




NATIONAL HUMAN GENOME RESEARCH INSTITUTE Division of Intramural Research




Current Topics in Genome Analysis 2014
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



Current Topics in Genome Analysis 2014
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

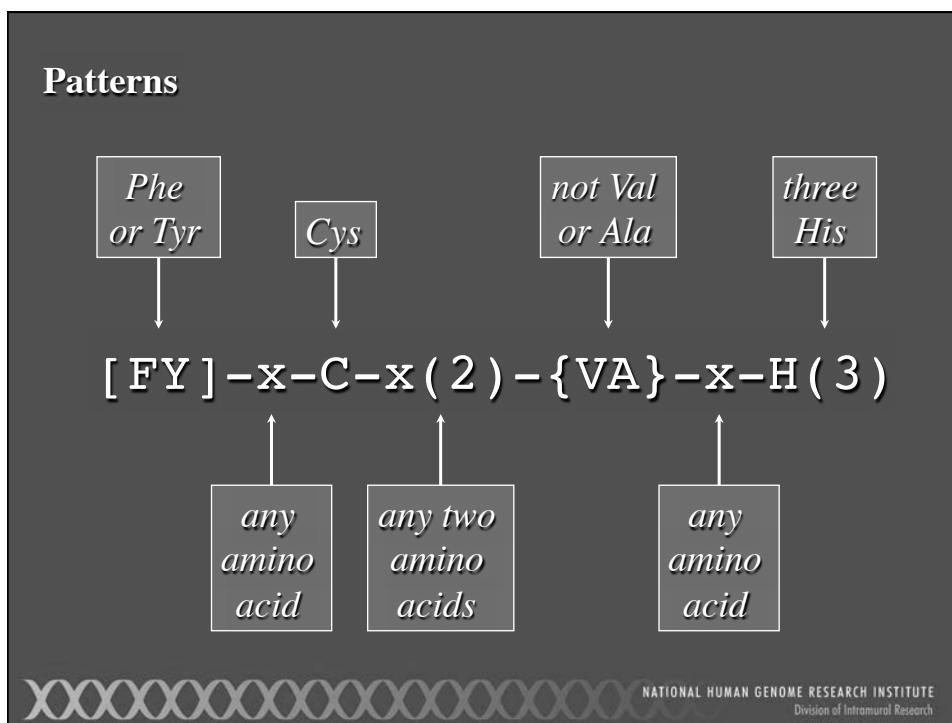
APHIIVATPG
GCEIVATPG
GVEICATPG
GVDILIGTTPG
RPHIIVATPG
KPHIIVATPG
KVQLIATPG
RPDIVIATPG
APHIIVGTPG
APHIIVGTPG
GCHVVIATPG
NQDIVVATPG

- 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	0	13	0	0	-12	13	0	0	-5	-3	-1	-2	23	2	-2	12	11	17	-31	-8	1
H	5	24	-12	29	25	-20	8	32	-9	9	-10	-9	22	7	30	10	0	4	-8	-20	-7	27
I	-1	-12	6	-13	-11	33	-12	-13	63	-11	40	29	-15	-9	-14	-15	-6	7	50	-17	8	-11
V	3	-11	1	-11	-9	22	-3	-11	46	-9	37	30	-13	-3	-9	-13	-6	6	50	-19	2	-8
V	5	-9	9	-9	-9	19	-1	-13	57	-9	35	26	-13	-2	-11	-13	-4	9	58	-29	0	-9
A	54	15	12	20	17	-24	44	-6	-4	-1	-11	-5	12	19	9	-13	21	19	9	-39	-20	10
T	40	20	20	20	20	-30	40	-10	20	20	-10	0	20	30	-10	-10	30	150	20	-60	-30	10
P	-11	0	7	0	0	-11	-11	-2	0	-16	-11	-2	-2	89	17	17	24	22	9	-50	-48	12
G	-70	-66	20	-76	-30	-68	150	-20	-30	-10	-50	-30	40	30	20	-30	60	40	20	-100	-70	30





Pfam

- Collection of multiple alignments of protein domains and conserved protein regions that 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

Finn et al., *Nucleic Acids Res.* 42: D222-D230, 2014

Pfam

- Pfam A
 - Based on *curated* multiple alignments (“seed alignment”)
 - HMMER used to find all detectable protein sequences belonging to the family (Eddy, 2011)
 - 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



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

Sequences Used in Examples

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

The screenshot shows the NHGRI website with a navigation bar and a main content area. The main content area is titled 'Current Topics in Genome Analysis 2014' and 'Protein and Nucleotide Sequences for Analysis'. It features a section for 'Current Topics in Genome Analysis 2014' with a 'Current Home' link. Below this, there are two sections: 'Current Topics in Genome Analysis 2014' and 'Protein and Nucleotide Sequences for Analysis'. The 'Protein and Nucleotide Sequences for Analysis' section contains a 'BLAST' section with a 'BLAST 2 Sequences' subsection. The 'BLAST 2 Sequences' subsection shows a 'Query' sequence and a 'Subject' sequence, both of which are protein sequences. The 'Query' sequence is a long string of amino acids, and the 'Subject' sequence is a long string of amino acids. The 'BLAST' section also includes a 'BLAST' subsection with a 'BLAST' subsection. The 'BLAST' subsection shows a 'Query' sequence and a 'Subject' sequence, both of which are protein sequences. The 'Query' sequence is a long string of amino acids, and the 'Subject' sequence is a long string of amino acids.



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

<http://pfam.sanger.ac.uk>

wellcome trust
sanger
institute

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam
keyword search

Pfam 27.0 (March 2013, 14831 families)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

QUICK LINKS

- SEQUENCE SEARCH**
- VIEW A PFAM FAMILY**
- VIEW A CLAN**
- VIEW A SEQUENCE**
- VIEW A STRUCTURE**
- KEYWORD SEARCH**
- JUMP TO**

YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...

- Analyze your protein sequence for Pfam matches
- View Pfam family annotation and alignments
- See groups of related families
- Look at the domain organisation of a protein sequence
- Find the domains on a PDB structure
- Query Pfam by keywords

Enter any type of accession or ID to jump to the page for a Pfam family or clan, UniProt sequence, PDB structure, etc.

Or view the help pages for more information

Recent Pfam blog posts

Short-term Pfam position available. (posted 7 February 2014)

We have just advertised a 9-month maternity cover position in Pfam. We are looking for a skilled Bioinformatician to help us take Pfam into its next phase of development as we become more integrated into the European Bioinformatics Institute (EMBL-EBI). Essential knowledge, skills and experience: Degree in Science with relevant experience Computer literacy (unix experience) [...]

Join Rfam, see the world (posted 31 January 2014)

<http://pfam.sanger.ac.uk>

wellcome trust
sanger
institute

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam
keyword search

Pfam 27.0 (March 2013, 14831 families)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

QUICK LINKS

- SEQUENCE SEARCH**
- VIEW A PFAM FAMILY**
- VIEW A CLAN**
- VIEW A SEQUENCE**
- VIEW A STRUCTURE**
- KEYWORD SEARCH**
- JUMP TO**

ANALYZE YOUR PROTEIN SEQUENCE FOR PFAM MATCHES

Paste your protein sequence here to find matching Pfam families.

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.

Recent Pfam blog posts

Short-term Pfam position available. (posted 7 February 2014)

We have just advertised a 9-month maternity cover position in Pfam. We are looking for a skilled Bioinformatician to help us take Pfam into its next phase of development as we become more integrated into the European Bioinformatics Institute (EMBL-EBI). Essential knowledge, skills and experience: Degree in Science with relevant experience Computer literacy (unix experience) [...]

Join Rfam, see the world (posted 31 January 2014)

Rfam is recruiting! We are currently recruiting an RNA informatician to join our team. We're looking for someone really enthusiastic about RNA and who's interested in working with Rfam as we move to

wellcome trust
sanger
institute

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam
keyword search Go

Search Pfam

Sequence search

Find Pfam families within your sequence of interest. Paste your **protein** or **DNA** sequence into the box below to have it searched for matching Pfam families. [More...](#)

Sequence

Protein sequence options

Cut-off ☐ Gathering threshold
☒ Use E-value

E-value

Search for PfamBs ☐ Note that we search only the 20,000 largest Pfam-B families

Submit [Reset](#) [Example protein sequence](#) [Example DNA sequence](#)

Comments or questions on the site? Send a mail to pfam-help@ebi.ac.uk. Our [cookie policy](#).
The Wellcome Trust

wellcome trust
sanger
institute

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam
keyword search Go

Sequence search results

[Show](#) the detailed description of this results page.

We found **3** Pfam-A matches to your search sequence (**1** significant and **2** insignificant). You did not choose to search for Pfam-B matches.

[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		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

Insignificant Pfam-A Matches

[Show or hide all alignments.](#)

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					
COG7	Golgi complex component 7 (COG7)	Family	CL0294	188	309	247	296	313	362	766	11.0	0.069	n/a	Show
Sec8_exocyst	Sec8 exocyst complex component specific	Domain	CL0295	246	286	249	277	42	70	142	13.3	0.047	n/a	Show

Comments or questions on the site? Send a mail to pfam-help@ebi.ac.uk. Our [cookie policy](#).
The Wellcome Trust

[illegible]

PFfam Family p450 (PF00067)

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Summary Cytochrome P450

Domain organisation

Alignments

HMM logo

Trees

Curation & model

Species

Interactions

Structures

Jump to... ↓

PROSITE PROSITE

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

Wikipedia: Cytochrome P450 Plam 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 EDSX motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with mitochondrial membranes, their general enzymatic function is to catalyse regioselective and stereospecific oxidation of non-volatile hydrocarbons at physiological temperatures [6].

