

joanne@msl.ubc.ca

wireless login:

mslguest

4myguest

# Laboratory Bioinformatics

Common tools, useful databases, and tricks of the trade for practical use in the laboratory.



[bioteach.ubc.ca/bioinfo2010](http://bioteach.ubc.ca/bioinfo2010)

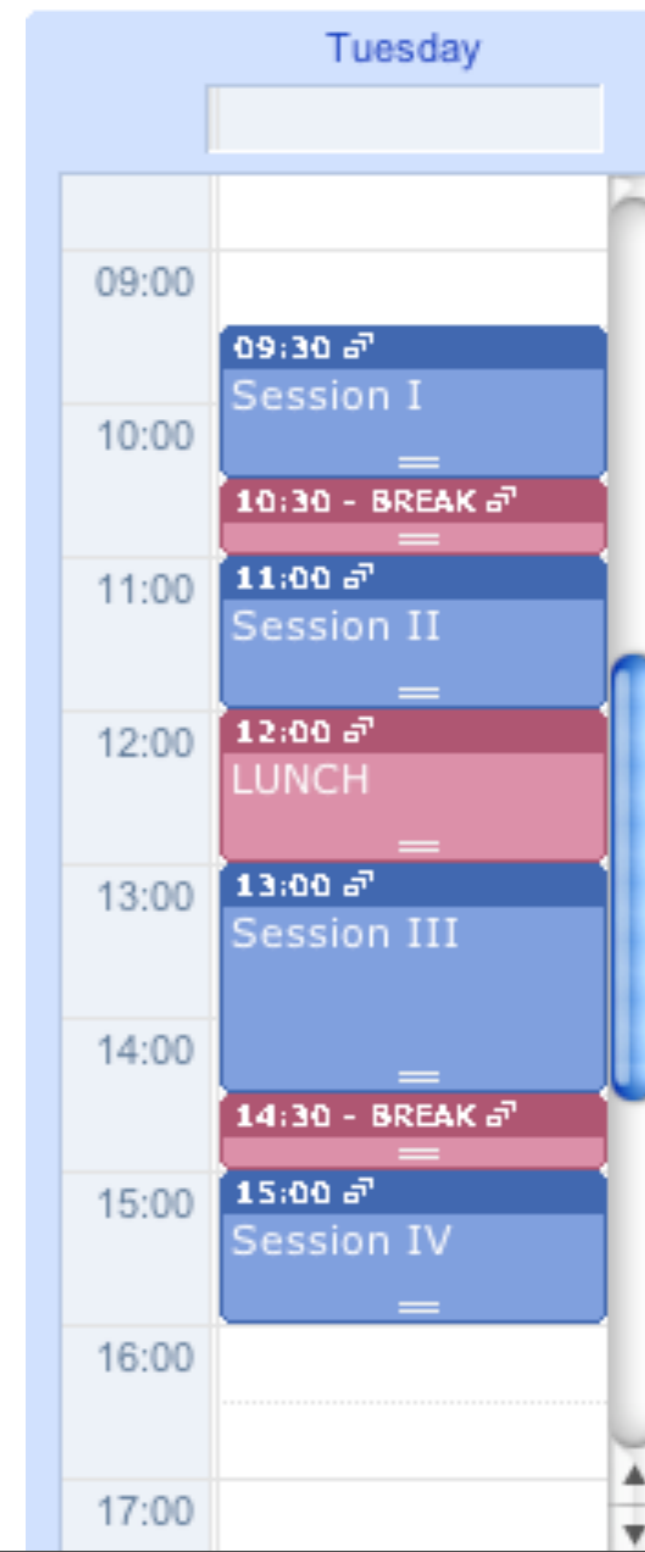
# Workshop Schedule

- Laptops, available here for your use 9am - 4:30pm

- wireless login

mslguest

4myguest





# Today's Topics

- **BLAST** - Finding Function by Sequence Similarity
- **GUIDED TOUR** - Advanced Tips & Tricks for Using BLAST
- **PRACTICAL EXERCISES** - The Jurassic Park Detective Story
- **COMMON TASKS** - Basic Search; Searching Sets of Sequences (multiple inputs; small custom databases); Primer Design

# BLAST

Finding Function By Sequence Similarity



# What do the Score and the e-value really mean?

- The quality of the alignment is represented by the **Score (S)**.

The score of an alignment is calculated as the sum of substitution and gap scores. Substitution scores are given by a look-up table (PAM, BLOSUM) whereas gap scores are assigned empirically .

- The significance of each alignment is computed as an **E value (E)**.

Expectation value. The number of different alignments with scores equivalent to or better than S that are expected to occur in a database search by chance. The lower the E value, the more significant the score.

# BLAST Algorithm

- Scoring of matches done using scoring matrices
- Sequences are split into words (default  $n=3$ )
  - Speed, computational efficiency
- BLAST algorithm extends the initial “seed” hit into an HSP
  - HSP = high scoring segment pair = Local optimal alignment

# How Does BLAST Really Work?

- The BLAST programs improved the overall speed of searches while retaining good sensitivity (important as databases continue to grow) by breaking the query and database sequences into fragments ("words"), and initially seeking matches between fragments.
- Word hits are then extended in either direction in an attempt to generate an alignment with a score exceeding the threshold of "S".



# BLAST Algorithm

Query Word ( $W = 3$ )

TLSHAWRLSNETDKRPFIEAERL**RDQ**HKKDYPEYKYQPRRRKNGKPGSSSEADAHSE



Determine neighborhood

<b>RDQ</b> 16	QDQ 12	EDQ 11	RDN 11	RDB 11	BDQ 10	RDP 10
RBQ 14	<b>REQ</b> 12	HDQ 11	RDD 11	ADQ 10	XDQ 10	RDT 10
RDZ 14	RDR 12	ZDQ 11	RDH 11	MDQ 10	RQQ 10	RDY 10
KDQ 13	RDK 12	RNQ 11	RDM 11	SDQ 10	RSQ 10	RDX 10
RDE 13	NDQ 11	RZQ 11	RDS 11	TDQ 10	RDA 10	DDQ 9 ...

# How Does BLAST Really Work?

- The BLAST programs improved the overall speed of searches while retaining good sensitivity (important as databases continue to grow) by breaking the query and database sequences into fragments ("words"), and initially seeking matches between fragments.
- Word hits are then extended in either direction in an attempt to generate an alignment with a score exceeding the threshold of "S".

# BLAST Algorithm

<b>RDQ</b> 16	QDQ 12	EDQ 11	RDN 11	RDB 11	BDQ 10	RDP 10
RBQ 14	<b>REQ</b> 12	HDQ 11	RDD 11	ADQ 10	XDQ 10	RDT 10
RDZ 14	RDR 12	ZDQ 11	RDH 11	MDQ 10	RQQ 10	RDY 10
KDQ 13	RDK 12	RNQ 11	RDM 11	SDQ 10	RSQ 10	RDX 10
RDE 13	NDQ 11	RZQ 11	RDS 11	TDQ 10	RDA 10	DDQ 9 ...

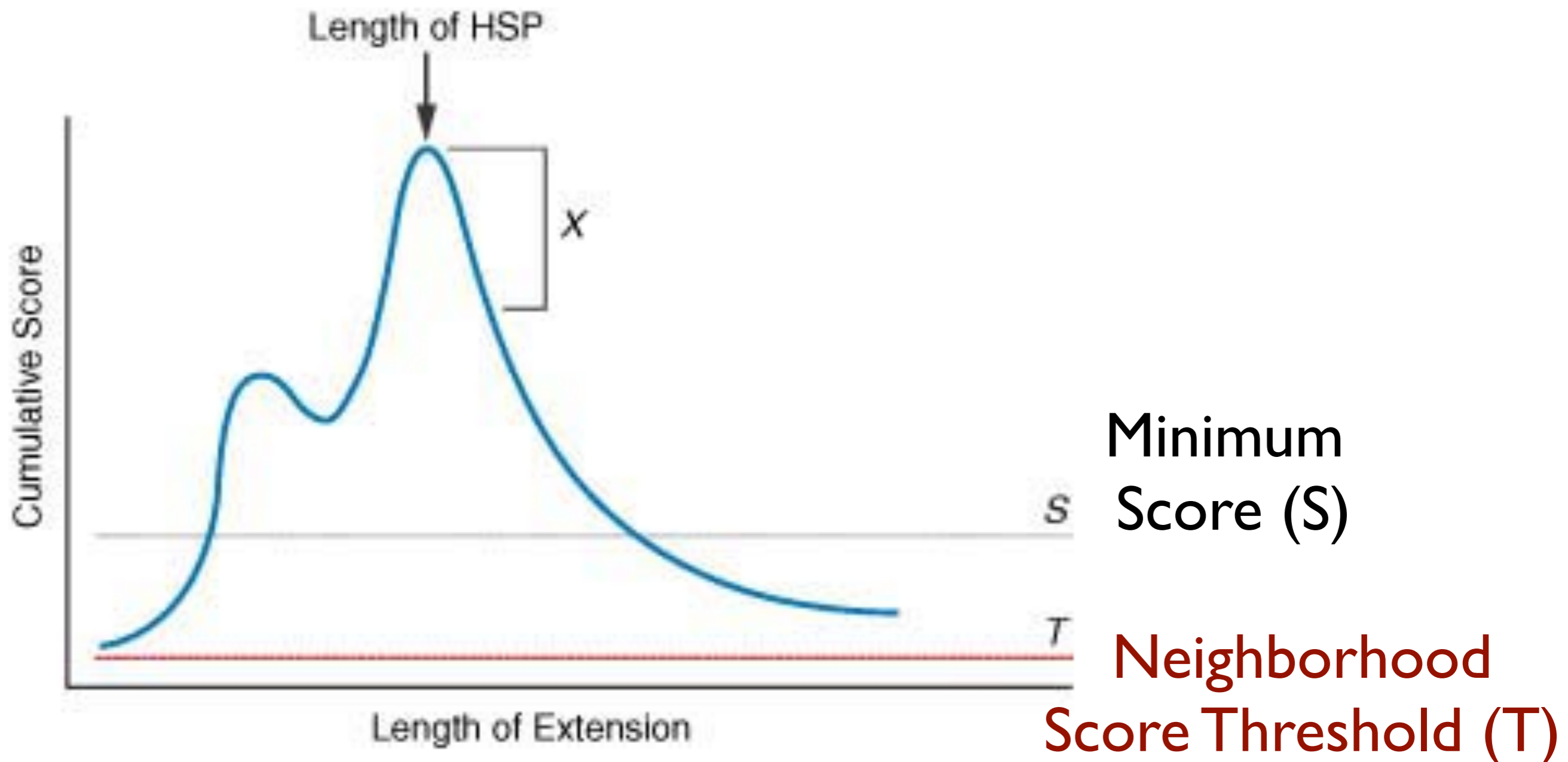
*Extension using neighborhood words greater than neighborhood score threshold ( $T = 11$ )*



```

Query: 1  TL SHAWRLSNETDKRPFIE TAERL RDQ HKKDYPEYKYQPRRRKNGKPGSSSEADAHSE 58
          TL   WRL N  +KRPF+E  AERLR+QHKKD+P+YKYQPRRRK+  K G S   D   +
Sbjct: 140 TLESGWRLNPGEKRPFVEGAERL REQ HKKDHPDYKYQPRRRKSVKNGQSEPEDGSEQ 197
  
```

# Extending the High Scoring Segment Pair (HSP)



> [gb|AAL08419.1](#) PTEN [Takifugu rubripes]  
Length=412

Score = 197 bits (501), Expect = 2e-49, Method: Composition-based stats.  
Identities = 95/100 (95%), Positives = 98/100 (98%), Gaps = 0/100 (0%)

```
Query 2 IVSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI 61
      +VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI
Sbjct 8 MVS RNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI 67

Query 62 YNLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFKQN 101
      YNLCAERHYD AKFNCRVAQYPPFEDHNPPQLELIKPF ++
Sbjct 68 YNLCAERHYDAAKFNCRVAQYPPFEDHNPPQLELIKPFCE D 107
```

Score = 83.6 bits (205), Expect = 4e-15, Method: Composition-based stats.  
Identities = 60/103 (58%), Positives = 68/103 (66%), Gaps = 32/103 (31%)

```
Query 99 KQNKMLKKDKMPHFVNTFFIPGPEEV-----D 126
      KQNKMK+KKDKMPHFVNTFFIPGPEE +
Sbjct 260 KQNKMMKKDKMPHFVNTFFIPGPEESRDKLENGAVNNADSQQGV P APGQGQPQSAECRE 319

Query 127 NDKEYLVLTLTkndldkankdkanRYFSPNFKVKLYFTKTVEE 169
      +D++YL+LTL+KND DKANKDKANRYFSPNFKVKL F+KTVEE
Sbjct 320 SDRDY LILTL SKNDRDKANKDKANRYFSPNFKVKLCFSKTVEE 362
```

> [gb|AAH93110.1](#) **UG** Ptenb protein [Danio rerio]  
Length=289

Score = 197 bits (500), Expect = 2e-49, Method: Composition-based stats.  
Identities = 95/99 (95%), Positives = 98/99 (98%), Gaps = 0/99 (0%)

```
Query 3 VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKIY 62
      VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHK+HYKIY
Sbjct 9 VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKDHYKIY 68

Query 63 NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFKQN 101
      NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPF ++
Sbjct 69 NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFCE D 107
```

# BLAST Algorithm

- Scoring of matches done using scoring matrices
- Sequences are split into words (default  $n=3$ )
  - Speed, computational efficiency
- BLAST algorithm extends the initial “seed” hit into an HSP
  - HSP = high scoring segment pair = Local optimal alignment

# Credits

- Materials for this presentation have been adapted from the following sources:

Bioinformatics: A practical guide to the analysis of genes and proteins

- Questions? Please contact:

Dr. Joanne Fox

Michael Smith Laboratories

[joanne@mssl.ubc.ca](mailto:joanne@mssl.ubc.ca)

# BLAST

GUIDED TOUR: Advanced Tips & Tricks for Using BLAST





# <http://blast.ncbi.nlm.nih.gov/>

▶ [NCBI/BLAST Home](#)

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

**New** Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#).

## BLAST Assembled Genomes

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

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

## Basic BLAST

Choose a BLAST program to run.

- [nucleotide blast](#) Search a **nucleotide** database using a **nucleotide** query  
*Algorithms: blastn, megablast, discontinuous megablast*
- [protein blast](#) Search **protein** database using a **protein** query  
*Algorithms: blastp, psi-blast, phi-blast*
- [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

### News

#### [BLAST 2.2.23 release](#)

A new version of the stand-alone applications is available.

Mon, 22 Mar 2010 15:00:00 EST

 [More BLAST news...](#)

### Tip of the Day

#### [How to do Batch BLAST jobs.](#)

BLAST makes it easy to examine a large group of potential gene candidates.

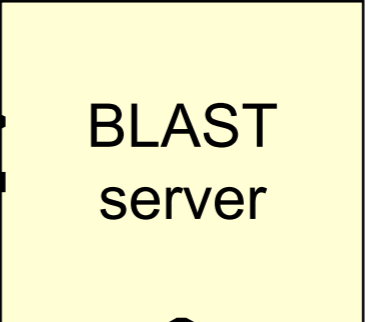
 [More tips...](#)

```
>gi|15237380|ref|NP_197163.1| myb family transcription factor (MYB43) [Arabidopsis thaliana]
MGRQPCCDKVGLKKGPTIEEDKKLINFILTNHGCCWRALPKLSGLLRGKSCRLRWINYLRPDLKRGLL
SEYEEQKVINLHAQLGNRWSKIASHLPGRTDNEIKNHWNTHIKKLRKMGIDPLTHKPLSEQEASQQAQG
RKKSLVPHDDKNPKDQQTQKDEQEQHLEQALEKNNTSVSGDGFIDEVPLLPHEILIDISSHHHHSN
DDNVNINTSKFTSPSSSSSTSSCISSVVPDGFSPKFFDEMEILDKWLSSDSSLGDDISKDGFNNSTV
DTMNLWDINDLSSLDMMFNEHDDGFIGNGGCRMLVLDQDSWTFDLL
```

**Submit Query**

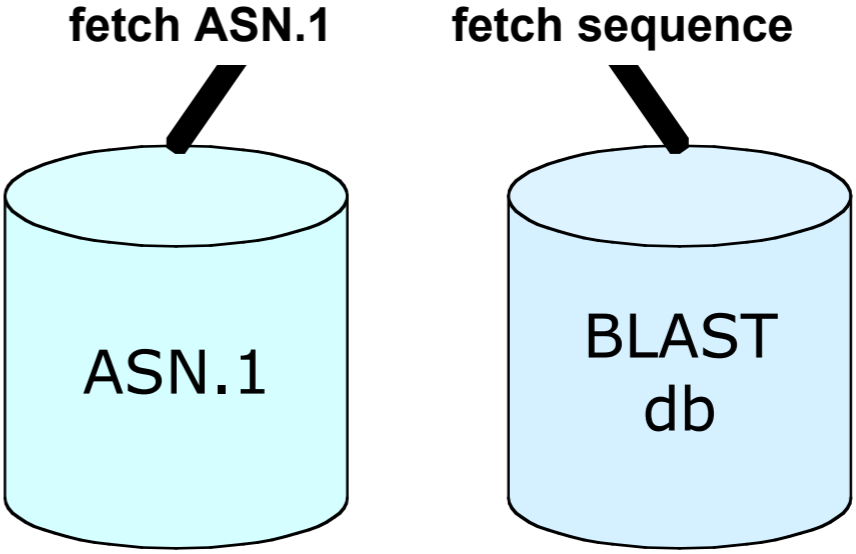
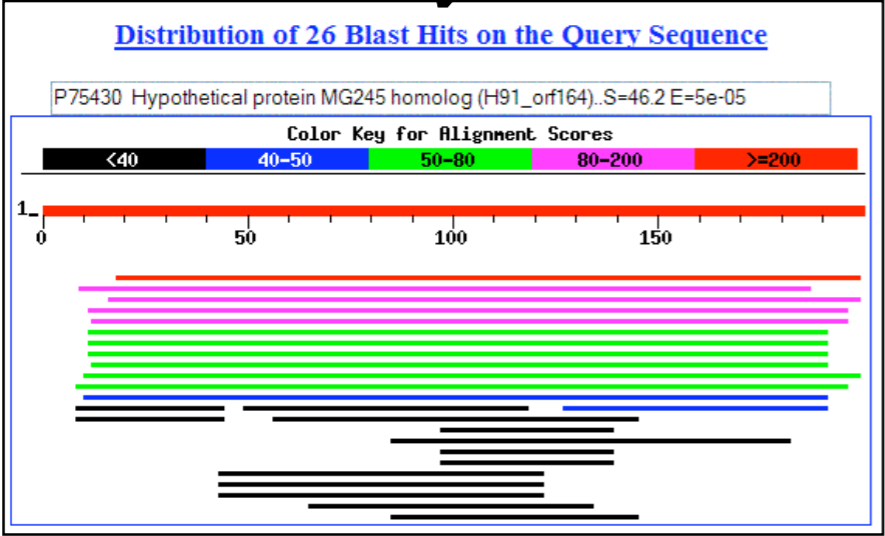


**Request Results**



**Return Formatted Results**

**Display Results**



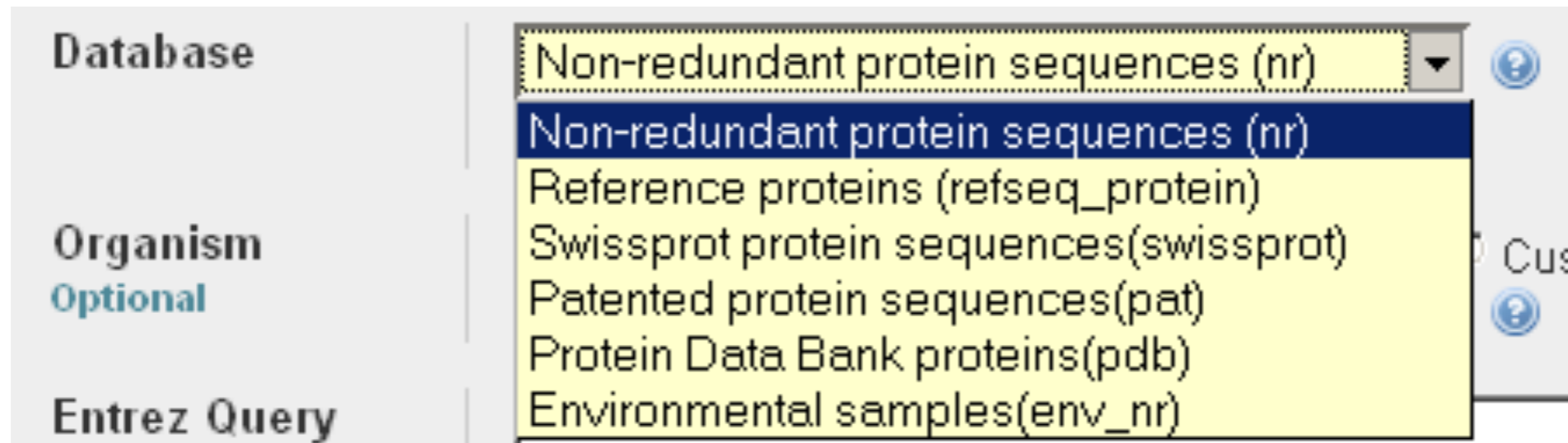
# Consider your research question ...

- Are you looking for a particular gene in a particular species?
- Are you looking for additional members of a protein family across all species?
- Are you looking to annotate genes in your species of interest?

# Know your reagents

- Changing your choice of database is changing your search space
- Database size affects the BLAST statistics
- Databases change rapidly and are updated frequently

# Protein Databases: nr



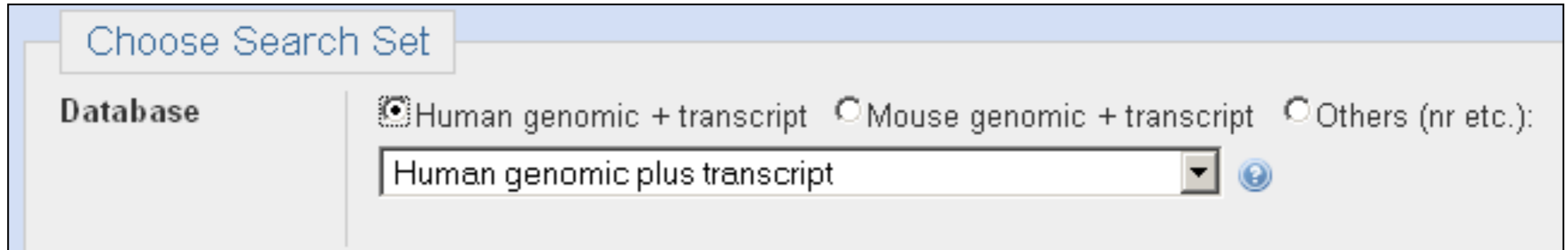
- nr (non-redundant protein sequences) default
  - GenBank CDS translations
  - Refseq Proteins
  - Outside Protein
    - PIR, Swiss-Prot, PRF
    - PDB (sequences from structures)
- pat protein patents
- env\_nr environmental samples

Services

blastp

blastx

# Nucleotide Databases: Human and Mouse



The screenshot shows a web interface titled "Choose Search Set". Under the "Database" section, there are three radio button options: "Human genomic + transcript" (which is selected), "Mouse genomic + transcript", and "Others (nr etc.):". Below these options is a dropdown menu currently displaying "Human genomic plus transcript" and a blue question mark icon.

- Human and mouse genomic + transcript default
- Separate sections in output for mRNA and genomic
- Direct links to Map Viewer for genomic sequences

Megablast, blastn service

# Nucleotide Databases: Traditional

Choose Search Set

Database

Nucleotide collection (nr/nt)

Nucleotide collection (nr/nt)

Reference mRNA sequences (refseq\_rna)

Reference genomic sequences (refseq\_genomic)

NCBI Genomes (chromosome)

Expressed sequence tags (est)

Non-human, non-mouse ESTs (est\_others)

Genomic survey sequences (gss)

High throughput genomic sequences (HTGS)

Patent sequences(pat)

Protein Data Bank (pdb)

Human ALU repeat elements (alu\_repeats)

Sequence tagged sites (dbsts)

Whole-genome shotgun reads (wgs)

Environmental samples (env\_nt)

Organism

Optional

Entrez Query

Optional

BLAST

Services

blastn

tblastn

tblastx

# Nucleotide Databases:

- **nr (nt)** Traditional GenBank
  - + RefSeq nucleotides
  - + PDB sequences
- **refseq\_rna**
- **refseq\_genomic** NC\_
- **NCBI genomes**
  - complete genomes
  - + chromosomes from RefSeq
- **est** expressed sequence tags
  - human + mouse, others
- **htgs** high throughput genomic
  - unfinished
- **gss** genome survey sequence
  - single-pass genomic data
- **pdb** protein data bank
  - derived from 3D structures
- **wgs**
  - whole genome shotgun
- **env\_nt**
  - environmental samples

Databases are mostly non-overlapping



# <http://blast.ncbi.nlm.nih.gov/>

The screenshot shows the NCBI BLAST website interface. At the top, there is a navigation bar with the BLAST logo and the text 'BLAST Basic Local Alignment Search Tool'. Below this are several tabs: 'Home', 'Recent Results', 'Saved Strategies', and 'Help'. The 'Help' tab is highlighted with a circular callout. To the right of the navigation bar, there is a 'My NCBI' section with a welcome message for 'joannealisonfox' and a 'Sign Out' link.

