
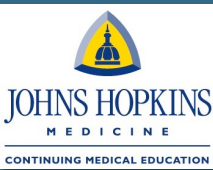



NATIONAL HUMAN GENOME RESEARCH INSTITUTE *Division of Intramural Research*




Current Topics in Genome Analysis 2016
Week 4: Biological Sequence Analysis II
Andy Baxevanis, Ph.D.

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



Current Topics in Genome Analysis 2016
Andy Baxevanis, Ph.D.
***No Relevant Financial Relationships with
Commercial Interests***

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research



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, CDD*
or 'many-against-one': *PSI-BLAST, DELTA-BLAST*



Profiles, Patterns, Motifs, and Domains



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

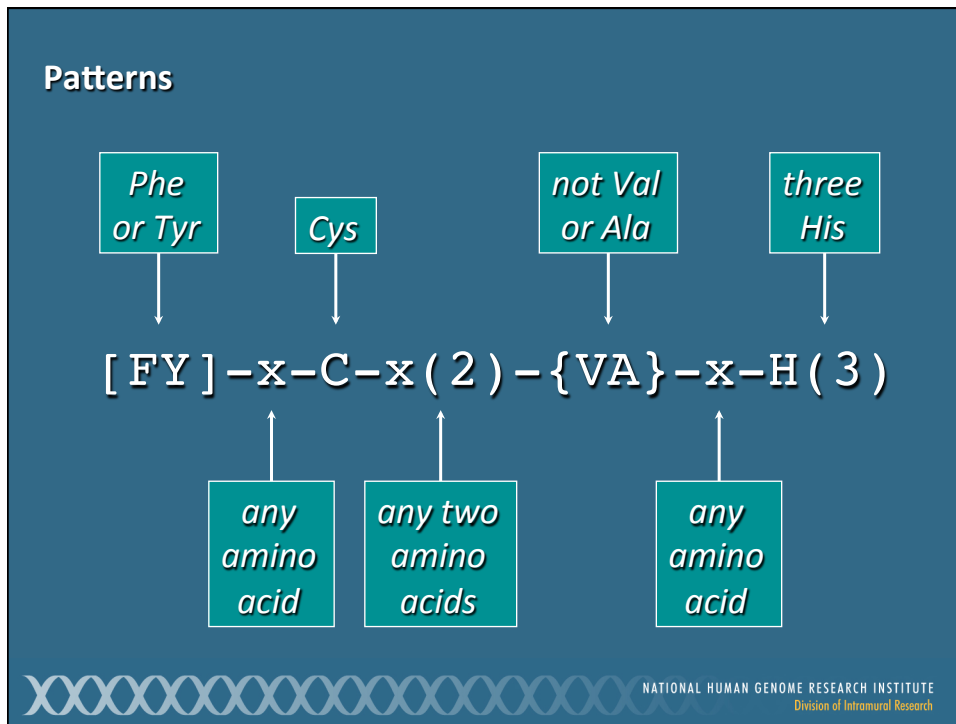
APHIIVATPG
 GCEIIVATPG
 GVEICIAATPG
 GVDILIGTTC
 RPHIIVATPG
 KPHIIATPG
 KVQLIIATPG
 RPDIVIAATPG
 APHIIVGTPG
 APHIIVGTPG
 GCHVVIATPG
 NQDIVVATPG

- 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	13	0	-12	13	0	0	-5	-5	-5	-4	15	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	-4	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	13	0	-12	13	0	0	-5	-5	-5	-4	15	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





Pfam

- Collection of multiple alignments of protein domains and conserved protein regions that probably have structural, functional, or evolutionary 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

Finn et al., *Nucleic Acids Res.* 44: D279-D285, 2016

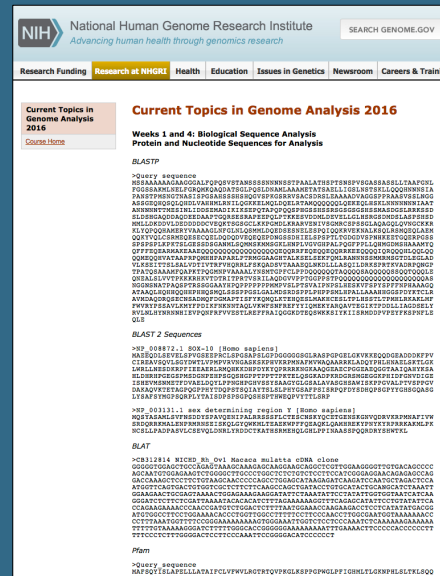
Pfam A

- Based on *curated* multiple alignments of known members of a protein family ('seed alignment')
 - Pfam definition of 'family': a collection of related protein regions
 - Based on reference proteomes (UniProtKB)
- HMMER used to find all detectable protein sequences belonging to the family
- New 'true members' of the family are then used to generate the 'full alignment' for the protein family
- Given the method used to construct the alignments, hits are highly likely to be true positives



Sequences Used in Examples

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



The screenshot displays the NIH website's 'Current Topics in Genome Analysis 2016' page. The header includes the NIH logo and navigation links for Research Funding, Research at NHGRI, Health, Education, Issues in Genetics, Newsroom, and Careers & Training. The main content area is titled 'Current Topics in Genome Analysis 2016' and 'Weeks 1 and 4: Biological Sequence Analysis Protein and Nucleotide Sequences for Analysis'. Below this, there are sections for 'BLASTP' and 'BLAST 2 Sequences'. The 'BLASTP' section shows a query sequence and its alignment with a reference sequence. The 'BLAST 2 Sequences' section shows two sequences being compared. The 'BLAT' section shows a query sequence and its alignment with a reference sequence. The footer includes the National Human Genome Research Institute logo and the text 'Division of Intramural Research'.



The screenshot shows the Pfam website homepage. At the top, there is a navigation bar with links for HOME, SEARCH, BROWSE ABOUT, FTP, and HELP. The Pfam logo is on the right, with a 'keyword search' button and a 'Go' button. Below the navigation bar, the main heading is 'Pfam 29.0 (December 2015, 16295 entries)'. A sub-heading states: 'The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). [More...](#)'

Below this, there are two columns of quick links. The left column lists: QUICK LINKS, SEQUENCE SEARCH, VIEW A PFAM ENTRY, VIEW A CLAN, VIEW A SEQUENCE, VIEW A STRUCTURE, KEYWORD SEARCH, and JUMP TO. The right column is titled 'YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...' and lists: Analyze your protein sequence for Pfam matches, View Pfam annotation and alignments, See groups of related entries, Look at the domain organisation of a protein sequence, Find the domains on a PDB structure, Query Pfam by keywords, and a 'JUMP TO' section with a text input field, 'Go' and 'Example' buttons, and instructions: 'Enter any type of accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc. Or view the [help](#) pages for more information'.

At the bottom, there is a 'Recent Pfam blog posts' section with a 'Hide this' button. The first post is 'Pfam 29.0 is now available' (posted 22 December 2015). The text of the post begins: 'Pfam 29.0, our second release of 2015, contains 16295 entries and 559 clans. We have made some major changes to our underlying sequence database and the data that are displayed on the website, which we've outlined below. Full details can be found in our Nucleic Acids Research paper, which is available here. The

This screenshot is similar to the first one, showing the Pfam website homepage. The navigation bar and main heading are identical. However, the 'SEQUENCE SEARCH' section is expanded. It includes a text input field for 'Paste your protein sequence here to find matching Pfam entries.' and 'Go' and 'Example' buttons. Below the input field, there is a note: 'This search will use an E-value of 1.0. You can set your own search parameters and perform a range of other searches [here](#).' The 'VIEW A SEQUENCE' link in the left column is highlighted with a red arrow pointing to the 'here' link in the search note.

The 'Recent Pfam blog posts' section is also present, with the first post being 'Pfam 29.0 is now available' (posted 22 December 2015). The text of the post is the same as in the first screenshot. A second post is visible: 'Moving to xfam.org' (posted 1 May 2014).

The screenshot shows the Pfam search page. The search form is titled "Sequence search" and contains a text area with a protein sequence. Below the sequence, there are radio buttons for "Cut-off" options: "Gathering threshold" and "Use E-value". The "Use E-value" option is selected, and the "E-value" is set to 1.0. A red box highlights the "E-value" input field. There are also "Submit", "Reset", "Example protein sequence", and "Example DNA sequence" buttons. The page footer includes contact information for EMBL-EBI and the European Molecular Biology Laboratory.

The screenshot shows the Pfam search results page. It displays a single match for the family p450, which is Cytochrome P450. The match is highlighted with a green bar. Below the match, there is a table of "Significant Pfam-A Matches" with columns for Family, Description, Entry type, Clan, Envelope, Alignment, HMM, HMM length, Bit score, E-value, Predicted active sites, and Show/hide alignment. The table shows one match for p450 with a bit score of 344.0 and an E-value of 1.1e-102. A red box highlights the "Significant Pfam-A Matches" section header.

Family	Description	Entry type	Clan	Envelope	Alignment	HMM	HMM length	Bit score	E-value	Predicted active sites	Show/hide alignment
				Start End	Start End	From To					
p450	Cytochrome P450	Domain	n/a	41 505	41 500	1 457	463	344.0	1.1e-102	n/a	Show

Sequence search results
[Hide](#) the detailed description of this results page.
 Below are the details of the matches that were found. We separate Pfam-A matches into two tables, containing the significant and insignificant matches. A significant match is one where the bits score is greater than or equal to the gathering threshold for the Pfam domain. Hits which do not start and end at the end points of the matching HMM are **highlighted**.
 The Pfam graphic below shows only the **significant** matches to your sequence. Clicking on any of the domains in the image will take you to a page of information about that domain.
 Pfam does not allow any amino-acid to match more than one Pfam-A family, unless the overlapping families are part of the same clan. In cases where two members of the same clan match the same region of a sequence, only one match is shown, that with the lowest E-value.
 A small proportion of sequences within the enzymatic Pfam families have had their active sites experimentally determined. Using a strict set of rules, chosen to reduce the rate of false positives, we transfer experimentally determined active site residue data from a sequence within the same Pfam family to your query sequence. These are shown as "Predicted active sites". Full details of Pfam active site prediction process can be found in [the accompanying paper](#).
 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 each table.
 This alignment row for each hit shows the alignment between your sequence and the matching HMM. The alignment fragment includes the following rows:
 #HMM: consensus of the HMM. Capital letters indicate the most conserved positions
 #MATCH: the match between the query sequence and the HMM. A '-' indicates a positive score which can be interpreted as a conservative substitution
 #PP: posterior probability. The degree of confidence in each individual aligned residue. 0 means 0-5%, 1 means 5-15% and so on; 9 means 85-95% and a '*' means 95-100% posterior probability
 #SEQ: query sequence. A '-' indicate deletions in the query sequence with respect to the HMM. Columns are coloured according to the posterior probability
 0% 100%
 You can bookmark this page and return to it later, but please use the URL that you can find in the "Search options" section below. Please note that old results may be removed after **one week**.
 We found **1** Pfam-A match to your search sequence (**all significant**)

[Show](#) the search options and sequence that you submitted.
[Return](#) to the search form to look for Pfam domains on a new sequence.