Literature references

- Graham-Lorance S, Amarnah B, White RE, Petersen JA, Simpson LR. Protein Sci 1995;4:1065-1060.: A three-dimensional model of aromatase cytochrome P450. PubMed:7549871 Q EPNC:7549871 Q
- Degtyarenko KN, Archakov AI., FEBS Lett 1993;332:1-6.: Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. PubMedID:6405421 Q EPNC:6405421 Q
- Nelson DR, Kamakatsi T, Waxman DJ, Guengerich PP, Estabrook RW, Feyereisen R, Gonzalez FJ, Coon MJ, Gunsalus IC, Gotoh O, et al.; DNA Cell Biol 1993;12:1-S1: 1-S1.: The P450 superfamily: update on new sequences, gene mapping, accession numbers, early trivial names of enzymes, and nomenclature. PubMed:7678494 Q EPNC:7678494 Q
- Guengerich PP., J Biol Chem 1991;266:10019-10022.: Reactions and significance of cytochrome P-450 enzymes. PubMed:2037557 Q EPNC:2037557 Q
- Nebert DW, Gonzalez FJ., Annu Rev Biochem 1987;56:945-993.: P450 genes: structure, evolution, and regulation. PubMed:3304150 Q EPNC:3304150 Q
- Werck-Reichert D, Feyereisen R., Genome Biol 2000;1:REVIEW53003.: Cytochromes P450: a success story. PubMed:11178272 Q EPNC:11178272 Q

External database links

HOMSTRAD:	p450 Q
PANDIT:	PF00067 Q
PRINTS:	PR00385 Q PR00359 Q PR00408 Q PR00463 Q PR00464 Q PR00465 Q

Example structure
PDB entry 2D8R: Crystal Structure of Alene oxide synthase
[View a different structure](#)
3DM Q

wellcome trust sanger

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam

Family: **p450 (PF00067)**

282 architectures 39592 sequences 2 interactions 2977 species 873 structures

Summary

Domain organisation

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

There are **33674** sequences with the following architecture: **p450**
Q99N15_MOUSE [Mus musculus (Mouse)] Protein Cyp450 (534 residues)

p450

Show all sequences with this architecture.

There are **2329** sequences with the following architecture: **p450 x 2**
Q27YCA_MYCB [Mycobacterium paratuberculosis] Putative uncharacterized protein (422 residues)

p450 **p450**

Show all sequences with this architecture.

There are **307** sequences with the following architecture: **p450, Flavodoxin_1, FAD_binding_1, NAD_binding_1**
Q8KUJ0_ACTTA [Actinostroma pretiosum subsp. aurumcolum] Cytochrome P450 (1005 residues)

p450 **p450** **p450** **p450**

Show all sequences with this architecture.

There are **88** sequences with the following architecture: **p450 x 3**
Q331R6_SACTO [Saccharomyces sp. KCTC 00418P] Cytochrome P-450 (401 residues)

p450 **p450** **p450**

Show all sequences with this architecture.

There are **86** sequences with the following architecture: **An_peroxidase, p450**
Q2TW67_ASPOK [Aspergillus oryzae (strain ATCC 42146 / RIB 40) (Yellow koji mold)] Peroxidase/orygenase (1147 residues)

An_peroxidase **p450**

Show all sequences with this architecture.

There are **71** sequences with the following architecture: **p450, FAD_binding_6, NAD_binding_1, Fer2**
A1U298_BURP1 [Burkholderia pseudomallei (strain SA911)] Cytochrome P450 (784 residues)

p450 **FAD_binding_6** **NAD_binding_1** **Fer2**

Show all sequences with this architecture.

There are **36** sequences with the following architecture: **An_peroxidase x 2, p450**
Q0C299_ASPTN [Aspergillus terreus (strain NIH 2524 / FGSC A1156)] Putative uncharacterized protein (1045 residues)

An_peroxidase **An_peroxidase** **p450**

Show all sequences with this architecture.

There are **19** sequences with the following architecture: **p450, KR**
Q331H2_BURP1 [Burkholderia pseudomallei (strain 1710b)] Cytochrome P450 family protein (1373 residues)

p450 **KR**

Show all sequences with this architecture.

There are **13** sequences with the following architecture: **p450 x 4**
C1X958_BSAF1 [Brachyostoma forbesi (Florida lancelet) (Amphioxus)] Putative uncharacterized protein (562 residues)

p450 **p450** **p450** **p450**

wellcome trust sanger

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam

Family: **p450 (PF00067)**

282 architectures 39592 sequences 2 interactions 2977 species 873 structures

Summary

Domain organisation

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 using the family HMM. We also generate alignments using four representative proteomes¹ (RP) sets, 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 (39592)	Representative proteomes				NCBI (39556)	Meta (2723)
			RP15 (3463)	RP35 (11134)	RP55 (16825)	RP75 (20665)		
Jalview	✓	✓	✓	✓	✓	✓	✓	✓
HTML	✓	—	—	—	—	—	✓	✓
PP/heatmap	X	—	—	—	—	—	✓	✓
Pfam viewer	✓	✓	X	X	X	X	X	X

¹Cannot generate PP/heatmap alignments for seeds; no PP data available.

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

Format an alignment

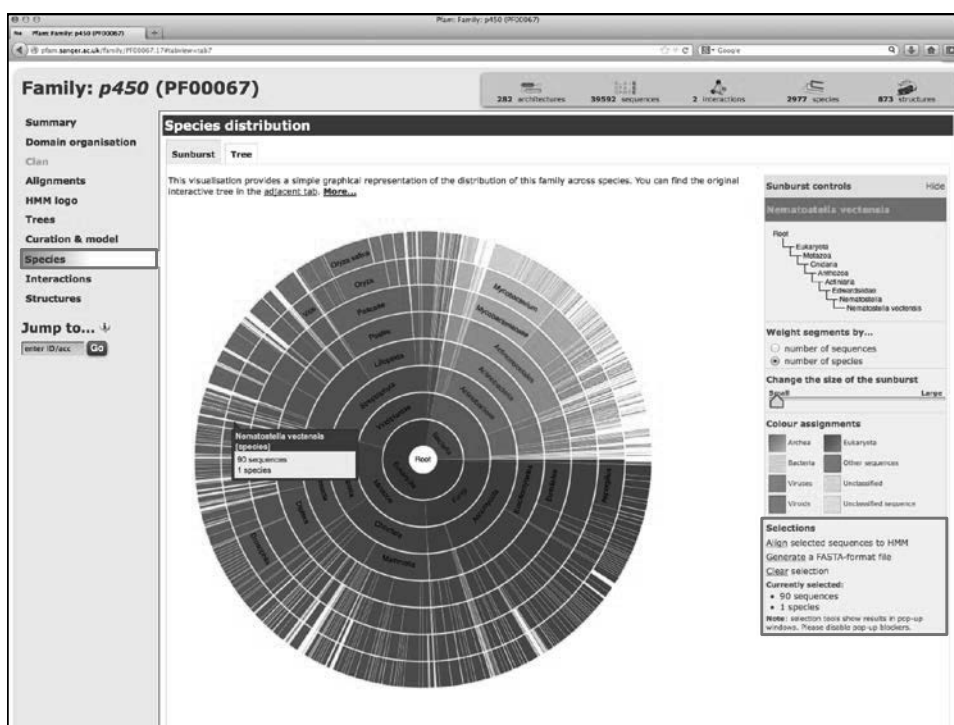
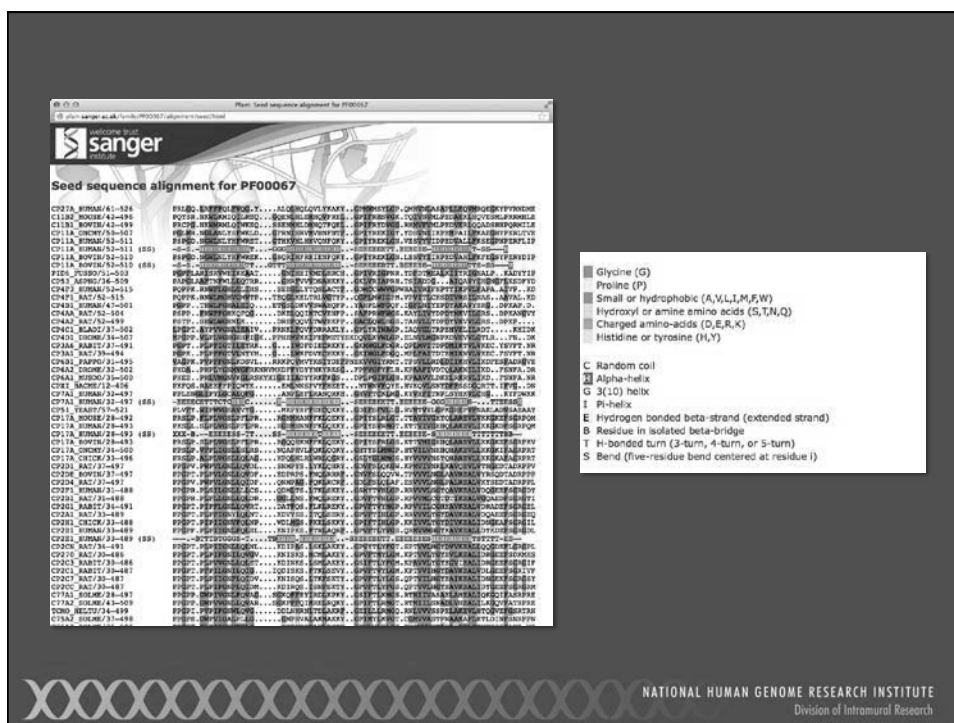
	Seed (50)	Full (39592)	Representative proteomes				NCBI (39556)	Meta (2723)
			RP15 (3463)	RP35 (11134)	RP55 (16825)	RP75 (20665)		
Alignment:	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Format:	Select							
Order:	<input type="radio"/> Tree <input type="radio"/> Alphabetical							
Sequence:	<input type="radio"/> Inserts lower case <input type="radio"/> All upper case							
Gaps:	Gaps as "-" or "." (mixed)							
Download/View:	<input 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 (39592)	Representative proteomes				NCBI (39556)	Meta (2723)
			RP15 (3463)	RP35 (11134)	RP55 (16825)	RP75 (20665)		
Alignment:	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>



wellcome trust sanger

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam Family: p450 (PF00067)