Below the navigation bar, the main content area is titled 'NCBI/BLAST/Help'. It contains a search box with the text 'Browse BLAST documentation.' and two columns of links:

- Getting Started**
  - [BLAST short course](#)
  - [BLAST program selection guide](#)
- Getting Help**
  - [Email blast-help](#)
  - [Mailing list](#)
- About BLAST**
  - [Frequently Asked Questions](#)
  - [NCBI Handbook: BLAST](#)
  - [The Statistics of Sequence Similarity Scores](#)
  - [NAR 2004 Web server issue](#)
  - [NAR 2006 Web server issue](#)
  - [BLAST glossary](#)
  - [References](#)
- BLAST information**
  - [Download BLAST Software and Databases](#)
  - [Developer information](#)

A callout box with the text 'Program Selection Guide' and a large arrow points to the 'BLAST program selection guide' link in the 'Getting Started' section.

At the bottom of the page, there is a 'BLAST News' section with a link to the 'BLAST News directory'.

### 3. Program Selection Tables

The appropriate selection of a BLAST program for a given search is influenced by the following three factors **1)** the nature of the query, **2)** the purpose of the search, and **3)** the database intended as the target of the search and its availability. The following tables provide recommendations on how to make this selection.

Table 3.1 Program Selection for Nucleotide Queries				
Length <sup>1</sup>	Database	Purpose	Program	Explanation
20 bp or longer  28 bp or above for megablast	<a href="#">Nucleotide</a>	Identify the query sequence	<a href="#">discontiguous megablast</a> , <a href="#">megablast</a> , or <a href="#">blastn</a>	<a href="#">Learn more ...</a>
		Find sequences similar to query sequence	<a href="#">discontiguous megablast</a> or <a href="#">blastn</a>	<a href="#">Learn more ...</a>
		Find similar sequence from the Trace archive	<a href="#">Trace megablast</a> , or <a href="#">Trace discontiguous megablast</a>	<a href="#">Learn more ...</a>
	<a href="#">Peptide</a>	Find similar proteins to translated query in a translated database	<a href="#">Translated BLAST (tblastx)</a>	<a href="#">Learn more ...</a>
	<a href="#">Peptide</a>	Find similar proteins to translated query in a protein database	<a href="#">Translated BLAST (blastx)</a>	<a href="#">Learn more ...</a>
7 - 20 bp	<a href="#">Nucleotide</a>	Find primer binding sites or map short contiguous motifs	<a href="#">Search for short, nearly exact matches</a>	<a href="#">Learn more ...</a>

**NOTE:**

<sup>1</sup> The cut-off is only a recommendation. For short queries, one is more likely to get matches if the "Search for short, nearly exact matches" page is used. Detailed discussion is in the [Section 4](#) below. With default setting, the shortest unambiguous query one can use is 11 for blastn and 28 for MEGABLAST.

Table 3.2 Program Selection for Protein Queries

Length <sup>1</sup>	Database	Purpose	Program	Explanation
15 residues or longer	<a href="#">Peptide</a>	Identify the query sequence or find protein sequences similar to the query	<a href="#">Standard Protein BLAST (blastp)</a>	<a href="#">Learn more</a> ...
		Find members of a protein family or build a custom position-specific score matrix	<a href="#">PSI-BLAST</a>	<a href="#">Learn more</a> ...
		Find proteins similar to the query around a given pattern	<a href="#">PHI-BLAST</a>	<a href="#">Learn more</a> ...
		Find conserved domains in the query	CD-search ( <a href="#">RPS-BLAST</a> )	<a href="#">Learn more</a> ...
		Find conserved domains in the query and identify other proteins with similar domain architectures	Conserved Domain Architecture Retrieval Tool ( <a href="#">CDART</a> )	<a href="#">Learn more</a> ...
	<a href="#">Nucleotide</a>	Find similar proteins in a translated nucleotide database	<a href="#">Translated BLAST (tblastn)</a>	<a href="#">Learn more</a> ...
5-15 residues	<a href="#">Peptide</a>	Search for peptide motifs	<a href="#">Search for short, nearly exact matches</a>	<a href="#">Learn more</a> ...

Note:

<sup>1</sup> The cut-off is only a recommendation. For short queries, one is more likely to get matches if the "Search for short, nearly exact matches" page is used. Detailed discussion is in [Section 4](#) below.

As genomic and other specialized sequence information is made available to the public, NCBI creates specialized BLAST pages for those sequences. The table below provides a general guide on how to select and use those special BLAST databases.

Table 3.3 Search against Organism Specific or Genome Databases <sup>1</sup>				
Query <sup>2</sup>	Database	Purpose	BLAST Pages to Use <sup>3</sup>	Explanation
Nucleotide: 20 or 28 bp and above  Protein: 15 residues and above	Human Genome	Map the query sequence	<a href="#">Human</a>	<a href="#">Learn more ...</a>
	Mouse Genome		<a href="#">Mouse</a>	<a href="#">Learn more ...</a>
	Rat Genome		<a href="#">Rat</a>	<a href="#">Learn more ...</a>
	Chimp, Cow, Dog, or Chicken Genome		<a href="#">Chimp</a> , or <a href="#">Cow</a> , <a href="#">Dog</a> , <a href="#">Chicken</a>	<a href="#">Learn more ...</a>
	Cat, Sheep, or Pig Genome	Determine the genomic structure	<a href="#">Cat</a> , <a href="#">Sheep</a> , or <a href="#">Pig</a>	<a href="#">Learn more ...</a>
	Zebrafish or Fugu (Pufferfish)	Identify novel genes	<a href="#">Zebrafish</a> or <a href="#">Fugu rubripes</a>	<a href="#">Learn more ...</a>
	Insects (flies and honeybees)		<a href="#">Insects</a>	<a href="#">Learn more ...</a>
	Nematodes (worms)	Find homologs	<a href="#">Nematodes</a>	<a href="#">Learn more ...</a>
	Plants		<a href="#">Plants</a>	<a href="#">Learn more ...</a>
	Fungi Genomes (including yeasts)	Other data mining	<a href="#">Fungi</a>	<a href="#">Learn more ...</a>
	Protozoa		<a href="#">Protozoa</a>	<a href="#">Learn more ...</a>
	Environmental Samples		<a href="#">Environmental Samples</a>	<a href="#">Learn more ...</a>
	Other Lower Eukaryotic Genomes		<a href="#">Other eukaryotes genomes</a>	<a href="#">Learn more ...</a>
Microbial Genomes	<a href="#">Microbial genomes</a>		<a href="#">Learn more ...</a>	

**NOTE:**

<sup>1</sup> Those pages access the genome database consisting of contig assemblies and other sequences specific to the organisms. Not all organisms listed here have genome assemblies available.

<sup>2</sup> Sequence length is only a suggestion. For most of the pages, the search parameters can be modified to enable searches with a short query by pasting additional options in the "Advanced Options" text box. For protein comparisons, -F F -e 20000 -W 2 should be used. For nucleotide comparison, use -F F -e 1000 -W 7. This also requires the uncheck of the megablast checkbox.

<sup>3</sup> Available databases and their contents are described in Section 5.

BLAST pages for special purposes are listed under Special and Meta sections. Their functions are described in Table 3.4 below.

Table 3.4 Function of Special BLAST Pages under Special/Meta Sections				
Query <sup>1</sup>	Database	Purpose	BLAST Page to Use	Explanation
Nucleotide: 11 bp or above Protein: 15 or above	- <sup>2</sup>	Compare two sequences directly	<a href="#">Align two sequences</a>	<a href="#">Learn more ...</a>
	Immunoglobulin sequences	Find matches to curated immunoglobulin sequences	<a href="#">igBLAST</a>	<a href="#">Learn more ...</a>
Nucleotide: 20 or 28 bp and above	UniVec	Screen for vector contamination	<a href="#">VecScreen</a>	<a href="#">Learn more ...</a>
	GEO	Find matches to sequences with MicroArray information	<a href="#">GEO BLAST</a>	<a href="#">Learn more ...</a>
	SNP	Find matches to human reference SNPs	<a href="#">SNP BLAST</a>	<a href="#">Learn more ...</a>
-	- <sup>3</sup>	To retrieve results for a search with its RID	<a href="#">Retrieve result for an RID</a>	<a href="#">Learn more ...</a>

Note:

<sup>1</sup> The query sequence length is only a suggestion. For most of the pages, the search parameters can be modified to enable better handling of short query by pasting additional options in the "Advanced Options" text box. For protein comparisons, -F F -e 20000 -W 2 should be used. For nucleotide comparison, use -F F -e 2000 -W 7.

<sup>2</sup> "Align two sequences" treats the second sequence as the database.

<sup>3</sup> Requires valid RIDs that are assigned within the past 24 hours.



► [NCBI/BLAST Home](#)

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

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

## BLAST Assembled Genomes

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

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

## Basic BLAST

Choose a BLAST program to run.

[nucleotide blast](#)

Search a **nucleotide** database using a **nucleotide** query  
*Algorithms: blastn, megablast, discontinuous megablast*

[protein blast](#)

Search **protein** database using a **protein** query  
*Algorithms: blastp, psi-blast, phi-blast*

[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

### News

[Old BLAST Web Pages to be deleted June 11th 2007](#)

As previously announced access to the old pages will be removed on June 11, 2007.

2007-06-01 12:15:00

[More BLAST news...](#)

### Tip of the Day

**How to use BLAST to find human sequences in a database that can be amplified with a particular primer pair.**

A frequent use of nucleotide-nucleotide BLAST is to check the specificity of oligonucleotides for hybridization in PCR. The goal is usually to make sure that the primers will give a unique product from the target genome or cDNA.

### Enter Query Sequence

Enter accession number, gi, or FASTA sequence

[Clear](#)

Query subrange

231571

**231571**

From

To

Or, upload file

[Browse...](#)

Job Title

Q02067:Achaete-scute homolog 1 (Mash-1)

Enter a descriptive title for your BLAST search

### Choose Search Set

Database

Swissprot protein sequences(swissprot)

Organism  
Optional

Enter organism name or id--completions will be suggested

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

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)

Choose a BLAST algorithm

Let's look at some of the options!

**BLAST**

Search database **swissprot** using **Blastp (protein-protein BLAST)**

Show results in a new window

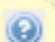
▼ [Algorithm parameters](#)

Note: Parameter values that differ from the default are highlighted in yellow

# Context Specific Help

Choose Search Set


**Database**

Swissprot protein sequences(swissprot) 

Select the sequence database to run searches against. No BLAST database contains all the sequences at NCBI. BLAST databases are organized by informational content (nr, RefSeq, etc.) or by sequencing technique (WGS, EST, etc.). [more...](#)


**Organism**  
**Optional**

Enter organism name or id--completions will be suggested

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

Select from the list or choose "Custom" to enter the name of an organism. The search will be restricted to the sequences in the database which are from the organism selected.

**Entrez Query**  
**Optional**

Enter an Entrez query to limit search 

You can use Entrez query syntax to search a subset of the selected BLAST database. This can be helpful to limit searches to molecule types, sequence lengths or to exclude organisms. [more...](#)



# Limiting Database: Organism

Organism  
Optional

Any  Human  *A.thaliana*  Mouse  Custom...

Search: bacter

- CFB group **bacteria** (taxid:976)
- GNS **bacteria** (taxid:200795)
- green sulfur **bacteria** (taxid:1090)
- Bacteria** (taxid:2)
- purple **bacteria** and relatives (taxid:1224)
- purple non-sulfur **bacteria** (taxid:1224)
- purple photosynthetic **bacteria** (taxid:1224)
- purple photosynthetic **bacteria** and relatives (taxid:1224)
- purple **bacteria** (taxid:1224)
- low G+C Gram-positive **bacteria** (taxid:1239)


taxa will be shown.

Organism autocomplete

# Limiting Database: Entrez Query

Entrez Query  
Optional

all[filter] NOT mammals[organism]

Enter an Entrez query to limit search 

all[filter] NOT mammals[organism]

gene\_in\_mitochondrion[Properties]  
2006:2007 [Modification Date]

Nucleotide

biomol\_mrna[Properties]

biomol\_genomic[Properties]

## BLAST

Search database **swissprot** using **Blastp (protein-protein BLAST)**

Show results in a new window

### Algorithm parameters

Note: Parameter values that differ from the default

#### General Parameters

Max target sequences

100

Select the maximum number of aligned sequences to display

Short queries

Automatically adjust parameters for short input sequences

Expect threshold

10

Word size

3

#### Scoring Parameters

Matrix

BLOSUM62

Gap Costs

Existence: 11 Extension: 1

# Algorithm parameters: Protein

The image shows a screenshot of a web-based interface for protein sequence alignment parameters. The interface is organized into four main sections: General Parameters, Scoring Parameters, Filters and Masking, and a partially visible section below. A yellow arrow labeled "Expand" points to the "Algorithm parameters" header. A dropdown menu for "Max target sequences" is open, showing options from 10 to 20000. Callout boxes provide context for several settings: "Max target sequences" (100) is noted as "May limit results"; "Expect threshold" (10) is noted as "Adjust to set stringency"; "Compositional adjustments" (Composition-based statistics) is noted as "Default statistics adjustment for compositional bias"; and the "Low complexity regions" filter is noted as "Off now by default. Conflicts with comp-based stats".

**Algorithm parameters** (Expand)

**General Parameters**

- Max target sequences:** 100 (May limit results)
- Short queries:**  Automatic
- Expect threshold:** 10 (Adjust to set stringency)
- Word size:** 3

**Scoring Parameters**

- Matrix:** BLOSUM62
- Gap Costs:** Existence: 11 Extension: 1
- Compositional adjustments:** Composition-based statistics (Default statistics adjustment for compositional bias)

**Filters and Masking**

- Filter:**  Low complexity regions (Off now by default. Conflicts with comp-based stats)
- Mask:**  Mask for lookup table only;  Mask lower case letters

# Automatic Short Sequence Adjustment

Job Title: Elvis Lives!

No putative conserved domains have been detected

Your search parameters were adjusted to search for a short input sequence.

WAITING

Request ID 1WSB0FX012

Status

Subr

Curre

Time

This p

e-value	200000
Word Size	2
Matrix	PAM30
Gap Costs	-9, -1
Comp Stats	Off
Low Comp Filter	Off

> [ref|ZP\\_01712014.1|](#) conserved hypothetical protein [Pseudomonas putida] Length=245


Score = 18.5 bits (36), Expect = 15305  
Identities = 5/5 (100%), Positives = 5/5 (100%), Gaps = 0/5 (0%)

Query 1 ELVIS 5  
          ELVIS  
Sbjct 126 ELVIS 130

> [ref|ZP\\_01712512.1|](#) Substrate-binding region of ABC-type glycine binding system [Pseudomonas putida GB-1] Length=342

Score = 18.5 bits (36), Expect = 15305  
Identities = 5/5 (100%), Positives = 5/5 (100%), Gaps = 0/5 (0%)

Query 1 ELVIS 5  
          ELVIS  
Sbjct 172 ELVIS 176

> [ref|XP\\_001366374.1|](#)  PREDICTED: similar to R7 binding protein [Mycobacterium tuberculosis H37Rv] Length=257

Score = 18.5 bits (36), Expect = 15305  
Identities = 5/5 (100%), Positives = 5/5 (100%), Gaps = 0/5 (0%)

Query 1 ELVIS 5  
          ELVIS  
Sbjct 69 ELVIS 73

> [ref|ZP\\_01711731.1|](#) GCN5-related N-acetyltransferase [Caldivirga marisnigri] Length=166

Score = 18.5 bits (36), Expect = 15305  
Identities = 5/5 (100%), Positives = 5/5 (100%), Gaps = 0/5 (0%)

Query 1 ELVIS 5  
          ELVIS  
Sbjct 20 ELVIS 24

## Enter Query Sequence

Enter accession number, gi, or FASTA sequence 

[Clear](#)


Query subrange 

```
>gi|231571|sp|Q02067|ASCL1_MOUSE Achaete-scute homolog 1  
(Mash-1)  
MESSGKMEAGAGQPPQPPQPPFLPPAACFFATAAAAAAAAAAAAAAQAQQQQPQAPPQQAPQLS  
GGGHKSAAKQDKRQRSSPELMRCKRRLNFSGFCYSLPQQQPAAVARRNERERNRVKLVNLG  
PNGAANKMSKVETLRSVQYIRALQQLLDEHDAVSAAFQAGVLSPTISPNYSNDLNSMAGS
```


From

To

Or, upload file




Job Title

Enter a descriptive title for your BLAST search 


## Choose Search Set

Database




Organism

Optional

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

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)

Choose a BLAST algorithm 

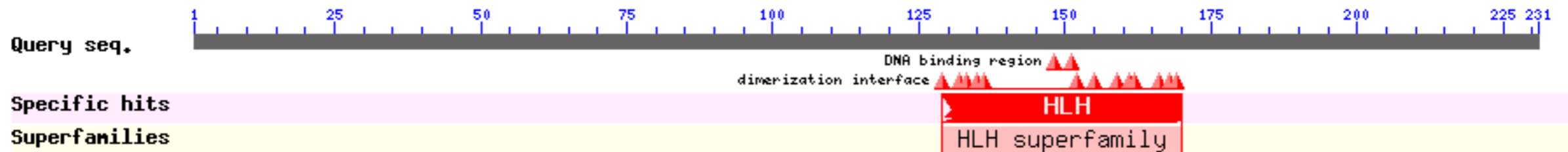
**BLAST**

Search database **swissprot** using **Blastp (protein-protein BLAST)**

Show results in a new window

**Job Title: Q02067:RecName: Full=Achaete-scute homolog...**

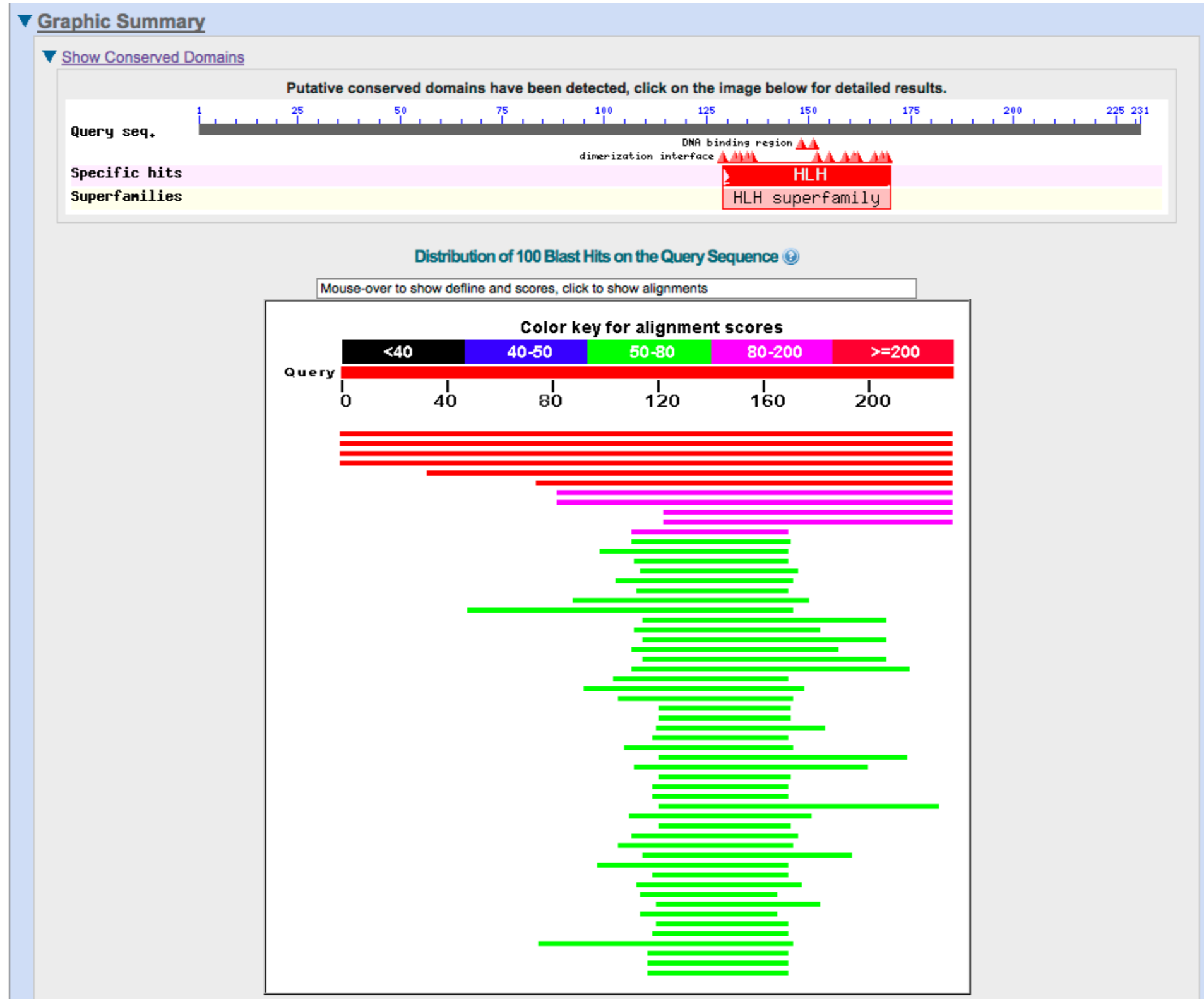
Putative conserved domains have been detected, click on the image below for detailed results.



Request ID	T9U0ZFN4011
Status	Searching
Submitted at	Thu Feb 12 22:25:19 2009
Current time	Thu Feb 12 22:25:26 2009
Time since submission	00:00:06

This page will be automatically updated in **78** seconds

# A graphical view





Full=Masn-1

Program BLASTP 2.2.19+ [Citation](#)

Molecule type amino acid

Query Length 231

Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#)[Graphic Summary](#)[Descriptions](#)

Sequences producing significant alignments:

			Score (Bits)	E Value	
<a href="#">sp Q02067.1 ASCL1</a>	MOUSE	RecName: Full=Achaete-scute homolog 1...	<a href="#">466</a>	4e-131	G
<a href="#">sp P19359.1 ASCL1</a>	RAT	RecName: Full=Achaete-scute homolog 1	<a href="#">347</a>	4e-95	G
<a href="#">sp P50553.2 ASCL1</a>	HUMAN	RecName: Full=Achaete-scute homolog 1...	<a href="#">332</a>	1e-90	G
<a href="#">sp Q90259.1 ASL1A</a>	DANRE	RecName: Full=Achaete-scute homolog 1...	<a href="#">298</a>	1e-80	G
<a href="#">sp Q06234.1 ASCL1</a>	XENLA	RecName: Full=Achaete-scute homolog 1	<a href="#">289</a>	9e-78	G
<a href="#">sp Q90260.1 ASL1B</a>	DANRE	RecName: Full=Achaete-scute homolog 1...	<a href="#">217</a>	3e-56	G
<a href="#">sp Q2EGB9.1 ASCL2</a>	BOVIN	RecName: Full=Achaete-scute homolog 2...	<a href="#">135</a>	1e-31	G
<a href="#">sp Q99929.2 ASCL2</a>	HUMAN	RecName: Full=Achaete-scute homolog 2...	<a href="#">124</a>	3e-28	G
<a href="#">sp P19360.1 ASCL2</a>	RAT	RecName: Full=Achaete-scute homolog 2; ...	<a href="#">106</a>	8e-23	G
<a href="#">sp O35885.2 ASCL2</a>	MOUSE	RecName: Full=Achaete-scute homolog 2...	<a href="#">103</a>	1e-21	G
<a href="#">sp Q7RTU5.2 ASCL5</a>	HUMAN	RecName: Full=Achaete-scute homolog 5	<a href="#">80.5</a>	6e-15	G
<a href="#">sp Q6XD76.1 ASCL4</a>	HUMAN	RecName: Full=Achaete-scute homolog 4...	<a href="#">78.2</a>	4e-14	G
<a href="#">sp Q9NQ33.2 ASCL3</a>	HUMAN	RecName: Full=Achaete-scute homolog 3...	<a href="#">75.9</a>	2e-13	G
<a href="#">sp Q9JJR7.1 ASCL3</a>	MOUSE	RecName: Full=Achaete-scute homolog 3...	<a href="#">75.1</a>	3e-13	G
<a href="#">sp P10083.1 AST5</a>	DROME	RecName: Full=Achaete-scute complex pr...	<a href="#">74.7</a>	3e-13	G
<a href="#">sp P10084.2 AST4</a>	DROME	RecName: Full=Achaete-scute complex pr...	<a href="#">71.6</a>	3e-12	G
<a href="#">sp Q10007.1 HLH6</a>	CAEEL	RecName: Full=Helix-loop-helix protein 6	<a href="#">64.3</a>	5e-10	G
<a href="#">sp P00334.2 AST3</a>	DROME	RecName: Full=Achaete-scute complex pr...	<a href="#">63.2</a>	1e-09	G

# Re-Format and/or Download your BLAST results

[Edit and Resubmit](#) [Save Search Strategies](#) [Formatting options](#) [Download](#)

[Reformat](#)

**Formatting options**

**Show** Alignment as HTML  Advanced View  Use old BLAST report format [Reset form to defaults](#)

**Alignment View** Pairwise

**Display**  Graphical Overview  Linkout  Sequence Retrieval  NCBI-gi

Masking Character: Lower Case Masking Color: Grey

**Limit results** Descriptions: 100 Graphical overview: 100 Alignments: 100

Organism Type common name, binomial, taxid, or group name. Only 20 top taxa will be shown.

