

Current Topics in Genome Analysis Spring 2005

Week 5 Biological Sequence Analysis II

Andy Baxevanis, Ph.D.



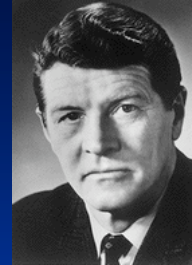
Overview

- 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

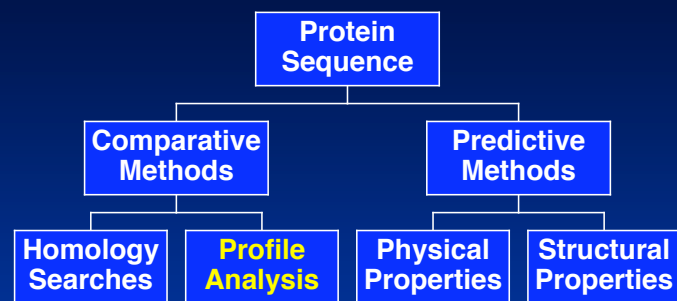


Protein Conformation

- 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

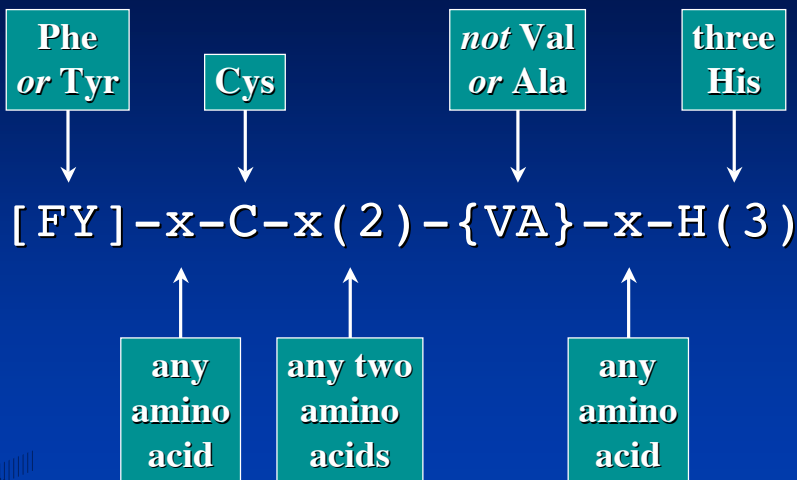
APHIIVATPG
 GCEIVIAATPG
 GVEICIAATPG
 GVDILIGATPG
 RPHIIVATPG
 KPHIIIAATPG
 KVQLIIATPG
 RPDIVIAATPG
 APHIIIVGTPG
 APHIIIVGTPG
 GCHVVIAATPG
 NQDIVVATPG

- 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	-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	20	-30	40	-10	20	20	-10	0	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	70	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>



myhits
Query Hub Result Help Database

Motif Scan user: anonymous
log in

Protein Sequence Input
Enter a protein sequence in RAW or FASTA or Swiss-Prot format or a db:AC or db:ID identifier

FASTA format

```
>NP_036673 cytochrome P450 [Rattus norvegicus]  
MAFSQYISLAPELLLATAI FCLVFNVLGRTQVFKGLSPGPGWLPFI  
YGVVIGIRIGSTFPVVLGSLNTI KQALVKGDDPGRGRLYSFTLITNGK  
DALSGTIRASDPTFVRYCYLLEEVYKKAHILIDFQGLMAEYQGEIYFNG  
KSEMLNLVYSSKDFENVTSGNADVDFPVLAYLNPALKRFYFNDFNFV  
DITCALFKSESNFKDGGLLIPEKIVNIVNDI FQAGFVTVTTAIFWSILL  
KDRQPLRSRPLQPLYLEAFILIEYRYTSFVPTIPISTTTDTSILNGHIF  
DQFVYKPERFLTNGATAIDWTLSEKQMLKGLGSRKCTGELPAGSRVFLYI
```

Clear input
Reset page

Motif scanning means finding all known motifs that occur in a sequence. This form lets you paste a protein sequence, select the collections of motifs to scan for, and launch the search. Some general [documentation](#) is available about the Prosite and Pfam collections of motifs. Another [document](#) deals with the interpretation of the match scores. You should consult the home pages of [Prosite](#) on EXPASY, [Pfam](#) and [InterPro](#) for additional information.

Warning! The scan might take a few minutes, thus if your proteins of interest are already in the sequence databases (see [list](#)), the [Query by Protein](#) form is much faster, and the [Protein Hub](#) provides a collection of tools that you might find useful.

Parameters

Database of motifs ([db description](#))

- PROSITE patterns
- PROSITE patterns (frequent match producers)
- PROSITE profiles
- Profile (more profiles)
- HAMAP profiles
- Pfam HMMs (local models)
- Pfam HMMs (global models)

search

Question or comment about this page.

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].

FT	MYHIT	449	458	pat:CYTOCHROME_P450 [1]
	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

HMNER2.0 [2.3.1] [documentation at Sanger Institute]

NAME p450
 ACC PF00067.10
 DESC Cytochrome P450
 LENG 499
 ALPH Amino
 RF no
 CS no
 MAP 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
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	-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	-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	-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

Pfam Protein families database of alignments and HMMs

Wellcome Trust Sanger Institute

p450

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)

Figure 1: 1mpw Oxidoreductase
 Molecular recognition in (+)- α -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](#)]

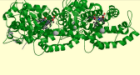
Alignment Seed (52) Full (3878)
 Format: Coloured alignment [View HMM logo](#)

Domain organisation View 13 representative architectures View architectures for 3878 proteins
 Zoom: 0.5 pixels/aa. [View Graphic](#)

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)

Key:

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

The SwissProt/PDB mapping was provided by HSD

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

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](#)

Q7P63 [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](#)

Q94HMS [oryza sativa (rice)] putative cytochrome p-450 like protein



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

Pfam: p450

Protein families database of alignments and HMMs

Wellcome Trust Sanger Institute

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 (+)- α -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 HSD

1akd Display pdb

For additional annotation, see the PROSITE document PDOC0081 | [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](#)

Browser: InterPro: Cytochrome P450 - Netscape
 URL: http://www.ebi.ac.uk/interpro/Entry?ac=IPR001128

EMBL-EBI European Bioinformatics Institute

EBI Home | About EBI | Research | Services | Toolbox | Databases | Downloads | Submissions

InterPro

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
 IPR02401: E-class P450, group I
 IPR02402: E-class P450, group II
 IPR02403: 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

InterPro: 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)
 SSF48284.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
 PfamMSD - [click here](#)

Database links PANDIT PF00067

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

InterPro: Cytochrome P450

COMe PRX000238
 Blocks IPR001128
 PROSITE doc PDOC00081

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
		11

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

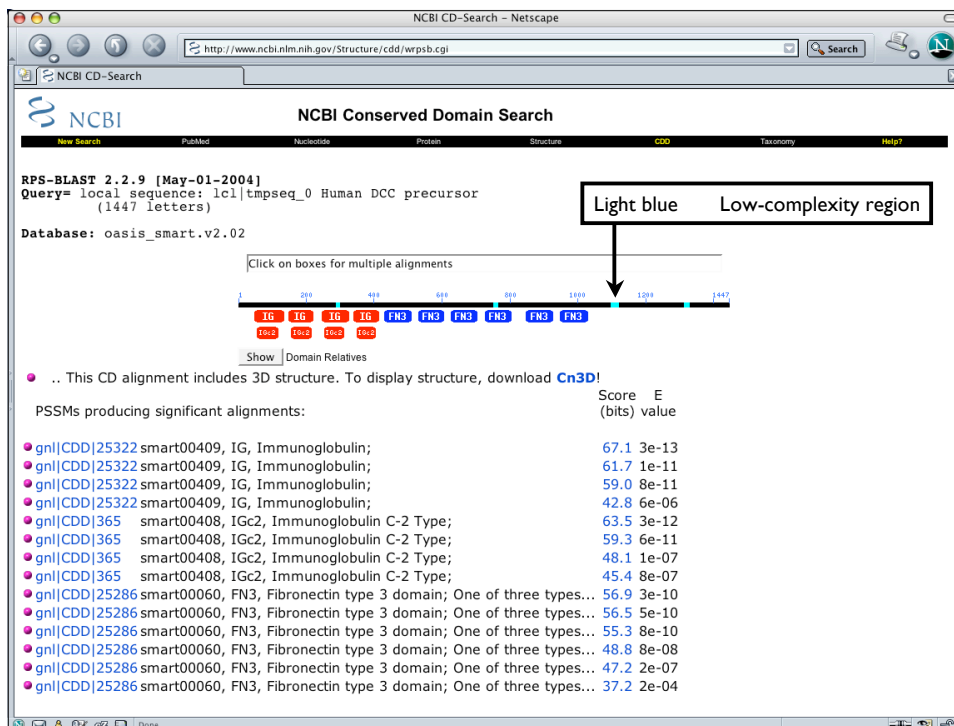
There is no significance to the placement of individual nodes on the circles

