


NATIONAL HUMAN GENOME RESEARCH INSTITUTE Division of Intramural Research




**Current Topics in Genome Analysis
Spring 2008**

Week 3: Biological Sequence Analysis II

Andy Baxevanis, Ph.D.

U.S. DEPARTMENT OF HEALTH AND HUMAN SERVICES | NATIONAL INSTITUTES OF HEALTH | genome.gov/DIR



Overview

- Week 2
 - Similarity vs. Homology
 - Global vs. Local Alignments
 - Scoring Matrices
 - BLAST
 - BLAT
- Week 3
 - Profiles, Patterns, Motifs, and Domains
 - Structures: VAST, Cn3D, and *de novo* Prediction
 - Multiple Sequence Alignment



Sequence Comparisons

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



Profiles

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



Profile Construction

APHIIVATPG
 GCEIIVATPG
 GVEICIATPG
 GVDILIGTGG
 RPHIIVATPG
 KPHIIVATPG
 KVQLIATPG
 RPDIVIVATPG
 APHIIIVGTPG
 APHIIIVGTPG
 GCHVVIVATPG
 NQDIVVATPG

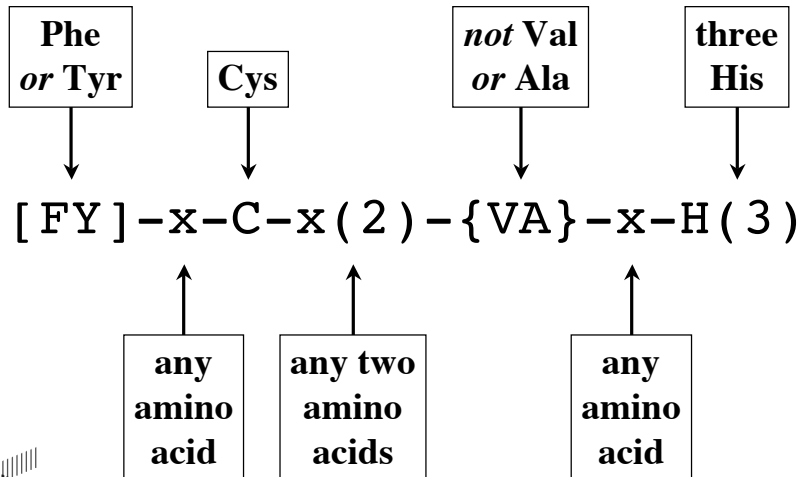
- Which residues are seen at each position?
- What is the frequency of observed residues?
- Which positions are conserved?
- Where can gaps be introduced?

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	9	10	9	0	12	10	0	0	0	0	0	0	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	10	9	10	9	0	12	10	0	0	0	0	0	0	23	2	-2	12	11	17	-31	-8	1
G	17	18	0	19	14	-22	31	0	-9	12	-15	-5	15	10	9	6	18	14	1	-15	-22	11



Patterns



Pfam

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



Pfam

- Pfam A
 - Based on *curated* multiple alignments (“seed alignment”)
 - Hidden Markov models (HMMs) used to find all detectable protein sequences belonging to the family
 - Given the method used to construct the alignments, hits are highly likely to be true positives
- Pfam B
 - Automatically generated from database searches
 - Deemed “lower quality”, but can be useful when no Pfam A family is identified



Sequences Used in Examples

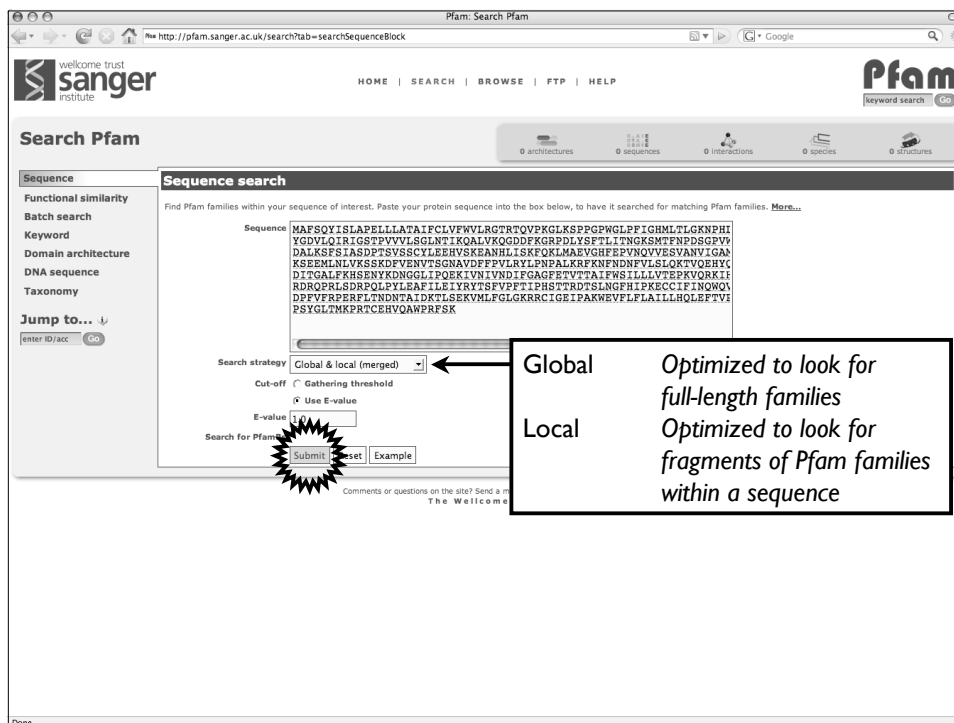
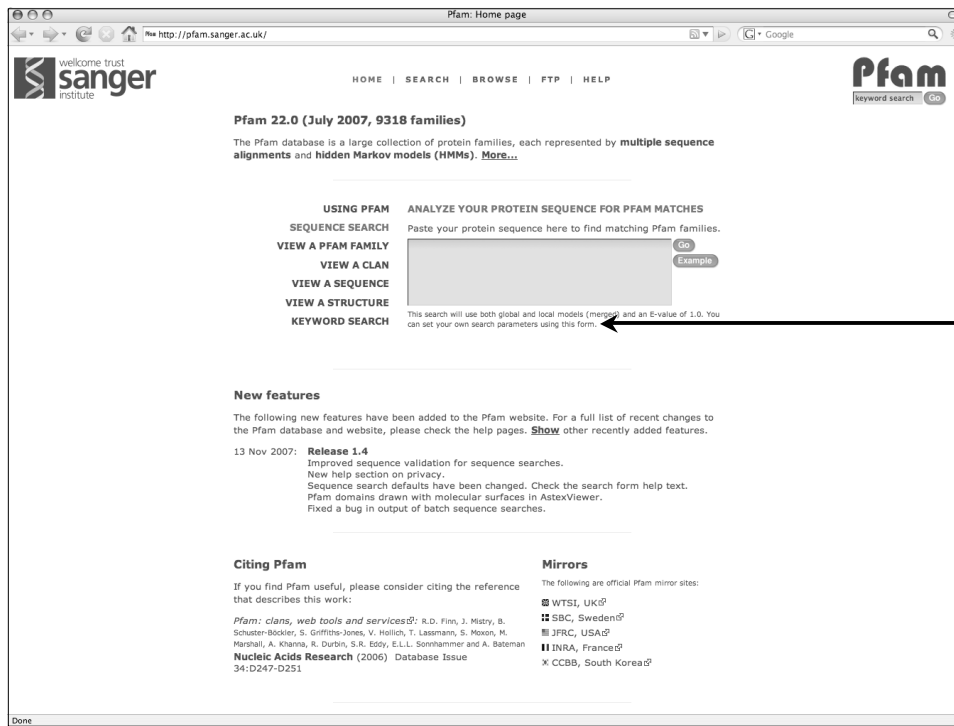
http://research.nhgri.nih.gov/teaching/seq_analysis.shtml



The screenshot shows a web browser displaying the NHGRI website. The page title is "Current Topics in Genome Analysis 2008". The main content area is titled "Weeks 2 and 3: Biological Sequence Analysis Protein and Nucleotide Sequences for Analysis". It contains several sections of sequence data:

- BLASTP**: A section with a query sequence and a long list of protein sequences.
- BLAST 2 Sequences**: A section with two sequences, one labeled "NP_038872.1 BOX-10 (Bbox eaglena)" and another labeled "NP_031311.1 sex determination region Y (Bbox eaglena)".
- MegaBLAST**: A section with a query sequence and a long list of DNA sequences.
- BLAT**: A section with a query sequence and a long list of DNA sequences.

The screenshot shows the Pfam website home page. The URL in the browser is <http://pfam.sanger.ac.uk>. The page features a navigation menu with "HOME", "SEARCH", "BROWSE", "FTP", and "HELP". A search bar is located in the top right corner. The main content area is titled "Pfam 22.0 (July 2007, 9318 families)" and includes a brief description of the database. Below this, there is a section titled "USING PFAM" with a list of links: "SEQUENCE SEARCH", "VIEW A PFAM FAMILY", "VIEW A CLAN", "VIEW A SEQUENCE", "VIEW A STRUCTURE", and "KEYWORD SEARCH". A black arrow points to the "VIEW A PFAM FAMILY" link. The page also includes a "New features" section with a list of updates and a "Citing Pfam" section with a list of references. The footer contains a "Mirrors" section with a list of official Pfam mirror sites.



Pfam: Sequence search results

wellcome trust sanger
 HOME | SEARCH | BROWSE | FTP | HELP

Sequence search results

We found 1 Pfam-A match to your search sequence. You did not choose to search for Pfam-B matches. The Pfam graphic below shows the arrangement of the domains on your search sequence. Clicking on any of the domains will take you to a page of information about that domain.

Below showing the details of the domains that were found. Rows containing significant hits are highlighted. Hits which do not start and end at the end points of the matching HMM are also highlighted.

For Pfam-A hits we show the alignments between your search sequence and the matching HMM. You can show individual alignments by clicking on the "Show" button in each row of the result table, or you can show all alignments using the links above the table. You can bookmark this page and return to it later, but please note that old results will be removed after **one week**. Return to the search form to look for Pfam domains on a new sequence.

Pfam-A Matches

Show or hide all alignments.

Pfam-A	Description	Entry type	Sequence		HMM		Bits score	E-value	Alignment mode	Show/hide alignment
			Start	End	From	To				
p450	Cytochrome P450	Domain	41	506	1	504	367.2	9.1e-113	fs	Show

#HM -->Ppgeptp1p1Gnl1q1grgrf1kdlhsvtklkkkYQpifllylQpkvPv1agpeavkveLkkbaefsggdeafytlkpgf1ghvifang.GerVqlRcfltpftrfmgkllk.....#fsgpgeardvevlkktagppg
 #MATCH Pppg +1P++Q++lg +nh +tkl++YG+++G++PvV1ag+ +k +L+kq++Eg+d +y+++ +gk + E +G+ W Rr+ ++sf + +++ ++ +e+ +ea+ L+ k+k +e g
 #SEQ PPGPGLPFIQMLTG-----KHPFLSITKLBQOYGVQLGIRIGSTPVV1VLSGLNTKQALVQKQDDFKGRPD---LVSFLLITNGKMTFNPDeGFWAARRRAGDALKSF61-ASDptvavscYL8HVSKKANHLISKFKRLMAYG

Comments or questions on the site? Send a mail to pfam-help@sanger.ac.uk
 The Wellcome Trust

Pfam: Family: p450 (PF00067)

70 architectures 8758 sequences 2 interactions 1045 species 148 structures

Family: p450 (PF00067)

Summary

Cytochrome P450 [Add annotation](#)

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

Literature references

- Werk-Richardt D, Feyereisen R, Genome Biol 2000;1:REVIEWS3083. Cytochromes P450: a success story. PUBMED:11178272
- Nebert DW, Gonzalez FJ, Annu Rev Biochem 1987;56:945-993. P450 genes: structure, evolution, and regulation. PUBMED:3304150
- Guengerich FP, J Biol Chem 1991;266:10019-10022. Reactions and significance of cytochrome P-450 enzymes. PUBMED:2037557
- Nelson DR, Kamabaki T, Waxman DJ, Guengerich FP, Estabrook RW, Feyereisen R, Gonzalez FJ, Coon MJ, Gunsalus IC, Gotoh O, et al., DNA Cell Biol 1993;12:1-51. The P450 superfamily: update on new sequences, gene mapping, accession numbers, early trivial names of enzymes, and nomenclature. PUBMED:7678949
- Deglyarenko KN, Archakov AI, FEBS Lett 1993;332:1-8. Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. PUBMED:8405421
- Graham-Lorence S, Amarnah B, White RE, Peterson JA, Simpson ER, Protein Sci 1995;4:1065-1080. A three-dimensional model of aromatase cytochrome P450. PUBMED:7549871

Interpro entry IPR001128

Cytochrome P450 enzymes are a superfamily of haem-containing mono-oxygenases that are found in all kingdoms of life, and which show extraordinary diversity in their reaction chemistry. In mammals, these proteins are found primarily in microsomes of hepatocytes and other cell types, where they oxidise steroids, fatty acids and xenobiotics, and are important for the detoxification and clearance of various compounds, as well as for hormone synthesis and breakdown, cholesterol synthesis and vitamin D metabolism. In plants, these proteins are important for the biosynthesis of several compounds such as hormones, defensive compounds and fatty acids. In bacteria, they are important for several metabolic processes, such as the biosynthesis of antibiotic erythromycin in *Saccharopolyspora erythraea*.

Cytochrome P450 enzymes use haem to oxidise their substrates, using protons derived from NADH or NADPH to split the oxygen so a single atom can be added to a substrate. They also require electrons, which they receive from a variety of redox partners. In certain cases, cytochrome P450 can be fused to its redox partner to produce a bi-functional protein, such as with P450BM-3 from *Bacillus megaterium* PUBMED:13023115, which has haem and flavin domains.