282 structures 26992 sequences 2 interactions 2977 domains 873 structures

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 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 mitochondrial 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;4:1065-1069. A three-dimensional model of aromatase cytochrome P450. PubMed:7549871 EPMC:7549871
- DeGyarenko KN, Archakov AI; FEBS Lett 1993;332:1-6. Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. PubMed:8405421 EPMC:8405421
- Neison 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 EPMC:7678494
- Guengerich FP; J Biol Chem 1991;266:10019-10022. Reactions and significance of cytochrome P-450 enzymes. PubMed:2037557 EPMC:2037557
- Nebert DW, Gonzalez FJ; Annu Rev Biochem 1987;56:945-993. P450 genes: structure, evolution, and regulation. PubMed:3304150 EPMC:3304150
- Wreck-Reichhart D, Feyereisen R; Genome Biol 2000;1:REVIEWS3003. Cytochromes P450: a success story. PubMed:11178272 EPMC:11178272

External database links

HOMSTRAD:	p450
PANDIT:	PF00067
PRINTS:	PR00385 PR00359 PR00408 PR00463 PR00464 PR00465

wellcome trust sanger

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam Family: p450 (PF00067)

282 structures 26992 sequences 2 interactions 2977 domains 873 structures

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 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 mitochondrial 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;4:1065-1069. A three-dimensional model of aromatase cytochrome P450. PubMed:7549871 EPMC:7549871
- DeGyarenko KN, Archakov AI; FEBS Lett 1993;332:1-6. Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. PubMed:8405421 EPMC:8405421
- Neison 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 EPMC:7678494
- Guengerich FP; J Biol Chem 1991;266:10019-10022. Reactions and significance of cytochrome P-450 enzymes. PubMed:2037557 EPMC:2037557
- Nebert DW, Gonzalez FJ; Annu Rev Biochem 1987;56:945-993. P450 genes: structure, evolution, and regulation. PubMed:3304150 EPMC:3304150
- Wreck-Reichhart D, Feyereisen R; Genome Biol 2000;1:REVIEWS3003. Cytochromes P450: a success story. PubMed:11178272 EPMC:11178272

External database links

HOMSTRAD:	p450
PANDIT:	PF00067
PRINTS:	PR00385 PR00359 PR00408 PR00463 PR00464 PR00465
PROSITE:	POC00081
Pseudofam:	PF00067
SCOP:	2cnp
SYSTEMS:	p450

Comments or questions on the site? Send a mail to pfam-help@sanger.ac.uk. Our cookie policy.

The Wellcome Trust

PROSITE documentation PDOC00081

Cytochrome P450 cysteine heme-iron ligand signature

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:

Daghyarenko 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]-[RKHPT]-[P]-C-[LVMFAP]-[GAD]
- C is the heme iron ligand
- Sequences in UniProtKB/Swiss-Prot known to belong to this class: 962
 - detected by PS00086: 922 (true positives)
 - undetected by PS00086: 70 (50 false negatives and 10 'borderline')
- 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 tree view 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 1BVY 1CBJ ... [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 (not Pfam B)
 - Simple Modular Architecture Research Tool (SMART)
 - COG (orthologous prokaryotic protein families)
 - KOG (eukaryotic equivalent of COG)
 - PRK ("protein clusters" of related protein RefSeq entries)
 - TIGRFAM

Marchler-Bauer *et al.*, *Nucleic Acids Res.* 41: D348-D352, 2013

Conserved Domain Database (CDD)

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



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

http://ncbi.nlm.nih.gov/Structure

NCBI Conserved Domain Search

Conserved Domains

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

Submit

Reset

OPTIONS

Search against database (D): CDD v3.11 - 45746 PSSMs

Expect Value (E) threshold: 0.010000

Apply low-complexity filter (F) ☐

Force live search (L) ☐

Maximum number of hits (H): 500

Result mode (M) ☒ Concise ☐ Standard ☐ Full

Retrieve previous CD-search result

Request ID: Retrieve

References:

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

NCBI Conserved Domain Search

Conserved Domains

Conserved domains on [cl|seqsig_6a6d6fc2d397e2ad4c3c4e6ae54c1b1]

NP_005206.1 deleted in colorectal carcinoma [Homo sapiens]

Graphical summary show options

Query seq.

Specific hits: **FN3** **FN3** **FN3** **FN3** **FN3** **FN3**

Superfamilies: **FN3 super** **FN3 super** **FN3 super** **FN3 super** **FN3 super** **FN3 super**

Multi-domains: **IGC2** **IGC2** **I-set** **I-set**

Neogenin_C

List of domain hits

Description	Pasid	Multi-dom	E-value
[H]g1_Neogenin[cd05722], First immunoglobulin (Ig)-like domain in neogenin and similar proteins; Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and similar proteins; Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and similar 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.	143199	no	6.80e-50
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	3.75e-16
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	1.17e-17
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	2.22e-16
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	2.04e-13
[H]g[cd00066], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	143166	yes	1.03e-10
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	3.63e-10
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	9.45e-09
[H]g super family[cd11963], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	264487	no	8.84e-35
[H]g super family[cd11963], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	264487	no	1.46e-14
[H]Neogenin_C[cd00066], Neogenin C-terminus: This family represents the C-terminus of eukaryotic neogenin precursor proteins, which ...	253838	yes	5.62e-143
[H]seq[cd00066], Immunoglobulin I-set domain;	254352	yes	3.34e-19
[H]seq[cd00066], Immunoglobulin I-set domain;	254352	yes	5.04e-19
[H]G2[cd00066], Immunoglobulin C-2 Type;	197706	yes	7.30e-16
[H]G2[cd00066], Immunoglobulin C-2 Type;	197706	yes	1.04e-08

References:

- 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

NCBI Conserved Domain Search

Conserved Domains

Conserved domains on [cl|seqsig_6a6d6fc2d397e2ad4c3c4e6ae54c1b1]

NP_005206.1 deleted in colorectal carcinoma [Homo sapiens]

Graphical summary show options

Query seq.

Specific hits: **FN3** **FN3** **FN3** **FN3** **FN3** **FN3**

Superfamilies: **FN3 super** **FN3 super** **FN3 super** **FN3 super** **FN3 super** **FN3 super**

Multi-domains: **IGC2** **IGC2** **I-set** **I-set**

Neogenin_C

List of domain hits

Description	Pasid	Multi-dom	E-value
[H]g1_Neogenin[cd05722], First immunoglobulin (Ig)-like domain in neogenin and similar proteins; Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and similar proteins; Ig1_Neogenin: first immunoglobulin (Ig)-like domain in neogenin and similar 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.	143199	no	6.80e-50
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	3.75e-16
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	1.17e-17
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	2.22e-16
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	2.04e-13
[H]g[cd00066], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	143166	yes	1.03e-10
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	3.63e-10
[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...	236020	no	9.45e-09
[H]g super family[cd11963], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	264487	no	8.84e-35
[H]g super family[cd11963], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...	264487	no	1.46e-14
[H]Neogenin_C[cd00066], Neogenin C-terminus: This family represents the C-terminus of eukaryotic neogenin precursor proteins, which ...	253838	yes	5.62e-143
[H]seq[cd00066], Immunoglobulin I-set domain;	254352	yes	3.34e-19
[H]seq[cd00066], Immunoglobulin I-set domain;	254352	yes	5.04e-19
[H]G2[cd00066], Immunoglobulin C-2 Type;	197706	yes	7.30e-16
[H]G2[cd00066], Immunoglobulin C-2 Type;	197706	yes	1.04e-08

Cd Length: 95 Bit Score: 173.43 E-value: 6.80e-50

seqsig_6a6d6fc2d397e2ad4c3c4e6ae54c1b1
cdid:cd05722

seqsig_6a6d6fc2d397e2ad4c3c4e6ae54c1b1
cdid:cd05722

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

[H]g[cd00066], Immunoglobulin domain; Ig: immunoglobulin (Ig) domain found in the Ig superfamily. The Ig superfamily is a ...

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

[H]FN3[cd00063], Fibronectin type 3 domain; One of three types of internal repeats found in the plasma ...

