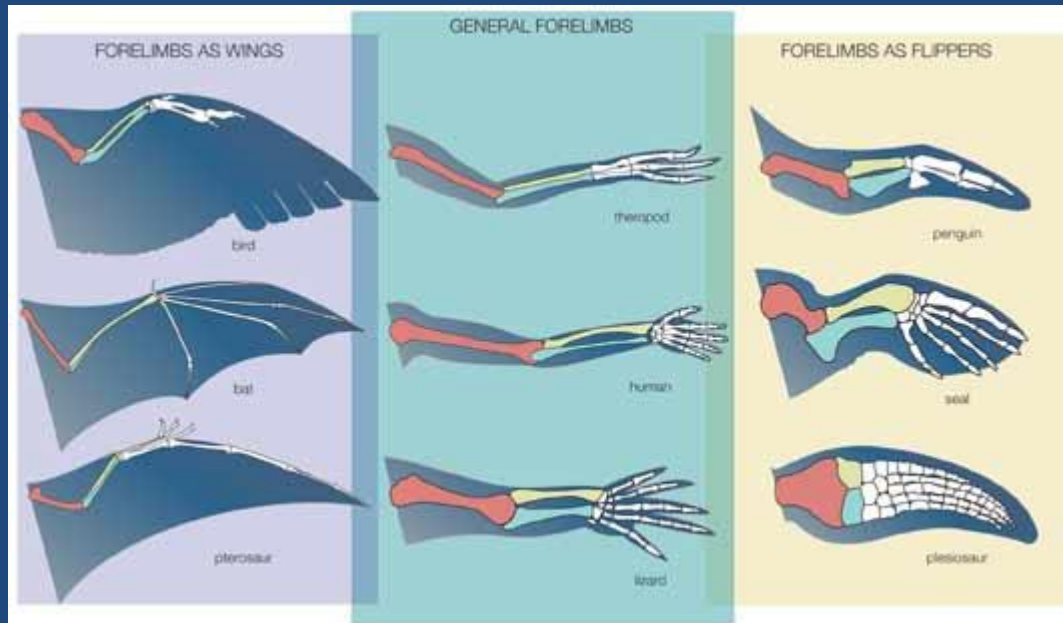


# Bioinformatics Tools for Sequence homology and alignment

# Homology

- Similarity between characters due to a common ancestry



# Sequence homology

- Similarity between sequences that results from a common ancestor

**VLSPAVKWAKVGAHAAGHG**

**VLSEAVLWAKVEADVAGHG**

◆ Basic assumption:

Sequence homology →  
similar structure/function

# Sequence alignment

**Alignment:** Comparing two (pairwise) or more (multiple) sequences. Searching for a series of identical or similar characters in the sequences.

# Homology

- Ortholog – homolog with similar function (via speciation)
- Paralog – homolog which arose from gene duplication

Common use:

**Orthologs** –  
2 homologs  
from **different**  
species

**Paralogs** –  
2 homologs  
within the  
**same** species

# How close?

- Rule of thumb:
- Proteins are homologous if 25% identical  
(length >100)
- DNA sequences are homologous if 70% identical



# Twilight zone

- < 20% identity in proteins – may be homologous and may not be....
- (Note that 5% identity will be obtained completely by chance!)



# Local vs. Global

- **Global alignment** – finds the best alignment across the **entire** two sequences.

```

ADLGAVFALCDRYFQ
|||||      |||||
ADLGRTON-CDRYYQ
    
```

- **Local alignment** – finds regions of similar **parts** of the sequences.

```

ADLG      CDRYFQ
|||||    |||||
ADLG      CDRYYQ
    
```

Global alignment: forces alignment in regions which differ

Local alignment will return only **regions** of good alignment



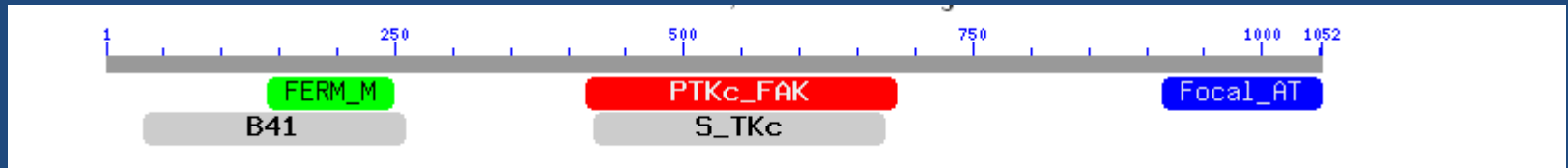
When global and when local?

# Global alignment

- PTK2 protein tyrosine kinase 2 of human and rhesus monkey

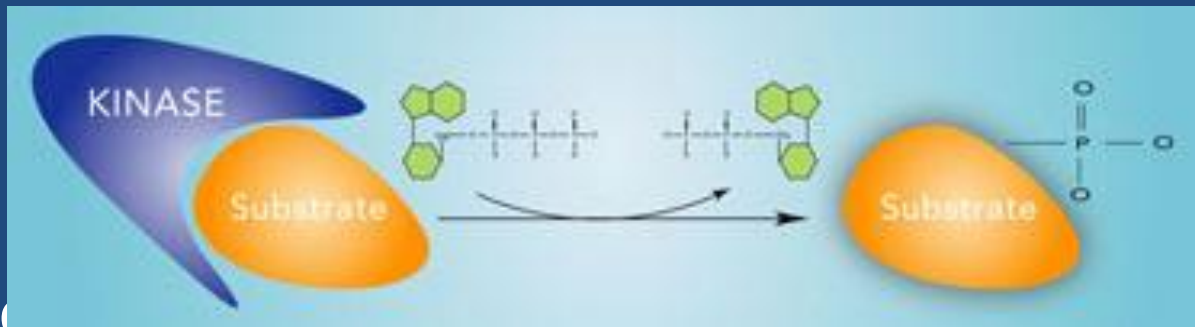
human	107	VREKYELAHPPEEWKYLRI RYLPKGFLNQFTEDKPTLNFFYQQVKS DYM	156
rhesus	151	VREKYELAHPPEEWKYLRI RYLPKGFLNQFTEDKPTLNFFYQQVKS DYM	200
human	157	LEIADQVDQEIALKLGCLEIRRSYWEMRGNALEKKS NYEVLEKDVGLKRF	206
rhesus	201	LEIADQVDQEIALKLGCLEIRRSYWEMRGNALEKKS NYEVLEKDVGLKRF	250
human	207	FPKSLDSVKAKTLRKL IQQTFRQFANLNREESILKFFEILSPVYRFDKE	256
rhesus	251	FPKSLDSVKAKTLRKL IQQTFRQFANLNREESILKFFEILSPVYRFDKE	300
human	257	CFKCALGSSWII SVELAIGPEEGISYLTDRGCNPTHLADFTQVQTIQYSN	306
rhesus	301	CFKCALGSSWII SVELAIGPEEGISYLTDRGCNPTHLADFTQVQTIQYSN	350
human	307	SEKDRKGM LQLKIAGAPEPLTVTAPSLTIAENMADLIDGYCRLVNGTSQ	356
rhesus	351	SEKDRKGM LQLKIAGAPEPLTVTAPSLTIAENMADLIDGYCRLVNGASQ	400
human	357	SFIIRPQKEGERALPSIPKLANSEKQGM RTHAVSVSETDDYAEIIDEEDT	406
rhesus	401	SFIIRPQKEGERALPSIPKLANSEKQGM RTHAVSVSETDDYAEIIDEEDT	450
human	407	YTMPSTRDYEIQRERIELGRCIGEGQFGDVHQGIYMS PENPALAVAIKTC	456
rhesus	451	YTMPSTRDYEIQRERIELGRCIGEGQFGDVHQGVYMS PENPALAVAIKTC	500
human	457	KNCTSDSVREKFLQEAL TMRQFD-HPHIVKLVITENPVWIIIMELCTLG	505
rhesus	501	KNCTSDSVREKFLQEAL TMRQFD-HPHIVKLVITENPVWIIIMELCTLG	550

# Protein tyrosine kinase domain



# Protein tyrosine kinase domain

- Human PTK2 and leukocyte tyrosine kinase
- Both function as tyrosine kinases, in completely different contexts