Organisms produce many different cytochrome P450 enzymes (at least 58 in humans), which together with alternative splicing can provide a wide array of enzymes with different substrate and tissue specificities. Individual cytochrome P450 proteins follow the nomenclature: CYP, followed by a number (family), then a letter (subfamily), and another number (protein); e.g. CYP3A4 is the fourth protein in family 3, subfamily A. In general, family members should share >40% identity, while subfamily members should share >55% identity.

Cytochrome P450 proteins can also be grouped by two different schemes. One scheme was based on a taxonomic split: class I (prokaryotic/mitochondrial) and class II (eukaryotic microsomes). The other scheme was based on the number of components in the system: class B (2-components) and class E (2-components). These classes merge to a certain degree. Most prokaryotes and mitochondria (and fungal CYP55) have 3-component systems (class I/class B) - a FAD-containing flavoprotein (NAD(P)H-dependent reductase), an iron-sulphur protein and P450. Most eukaryotic microsomes have 2-component systems (class II/class E) - NAD(P)H reductase (FAD and FMN-containing flavoprotein) and P450. There are exceptions to this scheme, such as 1-component systems that resemble class E enzymes PUBMED:16042601, PUBMED:15128046, PUBMED:8637843. The class E enzymes can be further subdivided into five sequence clusters, groups I-V, each of which may contain more than one cytochrome P450 family (eg, CYP1 and CYP2 are both found in group I). The divergence of the cytochrome P450 superfamily into B- and E-classes, and further divergence into stable clusters within the E-class, appears to be very ancient, occurring before the appearance of eukaryotes.

More information about these proteins can be found at Protein of the Month: Cytochrome P450 PUBMED..

Gene Ontology

Molecular function	heme binding (GO:0020037)
Molecular function	iron ion binding (GO:0055059)
Biological process	electron transport (GO:0006119)
Molecular function	monooxygenase activity (GO:0004497)

External database links

InterPro: IPR001128 Cytochrome P450

EMBL-EBI | EBI Home | All Databases | Enter Text Here | Go | Reset | Give us feedback

Jump to: InterProScan | Databases | Documentation | FTP site | Help | Advanced search

InterPro: IPR001128 Cytochrome P450

Protein matches

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

Architectures
 Accession List

UniProtKB Matches: 10695 proteins

Accession: IPR001128 Cyt_P450

Type: Family

Database	ID	Name
Gene3D	G3DSA:1.10.630.10	Cyt_P450
Pfam	PF00067	p450
PRINTS	PR00385	P450
PROSITE pattern	PS00086	CYTOCHROME_P450 7851
PANTHER	PTHR19383	Cyt_P450 10317
SuperFamily	SSF48264	Cytochrome_P450 10420

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

InterPro Relationships

Children

- IPR002397 Cytochrome P450, B-class
- IPR002401 Cytochrome P450, mitochondrial
- IPR002402 Cytochrome P450, E-class, group I
- IPR02403 Cytochrome P450, E-class, group II
- IPR02404 Cytochrome P450, E-class, group IV

GO Term annotation

Process: GO:0008118 electron transport
 GO:0004497 monoxygenase activity
 Function: GO:0005506 iron ion binding
 GO:0003037 heme binding

InterPro annotation

Abstract: Cytochrome P450 enzymes are a superfamily of haem-containing mono-oxygenases that are found in all kingdoms of life, and which show extraordinary diversity in their reaction chemistry. In mammals, these proteins are found primarily in microsomes of hepatocytes and other cell types, where they oxidise steroids, fatty acids and xenobiotics, and are important for the detoxification and clearance of various compounds, as well as for hormone synthesis and breakdown, cholesterol synthesis and vitamin D metabolism. In plants, these proteins are important for the biosynthesis of several compounds such as hormones, defensive compounds and fatty acids. In bacteria, they are important for several metabolic processes, such as the biosynthesis of antibiotic erythromycin in *Saccharopolyspora erythraea*.

Cytochrome P450 enzymes use haem to oxidise their substrates, using protons derived from NADH or NADPH to split the oxygen so a single atom can be added to a substrate. They also require electrons, which they receive from a variety of redox partners. In certain cases, cytochrome P450 can be fused to its redox partner to produce a bi-functional protein, such as with P450BM-3 from *Bacillus megaterium* [1], which has haem and flavin domains.

Organisms produce many different cytochrome P450 enzymes (at least 58 in humans), which together with alternative splicing can provide a wide array of enzymes with different substrate and tissue specificities. Individual cytochrome P450 proteins follow the nomenclature: CYP, followed by a number (family), then a letter (subfamily), and another number (protein); e.g. CYP3A4 is the fourth protein in family 3, subfamily A. In general, family members should share >40% identity, while subfamily members should share >55% identity.

Cytochrome P450 proteins can also be grouped by two different schemes. One scheme was based on a taxonomic split: class I (prokaryotic/mitochondrial) and class II (eukaryotic microsomes). The other scheme was based on the number of components in the system: class B (3-components) and class E (2-components). These classes merge to a certain degree. Most prokaryotes and mitochondria (and fungal CYP5) have 3-component systems (class I/class B) - a FAD-containing flavoprotein (NAD(P)H-dependent reductase), an

InterPro: IPR001128 Cytochrome P450

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

Database links: COME: PR000236
 PANTHER: PF00067
 PROSITE doc: PDOC00081
 Enzyme: EC:1.14
 MSDsite: PS00086

Taxonomic coverage

Group	Count
Unclassified	2
Virus	14
Archaea	2156
Bacteria	66
Cyanobacteria	66
Synechocystis PCC 6803	1
Oryza sativa (Rice)	3252
Arabidopsis thaliana	422
Green Plants	3132
Plastid Group	3214
Human	3214
Mouse	3214
Other Eukaryotes	54
Eukaryota	8542

Overlapping InterPro entries

IPR001128	Numbers of overlapping proteins	Average numbers of overlapping amino acids
IPR002397	9130 1566 0	N/A
IPR002399	10664 32 0	N/A
IPR002401	4250 6446 0	N/A
IPR002402	10596 100 0	N/A
IPR002403	9405 1291 0	N/A
IPR002404	10681 15 0	N/A
IPR002405	10483 213 0	N/A
IPR002406	10472 224 0	N/A
IPR002407	10619 77 0	N/A
IPR002408	10643 53 0	N/A
IPR002409	10553 143 0	N/A
IPR002410	10657 39 0	N/A
IPR002411	10658 38 0	N/A
IPR002412	10555 141 0	N/A

Example proteins

Q00158 Cytochrome P450 3A25 (EC 1.14.14.1) (CYP11A25)

Q17824 Putative cytochrome P450 cyp-13B1 (EC 1.14.-.-)

O46051 Probable cytochrome P450 4d14 (EC 1.14.-.-) (CYP1VD14)

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

InterPro: IPR001128 Cytochrome P450

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

Example proteins

- Q09158 Cytochrome P450 3A26 (EC 1.14.14.1) (CYP3A26)
- Q17824 Putative cytochrome P450 cyp-13B1 (EC 1.14.-.)
- Q46051 Probable cytochrome P450 4d14 (EC 1.14.-.) (CYP1VD14)
- P05177 Cytochrome P450 1A2 (EC 1.14.14.1) (CYP1A2) (P450-P3) (P3)450 (P450 4)
- P10614 Cytochrome P450 51 (EC 1.14.13.70) (CYP1) (P450-L1A1) (Sterol 14-alpha demethylase) (Lanosterol 14-alpha demethylase) (P450-14DM)

More proteins

- IPR001128 Cytochrome P450
- IPR008086 Cytochrome P450, E-class, group I, CYP1
- IPR008072 Cytochrome P450, E-class, CYP3A
- IPR002403 Cytochrome P450, E-class, group IV
- IPR002401 Cytochrome P450, E-class, group I
- ModBase
- SWISS-MODEL
- PDB Chain

Publications

- Munro A.W., Girvan H.M., McLean K.J. Cytochrome P450-reductase fusion enzymes. *Biochim. Biophys. Acta* 2006 [PubMed: 17023115]
- McLean K.J., Sabri M., Marshall K.R., Lawson R.J., Lewis D.G., Clift D., Balding P.R., Dunford A.J., Warman A.J., McVey J.P., Quinn A.M., Sutcliffe M.J., Scrutton N.S., Munro A.W. Biodiversity of cytochrome P450 reductase systems. *Biochem. Soc. Trans.* 33 796-801 2005 [PubMed: 16042901]
- Neilson D.R., Zaidin D.G., Hoffman S.M., Mallik S.J., Wein H.M., Nebert D.W. Comparison of cytochrome P450 (CYP) genes from the mouse and human genomes, including nomenclature recommendations for genes, pseudogenes and alternative-splice variants. *Pharmacogenetics* 14 1-16 2004 [PubMed: 15128046]
- Dedyer R.N.K. Structural domains of P450-containing monooxygenase systems. *Protein Eng.* 8 737-47 1996 [PubMed: 8837843]
- McDowell J. Protein of the Month - Cytochrome P450. 2006 http://www.ebi.ac.uk/interpro/iptm/2006_10/Page1.htm

Additional Reading

- Chen R., Fuhndou S., Su F., Zhang L., Takaya N., Shoun H. Structural evidence for direct histidine transfer from NADH to cytochrome P450d. http://www.ebi.ac.uk/interpro/iptm/2006_10/Page1.htm

Current Protocols in Bioinformatics

CPBI Unit 2.5 Pfam

CPBI Unit 2.7 InterPro

Identifying Protein Domains with the Pfam Database

UNIT 2.5

Hundreds of thousands of protein sequences are now known and the display of data shows an upsurge in diversity. The sequence analysis of proteins may seem like a peripheral task. However, the majority of protein sequences appear to fall into a few thousand protein families (Chothia, 1992). Now these families are representative of proteins at the domain level. These domains are discrete structural units that are functionally linked and define protein function. Pfam is a database of protein domain families (Coleman et al., 1997; Bateman et al., 2002), which is a highly recognizable multiple-sequence alignment and profile hidden Markov model (HMM) database. In addition, each family has associated annotations, literature references, and links to other databases. The entries in Pfam are available via the Web and in flat-file format.

The use of Pfam by molecular biologists as a protein information resource and analysis tool is widespread. Current sequencing projects, including human genome, the International Protein Consortium (IPC) for large-scale functional annotation of proteins, and other protein genome projects, all utilize a single Pfam as a functional domain resource. Pfam is used for techniques such as secondary structure prediction, fold recognition, and phylogenetic analysis, and for site and mutation design. This unit contains detailed information on how to access and utilize the information present in the Pfam database, namely the families, multiple alignments, and annotations. Details on using Pfam include the two Basic Protocols and locally the Alternative Protocol (if you are interested). The unit also includes a brief discussion of analyzing genomes (DNA) with Pfam, the Alternative Protocol 2. Information on underlying families is also presented (see Appendix of the end of this unit).

ANALYZING A PROTEIN SEQUENCE WITH Pfam VIA THE WEB

A primary use of the Pfam database is to determine what domains are in a protein of interest. This section describes several approaches for carrying out this analysis. In a typical identification you first enter a protein or nucleotide sequence into the search box. In a typical identification you first enter a protein or nucleotide sequence into the search box. In a typical identification you first enter a protein or nucleotide sequence into the search box.

NECESSARY REVISIONS

Hardware

Workstation with network connection

Software

Internet-capable browser (e.g., Netscape 4.0, Internet Explorer 4.0)

Files

Protein sequence of interest in appropriate format (e.g., FASTA, Genbank, etc.)

The complete sequence used for the analysis is available at the Current Protocols Web site (<http://www.currentprotocols.com/>).

Contributed by Robert Fux, Santa Clotilde, Spain, and John Bateman, Current Protocols in Bioinformatics (UNIT 2.5). © 2006 Copyright 2006 by John Wiley & Sons, Inc.



The InterPro Database and Tools for Protein Domain Analysis

UNIT 2.7

InterPro (Apicler et al., 2001a) is an integrated documentation resource of protein families, domains, and functional sites that ratifies data from the major protein signature databases: Pfam (Bateman et al., 2002), PROSITE (Gasteiger et al., 2002), PRINTS (Atkinson et al., 2002), ProDom (Coppe et al., 2000), SMART (Hansen et al., 2002), IPD (Sapich et al., 2000) and PIRAM (Jiang et al., 2001). These databases, summarized in Table 2.7.1 along with other useful sites (also see InterPro Significance), are a source of methods for sequence identification, including regular expressions, profile, signature, matrix, hidden Markov models (HMMs), and sequence-clustering algorithms. In InterPro, signatures describing the same domain, family, motif, or post-translational modification are grouped into single entries with unique accession numbers. Entries that are subsets of others are annotated by parent-child (family and subfamily) or contained-by (domain and superfamily) relationships. The associated reference package, InterProScan (Bateman and Apicler, 2001), combines the search methods from each of the member databases into one package and provides an output with all results in a single format, namely HTML, XML, or text. The software has been used for annotation of numerous genomes (see the Introduction, especially Section The International Human Genome Consortium, 2001).

InterPro provides a one-stop shop for protein sequence classification. Using the user-friendly interface, users can search for protein signatures, classification, using the user-friendly interface. The high-quality search results provide useful information on each domain entry. The protein family, subfamily, and motifs are supplied in the UniProt (Bateman, 2001), The Gene Ontology Consortium, 2001) facilitating automatic mapping of large sets of proteins (GO). The four search methods (profile, signature, matrix, HMM) are available in InterProScan via a Web tool Basic Protocol 2. In addition, details on harnessing InterPro for an alternative, a workflow requiring how to install InterProScan for local running is also described (see Alternative Protocol 2). In addition, details on harnessing InterPro families and domains of interest using the InterPro Web server (see Basic Protocol 2) and sequence retrieval system (see Alternative Protocols 3 and 4) are provided. Where necessary, the protocol is illustrated with examples.