Significant Pfam-A Matches
 Show or [hide](#) all alignments.

Family	Description	Entry type	Clan	Envelope	
				Start	End
p450	Cytochrome P450	Domain	n/a	41	

#HMM: Fpqpptlpvgnllqgxeelhevlsklkkyypifrlklygkpvvvlagvskvevlkkqeeqfprdeallatarkpkfkgvlfang.ekvklrftptlrf.....qklleelvehseelveklrkagselelditell
 #MATCH: Ppqp lp+++l+lg+++b l+k++++g++++s+pvvvlag + *k++l+kq++f+grpd ++ ++gk++ [+ + w Rr + +l sT + + lee V +eaa+ l+ k+k+ e +++++
 #PP: 88939*****g*****7755555...58899988875555*****999998867799*****
 #SEQ: **PKDQKLPFZDQMLTFLR** **ENRFLSLKLSQYQDQVFLTQSPVYVYDLSLSTKQALVYQGDQFQKRLVYGRS** **LDKESSTLPPDQVAAARRAQDALKQSLAASDPLVSRLLTAVYKXIAHLLTDTLTKAKAVQHLVSTPVPVYV**

Comments or questions on the site? Send a mail to pfam-help@ebi.ac.uk.
 European Molecular Biology Laboratory

EMBL-EBI **Pfam**
 HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Family: p450 (PF00067)

455 architectures 1973 sequences 4 interactions 929 species 1275 structures

Summary: Cytochrome P450
 Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

Wikipedia: [Cytochrome P450](#) **Pfam** [InterPro](#)

This tab holds the annotation information that is stored in the Pfam database. As we move to using Wikipedia as our main source of annotation, the contents of this tab will be gradually replaced by the Wikipedia tab.

Cytochrome P450 [Provide feedback](#)

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

Literature references

- Graham-Lorence S, Amareh B, White RE, Peterson JA, Simpson ER; , Protein Sci 1995;4:1065-1080.: A three-dimensional model of aromatase cytochrome P450. [PubMed:7549871](#) [EMBL:7549871](#)
- Degtyarenko KN, Archakov AI; , FEBS Lett 1993;332:1-8.: Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. [PubMed:8405421](#) [EMBL:8405421](#)
- Nelson DR, Kamatani 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:7678494](#) [EMBL:7678494](#)
- Guengerich FP; , J Biol Chem 1991;266:10019-10022.: Reactions and significance of cytochrome P-450 enzymes. [PubMed:2037557](#) [EMBL:2037557](#)
- Nebert DW, Gonzalez FJ; , Annu Rev Biochem 1987;56:945-993.: P450 genes: structure, evolution, and regulation. [PubMed:3304150](#) [EMBL:3304150](#)
- Werck-Reichhart D, Feyereisen R; , Genome Biol 2000;1:REVIEWS3003.: Cytochromes P450: a success story. [PubMed:11178272](#) [EMBL:11178272](#)

External database links

HOMSTRAD: [p450](#)

PRINTS: [PR00385](#) [PR00359](#) [PR00408](#) [PR00463](#) [PR00464](#) [PR00465](#)

PROSITE: [PDOC0081](#)

SCOP: [Zcmp](#)

Example structure
[PDB entry 4C9P](#): Structure of camphor bound P450 mutant of CYP101D1
 View a different structure: [4C9P](#)

Family: p450 (PF00067)

455 architectures 41973 sequences 4 interactions 929 species 1275 structures

Domain organisation

Below is a listing of the unique domain organisations or architectures in which this domain is found. [More...](#)

There are **34971 sequences with the following architecture: p450**
 WSPB74_SHEEP [Ovis aries (Sheep)] Uncharacterized protein (ECO:0000313|Ensembl:ENSDARP0000007685) (494 residues)

There are **2629 sequences with the following architecture: p450 x 2**
 W4Z265_STRPU [Strongylocentrotus purpuratus (Purple sea urchin)] Uncharacterized protein (ECO:0000313|Ensembl:Metazoa:SPLU_026477-tr) (575 residues)

There are **252 sequences with the following architecture: p450, Flavodoxin_1, FAD_binding_1, NAD_binding_1**
 Q89R90_BRADU [Bradyrhizobium diazoefficiens (strain JCM 10833 / IAM 13628 / NBRC 14792 / USDA 110)] Bir2882 protein (ECO:0000313|EMBL:BAC48147.1) (1078 residues)

There are **152 sequences with the following architecture: p450 x 3**
 M4E506_BRAR2 [Brassica rapa subsp. pekinensis (Chinese cabbage) (Brassica pekinensis)] Uncharacterized protein (ECO:0000313|Ensembl:Plants:Bra039327.1-P) (983 residues)

There are **93 sequences with the following architecture: An_peroxidase x 2, p450**
 W7MZF6_GIBM7 [Gibberella moniliformis (strain M3125 / FGSC 7600) (Maize ear and stalk rot fungus) (Fusarium verticillioides)] Prostaglandin-endoperoxide synthase 1 (ECO:0000313|EMBL:EWG53184.1) (1101 residues)

There are **68 sequences with the following architecture: An_peroxidase, p450**
 J4GN3_FIBRA [Fibroporia radiculosa (strain TFFH 294) (Brown rot fungus) (Antrodia radiculosa)] Uncharacterized protein (ECO:0000313|EMBL:CCL99225.1) (1228 residues)

There are **28 sequences with the following architecture: p450 x 4**
 V4L6B3_EUTSA [Eutrema saalsugineum (Saltwater cress) (Sisymbrium saalsugineum)] Uncharacterized protein (ECO:0000313|EMBL:ESQ39189.1) (1387 residues)

There are **15 sequences with the following architecture: p450, Fer2**
 X5EDG3_9CORY [Corynebacterium glyciniphilum AJ 3170] Cytochrome P450 (ECO:0000313|EMBL:AHW64656.1) (774 residues)

Family: p450 (PF00067)

455 architectures 41973 sequences 4 interactions 929 species 1275 structures

Alignments

We store a range of different sequence alignments for families. As well as the seed alignment from which the family is built, we provide the full alignment, generated by searching the sequence database (reference proteomes) using the family HMM. We also generate alignments using four representative proteomes (RP) sets, the UniProtKB sequence database, the NCBI sequence database, and our metagenomics sequence database. [More...](#)

View options

We make a range of alignments for each Pfam-A family. You can see a description of each [above](#). You can view these alignments in various ways but please note that some types of alignment are never generated while others may not be available for all families, most commonly because the alignments are too large to handle.

	Seed (50)	Full (41973)	Representative proteomes				UniProt (105935)	NCBI (141176)	Meta (2644)
			RP15 (9588)	RP35 (24393)	RP55 (37142)	RP75 (44573)			
Jalview	✓	✓	✓	✓	✓	✓	✓	✓	✓
HTML	✓	✓	✗	✗	✗	✗	✗	✗	✗
PP/heatmap	✗	✗	✗	✗	✗	✗	✗	✗	✗

✗ Cannot generate PP/Heatmap alignments for seeds; no PP data available

Key: ✓ available, ✗ not generated, — not available.

Format an alignment

	Seed (50)	Full (41973)	Representative proteomes				UniProt (105935)	NCBI (141176)	Meta (2644)
			RP15 (9588)	RP35 (24393)	RP55 (37142)	RP75 (44573)			
Alignment:	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Format:	Setex								
Order:	<input checked="" type="radio"/> Tree <input type="radio"/> Alphabetical								
Sequence:	<input checked="" type="radio"/> Inserts lower case <input type="radio"/> All upper case								
Gaps:	<input type="radio"/> Gaps as "." or "-" (mixed) <input checked="" type="radio"/> Gaps as "-" (lower) <input type="radio"/> Gaps as "." (higher)								
Download/view:	<input checked="" type="radio"/> Download <input type="radio"/> View								

Generate

Download options

We make all of our alignments available in Stockholm format. You can download them here as raw, plain text files or as gzip-compressed files.

	Seed (50)	Full (41973)	Representative proteomes				UniProt (105935)	NCBI (141176)	Meta (2644)
			RP15 (9588)	RP35 (24393)	RP55 (37142)	RP75 (44573)			
Raw Stockholm	✓	✓	✓	✓	✓	✓	—	—	✓

EMBL-EBI
 Seed sequence alignment for PF00067

Legend:

- Glycine (G)
- Proline (P)
- Small or hydrophobic (A,V,L,I,M,F,W)
- Hydroxyl or amine amino acids (S,T,N,Q)
- Charged amino-acids (D,E,R,K)
- Histidine or tyrosine (H,Y)
- Random coil
- Alpha-helix
- G 3(10) helix
- I Pi-helix
- Hydrogen bonded beta-strand (extended strand)
- Residue in isolated beta-bridge
- H-bonded turn (3-turn, 4-turn, or 5-turn)
- Bend (five-residue bend centered at residue i)

EMBL-EBI Pfam
 Family: p450 (PF00067)

455 architectures 41973 sequences 4 interactions 929 species 1275 structures

Species distribution

Sunburst controls: Hide

Nematostella vectensis

Weight segments by...
 number of sequences
 number of species

Change the size of the sunburst
 Small Large

Colour assignments

- Archaea
- Bacteria
- Viruses
- Eukaryota
- Other sequences
- Unclassified
- Unclassified sequence

Selections

- Align selected sequences to HMM
- Generate a FASTA-format file
- Clear selection
- Currently selected:
 - + 90 sequences
 - + 1 species
- Note: selection tools show results in pop-up windows. Please disable pop-up blockers.

Family: **p450 (PF00067)**

Summary: **Cytochrome P450**

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

Wikipedia: [Cytochrome P450](#) Pfam InterPro

This tab holds the annotation information that is stored in the Pfam database. As we move to using Wikipedia as our main source of annotation, the contents of this tab will be gradually replaced by the Wikipedia tab.