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: cl11960
- BioSystems: 310 links

Statistics

- PSSM-Id: 143199
- View PSSM: cd05722
- Aligned: 7 rows
- Threshold/BitsScore: 148.395
- Threshold/SettingGk: 148277558
- Created: 27-Sep-2007
- Updated: 17-Jan-2013

Structure

- Interactive View
- Aligned Rows: All 7 rows
- Download Cn3D

PubMed References

- Neogenin: one receptor, many functions. *Int J Biochem. Cell Biol.* 2007; 39(5):674-678
- 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.* 1964 Sep 30; 24(4):305-320
- The immunoglobulin superfamily: an insight on its tissue, species, and functional diversity. *J. Mol. Evol.* 1998 Apr; 46(4):385-407
- Evolution of antigen binding receptors. *Annu. Rev. Immunol.* 1996; 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 cl11960.

cd05722 Sequence Cluster

Sub-family Hierarchy

- cd05722 Ig1_Neogenin
- cd05723 Ig4_Neogenin
- cd05724 Ig2_Robo
- cd05725 Ig3_Robo
- cd05726 Ig4_Robo
- cd05727 Ig2_Conf.actin-2-like
- cd05728 Ig4_Conf.actin-2-like
- cd05729 Ig2_FGFRL1-like
- cd05856 Ig2_FGFRL1-like
- cd05857 Ig2_FGFRL1-like
- cd05730 Ig3_NCRN-1-like
- cd05731 Ig3_L1-NCRN-like
- cd05876 Ig3_L1-NCRN-like
- cd05732 Ig5_NCRN-1-like
- cd05863 Ig5_NCRN-1-like
- cd05870 Ig5_NCRN-2-like
- cd05733 Ig6_L1-NCRN-like
- cd05874 TcR_M-FRM

Created: 27-Sep-2007
Updated: 17-Jan-2013

Structure

- Interactive View
- Aligned Rows: All 7 rows
- 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

Accession	Position	Sequence
gi 62204258	35	WSTEPDZLA [5] VLLNCSVRS [3] AKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 111
gi 110645196	46	YFLEPVDVTE [5] AVLNCSATA [3] PKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 124
gi 113675978	28	FFLEPVDVTE [5] VLLDCQMG [3] IGINWLNQVITE [6] LSNGLISVYVSR [1] NKPDEGVTQCV 101
gi 148277558	30	WSTEPDZLA [5] VLLNCSVRS [3] AKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 109
gi 1169233	41	WSTEPDZLA [5] VLLNCSVRS [3] AKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 118
gi 10720134	20	WSTEPDZLA [5] VLLNCSVRS [3] AKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 96
gi 147903889	41	WSTEPDZLA [5] VLLNCSVRS [3] AKIENKXDSFLS [8] LADGSLISVYVSR [1] NKPDEGVTQCV 118
gi 62204258	112	ATI [3] QTIVSRATLV 129
gi 110645196	125	ATV [3] QTIVSRATLV 142
gi 113675978	102	AON [2] QTIVSRATLV 118
gi 148277558	106	AON [2] QTIVSRATLV 122
gi 1169233	119	ATV [3] QTIVSRATLV 136
gi 10720134	97	ATV [3] QTIVSRATLV 114
gi 147903889	119	ATV [3] QTIVSRATLV 136

Citing CDD

Marchler-Bauer A et al. (2013), "CDD: conserved domains and protein three-dimensional structure," *Nucleic Acids Res.* 41(D1):D384-52.

Disclaimer | Privacy statement | Accessibility

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*

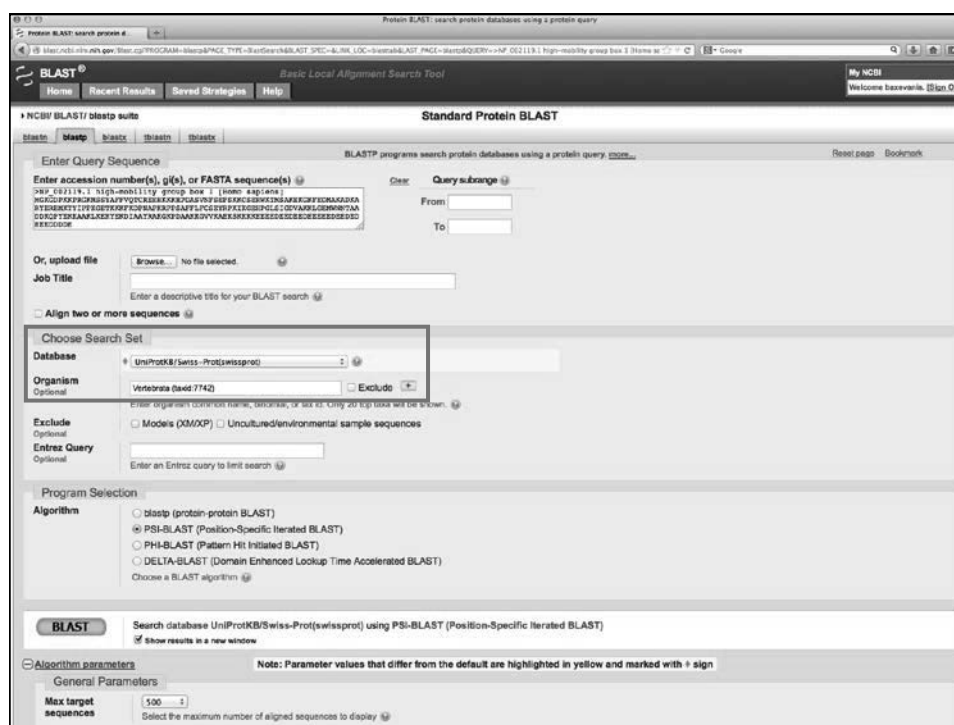
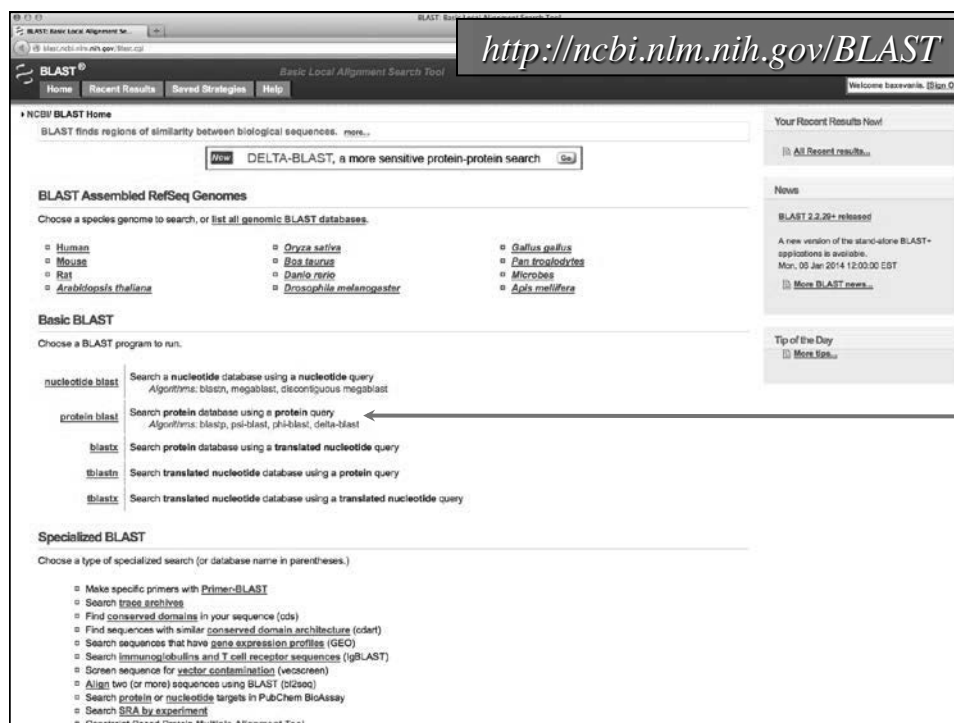


PSI-BLAST

- 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





Swiss-Prot

- *Goal*: Provide a single reference sequence for each protein sequence
- Distinguishing Features
 - Non-redundancy
 - Ongoing curation by EBI staff and *external experts*
 - Expert annotation includes editing/updates of
 - KW Keyword lines
 - CC Comment lines (the “executive summary”)
 - FT Feature table
 - Distinct accession series
 - [OPQ] 1 2 3 4 5



Protein BLAST: search protein databases using a protein query

Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI BLAST: blastp suite

Standard Protein BLAST

blastp blastx tblastn tblastx

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

Enter Query Sequence

Enter accession number(s), g(x), or FASTA sequence(s) [?](#)

or Query subrange [?](#)

From

To

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

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database [+](#) UnProtKB/Swiss-Prot[swissprot] [-](#) [?](#)

Organism [?](#) Vertebrate (taxid:7742) [-](#) Exclude [+](#)

Optional Enter organism common name, taxid, or tax id. Only 20 top taxa will be shown. [?](#)

Exclude [?](#) ☐ Models (MOLPROP) ☐ Uncultured/environmental sample sequences

Optional

Entrez Query [?](#)

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 [?](#)

BLAST Search database UniProtKB/Swiss-Prot[swissprot] using PSI-BLAST (Position-Specific Iterated BLAST)