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>



NHGRI Current Topics in Genome Analysis 2005
 Biological Sequence Analysis II



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

Light blue Low-complexity region

.. This CD alignment includes 3D structure. To display structure, download [Cn3D!](#)

PSSMs producing significant alignments:

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

gn|CDD|25322, smart00409, IG, Immunoglobulin;

CD-Length = 86 residues, 98.8% aligned
 Score = 67.1 bits (163), Expect = 3e-13

Query: 337 PSNLYAYESMDIEFECTVSGKPPVPTVNMKN-CDVVIPSDYFQIVGGSN---LRILGVVK 392
 Sbjct: 1 PPSVTVKEGESVTLSCAESGNPPPEVTWYKGGKLLAYSGRFSVSRSGNSTLTISNVTP 60

Query: 393 SDEGFYQCAENEAGNAQTS AOLIV 417
 Sbjct: 61 EDSGTTC AATNSSGSASSGTTLTV 85

gn|CDD|25322, smart00409, IG, Immunoglobulin;

CD-Length = 86 residues, 91.9% aligned
 Score = 61.7 bits (149), Expect = 1e-11

Query: 147 ESVTAFMGDTVLLKCEVIGEPMPITHWQKNQDLTPIPGDSRVVLPSPG---ALQISR LQ 203
 Sbjct: 2 PPSVTVKEGESVTLSCAESGNPPPEVTWYK--OGKLLAYSGRFSVSRSGNSTLTISNVTP 59

Query: 204 PGDIGYRCSARNPASSRTGN 224
 Sbjct: 60 PEDSGTYTCAATNSSGSASSG 80

gn|CDD|25322, smart00409, IG, Immunoglobulin;

CD-Length = 86 residues, 100.0% aligned
 Score = 59.0 bits (142), Expect = 8e-11

Query: 246 PSNVAIEGKDAVLECCVGYPPSPFWLARGEEVQLRSKYY---SLLGGSNLLISWTD 302
 Sbjct: 1 PPSVTVKEGESVTLSCAESGNPPPEVTWYKGGKLLAYSGRFSVSRSGNSTLTISNVTP 60

Query: 303 DDSGMYTCVVYTKNENISASAE LTVL 328
 Sbjct: 61 EDSGTTC AATNSSGSASSGTTLTVL 86

gn|CDD|25322, smart00409, IG, Immunoglobulin;

CD-Length = 86 residues, 98.8% aligned
 Score = 42.8 bits (100), Expect = 6e-06

NHGRI Current Topics in Genome Analysis 2005
 Biological Sequence Analysis II

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

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

NCBI CDD smart00409 - Netscape

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

NCBI CDD smart00409

Conserved Domains

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

smart00409.10 Immunoglobulin; **IG**

Links: 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 PPSVTVKEGESVTLSCAEG.[1].PPPEVTW YK.[2].GKLL.[6].SVSR.[3].NSTLTISNVTPE.[2].63
 IMCP_H 7 SGGGLVQPGGSLRLSCATSG.[3].SDFYMEW VR.[6].LEWI.[22].IVSR.[5].ILYLQMNALRAE.[2].93
 1GYA 6 ALETWALGQDINLDIPSPQ.[3].DIDDIKW EK.[3].KKKI.[16].KLFK NGTLKIKHLKTD.[2].78
 1ZXO 9 PKKLAVEPKGSLEVNCSSTC.[1].QPEVGLL ET.[1].LNKI LLDE.[3].WKHYLVSNISHD 62
 g1 399208 25 QRLLIVANRTATLVNCYTY.[4].KEFRASL HK.[4].AVEV.[20].RGIH.[3].KVIFNLWNMSAS.[2].106

NCBI CDD smart00409 - Netscape

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

NCBI CDD smart00409

Show Structure

Program: Cn3D

Drawing: Virtual Bonds

(download Cn3D)

[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 PPSVTVEGESVTLSCAEG.[1].PPPEVTW YK.[2].GKLL.[6].SVSR.[3].NSTLTISNVTPE.[2]. 63
IMCP_H 7 SGGGLVQPGGSLRLSCATSG.[3].SDFYMEW VR.[6].LEWI.[22].IVSR.[5].ILYLQNALRAE.[2]. 93
1GYA 6 ALETWALGQDINLDIPSTQ.[3].DIDDIRK EK.[3].KKKI.[16].KLFR NGTLKIKHLKTD.[2]. 78
1ZXO 9 PKKLAVPEKGSLEVNCSTTC.[1].QPEVGLL ET.[1].LNKI LLDE.[3].WRHYLVSNISHD 62
g1 399208 25 QRPLLVANRTATLVCMYYT.[4].KEFRASL HK.[4].AVEV.[20].RGLH.[3].KVIFLWMMSSA.[2]. 106
g1 461714 216 SNTFYAREGDQVFSFPLSF.[2].ENLVGEL RW.[9].LWIS.[19].QMKR.[2].PLRFIIPQVLSR.[2]. 298
g1 6166597 64 PQGGTVKVGEDITPIAKVKA.[6].PTIKWFK.[1].KW.[6].AGKH.[7].ERHS.[3].TFEMQIRAKDN.[2]. 137
g1 729801 435 QRTQYGLVGDITARIECFASS.[3].ARHVSWT FN.[1].QEIS.[7].SILV.[7].KSTLIIRDSQAY.[2]. 503
g1 124310 243 LRTISASLGSRLTIPCKVFL.[4].PLTMTLW WT.[1].NDTH.[18].SENN.[4].EVLPIFDVPTRE.[3]. 321
g1 1709202 281 REGETMSLGRKVVITPEIKH FOPEIRW YR.[1].GVPL.[6].QTLW.[3].RATLTFSHLNKE.[2]. 341

consensus 64 GTYTCAT.[2].SGSASS GTTLTVL 86
IMCP_H 94 AIYYCARN.[6].YFDVWG.[1].GTTVTVS 121
1GYA 79 DIYKVSII.[4].KNVLEK IFDLKIQ 103
1ZXO 63 TVLQCHFT.[2].GKQESM NSNVSIV 85
g1 399208 107 DIYFCIE.[7].VYNEKS.[1].GTVIHVR 135
g1 461714 299 GSGILTLN.[2].KGTLYQ EVNLVVM 321
g1 6166597 138 GNVRCVET.[2].DRFDSC SFDFEVH 160
g1 729801 504 GRYNCTVV.[2].YGNQVA EIQLOAK 526
g1 124310 322 MDFKCVVH NTLSTFO TLRTTVK 342
g1 1709202 342 GLYTIKVR MGEYVE QYSAYVF 362
    