PROTEIN SEQUENCE CLASSIFICATION USING INTERPROSCAN VIA THE INTERNET

A core protein sequence may be run against all the protein signatures in InterPro to identify protein domains or domains within protein families. The sequence-based search, InterProScan, may be provided by the InterPro database, including Searchable for PROSITE patterns, files for PROSITE profiles (Bateman et al., 2002), Intermap (Eddy, 1998), AnnotateProteinData for Pfam (Bateman et al., 2002), SMART (Pfam, 2002), and the InterProScan (Bateman, 2001) interface. InterProScan (Bateman, 2001) is available via a Web tool Basic Protocol 2. In addition, details on harnessing InterPro for an alternative, a workflow requiring how to install InterProScan for local running is also described (see Alternative Protocol 2). In addition, details on harnessing InterPro families and domains of interest using the InterPro Web server (see Basic Protocol 2) and sequence retrieval system (see Alternative Protocols 3 and 4) are provided. Where necessary, the protocol is illustrated with examples.

NECESSARY REVISIONS

Hardware

Workstation with network connection

Software

Internet-capable browser (e.g., Netscape 4.0, Internet Explorer 4.0)

Files

Protein sequence of interest in appropriate format (e.g., FASTA, Genbank, etc.)

The complete sequence used for the analysis is available at the Current Protocols Web site (<http://www.currentprotocols.com/>).

Contributed by Mark J. Martin and John Bateman, Current Protocols in Bioinformatics (UNIT 2.7). © 2006 Copyright 2006 by John Wiley & Sons, Inc.

<http://nihlibrary.nih.gov>
Search "Online Journals" for "Current Protocols in Bioinformatics"

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 Pfam or InterPro



NCBI Conserved Domain Database (CDD)

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

Domains SH3 SH2 SH3 SH3

Search across Entrez databases:

CDTree ^{NEW} **A Conserved Domain Database and Search Service, v2.13**

CDD help

NCBI Handbook

CD-Search

CDART

Pfam

SMART

COG

Find CDS

In Entrez:

Structure

MMDB

Cn3D

VAST

Research

CDD FTP site

Last Revised 11/15/07

Submit Query Search Database: CDD v2.13 - 24083 PSSMs

Enter a **Protein** query as Accession, GI, or Sequence in FASTA format:

```
>NP_005206.1 deleted in colorectal carcinoma [Homo sapiens]
MENSLRGVVWPKLAFVLFQASLLSAHLQVTFQIKAPTALRFLSEPSDAVTMRGCVLLDCSAESDRGVP
VIKWRKDGIIHALGMDERKQQLSNGSLLIQNLHSHRHHKPDGLYCEASLDGSGSIIISRTAKVAVAGPL
RFLSQTESVIAEMQDVIILKCEVIGEPMTIHWKQQDLTFPGDSRVVLFSGALQISRLQFGDIGIY
```

Read about the FASTA format description. [Click here for advanced options.](#)

Computational biologists define conserved domains based on recurring sequence patterns or motifs. The un-curated section of CDD contains domains imported from SMART, Pfam and COGs. The source databases also provide descriptions and links to citations. Because conserved domains correspond to compact structural units, CDS are linked to 3D structure when possible. The NCBI-curated section of CDD attempts to group ancient domains related by common descent into family hierarchies.

To identify conserved domains in a protein sequence, the CD-Search service uses the reverse position-specific BLAST algorithm. The query sequence is compared to a position-specific score matrix prepared from the underlying conserved domain alignment. Hits may be displayed as a pairwise alignments of the query sequence with representative domain sequences, or as multiple alignments. CD-Search now is run by default in parallel with protein BLAST searches. Although the user waits for the BLAST queue to further process the request, the domain architecture of the query may already be studied.

Run CDART, the Conserved Domain Architecture Retrieval Tool, to search for proteins with similar domain architectures. CDART uses pre-computed CD-Search results to quickly identify proteins with a set of domains similar to that of the query.

Read more about CDD:

Marchler-Bauer A, Anderson JB, Cherukuri PF, DeWeese-Scott C, Geer LY, Gwade M, He S, Hurwitz DI, Jackson JD, Ke Z, Lanczycki C, 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. CDD: a Conserved Domain Database for protein classification. *Nucleic Acids Res.* 2005;33 Database Issue:D192-6. [Abstract] [Full Text]

Marchler-Bauer A, Bryant SH. CD-Search: protein domain annotations on the fly. *Nucleic Acids Res.* 2004;32(Web Server Issue):W327-31. [Abstract] [Full Text]

Marchler-Bauer A, 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 GH, Mazumder B, Niolekaya AN, Panchenko AR, Rao BS, Shoemaker BA, Simonyan V, Song JS, Thissen PA, Vasudevan S, Wang Y, Yamashita RA, Yin JJ, Bryant SH. CDD: a curated Entrez database of conserved domain alignments. *Nucleic Acids Res.* 2003;31:383-7. [Abstract] [Full Text][Terms]

Marchler-Bauer A, Panchenko AR, Shoemaker BA, Thissen PA, Geer LY, and Bryant SH. CDD: a database of conserved domain alignments with links to domain three-dimensional structure. *Nucleic Acids Res.* 2002;30:281-3. [Abstract] [Full Text]

Citing CDD: Marchler-Bauer A, Anderson JB, Cherukuri PF, DeWeese-Scott C, Geer LY, Gwade M, He S, Hurwitz DI, Jackson JD, Ke Z, Lanczycki CJ, Liebert CA, Liu C, Lu F, Marchler GH, Mullokandov M, Shoemaker BA. [Done](#)

The screenshot shows the NCBI Conserved Domain Search results page. At the top, there's a navigation bar with 'NewSearch', 'PubMed', 'Nucleotide', 'Protein', 'Structure', 'Taxonomy', and 'Help'. The 'Query sequence' is '[(local sequence)|cd|1]'. Below this is a diagram of a protein sequence with domains highlighted: Icam, FN3, FN3, FN3, FN3, FN3, and Neogenin_C. The 'Descriptions' table lists several entries, with the first entry for 'Icam' highlighted in blue. An arrow points from the left margin to this first entry.

Title	Pssmid	Multi-Dom	E-value
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	3e-15
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	8e-13
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	1e-12
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	9e-12
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	3e-11
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	3e-10
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	6e-09
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	8e-07
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	1e-06
hpfam06583, Neogenin_C, Neogenin C-terminus. This family represents the C-terminus of e...	87114	No	3e-96
hpfam07686, V-set, Immunoglobulin V-set domain. This domain is found in antibodies as w...	87333	Yes	2e-04

This screenshot shows the 'Full Result' view for the Icam domain. It includes a sequence alignment diagram with domains: Icam, FN3, FN3, FN3, FN3, FN3, and Neogenin_C. The 'Descriptions' table is expanded to show details for the first entry (PSSMID 28983):

Title: hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...

Pssmid: 28983

Multi-Dom: No

E-value: 3e-15

CD Length: 89, **Pct. Aligned:** 100, **Bit Score:** 79.775909, **E-value:** 3e-15

```

1          330  *P*P*F*L*N*E*P*S*Y*E*S*E*D*E*E*F*E*C*T*V*S*G*K*V*P*V*V*N*N*K*G*Q*W*V*P*S*P*--Y*P*Q*V*G*G*S*N*L*R*L*L*V*V*K*S*D*E*G*F*Y*Q*V*A*E*N*E* 407
cd00931    1  P*P*T*Q*K*P*P*P*D*V*V*G*G*E*D*V*T*L*E*C*R*A*S*G*N*P*P*P*I*T*W*L*K*N*K*P*L*S*L*L*D*P*Y*T*V*I*D*N*N*G*T*L*I*T*S*N*V*T*R*E*D*A*G*T*Y*T*C*V*A*T*S*G  80
    
```

The table below the alignment lists other related domains:

hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	8e-13
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	1e-12
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	9e-12
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	3e-11
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	3e-10
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	6e-09
hcd00063, FN3, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	8e-07
hcd00931, ICam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	1e-06
hpfam06583, Neogenin_C, Neogenin C-terminus. This family represents the C-terminus of e...	87114	No	3e-96
hpfam07686, V-set, Immunoglobulin V-set domain. This domain is found in antibodies as w...	87333	Yes	2e-04

NCBI CDD cd00931

http://www.ncbi.nlm.nih.gov/Structure/cdd/cddsrv.cgi?uid=cd00931

Conserved Domains

cd00931: IGCam

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

Links

Source: Smart
 Taxonomy: Bilateria
 PubMed: 3 links
 Protein: Related Protein
 Related Structure
 Architectures
 Representatives
 Related CDS: 8 links

Statistics

PSSM-Id: 28983
 View PSSM: cd00931
 Aligned: 41 rows
 Status: curated CD
 Created: 1-Nov-2000
 Updated: 10-Jan-2006

Structure

Structure View

Program: Cn3D
 Drawing: All Atoms
 Aligned Rows: up to 10
 Download Cn3D

Hierarchy

Interactive Display
 Display: cd00931 branch
 Download CDTree

FGF/FGF-Receptor

Feature 1: FGF/FGF-Receptor Interaction

Evidence:

- Structure: 1EV2: Receptor domain (chain F) contacts FGF (chain B)

View structure with Cn3D

Download Ch3D for Viewing 3D Structure

Scroll to Sequence Alignment Display

cd00931 is part of a hierarchy of related CD models.
 Use the graphical representation to navigate this hierarchy.

cd00931 Sequence Cluster

Zoom Out Detailed View

Sub-family Hierarchy

Interactive Display with CDTree

- cd00936 IGC
- cd00938 IGC
- cd00939 IGV
- cd00931 IGCam

NCBI CDD cd00931

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

Drawing: All Atoms
 Aligned Rows: up to 10
 Download Cn3D

Hierarchy

Interactive Display
 Display: cd00931 branch
 Download CDTree

Other Related Conserved Domains

Start#1816 Start#1818 Start#1819 Start#2015 Start#2017 Start#7073

Sequence Alignment

Reformat Format: Hypertext Row Display: up to 10 Color Bks: 2.0 bit Type Selection: Top listed sequences

Feature 1

IB1H_A 309 PKYeqK---pkVIVVkgp---gDTIPCKVGLpaFNWVSHnakp---lsgRATVTD---SGLVIKGVKs-gDK 372
 IC56_A 295 PMLdV---LrTEADlg---dLRNSCVASGkprPAVRWRldgqpl---asqRlEVSG---GELRFSLVl-eDS 358
 IF6C_A 42 PTFtk---LlMDVYkg---aAATFCVYDygpEVWFKdmgp---kssLTIQYDaaqgCGLTSECVK---GDD 110
 LV2_P 8 PFWntokokmrLRAVPa---kTVFKPAGGmpPTDRGLKngkeLkqehrlgCYVVRng---hMLLMSVVP-sDK 80
 IC56_A 5 PVVteq---paHTLFPgaaseKVTLZCARAappATTRKngtLk-kgpdrKTRLVA---GDLVSNPVAkda 74
 INCI_A 2 VWFpa---pQRFKkg---dNVVYDYNlspRTIKWkgrdri-LkkoVETId---NlQIRKIKs-tDE 68
 g1 20455467 332 PFWlgk---pqsHLXGp---eTARLDQVQcpgpEVWIRngmsLkvnkdgKYRIEQ---SGLLNSVQp-sDT 398
 g1 14286138 339 PTFk---paLRARAd---EVVFCRAkqepKLSWlHnglqgptqgRVTVD---NTRILNVA-gPT 405
 g1 3334268 209 PAFIm---pgSPNAAsrgpENTSCASGepPALSNRngkll---epmKILIKG---NTELTVNIIIn-sDG 276
 g1 14286138 432 PTLbaa---paVSTVDg---rNVTIKRVNgspkFLVWRlRaanw---ltgRTRVQA---GDLIQDVTf-sDA 495

Feature 1

IB1H_A 373 GYGCRATne---hgDKYF-ETLQ 393
 IC56_A 359 GMYQVARKk---hgTVVA-SAEIT 379
 IF6C_A 111 AKYTCVAVns---lGEAC-SAEEL 131
 LV2_P 81 GNYTCVVENs---yGSIMH-TYELD 101
 IC56_A 75 GSYCVAVTns---rYTVVS-SAGLR 96
 INCI_A 69 GNYTCVSRLLlscqgINPK-DQIT 91
 g1 20455467 399 NVTQCARng---hGLL-NATY 419
 g1 14286138 406 GNYTCVAVns---lQVY-DYELM 426
 g1 3334268 277 GPVCRATK---agEDEK-QALFQ 297
 g1 14286138 496 GKYTCYQnk---fGEIQA-DGSLV 516

Citing CDD

Wanchler Bauer A, Anderson JB, Derynshire MK, DeLorenzo Scott C, Gonzalez NB, Gwartz M, Hsu L, He S, Hummel D, Jackson JD, Ke Z, Klybi D, Lenczycki C, Liebert CA, Liu C, Lu F, Lu S, Marcher GH, Mukandov M, Song JS, Thakki N, Yamashita RA, Yin JJ, Zhang D, Bryant SH. (2007) CDD: a conserved domain database for interactive domain family analysis. *Nucleic Acids Res.* 35: D237-43

Disclaimer: Privacy statement | Accessibility

NCBI Conserved Domain Search

Query sequence: [(local sequence)|cd|1]

Concise Result Full Result Show Search Information

Click on the colored bar for a conserved domain to view your query sequence within the multiple sequence alignment for that domain. To see only the sequences used to generate the domain, click on its PSSMID in the tabular summary.