Cytochrome P450 [Provide feedback](#)

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

Literature references

- Graham-Lorence S, Amarnath B, White RE, Peterson JA, Simpson ER; , Protein Sci 1995;4:1065-1080.: A three-dimensional model of aromatase cytochrome P450. [PubMed:7549871](#) [EMBL:7549871](#)
- Deglyarenko KN, Archakov AJ; , FEBS Lett 1993;332:1-8.: Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. [PubMed:8405421](#) [EMBL:8405421](#)
- Nelson DR, Kamatani 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:7678494](#) [EMBL:7678494](#)
- Guengerich FP; , J Biol Chem 1991;266:10019-10022.: Reactions and significance of cytochrome P-450 enzymes. [PubMed:2037557](#) [EMBL:2037557](#)
- Neibert DW, Gonzalez FJ; , Annu Rev Biochem 1987;56:945-993.: P450 genes: structure, evolution, and regulation. [PubMed:3304150](#) [EMBL:3304150](#)
- Werk-Reichert D, Feyereisen R; , Genome Biol 2000;1:REVIEWS3003.: Cytochromes P450: a success story. [PubMed:11178272](#) [EMBL:11178272](#)

External database links

HOMSTRAD: [p450](#)

PRINTS: [PR00385](#) [PR00359](#) [PR00408](#) [PR00463](#) [PR00464](#) [PR00465](#)

PROSITE: [PDOC00081](#)

SCOP: [2zpp](#)

Comments or questions on the site? Send a mail to pfam-help@ebi.ac.uk
 European Molecular Biology Laboratory

PROSITE documentation **PDOC00081**

Cytochrome P450 cysteine heme-iron ligand signature

Description Technical section References Copyright Miscellaneous

Description

Cytochrome P450's [1,2,3,E1] are a group of enzymes involved in the oxidative metabolism of a high number of natural compounds (such as steroids, fatty acids, prostaglandins, leukotrienes, etc) as well as drugs, carcinogens and mutagens. Based on sequence similarities, P450's have been classified into about forty different families [4,5]. P450's are proteins of 400 to 530 amino acids; the only exception is *Bacillus BM-3* (CYP102) which is a protein of 1048 residues that contains a N-terminal P450 domain followed by a reductase domain. P450's are heme proteins. A conserved cysteine residue in the C-terminal part of P450's is involved in binding the heme iron in the fifth coordination site. From a region around this residue, we developed a ten residue signature specific to P450's.

Note:

The term 'cytochrome' P450, while commonly used, is incorrect as P450 are not electron-transfer proteins; the appropriate name is P450 heme- thiolate proteins'.

Expert(s) to contact by email:
[Deglyarenko K.N.](mailto:Deglyarenko.K.N.)

Last update:
 December 2004 / Pattern and text revised.

Technical section

PROSITE method (with tools and information) covered by this documentation:

CYTOCHROME_P450, PS00086; Cytochrome P450 cysteine heme-iron ligand signature (PATTERN)

- Consensus pattern:
 [FW][SGNH]-x-[GD]-[F]-[RKHP]-[P]-C-[LIVMFAP]-[GAD]
 C is the heme iron ligand
- Sequences in UniProtKB/Swiss-Prot known to belong to this class: 1077
 - detected by PS00086: 904 (true positives)
 - undetected by PS00086: 83 (73 false negatives and 10 'parialals')
- Other sequence(s) in UniProtKB/Swiss-Prot detected by PS00086:
 46 false positives.
- Retrieve an alignment of UniProtKB/Swiss-Prot true positive hits:
 Clustal format, color, condensed view / Clustal format, color / Clustal format, plain text / Fasta format
- Retrieve the sequence logo from the alignment
- Taxonomic distribution of all UniProtKB (Swiss-Prot + TrEMBL) entries matching PS00086
- Retrieve a list of all UniProtKB (Swiss-Prot + TrEMBL) entries matching PS00086
- Scan UniProtKB (Swiss-Prot and/or TrEMBL) entries against PS00086
- View ligand binding statistics of PS00086
- Matching PDB structures: 1AKD 1BU7 1BVI 1C8J ... [ALL]

Conserved Domain Database (CDD)

- Identify conserved domains in a protein sequence
- Incorporates three-dimensional structural information to define domain boundaries and refine alignments
- Source data derived from:
 - Pfam A
 - Simple Modular Architecture Research Tool (SMART)
 - COG (orthologous prokaryotic protein families)
 - PRK ('protein clusters' of related protein RefSeq entries)
 - TIGRFAM

Marchler-Bauer et al., Nucleic Acids Res. 43: D222-D226, 2015



Conserved Domain Database (CDD)

- CD-Search performed using RPS-BLAST
- Query sequence is used to search a database of pre-calculated position-specific scoring matrices
- *Not* the same method used by Pfam



Conserved Domains
<http://ncbi.nlm.nih.gov/Structure>

Search for Conserved Domains within a protein or coding nucleotide sequence

NEW! Use **Batch CD-search** to submit multiple query proteins at once!

Enter protein or nucleotide query as accession, gi, or sequence in FASTA format

SNP_005206.1 deleted in colorectal carcinoma (Homo sapiens)
 MENSLEKCVWYKLAFLVFGASLISLHQLVIGFQIKAFALRFLSEEDAVYMRGNNLDDCSAESRQV
 VIKKDKGHLALGMDERKQKLSNGSLLIQNLHSHRHHKDEGLVQCEASLIGSGSIIISRTAKVAVAGL
 RFLSQTSYPTAFMGDTVLLKCEVIGEPMPTEWQKNOQDLTPIGDSRVVLPFGALQISRLQPGDIGY
 NCBANPASPSTGNEARVLLSDPGLHQLGFLQNSVVAIEGKAVLCCVQYFPFPTWLAGSEVI
 QLAKSYLLGSGLLINVDVDDSGWYCVTYKNEISASAELVLPVFPFHNHNSLXAYESMDIEF
 ECTVSGKPVPTVNMKNGDVVIPSDFYQVIGSSNLRILGVVKSDEGFYQCAENAGNAQTSQQLVFPK
 AIPSSSVLPSAPRDPVPLVSSRFRLSNRPFAEAKGNIQTTFVFSRQGNRREALITQPSGLQTLVG
 NLKPEAMTFYVAVNGKQSSGPIKATQPELQVFPVFNQVSTSPFILLTWEPANANGPVGG
 YRLCFYVSTGKQZNI:EVGDSYKLEGLKFTFYSRLFLAYNRVGPVSTDDITVVLSDVSPAFQVNS
 LEVNSRSIKVSNLPPFGTQNGFI:GYKIRRRKTRRGMETLFPNNLWYLFGLKQSGYSFQVSA
 VNGTQPSNWTAEIPENDLDEGVFQDPSRLVROPQNCI:IMSWTFLMNIIVRGIIGYGVSPYAE
 FVWDSKQYYSI:ERLESSESHVILKAFNAGSVPFESATTSIDTDFDVPVYLLDQFPFVPL
 STMLFPVQVALTEADVYVNSAONSVEKNQKTSSEVRLTYVNRVTSFASAKYKSEDTTSLSYATGL
 KFMVYFVSVMVNRASSTWSMTAAHTTVEAAPTSAKPDVTI:TRGKPRVIVSWQPFLEANGKITAY
 ILPFTYLDKNIPIDQIMETISGDLTHQIMOLNLDYMYFRIQARNKSGVGLSDPILFRYLKVEHPDKM

OPTIONS

Search against database: CDD v3.14 - 47363 PSSMs

Expect Value threshold: 0.010000

Apply low-complexity filter:

Composition based statistics adjustment:

Force live search:

Rescue borderline hits: Suppress weak overlapping hits:

Maximum number of hits: 500

Result mode: Concise Standard Full

Retrieve previous CD-search result

Request ID: Retrieve

References:

- Marchler-Bauer A et al. (2015), "CDD: NCBI's conserved domain database.", *Nucleic Acids Res.*43(D)222-6.
- Marchler-Bauer A et al. (2011), "CDD: a Conserved Domain Database for the functional annotation of proteins.", *Nucleic Acids Res.*39(D)225-9.
- Marchler-Bauer A et al. (2009), "CDD: specific functional annotation with the Conserved Domain Database.", *Nucleic Acids Res.*37(D)205-10.
- Marchler-Bauer A, Bryant SH (2004), "CD-Search: protein domain annotations on the fly.", *Nucleic Acids Res.*32(W)327-331.

[Help](#) | [Disclaimer](#) | [Write to the Help Desk](#)
 NCBI | NLM | NIH

Conserved domains on [cl|seqsig_MENSL_43d16cb872e3ad4b6afc9580b64484e]

NP_005206.1 deleted in colorectal carcinoma (Homo sapiens)

Graphical summary | Zoom to residue level | show extra options

Query seq. [250 500 750 1000 1250 1500 1647]

Inter-domain contacts: Cytokine receptor motif

Specific hits: FN3, FN3, FN3, FN3, FN3, FN3

Superfamilies: Ig superfamily, I-set

Multi-domains: IGc2, I-set, Neogenin_C

List of domain hits

Name	Accession	Description	Interval	E-value
[H] Ig1_Neogenin	cd05722	First immunoglobulin (Ig)-like domain in neogenin and similar proteins; Ig1_Neogenin: first ...	41-136	4.80e-51
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	528-617	4.17e-19
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	429-521	1.33e-18
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	625-715	2.66e-17
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	946-1041	2.89e-14
[H] Ig	cd00096	Immunoglobulin domain; Ig; immunoglobulin (Ig) domain found in the Ig superfamily. The Ig ...	157-222	1.53e-11
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	846-939	6.70e-11
[H] FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	728-814	2.03e-09
[H] Ig super family	d11960	Immunoglobulin domain; Ig; immunoglobulin (Ig) domain found in the Ig superfamily. The Ig ...	347-417	6.26e-36
[H] Ig super family	d11960	Immunoglobulin domain; Ig; immunoglobulin (Ig) domain found in the Ig superfamily. The Ig ...	244-327	1.84e-15
[H] Neogenin_C	pfam06583	Neogenin C-terminus; This family represents the C-terminus of eukaryotic neogenin precursor ...	1146-1445	4.71e-144
[H] I-set	pfam07679	Immunoglobulin I-set domain;	331-417	3.39e-20
[H] I-set	pfam07679	Immunoglobulin I-set domain;	241-327	5.30e-20
[H] IGc2	smart00408	Immunoglobulin C-2 Type;	153-219	7.58e-17
[H] IGc2	smart00408	Immunoglobulin C-2 Type;	54-120	1.63e-09

References:

- Marchler-Bauer A et al. (2015), "CDD: NCBI's conserved domain database.", *Nucleic Acids Res.*43(D)222-6.
- Marchler-Bauer A et al. (2011), "CDD: a Conserved Domain Database for the functional annotation of proteins.", *Nucleic Acids Res.*39(D)225-9.
- Marchler-Bauer A et al. (2009), "CDD: specific functional annotation with the Conserved Domain Database.", *Nucleic Acids Res.*37(D)205-10.