- Ancient duplication

# Global alignment of PTK and LTK



# Local alignment of PTK and LTK

human_ptk2	343	LIDGYCRLVNGTSQSFIIRPQKE----GERALPSIPKLANSEKQGMRTHA	388
		:        : . .     . :..  :     . .     :     ...: . .	
human_LTK	439	LL-----MVCGV---LILVKQKKWQGLQEMRLPS-PEL---ELSKLRTSA	476
human_ptk2	389	VSVSETDDYAEI-IDEEDTYTMPSTRDYEQERERIELGRCIGEGQFGDVH	437
		:...:.. .:: :...::: . ....  :...: . .:. . .  : :	
human_LTK	477	IRTAPNPYYCQVGLGPAQSWPLPPGVT-EVSPANVTLLRALGHGAFGEVY	525
human_ptk2	438	QG--IYMSPENPALAWAIKTCKNCTSDSVREKFLQEALTMROFDHPHIVK	485
		:     ...:.. .     ..... ..... .    :..: . .:. :	
human_LTK	526	EGLVIGLPGDSSPLQVAIKTLPELCSPQDELDFLMEALIISKFRHQNIWR	575
human_ptk2	486	LIGV-ITENPVWIIMELCTLGELRSFLQVRKYSLD-----LASLILYAY	528
		.: : :... . .  : :.. .  : : : : : : : : : : : : : : : : :	
human_LTK	576	CVGLSLRATPRLILLELMSGDMKSFLRHSRPHLGQPSPLVWRDLLQLAQ	625
human_ptk2	529	QLSTALAYLESKRFBVHRDIAARNVLVS---SNDCVKLGD FGLSRYMEDST	575
		:...:..   ... :       .  :     :... :    : . :..:..:	
human_LTK	626	DIAQGCHYLEENHFIHRDIAARNCLLSCAGPSRVAKIGDFGMARDIYRAS	675
human_ptk2	576	YY-KASKGKLPKWMAPESINFRRFTSASDVWMFGVCMWEILMHGVKPFQ	624
		:...:..  : .. .  :..... ..: . .   : .. . . . . :.	
human_LTK	676	YYRRGDRALLPVKWPPEAFLEGIFTSKTD SWSFGVLLWEIFSLGYMPYP	725
human_ptk2	625	GVKMNADVIGRIENGERLPMPNCPPTLYSLMTCWAYDPSRRPRFTE---	671
		.. .  : :..:.. .  :.. .  ..: .  : : : : : : : : : : : : :	
human_LTK	726	GRTNQEVLD FVVGGRMDPPRGCPGVYRIMTQCWQHEPELRPSFASILE	775

# Searching databases

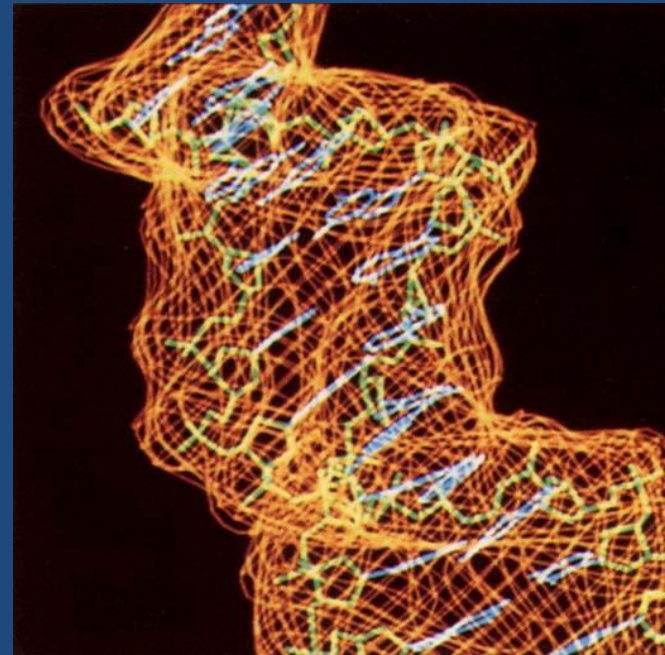
# Searching a database

- Using a sequence as the query to find **homologous** sequences in the database



# DNA or protein?

- For coding sequences, we can use the DNA sequence or the protein sequence to search for similar sequences.
- Which is preferable?



# Protein is better!

- Selection (and hence conservation) works on the protein level:

CTTTCA = Leu-Ser

TTGAGT = Leu-Ser

# Query type

- ◆ Nucleotides: a four letter alphabet
- ◆ Amino acids: a twenty letter alphabet



- Two random DNA sequences will share on average 25% of identity
- Two random protein sequences will share on average 5% of identity

# Conclusions

- Using the amino acid sequence is preferable for homology search
- Why use a nucleotide sequence after all?
- No ORF found, e.g. newly sequenced genome
- No similar protein sequences were found
- Specific DNA databases are available (EST)

# Some terminology

- **Query sequence** - the sequence with which we are searching
- **Hit** – a sequence found in the database, suspected as homologous

# How do we search a database?

- Assume we perform pairwise alignment of the query against all the sequences in the database
- Exact pairwise alignment is  $O(mn) \approx O(n^2)$   
( $m$  – length of sequence 1,  
 $n$  – length of sequence 2)

# How much time will it take?

- $O(n^2)$  computations per search.
- Assume  $n=200$ , so we have **40,000** computations per search
- Size of database - **~60 million entries**
- **$2.4 \times 10^{12}$**  computations for each sequence search we perform!
- Assume each computation takes  $10^{-6}$  seconds  $\rightarrow$  24,000 seconds  $\approx$  **6.66 hours for each sequence search**
- **150,000** searches (at least!!) are performed per day

# Conclusion

- Using the exact comparison pairwise alignment algorithm between query and all DB entries – too slow





# Heuristic

- **Definition:** a heuristic is a design to solve a problem that does not provide exact solution (but is not too bad) and reduces the time complexity of the exact solution

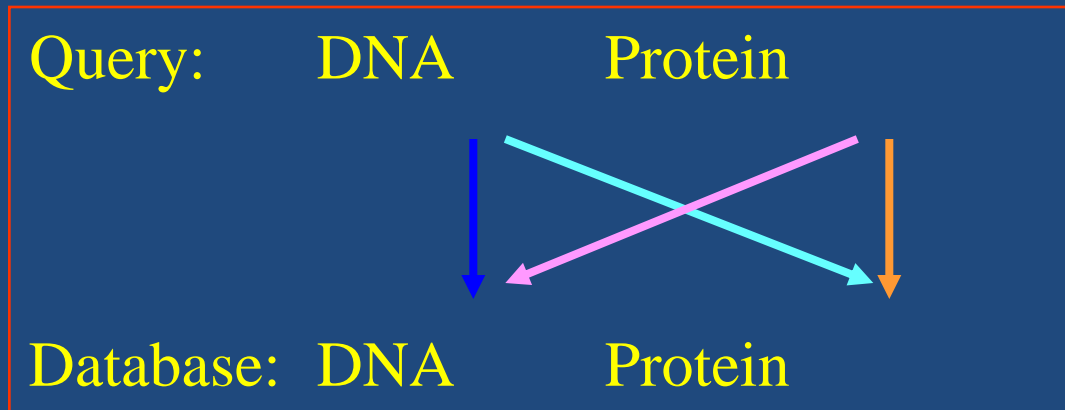
# BLAST

- BLAST - Basic Local Alignment and Search Tool
- A heuristic for searching a database for similar sequences



# DNA or Protein

- All types of searches are possible.



**blastn** – nuc vs. nuc

**blastp** – prot vs. prot

**blastx** – translated query vs. protein database

**tblastn** – protein vs. translated nuc. DB

**tblastx** – translated query vs. translated database

Translated  
databases:

trEMBL

genPept

# BLAST - underlying hypothesis

- **The underlying hypothesis:** when two sequences are similar there are **short ungapped regions of high similarity** between the two
- **The heuristic:**
  1. Discard irrelevant sequences
  2. Perform exact **local** alignment with remaining sequences

# How do we discard irrelevant sequences quickly?

- Divide the **database** into **words** of length  $w$  ( $w = 3$  for protein and  $w = 7$  for DNA)
- Save the words in a look-up table that can be searched quickly

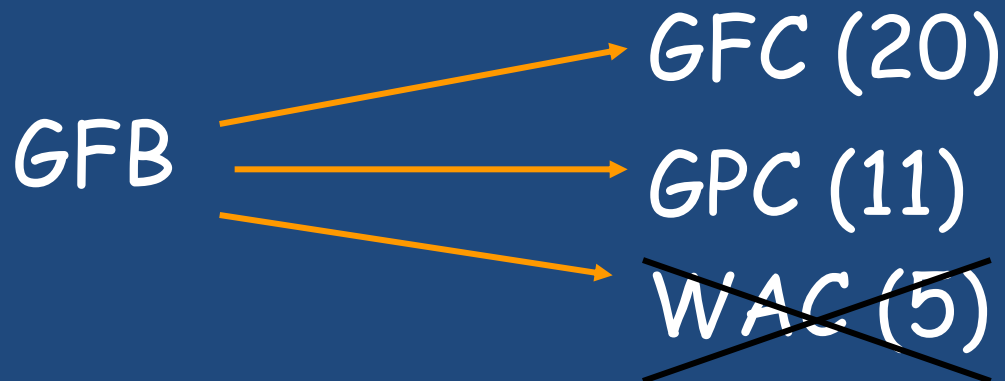


# BLAST: discarding sequences

- When the user gives a query sequence, divide it also into words
- Search the **database** for consecutive neighbor words