Entrez query:

Expect Min: Expect Max:

**Format for**  PSI-BLAST with inclusion threshold:

**Download**

Alignment					Search Strategies	Bioseq
<a href="#">Text</a>	<a href="#">XML</a>	<a href="#">ASN.1</a>	<a href="#">Hit Table(text)</a>	<a href="#">Hit Table(csv)</a>	<a href="#">ASN.1</a>	<a href="#">ASN.1</a>



# BLAST Alignments

```
>| sp|P20389|MYC2_MARMO N-myc 2 proto-oncogene protein  
Length=454
```

```
Score = 35.8 bits (81), Expect = 0.14, Method: Composition-based stats.  
Identities = 22/52 (42%), Positives = 30/52 (57%), Gaps = 4/52 (7%)
```

```
Query 133 FATLREHVPNGAANKKMSKVETLRSVQYIRALQ----QLLDEHDAVSAAFQ 180  
      F TLR+HVP      N+K +KV L+ A +Y+ LQ      QLL E + + A Q  
Sbjct 391 FTTLRDHVPPELVKNEKAAKVVILKKACEYVHYLQAKEHQLLMEKEKLOARQQ 442
```

Identical match

positive score  
(conservative)

gap

Negative or zero

# BLAST Alignments

> [\[sp|P04198|MYCN HUMAN\]](#) **G** N-myc proto-oncogene protein  
Length=464

Score = 35.4 bits (80), Expect = 0.025, Method: Composition-based stats.  
Identities = 22/52 (42%), Positives = 31/52 (59%), Gaps = 4/52 (7%)

```
Query 133 FATLREHVPNGAANKKMSKVETLRSVQYIRALQ----QLLDEHDAVSAAFQ 180
          F TLR+HVP      N+K +KV  L+ A +Y+ +LQ      QLL E + + A  Q
Sbjct 401 FLTLRDHVPELVKNEKAAKVVILKKATEYVHSLQAEEHQLLLEKEKQLQARQQ 452
```

> [\[sp|Q02363|ID2 HUMAN\]](#) **G** DNA-binding protein inhibitor ID-2 (Inhibitor of DNA binding 2)  
Length=134

Score = 35.4 bits (80), Expect = 0.025, Method: Composition-based stats.  
Identities = 19/47 (40%), Positives = 29/47 (61%), Gaps = 0/47 (0%)

```
Query 129 VNLGFATLREHVPNGAANKKMSKVETLRSVQYIRALQQLLDEHDAV 175
          +N ++ L+E VP+      NKK+SK+E L+ + YI  LQ  LD H  +
Sbjct 39  MNDCYSKCLKELVPSIPQNKVKSKMEILQHVIDYILDQLIALDSHPTI 85
```

> [\[sp|P12980|LYL1 HUMAN\]](#) **G** Protein lyl-1 (Lymphoblastic leukemia-derived sequence 1)  
Length=267

Score = 35.4 bits (80), Expect = 0.025, Method: Composition-based stats.  
Identities = 22/50 (44%), Positives = 31/50 (62%), Gaps = 0/50 (0%)

```
Query 129 VNLGFATLREHVPNGAANKKMSKVETLRSVQYIRALQQLLDEHDAVSAA 178
          VN  FA LR+ +P      ++K+SK E LR A++YI  L +LL +  A  AA
Sbjct 153 VNGAFAELRKLLPHTPPDRKLSKNEVLRRLAMKYIGFLVRLLRDQAAALAA 202
```

- **Similarity**

The extent to which nucleotide or protein sequences are related. The extent of similarity between two sequences can be based on percent sequence identity and/or conservation. In BLAST similarity refers to a positive matrix score.

- **Identity**

The extent to which two (nucleotide or amino acid) sequences are invariant.

- **Homology**

Similarity attributed to descent from a common ancestor.

It is your responsibility as an informed bioinformatician to use these terms correctly: A sequence is either homologous or not. Don't use % with this term!

# Sorting BLAST by Taxonomy

BLAST

Basic Local Alignment Search Tool

My NCBI

[Home](#)

[Recent Results](#)

[Saved Strategies](#)

[Help](#)

[\[Sign In\]](#) [\[Regis](#)

NCBI/ BLAST/ blastp suite/ Formatting Results - T9U0ZFN4011

[Edit and Resubmit](#) [Save Search Strategies](#) [▶ Formatting options](#) [▶ Download](#)

Q02067:RecName: Full=Achaete-scute homolog...

**Query ID** gi|231571|sp|Q02067.1|ASCL1\_MOUSE  
**Description** RecName: Full=Achaete-scute homolog 1; AltName:  
Full=Mash-1  
**Molecule type** amino acid  
**Query Length** 231

**Database Name** swissprot  
**Description** Non-redundant SwissProt sequences  
**Program** BLASTP 2.2.19+ [▶ Citation](#)

Other reports: [▶ Search Summary](#) [\[Taxonomy reports\]](#) [▶ Distance tree of results\]](#)

▶ [Graphic Summary](#)

▼ [Descriptions](#)

Sequences producing significant alignments:			Score (Bits)	E Value	
<a href="#">sp Q02067.1 ASCL1_MOUSE</a>	RecName: Full=Achaete-scute homolog 1...		<a href="#">466</a>	4e-131	<a href="#">G</a>
<a href="#">sp P19359.1 ASCL1_RAT</a>	RecName: Full=Achaete-scute homolog 1		<a href="#">347</a>	4e-95	<a href="#">G</a>
<a href="#">sp P50553.2 ASCL1_HUMAN</a>	RecName: Full=Achaete-scute homolog 1...		<a href="#">332</a>	1e-90	<a href="#">G</a>
<a href="#">sp Q90259.1 ASL1A_DANRE</a>	RecName: Full=Achaete-scute homolog 1...		<a href="#">298</a>	1e-80	<a href="#">G</a>
<a href="#">sp Q06234.1 ASCL1_XENLA</a>	RecName: Full=Achaete-scute homolog 1		<a href="#">289</a>	9e-78	<a href="#">G</a>



Job Title: gi|231571 (231 letters)

Show Conserved Domains

Tax BLAST Report

Index

- Lineage Report
Organism Report
Taxonomy Report
Help

Lineage Report

Table with 4 columns: Taxonomy (e.g., Bilateria, Coelomata, Euteleostomi), Count, Hits (e.g., 22 hits), and Hit Names (e.g., Achaete-scute homolog 1 (Mash-1)).

Organism Report

Table with 4 columns: Organism (Mus musculus), Accession (sp\_Q02067), Gene (ASCL1 MOUSE), and Statistics (466 hits, 4e-131).

# Distance Tree of Results

Tree view for rid: **T9U0ZFN4011**, query ID: **sp|Q02067.1**, database: **swissprot**

This tree was produced using BLAST pairwise alignments. [more...](#)

BLAST computes a pairwise alignment between a query and the database sequences searched. It does not explicitly compute an alignment between the different database sequences (i.e., does not perform a multiple alignment). For purposes of this sequence tree presentation an implicit alignment between the database sequences is constructed, based upon the alignment of those (database) sequences to the query. It may often occur that two database sequences align to different parts of the query, so that they barely overlap each other or do not overlap at all. In that case it is not possible to calculate a distance between these two sequences and only the higher scoring sequence is included in the tree.

Tree method: **Fast Minimum Evolution** | Max Seq Difference: **0.85** | Distance: **Grishin (protein)** | **Reset** | **Download**

## Tree Method:

Algorithm used to produce a tree from given distances (or dissimilarities) between sequences. Available options:

- 1) Fast Minimum Evolution (*Desper R and Gascuel O, Mol Biol Evol 21:587-98, 2004*)
- 2) Neighbor Joining (*Saitou N and Nei M, Mol Biol Evol, 4:406-25, 2004*)

**Note:** Both algorithms produce un-rooted tree such as ones shown as *radial* or *force* in the tabs below. The rooted trees are created by placing a root in the middle of the longest edge.

read more in context specific help menus

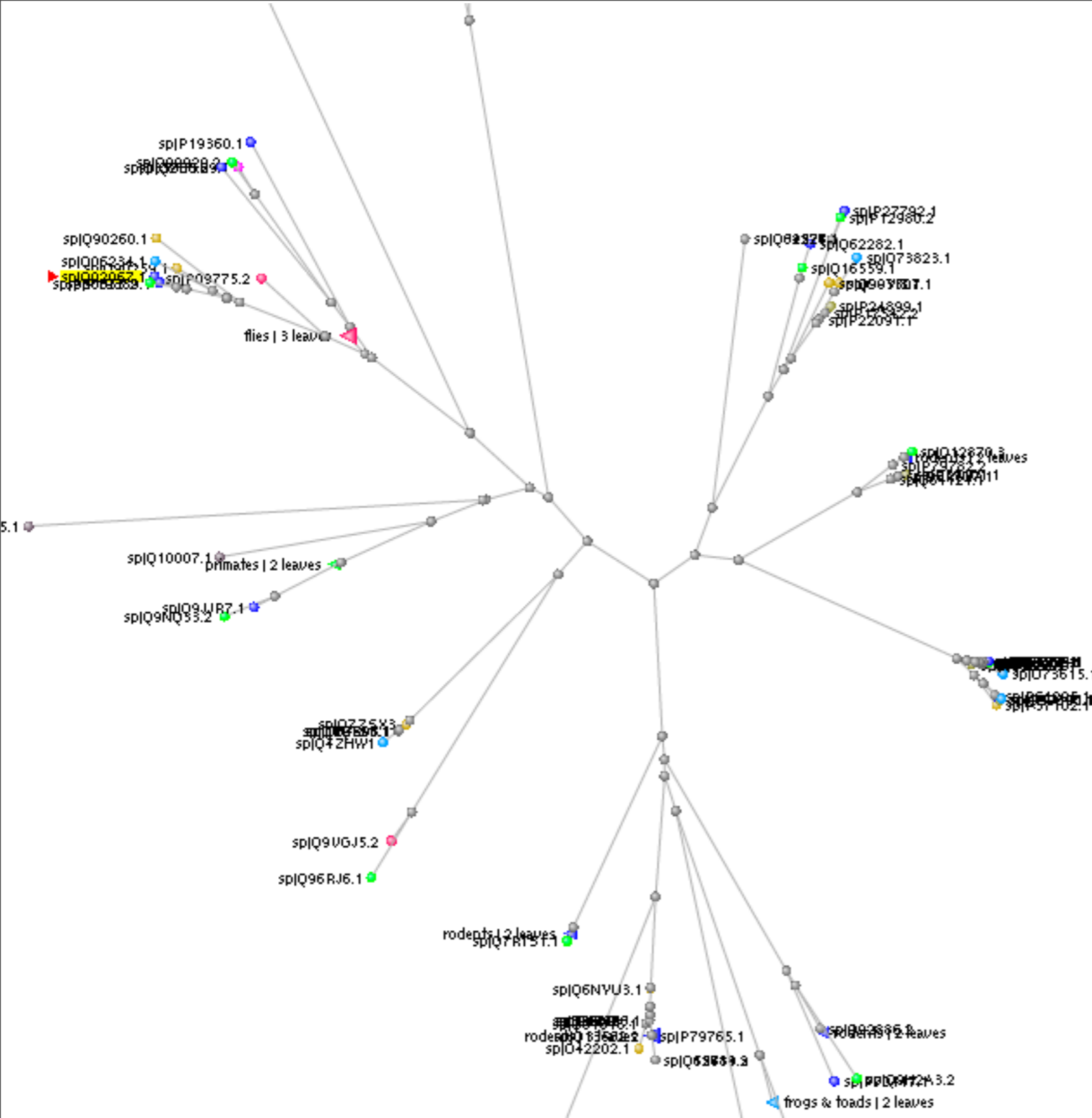
This distance tree is based on BLAST pairwise alignments

≠ phylogenetic tree

**rectangle** | **slanted** | **radial** | **force** |  Show distance | Mouse over an internal node for a







- ✓ Rectangle: rectangular shaped rooted tree, where root is placed in the longest edge
- ✓ Slanted: similar to rectangle, but with triangular tree shape
- ✓ Radial: un-rooted tree
- ✓ Force: similar to radial, where nodes are pushed away from one another for better presentation.

# New Alignment of Multiple Sequences



BLAST

Basic Local Alignment Search Tool

My NCBI

Welcome joannealison

[Home](#)

[Recent Results](#)

[Saved Strategies](#)

[Help](#)

NCBI/ BLAST/ blastp suite/ Formatting Results - ZV7DHV5E01N

[Edit and Resubmit](#) [Save Search Strategies](#) [Formatting options](#) [Download](#)

sp|Q02067| (231 letters)

**Query ID** [gi|231571|sp|Q02067.1|ASCL1\\_MOUSE](#)  
**Description** RecName: Full=Achaete-scute homolog 1; Short=ASH-1;  
Short=mASH-1; Short=mASH1 >gi|193876|gb|AAA37780.1|  
helix-loop-helix protein [Mus musculus]  
>gi|15131817|gb|AAK84426.1| achaete-scute complex  
homolog-like 1 [Mus musculus]

**Molecule type** amino acid

**Query Length** 231

Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Multiple alignment](#)

**Database Name** swissprot

**Description** Non-redundant SwissProt sequences

**Program** BLASTP 2.2.23+ [Citation](#)

[Graphic Summary](#)

[Descriptions](#)

[Alignments](#)

Select All

[Get selected sequences](#)

[Distance tree of results](#)

[Multiple alignment](#)

>  [sp|Q02067.1|ASCL1\\_MOUSE](#) RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=mASH-1;  
Short=mASH1

[gb|AAA37780.1|](#) helix-loop-helix protein [Mus musculus]

[gb|AAK84426.1|](#) achaete-scute complex homolog-like 1 [Mus musculus]  
Length=231

# BLAST results sent to COBALT

generates MSA

## Cobalt Results - sp|Q02067| (231 letters) - Cobalt RID ZV85B6BE212 (100 seqs)

Descriptions Select All Re-align Alignment parameters

Legend for links to other resources: U UniGene E GEO G Gene S Structure M Map Viewer

Accession	Description	Links
<a href="#">Q02067.1</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=mASH-1; Short=mASH1 >gi 193876 gb AAA37780.1  helix-loop-heli	<a href="#">G</a>
<a href="#">Q02067.1</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=mASH-1; Short=mASH1 >gi 193876 gb AAA37780.1  helix-loop-heli	<a href="#">G</a>
<a href="#">P19359.1</a>	RecName: Full=Achaete-scute homolog 1 >ref NP_071779.1  achaete-scute homolog 1 [Rattus norvegicus] >emb CAA37760.1  unn	<a href="#">G</a>
<a href="#">P50553.2</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=hASH1; AltName: Full=Class A basic helix-loop-helix protein 46; Shr	<a href="#">G</a>
<a href="#">Q90259.1</a>	RecName: Full=Achaete-scute homolog 1a; Short=Zash-1a; AltName: Full=Pituitary-absent protein >ref NP_571294.1  achaete-scute	<a href="#">G</a>
<a href="#">Q06234.1</a>	RecName: Full=Achaete-scute homolog 1 >ref NP_001079247.1  achaete-scute complex homolog 1 [Xenopus laevis] >gb AAA4964	<a href="#">G</a>
<a href="#">Q90260.1</a>	RecName: Full=Achaete-scute homolog 1b; Short=Zash-1b >ref NP_571306.1  achaete-scute homolog 1b [Danio rerio] >gb AAA788	<a href="#">G</a>
<a href="#">Q2EGB9.1</a>	RecName: Full=Achaete-scute homolog 2; AltName: Full=Mash2 >gb ABD39719.1  achaete scute-like protein 2 [Bos taurus] >gb AA	<a href="#">G</a>
<a href="#">Q99929.2</a>	RecName: Full=Achaete-scute homolog 2; Short=ASH-2; Short=hASH2; AltName: Full=Mash2; AltName: Full=Class A basic helix-lo	<a href="#">G</a>
<a href="#">P19360.1</a>	RecName: Full=Achaete-scute homolog 2; AltName: Full=Mash2 >ref NP_113691.1  achaete-scute homolog 2 [Rattus norvegicus] >	<a href="#">G</a>
<a href="#">Q35885.2</a>	RecName: Full=Achaete-scute homolog 2; Short=ASH-2; Short=mASH-2; Short=mASH2 >ref NP_032580.2  achaete-scute homolog	<a href="#">G</a>
<a href="#">Q7RTU5.2</a>	RecName: Full=Achaete-scute homolog 5; Short=ASH-5; Short=hASH5; AltName: Full=Class A basic helix-loop-helix protein 47; Shr	<a href="#">G</a>
<a href="#">Q6XD76.1</a>	RecName: Full=Achaete-scute homolog 4; Short=ASH-4; Short=hASH4; AltName: Full=Achaete-scute-like protein 4; AltName: Full=	<a href="#">G</a>
<a href="#">Q9NQ33.2</a>	RecName: Full=Achaete-scute homolog 3; Short=ASH-3; Short=hASH3; AltName: Full=Class A basic helix-loop-helix protein 42; Shr	<a href="#">G</a>
<a href="#">Q9JJR7.1</a>	RecName: Full=Achaete-scute homolog 3; Short=ASH-3; Short=mASH-3; Short=mASH3; AltName: Full=bHLH transcriptional regula	<a href="#">G</a>
<a href="#">P10083.1</a>	RecName: Full=Achaete-scute complex protein T5; AltName: Full=Protein achaete >ref NP_476824.1  achaete [Drosophila melanog	<a href="#">G</a>
<a href="#">P10084.2</a>	RecName: Full=Achaete-scute complex protein T4; AltName: Full=Protein scute >ref NP_476803.1  scute [Drosophila melanogaster]	<a href="#">G</a>
<a href="#">Q10007.1</a>	RecName: Full=Helix-loop-helix protein 6 >ref NP_496070.1  Helix Loop Helix family member (hlh-6) [Caenorhabditis elegans] >emb	<a href="#">G</a>
<a href="#">P09774.2</a>	RecName: Full=Achaete-scute complex protein T3; AltName: Full=Protein lethal of scute; Short=Lethal of sc >ref NP_476623.1  leth	<a href="#">G</a>
<a href="#">Q10574.2</a>	RecName: Full=Protein lin-32; AltName: Full=Abnormal cell lineage protein 32 >ref NP_508410.2  abnormal cell LINeage family men	<a href="#">G</a>

```

-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 129
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 128
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 129
-----MKRRRR----LRSDAEMQQ----LRQ-AANVRERRRMQSINDAFEGLRSHIP 147
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 132
-----Q---AGNCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111
-----RRGPKK----KKMTKARMQRFKMRRM-KANARERNRMHGLNDALESRLKRVVP 124
-----RRASSG----A-G----PVVVVRQRQ-AANARERDRTQSVNTAFTALRTLIP 98
-----S---PGRLE---ALGG-----RLPRRKG-SGPKKERRRTEINSAFELRECIPI 122
-----RRRRPG----PSGPGGRRDSSIQRRL-ESNERERQRMHKLNNAFQALREVIP 103
-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 101
-----LKRERRR---MRSEVEMQQ----LRQ-AANVRERRRMQSINDAFEGLRSHIP 143
-----RRASNG----A-G----PVVVVRQRQ-AANARERDRTQSVNTAFTALRTLIP 98
-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 113
-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 122

```

<a href="#">Q0VCE2</a>	84	-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 125
<a href="#">Q90691</a>	69	-----G---AGRLE---ALSG-----RLGRRKGVGGPKKERRRTEINSAFELRECIPI 111
<a href="#">Q91616</a>	84	-----RRGPKK----KKMTKARVERFKVRRM-KANARERNRMHGLNDALDSLKRVVP 130
<a href="#">Q6QHK4</a>	60	-----L-QL-----VLERRR-VANAKERERIKNLNRGFARLKLALVP 93
<a href="#">P17542</a>	142	QPLASLGSGFFGEPDAPPMFTNNRVKRRSPYE----MEITDGPHT-KVVRRI-FTNSRERWRQNVNGAFELRKLIP 215
<a href="#">P24699</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111
<a href="#">Q91154</a>	63	-----Q---AGHCL--LWACKACKRKSSTTDRRK-AATMRERRRLKKNVNSAFETLKRCTT 111
<a href="#">P22091</a>	142	QPLASLGSGFFGEPDAPPMFTNNRVKRRSPYE----MEISDGPHT-KVVRRI-FTNSRERWRQNVNGAFELRKLIP 215
<a href="#">P57100</a>	81	-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 122
<a href="#">Q63689</a>	103	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 149
<a href="#">Q62414</a>	104	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 150
<a href="#">P70447</a>	94	-----RAVSRG----AKTAETVQRIKTRRL-KANNRERNRMHNLNAALDALREVLP 140
<a href="#">Q15784</a>	103	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 149
<a href="#">P17667</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFDTLKRCTT 111
<a href="#">Q55208</a>	54	-----L-HL-----VLERRR-VANAKERERIKNLNRGFARLKLALVP 87
<a href="#">P13349</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111

More details in Papadopoulos JS and Agarwala R, Bioinformatics 23:1073-79, 2007 (PMID: 17332019)

# Phylogenetic Tree View - based on COBALT multiple alignment

**COBALT** *Phylogenetic Tree View*

This tree is based on COBALT multiple alignment [more...](#)

## Phylo Tree View for 100 sequences: Cobalt RID ZV85B6BE212

Tree method: **Fast Minimum Evolution** | Max Seq Difference: **0.85** | Distance: **Grishin (protein)** | **Reset** | **Download** in **Newick Format**

**rectangle** | **slanted** | **radial** | **force** |  **Show distance** | **Mouse over an internal node for a subtree or alignment** | [Hide Color Map](#) | [Show removed sequences](#)

**Sequence Label**  
Sequence Title (if available)

**Collapse Mode** **Blast Name**

Blast names color map	
<span style="color: green;">■</span>	primates
<span style="color: blue;">■</span>	rodents
<span style="color: yellow;">■</span>	bony fishes
<span style="color: orange;">■</span>	birds
<span style="color: cyan;">■</span>	frogs & toads
<span style="color: magenta;">■</span>	even-toed ungulates
<span style="color: darkgreen;">■</span>	salamanders
<span style="color: grey;">■</span>	nematodes
<span style="color: pink;">■</span>	flies

Key sequence labels visible in the tree include:  
 - RecName: Full=Factor in the germline alpha; Sho...  
 - RecName: Full=Factor in the germline alpha; Short=FIGal...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 1 homolog; Short=TAL-1; AltName: Full=Stem cell p...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 1; Short=TAL-1; AltName: Full=Stem cell protein;...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 1 homolog; Short=TAL-1; AltName: Full=Stem cell...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 1; Short=TAL-1; AltName: Full=Stem cell leuke...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 2; Short=TAL-2; AltName: Full=Class A basic helix-loop-helix protein 19; Short...  
 - RecName: Full=T-cell acute lymphocytic leukemia protein 2 homolog; Short=TAL-2  
 - RecName: Full=Protein Iyf-1; AltName: Full=Lymphoblastic leukemia-derived sequence 1; AltName: Full=Class...  
 - RecName: Full=Heart and neural crest derivatives-expressed prote...  
 - RecName: Full=Heart and neural crest derivatives-expressed prote...  
 - RecName: Full=Heart and neural crest derivatives-expressed prot...  
 - RecName: Full=Basic helix-loop-helix transcription factor scleraxis  
 - RecName: Full=Basic helix-loop-helix transcription factor scleraxis; AltName: Full=Class A basic helix-loop...  
 - RecName: Full=Basic helix-loop-helix transcription factor scleraxis  
 - RecName: Full=Helix-loop-helix protein 2; Short=HEN-2; AltName: Full=Nescent helix loop helix 2; Short=NSCL-2; AltNam...  
 - RecName: Full=Helix-loop-helix protein 2; Short=HEN-2; AltName: Full=Nescent helix loop helix 2; Short=NSCL-2  
 - RecName: Full=Helix-loop-helix protein 1; Short=HEN-1; AltName: Full=Nescent helix loop helix 1; Short=NSCL-1; Al...  
 - RecName: Full=Helix-loop-helix protein 1; Short=HEN-1; AltName: Full=Nescent helix loop helix 1; Short=NSCL-1  
 - RecName: Full=Myogenic factor 5; Short=M...  
 - RecName: Full=Myogenic factor 5; Short=Myf...  
 - RecName: Full=Myogenic factor 5; Short=Myf...  
 - RecName: Full=Myogenic factor 5; Short=Shor...  
 - RecName: Full=Myogenic factor 5; Short=M...  
 - RecName: Full=Class A basic helix-loop-helix protein 15; Short=bH...  
 - RecName: Full=Neurogenic differentiation factor 1; Short=NeuroD1; Sh...  
 - RecName: Full=Neurogenic differentiation factor 1; Short=NeuroD1; Short...  
 - RecName: Full=Neurogenic differentiation factor 1; Short=NeuroD1  
 - RecName: Full=Neurogenic differentiation factor 1; Short=NeuroD1  
 - RecName: Full=Neurogenic differentiation factor 6-A; Short=Ne...  
 - RecName: Full=Neurogenic differentiation factor 6; Short=NeuroD6; A...  
 - RecName: Full=Neurogenic differentiation factor 6; Short=NeuroD6; Al...  
 - RecName: Full=Neurogenic differentiation factor 6; Short=NeuroD6  
 - RecName: Full=Neurogenin-2; Short=NGN-2; Al...  
 - RecName: Full=Neurogenin-1; Short=NGN-1; AltName: ...  
 - RecName: Full=Neurogenin-3; Short=NGN-3; AltNam...  
 - RecName: Full=Protein atonal homolog 7; AltName: Full=Helix-loop-helix protein zATH-5; Sho

# Nucleotide BLAST

▶ [NCBI/BLAST Home](#)

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

**New** Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#).

## BLAST Assembled Genomes

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

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

## Basic BLAST

Choose a BLAST program to run.

