

# **EBI web resources III: Web-based tools in Europe (EBI, ExPASy, DTU)**

Yanbin Yin  
Spring 2013

# Homework assignment 3

1. Download <http://cys.bios.niu.edu/yyin/teach/PBB/purdue.cellwall.list.lignin.fa> to your computer
2. Select a C3H protein and a F5H protein from the above file and calculate the sequence identity between them using the Water server at EBI.
3. Perform a multiple sequence alignment using MAFFT with all FASTA sequences in the file
4. Built a phylogeny with the alignment using the "**A la Carte**" mode at <http://www.phylogeny.fr/>
5. Build another phylogeny starting from the unaligned sequences using the "**one-click**" mode at <http://www.phylogeny.fr/>; if you encounter any error reports, try to figure out why and how to solve it (hint: skip the Gblocks step).

Write a report (in word or ppt) to include all the operations, screen shots and the final phylogenies from step 3 and 4.

Due on Feb 21 (send by email)

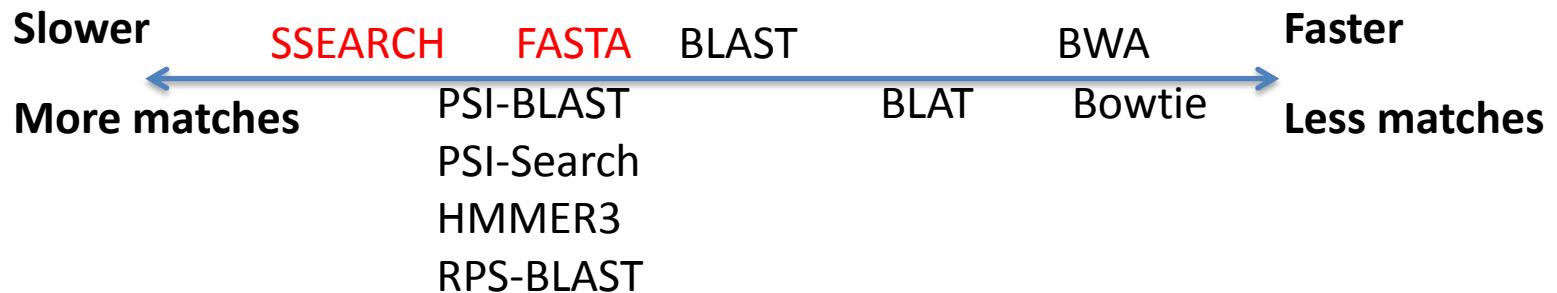
Office hour:

Tue, Thu and Fri 2-4pm, MO325A  
Or email: yyin@niu.edu

# Outline

- Hands on Practice!

# Pairwise alignment tools



- Tools Home
- Tools A-Z
- ID Mapping
- Literature
- Microarray Analysis
- Protein Functional Analysis
- Proteomic Services
- Sequence Analysis
- Similarity & Homology
- Structural Analysis
- Tools - Miscellaneous
- Web Services

- Databases
- Downloads

EBI &gt; Tools

<http://www.ebi.ac.uk/Tools/>

## Tools at the EBI

We provide a comprehensive range of bioinformatics tools.

This page shows a selection of those that are used most frequently. These include tools for the analysis and comparison of nucleotide and protein sequences, data from functional genomics experiments, text mining of the scientific literature and tools for determination and visualisation of macromolecular structures.

All these tools can be accessed over the web and most provide [Web Services](#) interfaces using SOAP or REST APIs.

### Nucleotide and Protein sequence searching

#### Nucleotide sequence searches

The sequence databases that can be searched with the tools outlined below include EMBL-Bank, Coding Sequences, immunoglobulins and High throughput cDNA:

- [ENA Search](#)
- [BLAST Nucleotide](#)
- [Fasta Genomes](#)
- [Search Genomes](#)

#### Protein sequence searches

The protein sequence databases available to search below include UniProtKB, sequences derived from macro molecular structures, immunoglobulins and sequences from patents:

- [BLAST Protein](#)
- [PSI-Search](#)
- [Fasta Proteomes](#)
- [Search Proteomes](#)

### Multiple Sequence Alignment

Alignment of three or more sequences to identify regions of conservation which may indicate functional constraints and infer evolutionary relationships.

- [Clustal Omega](#)

### Pairwise Sequence Alignments

Alignment of two sequences to identify regions where the sequence is conserved and conversely regions where the sequence is not conserved.

- [Needle](#)
- [J Align](#)

## Sequence Similarity Searching

**Sequence Similarity Searching** is a method of searching sequence databases by using alignment to a query sequence. By statistically assessing how well database and query sequences match one can infer homology and transfer information to the query sequence.

The tools can be launched with different form pre-sets using the buttons - these can be changed on the tool page as well.