Descriptions

Title	Pssmid	Multi-Dom	E-value
h:cd00931, IGcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	3e-15
cd00931, IGcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members are components of neural cell adhesion molecules (N-CAM L1), Fasciclin II and the insect immune protein Hemolin. The subfamily also includes receptor domains such as the extracellular ligand binding domain of Fibroblast Growth Factor Receptor 2. Members are phylogenetically diverse, occurring throughout metazoa, and are not components of the adaptive immune system molecules found in jawed vertebrates. A predominant feature of most Ig domains is a disulfide bridge connecting 2 beta-sheets with a Trp packing against the disulfide bond..			
CD Length: 89, Pct. Aligned: 100, Bit Score: 79.775909, E-value: 3e-15			
<pre> 10 20 30 40 50 60 70 80 1 PPWFLNEDSNLYAYESMDIEPECTVSGKRVPTVNMKNGDVIIPSD--YFQIVGGSNLRILGVVRSDEGFVQVAENEG 407 cd00931 1 PTFQKPPPTVAVGGEDVLECRASGNPPTITWLKNGKPLSLDgYrTVLDNNGTLTISNVTKEDAGTYTCVATNSAG 80 1 408 NAGTSAQLI 416 cd00931 81 GASASARLT 89 </pre>			
h:cd00063, FNS, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	9e-13
h:cd00063, FNS, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	1e-12
h:cd00931, IGcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	9e-12
h:cd00063, FNS, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	3e-11
h:cd00931, IGcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	3e-10
h:cd00063, FNS, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	6e-09
h:cd00063, FNS, Fibronectin type 3 domain; One of three types of internal repeats found ...	28945	No	9e-07
h:cd00931, IGcam, Immunoglobulin domain cell adhesion molecule (cam) subfamily; members ...	28983	No	1e-06
h:pfam06583, Neogenin_C, Neogenin C-terminus. This family represents the C-terminus of e...	87114	No	3e-96
h:pfam07686, V-set, Immunoglobulin V-set domain. This domain is found in antibodies as w...	87333	Yes	2e-04

Search for similar domain architectures

NCBI DART

CDART: Conserved Domain Architecture Retrieval Tool

About CDART

Query: I-set, Neogenin_C, FNS

Similar domain architectures:

- 99 Sequences: Caldesmon, neogenin, neogenin 1
- 6 Sequences: neural cell adhesion
- 2 Sequences: G33318-ac, lexafogf
- 30, 315479 Sequences: GalT, Eutheria, PROCT60, alu1a1
- 2 Sequences: EGF, LaminA/Progerin, Variable Kinase, VEGF_Like
- 2 Sequences: Hsaurin, Hsaurin
- 3 Sequences: Putrescine, tyrosine kinase re
- 2 Sequences: Eutheria, LPRCT
- 2 Sequences: tyrosine kinase re

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

Subset by Taxonomy

Subset by selected domains:

- cd00063 Fibronectin type 3 domain; One of three types of ... smart00060 pfam00041 includes: cd00160 Guanine nucleotide exchange factor for Rho/Rac/Cd... includes: smart00325 pfam00621
- cd00180 Serine/Threonine protein kinases, catalytic domai... cd00174 COG0510 COG2187 COG2334 COG3001 COG3173 COG3178 COG3231 COG3570 COG3642 COG4857 smart00090 smart00219 smart00326 smart00587 smart00750 pfam03109 pfam03881 pfam04655 pfam07914 PRK09902 PRK10271 PRK12396 pfam06293 PRK01723 PRK04750 PRK09550 PRK11768 pfam00018 pfam00069 pfam01163 pfam01633 pfam01636 pfam02958 pfam07653 pfam07714 cd05119 cd05144 cd05145 cd05146 cd05147 cd00192 cd05032 cd05033 cd05034 cd05035 cd05036 cd05037 cd05038 cd05039 cd05040 cd05041 cd05042

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



BLAST: Basic Local Alignment and Search Tool

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

BLAST finds regions of similarity between biological sequences. [more...](#)

[Learn more](#) about how to use the new BLAST design

BLAST Assembled Genomes

Choose a species genome to search, or [list all genomic BLAST databases](#).

<input type="checkbox"/> Human	<input type="checkbox"/> <i>Oryza sativa</i>	<input type="checkbox"/> <i>Gallus gallus</i>
<input type="checkbox"/> Mouse	<input type="checkbox"/> <i>Bos taurus</i>	<input type="checkbox"/> <i>Fam. troglodytes</i>
<input type="checkbox"/> Rat	<input type="checkbox"/> <i>Danio rerio</i>	<input type="checkbox"/> <i>Microbes</i>
<input type="checkbox"/> <i>Arabidopsis thaliana</i>	<input type="checkbox"/> <i>Drosophila melanogaster</i>	<input type="checkbox"/> <i>Apis mellifera</i>

Basic BLAST

Choose a BLAST program to run.

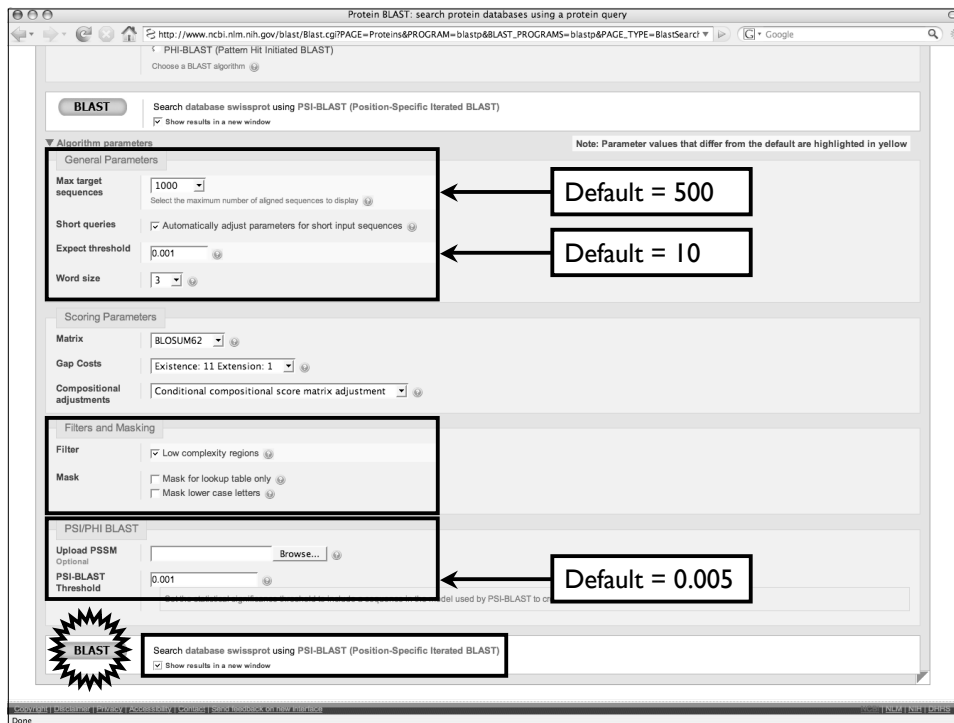
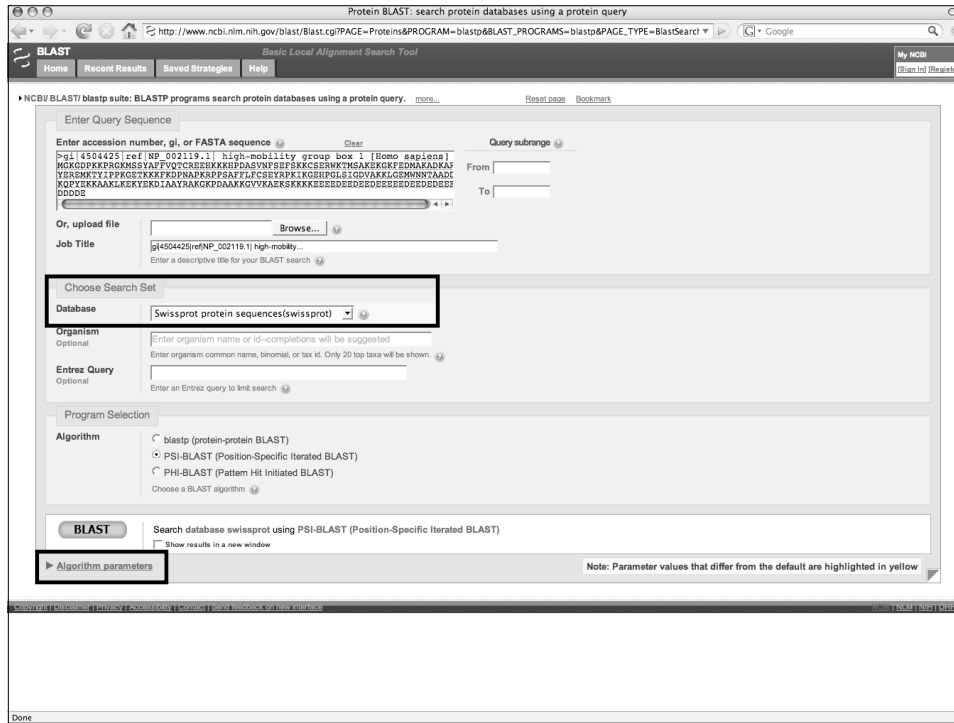
nucleotide blast	Search a nucleotide database using a nucleotide query <small>Algorithms: blastn, megablast, discontinuous megablast</small>
protein blast	Search protein database using a protein query <small>Algorithms: blastp, psi-blast, phi-blast</small>
blastx	Search protein database using a translated nucleotide query
tblastn	Search translated nucleotide database using a protein query
tblastx	Search translated nucleotide database using a translated nucleotide query

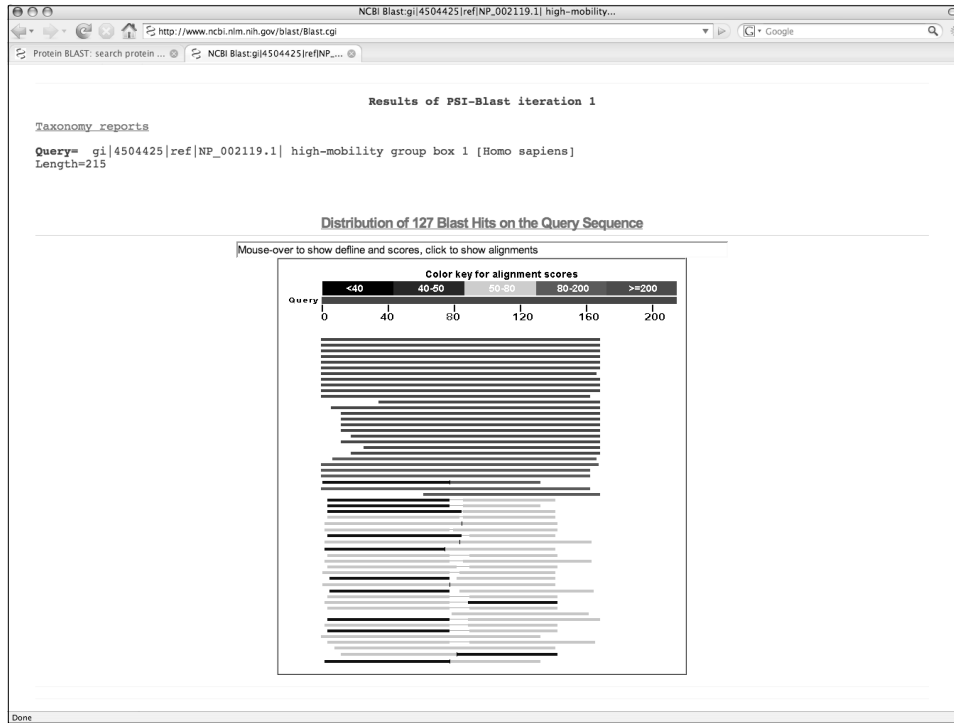
Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

- Search [trace archives](#)
- Find [conserved domains](#) in your sequence (cdt)
- Find sequences with [similar conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins](#) (igBLAST)
- Search for [SNPs](#) (snp)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align](#) two sequences using BLAST (bl2seq)

Done





Legend:

- NEW - means that the alignment score was below the threshold on the previous iteration
- OLD - means that the alignment was checked on the previous iteration