```

Citing CDD: Marchler-Bauer A, Anderson JB, Cherkuri 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, Mullokandov 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 CD-Search - Netscape

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

NCBI CD-Search

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

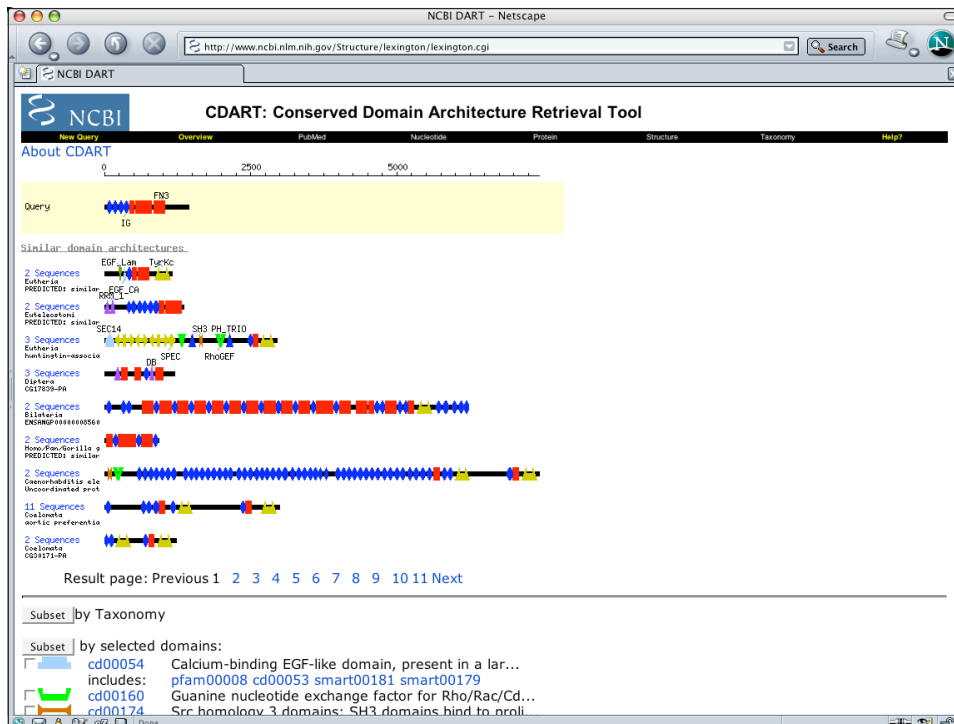
IG IG IG IG FN3 FN3 FN3 FN3 FN3 FN3

Show Domain Relatives

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

PSSMs producing significant alignments:

Accession	Domain	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

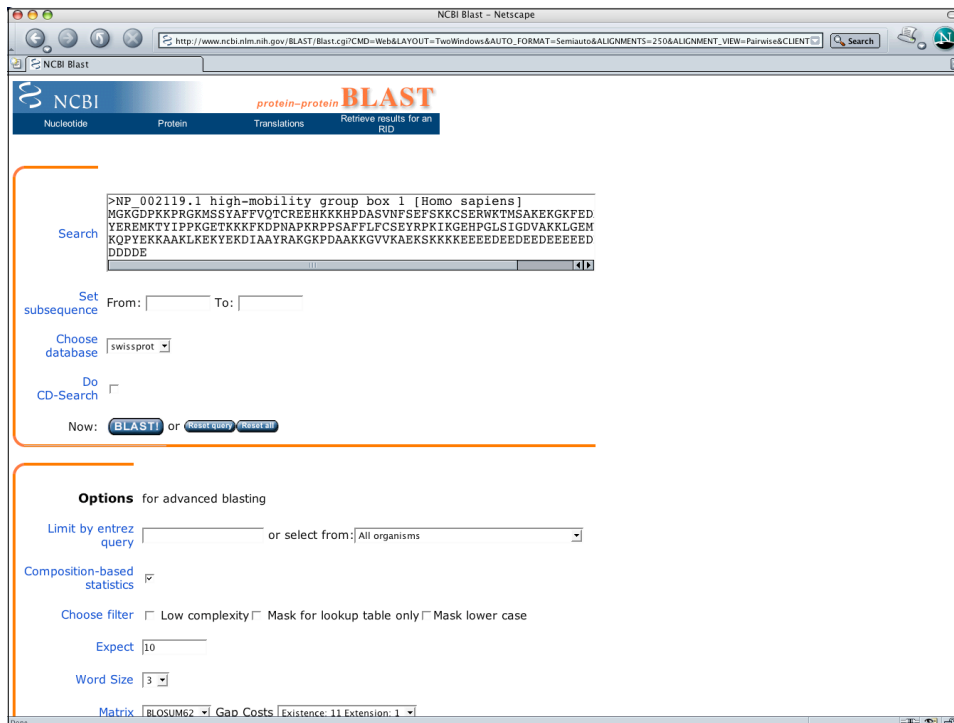
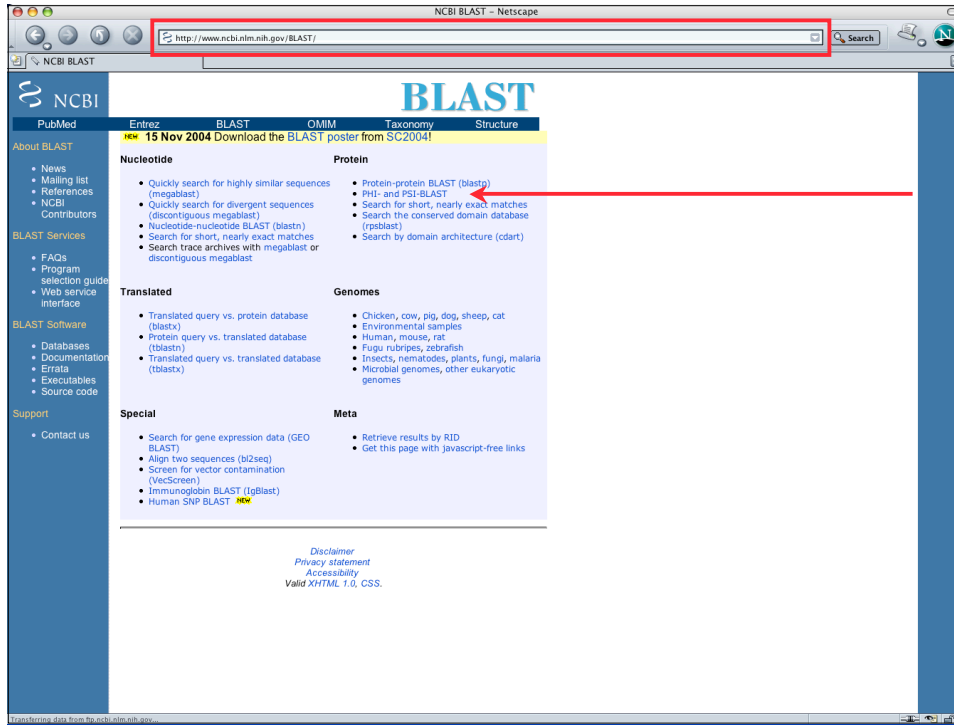


PSI-BLAST

- 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



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-gi 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

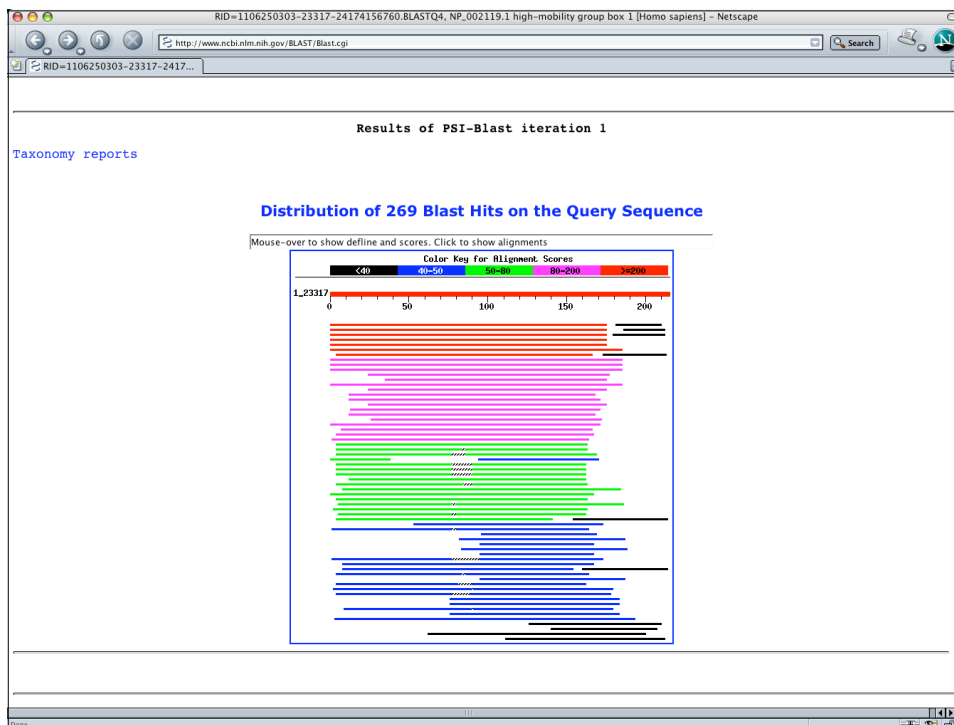
Expect value range:

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

Autoformat Semi-auto

BLAST!

Get the URL with preset values?



NHGRI Current Topics in Genome Analysis 2005
 Biological Sequence Analysis II

RID=1106250303-23317-2417... high-mobility group box 1 [Homo sapiens] - Netscape

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

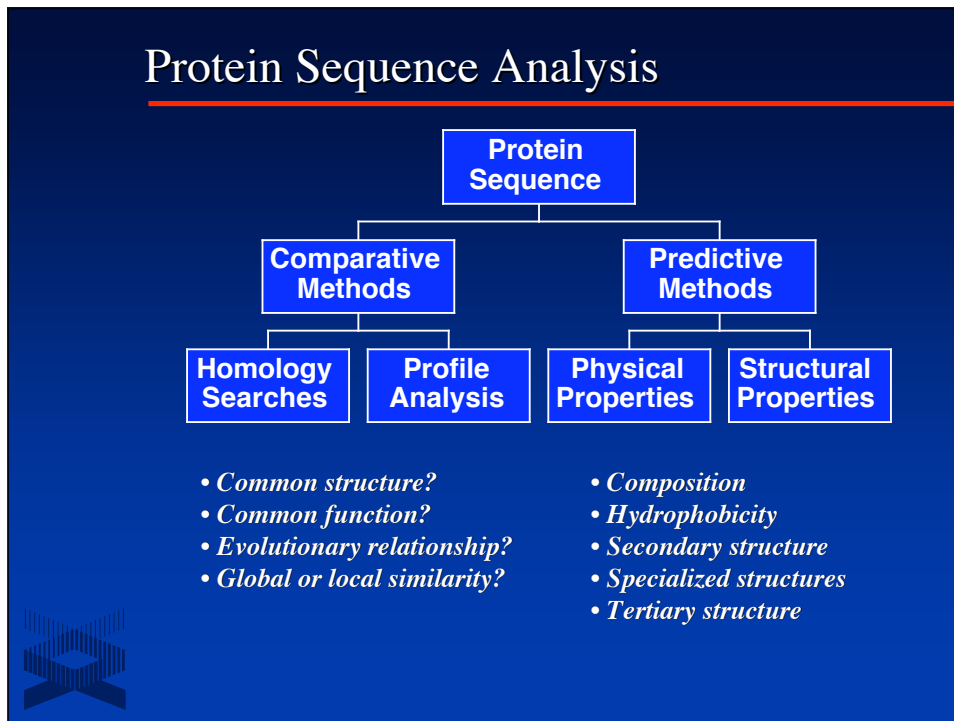
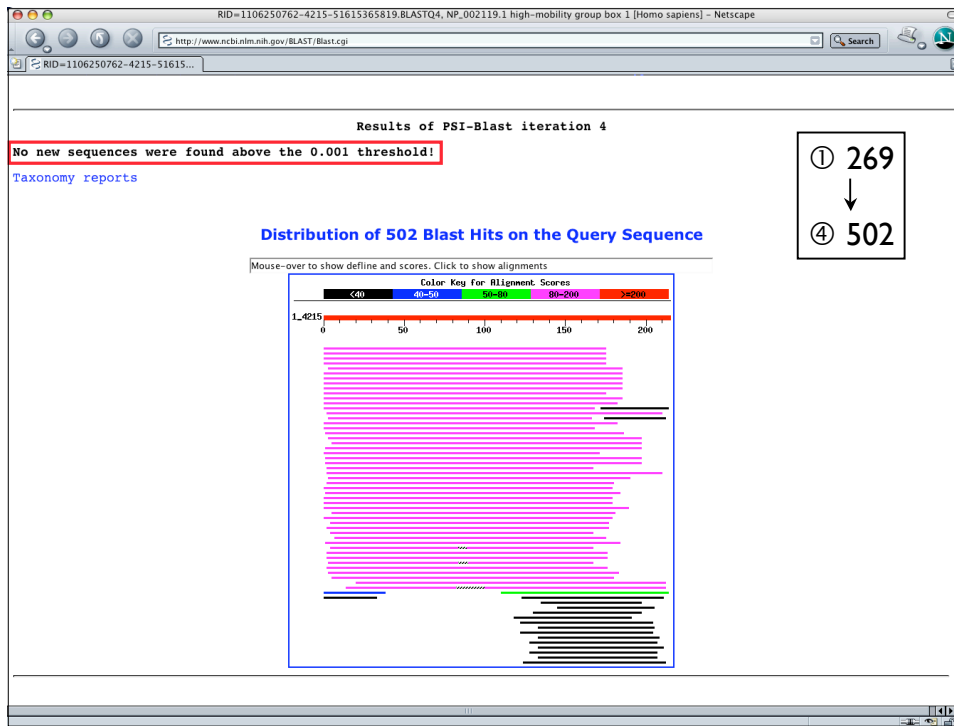
Sequences producing significant alignments:

	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 Q09UGV6 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 HMG2_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

