0006118)

**Figure 1: 1mpw**  
**Oxidoreductase**  
 Molecular recognition in (+)- $\alpha$ -pinene oxidation by cytochrome p450cam

**Key:**

Domain	Chain	Start Residue	End Residue
p450	A	13	399
p450	B	13	399

The SwissProt/PDB mapping was provided by MSD

1akd [Display pdb](#)

For additional annotation, see the PROSITE document PDOC00081 [ExPasy|SRS-UK|SRS-USA]

**Alignment**

Seed (52) Full (3878)

Format: Coloured alignment

[Get alignment](#) [View HMM logo](#)

Further alignment options here  
 Help relating to Pfam alignments here

**Domain organisation**

View 13 representative architectures  
 View architectures for 3878 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 (52) Full (3878)

[Download tree](#) [ATV Applet](#)

The trees were generated using Quicktree  
 To find out more about ATV phylogenetic tree-viewer click here

Pfam: Distinct architecture for all p450 domain proteins - Netscape

http://www.sanger.ac.uk/cgi-bin/Pfam/getallproteins.p?name=p450&acc=PF00067&verbose=true&type=full&domain\_view=arc

Pfam: Distinct architecture for all p450 domain proteins

Protein families database of alignments and HMMs

Wellcome Trust Sanger Institute


Home Search by Browse by ftp IPfam Help

**Distinct architecture for all p450 domain proteins**

This family may contain **overlapping domains**, to change the graphical view click here

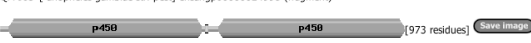
**3460 proteins with p450 architecture** [View](#)

C4GF\_DROME [ drosophila melanogaster (fruit fly)] cytochrome p450 4g15 (ec 1.14.-.-) (cyp1g15)




**192 proteins with p450, p450 architecture** [View](#)

Q7PJ63 [ anopheles gambiae str. pest] ensangp0000024998 (fragment)




**14 proteins with p450, Flavodoxin\_1, FAD\_binding\_1, NAD\_binding\_1 architecture** [View](#)

Q9HGE0 [ gibberella moniliformis] fum6p




**3 proteins with p450, p450, p450 architecture** [View](#)

Q6U7Q8 [ cryptococcus neoformans var. grubii h99] cytochrome p450 lanosterol 14a-demethylase (ec 1.14.13.70)



**2 proteins with Peptidase\_C48, p450 architecture** [View](#)

Q94HM5 [ oryza sativa (rice)] putative cytochrome p-450 like protein



Accession number: PF00067

### Cytochrome P450

Cytochrome P450s are involved in the oxidative degradation of various compounds. Particularly well known for their role in the degradation of environmental toxins and mutagens. Structure is mostly alpha, and binds a heme cofactor.

**NEW!** 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 P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class MEDLINE:7678494.

**QuickGO**

**PROCESS :** electron transport (GO:0006118)

**Figure 1: 1mpw**  
**Oxidoreductase**  
 Molecular recognition in (+)- $\alpha$ -pinene oxidation by cytochrome p450cam

**Key:**

Domain	Chain	Start Residue	End Residue
p450	A	13	399
p450	B	13	399

The SwissProt/PDB mapping was provided by MSD

1akd Display pdb

For additional annotation, see the PROSITE document PDOC00081 [ExPasy|SRS-UK|SRS-USA]

**Alignment**

Seed (52) Full (3878)

Format: Coloured alignment

Get alignment View HMM logo

Further alignment options here

**Domain organisation**

View 13 representative architectures

View architectures for 3878 proteins

Zoom 0.5 pixels/aa

View Graphic

InterPro: Cytochrome P450

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 Cytochrome P450

IPR001128 Cytochrome\_P450

Matches: 4047 proteins. View matches: Please be aware that match views for entries matching more than 1000 proteins may be slow.

Overview: sorted by AC, sorted by name, of known structure, grouped by taxonomy

Detailed: sorted by AC, sorted by name, of known structure

Table: For all matching proteins, of known structure, Architectures

**Name** [?]: Cytochrome P450

**Signatures** [?]: PF00067:p450 (3834 proteins), PR00385:P450 (2932 proteins), PS00086:CYTOCHROME\_P450 (3175 proteins), SSF48264:Cytochrome\_P450 (4012 proteins)

**Type** [?]: Family

**Dates** [?]: 1999-10-08 17:07:25.0 (created), 2000-02-17 17:11:42.0 (modified)

**Children** [?]: IPR002397: B-class P450, IPR002399: Mitochondrial P450, IPR002401: E-class P450, group I, IPR002402: E-class P450, group II, IPR002403: E-class P450, group IV

**Process** [?]: electron transport (GO:0006118)

**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** [?]: SCOP a\_104.1.1, CATH 1.10.630.10, PDB/MSD - click here

**Database links** [?]: PANDIT PF00067

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

NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis II

**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

**Children** (tree):  
 IPR002397: B-class P450  
 IPR002399: Mitochondrial P450  
 IPR002401: E-class P450, group I  
 IPR002402: E-class P450, group II  
 IPR002403: E-class P450, group IV

**Taxonomy**

4	Saccharomyces cerevisiae	Unclassified
384	Fungi	Virus
77	Caenorhabditis elegans	Archaea
127	Nematoda	Bacteria
2165	Metazoa	Cyanobacteria
112	Fruit Fly	Synechocystis PCC 6803
599	Arthropoda	Rice spp.
1041	Chordata	Arabidopsis thaliana
154	Mouse	Green Plants
168	Human	Plastid Group
3439	Eukaryota	Other Eukaryotes

**Examples**

- Q64459 CP3B\_MOUSE
- P12938 CPD5\_RAT
- P33267 CFZ2\_MOUSE
- P30612 CP5P\_CANTR
- P21595 CP56\_YEAST
- P26911 CPXH\_STRGR

**More proteins...**

- IPR001128 Cytochrome P450
- IPR002397 B-class P450
- IPR002401 E-class P450, group I
- IPR002974 P450, CYP52
- IPR008069 E-class P450, CYP2D
- IPR008072 E-class P450, CYP3A

**Center**  
 Inner circles  
 Outer circles

**Tree root**  
 Tree nodes  
 Representative model organisms

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

## 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
- <http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>



NCBI CD-Search - Netscape  
 http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi

NCBI Conserved Domain Search

New Search | PubMed | Nucleotide | Protein | Structure | CDD | Taxonomy | Help?

RPS-BLAST 2.2.9 [May-01-2004]  
 Query= local sequence: lcl|tmpseq\_0 Human DCC precursor  
 (1447 letters)

Database: oasis\_smart.v2.02

Click on boxes for multiple alignments

1 210 410 610 810 1010 1210 1447

IG IG IG IG FN3 FN3 FN3 FN3 FN3 FN3

1010 1010 1010 1010

Show | Domain Relatives

.. This CD alignment includes 3D structure. To display structure, download **Cn3D!**

PSSMs producing significant alignments:

	Score	E value
gnl CDD 25322 smart00409, IG, Immunoglobulin;	67.1	3e-13
gnl CDD 25322 smart00409, IG, Immunoglobulin;	61.7	1e-11
gnl CDD 25322 smart00409, IG, Immunoglobulin;	59.0	8e-11
gnl CDD 25322 smart00409, IG, Immunoglobulin;	42.8	6e-06
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	63.5	3e-12
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	59.3	6e-11
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	48.1	1e-07
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	45.4	8e-07
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.9	3e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.5	5e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	55.3	8e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	48.8	8e-08
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	47.2	2e-07
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	37.2	2e-04

NCBI CD-Search - Netscape  
 http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi

NCBI CD-Search

- gnl|CDD|25322, smart00409, IG, Immunoglobulin;  
 CD-Length = 86 residues, 98.8% aligned  
 Score = 67.1 bits (163), Expect = 3e-13  
 Query: 337 PSNLYAYESMDIEFECTVSCPKPPTVNMKN-GDVVIPSDFQIVGGSN---LRILGVVK 392  
 Sbjct: 1 PPSVTVKEGESVTLSCASGNPPPEVTWYKGGKLLAYSGRFSVSRSGNSTLTISNVTP 60
- gnl|CDD|25322, smart00409, IG, Immunoglobulin;  
 CD-Length = 86 residues, 91.9% aligned  
 Score = 61.7 bits (149), Expect = 1e-11  
 Query: 147 ESVTFMGDTVLLKCEVIGEPMPETHWQKNQDLTPIPGDSRVVLPSPG---ALQISRLQ 203  
 Sbjct: 2 PPSVTVKEGESVTLSCASGNPPPEVTWYK--OGGKLLAYSGRFSVSRSGNSTLTISNVTP 59
- gnl|CDD|25322, smart00409, IG, Immunoglobulin;  
 CD-Length = 86 residues, 100.0% aligned  
 Score = 59.0 bits (142), Expect = 8e-11  
 Query: 246 PSNVVAIEGKDAVLECCVGYPPSPFTWLRGEEVIQLRSKYY---SLGGSNLLISWTD 302  
 Sbjct: 1 PPSVTVKEGESVTLSCASGNPPPEVTWYKGGKLLAYSGRFSVSRSGNSTLTISNVTP 60
- gnl|CDD|25322, smart00409, IG, Immunoglobulin;  
 CD-Length = 86 residues, 98.8% aligned  
 Score = 42.8 bits (100), Expect = 6e-06  
 Query: 303 DDSGMYTCVVYKKNISASAEITVL 328  
 Sbjct: 61 EDSGTYTCAATNSSGSASSGTTLTVL 86

NCBI CD-Search - Netscape

http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi

NCBI CD-Search

NCBI Conserved Domain Search

New Search | PubMed | Nucleotide | Protein | Structure | CDD | Taxonomy | Help?