[nucleotide blast](#)

Search a **nucleotide** database using a **nucleotide** query  
*Algorithms: blastn, megablast, discontinuous megablast*

[protein blast](#)

Search **protein** database using a **protein** query  
*Algorithms: blastp, psi-blast, phi-blast*

[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

### News

#### [BLAST 2.2.23 release](#)

A new version of the stand-alone applications is available.

Mon, 22 Mar 2010 15:00:00 EST

[More BLAST news...](#)

### Tip of the Day

#### [How to do Batch BLAST jobs.](#)

BLAST makes it easy to examine a large group of potential gene candidates.

[More tips...](#)

# Where is homolog located in human?

► [NCBI/BLAST/blastn suite](#): BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

[Reset page](#) [Bookmark](#)

### Enter Query Sequence

Enter accession number, gi, or FASTA  [Clear](#)

```
>Crab eating macaque CDC20 mRNA
AGCGGAGAGTTTAAGAGGCGTAAGCGAGGCGTGTTAAACCCGGTCGGAAGTGC AACTTGCTC
ACGGGCTCCGCAGGCACCAACTGCAAGGACCCCTCCCGCTGCGGGCGTTCCCATGGCACAAT
GAGAGTGACCTGCACTCGCTGCTTCAGCTGGATGCACCCATCCCCAATGCACCCCTGCGCG
GCAAAGCCAAGGAAGCCTCAGGCCCGGCCCTCACCATGCGGGCCGCCAACCGATCCCAC
```

Or, upload file  [Browse...](#)

Job Title   
Enter a descriptive title for your BLAST search

Query subrange  
From   
To

### Choose Search Set

Database  Human genomic + transcript  Mouse genomic + transcript  Others (nr etc.):

Entrez Query **Optional**   
Enter an Entrez query to limit search

# Algorithm parameters: Nucleotide

**Algorithm parameters**

**General Parameters**

- Max target sequences: 100
- Short queries:  Automatically adjust word size for short input sequences
- Expect threshold: 10
- Word size: 11

**Scoring Parameters**

- Match/Mismatch Scores: 2,-3
- Gap Costs: Existence: 5 Extension: 2

**Filters and Masking**

- Filter:  Low complexity regions,  Species-specific repeats for Human
- Mask:  Mask for lookup table only,  Mask lower case letters

**Species-specific repeats dropdown:**

- Human
- Human
- Rodents
- Arabidopsis
- Rice
- Mammals
- Fungi
- C. elegans
- A. gambiae
- Zebrafish
- Fruit fly

**Callout 1 (General Parameters):**

- Prevents starting alignment in masked region
- Allows extensions through masked regions

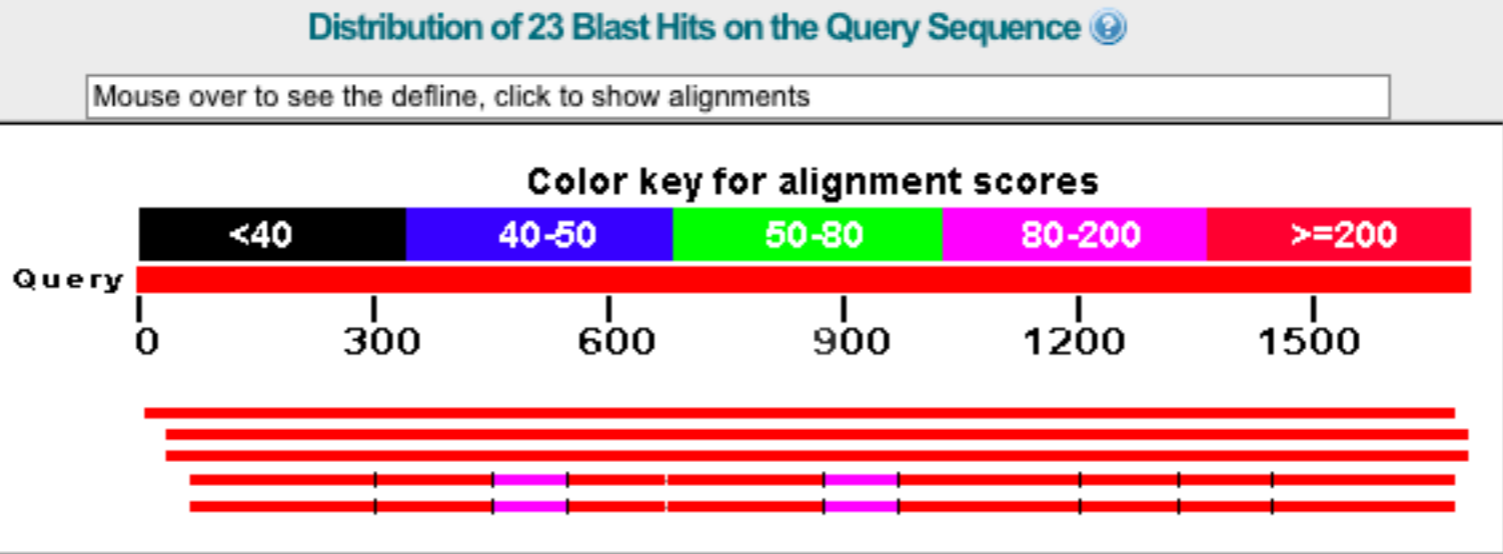
**Callout 2 (Species-specific repeats):**

- Masks LC sequence (simple repeats)
- Masks species-specific interspersed repeats
- Essential for genomic query sequences

**blastn**

# Sortable Results

Separate Sections for Transcript and Genome



## Descriptions

Legend for links to other resources: **U** UniGene **E** GEO **G** Gene **S** Structure **M** Map Viewer

### Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
<b>Transcripts</b>							
<a href="#">NM_001255.2</a>	Homo sapiens cell division cycle 20 homolog ( <i>S. cerevisiae</i> ) ( <i>CDC20</i> ),	<a href="#">2839</a>	2839	98%	0.0	97%	<b>GM</b>
<b>Genomic sequences</b> [ <a href="#">show first</a> ]							
<a href="#">NT_008470.19</a>	Homo sapiens chromosome 9 genomic contig, GRCh37 reference prim	<a href="#">2673</a>	2673	97%	0.0	95%	
<a href="#">NW_001839222.1</a>	Homo sapiens chromosome 9 genomic contig, alternate assembly (bas	<a href="#">2649</a>	2649	97%	0.0	95%	
<a href="#">NT_032977.9</a>	Homo sapiens chromosome 1 genomic contig, GRCh37 reference prim	<a href="#">411</a>	2853	94%	9e-112	100%	
<a href="#">NW_001838578.2</a>	Homo sapiens chromosome 1 genomic contig, alternate assembly (bas	<a href="#">411</a>	2853	94%	9e-112	100%	

Pseudogene on Chromosome 9

Functional Gene on Chromosome 1



# Total Score: All Segments

## Descriptions

Legend for links to other resources: **U** UniGene **E** GEO **G** Gene **S** Structure **M** Map Viewer

### Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
<b>Transcripts</b>							
<a href="#">NM_001255.2</a>	Homo sapiens cell division cycle 20 homolog (S. cerevisiae) (CDC20),	<a href="#">2839</a>	2839	98%	0.0	97%	<b>GM</b>
<b>Genomic sequences</b> [ <a href="#">show first</a> ]							
<a href="#">NW_001838578.2</a>	Homo sapiens chromosome 1 genomic contig, alternate assembly (bas	<a href="#">411</a>	2853	94%	9e-112	100%	
<a href="#">NT_032977.9</a>	Homo sapiens chromosome 1 genomic contig, GRCh37 reference prim	<a href="#">411</a>	2853	94%	9e-112	100%	
<a href="#">NT_008470.19</a>	Homo sapiens chromosome 9 genomic contig, GRCh37 reference prim	<a href="#">2673</a>	2673	97%	0.0	95%	
<a href="#">NW_001839222.1</a>	Homo sapiens chromosome 9 genomic contig, alternate assembly (bas	<a href="#">2649</a>	2649	97%	0.0	95%	

Functional Gene  
Now First

# Sorting in Exon Order

```
>  ref|NT\_032977.8|Hs1\_33153  Homo sapiens chromosome 1 genomic contig, reference assembly  
Length=73835825
```

Sort alignments for this subject sequence by:

[E value](#) [Score](#) [Percent identity](#)  
[Query start position](#) [Subject start position](#)

Features in this part of subject sequence:

[cell division cycle 20](#)

Score = 428 bits (216), Expect = 9e-117  
Identities = 231/236 (97%), Gaps = 0/236 (0%)  
Strand=Plus/Plus

```
Query  965      CTCCAGTGGTTCACGTTCTGGCCACATCCACCACCATGATGTTCGGGTAGCAGAACACCA  1024  
      |||  
Sbjct 13798316  CTCCAGTGGTTCACGTTCTGGCCACATCCACCACCATGATGTTCGGGTAGCAGAACACCA  13798375  
  
Query  1025     TGTGGCTACACTGAGTGGCCACAGCCAGGAAGTGTGTGGGCTGCGCTGGGCCCCAGATGG  1084  
      |||  
Sbjct 13798376  TGTGGCCACACTGAGTGGCCACAGCCAGGAAGTGTGTGGGCTGCGCTGGGCCCCAGATGG  13798435  
  
      TGTGGCCTAGCGCTCCTGG  1144  
      |||  
      TGTGGCCTAGTGCTCCTGG  13798495  
  
      AAGGGGCTGTCAAGG  1200  
      |||  
      AAGGGGCTGTCAAGG  13798551
```

**Default Sorting Order: Score**  
**Longest exon usually first**

# Sorting in Exon Order

```
>  ref|NT\_032977.8|Hs1\_33153  Homo sapiens chromosome 1 genomic contig, reference assembly
Length=73835825
```

Sort alignments for this subject sequence by:

E value Score Percent identity  
Query start position Subject start position

Features in this part of subject sequence:  
6169 bp at 5' side: myeloproliferative leukemia virus oncogene  
223 bp at 3' side: cell division cycle 20

Score = 42    Score = 89.7 bits (45),    Expect = 1e-14  
Identities = 51/53 (96%),    Gaps = 0/53 (0%)  
Strand=Plus/Plus

```
Query 965 Query 1      AGCGGAGAGTTTAAGAGGCGTAAGCGAGGCGTGTTAAACCCGGTTCGGAACTGC 53
          |||
Sbjct 13798 Sbjct 13796530 AGCGGAGAGTTTAAGAGGCGTAAGCCAGGCGTGTTAAAGCCGGTTCGGAACTGC 13796582
```

Query 1025  
Sbjct 13798  
Features in this part of subject sequence:  
cell division cycle 20

Score = 412 bits (208),    Expect = 5e-112  
Identities = 226/232 (97%),    Gaps = 0/232 (0%)  
Strand=Plus/Plus

```
Query 73      GGGCTCCGCAGGCACCAACTGCAAGGACCCCTCCCGCTGCGGGCGTTCCCATGGCACAAT 132
          |||
Sbjct 13796755 GGGCTCCGTAGGCACCAACTGCAAGGACCCCTCCCCCTGCGGGCGCTCCCATGGCACAGT 13796814

Query 133     TCGCGTTCGAGAGTGACCTGCACTCGCTGCTTCAGCTGGATGCACCCATCCCCAATGCAC 192
          |||
Sbjct 13796815 TCGCGTTCGAGAGTGACCTGCACTCGCTGCTTCAGCTGGATGCACCCATCCCCAATGCAC 13796874
```

Default  
Long

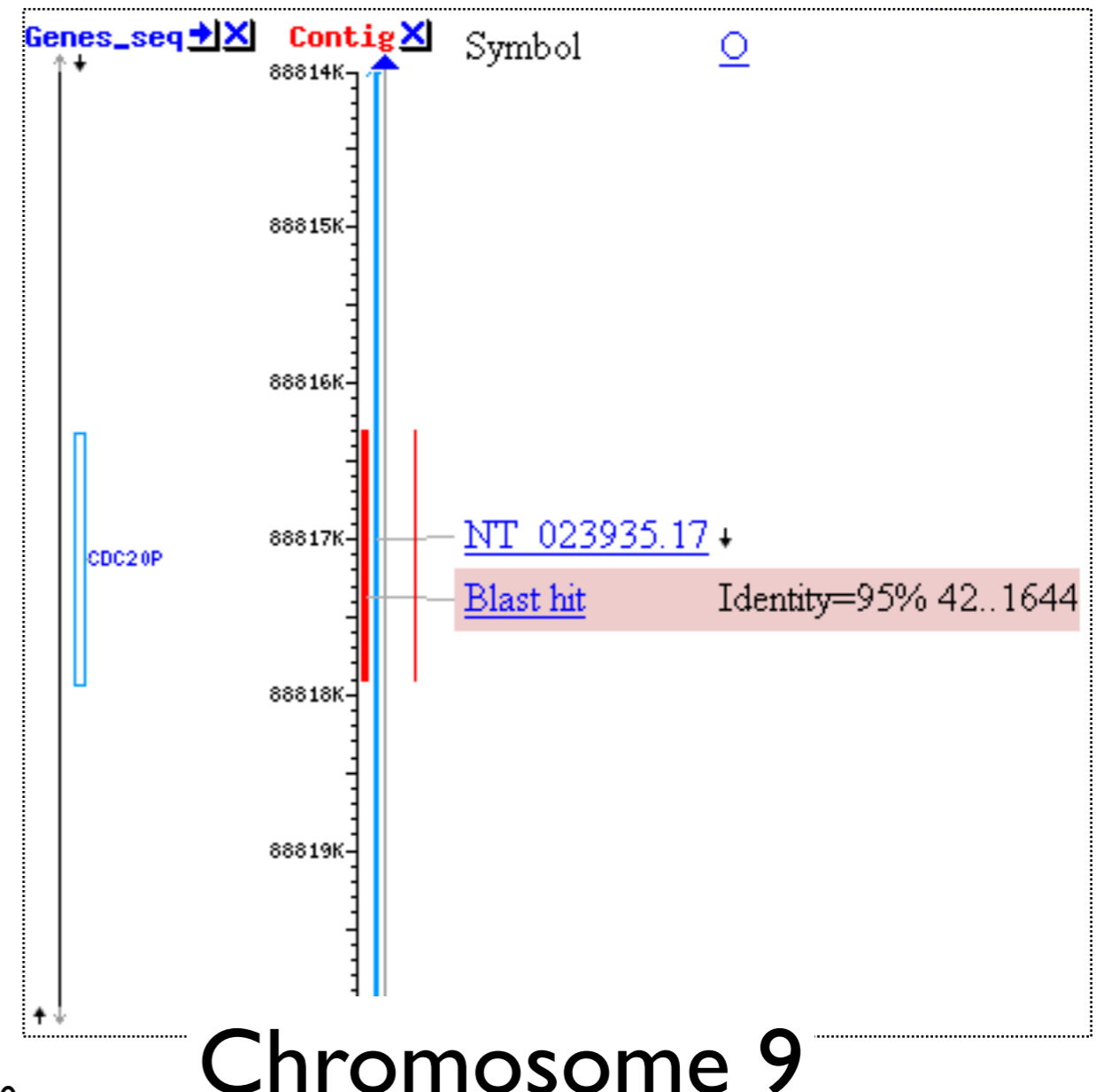
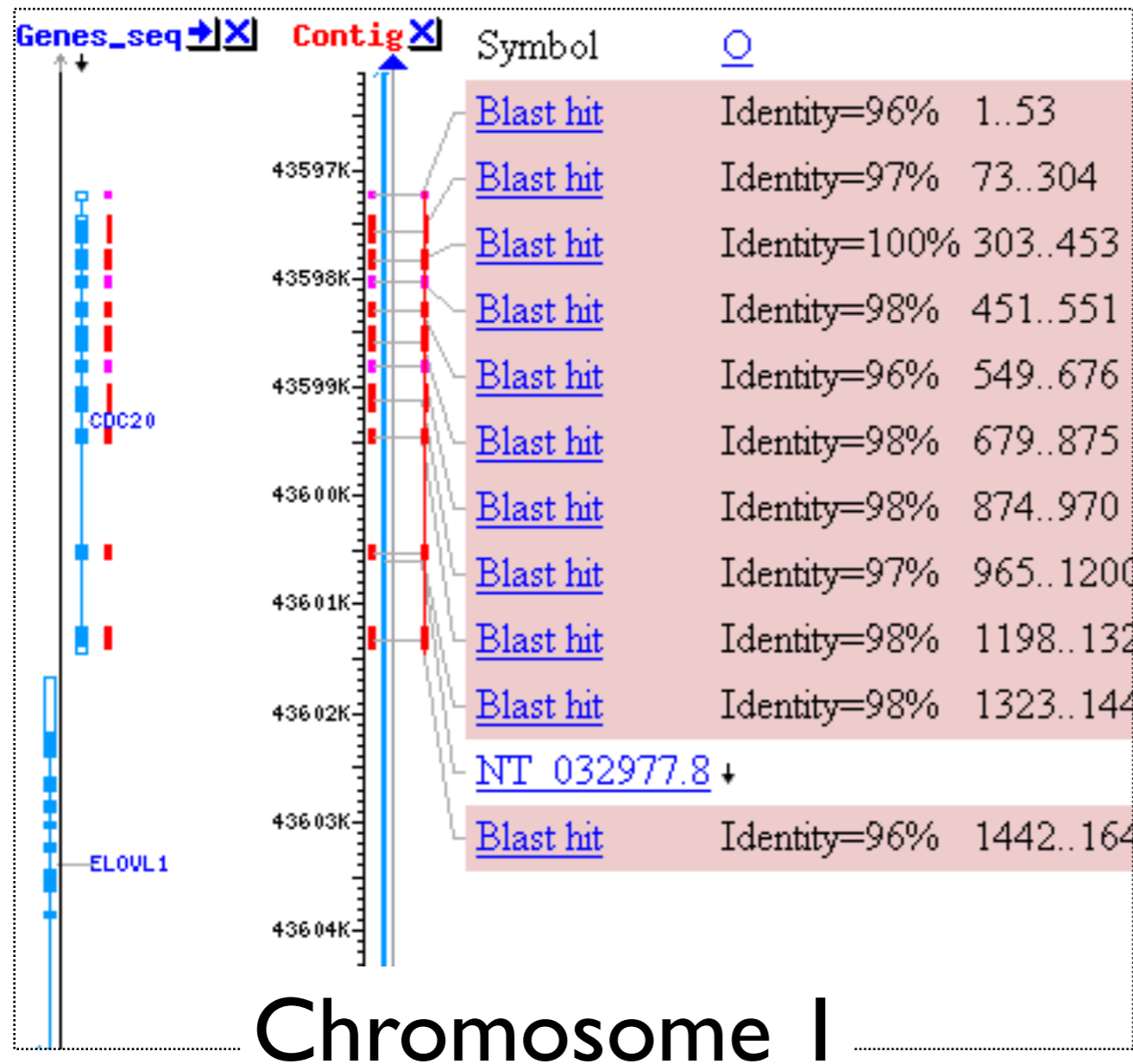
Query start  
position  
Exons in order

# Links to Genome View

**Query ID** [gi|67968779|dbj|AB168636.1|](#)  
**Description** Macaca fascicularis testis cDNA clone: QtsA-13692, similar to human CDC20 cell division cycle 20 homolog (S. cerevisiae) (CDC20), mRNA, RefSeq: NM\_001255.1  
**Molecule type** nucleic acid  
**Query Length** 1696

**Database Name** 3 databases  
**Description** [▶ See details](#)  
**Program** BLASTN 2.2.23+ [▶ Citation](#)

Other reports: [▶ Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#) [\[genome view\]](#)



# Recent and Saved Strategies

**BLAST** Basic Local Alignment Search Tool

Home Recent Results **Saved Strategies** Help

My NCBI  
Welcome joannealisonfox. [Sign Out]

► NCBI/ BLAST/ Recent Results  
Links to your unexpired BLAST jobs appear below. [more...](#)

Lookup BLAST Job

Request ID:  Go

Your Recent Results  
(Click headers to sort columns)

<a href="#">Submitted at</a>	<a href="#">Request ID</a>	<a href="#">Status</a>	<a href="#">Program</a>	<a href="#">Title</a>	<a href="#">Qlength</a>	<a href="#">Database</a>	<a href="#">Expires at</a>		
09-26 18:40	<a href="#">FNRZKDEZ012</a>	Done	blastp	Q02067:Achaete-scute homolog 1 (Mash-1)	231	swissprot	09-28 06:40	<a href="#">save</a>	✘
09-26 18:20	<a href="#">FNPT3VP9015</a>	Done	blastp	unknown protein - predict two seperate HSPs	169	nr	09-28 06:20	<a href="#">save</a>	✘
09-26 15:09	<a href="#">FNBKFCA3014</a>	Done	blastx	DinoDNA from THE LOST WORLD p. 135	1435	nr	09-28 03:09	<a href="#">save</a>	✘
09-26 14:57	<a href="#">FNAXJ9F4015</a>	Done	blastn	DinoDNA from JURASSIC PARK p. 103 nt 1-1200	1200	nr	09-28 02:57	<a href="#">save</a>	✘
09-26 12:43	<a href="#">FN31TZK015</a>	Done	megablast	dbj AB168636  (1696 letters)	1696	Human G+T	09-28 00:43	<a href="#">save</a>	✘

Login to My NCBI to save search strategies

# Specialized BLAST pages

## BLAST Assembled Genomes

---

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

- ▣ [Human](#)
- ▣ [Mouse](#)
- ▣ [Rat](#)
- ▣ [Arabidopsis thaliana](#)
- ▣ [Oryza sativa](#)
- ▣ [Bos taurus](#)
- ▣ [Danio rerio](#)
- ▣ [Drosophila melanogaster](#)
- ▣ [Gallus gallus](#)
- ▣ [Pan troglodytes](#)
- ▣ [Microbes](#)
- ▣ [Apis mellifera](#)

## Specialized BLAST

---

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

- ▣ Make specific primers with [Primer-BLAST](#)
- ▣ Search [trace archives](#)
- ▣ 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](#) (IgBLAST)
- ▣ Search for [SNPs](#) (snp)
- ▣ Screen sequence for [vector contamination](#) (vecscreen)
- ▣ [Align](#) two (or more) sequences using BLAST (bl2seq)
- ▣ Search [protein](#) or [nucleotide](#) targets in PubChem BioAssay
- ▣ Search SRA [transcript and genomic libraries](#)
- ▣ Constraint Based Protein [Multiple Alignment Tool](#)
- ▣ Needleman-Wunsch [Global Sequence Alignment Tool](#)

# Service Addresses

- ***General Help***      `info@ncbi.nlm.nih.gov`
- ***BLAST***              `blast-help@ncbi.nlm.nih.gov`

# BLAST

PRACTICAL EXERCISE: The Jurassic Park Detective Story





navigate to:  
 bioteach.ubc.ca/bioinfo2010

AMBL | The Educational Facilities of the Michael Smith Labs

# AMBL

**LABORATORY BIOINFORMATICS**

LABORATORY BIOINFORMATICS WORKSHOP, FEBRUARY 16-18TH, 2009

This workshop will focus on bioinformatics techniques for practical use in the laboratory. Hands-on exercises for retrieving data, primer design, BLAST searching, and genomics data navigation will be covered. Primarily aimed at researchers who are new to the area, or familiar but require a quick updating, where content covered can be tailored to laboratory needs.

joanne@msl.ubc.ca

**Laboratory Bioinformatics**  
 Common tools, useful databases, and tricks of the trade for practical use in the laboratory.