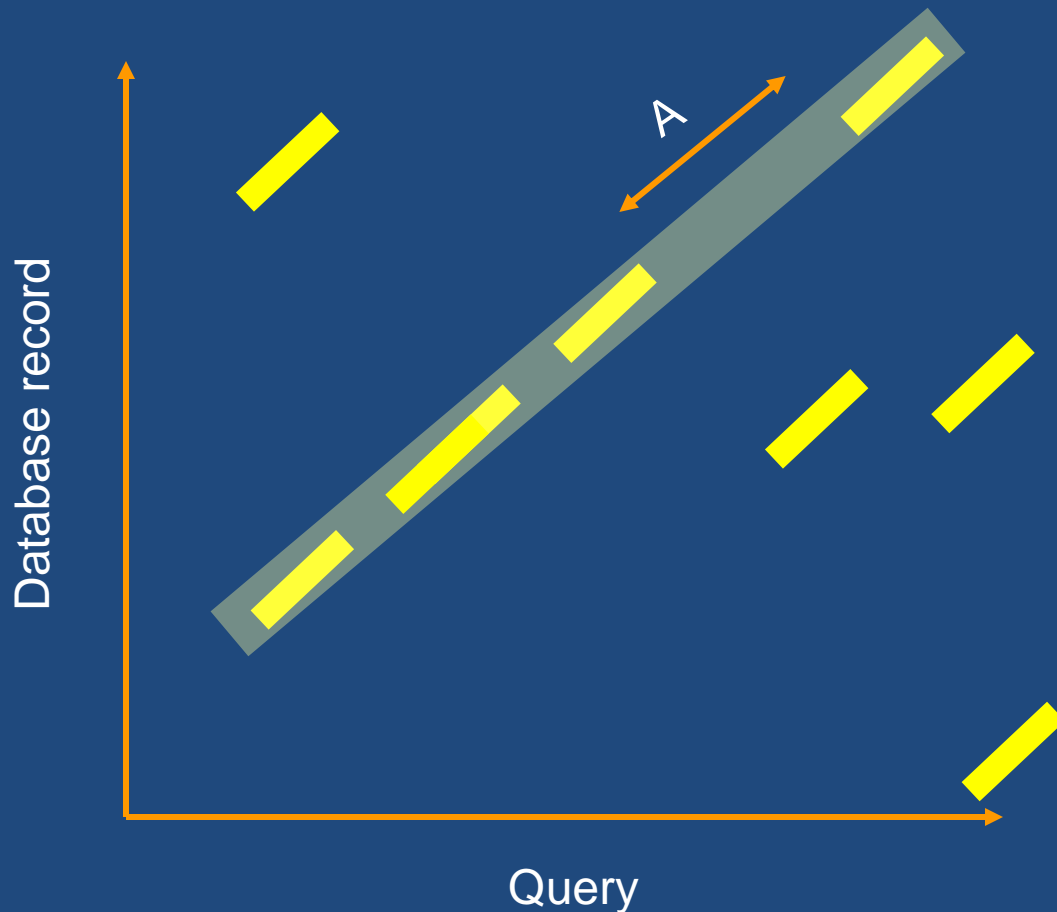
# Neighbour words

- **neighbor words** are defined according to a scoring matrix (e.g. BLOSUM62 for proteins) with a certain cutoff level



# Search for consecutive words

Neighbor word



Look for a seed: hits on the same diagonal which can be connected

At least 2 hits on the same diagonal with distance which is smaller than a predetermined cutoff

This is the filtering stage – many unrelated hits are filtered, saving lots of time!



# The result – local alignment

- The result of BLAST will be a series of **local alignments** between the query and the different hits found

# E-value

- The number of times we will theoretically find an alignment with a score  $\geq Y$  of a random sequence vs. a random database

Theoretically,  
we could trust  
any result  
with an  
E-value  $\leq 1$

In practice – BLAST uses estimations.  
E-values of  $10^{-4}$  and lower indicate a  
significant homology.  
E-values between  $10^{-4}$  and  $10^{-2}$  should  
be checked (similar domains, maybe  
non-homologous).  
E-values between  $10^{-2}$  and 1 are  
suspicious...

# Filtering low complexity

- **Low complexity regions** : e.g., Proline rich areas (in protein), Alu repeats (in DNA)
- Regions of low complexity generate high score of alignment BUT – this does not indicate homology

# Solution

- In BLAST there is masking of low-complexity regions in the query sequence (such regions are represented as XXXXX in query)

# BLAST Programs

- You can do it on line:

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

The screenshot shows the NCBI BLAST website interface. At the top, there is a navigation bar with "BLAST" and "Basic Local Alignment Search Tool" in the title, and "Home", "Recent Results", "Saved Strategies", and "Help" as menu items. Below the navigation bar, there is a section for "NCBI/BLAST Home" with a description: "BLAST finds regions of similarity between biological sequences. [more...](#)". A red box highlights a new feature: "New Aligning Multiple Protein Sequences? Try the COBALT Multiple Alignment Tool. [Go](#)".

The main content area is divided into several sections:

- BLAST Assembled Genomes**: "Choose a species genome to search, or [list all genomic BLAST databases](#)." This section lists various species with checkboxes: Human, Mouse, Rat, Arabidopsis thaliana, Oryza sativa, Bos taurus, Danio rerio, Drosophila melanogaster, Gallus gallus, Pan troglodytes, Microbes, and Apis mellifera.
- Basic BLAST**: "Choose a BLAST program to run." This section lists several programs with descriptions and algorithms:
  - [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)." This section includes a checkbox for "Make specific primers with [Primer-BLAST](#)".

```
[shzhang@usa]$ ./formatdb -help
```

```
formatdb 2.2.15 arguments:
```

```
-t Title for database file [String] Optional
```

```
-i Input file(s) for formatting [File In] Optional
```

```
-l Logfile name: [File Out] Optional
```

```
default = formatdb.log
```

```
-p Type of file
```

```
    T - protein
```

```
    F - nucleotide [T/F] Optional
```

```
default = T
```

1. Check whether BLAST is in your path  
> which blastall
  
2. Target sequences should be formatted before it's searched against.
  - a. Copy E.Coli protein sequences (NC\_00913.faa)
  
  - b. Now perform 'formatdb' in the BLAST directory  
>formatdb -i NC\_000913.faa -n EColi -p T
  - c. You will see these files created in the same directory.  
EColi.pin, EColi.psq, EColi.phr, formatdb.log
  
3. Let's perform a simple BLAST of "proteinSeq1.txt"
  - a. Copy the "proteinSeq1.txt" into the BLAST directory.
  
  - b. >blastall -p blastp -d EColi -i proteinSeq1.txt -o proteinSeq1.out  
blastall -p blastp -d EColi -i proteinSeq1.txt
  
4. Change the following options
  - A. -e : expectation value (Default: 10)
  - B. -m : alignment view option (Default: 0)
  - C. -b : Number of database sequences to show alignments (Default: 250)
  - D. -v : Number of database sequences to show one-line descriptor (Default: 500)
  - E. -g : Perform gapped alignment (Default: T)
  - F. -M : Scoring Matrix (Default: BLOSUM62)
  
5. There are many options you can adjust. Simply run blastall without any option.
  
6. Try to make BLAST print out result in html (with -T T)  
>blastall -p blastp -d EColi -i proteinSeq1.txt -o //index.html -T T

Find    [e.g. [D1S2806](#), [AP000869](#), [cancer](#)]

powered by  
**COMPAQ** NonStop™

## Ensembl BLAST Server

### RETRIEVE BLAST RESULTS

Enter the blast retrieval ID:

### SUBMIT A BLAST QUERY

Paste your sequence here in FASTA or plain text format.

```
mkwwwallll aawaaardc rvssfvken fdkarfsgtw yamakkdpeg lfqdnivae
fsvdetgqms atakgrvll nnwdvcadm vgtfdtedpa kfkmywgva sflqkgn dc
wivdtdydyt avqyscrlln ldgtcadsys mfsrdpngl ppeaqkivrq rqeelclarq
yrlivhngyc dgrsernll
```

OR select the sequence file you wish to search

### BLAST OPTIONS

Database

Executable

Report  alignments.

Mask repetitive sequences using Repeatmasker.

[Filter](#) low complexity regions.

Display histogram of score statistics.