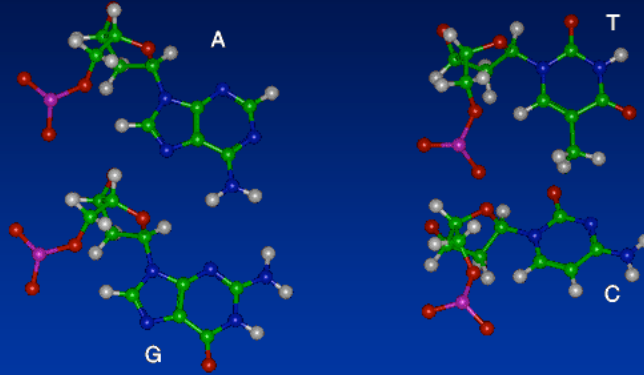
RID=1106250303-23317-2417... high-mobility group box 1 [Homo sapiens] - Netscape

Sequences with E-value WORSE than threshold

✔ 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
✔ 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

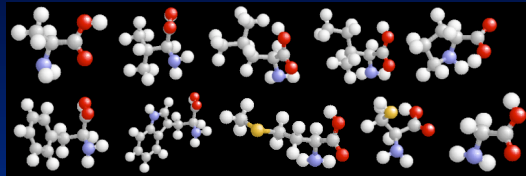


Information Landscape



Information Landscape

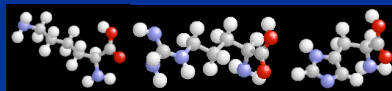
Nonpolar



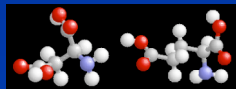
Polar Neutral



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

MNGEADCPDLEMAAPKQDRWSQEDMLTLLCEMKNLPSNDSSKFKTESHMDWEKVAFKDFSGDMCKL
 KWVEISNEVRKFRITLTELILDAQEHVKNPYKGGKLLKHPDFPKKPLTPYFRFFMEKRAKYAKLHPEM...

↓ 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*



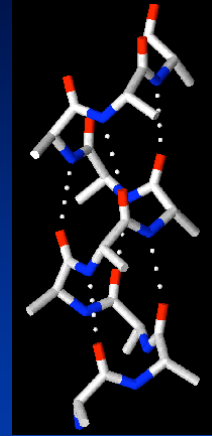
Secondary Structure Prediction

- 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



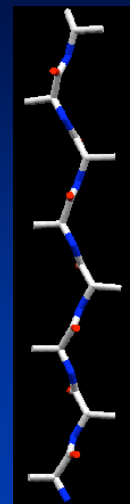
Alpha-helix

- 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



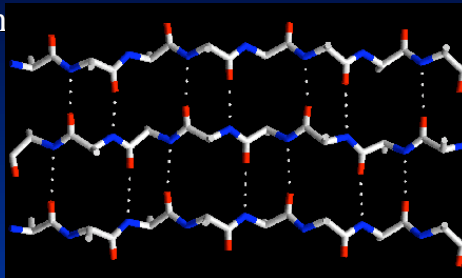
Beta-strand

- 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



α

Cyt c

Globins
 Orthogonal
 EF-hand
 Up-Down
 Cytochrome

β

CD4

Orthogonal
 Super-barrel
 Greek key
 Sandwich
 Jelly roll

$\alpha+\beta$

*Staph
 nuclease*

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

α/β

*Triose
 phosphate
 isomerase*

TIM barrel
 Doubly-wound

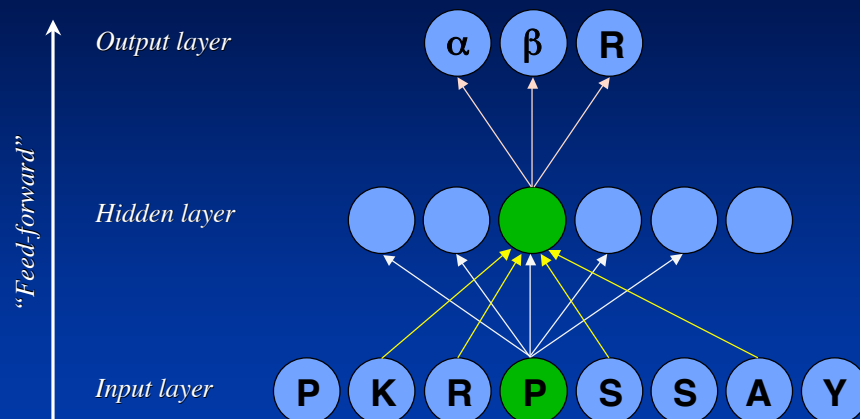


Neural Networks

- 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
 - Web <http://www.cmpharm.ucsf.edu/~nomi/nnpredict.html>



nnpredict Query