Run PSI-Blast iteration 2

Hit list size 1000

Distance tree of results NEW

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

Accession	Organism	Description	Score (Bits)	E Value
NEW <input type="checkbox"/> sp P09429 HMGB1_HUMAN	HUMAN	High mobility group protein B1 (High mo...	310	2e-84 C
NEW <input type="checkbox"/> sp P10103 HMGB1_BOVIN	BOVIN	High mobility group protein B1 (High mo...	310	2e-84 C
NEW <input type="checkbox"/> sp P63159 HMGB1_RAT	RAT	High mobility group protein B1 (High mobi...	310	2e-84 C
NEW <input type="checkbox"/> sp P12682 HMGB1_PIG	PIG	High mobility group protein B1 (High mobi...	308	9e-84 C
NEW <input type="checkbox"/> sp O9DGV6 HMGB1X_HUMAN	HUMAN	High mobility group protein 1-like 10 (HMG	290	2e-78 C
NEW <input type="checkbox"/> sp P26584 HMGB2_CHICK	CHICK	High mobility group protein B2 (High mo...	257	2e-68 C
NEW <input type="checkbox"/> sp P07746 HMGT_ONCMY	ONCMY	High mobility group-T protein (HMG-T) (HMG-	257	3e-68 C
NEW <input type="checkbox"/> sp P26583 HMGB2_HUMAN	HUMAN	High mobility group protein B2 (High mo...	252	6e-67 C
NEW <input type="checkbox"/> sp P52925 HMGB2_RAT	RAT	High mobility group protein B2 (High mobi...	251	1e-66 C
NEW <input type="checkbox"/> sp P30681 HMGB2_MOUSE	MOUSE	High mobility group protein B2 (High mo...	249	8e-66 C
NEW <input type="checkbox"/> sp P17741 HMGB2_PIG	PIG	High mobility group protein B2 (High mobi...	245	8e-65 C
NEW <input type="checkbox"/> sp P07156 HMGB1_CRIGR	CRIGR	High mobility group protein B1 (High mo...	239	4e-63 C
NEW <input type="checkbox"/> sp P23497 SP100_HUMAN	HUMAN	Nuclear autoantigen Sp-100 (Speckled 10...	211	1e-54 C
NEW <input type="checkbox"/> sp P40618 HMGB3_CHICK	CHICK	High mobility group protein B3 (High mo...	211	2e-54 C
NEW <input type="checkbox"/> sp O54879 HMGB3_MOUSE	MOUSE	High mobility group protein B3 (High mo...	210	3e-54 C
NEW <input type="checkbox"/> sp O32L31 HMGB1_BOVIN	BOVIN	High mobility group protein B3	209	7e-54 C
NEW <input type="checkbox"/> sp O15347 HMGB3_HUMAN	HUMAN	High mobility group protein B3 (High mo...	208	1e-53 C
NEW <input type="checkbox"/> sp O9N1Q6 SP100_GORGO	GORGO	Nuclear autoantigen Sp-100 (Speckled 10...	207	2e-53 C
NEW <input type="checkbox"/> sp P36194 HMGB1_CHICK	CHICK	High mobility group protein B1 (High mo...	203	5e-52 C
NEW <input type="checkbox"/> sp O9N1Q5 SP100_HYLLA	HYLLA	Nuclear autoantigen Sp-100 (Speckled 10...	201	2e-51 C
NEW <input type="checkbox"/> sp O9N1Q7 SP100_PANTR	PANTR	Nuclear autoantigen Sp-100 (Speckled 10...	201	2e-51 C
NEW <input type="checkbox"/> sp O24537 HMGB2_DROME	DROME	High mobility group protein DSP1 (Protein d	176	6e-44 C
NEW <input type="checkbox"/> sp P40644 HMGB4_STRPU	STRPU	High mobility group protein 1 homolog	152	1e-36 C
NEW <input type="checkbox"/> sp O32L34 HMGB4_BOVIN	BOVIN	High mobility group protein B4	134	3e-31 C
NEW <input type="checkbox"/> sp O8WW32 HMGB4_HUMAN	HUMAN	High mobility group protein B4	129	9e-30 C

NCBI Blast: gi|4504425|ref|NP_002119.1| high-mobility...

Protein BLAST: search protein ... NCBI Blast: gi|4504425|ref|NP_002119.1| high-mobility...

Accession	Description	E-value	Date
sp V22229 SSRP1_MOUSE	FACT complex subunit SSRP1 (Facilitates...	32.1	28-04
sp Q32169.1 HM20B_BOVIN	SWI/SNF-related matrix-associated act...	45.1	28-04
sp Q92104.1 HM20B_MOUSE	SWI/SNF-related matrix-associated act...	45.1	28-04
sp Q6DLJ5 HM20A_XENTR	High mobility group protein 20A (HMG bo...	45.1	28-04
sp Q90941 PBL_CHICK	Protein polybromo-1	45.1	28-04
sp Q6AZF8 HM20A_XENLA	High mobility group protein 20A (HMG bo...	44.3	48-04
sp P40623 HMG1B_CHITE	Mobility group protein 1B	43.9	48-04
sp P40622 HMG1A_CHITE	Mobility group protein 1A	43.5	68-04
sp Q9LEF5 SSRP1_MAIZE	FACT complex subunit SSRP1 (Facilitates...	43.5	68-04
sp Q912W1 TFAM_RAT	Transcription factor A, mitochondrial precurs...	43.5	78-04
sp Q52KFA HM20A_CHICK	High mobility group protein 20A (HMG bo...	43.5	78-04
sp Q5D144 TFAM_PIG	Transcription factor A, mitochondrial precurs...	43.1	78-04
sp Q9USU7.1 YHBB_SCHPO	HMG box-containing protein C28F2.11	42.1	88-04

Run PSI-Blast iteration 2

Alignments

Get selected sequences | Select all | Deselect all | Distance tree of results

```

>|_sp|P09429|HMGB1_HUMAN High mobility group protein B1 (High mobility group protein 1)
(HMG-1)
sp|Q6YKA4|HMGB1_CANFA High mobility group protein B1 (High mobility group protein 1)
(HMG-1)
sp|Q4RB44|HMGB1_MACFA High mobility group protein B1 (High mobility group protein 1)
(HMG-1)
sp|Q08TE6|HMGB1_HORSE High mobility group protein B1 (High mobility group protein 1)
(HMG-1)
Length=215
GENE ID: 3146 HMGB1 | high-mobility group box 1 [Homo sapiens]
(Over 100 PubMed links)
Score = 310 bits (795), Expect = 2e-84, Method: Compositional matrix adjust.
Identities = 169/169 (100%), Positives = 169/169 (100%), Gaps = 0/169 (0%)
Query 1  MGRGDPKPKFRGRMSSYAFFVQTCREEHKKHPDASVNFSEFSKCKSERWKTMSAKKRGKF 60
          MGRGDPKPKFRGRMSSYAFFVQTCREEHKKHPDASVNFSEFSKCKSERWKTMSAKKRGKF
Sbjct 1  MGRGDPKPKFRGRMSSYAFFVQTCREEHKKHPDASVNFSEFSKCKSERWKTMSAKKRGKF 60
    
```

NCBI Blast: gi|4504425|ref|NP_002119.1| high-mobility...

Protein BLAST: search protein ... NCBI Blast: gi|4504425|ref|NP_002119.1| high-mobility...

Results of PSI-Blast iteration 5

No new sequences were found above the 0.001 threshold!

Taxonomy reports

Query= gi|4504425|ref|NP_002119.1| high-mobility group box 1 [Homo sapiens]
 Length=215

Distribution of 183 Blast Hits on the Query Sequence

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

Color key for alignment scores

Score Range	Color
<40	Black
40-50	Dark Grey
50-60	Medium Grey
60-80	Light Grey
80-200	White
>=200	Black

Query 0 40 80 120 160 200

127
↓
183

Overview

- Week 2
 - Similarity vs. Homology
 - Global vs. Local Alignments
 - Scoring Matrices
 - BLAST
 - BLAT
- Week 3
 - Profiles, Patterns, Motifs, and Domains
 - Structures: VAST, Cn3D, and *de novo* Prediction
 - Multiple Sequence Alignment



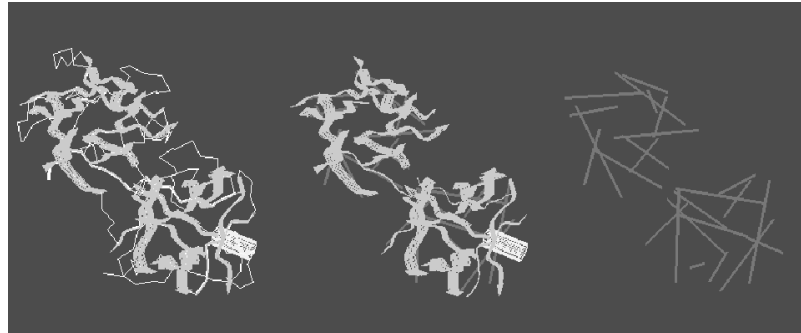
Predicting Tertiary Structure

- 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



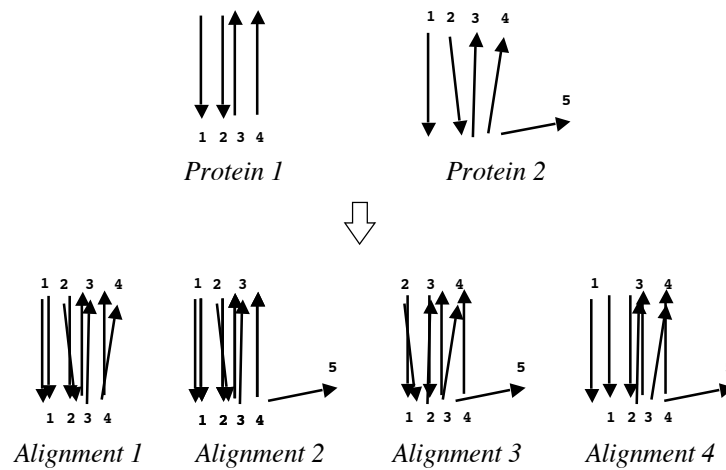
VAST Structure Comparison

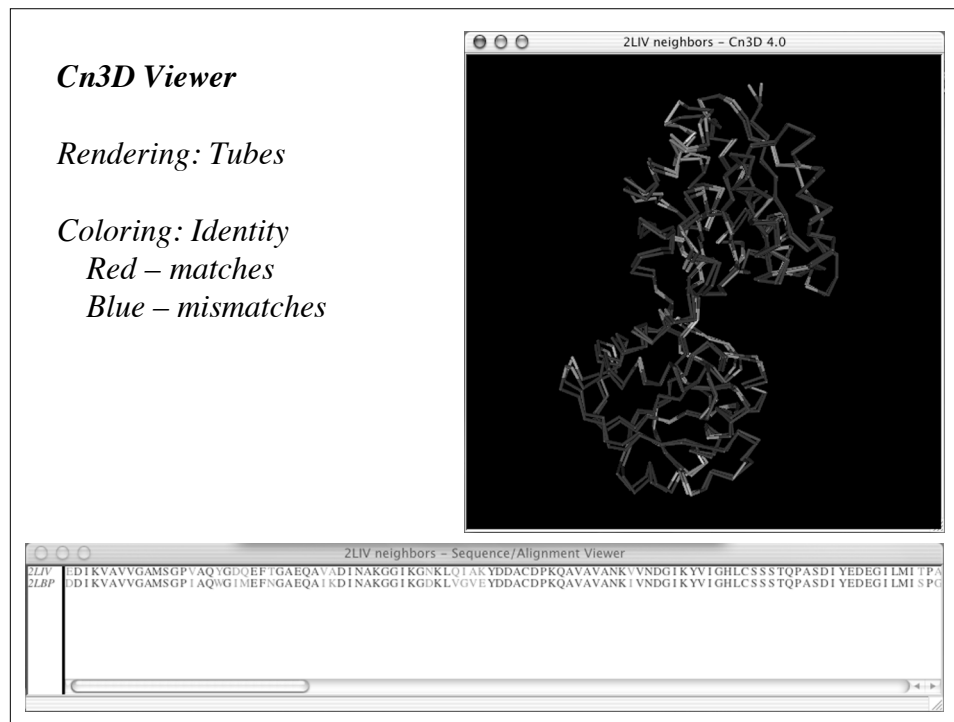
Step 1: Construct vectors for secondary structure elements



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



The screenshot shows the NCBI homepage with the URL <http://www.ncbi.nlm.nih.gov>. The page features a search bar at the top with the text "Search Structure" and "for 2LIV". Below the search bar, there are several navigation menus and sections:

- NCBI National Center for Biotechnology Information**: National Library of Medicine, National Institutes of Health.
- Navigation Menus**: PubMed, All Databases, BLAST, OMIM, Books, TaxBrowser, Structure.
- SITE MAP**: Alphabetical List, Resource Guide, About NCBI, GenBank, Literature databases, Molecular databases, Genomic biology, Tools, Research at NCBI, Software.
- What does NCBI do?**: Established in 1988 as a national resource for molecular biology information. NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. More...
- Hot Spots**: Assembly Archive, Clusters of orthologous groups, Coffee Break, Genes & Disease, NCBI Handbook, Electronic PCR, Entrez Home, Entrez Tools, Gene expression omnibus (GEO), Human genome resources, Influenza Virus Resource, Map Viewer, dbMHC, Mouse genome resources, My NCBI, ORF finder.
- GenBank vs. RefSeq**: Confused about the distinctions between GenBank, RefSeq, TPA and Uniprot? Click here for a brief description of the databases and their differences.
- New dbGaP**: NCBI's dbGaP Genome Wide Association Database. NCBI's dbGaP (Database of Genotype and Phenotype) provides data from Genome Wide Association (GWA) studies. The resource is intended to help elucidate the link between genes and disease. For each study, users have access to detailed information about the phenotypic variables measured and pre-computed associations between subjects' phenotypes and genotypes. Click here to read the press release. To read more about GWA projects, see NCBI's GWA resource page.
- PubMed Central**: An archive of biomedical and life sciences journals. Free fulltext. Over 1,100,000 articles from over 340 journals. Linked to PubMed and fully searchable. Use of PubMed Central requires no registration or fee. Access it from any computer with an internet connection.

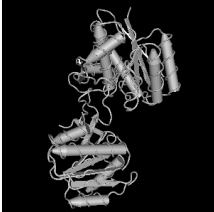
The screenshot shows the NCBI Structure search results page for the query "2LIV". The URL is <http://www.ncbi.nlm.nih.gov/sites/entrez?db=structure&cmd=search&term=2LIV>. The page displays the search results for the query "2LIV" and includes the following information:

- Search Results**: All: 1 Bacterial: 1 Eukaryotic: 0 Ligand: 0 NMR: 0 X-ray: 1
- Display**: Summary, Show 20, Sort by, Send to, Download Cn3D
- Results**: 1: 2LIV. Periplasmic Binding Protein Structure And Function. Refined X-Ray Structures Of The LeucineISOLEUCINEVALINE-Binding Protein And Its Complex With Leucine [mmdbid:58084].
- Navigation**: Limits, Preview/Index, History, Clipboard, Details.
- Left Sidebar**: About Entrez, Entrez Structure, Structure Research, MMDB, CDD, PDBaest, Cn3D, VAST, VAST Search, Research Structure Group research projects.
- Right Sidebar**: My NCBI, Sign In, Register.
- Footer**: Write to the Help Desk, NCBI | NLM | NIH, Department of Health & Human Services, Privacy Statement | Freedom of Information Act | Disclaimer.