RPS-BLAST 2.2.9 [May-01-2004]  
 Query= local sequence: lcl|tmpseq\_0 Human DCC precursor  
 (1447 letters)

Database: oasis\_smart.v2.02

Click on boxes for multiple alignments

1 210 410 610 810 1010 1210 1447

IG IG IG IG FN3 FN3 FN3 FN3 FN3 FN3

Show | Domain Relatives

.. This CD alignment includes 3D structure. To display structure, download **Cn3D!**

PSSMs producing significant alignments:

	Score	E value
gnl CDD 25322 smart00409, IG, Immunoglobulin;	67.1	3e-13
gnl CDD 25322 smart00409, IG, Immunoglobulin;	61.7	1e-11
gnl CDD 25322 smart00409, IG, Immunoglobulin;	59.0	8e-11
gnl CDD 25322 smart00409, IG, Immunoglobulin;	42.8	6e-06
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	63.5	3e-12
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	59.3	6e-11
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	48.1	1e-07
gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	45.4	8e-07
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.9	3e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.5	5e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	55.3	8e-10
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	48.8	8e-08
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	47.2	2e-07
gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	37.2	2e-04

NCBI CDD smart00409 - Netscape

http://www.ncbi.nlm.nih.gov/Structure/cdd/cddsvr.cgi?uid=smart00409&version=v2.02

NCBI CDD smart00409

Conserved Domains

HOME | SEARCH | SITE MAP | Entrez | CDD | Structure | Protein | Help

smart00409.10 Immunoglobulin; IG

Source: Smart  
 Taxonomy: root  
 Proteins: smart00409 related  
 Related CD: 7 links

Statistics:

PSSM-Id: 25322  
 Aligned: 472 rows  
 PSSM: 86 columns  
 Status: Alignment from source  
 Created: 12-Dec-2003  
 Updated: 12-Dec-2003

Structure:

Show Structure

Program: Cn3D

Drawing: Virtual Bonds

(download Cn3D)

This domain model appears to be related to other CDs:

[mouse over cd tag to display the number of PSSM pairs and cd name]

Show Alignment

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

Type Selection: the most diverse members

consensus	1	PPSVTVKEGESVTLSCASG.[1].PPPEVTW	YK.[2].GKLL.[6].SVSR.[3].NSTLTISNVTPE.[2].63
1MCP_H	7	SGGGLVQPQGSRLRSCATSG.[3].SDFYMEW	VR.[6].LEWI.[22].IVSR.[5].ILYLQMNALRAE.[2].93
1GYA	6	ALETWALGQDINLDIPSFQ.[3].DIDDIKW	EK.[3].KKKI.[16].KLFK
1ZXO	9	PKKLAVEPKGSLEVNCSTTC.[1].QPEVGGI	ET.[1].LNKI
g1_399208	25	QRPLLIIVANRTATLVYNYTY.[4].KEFRASL	HK.[4].AVEV.[20].RGIH.[3].KVIFNLWNMSAS.[2].106

NCBI CDD smart00409

Program: Cn3D  
 Drawing: Virtual Bonds  
 (download Cn3D)

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

Type Selection: the most diverse members

consensus	1 PPSVTVKEGESVTLSCAASG.[1].PPPEVTW	YK.[2].GKLL.[6].SVSR.[3].NSTLTISNVTPE.[2].63
IMCP_H	7 SGGGLVQPGGSLRLSCATSG.[3].SDPYMEW	VR.[6].LEWI.[22].IVSR.[5].ILYLQMNALRAE.[2].93
1GYA	6 ALETWALGQDINLDIPSFQ.[3].DIDDIRK	EK.[3].KKKI.[16].KLFFK NGTLKIKHLKTD.[2].78
1ZXQ	9 PKKLAVPEKGSLEVNCSSTC.[1].QPEVGGI	ET.[1].LNKI LLDE.[3].WKHYLVSNISHD 62
gi 399208	25 QRPLLVANRATILVCMYYI.[4].KEFRASL	HK.[4].AVEV.[20].RGLH.[3].KVIFMLWMSAS.[2].106
gi 461714	216 SNTFYAREGDQVEFSPLSF.[2].ENLVGEL	RW.[9].LWIS.[19].QMKE.[2].PLRFTIPOVLSR.[2].298
gi 6166597	64 PGGTVKVGEDITPIAKVKA.[6].PTIKWFK.[1].KW.[6].AGKH.[7].ERHS.[3].TFEMQIIRAKDN.[2].137	
gi 729801	435 QRTQYGLVGDTRARIECFASS.[3].ARHVSWT	FN.[1].QEIS.[7].SILV.[7].KSTLIIRDSQAY.[2].503
gi 124310	243 LRTISASLGSRLTIPCKVEL.[4].PLTMTLW	WT.[1].NDTH.[18].SENN.[4].EVPILFDPVTR.[3].321
gi 1709202	281 REGETMSLGCRRVITPEIKH	FOPEIRW YR.[1].GVPL.[6].QTLW.[3].RATLTFSHLNKE.[2].341

Citing CDD: Marchler-Bauer A, Anderson JB, Cherukuri PF, DeWeese-Scott C, Geer LY, Gwadz M, He S, Hurwitz DJ, Jackson JD, Ke Z, Lanczycki CJ, Liebert CA, Liu C, Lu F, Marchler GH, Mulloikandov M, Shoemaker BA, Simonyan V, Song JS, Thissen PA, Yamashita RA, Yin JJ, Zhang D, Bryant SH (2005), "CDD: a Conserved Domain Database for protein classification.", *Nucleic Acids Res.* 33: D192-6

NCBI Conserved Domain Search

RPS-BLAST 2.2.9 [May-01-2004]  
 Query= local sequence: lcl|tmpseq\_0 Human DCC precursor  
 (1447 letters)

Database: oasis\_smart.v2.02

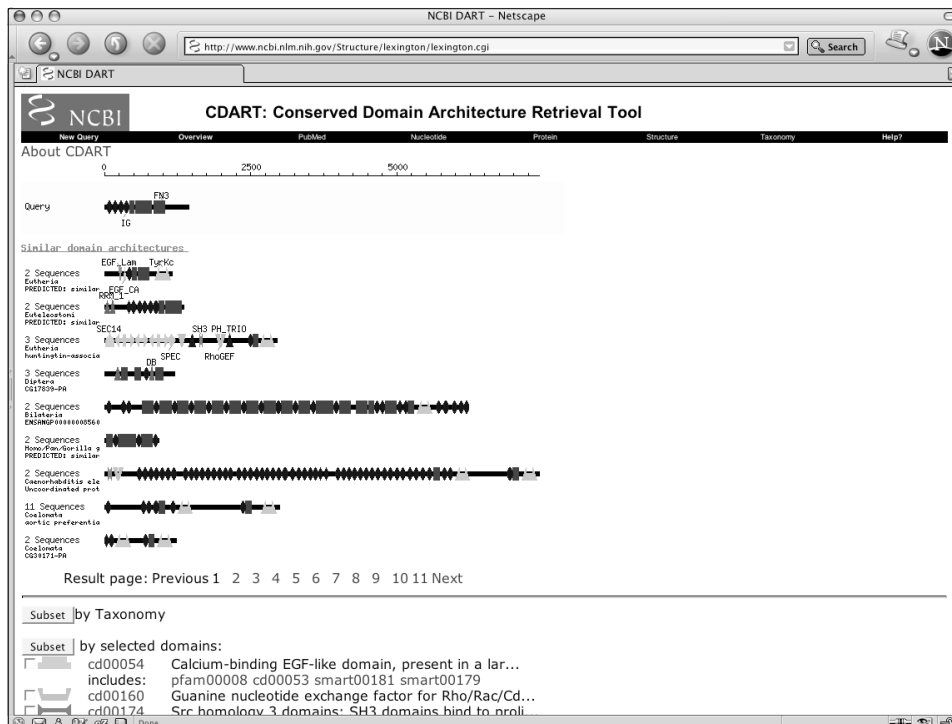
Click on boxes for multiple alignments

IG IG IG IG FN3 FN3 FN3 FN3 FN3 FN3

Show Domain Relatives

.. This CD alignment includes 3D structure. To display structure, download Cn3D!