Conserved Domains
 NCBI Conserved Domain Search

NP_005206.1 deleted in colorectal carcinoma (Homo sapiens)

Graphical summary Zoom to residue level show extra options

Query seq. 250 500 750 1000 1250 1497

Specific hits: Ig1-like, Ig super, FN3, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super

Superfamilies: Ig1-like, Ig super, FN3, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super

Multi-domains: Ig1-like, Ig super, FN3, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super, FN3 super

List of domain hits

Name	Accession	Description	Interval	E-value
Ig1_Neogenin	cd05722	First immunoglobulin (Ig)-like domain in neogenin and similar proteins. Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and related proteins. Neogenin is a cell surface protein which is expressed in the developing nervous system of vertebrate embryos in the growing nerve cells. It is also expressed in other embryonic tissues, and may play a general role in developmental processes such as cell migration, cell-cell recognition, and tissue growth regulation. Included in this group is the tumor suppressor protein DCC, which is deleted in colorectal carcinoma. DCC and neogenin each have four Ig-like domains followed by six fibronectin type III domains, a transmembrane domain, and an intracellular domain.	41-136	4.80e-51
FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	528-617	4.17e-19
FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	429-521	1.33e-18
FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	625-715	2.66e-17
FN3	cd00063	Fibronectin type 3 domain; One of three types of internal repeats found in the plasma protein ...	946-1041	2.89e-14

cd05722 Sequence Cluster

Sub-family Hierarchy

Conserved Protein Domain Family
 Ig1_Neogenin

cd05722: Ig1_Neogenin

First immunoglobulin (Ig)-like domain in neogenin and similar proteins
 Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and related proteins. Neogenin is a cell surface protein which is expressed in the developing nervous system of vertebrate embryos in the growing nerve cells. It is also expressed in other embryonic tissues, and may play a general role in developmental processes such as cell migration, cell-cell recognition, and tissue growth regulation. Included in this group is the tumor suppressor protein DCC, which is deleted in colorectal carcinoma. DCC and neogenin each have four Ig-like domains followed by six fibronectin type III domains, a transmembrane domain, and an intracellular domain.

Links

- Source: cd00096
- Taxonomy: Euteleostomi
- PubMed: 6 links
- Book: 2 links
- Protein: Representatives, Specific Protein, Related Protein, Related Structure, Architectures
- Superfamily: c11960
- BioSystems: 369 links

Statistics

- PSSM-Id: 143199
- View PSSM: cd05722
- Aligned: 7 rows
- ThresholdBitScore: 148.395
- ThresholdSettingG: 148277558
- Created: 27-Sep-2007
- Updated: 17-Jan-2013

PubMed References

- Neogenin: one receptor, many functions. *Int. J. Biochem. Cell Biol.* 2007; 39(5):874-878
- Neogenin, an avian cell surface protein expressed during terminal neuronal differentiation, is closely related to the human tumor suppressor molecule deleted in colorectal cancer. *J. Cell Biol.* 1994 Dec; 127(6):2009-2020
- Molecular characterization of human neogenin, a DCC-related protein, and the mapping of its gene (NEO1) to chromosomal position 15q22.3-q23. *Genomics* 1997 May 1; 41(3):414-421
- The immunoglobulin fold. Structural classification, sequence patterns and common core. *J. Mol. Biol.* 1994 Sep 30; 242(4):309-320
- The immunoglobulin superfamily: an insight on its tssular, species, and functional diversity. *J. Mol. Evol.* 1998 Apr; 46(4):369-400
- Evolution of antigen binding receptors. *Annu. Rev. Immunol.* 1999; 17:109-147

cd05722 is part of a hierarchy of related CD models. Use the graphical representation to navigate this hierarchy. cd05722 is a member of the superfamily c11960.

cd05722 Sequence Cluster

Sub-family Hierarchy

Interactive Display with CDTree

Download Cn3D

Hierarchy

Interactive Display

Display: cd05722 Branch

Download CDTree

LinkOut - more resources

Sequence Alignment

Reformat Format: Compact Hypertext Row Display: All 7 rows Color Bits: 2.0 bit Type Selection: top listed sequences

gi 62204258	35	WFSTEPSDTLA.[5].VLLNCSVHS.[3].AKIEWKKGDFLSL.[8].LADGSLLSVVVHSK.[1].NKPDEGVYQCV	111
gi 110645196	48	YFLTEPVDVPT.[5].AVLNCSAYA.[3].PKIEWKKGDTFLNL.[8].LPDGSLLITSVVHSK.[1].NKPDEGVYQCV	124
gi 113675978	28	FFIKPEPDIPTA.[5].VVLDCQARG.[3].GIRWLRNGVETTE.[6].LSNGSLLISVSVSRK	DKDDEGFYQCL 101
gi 148277558	30	SFPLESPDIIA.[5].LMLHCQVEG.[3].ISTQWRRSALVQE.[6].FTNGSLLITHFQKIK.[2].GSSIDEGDYECI	105
gi 1169233	41	RFLSEPSDAVT.[5].VLLDCSAES.[4].PVIKWKKGDIHLAL.[8].LSNGSLLIQNILHSR.[1].HKPDEGLYQCE	118
gi 10720134	20	YFLVEPMDLIS.[5].VIMNCSVYC.[3].PKIEWKKGDTLLNL.[8].LPDGSLLINSVVHSK.[1].NKPDEGVYQCV	96
gi 147903889	41	WFLSEPSDAVT.[5].VVLNCSAQS.[4].PIIKWKKGDTLLNL.[8].LPDGSFLIQNVHSR.[1].HRPDEGVYQCE	118

Citing CDD

Marchler-Bauer A et al. (2015), "CDD: NCBI's conserved domain database.", *Nucleic Acids Res.* 43(Database issue):D222-6.

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, CDD*
 - or 'many-against-one': *PSI-BLAST, DELTA-BLAST*

PSI-BLAST

- Position-Specific Iterated BLAST search
- Used to identify distantly related sequences that are possibly missed during a standard 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

Altschul et al., *Nucleic Acids Res.* 25: 3389-3402, 1997

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help Sign In (Registered)

NCBI/BLAST Home

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

BLAST Assembled Genomes

Find Genomic BLAST pages:

[Enter organism name or id--completions will be suggested] **GO**

<input type="checkbox"/> Human	<input type="checkbox"/> Rabbit	<input type="checkbox"/> Zebrafish
<input type="checkbox"/> Mouse	<input type="checkbox"/> Chimp	<input type="checkbox"/> Clawed frog
<input type="checkbox"/> Rat	<input type="checkbox"/> Guinea pig	<input type="checkbox"/> Arabidopsis
<input type="checkbox"/> Cow	<input type="checkbox"/> Fruit fly	<input type="checkbox"/> Rice
<input type="checkbox"/> Pig	<input type="checkbox"/> Honey bee	<input type="checkbox"/> Yeast
<input type="checkbox"/> Dog	<input type="checkbox"/> Chicken	<input type="checkbox"/> Microbes

Basic BLAST

Choose a BLAST program to run.

nucleotide blast	Search a nucleotide database using a nucleotide query <i>Algorithms: blastn, megablast, discontinuous megablast</i>
protein blast	Search protein database using a protein query <i>Algorithms: blastp, psi-blast, phi-blast, delta-blast</i>
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.)

- Get faster protein results with a graphical view using [SmartBLAST](#)
- Make specific primers with [Primer-BLAST](#)
- Cluster multiple sequences together with their database neighbors using [MOLE-BLAST](#)
- Find [conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IgBLAST)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align two \(or more\) sequences using BLAST \(bj2app\)](#)

Your Recent Results [New!](#)

[All Recent results...](#)

News

[Searching Whole Genome Shotgun sequences](#)

It is now much easier to search WGS (Whole Genome Shotgun) with stand-alone BLAST on your own computer.

Wed, 20 Jan 2016 10:00:00 EST

[More BLAST news...](#)

Tip of the Day

[Use Genomic BLAST to see the genomic context](#)

If you are interested in the evolution of a particular gene or gene family it is often interesting to examine the intro-exon structure even across species.

[More tips...](#)

Protein BLAST: search prot...
blast.ncbi.nlm.nih.gov/blast.cgi

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Registered]

NCBI/ BLAST/ blastp suite Standard Protein BLAST

blastn blastp blastx tblastn tblastx

Enter Query Sequence

BLASTP programs search protein databases using a protein query. [more...](#) [Reset page](#) [Bookmark](#)

Enter accession number(s), gi(s), or FASTA sequence(s) Clear [Query subrange](#)

From To

Or, upload file No file selected.

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database

Organism Exclude

Optional Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Optional

Entrez Query [YouTube](#) [Create custom database](#)

Optional Enter an Entrez query to limit search

Program Selection

Algorithm

blastp (protein-protein BLAST)

PSI-BLAST (Position-Specific Iterated BLAST)

PHI-BLAST (Pattern Hit Initiated BLAST)

DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm

Search database Reference proteins (refseq_protein) using PSI-BLAST (Position-Specific Iterated BLAST)

Show results in a new window

[Algorithm parameters](#) **Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign**

Protein BLAST: search prot...
blast.ncbi.nlm.nih.gov/blast.cgi

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Registered]

NCBI/ BLAST/ blastp suite Standard Protein BLAST

blastn blastp blastx tblastn tblastx

Enter Query Sequence

BLASTP programs search protein databases using a protein query. [more...](#) [Reset page](#) [Bookmark](#)

Enter accession number(s), gi(s), or FASTA sequence(s) Clear [Query subrange](#)

From To

Or, upload file No file selected.

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database

Organism Exclude

Optional Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Optional

Entrez Query [YouTube](#) [Create custom database](#)

Optional Enter an Entrez query to limit search

Program Selection

Algorithm

blastp (protein-protein BLAST)

PSI-BLAST (Position-Specific Iterated BLAST)

PHI-BLAST (Pattern Hit Initiated BLAST)

DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm

Search database Reference proteins (refseq_protein) using PSI-BLAST (Position-Specific Iterated BLAST)

Show results in a new window

[Algorithm parameters](#) **Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign**

The screenshot shows the 'Algorithm parameters' section of the NCBI BLAST interface. Under 'General Parameters', the 'Expect threshold' is highlighted in yellow and has a callout box pointing to it with the text 'Default = 10'. The 'PSI-BLAST Threshold' is also highlighted in yellow and has a callout box pointing to it with the text 'Default = 0.005'. Other parameters like 'Max target sequences' (500), 'Word size' (3), and 'Filters and Masking' (Low complexity regions checked) are also visible.