```
option: a/b
>flavodoxin - Anacystis nidulans
AKIGLFYGTQGTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIY
DDLDSVNFQGGKVAIFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQPDLTKNRIKTWVSQKSEFGL
```

↓ *α/β folding class*

Tertiary structure class: alpha/beta

```
Sequence:
AKIGLFYGTQGTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVG
ELQSDWEGIYDDLDSVNFQGGKVAIFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYW
PIEGYDFNESKAVRNNQFVGLAIDEDNQPDLTKNRIKTWVSQKSEFGL
```

```
Secondary structure prediction (H = helix, E = strand, - = no prediction):
---EEE-----EEHHHHHHH-----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
DDLDSVNFQGGKVAIFYGAGDQVGYSDNFQDAMGILEEKISSLSGQTWGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQDPLTKNRITWVSQLKSEFGL
```



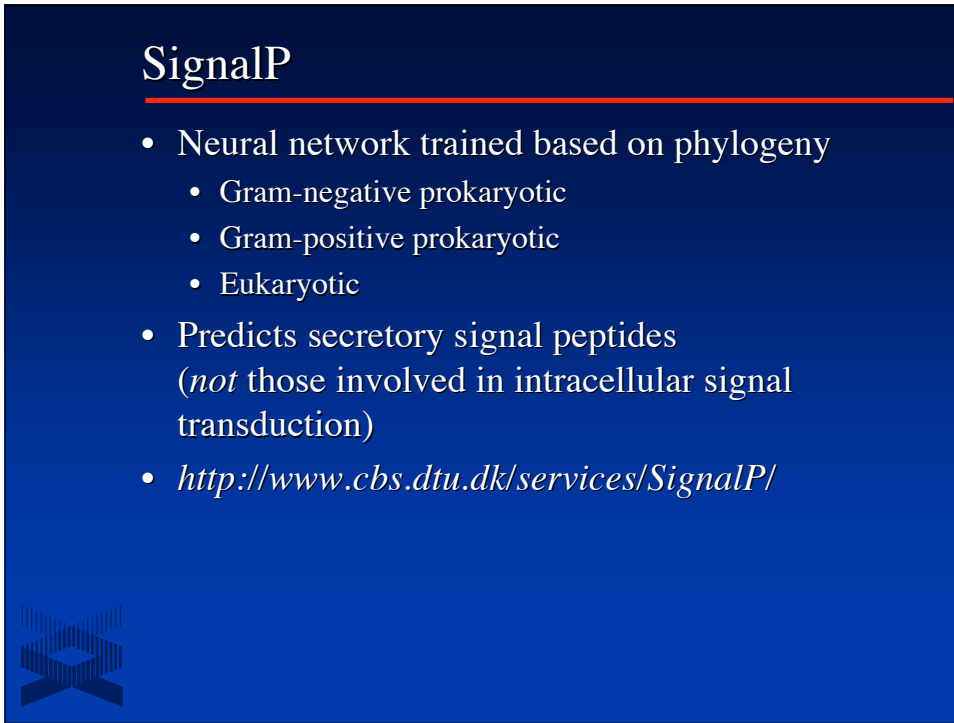
```
PROF results (normal)
.....1.....2.....3.....4.....5.....6.....7.....8.....9.....1
AA      AKIGLFYGTQTGVTTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIY
OBS_sec EEEEEEE HHHHHHHHHHHH EEEEE EEEEE HHHHHHHHHH EEEEEEE HHHH
Rel_sec 92789984267626888889888731068425442235670001225537887203476665421678888863045688863788841466443311446

0.1.....11.1.....12.1.....13.1.....14.1.....15.1.....16.1.....
AA      AMGILEEKISSLSGQTWGYWPIEGYDFNESKAVRNNQFVGLAIDEDNQDPLTKNRITWVSQLKSEFGL
OBS_sec HHHHHHHHHH EEEEE EEEEE HHHHHHHHHHHH
Rel_sec 788899888740782542202456544533100178268876426664211178899999888754289
```

- 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)
MGKISSLPQLFKCCDFLKVKHTMSSSHLYLALCLLFTSSATAGPELLOGAELVDALQVY
FFYFNKPTGYGSSRRAPQTGIVDECCFRSCDLRRLREMYCAPLKPARSARSVRAQRHTDMPKTKYV
```

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)