PSSMs producing significant alignments:	Score E (bits) value
● gnl CDD 25322 smart00409, IG, Immunoglobulin;	67.1 3e-13
● gnl CDD 25322 smart00409, IG, Immunoglobulin;	61.7 1e-11
● gnl CDD 25322 smart00409, IG, Immunoglobulin;	59.0 8e-11
● gnl CDD 25322 smart00409, IG, Immunoglobulin;	42.8 6e-06
● gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	63.5 3e-12
● gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	59.3 6e-11
● gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	48.1 1e-07
● gnl CDD 365 smart00408, IGc2, Immunoglobulin C-2 Type;	45.4 8e-07
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.9 3e-10
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	56.5 5e-10
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	55.3 8e-10
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	48.8 8e-08
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	47.2 2e-07
● gnl CDD 25286 smart00060, FN3, Fibronectin type 3 domain; One of three types...	37.2 2e-04



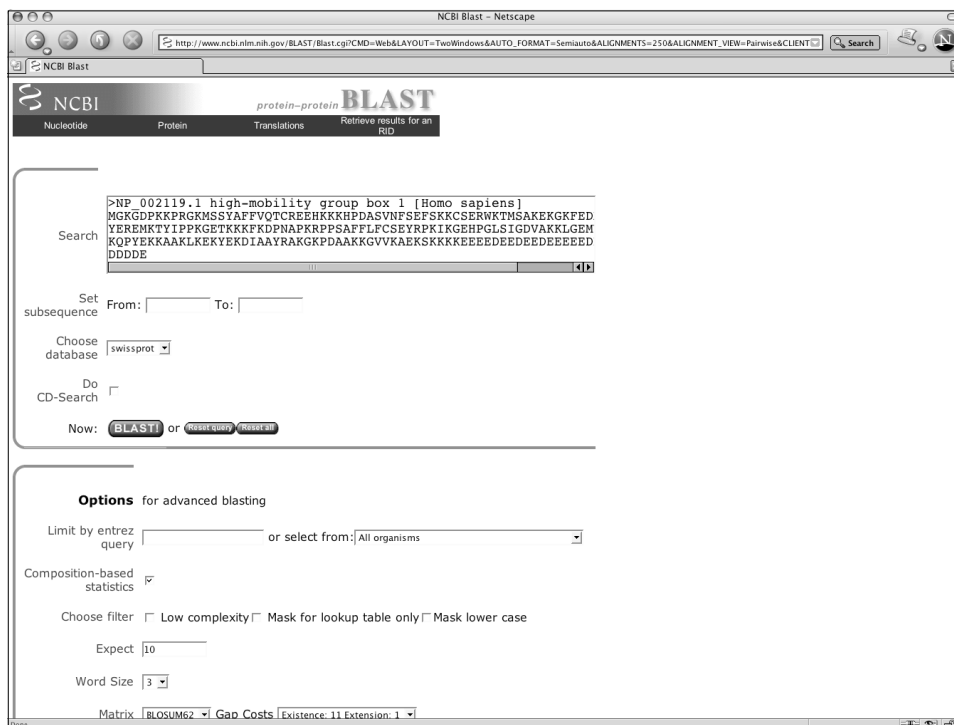
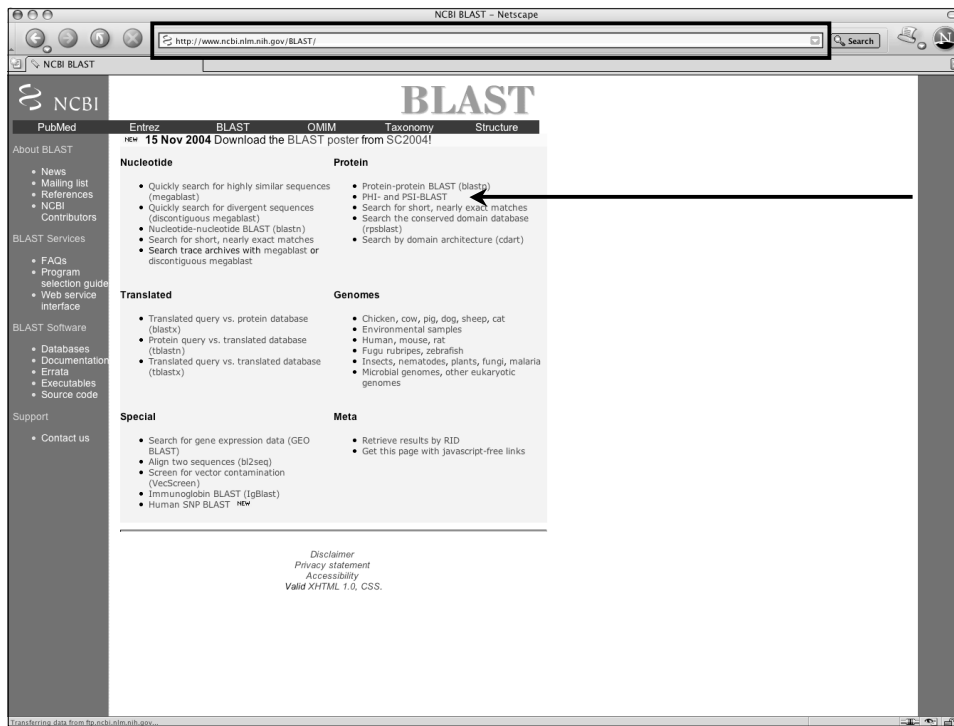
## PSI-BLAST

- 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





NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis II



NHGRI Current Topics in Genome Analysis 2005  
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=...

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 ]

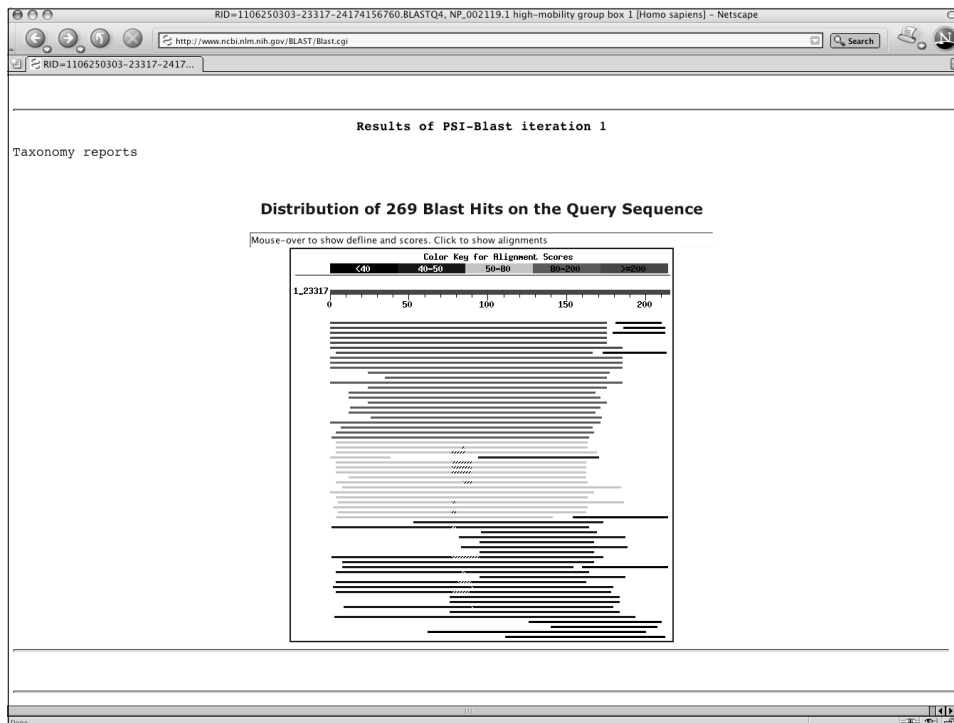
Expect value range: [ ] [ ]

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

Autoformat [ Semi-auto ]

**BLAST!**

Get WWW URL with preset values?



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis II

RID=1106250303-23317-24174156760.BLASTQ4, NP\_002119.1 high-mobility group box 1 [Homo sapiens] - Netscape

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

RID=1106250303-23317-2417...