### ADVANCED BLAST OPTIONS

Matrix  Expect (E)

Descriptions  HSP score

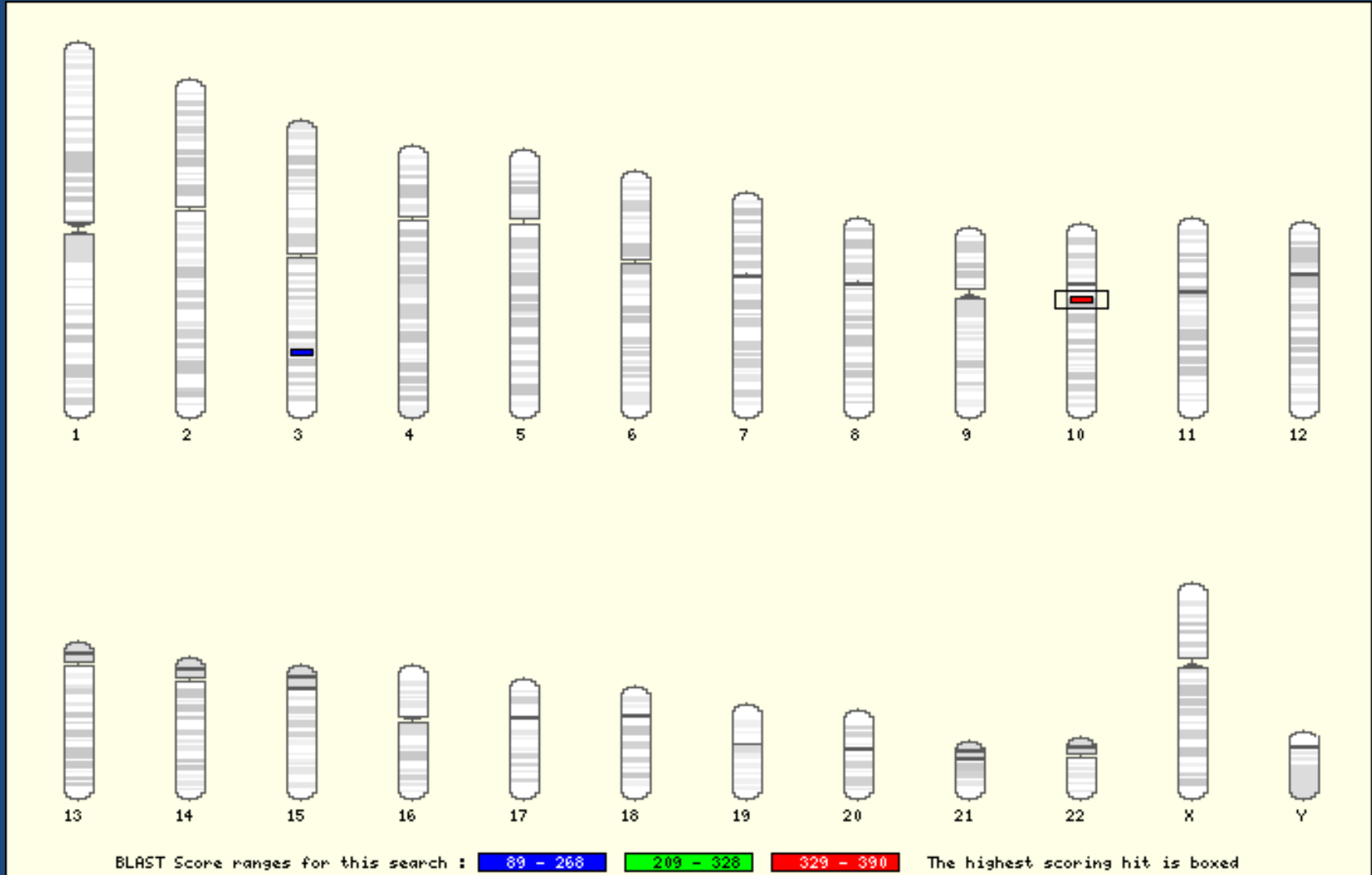
Sort results by  Filter type

Genetic Code  (blastx only)

**other options**  (not validated)



# Ensembl BLAST output includes an ideogram



[Program:](#)

[Database:](#)

[Alignments:](#)

- Porcine (*Sus scrofa*)
- African clawed frog (*Xenopus laevis*)
- Zebrafish (*Danio rerio*)
- 
- Arabidopsis thaliana
- Barley (*Hordeum vulgare*)
- Ice Plant (*Mesembryanthemum crystallinum*)
- Maize (*Zea mays*)
- Medicago truncatula
- Potato (*Solanum tuberosum*)

**Upload a file containing a sequence OR paste it into the textbox:**

(Note: If both are entered, the file will be ignored.)

Enter the name of the file containing a sequence in [FASTA](#) or raw format:

Enter your sequence in [FASTA](#) or raw format:

```
MKWVWALLLLAAWAAAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFVSDETGQMS
ATAKGRVRLNNDVCADMVGTFTDTEPAKFKMKYWGVAFLQKGNDDHWIVD TDYDTYAVQYSCRLLN
LDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEELCLARQYRLIVHNGYCDGRSERLL
```

Sequence identifier (for FASTA):

raw sequence entered:

- blosum62
- blosum100
- blosum30
- blosum35
- blosum40
- blosum45
- blosum50
- blosum55
- blosum60
- blosum65
- blosum70
- blosum75
- blosum80
- blosum85
- blosum90
- blosumn
- dayhoff
- gonnet
- identity
- match

**Options:**

[Matrix:](#)

[Filter:](#)

[Expect:](#)

[Cutoff:](#)

[Strand:](#)

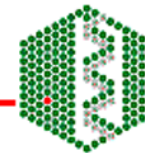
[Descriptions:](#)

[Wordlength](#) (for blastn only):

[Echofilter](#)

[Graphical Overview](#)

[Ignore Hypotheticals](#)



## WU-Blast2

[Help](#)[Tools](#)[EBI Home](#)[RUN BLAST](#)[RESET FORM](#)

<a href="#">YOUR EMAIL</a>	<a href="#">SEARCH TITLE</a>	<a href="#">RESULTS</a>	<a href="#">DATABASE</a>	<a href="#">PROGRAM</a>
<input type="text"/>	<input type="text" value="Sequence"/>	<input type="text" value="interactive"/>	<input type="text" value="swal"/>	<input type="text" value="WU-blastp"/>
<a href="#">MATRIX</a>	<a href="#">DNA STRAND</a>	<a href="#">EXP. THR</a>	<a href="#">FILTER</a>	<a href="#">VIEW FILTER</a>
<input type="text" value="blosum62"/>	<input type="text" value="none"/>	<input type="text" value="default"/>	<input type="text" value="none"/>	<input type="text" value="no"/>
<a href="#">HISTOGRAM</a>	<a href="#">STATS</a>	<a href="#">SORT</a>	<a href="#">SCORES</a>	<a href="#">ALIGNMENTS</a>
<input type="text" value="no"/>	<input type="text" value="sump"/>	<input type="text" value="pvalue"/>	<input type="text" value="default"/>	<input type="text" value="default"/>

[Enter or Paste](#) a  [Sequence](#) in any format:

[Upload a file:](#)

[Browse...](#)[RUN BLAST](#)[RESET FORM](#)

This document was last modified on : Thursday, July 05, 2001 10:52:03

Comments or suggestions [support@ebi.ac.uk](mailto:support@ebi.ac.uk)

© EBI 2000

If you plan to use these services during a course please contact us using the email above.

# BLAST-related tools for genomic DNA

---

Recently developed tools include:

- MegaBLAST at NCBI.
- BLAT (BLAST-like alignment tool). BLAT parses an entire genomic DNA database into words (11mers), then searches them against a query. Thus it is a mirror image of the BLAST strategy. See <http://genome.ucsc.edu>
- SSAHA at Ensembl uses a similar strategy as BLAT. See <http://www.ensembl.org>

To access BLAT, visit <http://genome.ucsc.edu>

**UCSC Genome Bioinformatics**

Genomes - Gene Sorter - Blat - PCR - Tables - FAQ - Help

**Genome Browser**

Gene Sorter

**Blat**

In Silico PCR

Table Browser

Utilities

Downloads

Release Log

Custom Tracks

ENCODE

**About the UCSC Genome Bioinformatics Site**

This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also shows the CFTR (cystic fibrosis) region in 13 species and provides a portal to the ENCODE project.

encourage you to explore these sequences with our tools. The Genome Browser zooms and scrolls over chromosomes, showing the work of annotators worldwide. The Gene Sorter shows expression, homology and other information on groups of genes that can be related in many ways. Blat quickly maps your sequence to the genome. The Table Browser provides convenient access to the underlying database.

**News** News Archives ►

**10 September 2004 - Tetraodon Genome Assembly in Genome Browser**

The Genoscope v7 *Tetraodon nigroviridis* genome assembly is now available in the UCSC Genome Browser and Blat server. This assembly, UCSC version tetNig1 dated Feb. 2004, is the result of a collaboration between [Genoscope](#) and the [Broad Institute](#) of MIT and Harvard.

The v7 assembly was constructed using the whole genome shotgun (WGS) approach, resulting in a sequence coverage of about 7.9X. The assembly contains 45,609 contigs and 25,773 scaffolds generated by the Arachne program and covers more than 90% of the genome.

“BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 20 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein blat on land vertebrates.” --BLAT website

## Human BLAT Search

# BLAT Search Genome

Genome:  Assembly:  Query type:  Sort output:  Output type:

Paste in a query sequence to find its location in the the genome. Multiple sequences may be searched at once if separated by a line starting with > followed by the sequence name.

```
>gi|8400727|ref|NM_006744.2| Homo sapiens retinol binding protein 4, plasma
(RBP4), mRNA
CGCTCGCCTCCCTCGCTCCACGCGCGCCCGGACGCGCGGCCAGGCTTGC GCGTGGTTCCCTCCCGGTG
GGCGGATTCTCTGGCAAGATGAAAGTGGGTGTGGGCGCTCTTGCTGTTGGCGCGTGGGCAGCGGCCGAGC
GCGACTGCCGAGTGAGCAGCTTCCGAGTCAAGGAGAACTTCGACAAGGCTCGCTTCTCTGGGACCGCTA
CGCCATGGCCAAGAAGGACCCCGAGGGCCTCTTCTGCAGGACAACATCGTTCGCGGAGTTCTCC
GAGACCGGCCAGATGAGCGCCACAGCCAAGGGCCGAGTCCGTCTTTTGAATAACTGGGACGTGTGGCGAG
ACATGGTGGGCACCTTCACAGACACCGAGGACCTGCCAAGTTCAAGATGAAAGTACTGGGGCGTAGCCTC
CTTTCTGCAGAAAGGAAATGATGACCACTGGATCGTCGACACAGACTACGACACGTATGCCGTACAGTAC
TCCTGCCGCTCCTGAACTCGATGGCACCTGTGCTGACAGCTACTCCTTCGTGTTTTCCCGGGACCCCA
ACGGCCTGCCCCAGAAAGCGCAGAAAGATTGTAAGGCAGCGGCAGGAGGAGCTGTGCTGGCCAGGCAGTA
CAGGCTGATCGTCCACAAACGGTTACTGCGATGGCAGATCAGAAAAGAAACCTTTTGTAGCAATATCAAGAA
TCTAGTTTCATCTGAGAACTTCTGATTAGCTCTCAGTCTTCAGCTCTATTTATCTTAGGAGTTTAATTG
CCCTTCTCTCCCATCTCCCTCAGTCCCATAAAACCTTCATTACACATAAAGATACACGTGGGGGTCA
```

Paste DNA or protein sequence here in the FASTA format

Rather than pasting a sequence, you can choose to upload a text file containing the sequence.

Upload sequence:

Only DNA sequences of 25,000 or fewer bases and protein or translated sequence of 5000 or fewer letters will be processed. Up to 25 sequences can be submitted at the same time. The total limit for multiple sequence submissions is 50,000 bases or 12,500 letters.

## About BLAT

BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 22 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein blat on land vertebrates.

BLAT is not BLAST. DNA BLAT works by keeping an index of the entire genome in memory. The index consists of all non-overlapping 11-mers except for those heavily involved in repeats. The index takes up a bit less than a gigabyte of RAM. The genome itself is not kept in memory, allowing BLAT to deliver high performance on a reasonably priced Linux box. The index is used to find areas of probable homology, which are then loaded into memory for a detailed alignment.



# Multiple sequence alignment



```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCTVSGTSFDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPSD--IAVEWWSNG--
```

Like pairwise alignment BUT compare  $n$  sequences instead of 2

Rows represent individual sequences  
Columns represent 'same' position

May be gaps in some sequences

# MSA & Evolution

MSA can give you a picture of the forces that shape evolution!

- Important amino acids or nucleotides are not “allowed” to mutate
- Less important positions change more easily

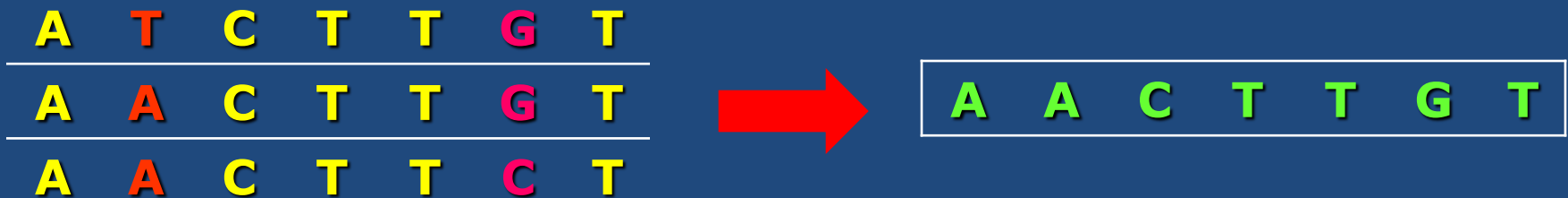
# Conserved positions

- Columns where all the sequences contain the same amino acids or nucleotides
- Important for the function or structure

```
VTISCTGSSSNIGAG-NHVKWYQQLPG  
VTISCTGSSSNIGS--ITVNWYQQLPG  
LRLSCTGSGFIFSS--YAMYWYQQAPG  
LSLTCTGSGTSFDD-QYYSTWYQQPPG
```

# Consensus Sequence

- The consensus sequence holds the most frequent character of the alignment at each column



# Profile

<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>T</b>	<b>G</b>	<b>T</b>	
<b>A</b>	<b>A</b>	<b>C</b>	<b>T</b>	<b>T</b>	<b>G</b>	<b>T</b>	→
<b>A</b>	<b>A</b>	<b>C</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>T</b>	→
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	
<b>A</b>	<b>1</b>	<b>0.67</b>	<b>0</b>	<b>0</b>	<b>.</b>	<b>.</b>	
<b>T</b>	<b>0</b>	<b>0.33</b>	<b>1</b>	<b>1</b>	<b>.</b>	<b>.</b>	
<b>C</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>.</b>	<b>.</b>	
<b>G</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>.</b>	<b>.</b>	

Profile =

PSSM – Position Specific Score Matrix

# Alignment methods

- Progressive alignment (Clustal)
- Iterative alignment (mafft, muscle)
- All methods today are an approximation strategy (**heuristic algorithm**), yield a possible alignment, but not necessarily the best one

# Progressive alignment

First step:



Compute the pairwise alignments for all against all (6 pairwise alignments) the similarities are stored in a table

	A	B	C	D
A				
B	11			
C	3	1		
D	2	2	10	

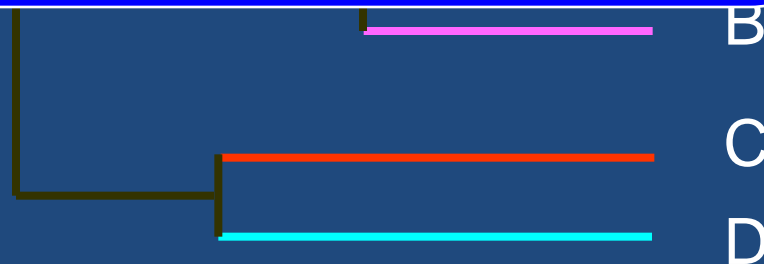
## Second step:

	A	B	C	D
A				
B	11			
C	3	1		
D				

cluster the sequences to create a tree  
(**guide tree**):

- Represents the order in which sequences are to be analyzed
- similar sequences are clustered together
- distant sequences are clustered together

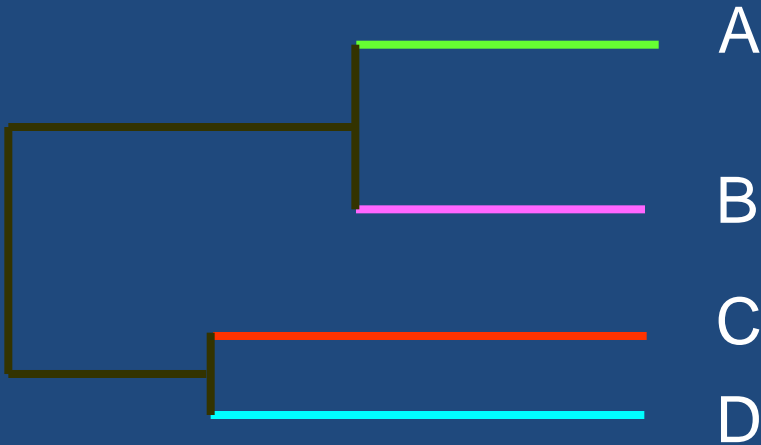
**The guide tree is imprecise and is NOT the tree which truly describes the relationship between the sequences!**