Structure Summary, 2LIV, 58084
 http://www.ncbi.nlm.nih.gov/Structure/mmdb/mmdbsrv.cgi?Dopt=s&id=58084

NCBI
Structure Summary
 MMDB

PubMed BLAST Structure Taxonomy OMIM Help? Cn3d



Reference: Sack JS, Saper MA, Quioco FA Periplasmic binding protein structure and function. Refined X-ray structures of the leucine/isoleucine/valine-binding protein and its complex with leucine *J. Mol. Biol.* v206, p.177-191
 All References

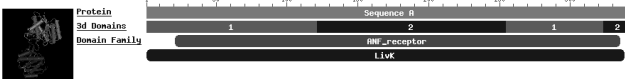
Description: Periplasmic Binding Protein Structure And Function. Refined X-Ray Structures Of The LeucineISOLEUCINEVALINE-Binding Protein And Its Complex With Leucine.

Deposition: 1989/4/10

Taxonomy: Escherichia coli
MMDB: 58084 **PDB:** 2LIV **Related Structures:** VAST

View options (Click image to view 3D structure)
 Download Cn3D!

Molecular components in the MMDB structure are listed below. The icons indicate macromolecular chains, 3D domains, protein classifications and ligands. Please hold the mouse over each icon for more information on the component. You may also click the thumbnails below to view corresponding chains and domains in Cn3D.



Done

Vast Neighbor Summary
 http://www.ncbi.nlm.nih.gov/Structure/vast/vastsrv.cgi?sdid=242528

NCBI
Related Structures
 VAST

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

VAST related structures for: MMDB 58084, 2LIV sequence A

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

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

View Sequence Alignment using Hypertext for Selected VAST related structures

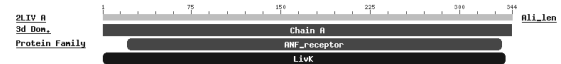
List All sequences subset, sorted by Vast E-value in Table

Advanced related structure search

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

Total related structures: 6424; 1 - 60 of 1134 representatives from the Medium redundancy subset displayed. Page: 1

Click to: Check All Uncheck All



Accession	Chain #	Residues
2LIV B	1	1-344
2LIV B	2	1-344
1EAT B	1	1-322
2E4Z B	1	1-317
1BP4 C	1	1-310
1JUN B	1	1-303
1Q00 B	1	1-291
2H46 B	1	1-290
1Z15 B 1	1	1-252
2LBP B 1	1	1-250
1G00 B	1	1-241
1ZHH B	1	1-241
2E4V B	1	1-239

Related Structures VAST

Published BLAST Structure Taxonomy OMIM Help? Cn3D

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

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

View Sequence Alignment using Hypertext for Selected VAST related structures

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

Advanced related structure search

1 - 60 of 6424 related structures displayed Page: 1

Click to: <input type="checkbox"/> Check All <input type="checkbox"/> Uncheck All	PDB C D	Ali. Len	Score	E_Val	Rmsd	%Id	MMDB Date	LHM	GSP	Description
<input type="checkbox"/>	1Z15 A	344	42.1	10e-48.8	1.3	99.7	10/2005	0.0	0.4	Crystal Structure Analysis Of Periplasmic LeuILEVAL-Binding Protein In Superopen Form
<input checked="" type="checkbox"/>	2LBP A	344	39.8	10e-44.6	0.9	79.1	10/2007	0.2	0.3	Structure Of The L-Leucine-Binding Protein Refined At 2.4 Angstroms Resolution And Comparison With The Leu(Slash)Ile(Slash)val-Binding Protein Structure
<input type="checkbox"/>	1USG A	343	40.1	10e-42.4	2.0	79.0	01/2004	0.2	0.6	L-Leucine-Binding Protein, Apo Form
<input type="checkbox"/>	1JDP B	302	29.8	10e-22.5	4.3	13.6	10/2001	6.2	1.5	Crystal Structure Of HormoneRECEPTOR COMPLEX
<input type="checkbox"/>	1YK0 A	314	29.7	10e-22.3	4.5	15.0	05/2006	6.0	1.5	Structure Of Natriuretic Peptide Receptor-C Complexed With Atrial Natriuretic Peptide
										Structure Of Natriuretic Pentide Receptor-C

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

Read the descriptions!

Cn3D Viewer

Rendering: Tubes

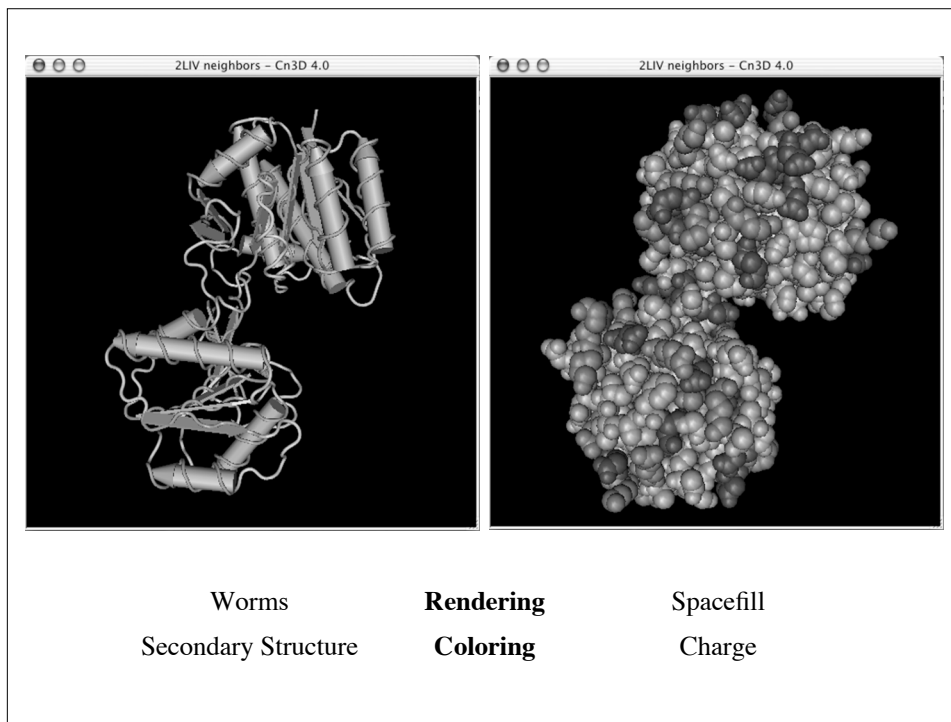
Coloring: Identity
Red – matches
Blue – mismatches

2LIV neighbors - Cn3D 4.0

2LIV neighbors - Sequence/Alignment Viewer

```

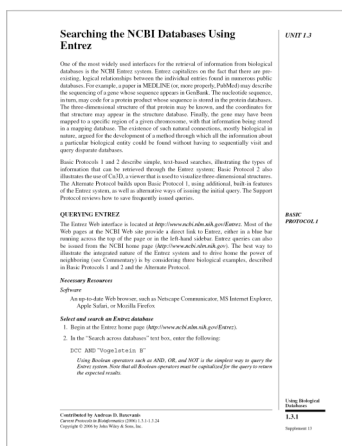
2LIV  | D I K V A V V G A M S G P V A Q Y G D Q E F T G A E Q A V A D I N A K G G I R G N K L O I A K Y D D A C D P K Q A V A V A N K I V N D G I K Y V I G H L C S S S T Q P A S D I Y E D E G I L M I T P A
2LBP  | D D I K V A V V G A M S G P I A Q W G I M E F N G A E Q A I K D I N A K G G I R G D K L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

Worms Secondary Structure Rendering Coloring Spacefill Charge

Current Protocols in Bioinformatics

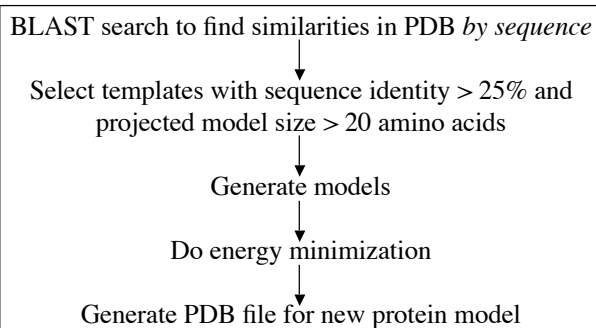
CPBI Unit 1.3 Entrez and Cn3D



<http://nihlibrary.nih.gov>
 Search "Online Journals" for "Current Protocols in Bioinformatics"

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
    
```

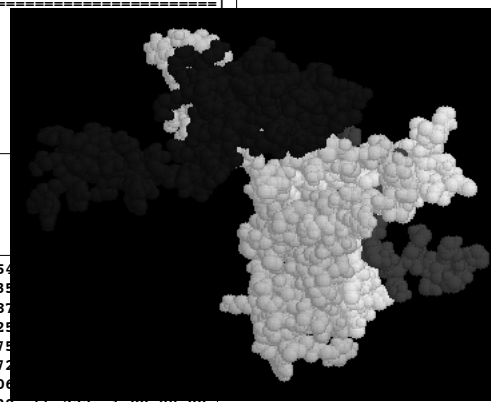
Target:

```

    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

- Modeller <http://www.salilab.org/modeller/>
- DeepView <http://us.expasy.org/spdbv/>
- WHAT IF <http://swift.cmbi.kun.nl>



Current Topics in Genome Analysis

Week 14
Tuesday, April 15, 2008

Protein Structure Analysis and Protein-Protein Interactions

*David Wishart, Ph.D.
Departments of Computing Science and
Biological Sciences
University of Alberta*



Overview

- Week 2
 - Similarity vs. Homology
 - Global vs. Local Alignments
 - Scoring Matrices
 - BLAST
 - BLAT
- Week 3
 - Profiles, Patterns, Motifs, and Domains
 - Structures: VAST, Cn3D, and *de novo* Prediction
 - Multiple Sequence Alignment



Why do multiple sequence alignments?

- Identify conserved regions, patterns, and domains
 - Experimental design
 - Predicting structure and function
 - Identifying new members of protein families
- Perform phylogenetic analysis
- Generate position-specific scoring matrices for subsequent searches (“many-against-one” or “one against many”)
- Bolster confidence in secondary structure predictions



Considerations

- Absolute sequence similarity
Create the alignment by lining up as many common characters as possible
- Conservation
Take into account residues that can substitute for one another and not adversely affect the function of the protein
- Structural similarity
Knowledge of the secondary or tertiary structure of the proteins being aligned can be used to fine-tune the alignment



General Guidelines

- As with most analyses, concentrate on the protein level rather than on the nucleotide level
 - More informative
 - Less prone to inaccurate alignment (“20 vs. 4”)
 - Can “translate back” to nucleotide sequences *after* doing the alignment



General Guidelines

- Use a reasonable number of sequences to avoid technical difficulties
 - **Global** alignment method: compute time increases exponentially as sequences are added to the set
 - Most alignment algorithms are ineffective on huge data sets (and may yield inaccurate alignments)
 - Phylogenetic studies resulting from inordinately large data sets are almost impossible
 - Good starting point: 10-15 sequences
 - Ballpark upper limit: 50 sequences



General Guidelines

- Selecting sequences for alignment
 - Sequences should be of about the same length
 - Use closely-related sequences to determine “required” amino acids
 - Use more divergent sequences to study evolutionary relationships
 - Good starting point: use sequences that are 30-70% similar to most of the other sequences in the data set
 - The most informative alignments result when the sequences in the data set are not “too similar”, but also not “too different”



General Guidelines

- Iterative process
 - Perform alignment on small set of sequences
 - Examine the quality of the alignment
 - If alignment good, can add new sequences to data set, then realign
 - If alignment not good, remove any sequences that result in the inclusion of long gaps, then realign



Interpretation

- Absolutely-conserved positions are *required* for proper structure and function
- Relatively well-conserved positions are able to tolerate limited amounts of change and not adversely affect the structure or function of the protein
- Non-conserved positions may “mutate freely,” and these mutations can possibly give rise to proteins with new functions



Interpretation

- Gap-free blocks probably correspond to regions of secondary structure
- Gap-rich blocks probably correspond to unstructured or loop regions



ClustalW2

- Automatic multiple alignment of nucleotide or amino acid sequences
- Implementations
 - Client versions
command-line text menu system, all platforms
 - Web-based version
<http://www.ebi.ac.uk/clustalw2>



Progressive Alignment

- Align two sequences at a time
- Gradually build up the multiple sequence alignment by merging larger and larger sub-alignments, clustering on the basis of similarity
- Uses protein scoring matrices and gap penalties to calculate alignments having the best score
- Major advantages of method
 - Very fast
 - Alignments generally of high quality



Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDVGGGEALGRLLVVYPWTQRRFFESFGDLST
>sequence B
VQLSGEAKAAVLALWDKVNVEEVGGGEALGRLLVVYPWTQRRFFDSFGDSLN
>sequence C
VLSPADKTNVKAANGKVGAAHAGEYGAEALERMFSLFPTTKTYFPHFDLSH
>sequence D
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSH
```



Progressive Alignment

1. Calculate a similarity score (percent identity) between every pair of sequences to drive the alignment

For N sequences, this requires the calculation of $[N \times (N - 1)] / 2$ pairwise alignments