Truncation

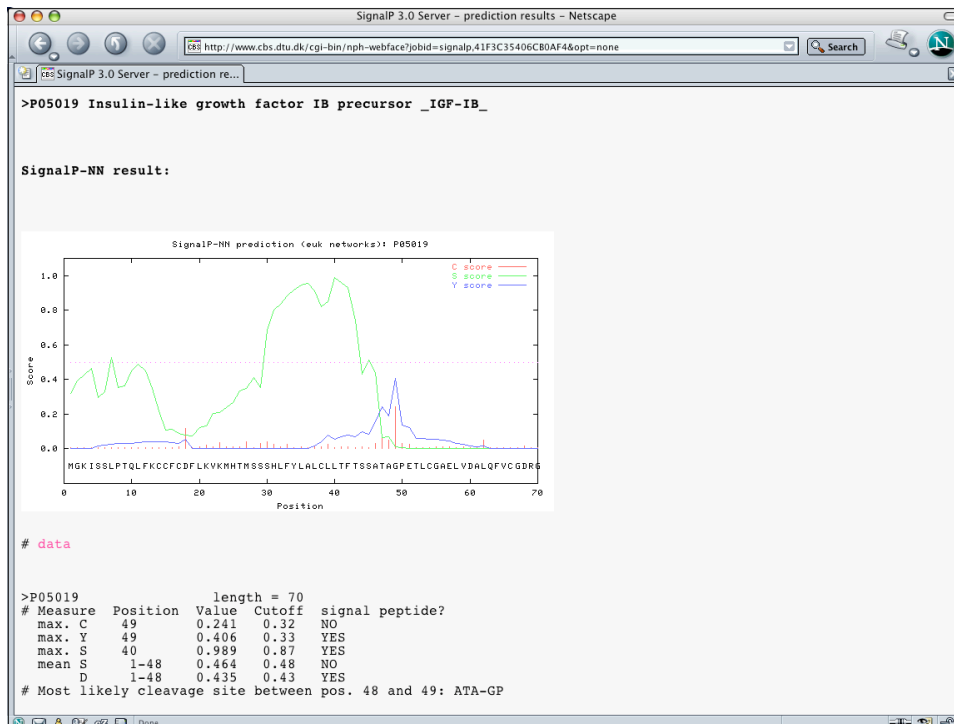
Truncate each sequence to max. 70 residues.

We recommend that only the N-terminal part of each protein sequence is submitted. Enter 0 (zero) to disable truncation.

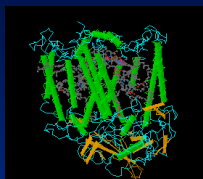
Submit | Clear fields

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.

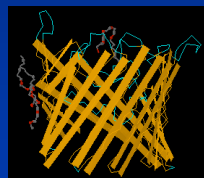


Transmembrane Classes



- 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)

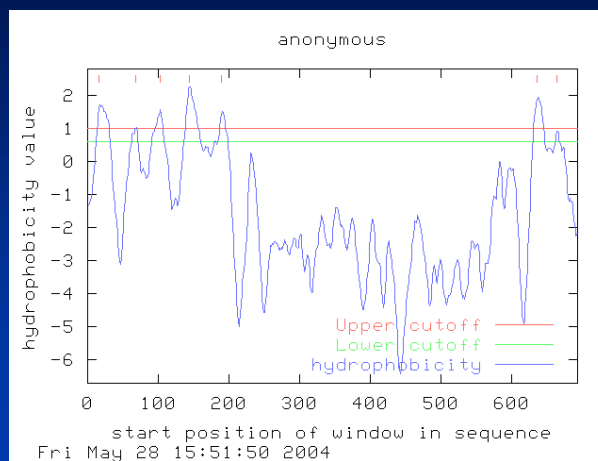


TopPred

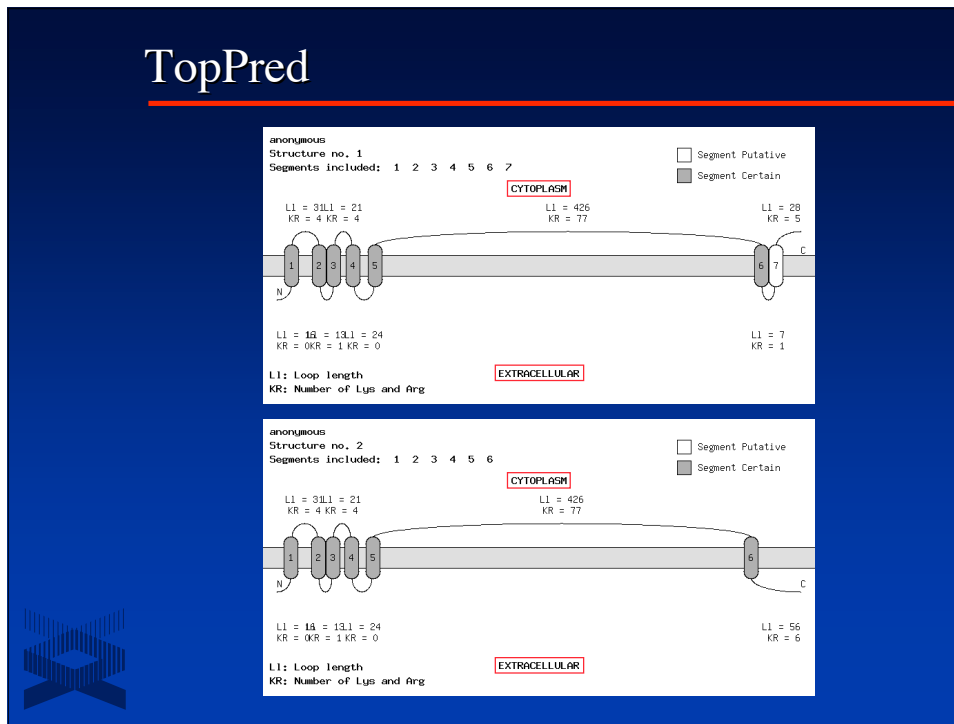
- Combines hydrophobicity analysis with the analysis of electrical charges
 - Calculates hydrophobicity profile
 - Hydrophobic-rich regions marked as “transmembrane”
 - Hydrophobic regions that fail to exceed a predefined cutoff are considered “putative transmembrane”
 - Topology prediction with and without putative helices
- Web-based search
 - <http://bioweb.pasteur.fr/seqanal/interfaces/toppred.html>