# Third step:

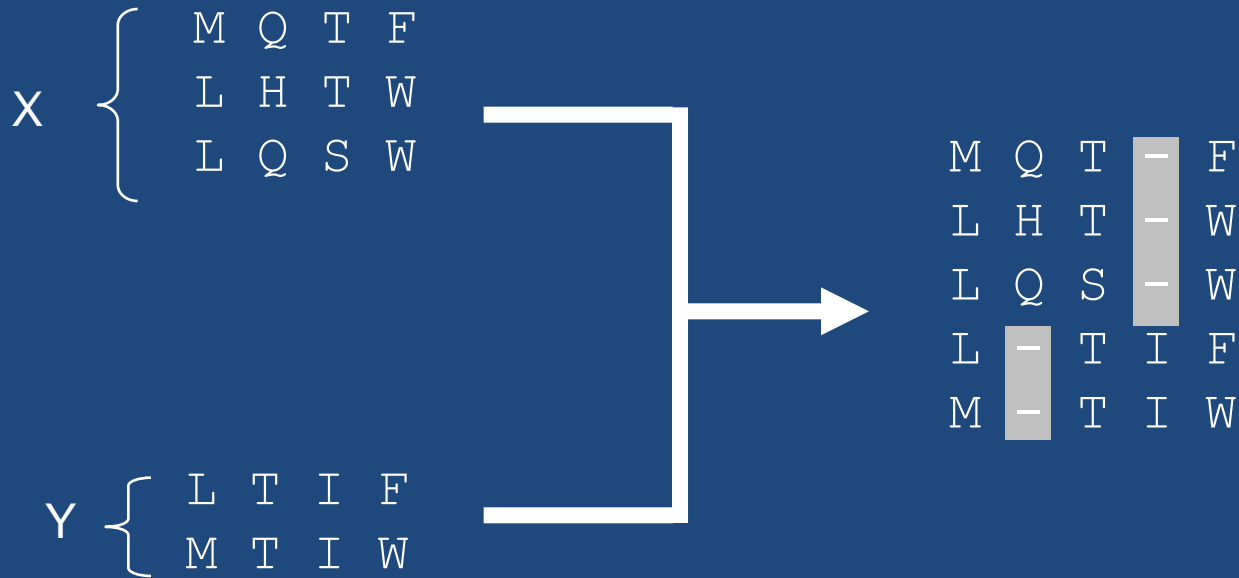
Align most similar pairs



Align the alignments as if each of them was a single sequence (replace with a single consensus sequence or use a profile)



# Alignment of alignments



# Iterative alignment

A  
B  
C  
D



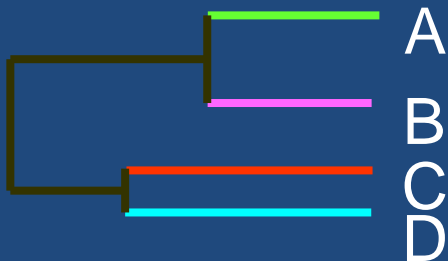
Pairwise distance  
table

	A	B	C	D
A				
B	11			
C	3	1		
D	2	2	10	

Iterate until the MSA  
doesn't change



Guide tree



MSA



# Online version

Multiple Sequence Alignment - CLUSTALW - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites History Mail Print

Address <http://align.genome.jp/> Go

## Multiple Sequence Alignment by CLUSTALW

CLUSTALW MAFFT PRRN

Help

**General Setting Parameters:**

Output Format:

Pairwise Alignment:  FAST/APPROXIMATE  SLOW/ACCURATE

Enter your sequences (with labels) below (copy & paste):  PROTEIN  DNA

Support Formats: FASTA (Pearson), NBRF/PIR, EMBL/Swiss Prot, GDE, CLUSTAL, and GCG/MSF

```
>CV523101_wheat
IARIFNTYGPRMCIDDDGRVVSNFVAQALRKEPLTVYGDGKQTRSFQYVSDLVEGLMRLME
GDHIGPFNLGNPGEFTMLELAKVVQDTIDPNARIEFRENTQDDPHKRKPDITRAKEQLGW
EPKIALRDGLPLMVTDFRKRIFGDQSAATATE
```

Or give the file name containing your query

More Detail Parameters...

**Pairwise Alignment Parameters:**

**For FAST/APPROXIMATE:**

K-tuple(word) size: , Window size: , Gap Penalty:

Number of Top Diagonals: , Scoring Method:

**For SLOW/ACCURATE:**

Gap Open Penalty: , Gap Extension Penalty:

Select Weight Matrix:

Done Internet

- >gi|115023|sp|P10425| MKKNTLLKVGLCVSLLGTTQFVSTISSVQAS  
QKVEQIVIKNETGTISISQLNKNVW VHTELGYFNGEAVPSNGLVLNTSKGL  
VLVDSSWDNKLTKELIEMVEKKFQKRVD VIITHAHADRIGGITALKERGIK  
AHSTALTAELAKKSGYEEPLGDLQVTNLKFGNTK VETFYPGKGHTEDNIV  
VWLPQYQILAGGCLVKSAEAKNLGNVADAYVNEWSTSIE NMLKRYRNINL  
VVPGHGKVGDKGLLLHTLDLLK >gi|115030|sp|P25910| MKTVFILIS  
MLFPVAVMAQKSVKISDDISITQLSDKVYTYVSLAEIEGWGMVPSNGM IVI  
NNHQAALLDTPINDAQTEMLVNWVTDSLHAKVTTFFIPNHHWGDCIGGLG  
YLQR KGVQSYANQMTIDLAKEKGLPVPEHGFTDSLTVSLDGMPLQCYLG  
GGHATDNIV VWLPTENILFGGCMLKDNQATSIGNISDADVTAWPKTLDK  
VKAKFPSARYVVPGH GDYGGTELIEHTKQIVNQYIESTSKP >gi|282554|  
pir||S25844  
MTVEVREVAEGVYAYEQAPGGWCVSNA GIVVGGDGALVVDTLSTIPRAR  
RLAEWV DKLAAGPGRTVVNTH  
FHGDHAFGNQVFAPGTRIIAHEDMRSAMVTTGLALTGLWP RVDWGEIEL  
RPPNVTFRDRLTLHVGERQVE  
LICVGAHTDHDV VVWLPEERVLFAGD VVMSGVTPFALFGSVAGTLAALD  
RLAELEPEVVGGHGPVAGP  
EVIDANRDYLRWV QRLAADAVDRRLTPLQAARRADLGAFAGLLDAERLVA  
NLHRAHEELLGGHV RDAM EI FAELVAYNGGQLPTCLA

# An output from ClustalW

## sequences have significant similarity

CLUSTAL W (1.82) multiple sequence alignment

```
gi|42542791|gb|AAH66228.1|   MSTAGKVIKCKAAVLWELKKPFSIEEVEVAPPKAHEVRIKMVAAGICRS- 49
gi|825623|emb|CAA39813.1|   MGTKGKVIKCKAAIAWEAGKPLCIEEVEVAPPKAHEVRIQIIATSLCHT- 49
gi|42738724|gb|AAS42652.1|  --MQNFVFRNPTKLIFGKGQ---LEQLKTEIPQFGKKVLLVYGGGSIKRN 45
      . *:: :: : : *::: * : : : . . . .
```

```
gi|42542791|gb|AAH66228.1|   ---DEHVSGNLV-TPLPVILGHEAAGIVESVGEVTTVKPG--DKVIPL 93
gi|825623|emb|CAA39813.1|   ---DASVIDSKFEGLAFPVIVGHEAAGIVESIGPGVTNVKPG--DKVIPL 94
gi|42738724|gb|AAS42652.1|   GIYDNVISILKDINAEVFELTGVEPNPRVSTVKKGIQICKDNGVEFILAV 95
      . * : : . . . : * * . *:: * : * .. ::::
```

```
gi|42542791|gb|AAH66228.1|   FTPQCGKCRICKNPESNYCLKN-DLGNPRG-----T 123
gi|825623|emb|CAA39813.1|   YAPLCRKCKFCLSPLTNLCGKISNLKSPASDQ-----QL 128
gi|42738724|gb|AAS42652.1|   GGGVIDCTKAI AAGSKYDGDVWDIVTKKAFASEALPFGTVLTLAATGSE 145
      . * . . :: . :: . . :: : : : : ::::
```

```
gi|42542791|gb|AAH66228.1|   LQDGTRRFTCSGKPIHHFVGVSTFSQYTVVDENAVAKIDAASPLEKVCLI 173
gi|825623|emb|CAA39813.1|   MEDKTSRFTCKGKPVYHFFGTSTFSQYTVVSDINLAKIDDDANLERVCLL 178
gi|42738724|gb|AAS42652.1|   MNAGSVITNWETNEKYGWGSPVTFPQFSILDPVHTASVPRDQTIYGMVDI 195
      :: : . . : : : . ** *::: . * : : : :
```

alcohol dehydrogenase, iron-containing [Bacillus cereus]

Class I alcohol dehydrogenase, gamma subunit [Homo sapiens]