Sequences	Alignments
4	6
10	45
25	300
50	1,225
100	4,950



Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDENVGGEALGRLLVVYPWTQRFFESFGDLST
>sequence B
VQLSGEELKAALVLALWDKVNEEEVGGGEALGRLLVVYPWTQRFFDSFGDSL
>sequence C
VLSPADKTNVKAANGKVGAAHAGEYGAEALERMFSLFPTTKTYFPHFDLSH
>sequence D
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSH
```

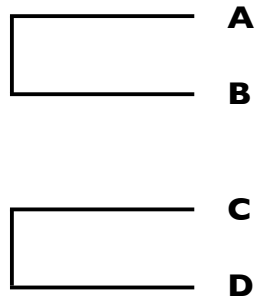
% ID	A	B	C	D
A	100			
B	80	100		
C	44	40	100	
D	40	40	92	100



Progressive Alignment

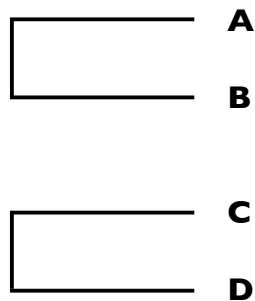
- Derive a dendrogram (guide tree) based on the pairwise comparisons (.dnd file)

Can infer from tree that A and B share greater similarity with each other than with C or D



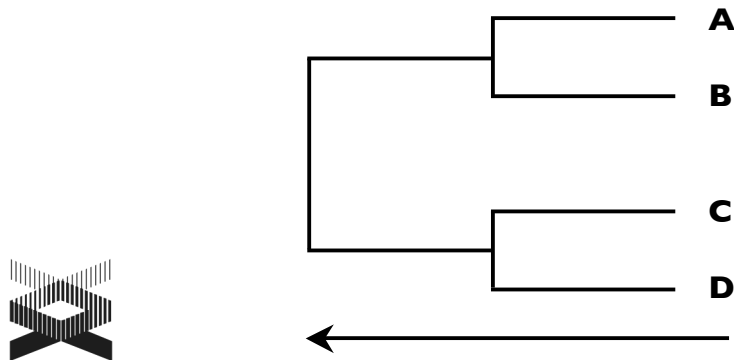
Progressive Alignment

- Align A with B → alignment AB (fixed)
- Align C with D → alignment CD (fixed)
- Represent alignments AB and CD as *single sequences*



Progressive Alignment

6. Align “sequence” AB with “sequence” CD
7. Continue following the branching order of the tree, from the tips to the root, merging each new pair of “sequences”



Progressive Alignment: Advantages

- Do “easier” alignments between highly-related sequences first
- Use information regarding conservation at each position to help with more difficult alignments between more distantly-related sequences later on in process



Progressive Alignment: Disadvantages

- If initial alignments are made on distantly related sequences, there may be errors in the initial alignments
- Once an alignment is “fixed”, it is not reconsidered, so any errors in the early alignments may propagate through subsequent alignments
- New version of ClustalW2 does provide a “remove first” iteration scheme to attempt to improve alignments



ClustalW2 Output

- Pairwise scores
- Multiple sequence alignment (.aln)
 - Alternative formats available:
GCG, Phylip, PIR, GDE



ClustalW2 Output

- Cladogram
 - Tree assumed to be an estimate of a phylogeny
 - Branches are of equal length
 - Cladograms show common ancestry, but do not provide an indication of the amount of “evolutionary time” separating taxa
- Phylogram
 - Tree that is assumed to be an estimate of phylogeny
 - Branch lengths proportional to the amount of inferred evolutionary change



ClustalW2 Conservation Patterns

- Conservation patterns in multiple sequence alignments usually follow the following rules:

[WYF]	Aromatics
[KRH]	Basic side chains (+)
[DE]	Acidic side chains (-)
[GP]	Ends of helices
[HS]	Catalytic sites
[C]	Cysteine cross-bridges



ClustalW2 Conservation Patterns

- Interpretation is *empirical* — there is no parallel to the *E*-values seen in BLAST searches to assess “significance”
 - * entirely conserved column
(want in at least 10% of positions)
 - “conserved”
(according to color table)
 - “semi-conserved”



ClustalW Colors

AVFPMILW	Red	Small
DE	Blue	Acidic
RK	Magenta	Basic
STYHCNGQ	Green	



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

ClustalW2

ClustalW2 is a general purpose multiple sequence alignment program for DNA or proteins. It produces biologically meaningful multiple sequence alignments of divergent sequences. It calculates the best match for the selected sequences, and lines them up so that the identities, similarities and differences can be seen. Evolutionary relationships can be seen via viewing Cladograms or Phylograms. New users, please read the FAQ.

Download Software

YOUR EMAIL: [] ALIGNMENT TITLE: Sequence RESULTS: interactive ALIGNMENT: full

KTUP (WORD SIZE): [] WINDOW LENGTH: [] SCORE TYPE: percent TOPDIAG: [] PAIRGAP: []

MATRIX: [] GAP OPEN: [] NO END GAPS: [] GAP EXTENSION: [] GAP DISTANCES: []

blosum ← PAM
 BLOSUM
 Gonnet (default)
 DNA Identity

ITERATION: [] NUMBER: []

alignment

OUTPUT FORMAT: [] OUTPUT ORDER: [] TREE TYPE: [] CORRECT DIST.: [] IGNORE GAPS: [] CLUSTERING: []

aln w/numbers | aligned | none | off | off | NJ

Enter or paste a set of sequences in any supported format:

```
>F05B_MOUSE Protein fosB
MFQAFPGDYDGSRCSSSPSAESQYLLSSVDSFGSPPTAAASQECAGLGEMPGSFV
ITTSQDLQWVQPLISSMAQSGQPLASQPPAVDYPDMPTSYSTPLGSAYSYTG
GGSTSTTGGVSRARARPRPRETTTEEEERKRVKREKLAACAKCNKJ
DRLOASTDLEEEKAELESEIABLQEKERLEFVLVAHPGCKIPYEEGGPGPPL
LPGSTSAKEDGFGLPPPPPLPFQSSRDAPPNLTA5LFTHSEVVLGDPFPV
TSSFVLTCPVSAFAGAORTSGSBQSDPLNSPSLLAL

>F05B_HUMAN Protein fosB
MFQAFPGDYDGSRCSSSPSAESQYLLSSVDSFGSPPTAAASQECAGLGEMPGSFV
```

Upload a file: [] Browse... [Run] [Reset]

If you plan to use these services during a course please contact us. Please read the FAQ before seeking help from our support staff.

Terms of Use EBI Funding Contact EBI European Bioinformatics Institute 2006-2007. EBI is an Outstation of the European Molecular Biology Laboratory.

ClustalW2

ClustalW2 is a general purpose multiple sequence alignment program for DNA or proteins. It produces biologically meaningful multiple sequence alignments of divergent sequences. It calculates the best match for the selected sequences, and lines them up so that the identities, similarities and differences can be seen. Evolutionary relationships can be seen via viewing Cladograms or Phylograms. New users, please read the FAQ.

Download Software

YOUR EMAIL: [] ALIGNMENT TITLE: Sequence RESULTS: interactive ALIGNMENT: full

KTUP (WORD SIZE): [] WINDOW LENGTH: [] SCORE TYPE: percent TOPDIAG: [] PAIRGAP: []

MATRIX: [] GAP OPEN: [] NO END GAPS: [] GAP EXTENSION: [] GAP DISTANCES: []

blosum

ITERATION: [] NUMBER: []

alignment ← Tree
 Alignment
 Default Iterations

Each step
 Final step
 3

OUTPUT FORMAT: [] OUTPUT ORDER: [] TREE TYPE: [] CORRECT DIST.: [] IGNORE GAPS: [] CLUSTERING: []

aln w/numbers | aligned | none | off | off | NJ

Enter or paste a set of sequences in any supported format:

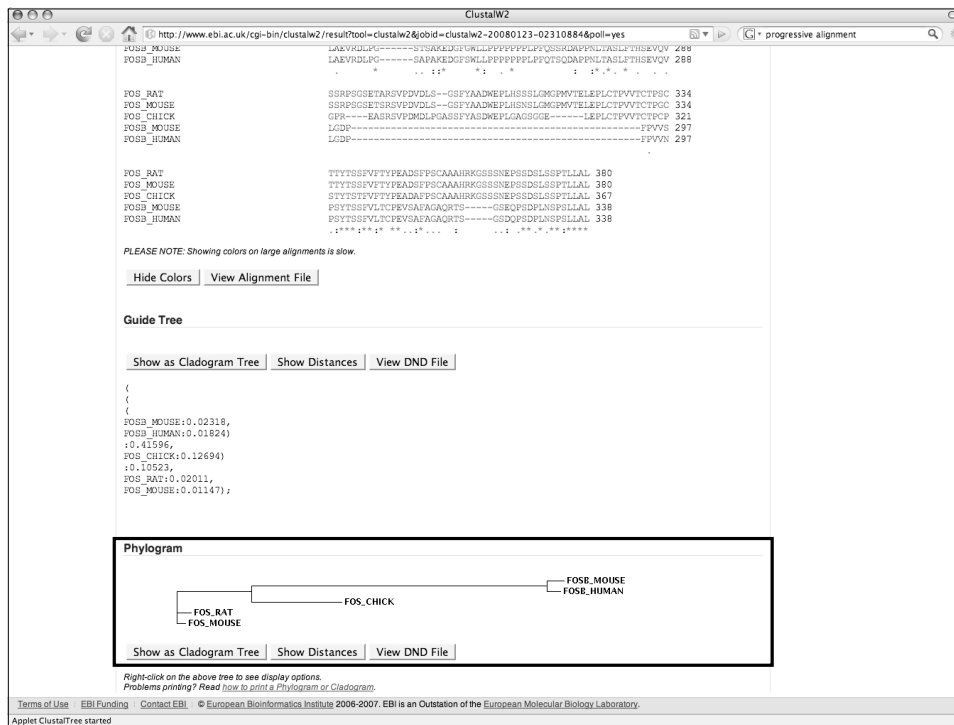
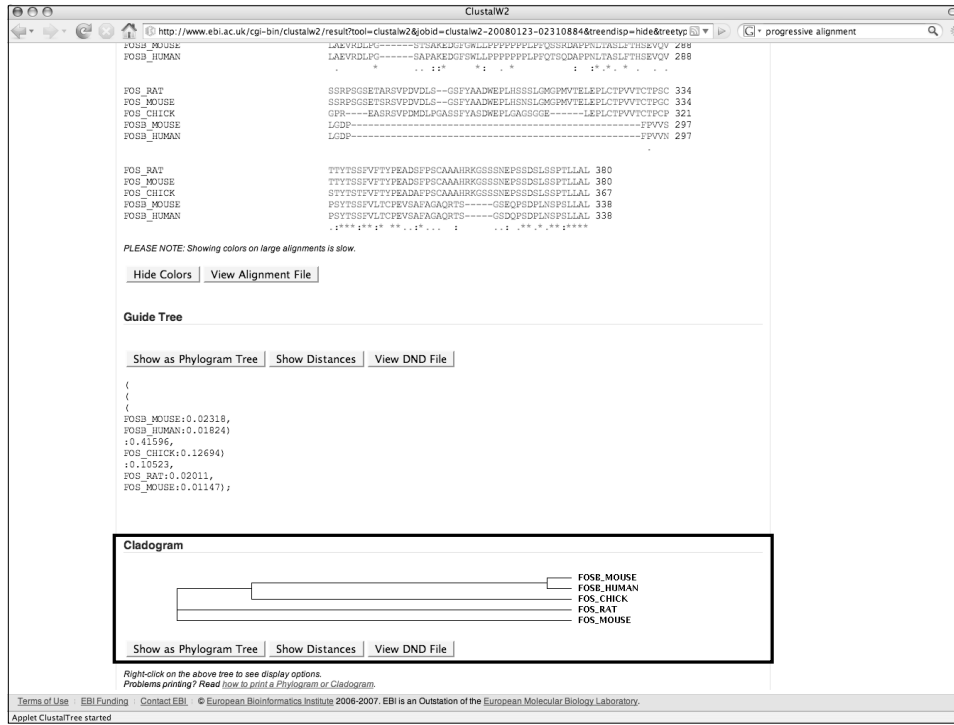
```
>F05B_MOUSE Protein fosB
MFQAFPGDYDGSRCSSSPSAESQYLLSSVDSFGSPPTAAASQECAGLGEMPGSFV
ITTSQDLQWVQPLISSMAQSGQPLASQPPAVDYPDMPTSYSTPLGSAYSYTG
GGSTSTTGGVSRARARPRPRETTTEEEERKRVKREKLAACAKCNKJ
DRLOASTDLEEEKAELESEIABLQEKERLEFVLVAHPGCKIPYEEGGPGPPL
LPGSTSAKEDGFGLPPPPPLPFQSSRDAPPNLTA5LFTHSEVVLGDPFPV
TSSFVLTCPVSAFAGAORTSGSBQSDPLNSPSLLAL