Written by AMBL  
 Edit

RESOURCES  
 UNIVERSITY

bioteach.ubc.ca/bioinfo2009

NCBI BLAST Home

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

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

### BLAST Assembled Genomes

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

- Human
- Mouse
- Rat
- Arabidopsis thaliana
- Oryza sativa
- Bos taurus
- Danio rerio
- Drosophila melanogaster
- Gallus gallus
- Pan troglodytes
- Microbes
- Apis mellifera

### Basic BLAST

Choose a BLAST program to run.

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

### Specialized BLAST

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

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

Let's compare our results



Get the sequences from the webpage and carry out BLAST searches

Can you identify the Dinosaur sequences?

Search #1:  
 Jurassic Park  
 sequence  
 use blastn

Search #2:  
 The Lost World  
 sequence  
 use blastx

Try some BLAST searches with  
your own sequence of interest...



Explore what happens when you  
change advanced parameters...

# Search #1 - blastn against nr



- Most common use of blastn
  - ✓ Sequence identification
  - ✓ Establish whether an exact match for a sequence is already present in the database

>|gi|157064989|gb|EU118176.1| Cloning vector pCM433, complete sequence  
Length=8081

Sort alignments for this subject sequence by:  
E value Score Percent identity  
Query start position Subject start position

Score = 437 bits (484), Expect = 4e-119  
Identities = 297/340 (87%), Gaps = 40/340 (11%)  
Strand=Plus/Plus

```
Query 1 GCGTTGCTGGCGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATCACAAAAATCGACGC 60
      |||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct 7309 GCGTTGCTGGCGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATCACAAAAATCGACGC 7368

Query 61 -----GGTGGCGAAACCCGACAGGACTATAAAGATAACCAGGCGTTTCCCCCTGGA 110
      |||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct 7369 TCAAGTCAGAGGTGGCGAAACCCGACAGGACTATAAAGATAACCAGGCGTTTCCCCCTGGA 7428

Query 111 AGCTCCCTCG-----TGTTCCGACCCTGCCGCTTACCGGATACCTGTCCGCCTTT 160
      ||||||||||| |||||||||||||||||||||||||||||||||||||||||||
Sbjct 7429 AGCTCCCTCGTGCCTCTCTCTGTTCCGACCCTGCCGCTTACCGGATACCTGTCCGCCTTT 7488

Query 161 CTCCCTTCGGGAAGCGTGGC-----TGCTCACGCTGTACCTATCTCAGTTCGGTG 210
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 7489 CTCCCTTCGGGAAGCGTGGCGCTTCTCATAGCTCACGCTGTAGGTATCTCAGTTCGGTG 7548

Query 211 TAGGTCGTTTCGCTCCAAGCTGGGCTGTGTG-----CCGTTTCAGCCCGACCGCTGC 260
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 7549 TAGGTCGTTTCGCTCCAAGCTGGGCTGTGTGCACGAACCCCGTTTCAGCCCGACCGCTGC 7608

Query 261 GCCTTATCCGGTAACTATCGTCTTGAGTCCAACCCGGTAA 300
      |||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct 7609 GCCTTATCCGGTAACTATCGTCTTGAGTCCAACCCGGTAA 7648
```

Score = 536 bits (594), Expect = 6e-149  
Identities = 360/410 (87%), Gaps = 50/410 (12%)  
Strand=Plus/Plus

```
Query 302 GTAGGACAGGTGCCGGCAGCGCTCTGGGTCATTTTCGGCGAGGACCGCTTTCGCTGGAG- 360
      |||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct 3591 GTAGGACAGGTGCCGGCAGCGCTCTGGGTCATTTTCGGCGAGGACCGCTTTCGCTGGAGC 3650

Query 361 -----ATCGGCCTGTGCTTGCAGGATTCGGAATCTTGACGCCCTCGCTCAAGCC 411
      |||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct 3651 GCGACGATGATCGGCCTGTGCTTGCAGGATTCGGAATCTTGACGCCCTCGCTCAAGCC 3710

Query 412 TTCGTCACT-----CCAAACGTTTCGGCGAGAAGCAGGCCATTATCGCCGGCATG 461
      ||||||||||| |||||||||||||||||||||||||||||||||||||||
Sbjct 3711 TTCGTCACTGGTCCCGCCACCAAACGTTTCGGCGAGAAGCAGGCCATTATCGCCGGCATG 3770

Query 462 GCGGCCGACGCGCTGGGCT-----GGCGTTCGCGACGCGAGGCTGGATGGCCTTC 511
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 3771 GCGGCCGACGCGCTGGGCTACGTCTTGCTGGCGTTCGCGACGCGAGGCTGGATGGCCTTC 3830

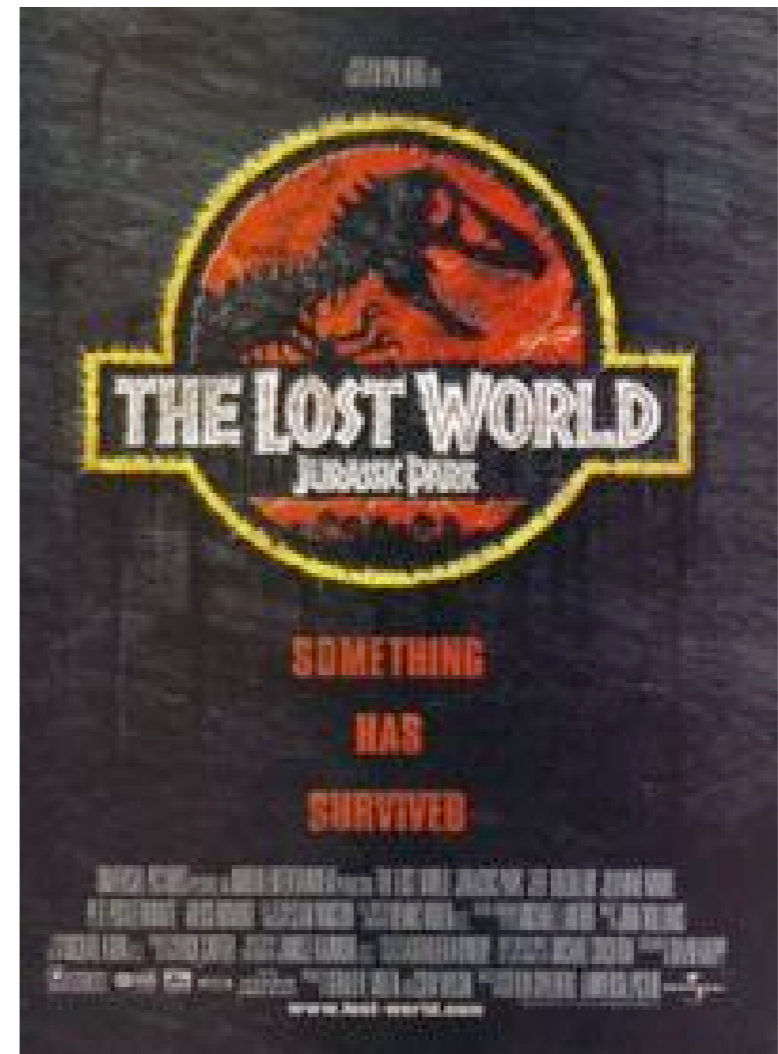
Query 512 CCCATTATGATTCTTCTCGCTTCCGGCG-----GCCCGCGTTGCAGGCCATGCTG 561
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 3831 CCCATTATGATTCTTCTCGCTTCCGGCGGCATCGGGATGCCCGCGTTGCAGGCCATGCTG 3890

Query 562 TCCAGGCAGGTAGATGACGACCATCAGGGACAGCTTCAA-----CGGCTCTTACC 611
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 3891 TCCAGGCAGGTAGATGACGACCATCAGGGACAGCTTCAAAGGATCGCTCGCGGCTCTTACC 3950

Query 612 AGCCTAACTTCGATCACTGGACCGCTGATCGTCACGGCGATTTATGCCGC 661
      ||||||||||||||||||||||| |||||||||||||||||||||||||||||||
Sbjct 3951 AGCCTAACTTCGATCACTGGACCGCTGATCGTCACGGCGATTTATGCCGC 4000
```

# Search #2 - blastx against nr

- Translating BLAST programs (blastx, tblastn, tblastx)
  - ✓ Look for similar proteins
  - ✓ Identify potential homologs in other species



```

> gi|45382623|ref|NP\_990795.1| UG erythroid-specific transcription factor eryfl [Gallus gallus]
gi|120955|sp|P17678|GATA1\_CHICK G Erythroid transcription factor (GATA-binding factor 1) (GATA-1)
(Eryfl) (NF-E1 DNA-binding protein) (NF-E1A)
gi|212629|gb|AAA49055.1| UG Eryfl protein
Length=304

Score = 366 bits (940), Expect = 2e-99
Identities = 304/318 (95%), Positives = 304/318 (95%), Gaps = 14/318 (4%)
Frame = +1

Query 121 MEFVALGGPDAGSPTFFPDeagafllgllggggerteaggl1aSYPPSGRVSLVPWADTGTLG 300
Sbjct 1 MEFVALGGPDAGSPTFFPDEAGAFLLGLGGGERTEAGGLLASYPSSGRVSLVPWADTGTLG 60

Query 301 TPQWVPPATQMEPPHYLEllqpprgsppphpsggpllp1ssgpppCEARECVMARKNCGAT 480
Sbjct 61 TPQWVPPATQMEPPHYLELLQPPRGSPPHPSGPLLPLSSGPPPCEARECV----NCGAT 116

Query 481 ATPLWRRDGTGHYLCNWASACGLYHRLNGQNRPLIRPKRLLVSKRAGTVCSHERENCQT 660
Sbjct 117 ATPLWRRDGTGHYLCN---ACGLYHRLNGQNRPLIRPKRLLVSKRAGTVCS----NCQT 169

Query 661 STTTLWRRSPMGDPVCNNIHACGLYYKLHQVNRPLTMRKDGIQTRNRKVSSKGGKRRPPG 840
Sbjct 170 STTTLWRRSPMGDPVCN---ACGLYYKLHQVNRPLTMRKDGIQTRNRKVSSKGGKRRPPG 226

Query 841 ggnpsatagggapmggggdpsmppppppppaaappQSDALYALGPVVLSGHFLPfgnsggf 1020
Sbjct 227 GGNPSATAGGGAPMGGGGDPSMPPPPPPPPAAAPPQSDALYALGPVVLSGHFLPFGNSGGF 286

Query 1021 fgggaggYTAPPGLSPQI 1074
Sbjct 287 FGGGAGGYTAPPGLSPQI 304

```

Mark was here, NIH

# BLAST

**COMMON TASKS - Basic Search; Searching Sets of Sequences (multiple inputs; small custom databases); Primer Design**



## A salmonid EST genomic study: genes, duplications, phylogeny and microarrays

Ben F Koop\*<sup>1,6</sup>, Kristian R von Schalburg<sup>1</sup>, Jong Leong<sup>1</sup>, Neil Walker<sup>1</sup>, Ryan Lieph<sup>1</sup>, Glenn A Cooper<sup>1</sup>, Adrienne Robb<sup>1</sup>, Marianne Beetz-Sargent<sup>1</sup>, Robert A Holt<sup>2</sup>, Richard Moore<sup>2</sup>, Sonal Brahmbhatt<sup>3</sup>, Jamie Rosner<sup>3</sup>, Caird E Rexroad III<sup>4</sup>, Colin R McGowan<sup>5</sup> and William S Davidson<sup>5</sup>

Address: <sup>1</sup>Centre for Biomedical Research, University of Victoria, Victoria, British Columbia, V8W 3N5, Canada, <sup>2</sup>Genome Sciences Centre, BC Cancer Agency, Vancouver, British Columbia, V5Z 4S6, Canada, <sup>3</sup>Prostate Centre, Vancouver, British Columbia, V6H 3Z6, Canada, <sup>4</sup>ARS, USDA, Natl Ctr Cool & Cold Water Aquaculture, Kearneysville, WV 25430, USA, <sup>5</sup>Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, British Columbia, V5A 1S6, Canada and <sup>6</sup>Department of Biology, University of Victoria, P.O. Box 3020, Victoria, British Columbia, V8W 3N5, Canada

Email: Ben F Koop\* - bkoop@uvic.ca; Kristian R von Schalburg - krvs@uvic.ca; Jong Leong - jong@uvic.ca; Neil Walker - nwalker@uvic.ca; Ryan Lieph - handsomryan@gmail.com; Glenn A Cooper - gac@uvic.ca; Adrienne Robb - arobb@uvic.ca; Marianne Beetz-Sargent - marianbs@uvic.ca; Robert A Holt - rholt@bcgsc.ca; Richard Moore - rmoore@bcgsc.ca; Sonal Brahmbhatt - Sonal.Brahmbhatt@vch.ca; Jamie Rosner - Jamie.Rosner@vch.ca; Caird E Rexroad - caird.rexroadIII@ARS.USDA.GOV; Colin R McGowan - cmcgowan@icywaters.com; William S Davidson - wdavidso@sfu.ca

\* Corresponding author

Published: 17 November 2008

Received: 13 June 2008

BMC Genomics 2008, 9:545 doi:10.1186/1471-2164-9-545

Accepted: 17 November 2008

This article is available from: <http://www.biomedcentral.com/1471-2164/9/545>

© 2008 Koop et al; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

### Abstract

**Background:** Salmonids are of interest because of their relatively recent genome duplication, and their extensive use in wild fisheries and aquaculture. A comprehensive gene list and a comparison of genes in some of the different species provide valuable genomic information for one of the most widely studied groups of fish.

**Results:** 298,304 expressed sequence tags (ESTs) from Atlantic salmon (69% of the total), 11,664 chinook, 10,813 sockeye, 10,051 brook trout, 10,975 grayling, 8,630 lake whitefish, and 3,624 northern pike ESTs were obtained in this study and have been deposited into the public databases. Contigs were built and putative full-length Atlantic salmon clones have been identified. A database containing ESTs, assemblies, consensus sequences, open reading frames, gene predictions and putative annotation is available. The overall similarity between Atlantic salmon ESTs and those of rainbow trout, chinook, sockeye, brook trout, grayling, lake whitefish, northern pike and rainbow smelt is 93.4, 94.2, 94.6, 94.4, 92.5, 91.7, 89.6, and 86.2% respectively. An analysis of 78 transcript sets show *Salmo* as a sister group to *Oncorhynchus* and *Salvelinus* within Salmoninae, and Thymallinae as a sister group to Salmoninae and Coregoninae within Salmonidae. Extensive gene duplication is consistent with a genome duplication in the common ancestor of salmonids. Using all of the available EST data, a new expanded salmonid cDNA microarray of 32,000 features was created. Cross-species hybridizations to this cDNA microarray indicate that this resource will be useful for studies of all 68 salmonid species.

**Conclusion:** An extensive collection and analysis of salmonid RNA putative transcripts indicate that Pacific salmon, Atlantic salmon and charr are 94–96% similar while the more distant whitefish, grayling, pike and smelt are 93, 92, 89 and 86% similar to salmon. The salmonid transcriptome reveals a complex history of gene duplication that is consistent with an ancestral salmonid genome duplication hypothesis. Genome resources, including a new 32 K microarray, provide valuable new tools to study salmonids.







NCBI/ BLAST/ blastx

blastn blastx **blastx** blastn tblastx

BLASTX search protein databases using a translated nucleotide query. [more...](#)

[Reset page](#) [Bookmark](#)

Enter Query Sequence

Enter accession number, gi, or FASTA sequence [?](#) [Clear](#)

```
GACACTCTTATTGCCATGACATTCAATTCTATAGTTGCCATTTTCTGTGTAGATTGATAATAAAATC
TTATATGCATTATGCAATCACGACTGTTGTTTACAGTGTACTCTGGAATTGTGTTATGCTCTCTCTTA
ATGGAATTATGTACCTTTCCATTCTATCTATACAAACCTTCAATAAACCTTTTCTGAACACAATTA
```

Query subrange [?](#)

From

To

Searching with Multiple Sequences as Input

Or, upload file

[Browse...](#) [?](#)

Genetic code

Standard (1) [?](#)

Job Title

8 sequences (gi|223585644|gb|GO065044.1|GO065044... [?](#)

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

Blast 2 sequences

Choose Search Set

Database

Reference proteins (refseq\_protein) [?](#)

Organism

Optional

Enter organism name or id--completions will be suggested

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

Entrez Query

Optional

Enter an Entrez query to limit search [?](#)

**BLAST**

Search **database refseq\_protein** using **Blastx** (search protein databases using a translated nucleotide query)

Show results in a new window

[Algorithm parameters](#)

Note: Parameter values that differ from the default are highlighted in yellow

Results for: 7:|cl|5773 gi|223585538|gb|GO065038.1|GO065038 EST\_ssal\_rgh\_1084502 ssalrgh mixed\_tissue full-length Salmo sal...(614bp)

Query ID Description

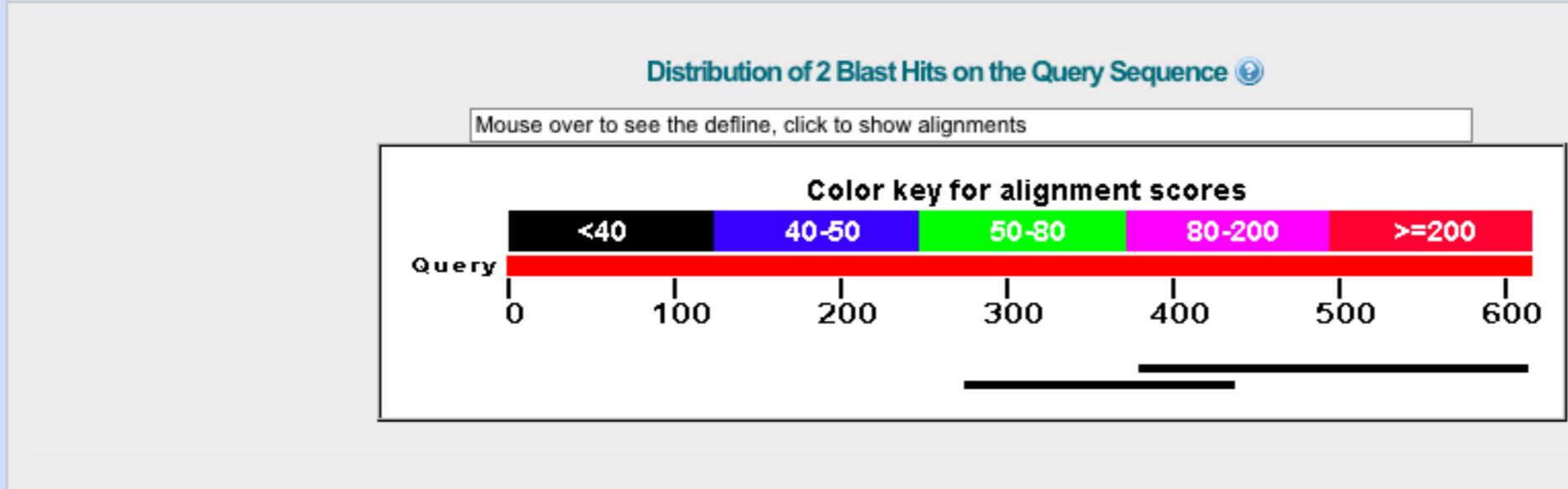
Molecule type Query Length

Other reports: [▶](#)

- 1:|cl|5767 gi|223585544|gb|GO065044.1|GO065044 EST\_ssal\_rgh\_1084509 ssalrgh mixed\_tissue full-length Salmo sal...(725bp)
- 2:|cl|5768 gi|223585543|gb|GO065043.1|GO065043 EST\_ssal\_rgh\_1079901 ssalrgh mixed\_tissue full-length Salmo sal...(897bp)
- 3:|cl|5769 gi|223585542|gb|GO065042.1|GO065042 EST\_ssal\_rgh\_1079900 ssalrgh mixed\_tissue full-length Salmo sal...(266bp)
- \*4:|cl|5770 gi|223585541|gb|GO065041.1|GO065041 EST\_ssal\_rgh\_1084506 ssalrgh mixed\_tissue full-length Salmo sal...(290bp)
- \*5:|cl|5771 gi|223585540|gb|GO065040.1|GO065040 EST\_ssal\_rgh\_1079898 ssalrgh mixed\_tissue full-length Salmo sal...(310bp)
- 6:|cl|5772 gi|223585539|gb|GO065039.1|GO065039 EST\_ssal\_rgh\_1084505 ssalrgh mixed\_tissue full-length Salmo sal...(432bp)
- 7:|cl|5773 gi|223585538|gb|GO065038.1|GO065038 EST\_ssal\_rgh\_1084502 ssalrgh mixed\_tissue full-length Salmo sal...(614bp)**
- 8:|cl|5774 gi|223585537|gb|GO065037.1|GO065037 EST\_ssal\_rgh\_1079894 ssalrgh mixed\_tissue full-length Salmo sal...(629bp)
- 9:|cl|5775 gi|223585536|gb|GO065036.1|GO065036 EST\_ssal\_rgh\_1079893 ssalrgh mixed\_tissue full-length Salmo sal...(884bp)
- 10:|cl|5776 gi|223585535|gb|GO065035.1|GO065035 EST\_ssal\_rgh\_1084500 ssalrgh mixed\_tissue full-length Salmo sal...(821bp)
- 11:|cl|5777 gi|223585534|gb|GO065034.1|GO065034 EST\_ssal\_rgh\_1079892 ssalrgh mixed\_tissue full-length Salmo sal...(791bp)

[What's this?](#)

▼ Graphic Summary



Results for:  
pull down list

► Descriptions

▼ Alignments  Select All [Get selected sequences](#)

```
> ref|YP\_934206.1 G hypothetical protein azo2703 [Azoarcus sp. BH72]
Length=774

GENE ID: 4607585\_azo2703 | hypothetical protein [Azoarcus sp. BH72]
(10 or fewer PubMed links)

Score = 35.0 bits (79), Expect = 4.3
Identities = 20/80 (25%), Positives = 36/80 (45%), Gaps = 2/80 (2%)
Frame = -2

Query 613 GEKPPQYPCNAAYSKL--DILILNGCQRHFKDIPAFYVNFVFCVHGEHETHWALTSIPR 440
      G++PP P + A + L D L+L +H+K A + + + G + W L P
Sbjct 557 GQRPPVTPLSRAEAGLPDDALVLAAFHQHYKITRASFALWMRLRLRGLPDALLWLLEGAPS 616

Query 439 WFKVISLK*HGNNIDPTSVC 380
      +S + + +DP +C
Sbjct 617 AMARLSQFARAHGVDPARLC 636
```



NCBI/ BLAST/ tblastn

blastn blastp blastx **tblastn** tblastx

Enter Query Sequence

TBLASTN search translated nucleotide subjects using a protein query. [more...](#)

Enter accession number, gi, or FASTA sequence

paste hbaa I sequence

>gi|47271417|ref|NP\_571332.2| hemoglobin alpha  
MSLSDTDKAVVKAIWAKISPKADEIGAEALARMLTVYPOTKTY  
E  
AVSKIDDLVGGLAALSELHAFKLRVDPANFKILSHNVIVVIAMLFPADFTPEVHVSVDKFFNNLALALS  
E  
KYR

To

Or, upload file

Job Title

gi|47271417|ref|NP\_571332.2| hemoglobin alpha...

Enter a descriptive title for your BLAST search

Align two or more sequences

Enter Subject Sequence

Enter accession number, gi, or FASTA sequence

>gi|223585544|gb|GO065044.1|GO065044 EST\_ssal\_rgh\_1084509\_ssalrgh  
mixed\_tissue\_full-length\_Salmo\_salar\_cDNA\_Salmo\_salar\_cDNA\_clone  
ssal\_rgh\_520\_381\_3', mRNA sequence  
AACTTGCAGCAAATACAAAAACAATAAATGATCAAACGAAACGTGACAAACAGTGACATGCCAAACAGG  
CAC  
CTACACAAAAACAAGATCCCACAAACCAGTGGGGAAATGGCTGCC

From

Or, upload file

Browse...