**Legend:**

- ✖ - 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

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

	Score	E Value
gi 123371 sp P12682 HMG1_PIG High mobility group protein 1 (HMG-1)	238	7e-63
gi 52783747 sp P63158 HMG1_MOUSE High mobility group protein 1 (...)	238	1e-62
gi 123367 sp P10103 HMG1_BOVIN High mobility group protein 1 (HM...	238	1e-62
gi 123369 sp P09429 HMG1_HUMAN High mobility group protein 1 (HM...	238	1e-62
gi 20138433 sp Q9UGV6 HM1X_HUMAN High mobility group protein 1-1...	229	5e-60
gi 123373 sp P26584 HMG2_CHICK High mobility group protein 2 (HM...	202	6e-52
gi 123382 sp P07746 HMGT_ONCMY High mobility group-T protein (HM...	201	1e-51
gi 1708260 sp P52925 HMG2_RAT High mobility group protein 2 (HMG-2)	194	2e-49
gi 123374 sp P26583 HMG2_HUMAN High mobility group protein 2 (HM...	193	2e-49
gi 1708259 sp P30681 HMG2_MOUSE High mobility group protein 2 (H...	193	2e-49
gi 13878931 sp P23497 SP10_HUMAN Nuclear autoantigen Sp-100 (Spe...	191	1e-48
gi 123368 sp P07156 HMG1_CRIGR High mobility group protein 1 (HM...	188	9e-48

RID=1106250303-23317-24174156760.BLASTQ4, NP\_002119.1 high-mobility group box 1 [Homo sapiens] - Netscape

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

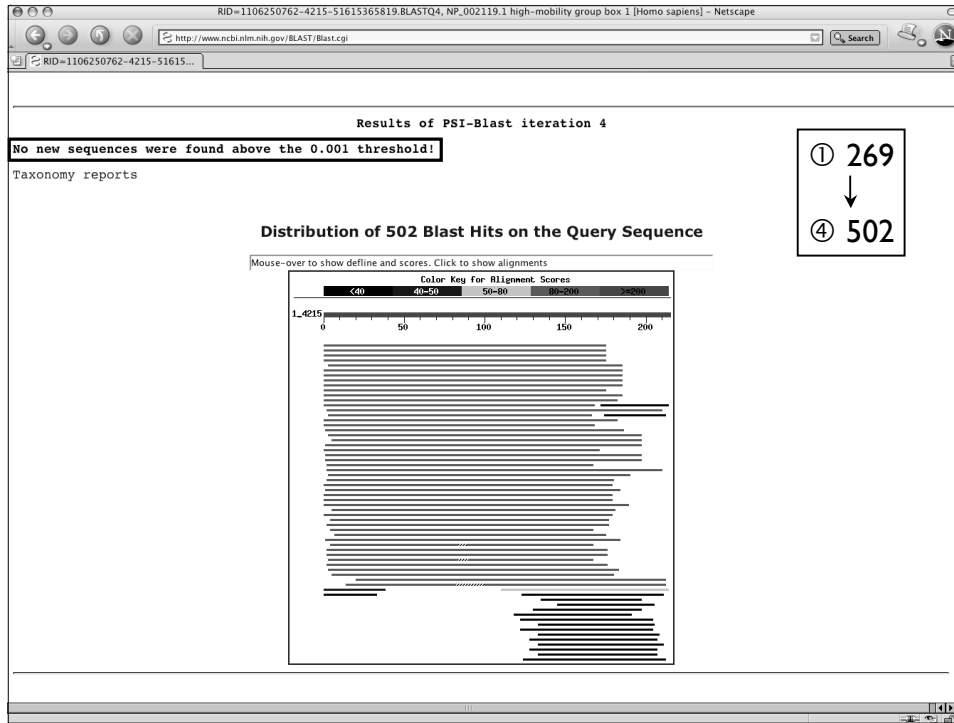
RID=1106250303-23317-2417...

gi 2495255 sp Q03435 NH10_YEAST Non-histone protein 10 (High mob...	42	8e-04
gi 6175054 sp P36389 SRY_HORSE Sex-determining region Y protein ...	42	8e-04
gi 22654148 sp Q91ZW1 TFAM_RAT Transcription factor A, mitochond...	42	8e-04
gi 6175076 sp Q04888 SX10_MOUSE Transcription factor SOX-10 (SOX...	42	8e-04
gi 6094380 sp O55170 SX10_RAT Transcription factor SOX-10	42	9e-04

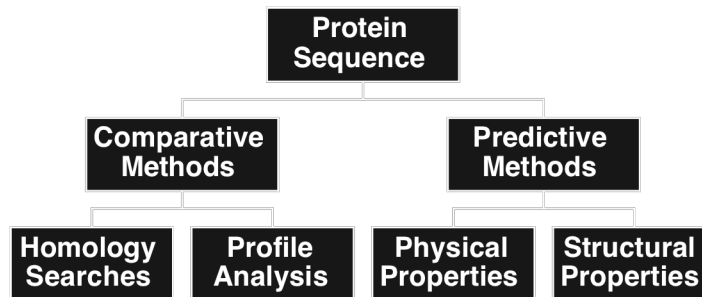
Run PSI-Blast iteration 2

Sequences with E-value WORSE than threshold

gi 6831689 sp O95416 SX14_HUMAN Transcription factor SOX-14 >gi ...	42	0.001
gi 2506521 sp P48434 SOX9_CHICK Transcription factor SOX-9	42	0.001
gi 24638225 sp Q9W7R6 SX14_CHICK Transcription factor SOX-14	42	0.001
gi 19862533 sp Q04892 SX14_MOUSE Transcription factor SOX-14	42	0.001
gi 1711465 sp P54231 SOX2_SHEEP Transcription factor SOX-2	42	0.001
gi 1351091 sp P48431 SOX2_HUMAN Transcription factor SOX-2	42	0.001
gi 3913481 sp Q24533 DICH_DROME SOX-domain protein dichaeete (Fis...	42	0.001
gi 12644266 sp P43267 SX15_MOUSE SOX-15 protein	42	0.001
gi 1723428 sp Q10241 CMB1_SCHPO Mismatch-binding protein cmb1	42	0.001
gi 6094324 sp P48432 SOX2_MOUSE Transcription factor SOX-2	42	0.001
gi 136654 sp P25977 UBF1_RAT Nucleolar transcription factor 1 (U...	42	0.001
gi 136652 sp P17480 UBF1_HUMAN Nucleolar transcription factor 1 ...	42	0.002
gi 729738 sp P40621 HMGL_WHEAT HMGL/2-like protein	41	0.002
gi 730136 sp P40632 NHP1_BABBO High mobility group protein homol...	41	0.002



## Protein Sequence Analysis



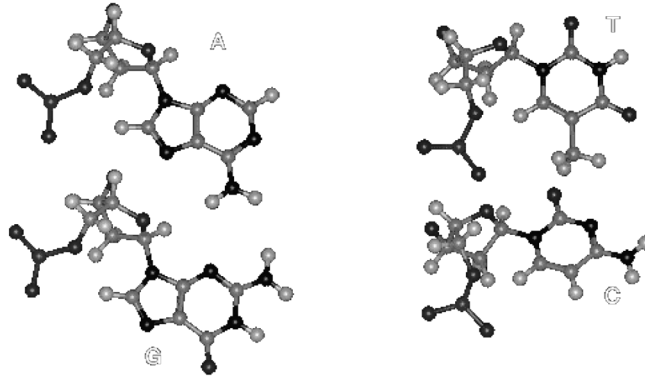
- *Common structure?*
- *Common function?*
- *Evolutionary relationship?*
- *Global or local similarity?*

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



## Information Landscape

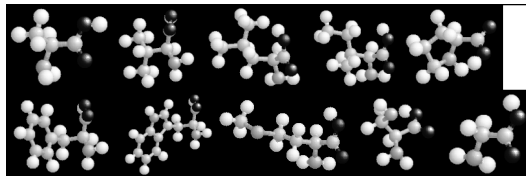
---



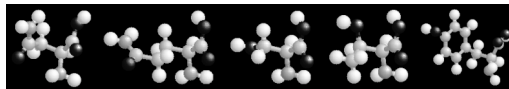
## Information Landscape

---

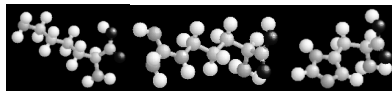
*Nonpolar*



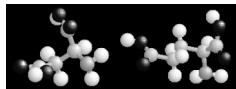
*Polar Neutral*



*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

MNGEADCPDLEMAAPKQDRWSQEDMLTLLCEMKNLPSNDSKFKTTESHMDWEKVAFKDFSGDMCKL  
 KWVEISNEVRKFRITLTELILDAQEHVKNPYKGGKLLKKHPDFPKKPLTPYFRFFMEKRAKYAKLHPEM...