The screenshot displays the results of a BLAST search for query NP_002119.1. The 'Graphic Summary' section shows a sequence alignment with two specific hits: 'HMG-UBF_HMG-box' and 'HMG-box superfamily'. Below this is a 'Distribution of 117 Blast Hits on the Query Sequence' chart. A color key for alignment scores is provided: <40 (black), 40-50 (blue), 50-80 (green), 80-200 (magenta), and >=200 (red). The chart shows a high density of hits with scores >=200, indicated by a thick red bar at the bottom.

NCBI Blast:NP_002119.1 N... x

blast.ncbi.nlm.nih.gov/Blast.cgi

Graphic Summary

Descriptions

Run PSI-Blast iteration 2 with max 500 Go

Sequences producing significant alignments with E-value BETTER than threshold

Select: All None Selected:0

Alignments Download GenPept Graphics Distance tree of results Multiple alignment

Description	Max score	Total score	Query cover	E value	Ident	Accession	Select for PSI blast	Used to build PSSM
<input type="checkbox"/> high mobility group protein B1 [Bos taurus]	310	310	78%	7e-106	100%	NP_788765.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Mus musculus]	310	310	78%	7e-106	100%	NP_034569.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Homo sapiens]	310	310	78%	7e-106	100%	NP_002119.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Sus scrofa]	308	308	78%	5e-105	99%	NP_001004034.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 1 like [Rattus norvegicus]	308	308	78%	7e-105	99%	NP_001102843.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Gallus gallus]	299	299	78%	2e-101	96%	NP_990233.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Xenopus tropicalis]	294	294	78%	2e-99	92%	NP_989226.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 1 [Xenopus laevis]	290	290	77%	6e-98	92%	NP_001080836.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein-1 [Xenopus laevis]	280	280	77%	7e-94	90%	NP_001081794.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Danio rerio]	266	266	77%	8e-89	87%	NP_001092721.2	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Danio rerio]	266	266	77%	2e-88	87%	NP_955849.2	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG-X protein [Xenopus laevis]	262	262	78%	5e-87	84%	NP_001079578.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> High mobility group-T protein [Salmo salar]	264	264	77%	5e-86	84%	NP_001140081.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Salmo salar]	259	259	77%	6e-86	84%	NP_001133101.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B2 [Gallus gallus]	257	257	78%	5e-85	85%	NP_990817.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group-T protein [Cnorchynchus mykiss]	257	257	77%	5e-85	83%	NP_001118186.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B2 [Homo sapiens]	252	252	78%	3e-83	86%	NP_002120.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B2 [Macaca fascicularis]	252	252	78%	4e-83	86%	NP_001271844.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B2 [Rattus norvegicus]	251	251	78%	1e-82	86%	NP_068883.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 2 [Xenopus laevis]	250	250	78%	4e-82	82%	NP_001079387.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>

NCBI Blast:NP_002119.1 N... x

blast.ncbi.nlm.nih.gov/Blast.cgi

<input type="checkbox"/> HMG domain-containing protein 4 [Homo sapiens]	47.8	47.8	21%	3e-05	43%	NP_001003681.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> uncharacterized protein LOC399067 [Xenopus laevis]	47.4	94.3	62%	3e-05	37%	NP_001083698.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily E member 1-related [Xer]	47.4	94.3	62%	3e-05	32%	NP_988941.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG domain-containing protein 4 [Bos taurus]	47.4	47.4	19%	3e-05	45%	NP_001095328.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> TOX high mobility group box family member 4-A [Xenopus laevis]	47.4	47.4	24%	3e-05	38%	NP_001086364.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> TOX high mobility group box family member 4 [Xenopus tropicalis]	47.4	47.4	24%	3e-05	38%	NP_001090624.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> TOX high mobility group box family member 4-B [Xenopus laevis]	47.4	47.4	24%	4e-05	38%	NP_001084977.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor CI-HMG20 [Ciona intestinalis]	46.2	46.2	23%	7e-05	39%	NP_001121587.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG domain-containing protein 4 [Danio rerio]	46.2	46.2	19%	7e-05	45%	NP_001120984.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily E member 1-related [Hoy]	45.4	89.7	62%	1e-04	37%	NP_006330.2	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG box-containing protein 4 [Xenopus laevis]	45.4	45.4	19%	1e-04	45%	NP_001082746.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor A_mitochondrial precursor [Rattus norvegicus]	45.1	88.2	64%	1e-04	29%	NP_112616.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily E member 1-related [Bor]	45.1	89.7	62%	1e-04	37%	NP_001033143.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily E member 1-related [Mu]	45.1	45.1	28%	1e-04	37%	NP_034570.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily E member 1-related [Rai]	45.1	87.8	62%	2e-04	37%	NP_001102201.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein 20A [Xenopus tropicalis]	45.1	45.1	34%	2e-04	31%	NP_001006760.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG box-containing protein 4 [Xenopus tropicalis]	45.1	45.1	19%	2e-04	45%	NP_001025555.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> protein polybromo-1 [Gallus gallus]	45.1	45.1	32%	3e-04	36%	NP_990496.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein 20A [Xenopus laevis]	44.3	44.3	34%	3e-04	31%	NP_001087141.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor A_mitochondrial precursor [Sus scrofa]	43.1	43.1	64%	5e-04	29%	NP_001123663.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein 20A [Gallus gallus]	43.5	43.5	34%	6e-04	31%	NP_001025565.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor protein [Ciona intestinalis]	43.5	43.5	33%	6e-04	36%	NP_001071666.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor protein [Ciona intestinalis]	43.1	43.1	33%	7e-04	36%	NP_001072029.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor A_mitochondrial [Esoc lucius]	42.7	42.7	34%	9e-04	27%	NP_001297981.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor protein [Ciona intestinalis]	43.1	43.1	21%	0.001	40%	NP_001071952.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor A_mitochondrial [Danio rerio]	42.7	42.7	22%	0.001	31%	NP_001070857.1	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Run PSI-Blast iteration 2 with max 500 Go

NCBI Blast:NP_002119.1 N... x

blast.ncbi.nlm.nih.gov/Blast.cgi

Graphic Summary

Descriptions

Run PSI-Blast iteration 3 with max: 500 Go

Sequences producing significant alignments with E-value BETTER than threshold

Select: All None Selected:0 Yellow: sequences scoring below threshold on previous iteration

Alignments Download GenPept Graphics Distance tree of results Multiple alignment

Description	Max score	Total score	Query cover	E value	Ident	Accession	Select for PSI blast	Used to build PSSM
<input type="checkbox"/> high mobility group protein B1 [Bos taurus]	250	250	78%	3e-82	100%	NP_788785.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Mus musculus]	250	250	78%	3e-82	100%	NP_034569.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Homo sapiens]	250	250	78%	3e-82	100%	NP_002119.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Sus scrofa]	250	250	78%	3e-82	99%	NP_001004034.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> High mobility group-T protein [Salmo salar]	254	254	77%	5e-82	84%	NP_001140081.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 1 like [Rattus norvegicus]	247	247	78%	3e-81	99%	NP_001102843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Gallus gallus]	246	246	78%	8e-81	96%	NP_990233.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Danio rerio]	236	236	77%	6e-77	87%	NP_001092721.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Xenopus tropicalis]	236	236	78%	6e-77	92%	NP_989226.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 1 [Xenopus laevis]	235	235	77%	2e-76	92%	NP_001080836.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Salmo salar]	234	234	77%	6e-76	84%	NP_001133101.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein-1 [Xenopus laevis]	232	232	77%	4e-75	90%	NP_001081794.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG-X protein [Xenopus laevis]	232	232	78%	4e-75	84%	NP_001079576.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B1 [Danio rerio]	230	230	77%	1e-74	87%	NP_955849.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B3 [Danio rerio]	230	230	77%	2e-74	66%	NP_001116308.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> High mobility group protein B3 [Salmo salar]	225	225	77%	2e-72	67%	NP_001133971.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B2 [Gallus gallus]	224	224	78%	3e-72	85%	NP_990817.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group box 3 [Callorhynchus milii]	223	223	77%	6e-72	77%	NP_001279444.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group protein B3 [Danio rerio]	223	223	75%	8e-72	68%	NP_001017769.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> high mobility group-T protein [Oncorhynchus mykiss]	223	223	77%	9e-72	83%	NP_001181868.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

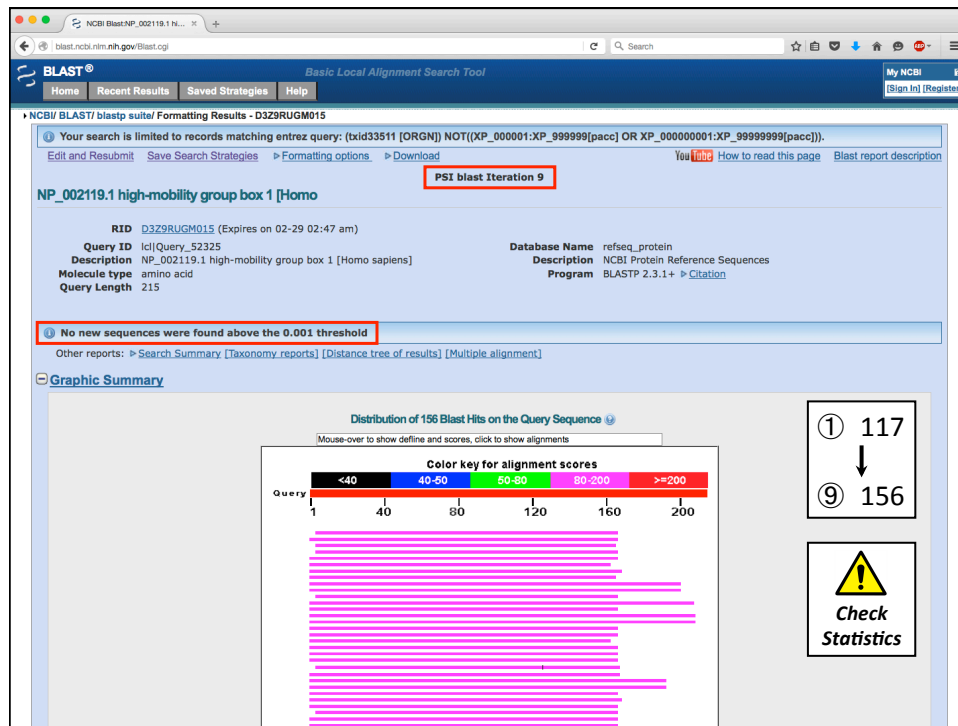
NCBI Blast:NP_002119.1 N... x

blast.ncbi.nlm.nih.gov/Blast.cgi

②... ③... ④...