### BLAST

**NCBI BLAST** ⓘ NCBI BLAST (blastall) is the most commonly used sequence similarity search tool. It uses heuristics to perform fast **local** alignment searches.



**WU-BLAST** ⓘ WU-BLAST is similar to NCBI BLAST but combines multiple parameter options into a simpler 'sensitivity' setting.



**PSI-BLAST** ⓘ PSI-BLAST allows users to construct and perform a BLAST search with a custom, position-specific, scoring matrix which can help find distant evolutionary relationships. PHI-BLAST functionality is also available to restrict results using patterns.



### FASTA

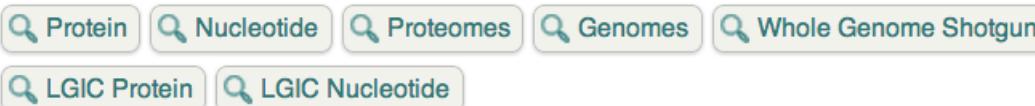
#### FASTA ⓘ

FASTA is another commonly used sequence similarity search tool which uses heuristics for fast **local** alignment searching.



#### SSEARCH ⓘ

SSEARCH is an optimal (as opposed to heuristics-based) **local** alignment search tool using the Smith-Waterman algorithm. Optimal searches guarantee you find the best alignment score for your given parameters.



#### PSI-Search ⓘ

PSI-Search combines the sensitivity of the Smith-Waterman search algorithm (SSEARCH)

**Help**

- [FASTA website](#)
- [Similar Applications](#)
- [Programmatic Access](#)
- Download**

- Database Information**
  - [UniProt](#)
  - [UniParc](#)

Tfastx: allow frame shift between codons

Tfasty: also allow frame shift within codons

Tolerate sequence errors

Good for finding pseudogenes

## FASTA/SSEARCH - Genomes Similarity Search

This tool provides sequence similarity searching against complete genomes databases using the FASTA suite of programs. FASTA provides a heuristic search with a nucleotide query. TFASTX and TFASTY translate the DNA database for searching with a protein query. Optimal searches are available with SSEARCH (local).

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Select your databases

##### GENOME DATABASES

1 Databank Selected	<input type="button" value="Clear Selection"/>
<b>Eukaryota</b> <ul style="list-style-type: none"> <li><input type="checkbox"/> Anopheles gambiae str. PEST</li> <li><input checked="" type="checkbox"/> Arabidopsis thaliana</li> <li><input type="checkbox"/> Ashbya gossypii ATCC 10895</li> <li><input type="checkbox"/> Aspergillus fumigatus Af293</li> </ul>	

##### OTHER TYPES

###### General

- Protein Databases
- Nucleotide Databases

###### Specialised

- Proteomes Databases
- WGS Databases
- LGIC Protein Databases

#### STEP 2 - Enter your input sequence

Enter or paste a **PROTEIN** sequence in any supported format:

```
>AT5G22740.1|AT5G22740.1|csIA
MDGVSPKFVLPETFDGVRMEITCQLGMIWELVKAPVIVPLLQLAVYICLLMSVMLLCERVYMGIVIVLVKLFWKPKDKRY
KFEPIHDDEELGSSNFPVVVLVQIPMFNEREVYKLSIGAACGLSWPSDRRLVIQVLDDSTDPTVKQMVEVECQRWASKGINIR
YQIRENRVGYKAGALKEGLKRSYVKHCEYYVIFDADFQPEPDFLRRSIPFLMHNPNALVQARWRFVNSDECLTRMQE
```

or Upload a file:  No file chosen

#### STEP 3 - Set your parameters

##### PROGRAM

TFASTX	▼
--------	---

EBI > Tools > Sequence Similarity Searching > FASTA

## FASTA Results

[Summary Table](#) [Tool Output](#) [Visual Output](#) [Submission Details](#) [Submit Another Job](#)

### Alignments

**Selection:** [Show Annotations](#) [Hide Annotations](#) [Show Alignments](#) [Hide Alignments](#)

[Download](#) in [fasta](#) format

[Clear Selection](#) [Select All](#) [Invert Selection](#)