**BLAST**

Search nucleotide sequence using Tblastn (search translated nucleotide subjects using a protein query)

Show results in a new window

# Search against small custom database

Blast 2 sequences

gi|47271417|ref|NP\_571332.2| hemoglobin alpha...

**Query ID** lcl|20148  
**Description** gi|47271417|ref|NP\_571332.2| hemoglobin alpha adult-1 [Danio rerio]  
**Molecule type** amino acid  
**Query Length** 143

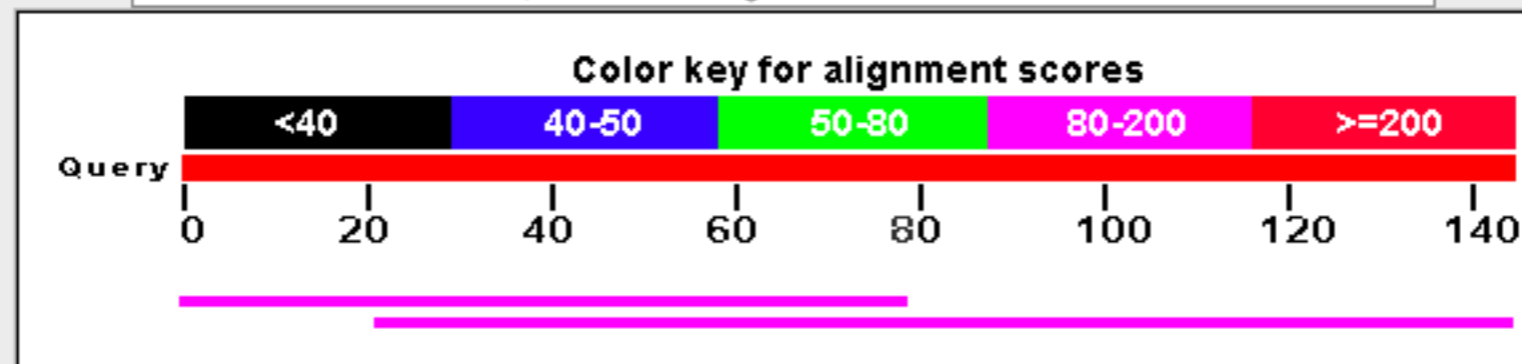
**Subject ID** 8 subjects  
**Description** [▶See details](#)  
**Molecule type** nucleic acid  
**Subject Length** n/a  
**Program** TBLASTN 2.2.19+ [▶Citation](#)

Other reports: [▶Search Summary](#) [\[Taxonomy reports\]](#)

## ▼ Graphic Summary

Distribution of 2 Blast Hits on the Query Sequence ⓘ

Mouse over to see the define, click to show alignments



## ▼ Descriptions

Sequences producing significant alignments:

					Score (Bits)	E Value
lcl 20152	gi 223585542	gb GO065042.1	GO065042	EST_ssal_rgh_10...	<u>116</u>	3e-31
lcl 20155	gi 223585539	gb GO065039.1	GO065039	EST_ssal_rgh_10...	<u>178</u>	4e-50

# BLAST tasks

## Basic BLAST

- ✓ Jurassic Park examples

## Batch BLAST searching

- ✓ Use Salmon ESTs as input

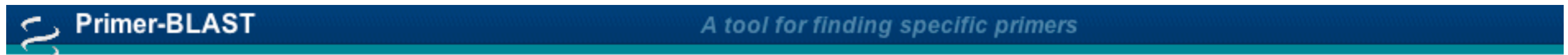
## Search against a small custom database

- ✓ Use “Align two or more sequences”

# Primer-BLAST

NCBI's Primer Designer and Specificity Checker

<http://www.ncbi.nlm.nih.gov/tools/primer-blast/>



▶ NCBI/Primer-BLAST: Finding primers specific to your PCR template (using Primer3 and BLAST). [more...](#) [Tips for finding specific primers](#)

**PCR Template** [Reset page](#) [Save search parameters](#)

Enter accession, gi, or FASTA sequence (A refseq record is preferred) [?](#) [Clear](#)

Or, upload FASTA file  no file selected

Range

Forward primer  From  To  [?](#) [Clear](#)

Reverse primer

**Primer Parameters**

Use my own forward primer (5'->3' on plus strand)  [?](#) [Clear](#)

Use my own reverse primer (5'->3' on minus strand)

PCR product size  Min  Max

# of primers to return

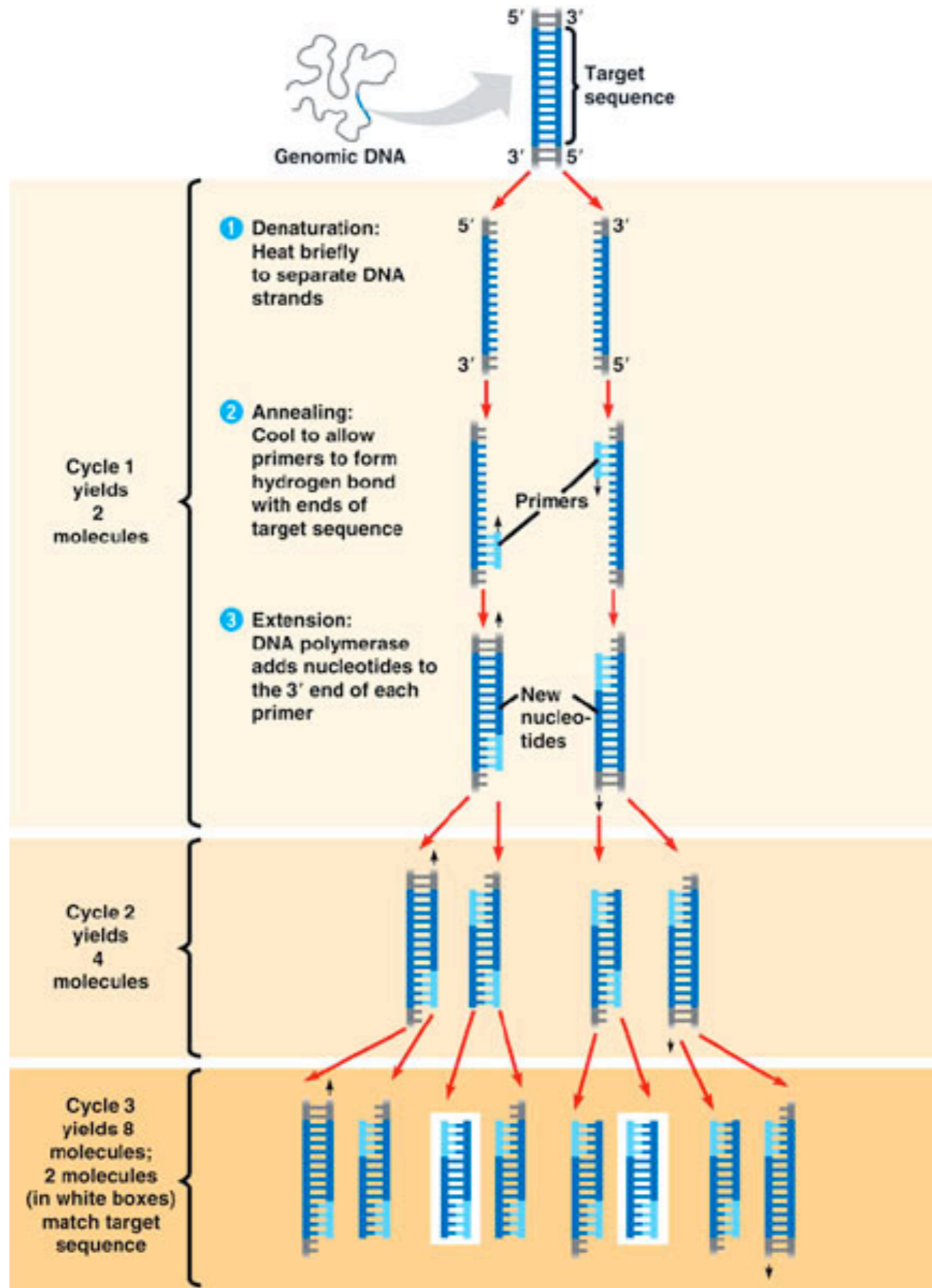
Primer melting temperatures (T<sub>m</sub>)  Min  Opt

**Primer Pair Specificity Checking Parameters**

Specificity check  Enable search for primer pairs specific to the intended PCR template [?](#)

Organism

offers integrated primer design with Primer3 & specificity check with custom BLAST





# Primer Design

## Balance:

- ✓ Specificity - frequency of mispriming
- ✓ Efficiency of Amplification - 2X increase

## Consider:

- primer length (18-24nt)
- primer  $T_m$  ( $>54^\circ\text{C}$ )
- 3' end (G or C)
- GC content (45-55%)
- primer dimers
- for cDNA - coding region; across intron/exon boundary

General Concepts for PCR Primer Design.  
Dieffenback CW, Lowe TM, Dveksler GS Genome Research  
3 (1993) S30-37 [PMID:8118394]

# Primer-BLAST input

designs primers specific to target template and unique in the target database

▶ [NCBI/ Primer-BLAST: Finding primers](#)

The screenshot displays the NCBI Primer-BLAST web interface. It is divided into two main sections: "PCR Template" and "Primer Parameters".

**PCR Template Section:**

- Label: "PCR Template" (highlighted with a blue box and an arrow pointing to the text "designs primers specific to target template and unique in the target database").
- Input field: "Enter accession, gi, or FASTA sequence (A refseq record is preferred)" with a "Clear" button.
- Range selection: "Range" section with "From" and "To" input boxes for "Forward primer" and "Reverse primer", and a "Clear" button.
- File upload: "Or, upload FASTA file" with a "Browse..." button.

**Primer Parameters Section:**

- Label: "Primer Parameters" (highlighted with a blue box and an arrow pointing to the text "can specify primer sequence(s), desired product size, T<sub>m</sub> ranges, T<sub>m</sub> difference (can be used with or without template)").
- Custom primers: "Use my own forward primer (5'→3' on plus strand)" and "Use my own reverse primer (5'→3' on minus strand)", each with an input field and a "Clear" button.
- Product size: "PCR product size" with "Min" (200) and "Max" (1000) input boxes.
- Number of primers: "# of primers to return" with an input box containing "10".
- Melting temperatures: "Primer melting temperatures (T<sub>m</sub>)" with "Min" (57.0), "Opt" (60.0), "Max" (63.0), and "Max T<sub>m</sub> difference" (3) input boxes.

can specify primer sequence(s), desired product size, T<sub>m</sub> ranges, T<sub>m</sub> difference (can be used with or without template)

# Primer-BLAST Specificity

By default human sequences are searched in specificity check

**Primer Pair Specificity Checking Parameters**

**Specificity check**

Enable search for primer pairs specific to the intended PCR template [?](#)

With this option on, the program will search the primers against the selected database and determine whether a primer pair can generate a PCR product on any targets in the database based on their matches to the targets and their orientations. The program will return, if possible, only primer pairs that do not generate a valid PCR product on unintended sequences and are therefore specific to the intended template. Note that the specificity is checked not only for the forward-reverse primer pair, but also for forward-forward as well as reverse-reverse primer pairs.

**Organism**

Homo sapiens

Enter an organism name, taxonomy id or select from the suggestion list as you type. [?](#)

**Database**

Refseq mRNA (refseq\_rna) [?](#)

**Primer specificity stringency**

At least  total mismatches to unintended targets, including  
at least  mismatches within the last  bps at the 3' end [?](#)

The larger the mismatches (especially those toward 3' end) are between primers and the unintended targets, the more specific the primer pair is to your template (i.e., it will be difficult to anneal to and amplify unintended targets). However, specifying a larger mismatch value may make it more difficult to find such specific primers. Try to lower the mismatch value in such case.

**Misprimed product size deviation**

[?](#)

**Splice variant handling**

Allow primer to amplify mRNA splice varia

[Get Primers](#)

Show results in a new window

custom BLAST; focus on 3' end to avoid mispriming

# Primer-BLAST Specificity

Four BLAST nucleotide databases available for searching

1. with refseq template, specific to splice variant

2. human, chimp, mouse, rat, cow, dog, chicken, zebrafish, fly, bee, Arabidopsis, rice

4. nr = database with widest coverage of organisms

3. all NC\_ (includes above) + other organisms, microbes

# Primer-BLAST Advanced

Adjustable settings from Primer3  
see Primer 3 Input Help:

<http://fokker.wi.mit.edu/primer3/input-help-040.htm>

▼ **Advanced parameters**

**Primer Pair Specificity Checking Parameters**

Blast max number of hit sequences: 250 (default)

Blast expect (E) value: 1000 (default)

Max primer pairs to screen: 3000 (default)

**Primer Parameters**

	Min	Opt	Max
PCR Product Tm			
Primer Size	15	20	27
Primer GC content (%)	20.0	80.0	
GC clamp	0		
Max self complementarity:	8.00		
Max 3' end complementarity:	3.00		
SNP handling	<input type="checkbox"/> Primer binding site may not contain known SNP		
Repeat filter	Automatic		
Low complexity filter	<input checked="" type="checkbox"/> Avoid low complexity region for primer selection		
Concentration of monovalent cations	50.0		
Concentration of divalent cations	0.0		
Concentration of dNTPs	0.0		
Salt correction formula:	Schildkraut and Lifson 1965		
Annealing Oligo Concentration	50.0		

Useful options specific to Primer-BLAST:

1. avoid regions that contain SNPs
2. avoid repetitive regions

# Primer-BLAST example



**Task #1:** Use Primer BLAST to design primers specific to the UNG1 splice variant, NM\_003362.

**Task #2:** Use Primer BLAST to design primers specific to the UNG2 splice variant, NM\_080911.

**Task #3:** Carry out a specificity check for one of your primer pairs. Will this primer pair (designed against the human UNG2 transcript, for example) also amplify transcripts from other primate species?

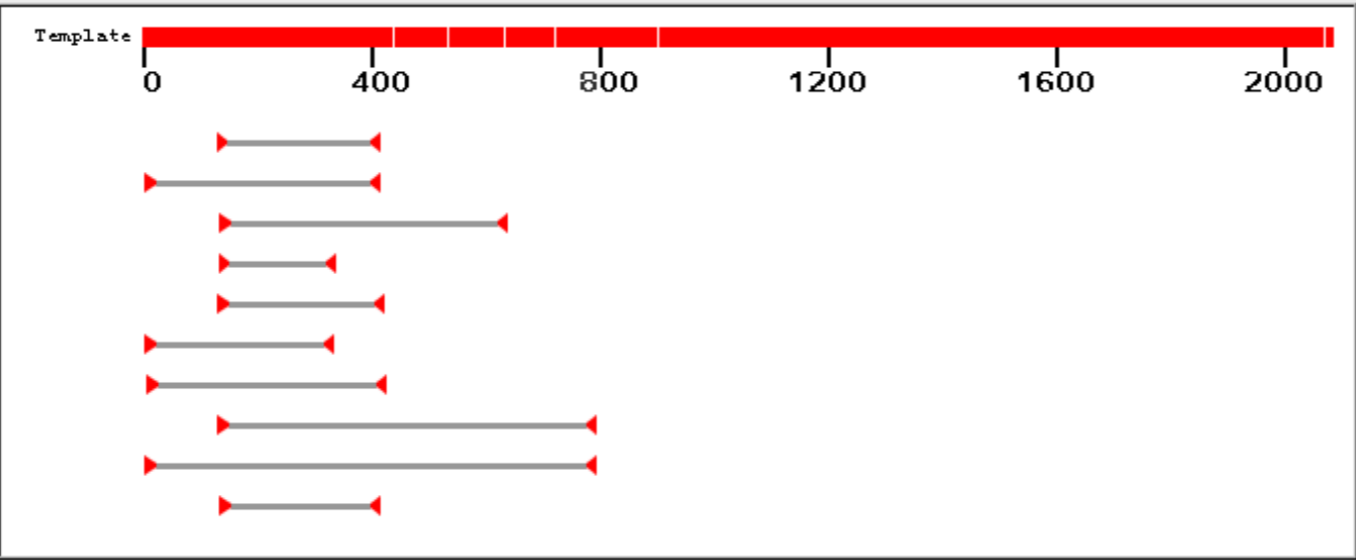
# Task #1: Use Primer BLAST to design primers that will identify the UNG1 splice variant.

**Input PCR template** [NM\\_003362.2](#) Homo sapiens uracil-DNA glycosylase (UNG), nuclear gene encoding mitochondrial protein variant 1, mRNA

**Range** 1 - 2081

**Specificity of primers** Primer pairs are specific to input template as no other targets were found in selected database: NCBI Tr Sequences (Organism limited to Homo sapiens)

▼ **Summary of primer pairs**



Note: a break in the template graph indicates the exon-exon junction

▼ **Detailed primer reports**

**Primer pair 1**

	Sequence (5'->3')	Strand on template	Length	Start	Stop	Tm	GC%
Forward primer	CTTCTGCCTTGGGCCGTGGG	Plus	20	134	153	59.97	70.00%
Reverse primer	<b>TCCCGAACTCCCCGCTGAGG</b>	Minus	20	420	401	59.97	70.00%
Product length	287						

**Products on intended target**

>[NM\\_003362.2](#) Homo sapiens uracil-DNA glycosylase (UNG), nuclear gene encoding mitochondrial protein, transcript variant 1, mRNA

enter NM\_003362  
as template

use all default  
settings

# Task #2: Use Primer BLAST to design primers specific to the UNG2 splice variant, NM\_080911.

enter NM\_080911 as template

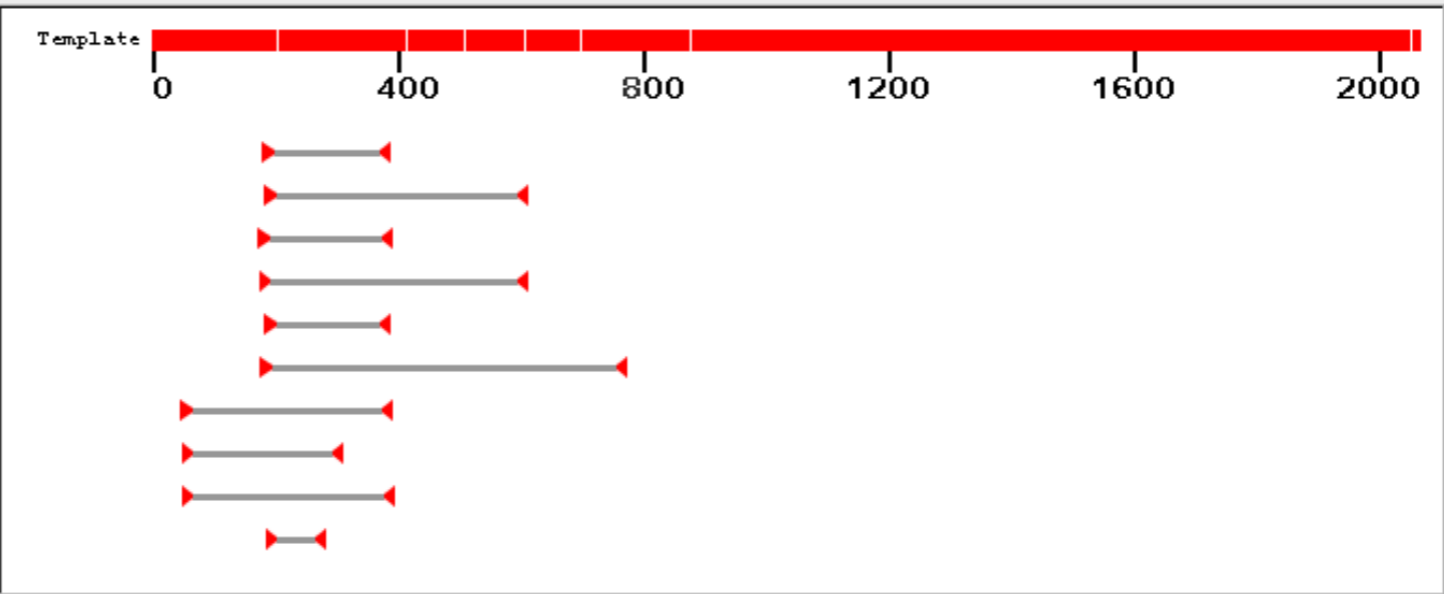
use all default settings

**Input PCR template** [NM\\_080911.1](#) Homo sapiens uracil-DNA glycosylase (UNG), transcript variant 2, mRNA

**Range** 1 - 2053

**Primer pairs** Primer pairs are specific to input template as no other targets were found in selected database: NCBI Transcript Reference Sequences (Organism limited to Homo sapiens)

**Number of primer pairs**



Note: a break in the template graph indicates the exon-exon junction

**▼ Detailed primer reports**

**Primer pair 1**

	Sequence (5'->3')	Strand on template	Length	Start	Stop	Tm	GC%
Forward primer	<a href="#">AGGAAAGCGGAGATGCGGCG</a>	Plus	20	183	202	59.84	65.00%
Reverse primer	TCCCGAACTCCCCGCTGAGG	Minus	20	392	373	59.97	70.00%
Product length	210						

**Products on intended target**

>[NM\\_080911.1](#) Homo sapiens uracil-DNA glycosylase (UNG), transcript variant 2, mRNA



**Task #3:** Carry out a specificity check for one of your primer pairs from Task #2. Will this primer pair (designed against the human UNG2 transcript) also amplify transcripts from other primate species?

no template

use my own:

forward primer

reverse primer

organism;  
specify primate

database;  
specify nr

**Primer pair 1**

	Sequence (5'->3')	Length	Tm	GC%
Forward primer	AGGAAAGCGGAGATGCGGCG	20	59.84	65.00%
Reverse primer	TCCCGAACTCCCCGCTGAGG	20	59.97	70.00%

**Products on target templates**

>[XM\\_002752978.1](#) PREDICTED: Callithrix jacchus uracil-DNA glycosylase-like (LOC100393193), mRNA

```
product length = 210
Forward primer 1  AGGAAAGCGGAGATGCGGCG  20
Template       113  .A..G.....  132

Reverse primer 1  TCCCGAACTCCCCGCTGAGG  20
Template       322  .....  303
```

>[AK313552.1](#) Homo sapiens cDNA, FLJ94113, Homo sapiens uracil-DNA glycosylase (UNG), nuclear gene encoding mitochondrial protein, transcript variant 2, mRNA

```
product length = 210
Forward primer 1  AGGAAAGCGGAGATGCGGCG  20
Template       196  .....  215

Reverse primer 1  TCCCGAACTCCCCGCTGAGG  20
Template       405  .....  386
```

>[NG\\_007284.1](#) Homo sapiens uracil-DNA glycosylase (UNG), RefSeqGene on chromosome 12

```
product length = 830
Forward primer 1  AGGAAAGCGGAGATGCGGCG  20
Template       5183  .....  5202

Reverse primer 1  TCCCGAACTCCCCGCTGAGG  20
Template       6012  .....  5993
```

>[AC193909.4](#) Pan troglodytes BAC clone CH251-19J18 from chromosome 12, complete sequence

```
product length = 831
Forward primer 1  AGGAAAGCGGAGATGCGGCG  20
Template       148574  .....  148555

Reverse primer 1  TCCCGAACTCCCCGCTGAGG  20
Template       147744  .....  147763
```

>[XM\\_509349.2](#) PREDICTED: Pan troglodytes uracil-DNA glycosylase, transcript variant 2 (UNG), mRNA

```
product length = 210
Forward primer 1  AGGAAAGCGGAGATGCGGCG  20
```

# Things you can do to maximize the chance of finding primers specific for your template.

- **Use refseq accession or GI (rather than the raw DNA sequence) as template whenever possible.** Even if you are only interested in part of the sequence, you can still use the accession or GI but you do need to specify the range (use forward primer "From" field for your sequence start position and reverse primer "To" field for your sequence stop position). The reason is that an accession or GI carries accurate information about its identity which allows primer-blast to better distinguish between intended template and off-targets.
- **Choose a non-redundant database (such as refseq\_rna or genome database).** The nr database contains redundant entries which can interfere with the process of finding specific primers.
- **Specify an organism** for database search if you are only amplifying DNA from a specific organism. Searching all organisms will be much slower and off-target priming from other organisms are irrelevant.