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

HelixWheel/HelixDraw

*Display protein sequence as a helical wheel*



## Secondary Structure Prediction

---

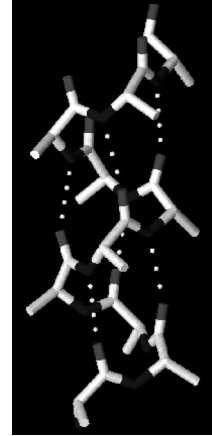
- Deduce the most likely position of alpha-helices and beta-strands
- Confirm structural or functional relationships when sequence similarity is weak
- Determine guidelines for rational selection of specific mutants for further laboratory study
- Basis for further structure-based studies



## Alpha-helix

---

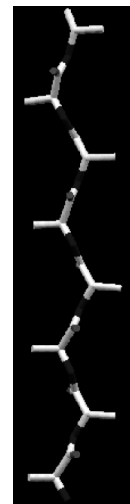
- Corkscrew
- Main chain forms backbone, side chains project out
- Hydrogen bonds between CO group at  $n$  and NH group at  $n+4$
- Helix-formers: Ala, Glu, Leu, Met
- Helix-breaker: Pro



## Beta-strand

---

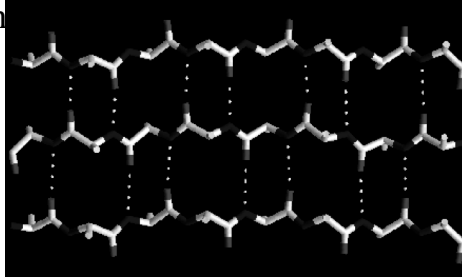
- Extended structure (“pleated”)
- Peptide bonds point in opposite directions
- Side chains point in opposite directions
- No hydrogen bonding *within* strand





## Beta-sheet

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



## Folding Classes



$\alpha$

*Cyt c*

Globins  
 Orthogonal  
 EF-hand  
 Up-Down  
 Cytochrome

$\beta$

*CD4*

Orthogonal  
 Super-barrel  
 Greek key  
 Sandwich  
 Jelly roll

$\alpha+\beta$

*Staph  
 nuclease*

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

$\alpha/\beta$

*Triose  
 phosphate  
 isomerase*

TIM barrel  
 Doubly-wound



## Neural Networks

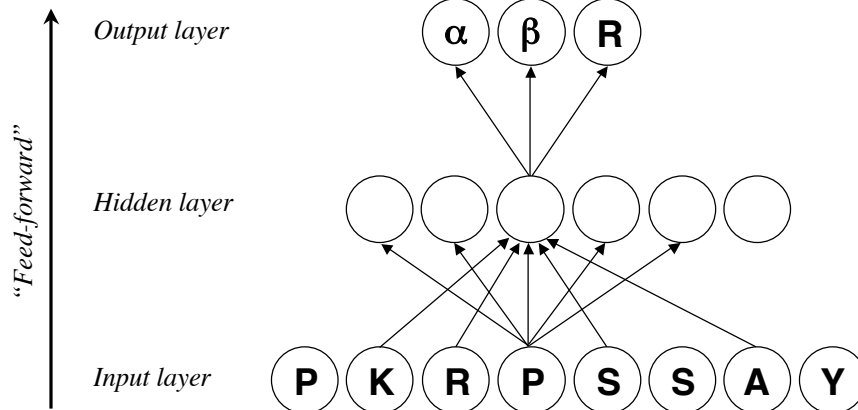
---

- 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



## Neural Networks

---



## nnpredict

---

- Neural network approach to making predictions  
(Kneller et al., 1990)
- Best-case accuracy > 65%
- Search engines
  - E-mail [nnpredict@celestes.ucsf.edu](mailto:nnpredict@celestes.ucsf.edu)
  - Web <http://www.cmpharm.ucsf.edu/~nomi/nnpredict.html>



## nnpredict Query

---

```
option: a/b
>flavodoxin - Anacystis nidulans
AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIY
DDLDSVNFQGGKVAIFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQPDLTKNRIKTWVSQLKSEFGL
```

↓ *α/β folding class*

Tertiary structure class: alpha/beta

```
Sequence:
AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVG
ELQSDWEGIYDDLDSVNFQGGKVAIFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYW
PIEGYDFNESKAVRNNQFVGLAIDEDNQPDLTKNRIKTWVSQLKSEFGL
```

```
Secondary structure prediction (H = helix, E = strand, - = no prediction):
---EEE-----EEEEHHHHH-----EEEH-----EEEE-----
-----HHHH-----EEEE-----H-----HHHHHHH-----E--E-
-E-----HH--E-----EHHHH-----
```



## 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%
- Search engines

<http://www.embl-heidelberg.de/predictprotein/>

<http://cubic.bioc.columbia.edu/predictprotein/>



```
>flavodoxin - Anacystis nidulans
AKIGLFYGTQTGVTTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIY
DDLDSVNFQGGKVA YFGAGDQVGYSDNFQDAMGILEEKISSLSGQT VGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQPDLTKNRIKTWVSQLKSEFGL
```



PROF results (normal)

```
.....1.....2.....3.....4.....5.....6.....7.....8.....9.....10
AA      AKIGLFYGTQTGVTTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIYDDLDSVNFQGGKVA YFGAGDQVGYSDNFQD
OBS_sec EEEEEEE HHHHHHHHHHHH EEEEE EEEEE HHHHHHHHHH EEEEE HHHH
Rel_sec 927899842676268888988873106842544223567000122553788720347666542167888886304568863788841466443311446

0.1,....11.1,....12.1,....13.1,....14.1,....15.1,....16.1,....
AA      AMGILEEKISSLSGQT VGYWPIEGYDFNESKAVRNNQFVGLAIDEDNQPDLTKNRIKTWVSQLKSEFGL
OBS_sec HHHHHHHHHH EEEEE EEEEE HHHHHHHHHHHH
PROF_sec 788899888740782542202456544533100178868876426664211178899999888754289
Rel_sec
```

- SWISS-PROT hits
- Multiple alignment
- PDB homologues

## 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)
- <http://www.cbs.dtu.dk/services/SignalP/>



SignalP 3.0 Server - new version

SignalP 3.0 server predicts the presence and location of signal peptide cleavage sites in amino acid sequences from different organisms: Gram-positive prokaryotes, Gram-negative prokaryotes, and eukaryotes. The method incorporates a prediction of cleavage sites and a signal peptide/non-signal peptide prediction based on a combination of several artificial neural networks and hidden Markov models.

View the version history of this server. All the previous versions are available on line, for comparison and reference.

Background Article abstracts Instructions Output format

**SUBMISSION**

Paste a single sequence or several sequences in FASTA format into the field below:

```
>P05019 Insulin-like growth factor IB precursor (IGF-IB)
MGKISSLPQLKCCDFLKVKHTMSSSHLPIALCLLFTTSSATAGPELGGAEIYVDALQVY
FFYFNKPTGYGSSRRAPQTGIVDECCFRSCDLRRLLEMYCAPLKPARSARSVRAQRHTDMPKTKYK
```

Submit a file in FASTA format directly from your local disk:  Browse...

**Organism group**

Eukaryotes  
 Gram-negative bacteria  
 Gram-positive bacteria

**Method**

Neural networks  
 Hidden Markov models  
 Both

**Graphics**

No graphics  
 GIF (inline)  
 GIF (inline) and EPS (as links)

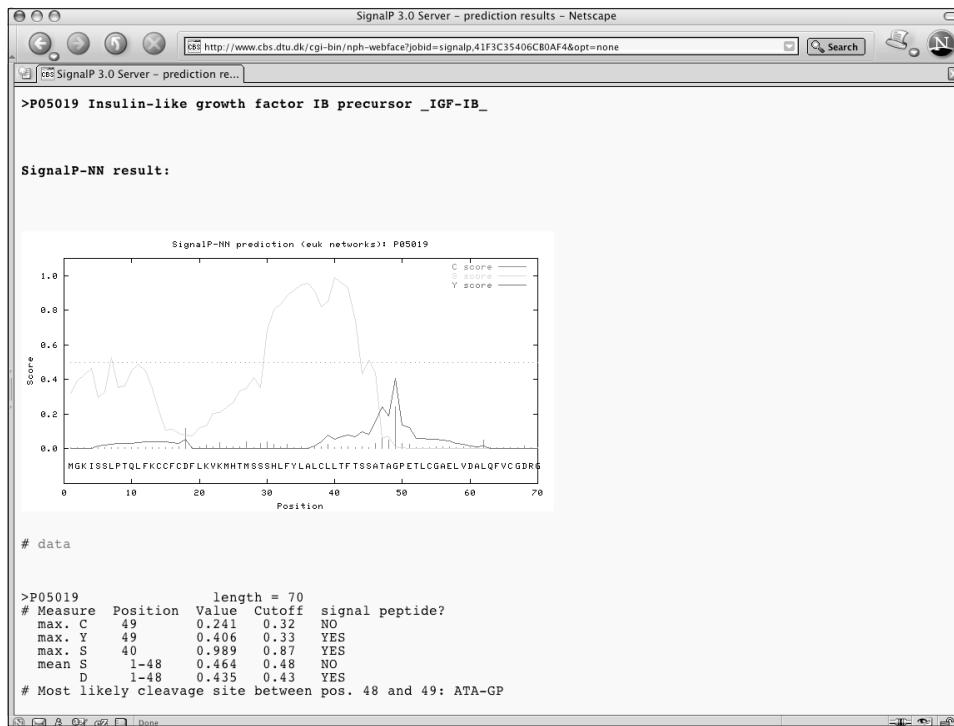
**Output format**

Standard  
 Full  
 Short (no graphics)

Submit | Clear fields

**Restrictions:**  
At most 2,000 sequences and 200,000 amino acids per submission; each sequence not more than 6,000 amino acids.