TopPred



TopPred

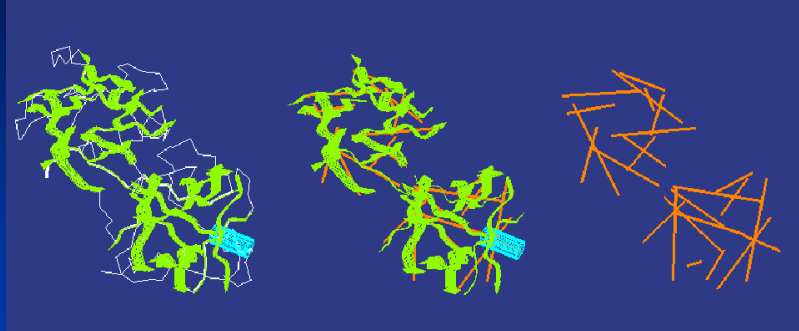


Predicting Tertiary Structure

- 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

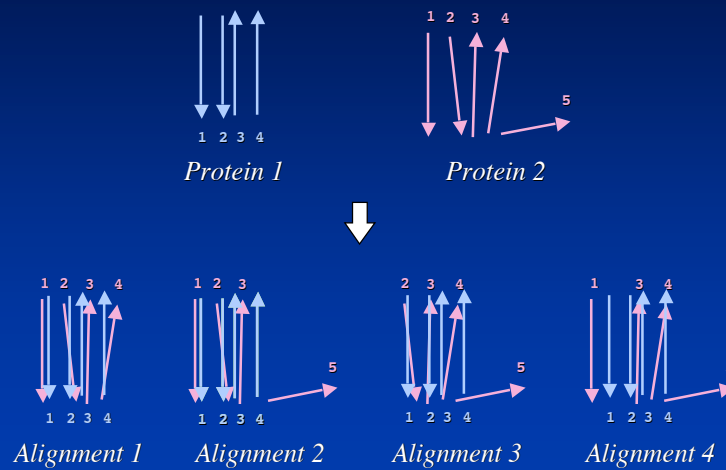
VAST Structure Comparison

Step 1: Construct vectors for secondary structure elements



VAST Structure Comparison

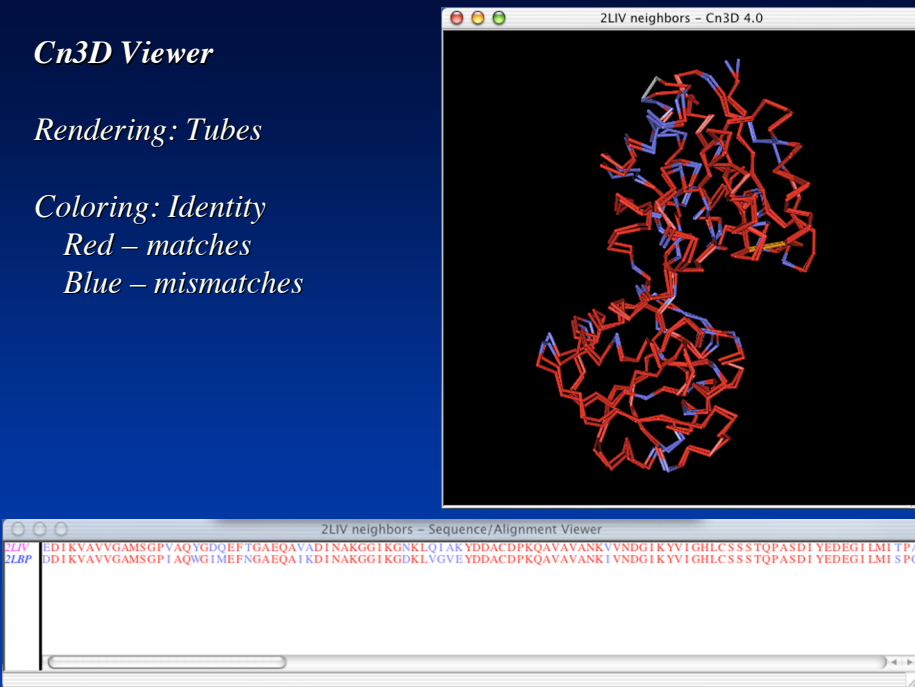
Step 2: Optimally align structure element vectors



Cn3D Viewer

Rendering: Tubes

Coloring: Identity
Red – matches
Blue – mismatches



2LIV neighbors - Cn3D 4.0

2LIV neighbors - Sequence/Alignment Viewer

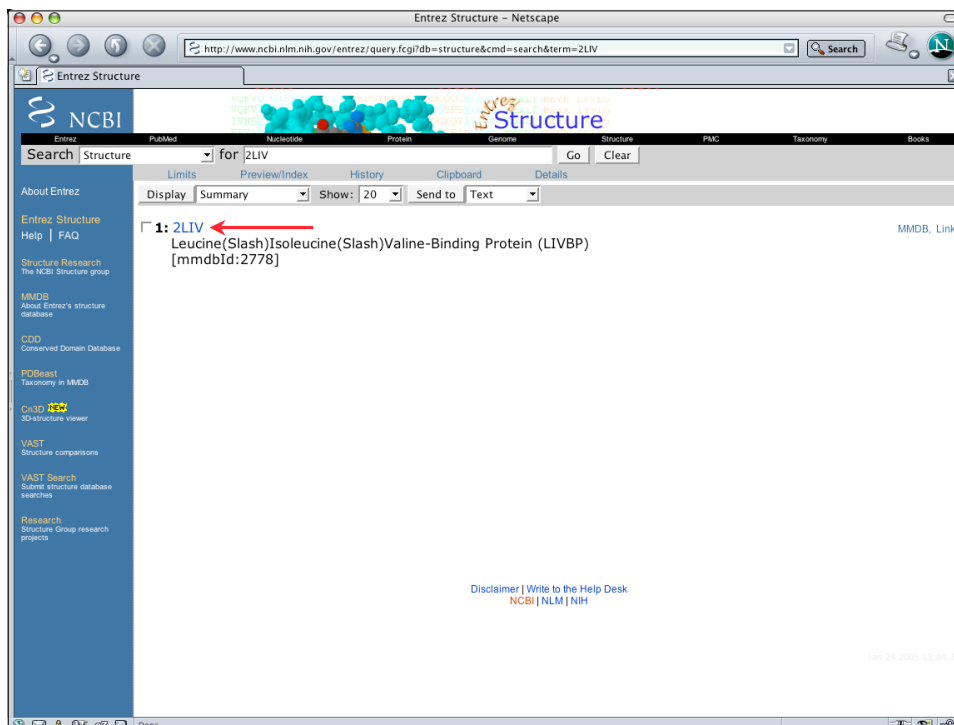
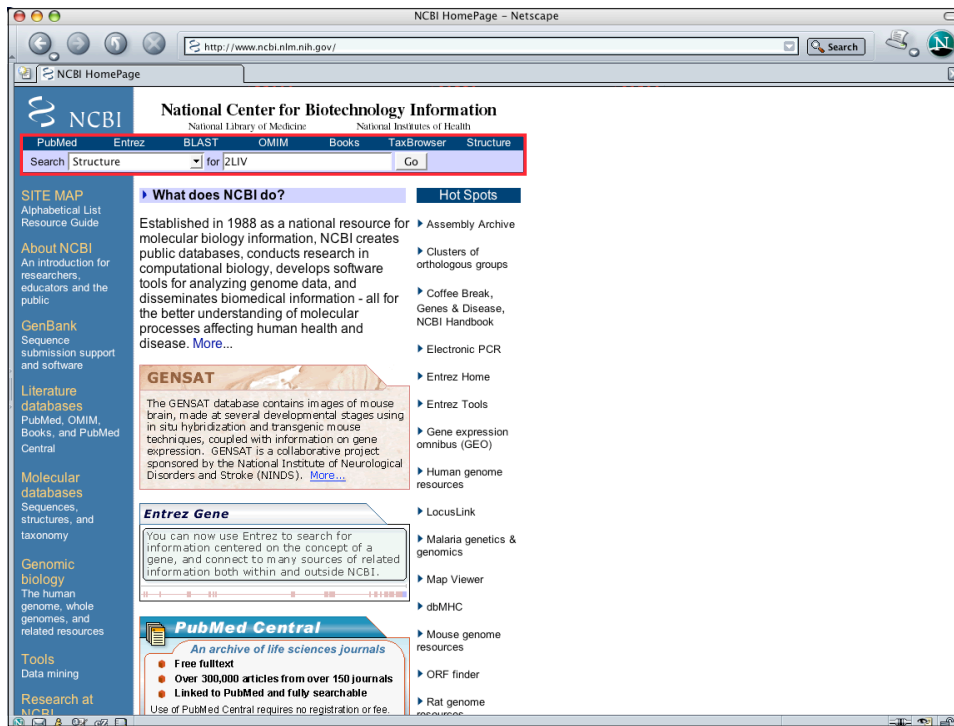
```
2LIV DDIKVAVVVGAMSGPVAQYGDQEFVIGAEQAVADINAKGGTRGNKLOIAKYDDACDPKQAVAVANKVVDGTRKVVIGHLCSSTQPPASDIYEDEGLMIITPA  
2LBP DDIKVAVVVGAMSGPIAQWGI MEFGAEQAIRKIDINAKGGIRGDKLVGVEYDDACDPKQAVAVANKVVDGTRKVVIGHLCSSTQPPASDIYEDEGLMISPG
```

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

The screenshot shows the NCBI MMDB Structure Summary page for the protein Leucine(Isoleucine)Valine-Binding Protein (LIVBP). The page includes a navigation menu with options like PubMed, BLAST, Structure, Taxonomy, OMIM, Help?, and Cn3d. The description states: "Description: Leucine(Isoleucine)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". Below this, there are controls for viewing the 3D structure, including a dropdown for "Best Model" and "with Cn3D". A protein domain diagram is shown with a red arrow pointing to the end of the chain at residue 344. The diagram labels "Protein Chain", "3d Domains" (1, 2, 1, 2), and "CDs" (RNF_receptor). At the bottom, there is a "Citing MMDB" section with a reference to Chen et al. (2003) and a disclaimer.