☒ Show results in a new window

Algorithm parameters

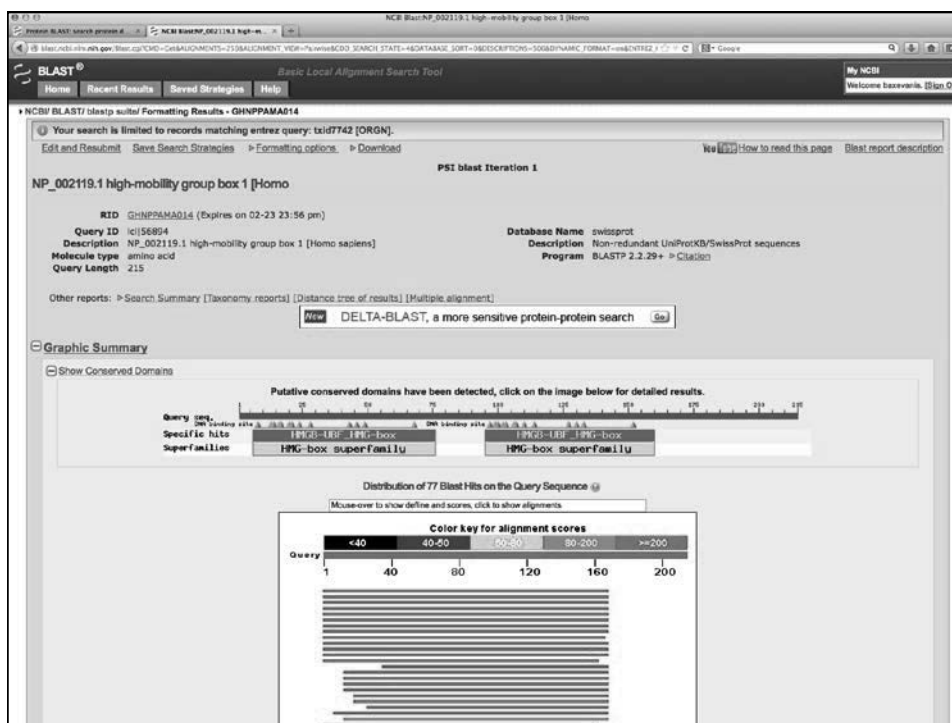
Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign

General Parameters

Max target sequences [100](#) [-](#) [+](#)

Select the maximum number of aligned sequences to display [?](#)

The screenshot shows the NCBI BLAST search interface. The search is performed against the UniProtKB/Swiss-Prot database using PSI-BLAST. The 'Expect threshold' is set to 0.001, with a callout indicating 'Default = 10'. The 'PSI-BLAST Threshold' is set to 0.001, with a callout indicating 'Default = 0.005'. The 'Matrix' is set to BLOSUM62. The 'Gap Costs' are set to Existence: 11, Extension: 1. The 'Filters and Masking' section is expanded, showing the 'Filter' set to 'Low complexity regions'. The 'Mask' section is also expanded, showing 'Mask for lookup table only' and 'Mask lower case letters' both checked. The 'BLAST' button is highlighted with a starburst.



Run PSI-Blast iteration 2 with max 500

Sequences producing significant alignments with E-value BETTER than threshold

Select: All None Selected 0

Alignments [Download](#) [GenPlot](#) [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"/>	RefName: Full-high mobility group protein B1, AltName: Full-high mobility group protein 1, ShortName: MG-1	310	310	78%	2e-106	100%	P10103.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-high mobility group protein B1, AltName: Full-Anahelelin, AltName: Full-Hesarin-binding protein p30, AltName: Full-His	310	310	78%	2e-106	100%	P36159.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-high mobility group protein B1, AltName: Full-Hisgag mobility group protein 1, ShortName: MG-1, GAG	310	310	78%	2e-106	100%	P36429.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-high mobility group protein B1, AltName: Full-Hisgag mobility group protein 1, ShortName: MG-1	308	308	78%	1e-105	99%	P12892.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-high mobility group protein B1, AltName: Full-Hisgag mobility group protein 1, ShortName: MG-1	299	299	78%	4e-102	96%	G9Y168.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Relative high mobility group protein B1-like 1, AltName: Full-Hisgag mobility group protein B1 pseudogene 1, AltName: Full	297	297	78%	2e-101	95%	G25760.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Relative high mobility group protein 1 like 10, ShortName: MG-1, 10	290	290	78%	1e-98	95%	G56526.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2	257	257	78%	1e-85	85%	P26584.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2	257	257	77%	1e-85	83%	P37745.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2, B2/10	252	252	78%	8e-84	86%	P59583.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2	251	251	78%	2e-83	86%	P26825.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2	249	249	78%	2e-82	86%	P30881.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B2, AltName: Full-Hisgag mobility group protein 2, ShortName: MG-2	245	245	75%	5e-81	87%	P17741.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B1, AltName: Full-Hisgag mobility group protein 1, ShortName: MG-1	239	239	82%	3e-79	100%	P37156.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B3, AltName: Full-Hisgag mobility group protein 3a, ShortName: MG-2a, AltName: Full-Hisgag mobility	210	210	73%	2e-67	75%	O56495.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B3, AltName: Full-Hisgag mobility group protein 3a, ShortName: MG-2a, AltName: Full-Hisgag mobility	209	209	73%	6e-67	75%	P40818.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B3	209	209	73%	6e-67	75%	O25131.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Hisgag mobility group protein B3, AltName: Full-Hisgag mobility group protein 3a, ShortName: MG-2a, AltName: Full-Hisgag mobility	208	208	73%	1e-66	75%	O15347.4	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Nuclear autoantigen Sp100, AltName: Full-Nuclear dot-associated Sp100 protein, AltName: Full-Spectro110 kDa	207	207	70%	5e-66	85%	Q9H126.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Nuclear autoantigen Sp100, AltName: Full-Nuclear dot-associated Sp100 protein, AltName: Full-Spectro110 kDa	201	201	70%	2e-63	85%	Q9H127.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Nuclear autoantigen Sp100, AltName: Full-Nuclear dot-associated Sp100 protein, AltName: Full-Spectro110 kDa	201	201	66%	2e-63	87%	Q9H126.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Nuclear autoantigen Sp100, AltName: Full-Nuclear dot-associated Sp100 protein, AltName: Full-Spectro110 kDa	211	211	73%	1e-62	82%	P27387.1	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	RefName: Full-Relative high mobility group protein B0-like protein	197	197	73%	2e-62	69%	P30636.1	<input type="checkbox"/>	<input type="checkbox"/>

[illegible]

NCBI BLAST search results - NCBI BLAST: NP_001151.1 high-mobility group box 1 (HMG1) -

Run PSI-Blast iteration 3 with max: 500

Sequences producing significant alignments with E-value BETTER than threshold

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

Alignments

Description	Max score	Total score	Query cover	E value	Ident	Accession	Select for PSI-Blast	Used to build PSI-Blast
ReName: Full=high mobility group protein B1; AName: Full=high mobility group protein 1; Short=HMG-1	242	242	78%	9e-80	99%	P12682.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B1; AName: Full=high mobility group protein 1; Short=HMG-1	242	242	78%	9e-80	100%	P10103.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B1; AName: Full=Ansholein; AName: Full=leapin-binding protein p30; AName: Full=high	242	242	78%	9e-80	100%	P31568.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B1; AName: Full=high mobility group protein 1; Short=HMG-1; sp=Q8Y5AM.3; MGB1_CAMP1	242	242	78%	9e-80	100%	P35429.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B1; AName: Full=high mobility group protein 1; Short=HMG-1	233	236	78%	3e-78	96%	Q2Y109.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=putative high mobility group protein B1-like 1; AName: Full=high mobility group protein B1 pseudogene 1; AName: Full=	235	236	78%	2e-77	95%	B2P293.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=putative high mobility group protein B1-like 10; Short=HMG-1.10	230	230	78%	7e-75	95%	Q5U0V6.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group T protein; Short=HMG-T; AName: Full=HMG-T; Short=HMG-1	211	211	77%	6e-68	83%	P37746.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B2; AName: Full=high mobility group protein 2; Short=HMG-2	211	211	78%	1e-67	85%	P25564.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B2; AName: Full=high mobility group protein 2; Short=HMG-2	206	206	78%	2e-65	86%	P25252.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B2; AName: Full=high mobility group protein 2; Short=HMG-2; sp=P40673.3; MGB2_BOVIN1	204	204	78%	4e-65	86%	P25563.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B2; AName: Full=high mobility group protein 2; Short=HMG-2	203	203	78%	2e-64	86%	P30681.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Nuclear autoantigen Sp-100; AName: Full=Nuclear dot-associated Sp100 protein; AName: Full=Sp100; 100 kDa	202	257	76%	8e-64	81%	Q9Y1Q6.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Nuclear autoantigen Sp-100; AName: Full=Nuclear dot-associated Sp100 protein; AName: Full=Sp100; 100 kDa	214	273	76%	1e-63	82%	P25472.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B3; AName: Full=high mobility group protein 2a; Short=HMG-2a; AName: Full=high mobility	200	200	78%	2e-63	76%	P5018.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B2; AName: Full=high mobility group protein 2; Short=HMG-2	200	200	75%	2e-63	87%	P17741.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B3	198	198	78%	8e-63	76%	Q32131.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B3; AName: Full=high mobility group protein 2a; Short=HMG-2a; AName: Full=high mobility	198	198	78%	1e-62	76%	O15347.4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B3; AName: Full=high mobility group protein 2a; Short=HMG-2a; AName: Full=high mobility	198	198	78%	1e-62	76%	O54753.3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Nuclear autoantigen Sp-100; AName: Full=Nuclear dot-associated Sp100 protein; AName: Full=Sp100; 100 kDa	196	249	76%	6e-62	80%	Q9Y1Q7.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=high mobility group protein B4	195	195	76%	2e-61	44%	Q32134.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=putative high mobility group protein B3-like protein	193	193	78%	8e-61	70%	P30684.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Nuclear autoantigen Sp-100; AName: Full=Nuclear dot-associated Sp100 protein; AName: Full=Sp100; 100 kDa	190	244	76%	4e-59	80%	Q9Y1Q8.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

NCBI BLAST search results - NCBI BLAST: NP_001151.1 high-mobility group box 1 (HMG1) -