<input type="checkbox"/> lymphoid enhancer factor 1 [Xenopus laevis]	46.9	46.9	32%	4e-05	21%	NP_001133001.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> lymphoid enhancer-binding factor 1 [Xenopus tropicalis]	46.9	46.9	32%	4e-05	21%	NP_001230763.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> lymphoid enhancer factor XLEF-1B [Xenopus laevis]	46.9	46.9	32%	4e-05	21%	NP_001092003.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> lymphoid enhancer-binding factor 1 [Xenopus laevis]	46.5	46.5	35%	5e-05	22%	NP_001082124.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> basic helix-loop-helix and HMG box domain-containing protein 1 [Homo sapiens]	46.9	46.9	36%	5e-05	22%	NP_001297053.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> PMS1 protein homolog 1 [Danio rerio]	46.1	46.1	21%	8e-05	32%	NP_958476.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> WD repeat and HMG-box DNA-binding protein 1 [Xenopus laevis]	46.1	46.1	52%	1e-04	26%	NP_001081495.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like protein [Saccoglossus kowalevskii]	45.7	45.7	27%	1e-04	29%	NP_001158464.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1 isoform 2 [Mus musculus]	45.3	45.3	27%	2e-04	28%	NP_033359.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1 [Rattus norvegicus]	45.3	45.3	27%	2e-04	28%	NP_001101335.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1 [Homo sapiens]	45.3	45.3	27%	2e-04	28%	NP_112573.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1 isoform 1 [Mus musculus]	45.3	45.3	27%	2e-04	28%	NP_001073290.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG protein TcfLef [Strongylocentrotus purpuratus]	45.3	45.3	37%	2e-04	19%	NP_999840.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1-A [Danio rerio]	44.9	44.9	30%	2e-04	28%	NP_571344.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> Wolf-Hirschhorn syndrome candidate 1 [Xenopus laevis]	44.9	44.9	28%	2e-04	23%	NP_001084939.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG box-containing protein 1 [Danio rerio]	44.9	44.9	23%	2e-04	29%	NP_001019602.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7 isoform 4 [Homo sapiens]	44.6	44.6	35%	2e-04	28%	NP_963965.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7 isoform 2 [Homo sapiens]	44.6	44.6	35%	2e-04	28%	NP_963963.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1 [Oryzias latipes]	44.9	44.9	30%	2e-04	26%	NP_001239177.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7 isoform 1 [Homo sapiens]	44.6	44.6	35%	2e-04	28%	NP_003193.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7 isoform 3 [Homo sapiens]	44.2	44.2	35%	3e-04	28%	NP_001128323.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1-B [Danio rerio]	44.2	44.2	27%	4e-04	27%	NP_571371.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> transcription factor 7-like 1-A [Xenopus laevis]	43.4	43.4	26%	6e-04	27%	NP_001081483.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG box domain-containing 3 [Xenopus laevis]	43.8	43.8	19%	7e-04	36%	NP_001089484.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/> HMG domain-containing protein 3 [Xenopus tropicalis]	43.8	43.8	19%	7e-04	36%	NP_001120640.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Run PSI-Blast iteration 3 with max: 500 Go



DELTA-BLAST

- Method different from that used by PSI-BLAST

Step 1: Align the query against conserved domains derived from CDD
Step 2: Compute PSSM
Step 3: Search sequence databases using PSSM as the query

- Intended to improve homology detection
- Produces high-quality alignments, even at low levels of sequence similarity
- Dependent on homologous relationships captured within CDD

Boratyn et al., *Biology Direct* 7: 12, 2012

Multiple Sequence Alignment: A Quick Primer



Why do multiple sequence alignments?

- Identify conserved regions, patterns, and domains
 - Experimental design
 - Predicting structure and function
 - Identifying new members of protein families
- Provide basis for:
 - Predicting secondary structure
 - Performing phylogenetic analyses, thereby determining evolutionary relationships (inferring homology)
 - Generating position-specific scoring matrices for use with sensitive sequence search methods



Overarching Considerations

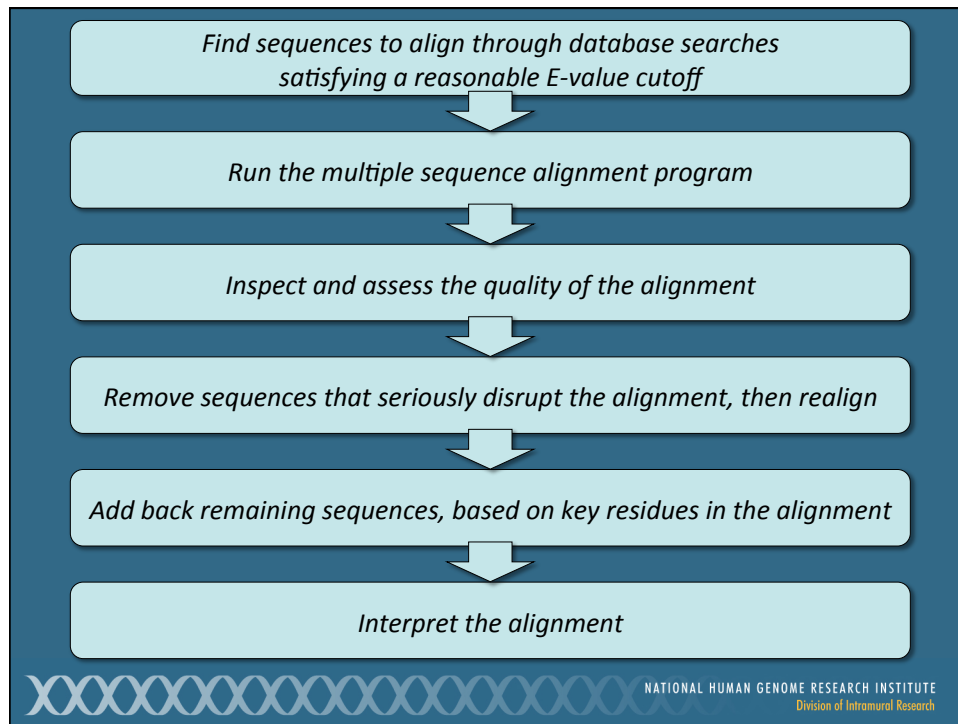
- Absolute sequence similarity
Create the alignment by lining up as many common characters as possible
- Conservation
Take into account residues that can substitute for one another and not adversely affect the function of the protein
- Structural similarity
Knowledge of the secondary or tertiary structure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



Protein vs. Nucleotide Multiple Sequence Alignments

- Concentrate on the protein level rather than on the nucleotide level
- Protein alignments tend to be more informative
- Less prone to inaccurate alignment ('20 vs. 4')
- Can 'translate back' to nucleotide sequences *after* doing the alignment





Selecting the Sequences

1. 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 can sometimes be intractable
 - Good starting point: 10-15 sequences
 - Ballpark upper limit: 50-100 sequences

Selecting the Sequences

2. Sequences should be of about the same length
3. Trim sequences down, so as to only use regions that have been deemed similar by either:
 - Pairwise search methods such as BLAST
 - Profile-based search methods such as PSI-BLAST



Selecting the Sequences

4. Consider the degree of similarity in the sequence set, *depending on what question is being asked*
 - Use closely-related sequences to determine 'required' (highly conserved) 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 dissimilar



Inspection: An Iterative Process

- Perform alignment on small set of sequences
- Examine the quality of the alignment, looking for:
 - Conservation of residues across alignment
 - Conservation of physicochemical properties
 - Relatively neat block-type structure
 - Excessive numbers of gaps
- If alignment is good, can add new sequences to data set, then realign
- If alignment is not good, remove any sequences that result in the inclusion of long gaps, then realign



Inspection: An Iterative Process

- Use visualization tools to identify 'key residues' and 'problem regions'
- Cross-check against 'expertly created' multiple sequence alignments available online
- Use any available information from solved X-ray or NMR structures to nail down structurally important regions and to assess where gaps can (or cannot) be tolerated



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



Clustal Omega

- Allows for automatic multiple alignment of nucleotide or amino acid sequences
- Aligns data sets quickly and easily
- Can align sequences against a pre-existing alignment (an 'external profile')
- Can bias the location of gaps, based on known structural information
- Works with Jalview, a Java applet for viewing and manipulating results

Sievers et al., Mol. Syst. Biol. 7: 539, 2011



Progressive Alignment

- Align two sequences at a time, starting with the two most related sequences
- Gradually build up the multiple sequence alignment by adding additional (less-related) sequences to the alignment
- Uses protein scoring matrices and gap penalties to calculate alignments having the best score
- Major advantages of method
 - Generally fast
 - Alignments generally of high quality



Clustal Omega Output

- Pairwise alignment scores
- Multiple sequence alignment
- Cladogram
 - Tree that is assumed to be an *estimate* of a phylogeny
 - Branches are of equal length
 - Cladograms can 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 a phylogeny
 - Branches are *not* of equal length
 - Branch lengths proportional to the amount of inferred evolutionary change



Clustal Omega 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



Clustal Omega Conservation Patterns

Interpretation is empirical — there is no parallel to the E-values seen in BLAST searches to assess statistical significance

- * entirely conserved column
(want in at least 10% of positions)
- ⋮ conserved
(strongly similar properties)
- semi-conserved
(weakly similar properties)



https://www.ebi.ac.uk/Tools/msa

Clustal Omega

Input form | Web services | Help & Documentation

Tools > Multiple Sequence Alignment > Clustal Omega

Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our [pairwise sequence alignment tools](#).

STEP 1 - Enter your input sequences

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

```
>FO5B_MOUSE Protein fosB
MFOAFPGDYDSGRSCSSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA
ITTSQDLQWLQPTLISMAQSQGQLASQPPAVDPYDMPGTSYSTPGLSAYSTGGASGS
GGPSTSTTTSGPVSARPARARPRHPRETLTPEEEKRVRRERNKJAAKCRNRRRELT
DRLCAETDQLEEKAELESJAELOKKEKLEFLVAKHKGCKIPYEEQGPGLAEVYD
LPGSTSAKEDGFGLWLPPLPPPLPFQSSRDAPPNLTASLFTHSEVQLGDPPFVVSPLY
TSSFVLTCPEVSFAQAQRTSGSEQSPDPLNSPSSLAL
```

Or, [upload](#) a file: [Browse...](#) No file selected.