**Confidentiality:**  
The sequences are kept confidential and will be deleted after processing.

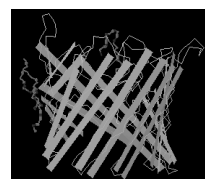


## Transmembrane Classes



- 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

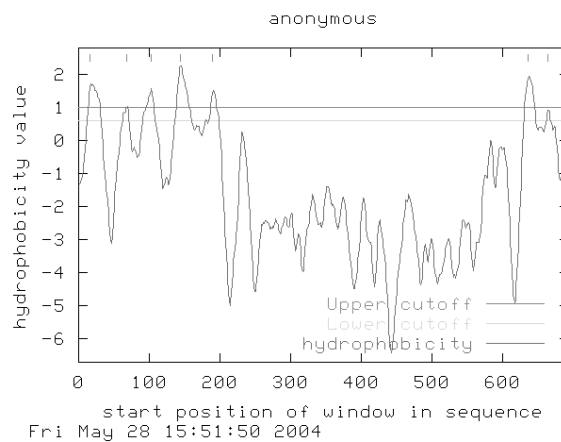
---

- 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
- Web-based search
  - <http://bioweb.pasteur.fr/seqanal/interfaces/toppred.html>

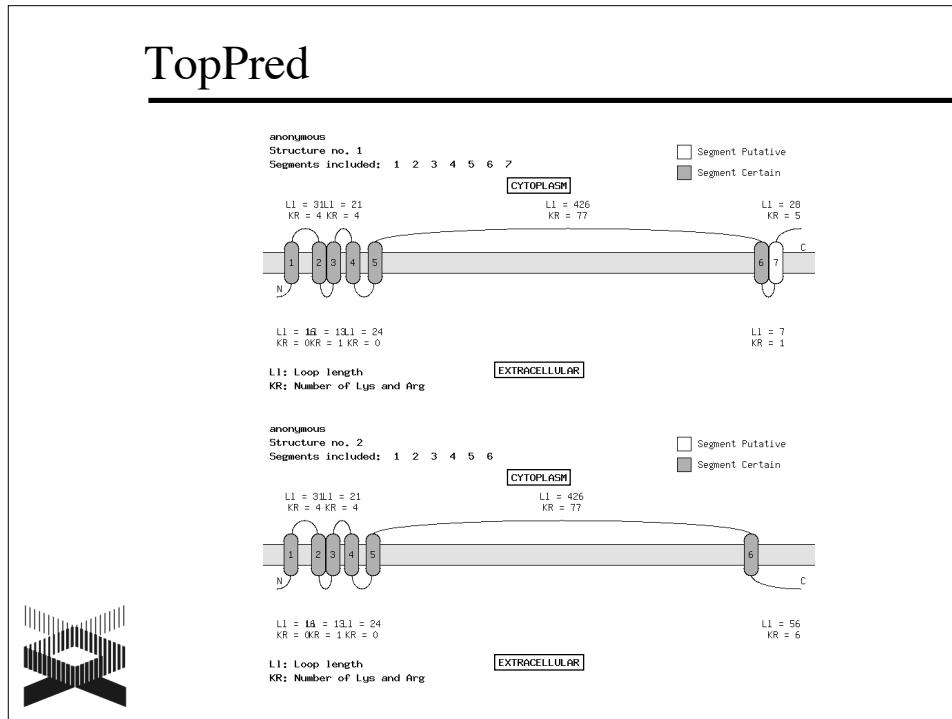


## TopPred

---



## TopPred



## 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



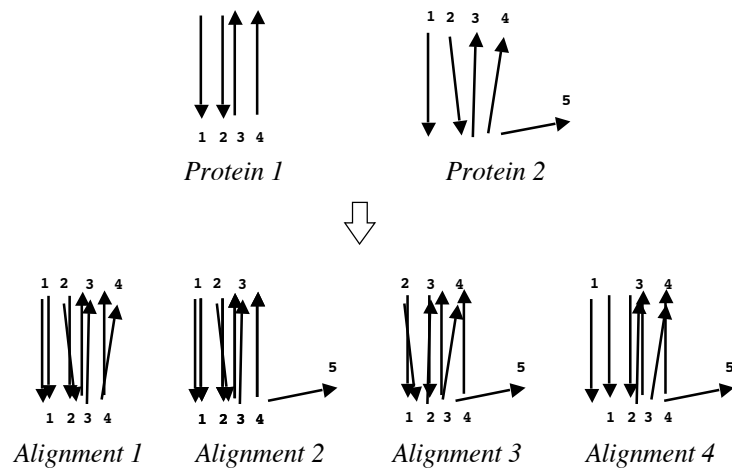
## VAST Structure Comparison

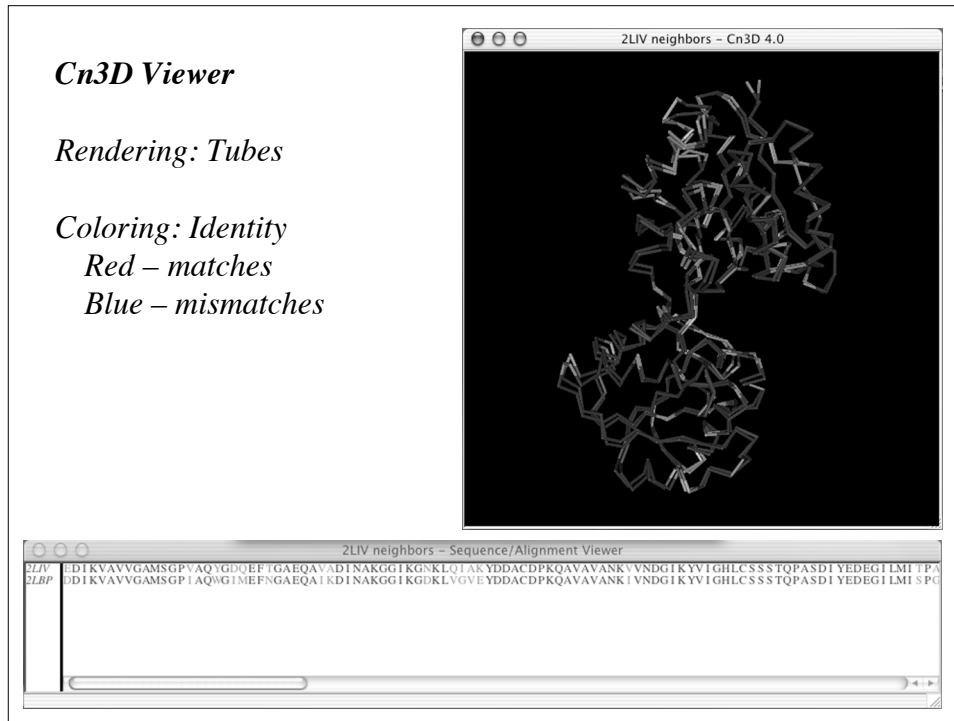
Step 1: Construct vectors for secondary structure elements