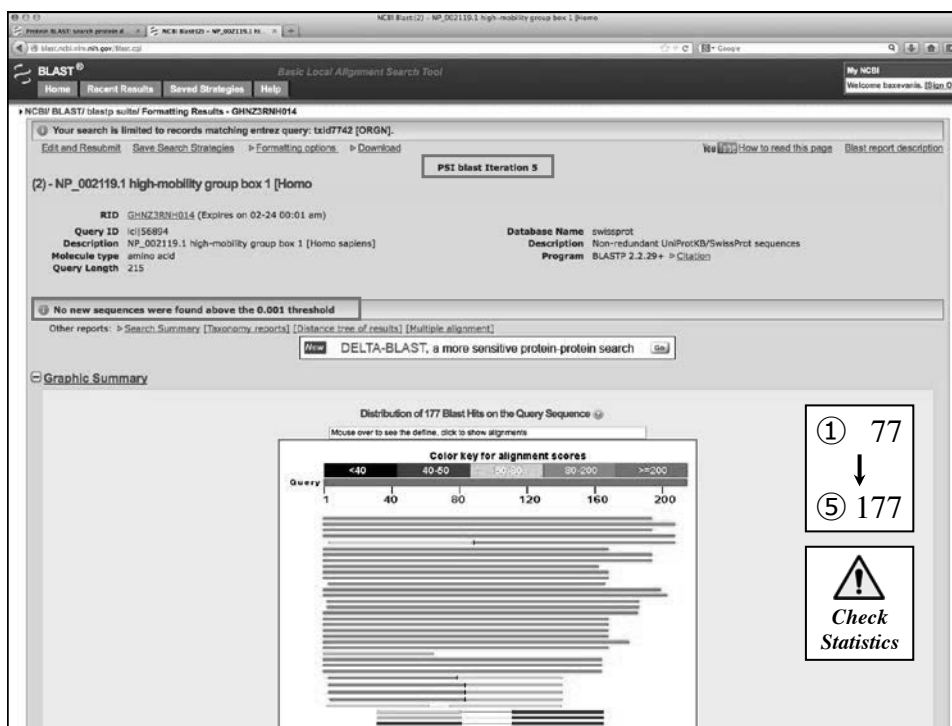
Run PSI-Blast iteration 3 with max: 500

Sequences producing significant alignments with E-value BETTER than threshold

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

Alignments

Description	Max score	Total score	Query cover	E value	Ident	Accession	Select for PSI-Blast	Used to build PSI-Blast
ReName: Full=TOX high-mobility group box family member 4; AName: Full=Endothelial Leukemia cell protein LCP1	82.5	152	68%	2e-17	38%	Q91842.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=TOX high-mobility group box family member 4; AName: Full=Endothelial Leukemia cell protein LCP1	82.5	152	68%	2e-17	38%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=TOX high-mobility group box family member 4	82.5	152	68%	2e-17	38%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=TOX high-mobility group box family member 4-B	81.3	150	68%	4e-17	38%	Q81980.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=TOX high-mobility group box family member 4	81.3	150	68%	4e-17	38%	A6Q2P5.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=TOX high-mobility group box family member 4-A	81.3	150	68%	4e-17	38%	Q8D48.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=FACT complex subunit SSRP1; AName: Full=Facilitates chromatin transcription complex subunit SSRP1; AName: Full=R	109	203	64%	1e-25	37%	Q05787.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=FACT complex subunit SSRP1; AName: Full=Facilitates chromatin transcription complex subunit SSRP1; AName: Full=R	106	203	64%	1e-25	36%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=FACT complex subunit SSRP1; AName: Full=Chromatin-specific transcription elongation factor B2 kDa subunit; AName: Full=	114	206	63%	4e-28	36%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	58.2	115	61%	5e-10	35%	Q05781.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.8	113	61%	9e-10	35%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor; sp=Q7/GF9.1; SRY_BOSGF; ReName: Full=Se	57.8	114	62%	9e-10	35%	Q7/GF7.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.8	114	62%	9e-10	35%	Q81922.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.8	114	62%	9e-10	35%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.5	114	61%	1e-09	35%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	56.3	111	62%	4e-09	35%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=SRV-related protein AMA2	39.3	39.3	22%	3e-04	35%	P43642.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=SRV-related protein AMA3	39.3	39.3	22%	4e-04	35%	P43643.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=SRV-related protein QH1	39.0	39.0	22%	5e-04	35%	P43645.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=FACT complex subunit SSRP1; AName: Full=Facilitates chromatin transcription complex subunit SSRP1; AName: Full=R	112	204	64%	1e-27	35%	Q05787.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=FACT complex subunit SSRP1; AName: Full=Facilitates chromatin transcription complex subunit SSRP1; AName: Full=R	112	211	64%	9e-28	34%	Q25843.2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=SRV-related protein QH3	40.1	40.1	23%	2e-04	34%	P43647.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=SRV-related protein QH1	38.6	38.6	23%	6e-04	34%	P43670.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Protein polyoma-1	77.9	141	60%	1e-15	34%	Q05941.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=MS1 protein homolog 1; AName: Full=DNA mismatch repair protein PM1	60.5	108	59%	4e-10	34%	P54272.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.8	114	62%	9e-10	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor; sp=Q68422.1; SRY_PHOPI; ReName: Full=Se	57.1	111	65%	1e-09	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.1	111	65%	1e-09	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.1	111	65%	1e-09	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor	57.1	111	65%	1e-09	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ReName: Full=Sex-determining region Y protein; AName: Full=Testis-determining factor; sp=Q68422.1; SRY_LAGG; ReName: Full=Se	57.1	111	65%	1e-09	33%	Q25843.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>



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

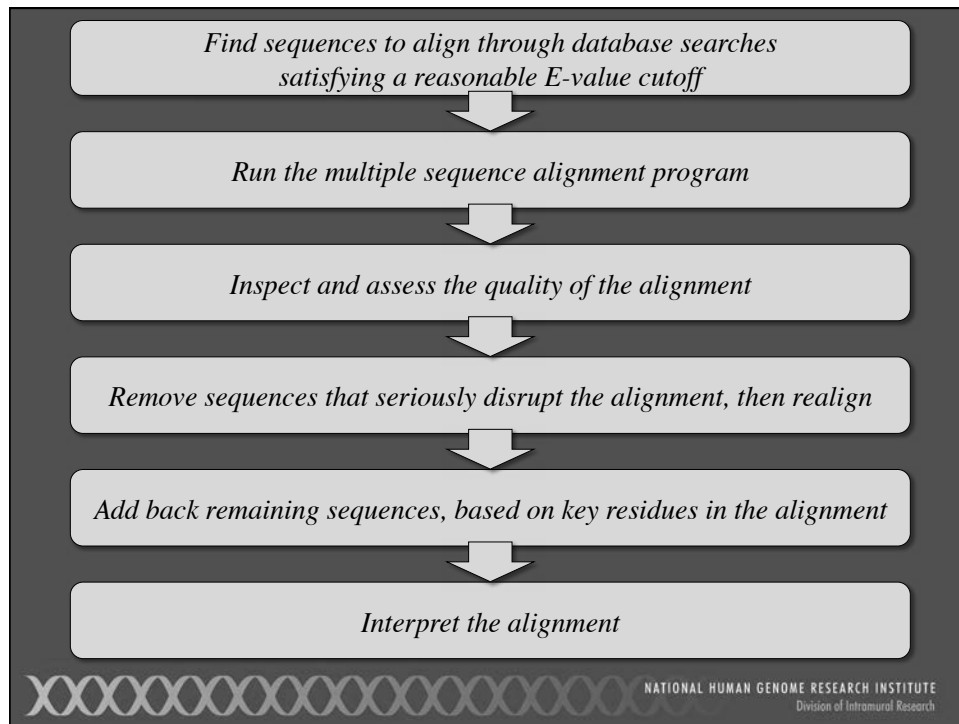
- Absolute sequence similarity
Create the alignment by lining up as many common characters as possible
- Conservation
Take into account residues that can substitute for one another and not adversely affect the function of the protein
- Structural similarity
Knowledge of the secondary or tertiary structure 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 are almost impossible
 - 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 (*e.g.*, BLAST)
 - Profile-based search methods (*e.g.*, 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” (*e.g.*, JalView)
- 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
- Can align data sets quickly and easily
- Can bias the location of gaps, based on known structural information
- Works with Jalview, 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



Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLST
>sequence B
VQLSGEEKAAVLALWDKVNNEEVGGEALGRLLVVYPWTQRFFDSFGDSLN
>sequence C
VLSPADKTNVKAANGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSH
>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



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLST
>sequence B
VQLSGEEKAAVLALWDKVNNEEVGGEALGRLLVVYPWTQRFFDSFGDSLN
>sequence C
VLSPADKTNVKAANGKVGAGHAGEYGAEALERMFLSFPTTKTYFPHFDLSH
>sequence D
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSH
```

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

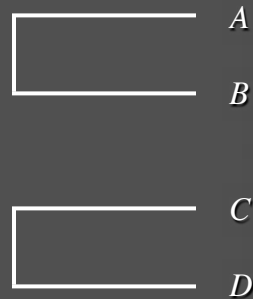


NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Progressive Alignment

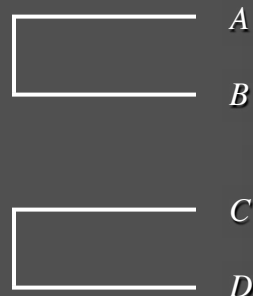
2. Derive a guide tree based on the pairwise comparisons

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



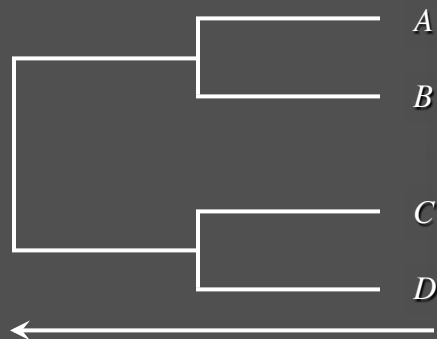
Progressive Alignment

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



Progressive Alignment

- Align “sequence” AB with “sequence” CD
- 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
- Clustal Omega does allow for guide tree iterations to hedge against errors introduced early in the alignment process (at the cost of increased compute time)



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 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 “significance”

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



Clustal Omega - Multiple Sequence Alignment

<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.

STEP 1 - Enter your input sequences

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