STEP 2 - Set your parameters

OUTPUT FORMAT: Clustal w/o numbers

The default settings will fulfill the needs of most users and, for that reason, are not visible.

More options... (Click here, if you want to view or change the default settings.)

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

Submit

If you plan to use these services during a course please [contact us](#).

Please read the [FAQ](#) before seeking help from our support staff.

https://www.ebi.ac.uk/Tools/msa

Clustal Omega

Input form | Web services | Help & Documentation

Tools > Multiple Sequence Alignment > Clustal Omega

Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our [pairwise sequence alignment tools](#).

STEP 1 - Enter your input sequences

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

```
>FO5B_MOUSE Protein fosB
MFOAFPGDYDSGRSCSSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA
ITTSQDLQWLQPTLISMAQSQGQLASQPPAVDPYDMPGTSYSTPGLSAYSTGGASGS
GGPSTSTTTSGPVSARPARARPRHPRETLTPEEEKRVRRERNKJAAKCRNRRRELT
DRLCAETDQLEEKAELESJAELOKKEKLEFLVAKHKGCKIPYEEQGPGLAEVYD
LPGSTSAKEDGFGLWLPPLPPPLPFQSSRDAPPNLTASLFTHSEVQLGDPPFVVSPLY
TSSFVLTCPEVSFAQAQRTSGSEQSPDPLNSPSSLAL
```

Or, [upload](#) a file: [Browse...](#) No file selected.

STEP 2 - Set your parameters

OUTPUT FORMAT: Clustal w/o numbers

DEALIGN INPUT SEQUENCES	MBED-LIKE CLUSTERING GUIDE-TREE	MBED-LIKE CLUSTERING ITERATION	NUMBER of COMBINED ITERATIONS
yes	yes	yes	default(0)
MAX GUIDE TREE ITERATIONS	MAX HMM ITERATIONS	ORDER	
default	default	input	

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

Submit

If you plan to use these services during a course please [contact us](#).

Alignments < Clustal Omega

Alignments Result Summary Phylogenetic Tree Submission Details

Download Alignment File Hide Colors Send to ClustalW2

CLUSTAL O (1.2.1) multiple sequence alignment

Residue	Colour	Property
AVFFMLW	RED	Small (small+ hydrophobic (incl.aromatic -Y))
DE	BLUE	Acidic
RK	MAGENTA	Basic - H
STYHCNGQ	GREEN	Hydroxyl + sulfhydryl + amine + G
Others	Grey	Unusual amino/imino acids etc

```

FOSB_MOUSE  -MPQAFPGDYDSGS--KCSSSPSA--ESQYLSVDSFG
FOSB_HUMAN  -MPQAFPGDYDSGS--KCSSSPSA--ESQYLSVDSFG
FOSB_CHICK  MMYQGFAGYEAFSSKCSSASFAGDSLTYYPSPADGFSMGSFVNSQDCTDLAVSSANF 60
FOSB_RAT    MMFSGFNADYEASSKCSSASFAGDSLTYHSPADGFSMGSFVNTQDFCADLVSSANF 60
FOSB_MOUSE  MMFSGFNADYEASSKCSSASFAGDSLTYHSPADGFSMGSFVNTQDFCADLVSSANF 60
          *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *::: *:::
FOSB_MOUSE  VPTVTAITTSQDLQWLQPTLISSMAQSQGQLASQPPAVDPYDMPGT--SYSTPGLS 110
FOSB_HUMAN  VPTVTAITTSQDLQWLQPTLISSMAQSQGQLASQPPAVDPYDMPGT--SYSTPGLS 110
FOSB_CHICK  VPTVTAITSPDLQWLQPTLISSVAPSQN---GHPYGVPAAPPAAYSRAVIL 112
FOSB_RAT    IPTVTAITSPDLQWLQPTLVSSVAPSQTA---PHPYGLPTFSTGAYARAGVVK 113
FOSB_MOUSE  IPTVTAITSPDLQWLQPTLVSSVAPSQTA---PHPYGLPTQSAGAYARAGVVK 113
          *****
FOSB_MOUSE  AYSTGGASGGGPGSTSTTSGPVSANPARAPRRPREETLTFEEEEKRVVREINLAAA 170
FOSB_HUMAN  GYSSGGASGGGPGSTSTTSGPVSANPARAPRRPREETLTFEEEEKRVVREINLAAA 170
FOSB_CHICK  KA---FG---GAGSIRGRGKVEQLSPEEERKRIKREINMAAA 151
FOSB_RAT    TM---SG---GAAQSIGRGRGKVEQLSPEEERKRIKREINMAAA 152
FOSB_MOUSE  TV---SG---GAAQSIGRGRGKVEQLSPEEERKRIKREINMAAA 152
          *****
FOSB_MOUSE  KCNRRRELDTDLQAEETDLEEEIAELESEIAELQNEKELEFLVLAHHPGCKIPYEEGP 230
FOSB_HUMAN  KCNRRRELDTDLQAEETDLEEEIAELESEIAELQNEKELEFLVLAHHPGCKIPYEEGP 230
FOSB_CHICK  KCNRRRELDTDLQAEETDLEEEISALQAEIANLLKEKELEFLVLAHHPACKMPLELRF 211
FOSB_RAT    KCNRRRELDTDLQAEETDLEEEISALQAEIANLLKEKELEFLVLAHHPACKIPNDLGF 212
FOSB_MOUSE  KCNRRRELDTDLQAEETDLEEEISALQAEIANLLKEKELEFLVLAHHPACKIPDDLGF 212
          *****
FOSB_MOUSE  GGGPLA-EVNDLPG---STSAKEDGFGWLLPPPPPPPLPFG--- 267
FOSB_HUMAN  GGGPLA-EVNDLPG---SAPAKEDGFGWLLPPPPPPPLPFG--- 267
FOSB_CHICK  SEELAAATLDLGA--PSPAAAEAFALPLMTAPPVVP--KEPSGGGLEKRAEPFDE 266
FOSB_RAT    PEEMSVTS-LDLTGGLPEATPESEEAFTLPLNDPEPKSLEPVNINMELKAEPFDD 271
FOSB_MOUSE  PEEMSVAS-LDLTGGLPEATPESEEAFTLPLNDPEPKSLEPVNINMELKAEPFDD 271
          *****
FOSB_MOUSE  ---SSVDAPPN-I-TASLFTHS---EVQVLGDPPF- 294
FOSB_HUMAN  ---TSQDAPPN-I-TASLFTHS---EVQVLGDPPF- 294
FOSB_CHICK  LLFSAGPR--EASRVSFMDLPGASSFYASDWEPLGAGSG--GELEPLCTPVVT 316
    
```

Phylogenetic Tree < Clustal

Alignments Result Summary **Phylogenetic Tree** Submission Details

Phylogenetic Tree

This is a Neighbour-joining tree without distance corrections.

Download Phylogenetic Tree File

```

(
(
(
FOSB_MOUSE:0.01854,
FOSB_HUMAN:0.02288)
:0.35561,
FOSB_CHICK:0.11070)
:0.11115,
FOSB_RAT:0.01948,
FOSB_MOUSE:0.01210);
    
```

Phylogram

Branch length: Cladogram Real

FOSB_MOUSE 0.01854
 FOSB_HUMAN 0.02288
 FOSB_CHICK 0.1107
 FOSB_RAT 0.01948
 FOSB_MOUSE 0.0121

EMBL-EBI Services Research Training Industry About us

News By topic Overview Overview Overview Overview
 Brochures By name (A-Z) Publications Train at EBI Members Area Leadership
 Contact us Help & Support Research groups Train outside EBI Workshops Funding
 Intranet Postdocs & PhDs Train online SME Forum Background

Phylogenetic Tree < Clustal ...

Alignments | **Result Summary** | **Phylogenetic Tree** | Submission Details

Phylogenetic Tree

This is a Neighbour-joining tree without distance corrections.

Download Phylogenetic Tree File

```
(
(
(
FOSB_MOUSE:0.01854,
FOSB_HUMAN:0.02288)
:0.35561,
FOS_CHICK:0.11070)
:0.11115,
FOS_RAT:0.01948,
FOS_MOUSE:0.01210);
```

Phylogram

Branch length: Cladogram Real

- FOSB_MOUSE 0.01854
- FOSB_HUMAN 0.02288
- FOS_CHICK 0.1107
- FOS_RAT 0.01948
- FOS_MOUSE 0.0121

EMBL-EBI | Services | Research | Training | Industry | About us

EMBL-EBI | Services | Research | Training | About us

Clustal Omega

Input form | Web services | Help & Documentation | Share | Feedback

Tools > Multiple Sequence Alignment > Clustal Omega

Results for job clustalo-I20160227-210614-0107-78029139-es

Alignments | **Result Summary** | Phylogenetic Tree | Submission Details

Input Sequences

clustalo-I20160227-210614-0107-78029139-es.input

Tool Output

clustalo-I20160227-210614-0107-78029139-es.output

Alignment in CLUSTAL format with base/residue numbering

clustalo-I20160227-210614-0107-78029139-es.clustal_num

Phylogenetic Tree

clustalo-I20160227-210614-0107-78029139-es.ph

Percent Identity Matrix

clustalo-I20160227-210614-0107-78029139-es.pim

Jalview

Start Jalview

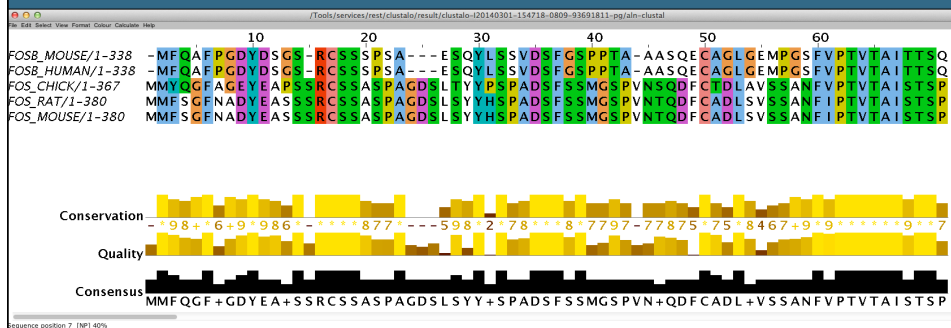
EMBL-EBI | Services | Research | Training | Industry | About us