>F05B_HUMAN Protein fosB
MFQAFPGDYDGSRCSSSPSAESQYLLSSVDSFGSPPTAAASQECAGLGEMPGSFV
```

Upload a file: [] Browse... [Run] [Reset]

If you plan to use these services during a course please contact us. Please read the FAQ before seeking help from our support staff.

Terms of Use EBI Funding Contact EBI European Bioinformatics Institute 2006-2007. EBI is an Outstation of the European Molecular Biology Laboratory.



Jalview

- Java applet available within ClustalW2 results
- Used to manually edit ClustalW2 alignments
- Color residues based on various properties
- Pairwise alignment of selected sequences
- Consensus sequence calculations
- Removal of redundant sequences
- Calculation of phylogenetic trees
- Color PostScript output



The screenshot shows the ClustalW2 web interface. The browser address bar shows the URL: <http://www.ebi.ac.uk/cgi-bin/clustalw2/result?tool=clustalw2&jobid=clustalw2-20080123-02310884&treeendp=hide&retype= progressive alignment>. The page title is "ClustalW2 Results".

Results of search

Number of sequences	5
Alignment score	1076
Sequence format	Pearson
Sequence type	aa
Jalview	Start Jalview
Output file	clustalw2-20080123-02310884.output
Alignment file	clustalw2-20080123-02310884.aln
Guide tree file	clustalw2-20080123-02310884.dnd
Your input file	clustalw2-20080123-02310884.inout

To save a result file right-click the file link in the above table and choose "Save Target As".
 If you cannot see the Jalview button, reload the page and check your browser settings to enable Java Applets.

Scores Table

Sort by: **Sequence Number** | View Output File

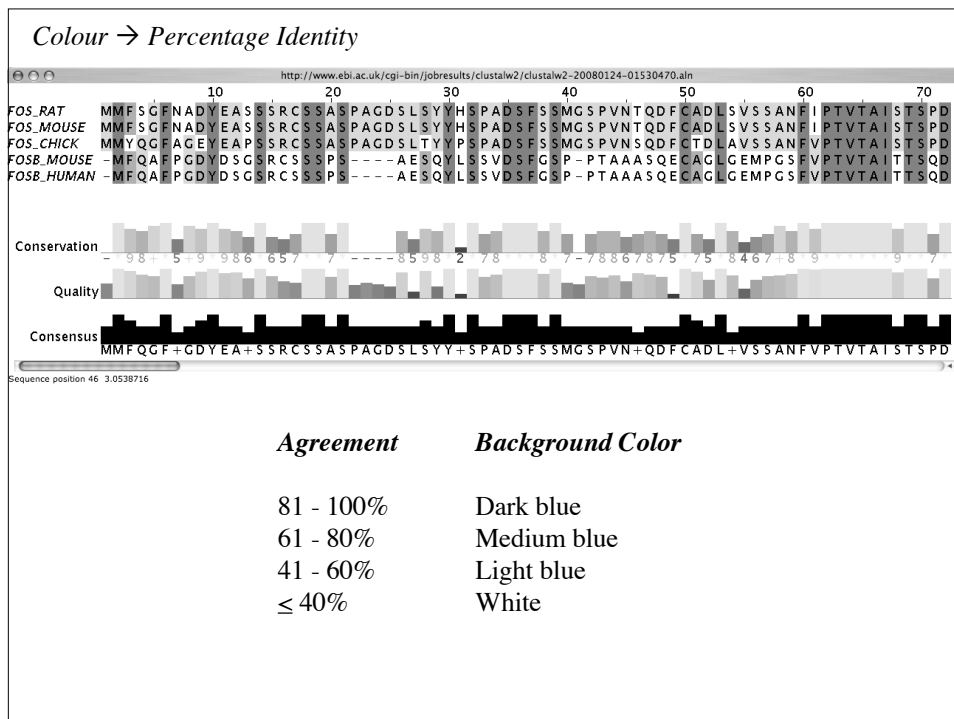
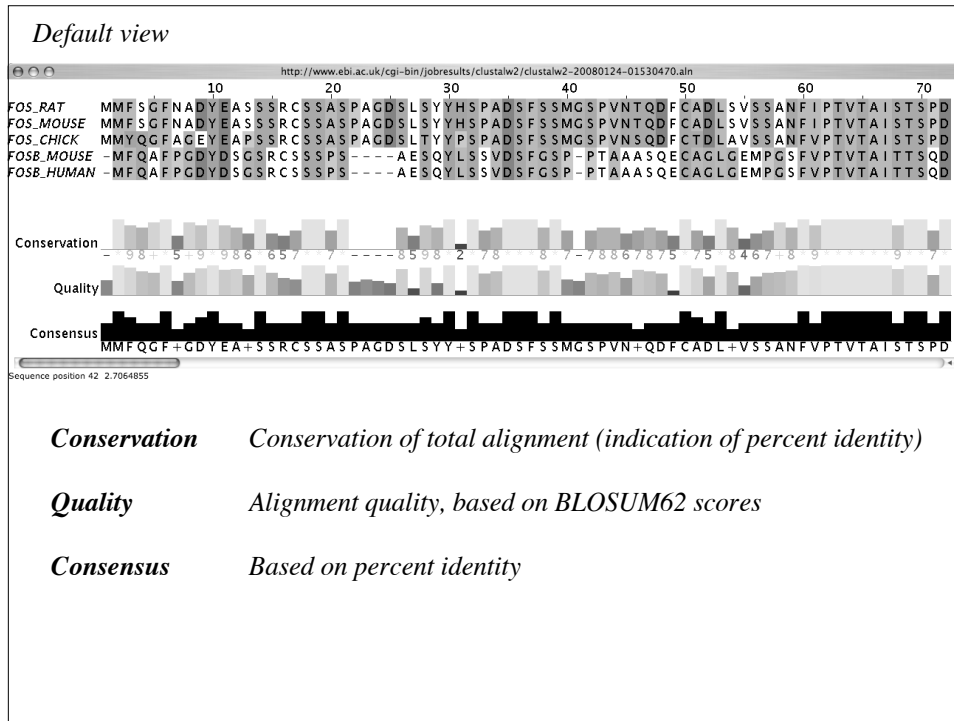
Seq#	Name	Len (aa)	Seq#	Name	Len (aa)	Score
1	POS_MOUSE	338	2	POS_HUMAN	338	95
1	POS_MOUSE	338	3	POS_CHICK	367	43
1	POS_MOUSE	338	4	POS_RAT	380	43
1	POS_MOUSE	338	5	POS_MOUSE	380	44
2	POS_HUMAN	338	3	POS_CHICK	367	43
2	POS_HUMAN	338	4	POS_RAT	380	43
2	POS_HUMAN	338	5	POS_MOUSE	380	45
3	POS_CHICK	367	4	POS_RAT	380	74
3	POS_CHICK	367	5	POS_MOUSE	380	75
4	POS_RAT	380	5	POS_MOUSE	380	96

PLEASE NOTE: Some scores may be missing from the above table if the alignment was done using multiple CPU mode. Please check the output.

Sort by: **Sequence Number** | View Output File

Alignment

Applet ClustalTree started



Current Protocols in Bioinformatics

CPBI Unit 2.3 ClustalW

CPBI Unit 3.8 T-Coffee

Multiple Sequence Alignment Using ClustalW and ClustalX

The Clustal programs are widely used for carrying out automatic, multiple alignment of sets of nucleotide or amino acid sequences. The most familiar version is ClustalW (Thompson et al., 1994), which uses a simple text menu system that is portable to most or less all computer systems. ClustalX (Thompson et al., 1997) features a graphical user interface and some powerful graphical options for editing the interpretation of alignments, and is the preferred version for Macintosh usage. ClustalW and ClustalX are developed in portable C, and the same source code/compilation system is used to build the software on different hardware, including Windows, Linux, and Solaris. The latest version for both programs was 1.8.1. The programs can both be run interactively, from the command line, or via a menu system on Mac or Windows. ClustalX, when run interactively, supports a rich command-line interface which allows for the most advanced use of the program. ClustalW can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script.

USING CLUSTALW AND CLUSTALX TO DO MULTIPLE ALIGNMENTS
 The programs ClustalW and ClustalX provide alternative user interfaces to the Clustal multiple alignment software. The alignment produced by the two programs are exactly the same, the only difference between ClustalW and ClustalX is the way in which the user interacts with the program. ClustalW uses a simple text menu system, while ClustalX uses a graphical user interface. Both programs are portable to most or less all computer systems. ClustalW and ClustalX are developed in portable C, and the same source code/compilation system is used to build the software on different hardware, including Windows, Linux, and Solaris. The latest version for both programs was 1.8.1. The programs can both be run interactively, from the command line, or via a menu system on Mac or Windows. ClustalX, when run interactively, supports a rich command-line interface which allows for the most advanced use of the program. ClustalW can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script. In the simplest usage (see Basic Protocol), the program can be run from a script.

Necessary Resources
Hardware:
 Unix (including Linux) workstation (e.g., Sun, Alpha, Silicon Graphics, PCs, PC with MS Windows, or Power Macintosh).
Software:
 ClustalW or ClustalX program (see Support Protocol).
Files:
 Sequences can be input to both ClustalW and ClustalX in one of seven file formats. All sequences must be in the same file. The format that is automatically recognized are: NBRF/PIR, DMSB/Protein-Protein, FASTA, Clustal, ClustalX, ClustalW, and GCG/Genbank. The sequences must be in the same file.

Contributed by John D. Thompson, Bodo A. Chichard, and David C. Higgins
 Copyright © 2005 by John Wiley & Sons, Inc.

UNIT 2.3
 BASIC
 PROTOCOL
 Receiving
 Feedback
 2.3.1

Computing Multiple Sequence/Structure Alignments with the T-Coffee Package

This unit describes how to assemble a multiple sequence alignment using the T-Coffee multiple sequence alignment package (Miyamoto et al., 2008). Although T-Coffee is often used to align closely related sequences and existing sequence alignments, T-Coffee is also much more flexible than most methods because it makes it possible to combine many alternative alignments into a single one based on an estimate of consistency between the alignments. The protocols below show how such a combination can be done, and how alternative alignment methods can be added to separate modules in the T-Coffee program.

This unit assumes that the user wants to align a set of sequences that are more or less homologous over their entire length (see Basic Protocol 1). These sequences may have been partitioned using appropriate database search strategy. Given such a data set, the user can assemble a multiple alignment in order to carry out family membership analysis (see Basic Protocol 2), structural modeling (see Basic Protocol 3), or phylogenetic modeling (see Support Protocol). This multiple sequence alignment may also be used to analyze the potential effect of a SNP, DNA hypermethylation, Single Nucleotide Polymorphism (see Support Protocol 4) on a specific sequence in a given member of the family.

Most of this unit assumes that the user is familiar with the Unix environment (without being a specialist or programmer, see entries in a glossary). The last section (see Appendix at the end of the unit) is a bit more demanding in terms of computer skills, but it should be relatively straightforward to acquire with a basic knowledge of the scripting language Perl.
 T-Coffee is more appropriate for generating high-quality alignments, but it is more demanding in terms of resources than other, simpler programs. Given a standard 2 GHz PC with 256 MB of memory, one can expect to align more than 100 sequences that are up to 2000 nucleotides long when using the default mode. This figure is simply an indication, since the memory requirement depends on the relationship of the sequences being considered (close sequences require less time and less memory). Given these limitations, it is often a good strategy to start with a small multiple sequence alignment method such as ClustalW (see Unit 2.3) in order to quickly identify potential problems within the data set, before relying on the system with T-Coffee. Nevertheless, the authors provide some alternative strategies that make it possible to bypass some of the limitations of T-Coffee regarding memory usage.
 In this unit, protocols are presented on how to use T-Coffee on a Unix/Linux environment, taking advantage of the rich command-line options of this program. We do not show who prefer a graphical Web interface, much of what is presented here can also be accomplished by using the Web-based T-Coffee interface (Poisson et al., 2003), available at <http://align.cse.cornell.edu/tcoffee/>. This service is provided by the community for the CNRS and HP servers. Other online versions of this software exist, and an alternative that is maintained by the T-Coffee home page (accessible from the T-Coffee servers at the aforementioned URL).
 JDFE: Investigators who are unfamiliar with the Unix environment are encouraged to read previous entries in this series.

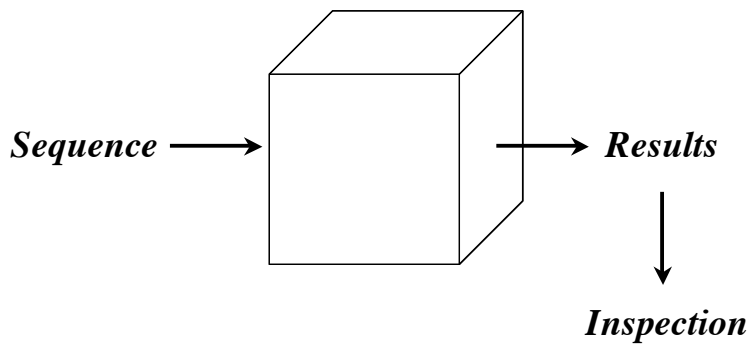
Contributed by Cedric Notredame and Karsten Willke
 Copyright © 2005 by John Wiley & Sons, Inc.

UNIT 3.8
 BASIC
 PROTOCOL
 Receiving
 Feedback
 3.8.1



<http://nihlibrary.nih.gov>
 Search "Online Journals" for "Current Protocols in Bioinformatics"

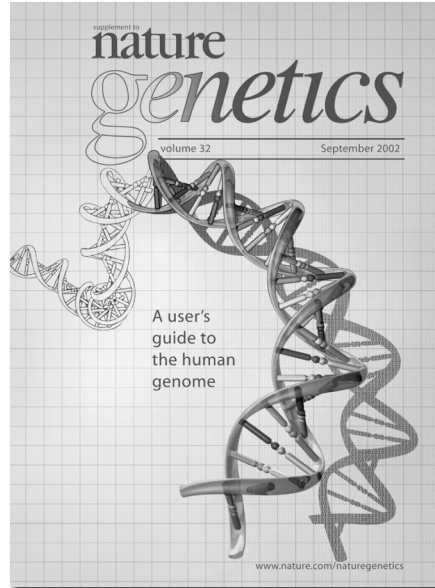
Understanding Analyses



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



Current Topics in Genome Analysis

Next Lecture:

Mining Data from Genome Browsers

Tyra Wolfsberg, Ph.D.

National Human Genome Research Institute

National Institutes of Health