```
>F0S8_MOUSE Protein fosB
MFQAFPGDYDGSRCSSPSAESCYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA
ITTSQDLWLWQPTLISSMAQSGQPLASQPPAVDYPDMPTSYSTPLGLSAYSTGASGS
GGPSTSTTTSGPVSARPARARPRRPREETLTPEEEKRRVRERKLAAMKCHNRHREL
DRLQAEITDQLEEKAELESEAELOKEKELEFLVLAHKPGCKPYEEGPGPLAEVRD
LPOSTSAKEDGFGWLPPPPPPPPPPGSSRDAPPLTASLFTHSEVGLGDPPFVVSPSY
TSSFVLTCPEVSFAFAGQRTSQSEQPSDPLNSPLLLAL
```

Or, upload a file: No file selected.

STEP 2 - Set your parameters

OUTPUT FORMAT

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

(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)

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

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

Clustal Omega - Multiple Sequence Alignment

Services Research Training About us

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.

STEP 1 - Enter your input sequences

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

```
>F0S8_MOUSE Protein fosB
MFQAFPGDYDGSRCSSPSAESCYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA
ITTSQDLWLWQPTLISSMAQSGQPLASQPPAVDYPDMPTSYSTPLGLSAYSTGASGS
GGPSTSTTTSGPVSARPARARPRRPREETLTPEEEKRRVRERKLAAMKCHNRHREL
DRLQAEITDQLEEKAELESEAELOKEKELEFLVLAHKPGCKPYEEGPGPLAEVRD
LPOSTSAKEDGFGWLPPPPPPPPPPGSSRDAPPLTASLFTHSEVGLGDPPFVVSPSY
TSSFVLTCPEVSFAFAGQRTSQSEQPSDPLNSPLLLAL
```

Or, upload a file: No file selected.

STEP 2 - Set your parameters

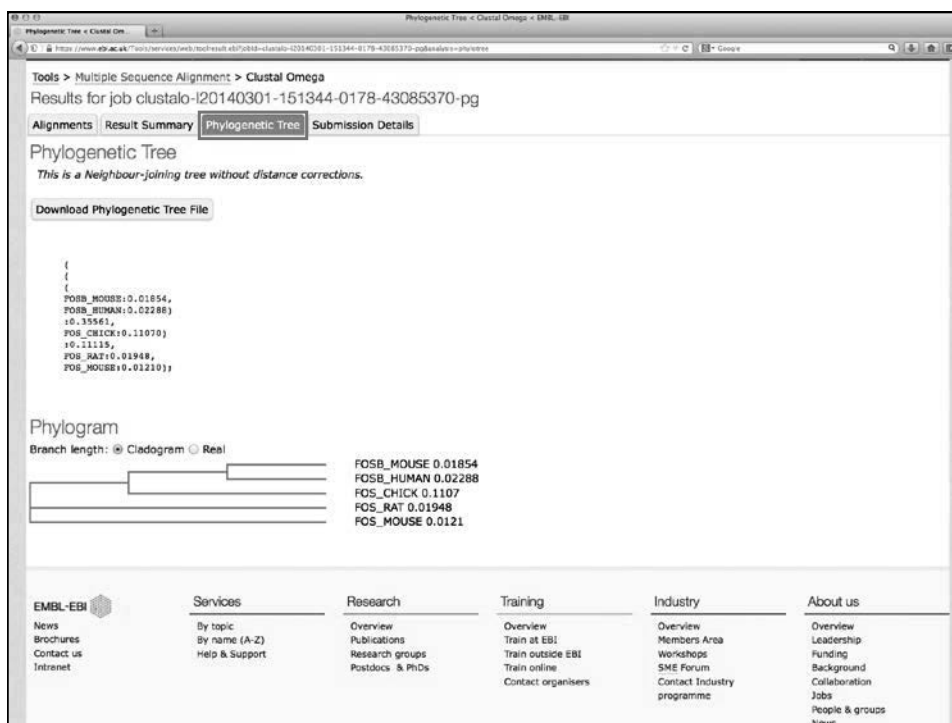
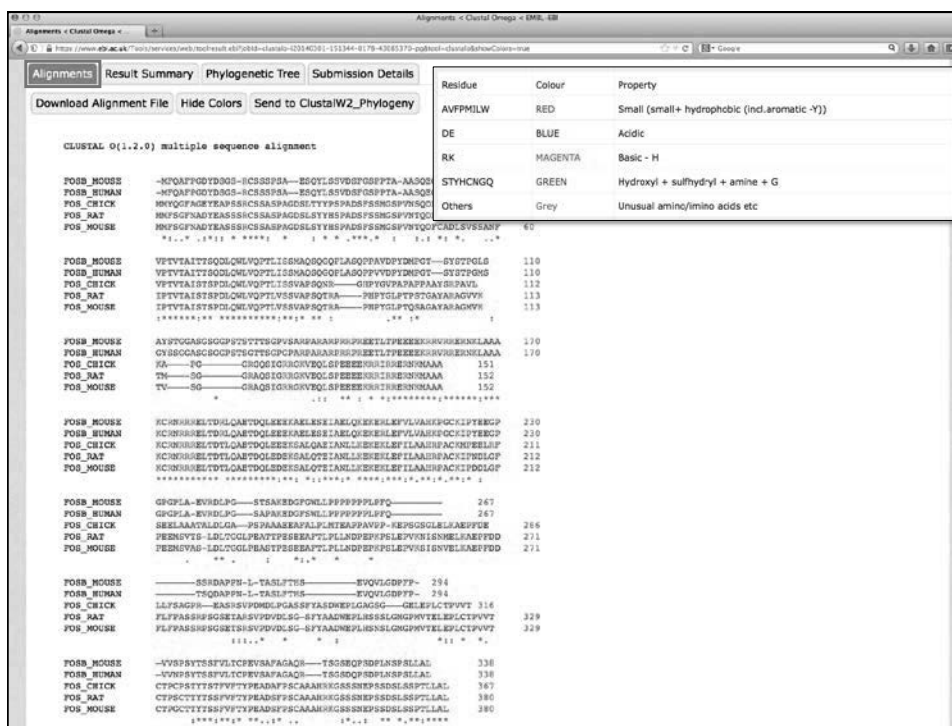
OUTPUT FORMAT

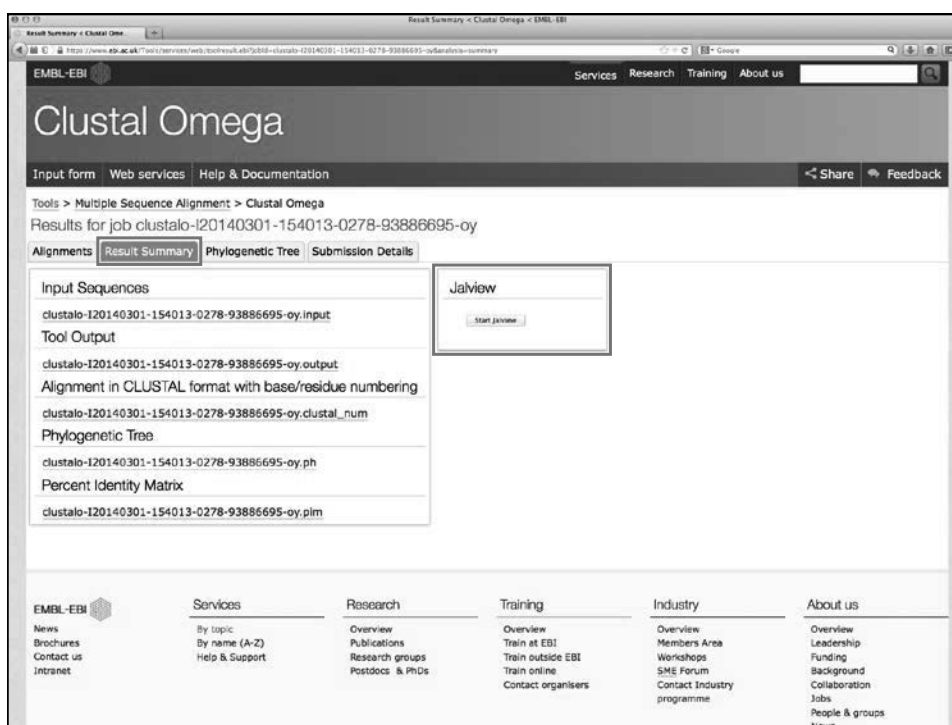
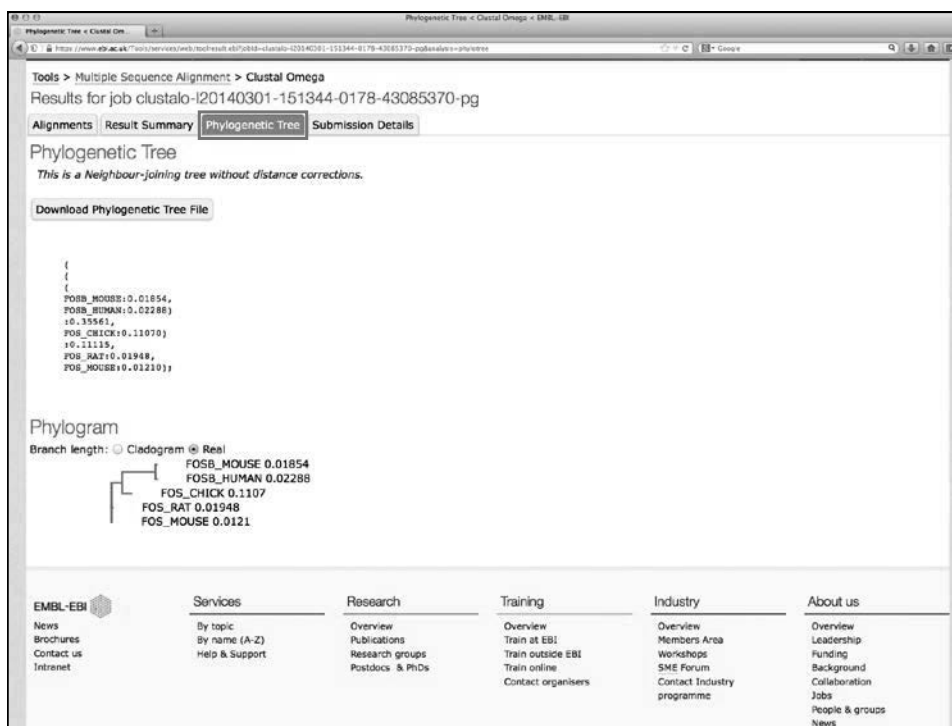
DEALIGN INPUT SEQUENCES	MBED-LIKE CLUSTERING GUIDE-TREE	MBED-LIKE CLUSTERING ITERATION	NUMBER OF COMBINED ITERATIONS
<input type="button" value="no"/>	<input type="button" value="yes"/>	<input type="button" value="yes"/>	<input type="button" value="default(0)"/>