Jalview

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

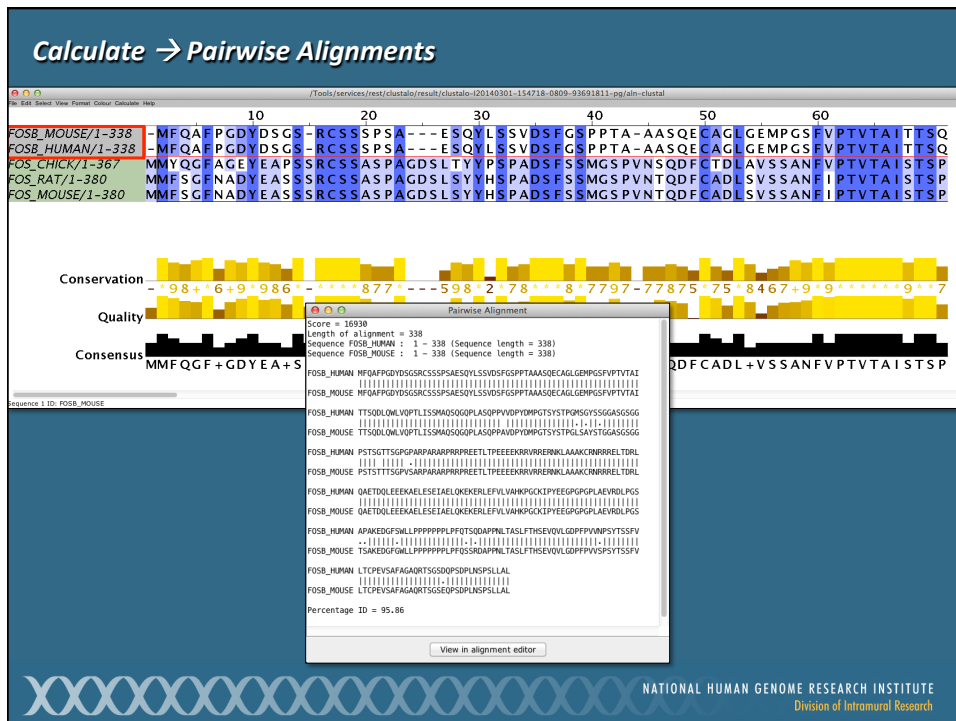
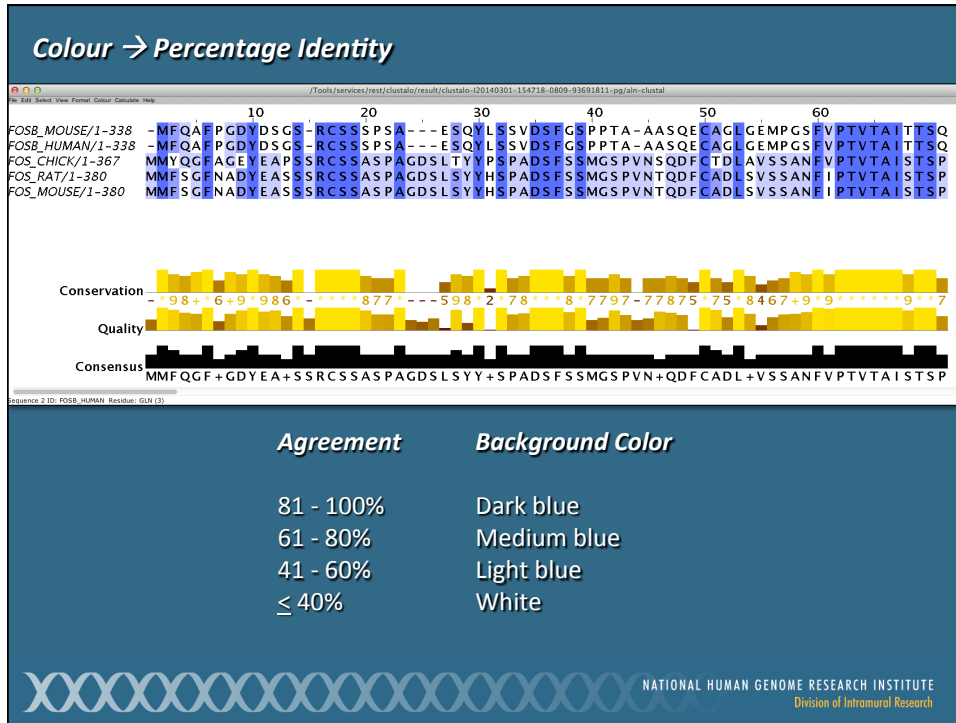


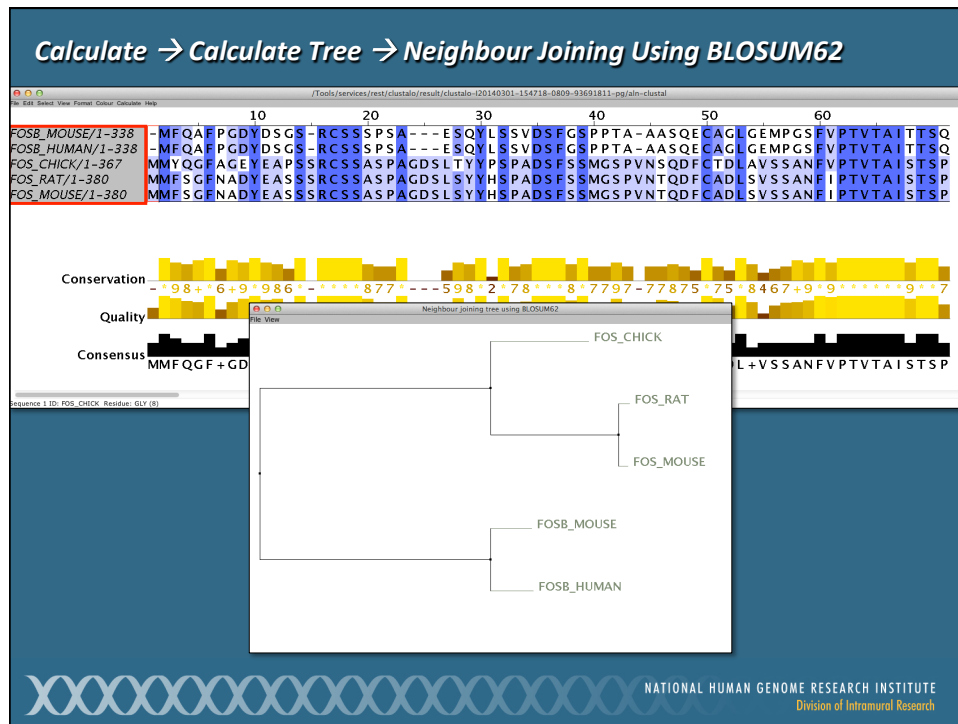
Default view



- Conservation** Conservation of total alignment (indication of percent identity)
- Quality** Alignment quality, based on BLOSUM scores
- Consensus** Based on percent identity





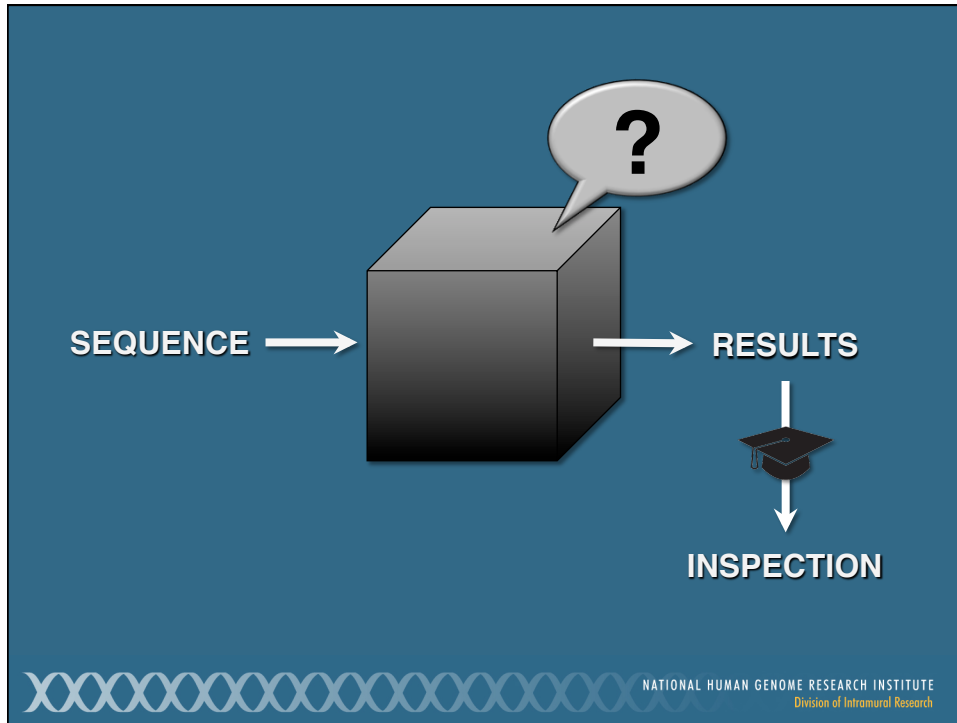


T-COFFEE

- Combines sequence, profile, and structural information
 - Protein structures
 - RNA secondary structures
- Specialized algorithm for aligning transmembrane proteins, non-coding RNAs, and homologous promoter regions
- Can combine output from other methods into a single 'master alignment'
- Freely available at <http://tcoffee.org>



Magis et al., *Methods Mol. Biol.* 1079: 117-129 (2014)



Online Training Resources

Suggested curriculum tracks tailored to individual needs

- Bioinformatic Analysis
- Data Mining
- Bioinformatics Tools
- Bioinformatics Systems
- Computational Biology

Searls, *PLoS Comput. Biol.* 8: e1002632, 2012
Searls, *PLoS Comput. Biol.* 10: e1003662, 2014

An Online Bioinformatics Curriculum
David B. Searls¹
Independent Consultant, Philadelphia, Pennsylvania, United States of America

A New Online Computational Biology Curriculum
David B. Searls¹
Independent Consultant, Philadelphia, Pennsylvania, United States of America

Current Topics in Genome Analysis 2016

Next Lecture
March 16, 2014

Regulatory and Epigenetic Landscapes of Mammalian Genomes

Laura Elnitski, Ph.D.
National Human Genome Research Institute
National Institutes of Health



NIH Intramural Research Program
Our Research Changes Lives

one program
many people
infinite possibilities

irp.nih.gov