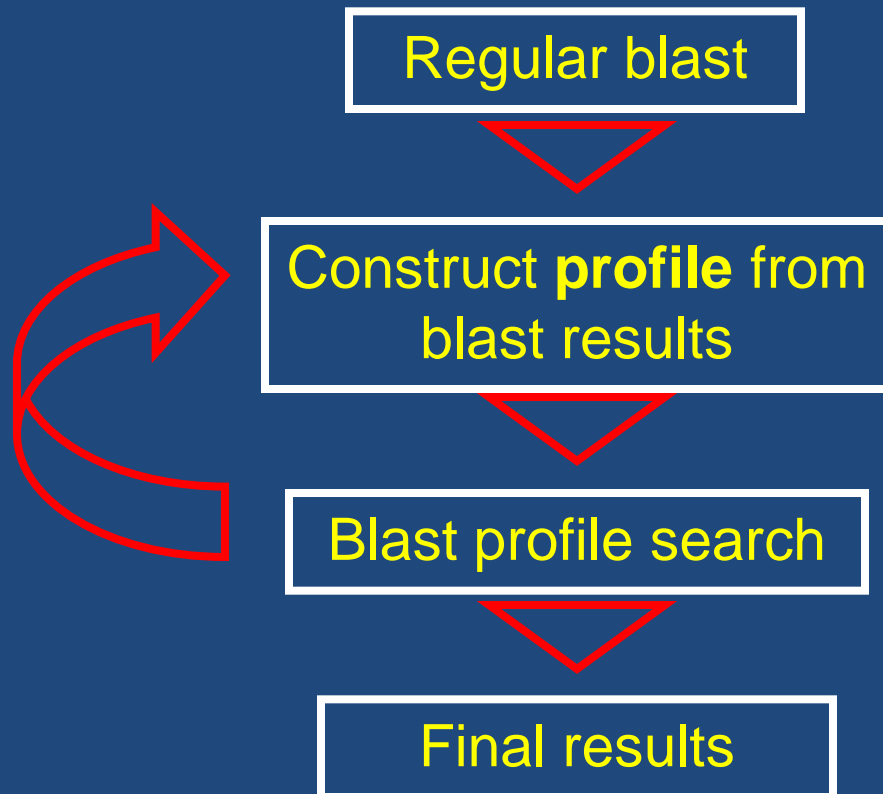
Different form of alcohol dehydrogenase [Homo sapiens]

# Searching for remote homologs

- Sometimes BLAST isn't enough.
- Large protein family, and BLAST only gives close members. We want more distant members
- PSI-BLAST

# PSI-BLAST

- Position Specific Iterated BLAST





# Position specific iterated BLAST: PSI-BLAST

---

The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a scoring matrix that is customized to your query.

# PSI-BLAST is performed in five steps

---

[1] Select a query and search it against a protein database

# PSI-BLAST is performed in five steps

---

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a “profile” or specialized position-specific scoring matrix (PSSM)

<a href="#">730496</a>	66	FTVDENGQMSATAKGRVRLFNWWDVCADMIGSFTDTEDEPAKFCKMRYWGVASFLQKGNDDH	125
<a href="#">200679</a>	63	FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDEPAKFCKMRYWGVASFLQKGNDDH	122
<a href="#">206589</a>	34	FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDEPAKFCKMRYWGVASFLQKGNDDH	93
<a href="#">2136812</a>	2	MSATAKGRVRLLSNWDVCADMVGTFTDTEDEPAKFCKMRYWGVASFLQKGNDDH	53
<a href="#">132408</a>	65	FKIEDNGKTTATAKGRVRILDKLELCANMVGTFTIETNDPAKYRMKYHGALAILERGLDDH	124
<a href="#">267584</a>	44	FSVDESGKVTATAHGRVILNWNWEMCANMFGTFEDTPDPAKFCKMRYWGAAASYLQKGNDDH	103
<a href="#">267585</a>	44	FSVDGSGKVTATAQGRVILNWNWEMCANMFGTFEDTPDPAKFCKMRYWGAAASYLQKGNDDH	103
<a href="#">8777608</a>	63	FTIHEDGAMTATAKGRVILNWNWEMCADMMATFETTPDPAKFRMRYWGAAASYLQKGNDDH	122
<a href="#">6687453</a>	60	FKVEEDGTMTATAIGRVILNWNWEMCANMFGTFEDTEDEPAKFCKMRYWGAAASYLQKGYDDH	119
<a href="#">10697027</a>	81	FKVQEDGTMTATATGRVILNWNWEMCANMFGTFEDTEEPARFKMRYWGAAASYLQKGYDDH	140
<a href="#">13645517</a>	1	MVGTFTDTEDEPAKFCKMRYWGVASFLQKGNDDH	32
<a href="#">13925316</a>	38	FSVDGSGKMTATAQGRVILNWNWEMCANMFGTFEDTPDPAKFCKMRYWGAAASYLQKGNDDH	97
<a href="#">131649</a>	65	YTVEEDGTMTASSKGRVKLFGFWVICADMAAQYTDPTTPAKMYMTYQGLASYLSSGGDNY	126

R,I,K

C

D,E,T

K,R,T

N,L,Y,G

# PSI-BLAST is performed in five steps

---

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a “profile” or specialized position-specific scoring matrix (PSSM)

[3] The PSSM is used as a query against the database

[4] PSI-BLAST estimates statistical significance (E values)

	✓	<a href="#">gi 6978523 ref NP_036909.1 </a>	apolipoprotein D [Rattus norvegicus]...	<a href="#">147</a>	4e-35
	✓	<a href="#">gi 1542847 dbj BAA13453.1 </a>	(D87752) alpha1-microglobulin/bikunin...	<a href="#">144</a>	6e-34
	✓	<a href="#">gi 619383 gb AAB32200.1 </a>	apolipoprotein D, apoD [human, plasma, ...	<a href="#">143</a>	8e-34
	✓	<a href="#">gi 5419892 emb CAB46489.1 </a>	(X02824) RBP (aa 101-172) [Homo sapiens]	<a href="#">139</a>	1e-32
	✓	<a href="#">gi 4502163 ref NP_001638.1 </a>	apolipoprotein D precursor [Homo sap...	<a href="#">138</a>	4e-32
	✓	<a href="#">gi 584763 sp P37153 APD_RABIT</a>	APOLIPOPROTEIN D PRECURSOR >gi 482...	<a href="#">134</a>	4e-31
	✓	<a href="#">gi 1703341 sp P51909 APD_CAVPO</a>	APOLIPOPROTEIN D PRECURSOR >gi 11...	<a href="#">133</a>	7e-31
	✓	<a href="#">gi 2895204 gb AAC02945.1 </a>	(AF025334) mutant retinol binding prot...	<a href="#">80</a>	9e-15
	✓	<a href="#">gi 1246096 gb AAB35919.1 </a>	(S80440) apolipoprotein D, apoD (C-ter...	<a href="#">77</a>	8e-14
	✓	<a href="#">gi 2895206 gb AAC02946.1 </a>	(AF025335) mutant retinol binding prot...	<a href="#">67</a>	8e-11
NEW	✓	<a href="#">gi 1346419 sp P49291 LAZA_SCHAM</a>	LAZARILLO PROTEIN PRECURSOR >gi ...	<a href="#">63</a>	1e-09
NEW	✓	<a href="#">gi 2506821 sp P00978 AMBP_BOVIN</a>	AMBP PROTEIN PRECURSOR [CONTAINS...	<a href="#">63</a>	2e-09
NEW	✓	<a href="#">gi 2497696 sp Q07456 AMBP_MOUSE</a>	AMBP PROTEIN PRECURSOR [CONTAINS...	<a href="#">63</a>	2e-09
NEW	✓	<a href="#">gi 6680684 ref NP_031469.1 </a>	alpha 1 microglobulin/bikunin [Mus m...	<a href="#">62</a>	2e-09
NEW	✓	<a href="#">gi 12836446 dbj BAB23659.1 </a>	(AK004907) putative [Mus musculus]	<a href="#">62</a>	3e-09
NEW	✓	<a href="#">gi 6978497 ref NP_037033.1 </a>	alpha-1 microglobulin/bikunin [Rattu...	<a href="#">62</a>	3e-09
NEW	✓	<a href="#">gi 2507586 sp P04366 AMBP_PIG</a>	AMBP PROTEIN PRECURSOR [CONTAINS: ...	<a href="#">61</a>	8e-09
NEW	✓	<a href="#">gi 1085207 pir  JC2556</a>	alpha-1-microglobulin/inter-alpha-trypsin...	<a href="#">60</a>	1e-08
NEW	✓	<a href="#">gi 2988354 dbj BAA25305.1 </a>	(AB006444) alpha-1-microglobulin/biku...	<a href="#">59</a>	2e-08
NEW	✓	<a href="#">gi 108233 pir  S13493</a>	alpha-1-microglobulin - pig	<a href="#">59</a>	2e-08
NEW	✓	<a href="#">gi 1882 emb CAA36306.1 </a>	(X52087) precursor codes for two protein...	<a href="#">59</a>	2e-08
NEW	✓	<a href="#">gi 9181923 gb AAF85707.1 AF276505_1</a>	(AF276505) neural Lazarillo ...	<a href="#">59</a>	3e-08
NEW	✓	<a href="#">gi 7296083 gb AAF51378.1 </a>	(AE003586) NLaz gene product [Drosophi...	<a href="#">58</a>	3e-08
NEW	✓	<a href="#">gi 117330 sp P80007 CRA2_HOMGA</a>	CRUSTACYANIN A2 SUBUNIT >gi 10275...	<a href="#">57</a>	8e-08
NEW	✓	<a href="#">gi 2497695 sp Q60559 AMBP_MESAU</a>	AMBP PROTEIN PRECURSOR [CONTAINS...	<a href="#">57</a>	1e-07
NEW	✓	<a href="#">gi 102968 pir  S22400</a>	insecticyanin A - tobacco hornworm >gi 971...	<a href="#">56</a>	1e-07
NEW	✓	<a href="#">gi 4502067 ref NP_001624.1 </a>	alpha-1-microglobulin/bikunin precu...	<a href="#">56</a>	2e-07
NEW	✓	<a href="#">gi 1146408 gb AAA85089.1 </a>	(L41641) gallerin [Galleria mellonella]	<a href="#">56</a>	2e-07
NEW	✓	<a href="#">gi 2497694 sp Q62577 AMBP_MERUN</a>	AMBP PROTEIN PRECURSOR [CONTAINS...	<a href="#">55</a>	3e-07
NEW	✓	<a href="#">gi 1213589 dbj BAA12075.1 </a>	(D83712) Prostaglandin D Synthase [Xe...	<a href="#">54</a>	5e-07
	✓	<a href="#">gi 539717 pir  A61233</a>	retinol-binding protein - cat (fragment)	<a href="#">54</a>	8e-07
NEW	✓	<a href="#">gi 266472 sp Q01584 LIPO_BUFMA</a>	LIPOCALIN PRECURSOR >gi 104284 pi...	<a href="#">53</a>	1e-06
	✓	<a href="#">gi 265042 gb AAB25283.1 </a>	retinol-binding protein, RBP (N-termina...	<a href="#">52</a>	3e-06
NEW	✓	<a href="#">gi 1079295 pir  S52354</a>	gene cpl-1 protein - African clawed frog ...	<a href="#">52</a>	3e-06
NEW	✓	<a href="#">gi 732003 sp P39281 BLC_ECOLI</a>	OUTER MEMBRANE LIPOPROTEIN BLC PRE...	<a href="#">51</a>	9e-06