MAX GUIDE TREE ITERATIONS	MAX HMM ITERATIONS	ORDER	
<input type="button" value="default"/>	<input type="button" value="default"/>	<input type="button" value="aligned"/>	

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)

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



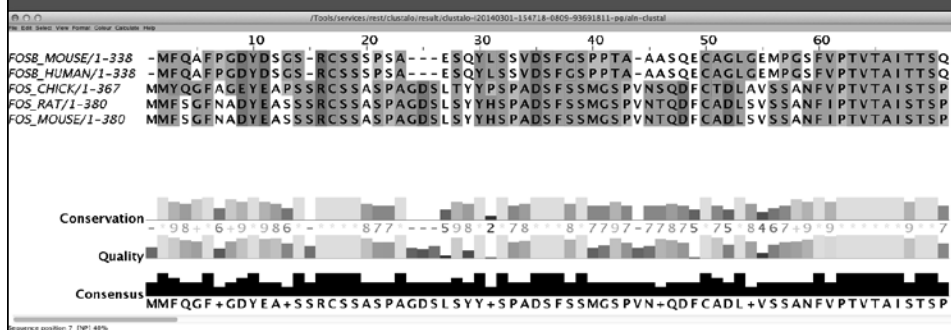


Jalview

- 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

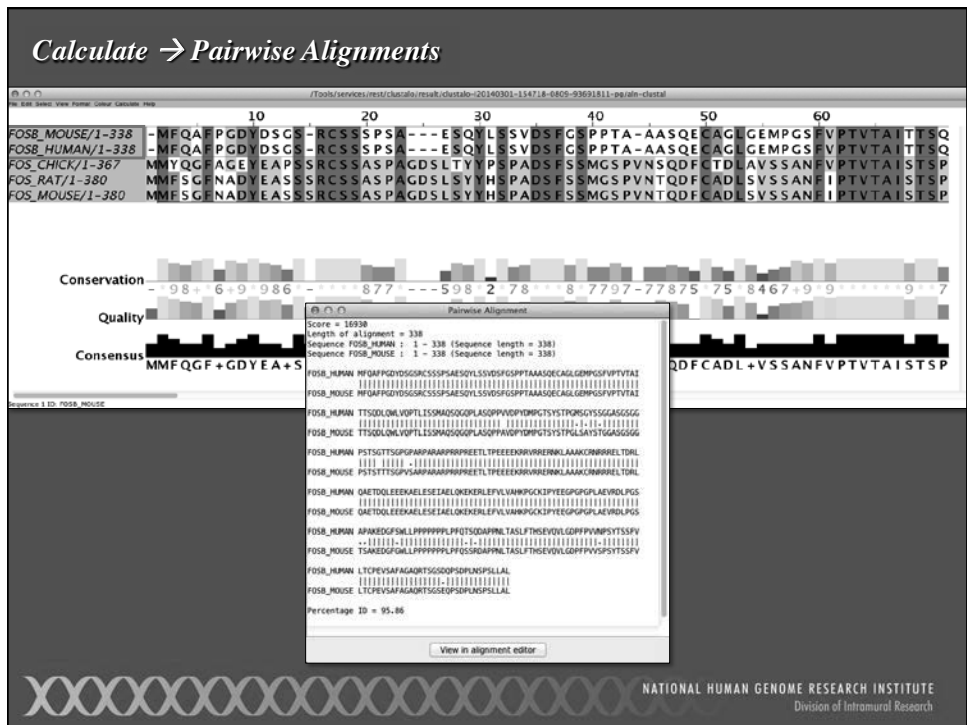


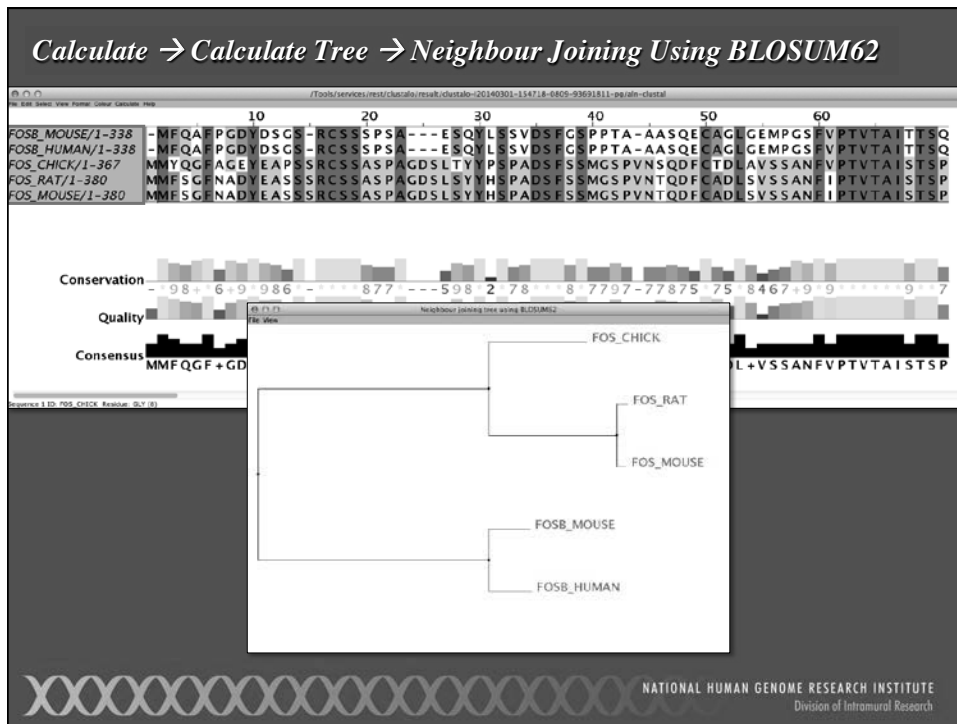
Default view



<i>Conservation</i>	Conservation of total alignment (indication of percent identity)
<i>Quality</i>	Alignment quality, based on BLOSUM scores
<i>Consensus</i>	Based on percent identity





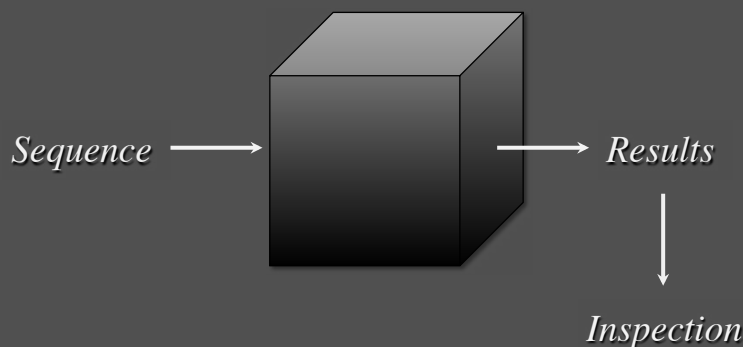


T-COFFEE

- 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)



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

[illegible]

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

Current Topics in Genome Analysis 2014

Next Lecture
March 19, 2014

Genome-Scale Sequence Analysis

Tyra Wolfsberg, Ph.D.
National Human Genome Research Institute
National Institutes of Health

A banner for the NIH Intramural Research Program. It features a series of five small, tilted photographs showing various researchers in lab coats working in a laboratory setting. Below the photos is the NIH logo (a stylized 'NIH' in a dark box) followed by the text "Intramural Research Program" and "Our Research Changes Lives". To the right of this, the text "one program many people infinite possibilities" is displayed, followed by the website "irp.nih.gov".

NIH Intramural Research Program
Our Research Changes Lives

one program
many people
infinite possibilities

irp.nih.gov