# Credits

- Materials for this presentation have been adapted with permission from the following NCBI HelpDesk course materials:

Field Guide Course Materials

Advanced Workshop for Bioinformatics Information Specialists

NCBI News

- NCBI BLAST

<http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>



# MSA

MSA = Multiple Sequence Alignments



# Examples

```

globin.aln
CLUSTAL 2.0.9 multiple sequence alignment

HBB_HUMAN      -----VHLTPEEKSAVTALWGKVNV--VDEVGGEALGRLLVYYPWTQRRFFESFGDLST
HBB_HORSE      -----VQLSGEEKAAVLALWVKVN--EEEVGGEALGRLLVYYPWTQRRFFDSFGDLSN
HBA_HUMAN      -----VLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFSLFPTTKTYFPHF--DLS-
HBA_HORSE      -----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF--DLS-
GLB5_PETMA     PIVDTGSAVPLSAAEKTIRSAAWAPVYSTYETSGVDILVKFFTSTPAAQEFFPKFKGLTT
MYG_PHYCA      -----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKDFRFKHLKT
LGB2_LUPLU     -----GALTESQAALVKSSWEEFNANIPKHTRHFFILVLEIAPAAKDLFSFLKGTSE
                *:  :  :  *  .      :  .:  *  :  *  :  .

HBB_HUMAN      PDAVMGNPKVKAHGKKVLAHFGKFTPPVQAAYQKVVAGVANALAHKYH-----
HBB_HORSE      PGAVMGNPKVKAHGKKVLSHFGGEGVHHLDN-----LKGTFEALSELHCDKLHVDPENFR
HBA_HUMAN      ----HGSAQVKGHGKKVADALTNVAHVDD-----MPNALSALSDLHAHKLRVDPVNFKL
HBA_HORSE      ----HGSAQVKAHGKKVGDALTLAVGHLD-----LPGALSNDLSDLHAHKLRVDPVNFKL
GLB5_PETMA     ADQLKKSADVRWHAERIINAVNDAVASMDDT--EKMSMKLRDLGKHAHSFQVDPQYFKV
MYG_PHYCA      EAEMKASEDLKKGVTVLTALGAILKKKGH-----HEAELKPLAQSHATKHKIPKYLEF
LGB2_LUPLU     VP--QNNPELQAHAGKVFKLVEEAAIQLVQVTVVVDATLKNLGSVHVSKG--VADAHFPV
                . .:  *  :  .      :  *  *  .  :  :  .

HBB_HUMAN      LGNVLVCVLAHFGKFTPPVQAAYQKVVAGVANALAHKYH-----
HBB_HORSE      LGNVLVVVLAHFGKDFPELQASYQKVVAGVANALAHKYH-----
HBA_HUMAN      LSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTISKYR-----
HBA_HORSE      LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTISKYR-----
GLB5_PETMA     LAAVIADTVAAG-----DAGFEKLSMICILLRSAY-----
MYG_PHYCA      ISEAIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
LGB2_LUPLU     VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA--
                :  :  .:  ...      .  :
    
```

ClustalX 2.0.9

Mode:  Font:

The screenshot shows the ClustalX 2.0.9 interface. On the left, a list of sequences is shown: HBB\_HUMAN, HBB\_HORSE, HBA\_HUMAN, HBA\_HORSE, GLB5\_PETMA, MYG\_PHYCA, and LGB2\_LUPLU. The main window displays a multiple sequence alignment of these sequences, with each column representing a position in the protein. The alignment is color-coded by amino acid type. Below the alignment, a horizontal bar indicates the conservation of each position, with a scale from 1 to 110. The bar shows higher conservation in certain regions, particularly around positions 10, 20, 50, 80, and 100.

# Multiple Sequence Alignment

```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCTVSGTSEDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPSD--IAVEWESNG--
```

The sole purpose of multiple sequence alignments is to place *homologous positions of homologous sequences* into the *same column.*

# Differences between MSA and BLAST?

## MSA

- global alignment method
  - Align complete sequence
- Assumes homology
- Complex gap penalties
- Slower
- Align protein-protein or nucleotide-nucleotide only

## BLAST

- local alignment method
  - Search for HSP
- Test for homology
- Simple gap penalties
- Fast
- Translated searches



# Clustal

- Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994)
- CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice.
- Nucleic Acids Research, 22:4673-4680.

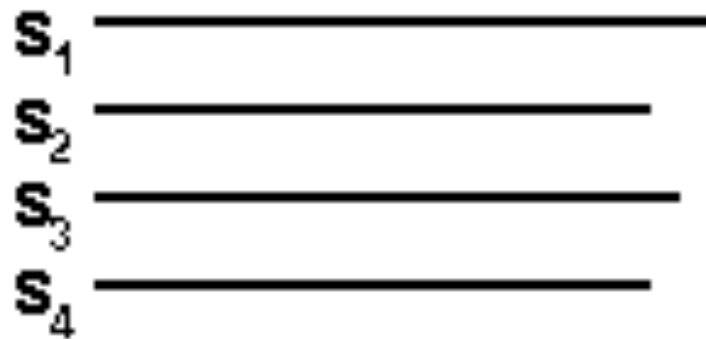
# CLUSTAL Algorithm Steps

1. Pairwise alignment of each sequence pair
  - Number of comparisons depends on how many sequences
2. Compute distance matrix
  - Percent non-identity between each alignment pair
  - Lower distance means more similar
3. Construct a sequence similarity tree
  - Cluster sequences according to distance (similarity)
4. Progressive alignment of sequences according to a tree

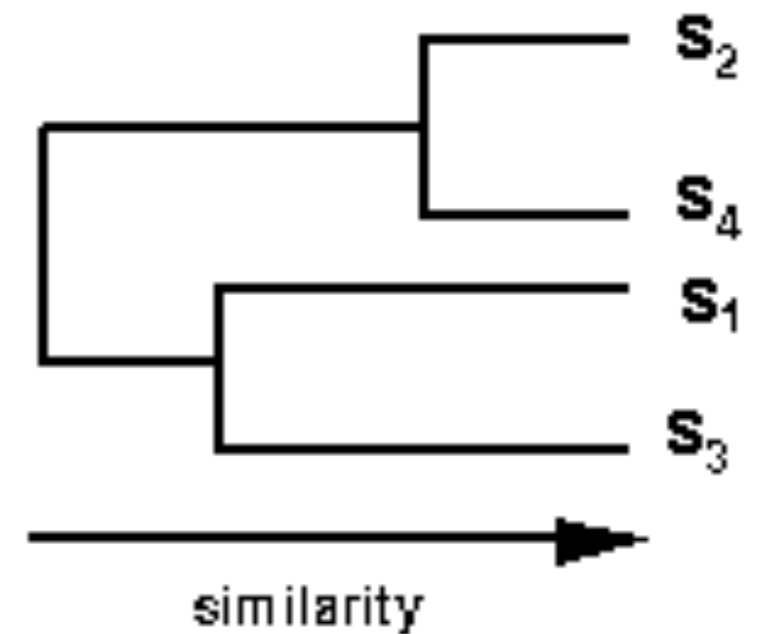
# How does the Clustal algorithm actually work?

## (A) Pairwise Alignment

Example - 4 sequences  $s_1$   $s_2$   $s_3$   $s_4$



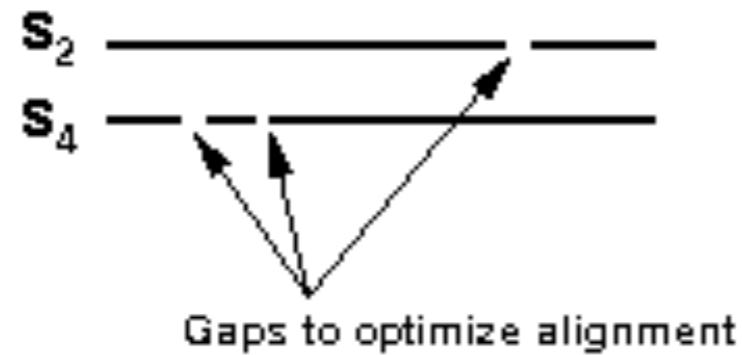
6 pairwise comparisons  
then cluster analysis



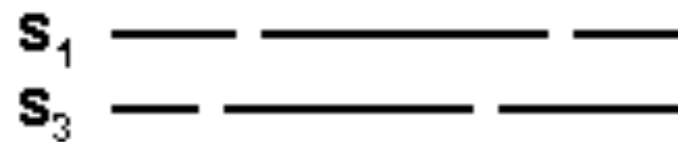
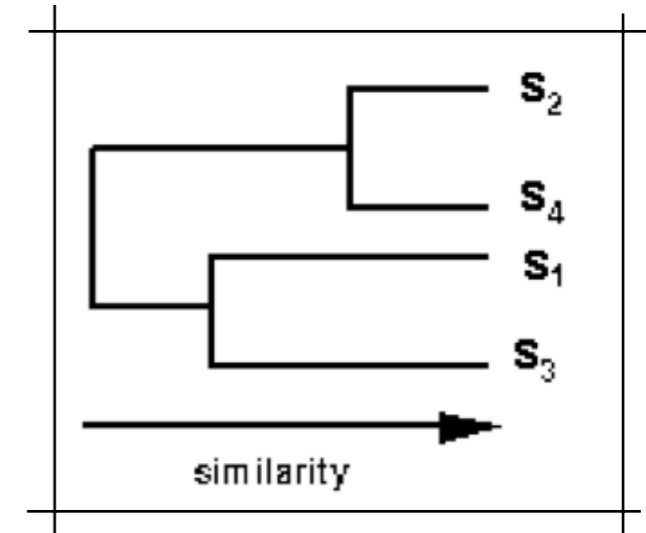
Which sequences would be aligned first?

# Steps in a Multiple Sequence Alignment continued ...

## (B) Multiple alignment following the tree from A

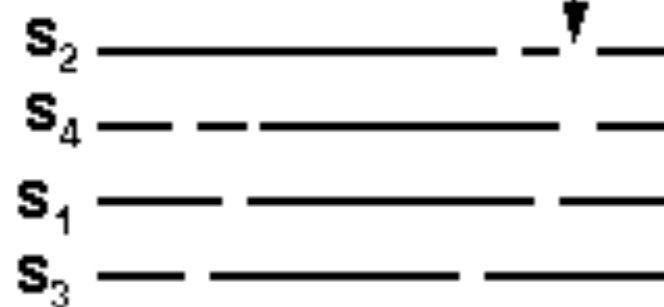


align most similar pair



align next most similar pair

New gap to optimize alignment of  $(s_2, s_4)$  with  $(s_1, s_3)$

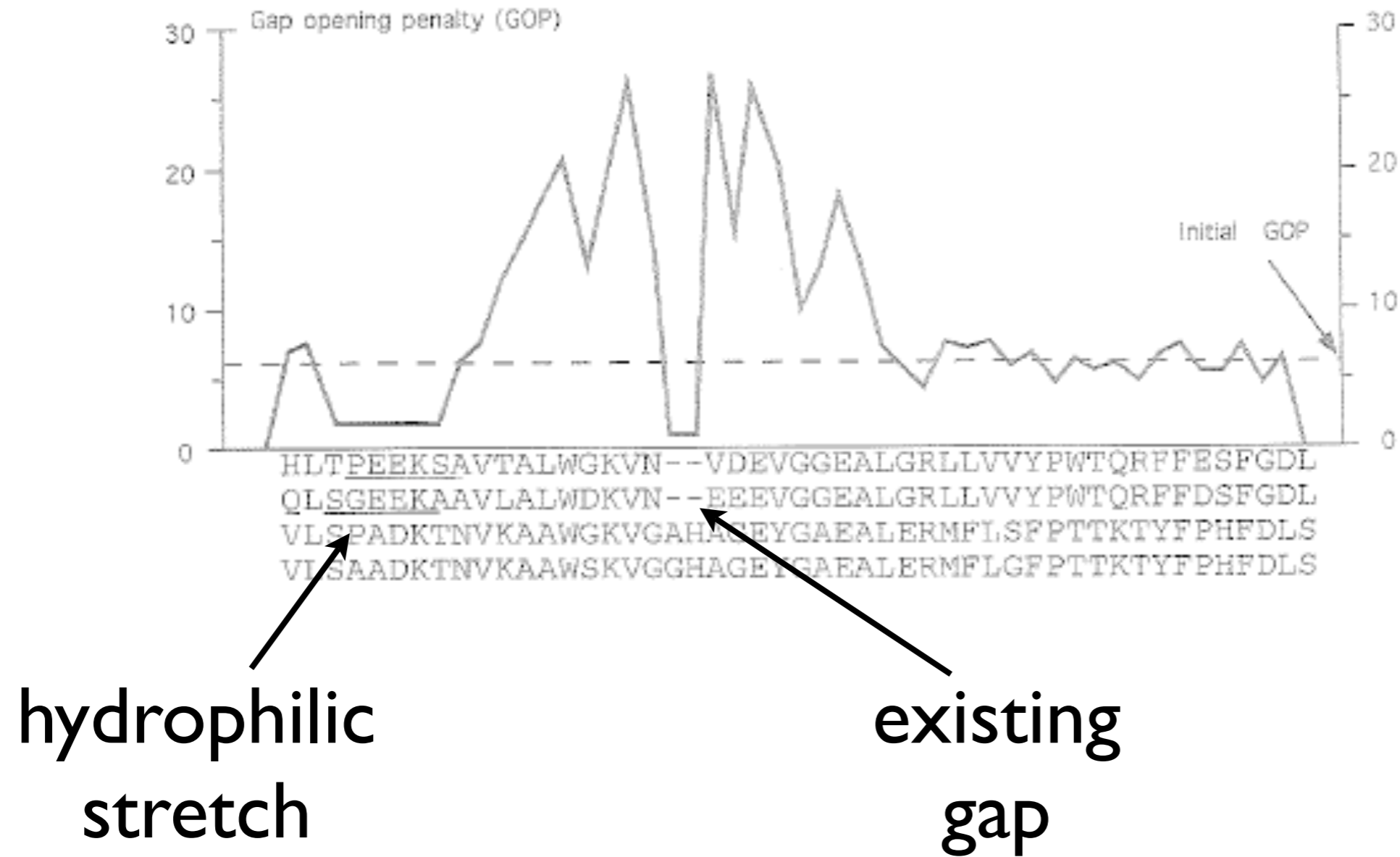


align alignments – preserve gaps

# Position Specific Gap Penalties

- There are two type of gap opening penalties: gap opening and gap extension
  - Determined empirically by user
- Decrease penalties where gaps already occurs
- Increase penalties in adjacent positions to where gap already occurs
  - Encourage extension of gaps in loop regions vs. introduction of new gaps
- Increase or decrease gap penalties according to amino acid type
  - Increase penalties in stretches of hydrophobic residues

# Gap Penalties Example



# Standard Multiple Sequence Alignment Approach

- Be as sure as possible that the sequences included are homologous
- Know as much as possible about the gene/protein in question before trying to create an alignment (secondary structure, domains etc..)
- Start with an automated alignment: preferably one that utilizes some evolutionary theory

**Which MSA tool/method should you use?**

# Comparing MSA Tools - Performance of aligning core regions on various benchmarks

Tool	Benchmark					Time
	1	2	3	4	5	
ProbCons	86.41	82.03	83.92	71.64	49.59	61h31
PCMA	85.75	80.37	90.01	69.76	46.27	11h57
MUSCLE	82.35	80.92	43.22	67.81	45.38	2h22
ClustalW	75.37	80.23	13.62	61.70	43.56	2h25
COBALT	84.44	84.40	88.13	67.05	50.50	8h54

% of letters aligned in reference alignment that are also aligned in the computed alignment



- [Help Index](#)
- [General Help](#)
- [Formats](#)
- [Gaps](#)
- [Matrix](#)
- [References](#)
- [ClustalW2 Help](#)
- [ClustalW2 FAQ](#)
- [Jalview Help](#)
- [Scores Table](#)
- [Alignment](#)
- [Guide Tree](#)
- [Colours](#)

■ [Similar Applications](#)

- [Align](#)
- [Kalign](#)
- [MAFFT](#)
- [MUSCLE](#)
- [T-Coffee](#)

■ [ClustalW Programmatic Access](#)

■ [www.clustal.org](http://www.clustal.org)

**Clustal Related Literature** 

Search for Clustal related literature in Medline...

[EBI](#) > [Tools](#) > [Sequence Analysis](#) > [ClustalW2](#)

## ClustalW2

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

[New users, please read the FAQ.](#)

>> [Download Software](#)



<b>YOUR EMAIL</b> <input type="text"/>	<b>ALIGNMENT TITLE</b> <input type="text" value="Sequence"/>	<b>RESULTS</b> <input type="button" value="interactive"/>	<b>ALIGNMENT</b> <input type="button" value="full"/>		
<b>KTUP (WORD SIZE)</b> <input type="button" value="def"/>	<b>WINDOW LENGTH</b> <input type="button" value="def"/>	<b>SCORE TYPE</b> <input type="button" value="percent"/>	<b>TOPDIAG</b> <input type="button" value="def"/>	<b>PAIRGAP</b> <input type="button" value="def"/>	
<b>MATRIX</b> <input type="button" value="def"/>	<b>GAP OPEN</b> <input type="button" value="def"/>	<b>NO END GAPS</b> <input type="button" value="yes"/>	<b>GAP EXTENSION</b> <input type="button" value="def"/>	<b>GAP DISTANCES</b> <input type="button" value="def"/>	
	<b>ITERATION</b> <input type="button" value="none"/>		<b>NUMITER</b> <input type="button" value="1"/>		
<b>OUTPUT</b>		<b>PHYLOGENETIC TREE</b>			
<b>OUTPUT FORMAT</b> <input type="button" value="aln w/numbers"/>	<b>OUTPUT ORDER</b> <input type="button" value="aligned"/>	<b>TREE TYPE</b> <input type="button" value="none"/>	<b>CORRECT DIST.</b> <input type="button" value="off"/>	<b>IGNORE GAPS</b> <input type="button" value="off"/>	<b>CLUSTERING</b> <input type="button" value="NJ"/>

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

# BLAST results sent to COBALT

generates MSA

## Cobalt Results - sp|Q02067| (231 letters) - Cobalt RID ZV85B6BE212 (100 seqs)

Descriptions Select All Re-align Alignment parameters

Legend for links to other resources: U UniGene E GEO G Gene S Structure M Map Viewer

Accession	Description	Links
<a href="#">Q02067.1</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=mASH-1; Short=mASH1 >gi 193876 gb AAA37780.1  helix-loop-heli	<a href="#">G</a>
<a href="#">Q02067.1</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=mASH-1; Short=mASH1 >gi 193876 gb AAA37780.1  helix-loop-heli	<a href="#">G</a>
<a href="#">P19359.1</a>	RecName: Full=Achaete-scute homolog 1 >ref NP_071779.1  achaete-scute homolog 1 [Rattus norvegicus] >emb CAA37760.1  unn	<a href="#">G</a>
<a href="#">P50553.2</a>	RecName: Full=Achaete-scute homolog 1; Short=ASH-1; Short=hASH1; AltName: Full=Class A basic helix-loop-helix protein 46; Shc	<a href="#">G</a>
<a href="#">Q90259.1</a>	RecName: Full=Achaete-scute homolog 1a; Short=Zash-1a; AltName: Full=Pituitary-absent protein >ref NP_571294.1  achaete-scute	<a href="#">G</a>
<a href="#">Q06234.1</a>	RecName: Full=Achaete-scute homolog 1 >ref NP_001079247.1  achaete-scute complex homolog 1 [Xenopus laevis] >gb AAA4964	<a href="#">G</a>
<a href="#">Q90260.1</a>	RecName: Full=Achaete-scute homolog 1b; Short=Zash-1b >ref NP_571306.1  achaete-scute homolog 1b [Danio rerio] >gb AAA788	<a href="#">G</a>
<a href="#">Q2EGB9.1</a>	RecName: Full=Achaete-scute homolog 2; AltName: Full=Mash2 >gb ABD39719.1  achaete scute-like protein 2 [Bos taurus] >gb AA	<a href="#">G</a>
<a href="#">Q99929.2</a>	RecName: Full=Achaete-scute homolog 2; Short=ASH-2; Short=hASH2; AltName: Full=Mash2; AltName: Full=Class A basic helix-lo	<a href="#">G</a>
<a href="#">P19360.1</a>	RecName: Full=Achaete-scute homolog 2; AltName: Full=Mash2 >ref NP_113691.1  achaete-scute homolog 2 [Rattus norvegicus] >	<a href="#">G</a>
<a href="#">Q35885.2</a>	RecName: Full=Achaete-scute homolog 2; Short=ASH-2; Short=mASH-2; Short=mASH2 >ref NP_032580.2  achaete-scute homolog	<a href="#">G</a>
<a href="#">Q7RTU5.2</a>	RecName: Full=Achaete-scute homolog 5; Short=ASH-5; Short=hASH5; AltName: Full=Class A basic helix-loop-helix protein 47; Shc	<a href="#">G</a>
<a href="#">Q6XD76.1</a>	RecName: Full=Achaete-scute homolog 4; Short=ASH-4; Short=hASH4; AltName: Full=Achaete-scute-like protein 4; AltName: Full=	<a href="#">G</a>
<a href="#">Q9NQ33.2</a>	RecName: Full=Achaete-scute homolog 3; Short=ASH-3; Short=hASH3; AltName: Full=Class A basic helix-loop-helix protein 42; Shc	<a href="#">G</a>
<a href="#">Q9JJR7.1</a>	RecName: Full=Achaete-scute homolog 3; Short=ASH-3; Short=mASH-3; Short=mASH3; AltName: Full=bHLH transcriptional regula	<a href="#">G</a>
<a href="#">P10083.1</a>	RecName: Full=Achaete-scute complex protein T5; AltName: Full=Protein achaete >ref NP_476824.1  achaete [Drosophila melanog	<a href="#">G</a>
<a href="#">P10084.2</a>	RecName: Full=Achaete-scute complex protein T4; AltName: Full=Protein scute >ref NP_476803.1  scute [Drosophila melanogaster]	<a href="#">G</a>
<a href="#">Q10007.1</a>	RecName: Full=Helix-loop-helix protein 6 >ref NP_496070.1  Helix Loop Helix family member (hlh-6) [Caenorhabditis elegans] >emb	<a href="#">G</a>
<a href="#">P09774.2</a>	RecName: Full=Achaete-scute complex protein T3; AltName: Full=Protein lethal of scute; Short=Lethal of sc >ref NP_476623.1  leth	<a href="#">G</a>
<a href="#">Q10574.2</a>	RecName: Full=Protein lin-32; AltName: Full=Abnormal cell lineage protein 32 >ref NP_508410.2  abnormal cell LINeage family men	<a href="#">G</a>

```

-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 129
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 128
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 129
-----MKRRRR----LRSDAEMQQ----LRQ-AANVRERRRMSINDAFEGLRSHIP 147
-----RRGPKK----KKMTKARLERFKLRRM-KANARERNRMHGLNAALDNLKRVVP 132
-----Q---AGNCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111
-----RRGPKK----KKMTKARMQRFKMRM-KANARERNRMHGLNDALESRLKRVVP 124
-----RRASSG----A-G----PVVVVRQRQ-AANARERDRTQSVNTAFTALRTLIP 98
-----S---PGRLE---ALGG-----RLPRRKG-SGPKKERRRTEINSAFELRECIPI 122
-----RRRRPG----PSGPGGRDSSIQRRL-ESNERERQRMHKLNNAFQALREVIP 103
-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 101
-----LKRERRR---MRSEVEMQQ----LRQ-AANVRERRRMSINDAFEGLRSHIP 143
-----RRASNG----A-G----PVVVVRQRQ-AANARERDRTQSVNTAFTALRTLIP 98
-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 113
-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 122

```

<a href="#">Q0VCE2</a>	84	-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 125
<a href="#">Q90691</a>	69	-----G---AGRLE---ALSG-----RLGRRKGVGGPKKERRRTEINSAFELRECIPI 111
<a href="#">Q91616</a>	84	-----RRGPKK----KKMTKARVERFKVRRM-KANARERNRMHGLNDALDSLKRVVP 130
<a href="#">Q6QHK4</a>	60	-----L-QL-----VLERRR-VANAKERERIKNLNRGFARLKLALVP 93
<a href="#">P17542</a>	142	QPLASLGSGFFGEPDAPPMFTNNRVKRRSPYE----MEITDGPHT-KVVRRI-FTNSRERWRQNVNGAFELRKLIP 215
<a href="#">P24699</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111
<a href="#">Q91154</a>	63	-----Q---AGHCL--LWACKACKRKSSTTDRRK-AATMRERRRLKKNVNSAFETLKRCTT 111
<a href="#">P22091</a>	142	QPLASLGSGFFGEPDAPPMFTNNRVKRRSPYE----MEISDGPHT-KVVRRI-FTNSRERWRQNVNGAFELRKLIP 215
<a href="#">P57100</a>	81	-----S---PGRLE---ALGG-----RLGRRKG-SGPKKERRRTEINSAFELRECIPI 122
<a href="#">Q63689</a>	103	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 149
<a href="#">Q62414</a>	104	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 150
<a href="#">P70447</a>	94	-----RAVSRG----AKTAETVQRIKTRRL-KANNRERNRMHNLNAALDALREVLP 140
<a href="#">Q15784</a>	103	-----KRGPKK----RKMTKARLERSKLRRQ-KANARERNRMHDLNAALDNLKRVVP 149
<a href="#">P17667</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFDTLKRCTT 111
<a href="#">Q55208</a>	54	-----L-HL-----VLERRR-VANAKERERIKNLNRGFARLKLALVP 87
<a href="#">P13349</a>	63	-----Q---AGHCL--MWACKACKRKSSTTDRRK-AATMRERRRLKKNVQAFETLKRCTT 111

More details in Papadopoulos JS and Agarwala R, Bioinformatics 23:1073-79, 2007 (PMID: 17332019)

# Phylogenetic Tree View - based on COBALT multiple alignment

**COBALT** *Phylogenetic Tree View*

This tree is based on COBALT multiple alignment [more...](#)

## Phylo Tree View for 100 sequences: Cobalt RID ZV85B6BE212

Tree method: Fast Minimum Evolution | Max Seq Difference: 0.85 | Distance: Grishin (protein) | Reset | Download in: Newick Format

rectangle | slanted | radial | force |  Show distance | Mouse over an internal node for a subtree or alignment | [Hide Color Map](#) | [Show removed sequences](#)

Sequence Label: Sequence Title (if available)

Collapse Mode: Blast Name

Blast names color map	
<span style="color: green;">■</span>	primates
<span style="color: blue;">■</span>	rodents
<span style="color: yellow;">■</span>	bony fishes
<span style="color: orange;">■</span>	birds
<span style="color: cyan;">■</span>	frogs & toads
<span style="color: magenta;">■</span>	even-toed ungulates
<span style="color: darkgreen;">■</span>	salamanders
<span style="color: grey;">■</span>	nematodes
<span style="color: pink;">■</span>	flies

- [Help Index](#)
- [General Help](#)
- [Formats](#)
- [Gaps](#)
- [Matrix](#)
- [References](#)
- [Muscle Help](#)
- [Jalview Help](#)

■ [Similar Applications](#)

- [Align](#)
- [ClustalW2](#)
- [Kalign](#)
- [MAFFT](#)
- [T-Coffee](#)

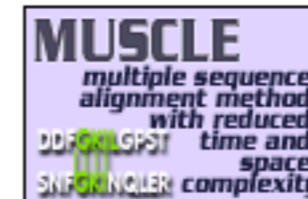
■ [Muscle Programmatic Access](#)

EBI > Tools > Sequence Analysis

## MUSCLE

MUSCLE stands for **M**ultiple **S**equence **C**omparison by **L**og-**E**xpectation. MUSCLE is claimed to achieve both better average accuracy and better speed than [ClustalW2](#) or [T-Coffee](#), depending on the chosen options.