# PSI-BLAST is performed in five steps

---

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a “profile” or specialized position-specific scoring matrix (PSSM)

[3] The PSSM is used as a query against the database

[4] PSI-BLAST estimates statistical significance (E values)

[5] Repeat steps [3] and [4] iteratively, typically 5 times. At each new search, a new profile is used as the query.

# Results of a PSI-BLAST search

---

		# hits	
<u>Iteration</u>	<u># hits</u>		<u>&gt; threshold</u>
1	104	49	
2	173	96	
3	236	178	
4	301	240	
5	344	283	
6	342	298	
7	378	310	
8	382	320	



# PSI-BLAST alignment of RBP and $\beta$ -lactoglobulin: iteration 1

Score = 46.2 bits (108), Expect = 2e-04

Identities = 40/150 (26%), Positives = 70/150 (46%), Gaps = 37/150 (24%)

Query: 27 VKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDETGQMSATAKGRVRLNNDVC 86  
V+ENFD ++ G WY + +K P + I A +S+ E G + K ++

Sbjct: 33 VQENFDVKKYLGRWYEI-EKIPASFEKGNCIQANYSLMENGNIENVLNK-----ELS 82

Query: 87 ADMVGTF-----TDTEDPAKFKMKYWGVASFLQKGNDDHWIVDTDYDTYAVQYSCR 137  
D GT ++ +PAK +++++ + +WI+ TDY+ YA+ YSC

Sbjct: 83 PD--GTMNQVKGEAKQSNVSEPAKLEVQFFPLMP-----PAPYWILATDYENYALVYSCT 135

Query: 138 ----LLNLDGTCADSYFVFSRDPNGLPPE 163  
L ++D + ++ R+P LPPE

Sbjct: 136 TFFWLFHVD-----FFWILGRNPY-LPPE 158

## PSI-BLAST alignment of RBP and $\beta$ -lactoglobulin: iteration 2

Score = 140 bits (353), Expect = 1e-32

Identities = 45/176 (25%), Positives = 78/176 (43%), Gaps = 33/176 (18%)

```
Query: 4   VWALLLLAAWAAAERDCRVSSF-----RVKENFDKARFSGTWYAMAKKDPEGLFLQD 55
          V L+ LA A      + +F          V+ENFD ++ G WY + +K P      +
Sbjct: 2   VTMLMFLATLAGLFTTAKGQNFHLGKCPSPVQENFDVKKYLGRWYEI-EKIPASFEKGN 60

Query: 56  NIVAEFSVDETGQMSATAKGRVRLNNDVDCADMV---GTFTDTEDPAKFKMKYWGVASF 112
          I A +S+ E G +      K      + D      + V      ++ +PAK +++++ +
Sbjct: 61  CIQANYSLMENGNIIEVLNKKEL-----SPDGTMNQVKGEAKQSNVSEPAKLEVQFFPL--- 112

Query: 113 LQKGNDDHWIVDTDYDTYAVQYSCR----LLNLDGTCADSYSEVFSRDPNGLPPEA 164
          +WI+ TDY+ YA+ YSC      L ++D      + ++ R+P LPPE
Sbjct: 113 --MPPAPYWILATDYENYALVYSCTTFFWLFHVD-----FFWILGRNPY-LPPET 159
```

## PSI-BLAST alignment of RBP and $\beta$ -lactoglobulin: iteration 3

Score = 159 bits (404), Expect = 1e-38

Identities = 41/170 (24%), Positives = 69/170 (40%), Gaps = 19/170 (11%)

Query: 3    WWALLLLLAWAAAERD-----CRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ 54

          V L+ LA A                   + S V+ENFD ++ G WY + K

Sbjct: 1    MVTMLMFLATLAGLFTTAKGQNFHLGKCPSPVQENFDVKKYLGRWYEIEKIPASFE-KG 59

Query: 55    DNIVAEFSVDETGQMSATAKGRVRLNNDVDCADMVGTFTDTEPAKFKMKYWGVASFLQ 114

          + I A +S+ E G +           K           V +           ++ +PAK +++++ +

Sbjct: 60    NCIQANYSLMENGNIIEVLNKELSPDGTMNQVKGE--AKQSNVSEPAKLEVQFFPL----- 112

Query: 115    KGNDDHWIVD TDYDTYAVQYSCRLNLDGTCADSYFVFSRDPNGLPPEA 164

          +WI+ TDY+ YA+ YSC                   + ++ R+P LPPE

Sbjct: 113    MPPAPYWILATDYENYALVYSCTTFFWL--FHVDFFWILGRNPY-LPPET 159

1

Score = 46.2 bits (108), Expect = 2e-04  
Identities = 40/150 (26%), Positives = 70/150 (46%), Gaps = 37/150 (24%)

Query: 27 VKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDETGQMSATAKGRVRLNNDVDC 86  
V+ENFD ++ G WY + +K P + I A +S+ E G + K ++

Sbjct: 33 VQENFDVKKYLGRWYEI-EKIPASFEKGNCIQANYSLMENGNI EVLNK-----ELS 82

Query: 87 ADMVGTF-----TDTEDPAKFKMKYWGVASFLQKGNDDHWIVD TDYDTYAVQYSCR 137  
D GT ++ +PAK +++++ + +WI+ TDY+ YA+ YSC

Sbjct: 83 PD--GTMNQVKGEAKQSNVSEPAKLEVQFFPLMP-----PAPYWILATDYENYALVYSCT 135

Query: 138 ----LLNLDGTCADSYFVFSRDPNGLPPE 163

L ++D + ++ R+P LPPE

Sbjct: 136 TFFWLFHVD-----FFWILGRNPY-LPPE 158

3

Score = 159 bits (404), Expect = 1e-38  
Identities = 41/170 (24%), Positives = 69/170 (40%), Gaps = 19/170 (11%)

Query: 3 WVWALLLLLAWAAAERD-----CRVSSFVRVKENFDKARFSGTWYAMAKKDPEGLFLQ 54  
V L+ LA A + S V+ENFD ++ G WY + K

Sbjct: 1 MVTMLMFLATLAGLFTTAKGQNFHLGKCPSPVQENFDVKKYLGRWYEIEKIPASFE-KG 59

Query: 55 DNIVAEFSVDETGQMSATAKGRVRLNNDVDCADMVGTFDTEDPAKFKMKYWGVASFLQ 114  
+ I A +S+ E G + K V + ++ +PAK +++++ +

Sbjct: 60 NCIQANYSLMENGNI EVLNKELSPDGTMNQVKGE--AKQSNVSEPAKLEVQFFPL----- 112

Query: 115 KGNDDHWIVD TDYDTYAVQYSCRLLNLDGTCADSYFVFSRDPNGLPPEA 164  
+WI+ TDY+ YA+ YSC + ++ R+P LPPE

Sbjct: 113 MPPAPYWILATDYENYALVYSCTTFFWL--FHVDFFWILGRNPY-LPPET 159