The screenshot shows the NCBI VAST Structure Neighbors page. The query is "MMDB 2778, 2LIV". The description is "Leucine(Isoleucine)Valine-Binding Protein (LIVBP)". The page features a search interface with a "View Alignment" section where "using Hypertext" is selected for "Selected" VAST neighbors. A red box highlights the "subset sorted by Vast P_value" dropdown. Below this, there are search fields for "MMDB or PDB ids" and "3D-Domain ids". The results section shows "3395 neighbors found. 60 out of 453 representatives from the Medium redundancy subset displayed." A table of results is shown, with the top entry being the query protein (2LIV) with a length of 344 residues. Other entries include 2LBP (344), 1EHT (332), 1OP4 (305), 1Q00 (256), 1G00 (253), 2LBP_1 (250), 2LBP_2 (227), 1Q00_B_1 (222), and 1G00 (220). Each entry has a corresponding domain diagram showing alignment with the query protein's domain structure.

VAST Summary - Netscape

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

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.

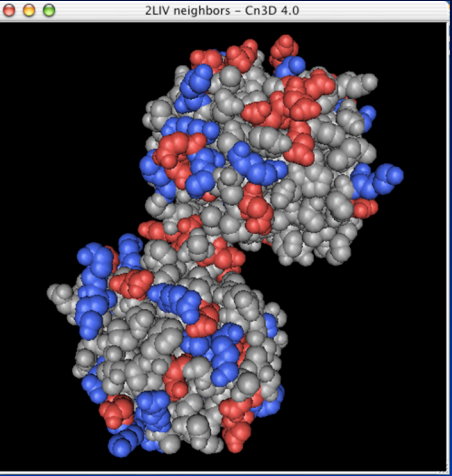
PDB	C	D	Ali.	Len.	SCORE	P-VAL	RMSD	%Id	MMDB	Date	Description
<input type="checkbox"/>	2LBP			344	39.8	10e-44.6	0.9	79.1	03/2001		Leucine-Binding Protein (LBP)
<input type="checkbox"/>	1USG	A		343	40.1	10e-42.4	2.0	79.0	01/2004		L-Leucine-Binding Protein, Apo Form
<input type="checkbox"/>	1JDP	B		310	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 ≤ 0.001
 and
 % Identity > 25
 over at least 20 residues

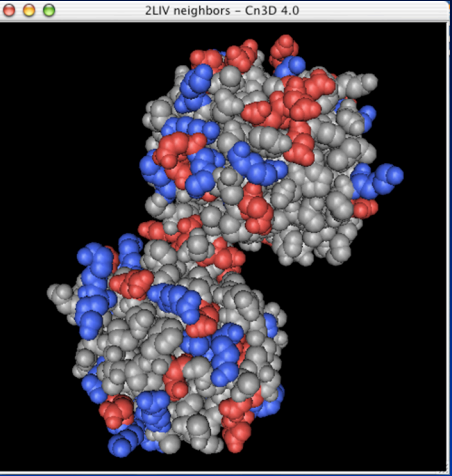
Read the descriptions!



Worms
Secondary Structure



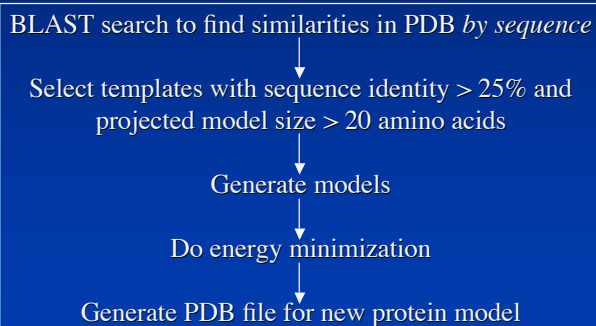
Rendering
Coloring



Spacefill
Charge

SWISS-MODEL

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



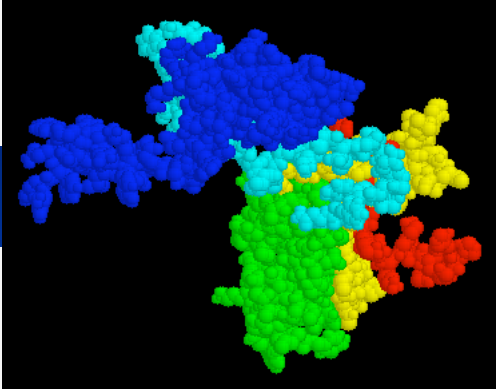
```

    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

```

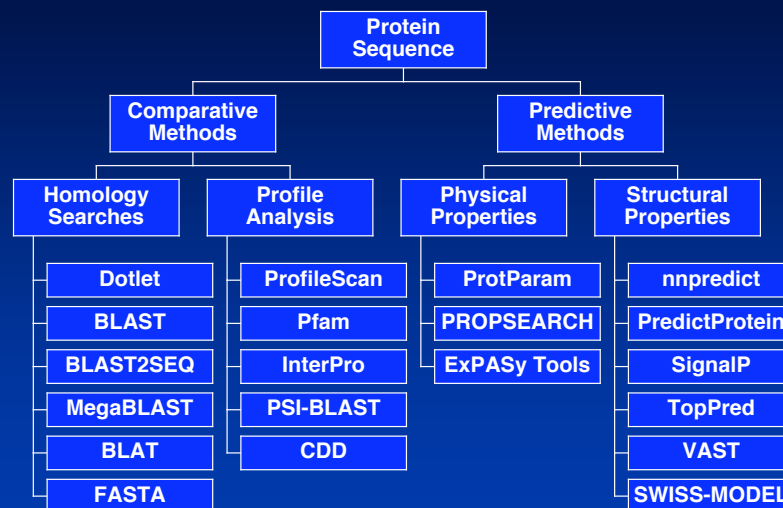
    ↗
    
```

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>

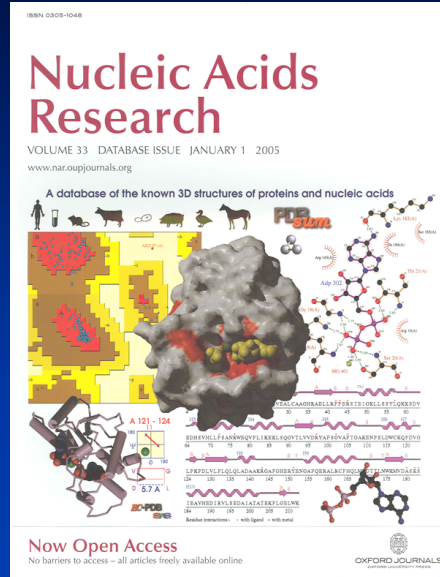


Protein Sequence Analysis

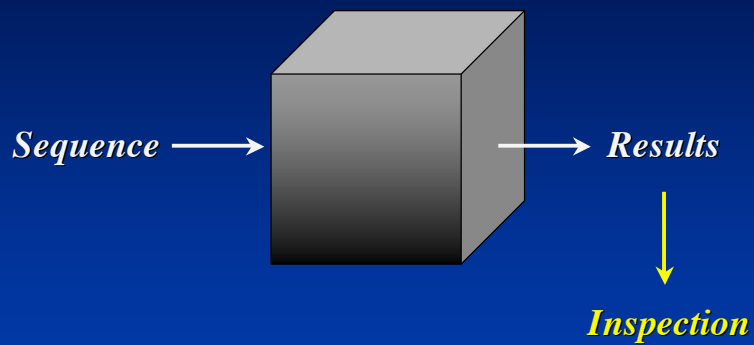


Annual NAR Database Issue

<http://nar.oupjournals.org>



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