 [Download Software](#)



<b>RESULTS</b> <input type="text" value="interactive"/>	<b>SEARCH TITLE</b> <input type="text" value="Sequence"/>	<b>YOUR EMAIL</b> <input type="text"/>
<b>OUTPUT FORMAT</b> <input type="text" value="FASTA"/>	<b>OUTPUT TREE</b> <input type="text" value="none"/>	<b>OUTPUT ORDER</b> <input type="text" value="aligned"/>

Enter or Paste a set of Sequences in any supported format:

Upload a file:  no file selected

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

# Standard Multiple Sequence Alignment Approach

Examine alignment:

- Are you confident that aligned residues/bases evolved from a common ancestor?
- Are domains of the proteins/predicted secondary structures, etc. aligning correctly?
- Are most indels outside of known motifs or secondary structure?
  - No? May need to edit sequences and redo...

# The Take Home Message

Why perform an MSA?

- Visualize trends between homologous sequences
  - Shared regions of homology
  - Regions unique to a sequence within a family
  - Consensus sequence
- As the first step in a phylogenetic analysis

# The Take Home Message

How does one perform an MSA?

- By hand: too hard!
- Automated alignment: Fast, but doesn't necessarily produce the "correct" alignment
- Developing methods of MSA is an active area of research

**Best approach = Automated alignment  
with manual editing**

# MSA

PRACTICAL EXERCISE: Comparing Sets of Protein Sequences





navigate to:  
[bioteach.ubc.ca/bioinfo2010](http://bioteach.ubc.ca/bioinfo2010)

We'll walk through  
 install + do MSA #1  
 together

# Clustal

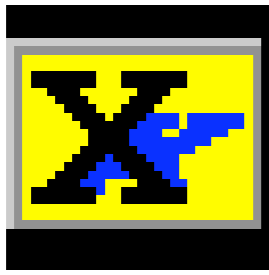
Install ClustalX on laptop

Use ClustalX to generate MSA

download program and install

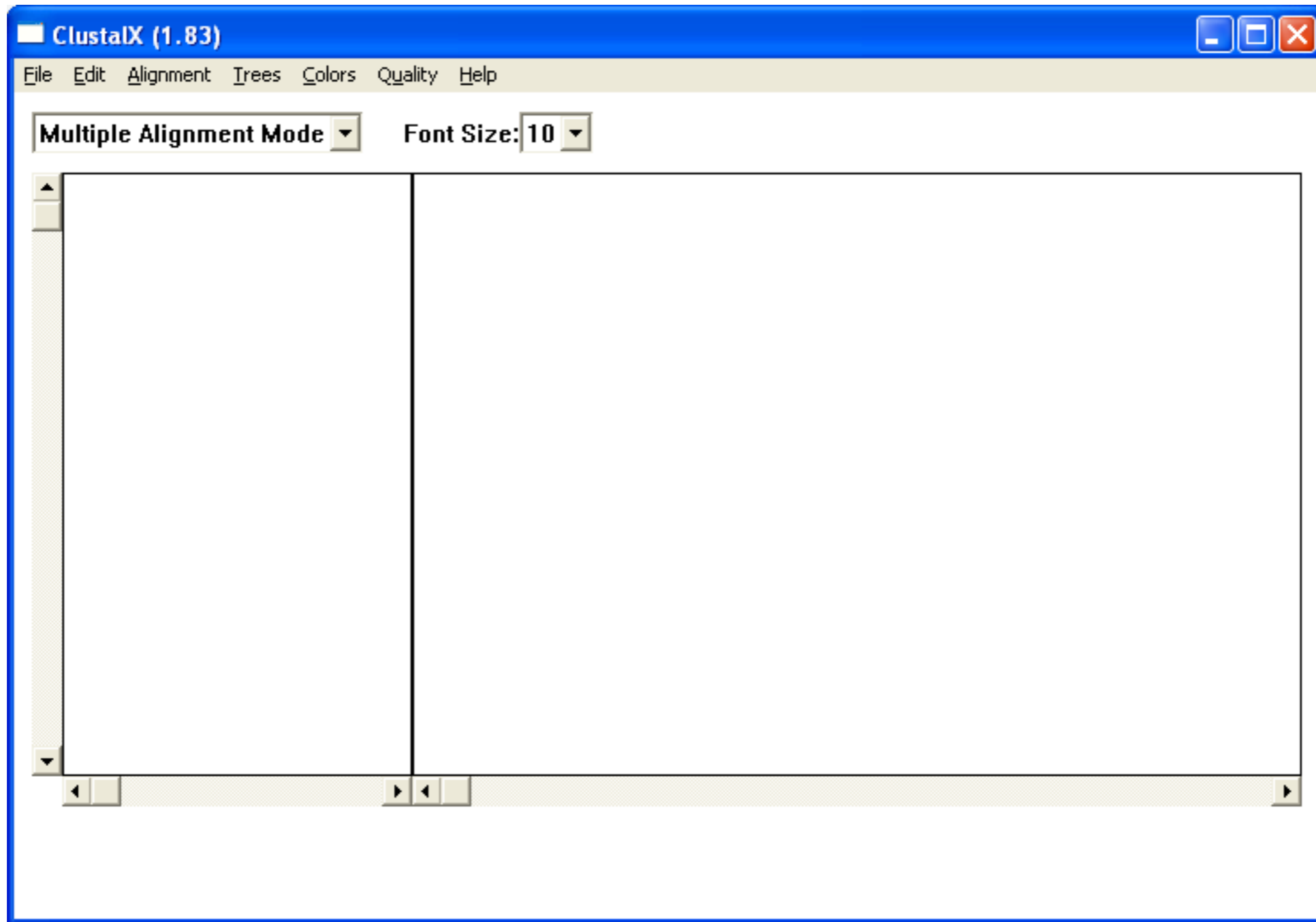
MSA #1: Use example sequences to generate alignment

MSA #2: Use your own sequences



Clustalx.exe

# Open ClustalX



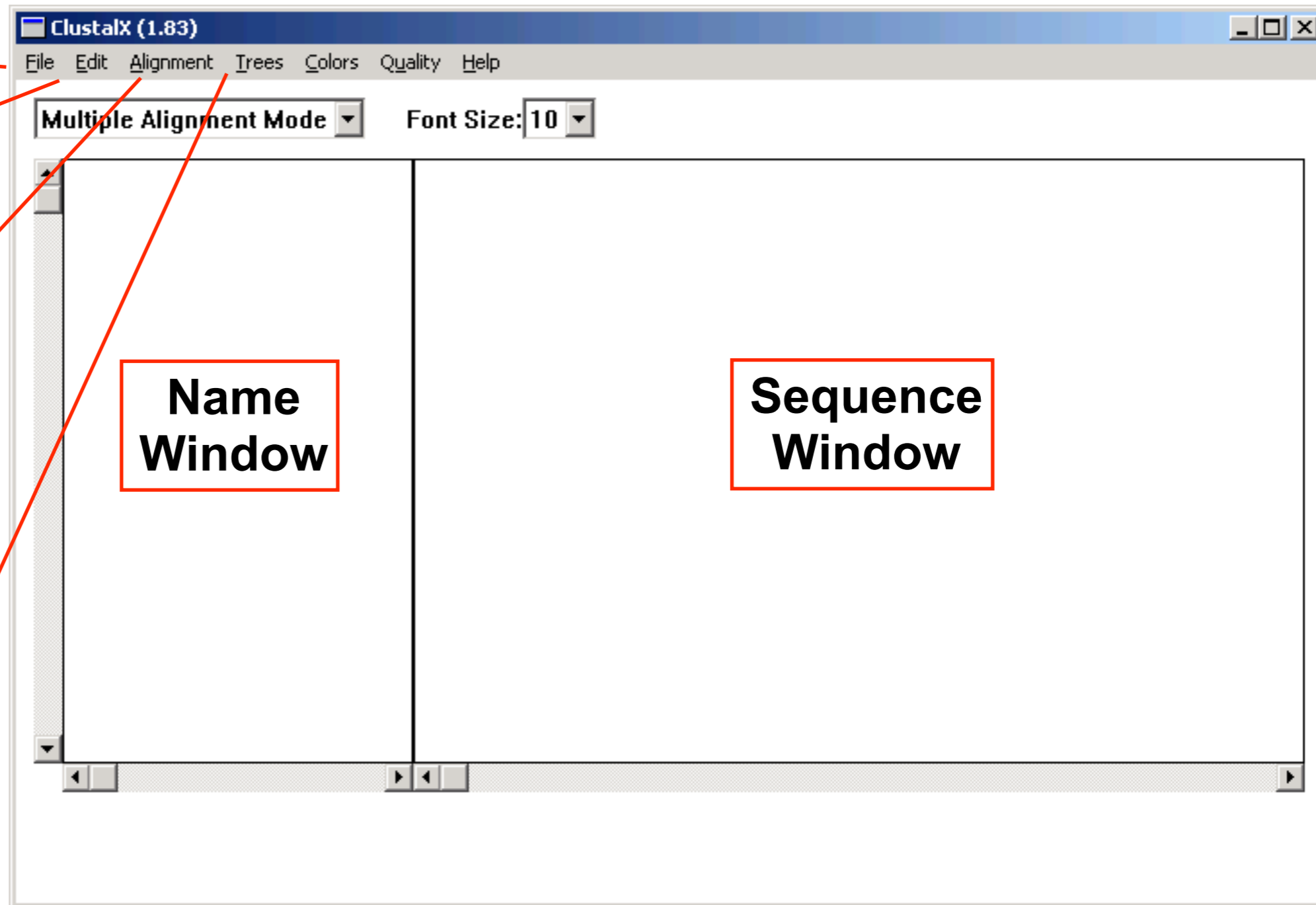
# Starting up ClustalX

File:  
-Load sequences

Edit:  
-Remove all gaps

Alignment:  
-Do complete alignment  
-Alignment parameters

Trees:  
-Bootstrapped NJ  
-Output format options



```

globin.pep - WordPad
File Edit View Insert Format Help eFax
[Icons]

>P1;HBB_HUMAN
Sw:Hbb_Human => HBB_HUMAN
      VHLTPEEKSA VTALWGKVNV DEVGGEALGR LLVYYPWTQR FFESFGDLST
      PDAVMGMPKV KAHGKKVLGA FSDGLAHLDM LKGTFAATLSE LHCDKLHVDP
      ENFRLLGNVL VCVLAHHFGK EFTPPVQAAAY QKVVAGVANA LAHKYH*
C;ID   HBB_HUMAN          STANDARD;          PRT;    146 AA.
C;AC   PO2023;
C;DT   21-JUL-1986 (REL. 01, CREATED)
C;DT   21-JUL-1986 (REL. 01, LAST SEQUENCE UPDATE)
C;DT   01-APR-1993 (REL. 25, LAST ANNOTATION UPDATE)
C;DE   HEMOGLOBIN BETA CHAIN. . . .

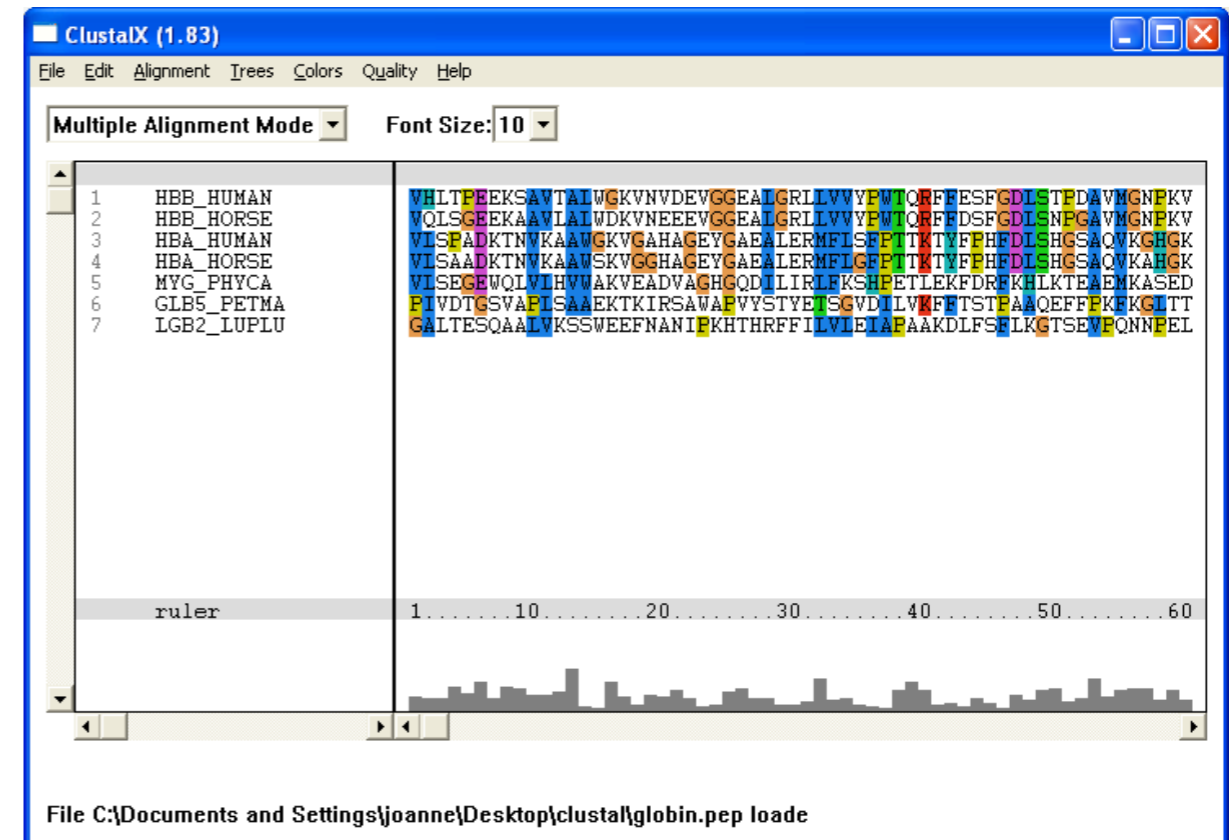
>P1;HBB_HORSE
Sw:Hbb_Horse => HBB_HORSE
      VQLSGEEKAA VLALWDKVNE EEVGGEALGR LLVYYPWTQR FFDSFGDLSN
      PGAVMGMPKV KAHGKKVLHS FGEGVHHLDM LKGTFAALSE LHCDKLHVDP
      ENFRLLGNVL VVVLARHFGK DFTPELQASY QKVVAGVANA LAHKYH*
C;ID   HBB_HORSE          STANDARD;          PRT;    146 AA.
C;AC   PO2062;
C;DT   21-JUL-1986 (REL. 01, CREATED)
C;DT   21-JUL-1986 (REL. 01, LAST SEQUENCE UPDATE)
C;DT   01-MAR-1992 (REL. 21, LAST ANNOTATION UPDATE)
C;DE   HEMOGLOBIN BETA CHAIN. . . .

>P1;HBA_HUMAN
Sw:Hba_Human => HBA_HUMAN
      VLSPADKTNV KAAWGKVGAAH AGEYGAEALE RMFLSFPTTK TYFPDFDLSH
      GSAQVKGHGK KVADALTNV AHVDDMPNAL SALSDLHAHK LRVDPVNFKL
      LSHCLLVTLA AHLPAEFTPA VHASLDKFLA SVSTVLTSKY R*
C;ID   HBA_HUMAN          STANDARD;          PRT;    141 AA.
C;AC   PO1922;
C;DT   21-JUL-1986 (REL. 01, CREATED)
C;DT   21-JUL-1986 (REL. 01, LAST SEQUENCE UPDATE)
C;DT   01-FEB-1994 (REL. 28, LAST ANNOTATION UPDATE)
C;DE   HEMOGLOBIN ALPHA CHAIN. . . .

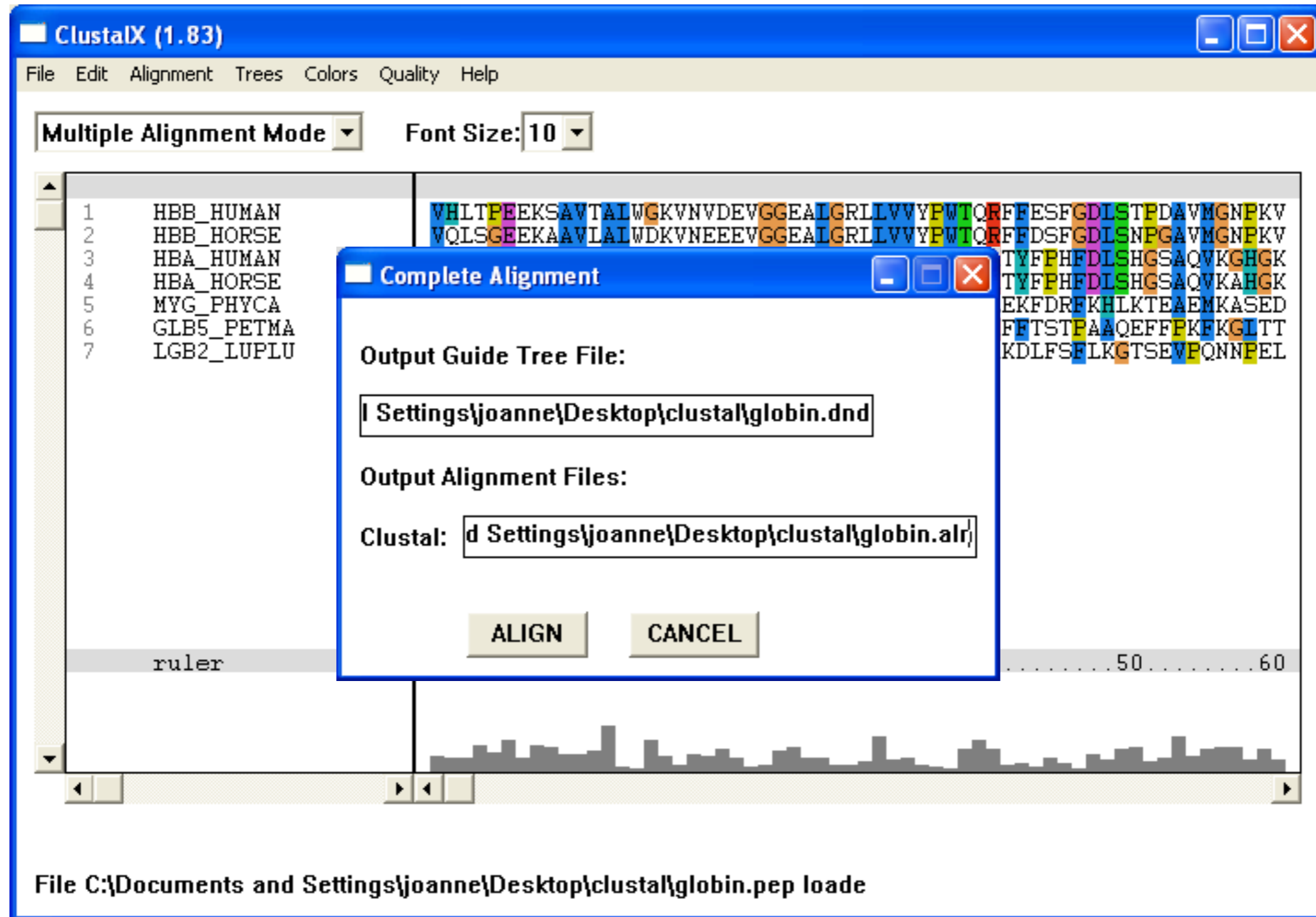
>P1;HBA_HORSE
Sw:Hba_Horse => HBA_HORSE
      VLSAADKTNV KAAWSKVGGAH AGEYGAEALE RMFLGFPTTK TYFPDFDLSH
      GSAQVKAHGK KVGDALTLAV GHLDDLPGAL SMLS DLHAHK LRVDPVNFKL
      LSHCLLSTLA VHLPNDFTPA VHASLDKFLS SVSTVLTSKY R*
C;ID   HBA_HORSE          STANDARD;          PRT;    141 AA.
C;AC   PO1958;
C;DT   21-JUL-1986 (REL. 01, CREATED)
C;DT   21-JUL-1986 (REL. 01, LAST SEQUENCE UPDATE)
C;DT   01-MAR-1992 (REL. 21, LAST ANNOTATION UPDATE)
C;DE   HEMOGLOBIN ALPHA CHAINS (SLOW AND FAST). . . .

```

# Load the sequences -globin.pep



# Alignment > Do Complete Alignment



also see: Alignment > Alignment Parameters

ClustalX (1.83)

File Edit Alignment Trees Colors Quality Help

Multiple Alignment Mode Font Size: 10

1	HBB_HUMAN	-----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYFPWTQRFFESEFGDLSI
2	HBB_HORSE	-----VQLSGEERAAVLALWDKVN--EEEVGGEALGRLLVVYFPWTQRFFDSFGDLSN
3	HBA_HUMAN	-----VLSPADKTNVKAAWGKVGAAHAGEYGAELERMFLSFPTTKTYFPHF-DLS-
4	HBA_HORSE	-----VLSAADKTNVKAAWSKVGGAHAGEYGAELERMFLGFPTTKTYFPHF-DLS-
5	GLB5_PETMA	FIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSQVDILVKFFTSTPAAQEFFPKFKGLTT
6	MYG_PHYCA	-----VISEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDREKHLKT
7	LGB2_LUPLU	-----GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSE

ruler 1.....10.....20.....30

CLUSTAL-Alignment file created []

Help

ALIGNMENT DISPLAY

The alignment is displayed on the screen with the sequence names on the left hand side. The sequence alignment is for display only, it cannot be edited here (except for changing the sequence order by cutting-and-pasting on the sequence names).

A ruler is displayed below the sequences, starting at 1 for the first residue position (residue numbers in the sequence input file are ignored).

A line above the alignment is used to mark strongly conserved positions. Three characters ('\*', ':' and '.') are used:

- '\*' indicates positions which have a single, fully conserved residue
- ':' indicates that one of the following 'strong' groups is fully conserved:-
  - STA
  - NEQK
  - NHQK
  - NDEQ
  - QHRK
  - MILV
  - MILF
  - HY
  - FYW
- '.' indicates that one of the following 'weaker' groups is fully conserved:-
  - CSA
  - ATV
  - SAG
  - STNK
  - STPA
  - SGND
  - SNDEQK
  - NDEQHK
  - NEQHRK
  - FVLIM
  - HFY

These are all the positively scoring groups that occur in the Gonnet Pam250 matrix. The strong and weak groups are defined as strong score >0.5 and weak

OK

see: Help > General

# Let's start at 9:30am

Genome Browsers

GEO - gene expression omnibus

Pathway Resources for Systems Biology