## VAST Structure Comparison

Step 2: Optimally align structure element vectors



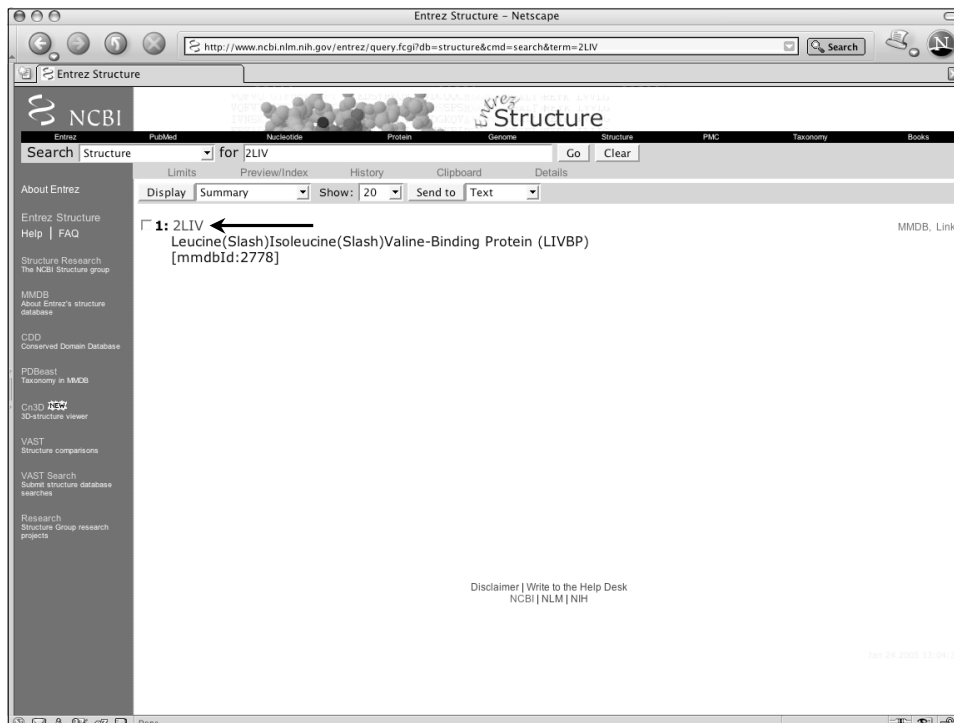
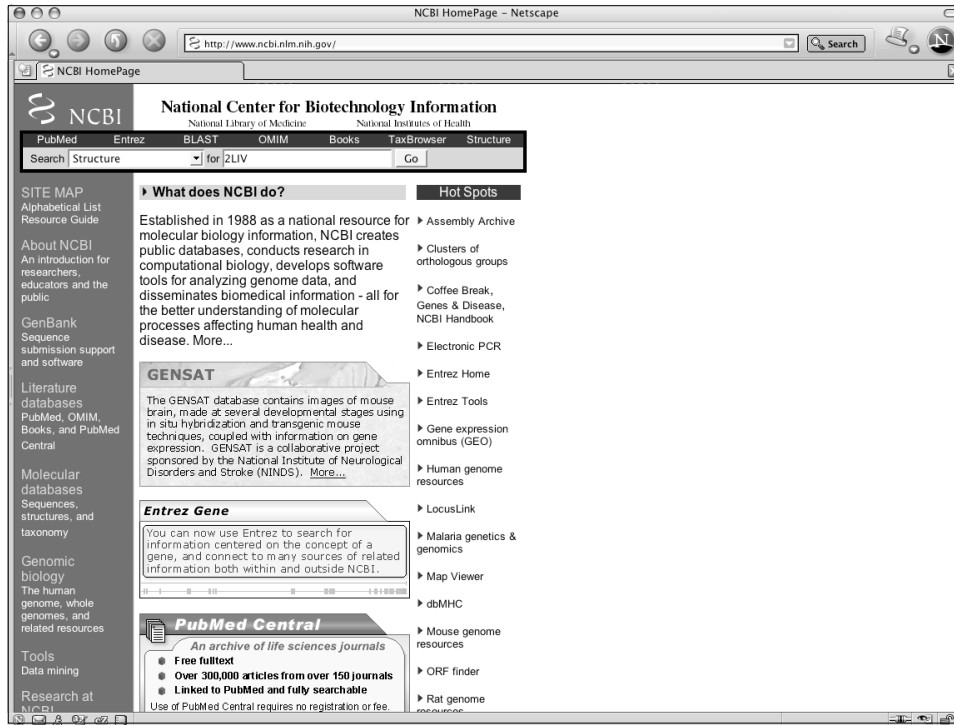


## VAST Shortcomings

- 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 2005  
 Biological Sequence Analysis II



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis II

Structure Summary - Netscape

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

**Description:** Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP).  
**Deposition:** J.S.Sack, M.A.Saper & F.A.Quiocho, 10-Apr-89  
**Taxonomy:** Escherichia coli  
**Reference:** PubMed **MMDB:** 2778 **PDB:** 2LIV

View 3D Structure of Best Model with Cn3D Display NEW Get Cn3D 4.1!

Protein Chain 1 50 100 150 200 250 300 344

3d Domains 1 2 1 2

CDs RNF\_receptor

Citing MMDB: Chen J, Anderson JB, DeWeese-Scott C, Fedorova ND, Geer LY, He S, Hurwitz DI, Jackson JD, Jacobs AR, Lanczycki CJ, Liebert CA, Liu C, Madej T, Marchler-Bauer A, Marchler GH, Mazumder R, Nikolskaya AN, Rao BS, Panchenko AR, Shoemaker BA, Simonyan V, Song JS, Thissen PA, Vasudevan S, Wang Y, Yamashita RA, Yin JJ, Bryant SH, "MMDB: Entrez's 3D-structure database". *Nucleic Acids Res.* 2003 Jan; 31(1): 474-7.

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

344 residues, click to see its structure neighbors

VAST Summary - Netscape

http://www.ncbi.nlm.nih.gov/Structure/vast/vastsvr.cgi?did=6728

NCBI **VAST** Structure Neighbors

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

**Query:** MMDB 2778, 2LIV  
**Description:** Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP)

View 3D Structure of All Atoms with Cn3D Display NEW Get Cn3D 4.1!

View Alignment using Hypertext for Selected VAST neighbors

List All sequences subset sorted by Vast\_P\_value page 1 in Table

Find MMDB or PDB ids: or 3D-Domain ids:

3395 neighbors found. 60 out of 453 representatives from the Medium redundancy subset displayed.

2LIV Chain 1 50 100 150 200 250 300 344 R11\_Res.

3d Domains RNF\_receptor

CDs

2LBP 344

1EMT B 332

1OP4 C 305

1Q00 B 256

1G00 B 253

2LBP 1 250

2LBP 2 227

1Q00 B 1 222

1G00 220

Aligned residues 284 to 299, click for sequence alignment

VAST Summary - Netscape

http://www.ncbi.nlm.nih.gov/Structure/vast/vaststrv.cgi?sid=6728&allfid=671001%2C4883801%2C4219601%2C3396101%2C

VAST Summary

NCBI VAST Structure Neighbors

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

Query: MMDB 2778, 2LIV  
 Description: Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP)

View 3D Structure of All Atoms with Cn3D Display Get Cn3D 4.1!

View Alignment using Hypertext for Selected VAST neighbors

List All sequences subset sorted by Vast P\_value page 1 in Table


Find MMDB or PDB ids: or 3D-Domain ids:

60 out of 3395 neighbors displayed.

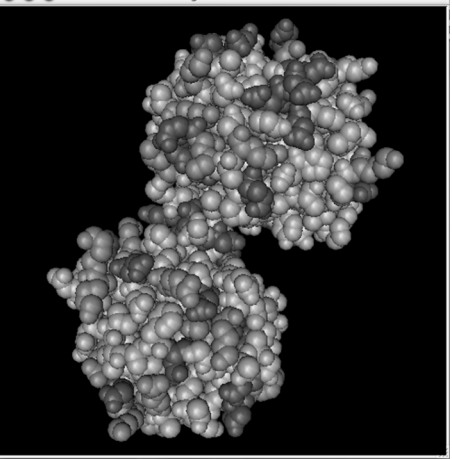
PDB	C	D	Ali.	Len.	SCORE	P-VAL	RMSD	%Id	MMDB	Date	Description
<input type="checkbox"/>	2LBP			344	39.8	10e-44.6	0.9	79.1	03/2001		Leucine-Binding Protein (LBP)
<input type="checkbox"/>	1USG	A		343	40.1	10e-42.4	2.0	79.0	01/2004		L-Leucine-Binding Protein, Apo Form
<input type="checkbox"/>	1JDP	B		310	29.9	10e-22.6	4.3	14.8	10/2001		Crystal Structure Of HormoneRECEPTOR COMPLEX
<input type="checkbox"/>	1ISS	B		330	30.3	10e-22.5	3.4	15.5	04/2002		Crystal Structure Of Metabotropic Glutamate Receptor Subtype 1 Complexed With An Antagonist
<input type="checkbox"/>	1JDP	A		322	29.8	10e-22.4	4.6	14.6	10/2001		Crystal Structure Of HormoneRECEPTOR COMPLEX

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

Read the descriptions!



2LIV neighbors - Cn3D 4.0



2LIV neighbors - Cn3D 4.0

Worms

Secondary Structure

Rendering

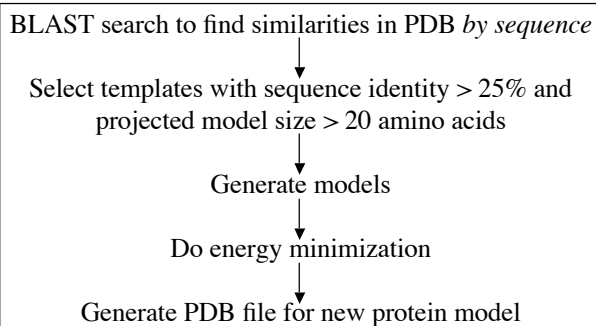
Coloring

Spacefill

Charge

## SWISS-MODEL

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



```

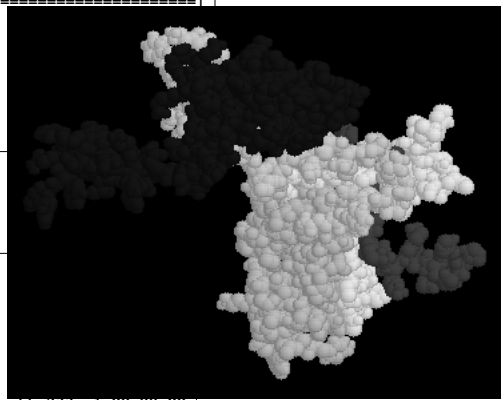
    21DJH.pdb: 42.77 % identity
    21DJG.pdb: 42.77 % identity
    11DJG.pdb: 42.22 % identity
    11QAS.pdb: 44.17 % identity
    11QAT.pdb: 43.52 % identity
    21QAT.pdb: 43.52 % identity
    21QAS.pdb: 43.52 % identity
    
```

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

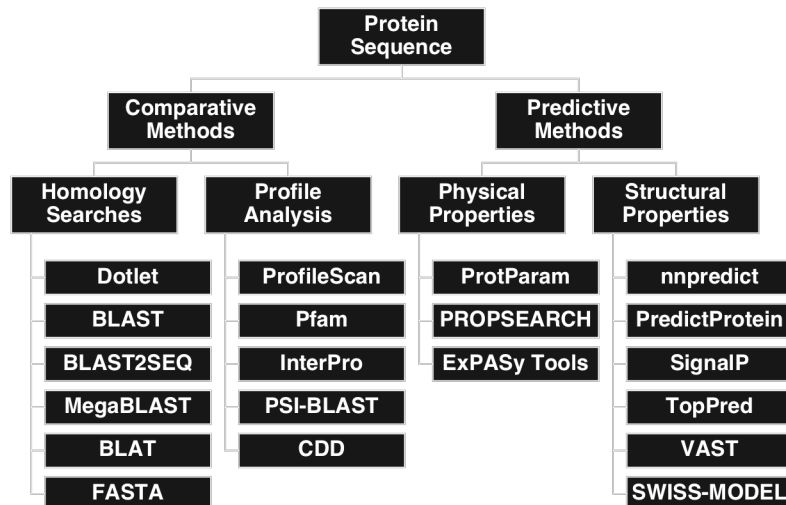


## Structural Modeling Software

- 3D-JIGSAW  
<http://www.bmm.icnet.uk/servers/3djigsaw>
- ESyPred3D  
<http://www.fundp.ac.be/urbm/bioinfo/esypred>
- MODELLER  
<http://www.salilab.org/modeller/modeller.html>
- Protinfo  
<http://protinfo.compbio.washington.edu>



## Protein Sequence Analysis



## Annual NAR Database Issue

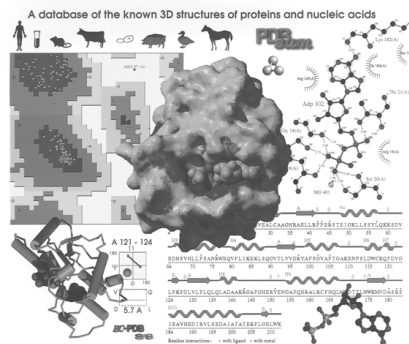
<http://nar.oupjournals.org>



ISSN 0305-1048

## Nucleic Acids Research

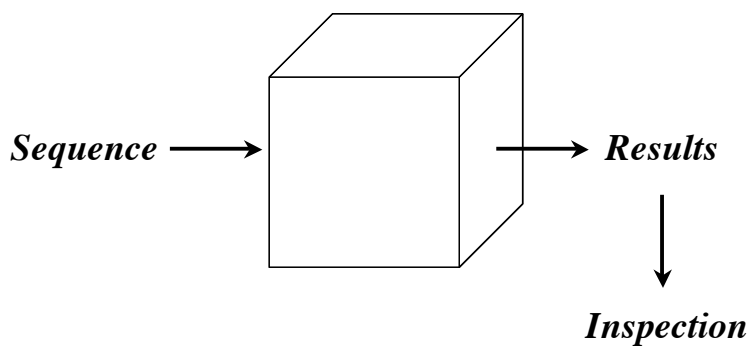
VOLUME 33 DATABASE ISSUE JANUARY 1 2005  
www.nar.oupjournals.org



Now Open Access  
No barriers to access - all articles freely available online



## Understanding Analyses



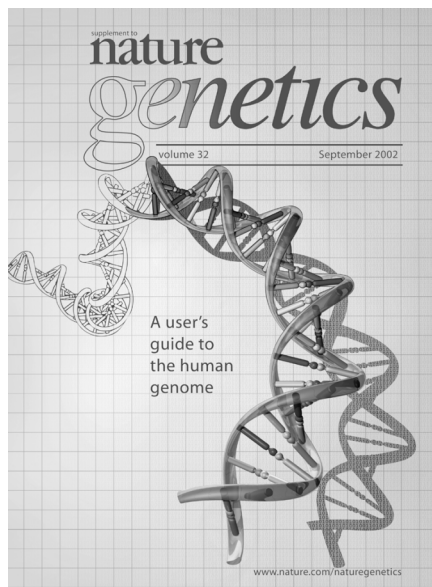


## A User's Guide to the Human Genome II

---

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

**Commentary:  
Keeping Biology  
in Mind**



## Further Reading

Apweiler, R., Attwood, T. K., Bairoch, A., Bateman, A., Birney, E., Biswas, M., Bucher, P., Cerutti, L., Corpet, F., Croning, M. D., et al. (2000). InterPro—an integrated documentation resource for protein families, domains and functional sites. *Bioinformatics* 16, 1145–1150. Describes an attempt to integrate several databases of motifs and patterns (including Pfam and PROSITE) into one comprehensive resource.

Ofran, Y. and Rost, B. (2005) Predictive methods using protein sequences. In *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins*, third edition (Baxevanis, A.D. and Ouellette, B.F.F., eds.), John Wiley and Sons. An overview of the methods used to generate pairwise sequence alignments and assess the biological significance of results.

Rost, B., Nair, R., Wrzeszczynski, K. O., and Ofran, Y. (2003). Automatic prediction of protein function. *CMLS Cell. Mol. Life Sci.* 60 1–14. More about automatic function prediction, with details about prediction of other aspects of functions and discussion of more methods, approaches, and prediction services.

Wolfsberg, T.G., Wetterstrand, K.A., Guyer, M., Collins, F.S., and Baxevanis, A.D. (2003) *A User's Guide to the Human Genome*, *Nature Genetics* 35, suppl. 1. Answers to frequently-asked questions about how to use genome sequence and annotation, written by researchers at the NHGRI to provide an introduction and guide for all genetics and bioinformatics researchers. <http://www.nature.com/ng/supplements/>

## References

*See also the annual Database Issue of Nucleic Acids Research, at <http://nar.oupjournals.org>, for updated versions of all papers describing individual database resources.*

Anfinsen, C. B. (1973). Principles that govern the folding of protein chains. *Science* 181, 223–230.

Bateman, A., Coin, L., Durbin, R., Finn, R. D., Hollich, V., Griffiths-Jones, S., Khanna, A., Marshall, M., Moxon, S., Sonnhammer, E. L., Studholme, D. J., Yeats, C., Eddy, S. R. (2004). The Pfam protein families database. *Nucl. Acids Res.* 32, D138–D141.

Chen, C. P., Kernytsky, A., and Rost, B. (2002). Transmembrane helix predictions, revisited. *Protein Sci.* 11: 2774–2791.

Cserzo, M., Wallin, E., Simon, I., von Heijne, G., and Elofsson, A. (1997). Prediction of transmembrane alpha-helices in prokaryotic membrane proteins: the dense alignment surface method. *Protein Eng.* 10, 673–676.

Falquet, L., Pagni, M., Bucher, P., Hulo, N., Sigrist, C. J., Hofmann, K., and Bairoch, A. (2002). The PROSITE database: its status in 2002. *Nucl. Acids Res.* 30, 235–238.

Hulo, N., Sigrist, C. J. A., Le Saux, V., Langenijk-Genevaux, P. S., Bordoli, L., Gattiker, A., De Castro, E., Bucher, P., and Bairoch, A. (2004). Recent improvements to the PROSITE database. *Nucl. Acids Res.* 32, D134–D137.

Rost, B., and Eyrich, V. (2001). EVA: large-scale analysis of secondary structure prediction. *Proteins Struct. Funct. Genet.* 45(Suppl 5), S192–S199.

Sonnhammer, E. L., Eddy, S. R., and Durbin, R. (1997). Pfam: a comprehensive database of protein domain families based on seed alignments. *Proteins Struct. Funct. Genet.* 28, 405–420.

von Heijne, G. (1992). Membrane protein structure prediction. *J. Mol. Biol.* 225, 487–494.