Align.	DB:ID	Source	Length	Score	Identities	Positives	E()
<input checked="" type="checkbox"/> 1	EM_PLN:AB006699	STD:Arabidopsis thaliana genomic DNA, chromosome 5, P1 clone:MDJ22. <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Literature ► Ontologies ► Protein sequences	77363	763	55.2	55.2	1.2E-44
<input checked="" type="checkbox"/> 2	EM_PLN:AP000415	STD:Arabidopsis thaliana genomic DNA, chromosome 3, P1 clone: MIG10. <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Literature ► Ontologies ► Protein sequences	31712	821	47.3	72.5	1.7E-33
<input checked="" type="checkbox"/> 3	EM_PLN:AC005990	STD:Arabidopsis thaliana chromosome 1 BAC F5O8 sequence, complete sequence. <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Macromolecular structures ► Ontologies ► Protein sequences	99923	1271	56.3	73.2	3.7E-32
<input checked="" type="checkbox"/> 4	EM_PLN:AC007945	STD:Genomic sequence for Arabidopsis thaliana BAC F28C11 from chromosome I, complete sequence. <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Ontologies ► Protein sequences	80472	1271	56.3	73.2	3.7E-32
<input checked="" type="checkbox"/> 5	EM_PLN:AL163832	STD:Arabidopsis thaliana DNA chromosome 3, BAC clone F27K19 <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Ontologies ► Protein sequences	95111	1041	44.9	64.2	2.2E-29
<input checked="" type="checkbox"/> 6	EM_PLN:AL162506	STD:Arabidopsis thaliana DNA chromosome 5, BAC clone F17C15 (ESSA project) <i>Cross-references and related information in:</i> ► Nucleotide sequences ► Genomes ► Protein families ► Macromolecular structures ► Ontologies ► Protein sequences	102897	814	58.5	74.0	7.2E-26
<input checked="" type="checkbox"/> 7	EM_PLN:AB005235	STD:Arabidopsis thaliana genomic DNA, chromosome 5, P1 clone:MED24.	75475	814	58.5	74.0	7.2E-

In the alignment, look for

/

\

\*

[Download](#) ▾ [GenBank](#) [Graphics](#) Sort by: E value ▾

▼ N

## Arabidopsis thaliana chromosome 1, complete sequence

Sequence ID: [ref|NC\\_003070.9|](#) Length: 30427671 Number of Matches: 13

BLAST: shorter alignment

Range 1: 8335055 to 8335714 [GenBank](#) [Graphics](#)

▼ Next Match

▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps	Frame
303 bits(777)	3e-149	Compositional matrix adjust.	140/220(64%)	169/220(76%)	28/220(12%)	+2

Features: [putative mannan synthase 3](#)  
[putative mannan synthase 3](#)

Query	228	VNSDECLLTRMQEMSLDYHFTVEQEVGSSTHAFFGFNGTAGIWRIAALNEAGGWKDRRTV	287
Sbjct	8335055	VNANECLMTRMQEMSLNYHFVAEQESGSSIHAFFGFNGTAGVWRIAALNEAGGWKDRRTV	8335234
Query	288	EDMDLAVRASLRGWKFLYLGDL-----QVKSELPSTF	319
Sbjct	8335235	EDMDLAVRA L GWKF+Y+ D+ QVK+ELPSTF	8335414
Query	320	RAFRFQQHRWSCGPANLFRKMVMEIVRNKKVRFWKKVYVIYSFFFVRKIIAHWVTFCFYC	379
Sbjct	8335415	+A+RFQQHRWSCGPANL+RKM MEI++NKKV WKK+Y+IY+FFF+RKI+ H TF FYC	8335594
Query	380	KAYRFQQHRWSCGPANLWRKMTMEILQNKVSAWKKLYIYNFFFIRKIVVHIFTVFY	
Sbjct	8335595	LILPTTVLFPELQVPKWATVYFPPTTITILNAIATPR*QHL 8335714	

Range 2: 8334112 to 8334624 [GenBank](#) [Graphics](#)

▼ Next Match

▲ Previous Match

▲ First Match

Score	Expect	Method	Identities	Positives	Gaps	Frame
137 bits(345)	3e-149	Compositional matrix adjust.	75/171(44%)	106/171(61%)	34/171(19%)	+1

Features: [putative mannan synthase 3](#)  
[putative mannan synthase 3](#)

Query	12	ETFDGV-RMEITGQLGMIWELVKAPVIVPLLQLAVYICLLMSVMLLCERVYMGIVIVLVK	70
Sbjct	8334112	+T DGV R I G++ IW+ + V +P+L+ V ICL+MS++L ERVYM IV+V VK	8334291
Query	71	DTTDGVVRSGIIGEIIYIWKQTRIFVFIPILKCLVTICLVMSSLFIERVYMSIVVVFK	
Sbjct	8334112	LFWKKPDKRYKFEPIHDDE-ELGSSNFPVVLVQIPMFNEREVY-----	112

## FASTA/SSEARCH - Proteome Similarity Search

This tool provides sequence similarity searching against complete **proteomes** databases using the FASTA suite of programs. FASTA provides a heuristic search with a protein query. FASTX and FASTY translate a DNA query. Optimal searches are available with SSEARCH (local).

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Select your databases

##### PROTEIN DATABASES

0 Databank Selected

Clear Selection

- ▶ Eukaryota
- ▶ Archaea
- ▶ Bacteria
- ▶ Phage

##### OTHER TYPES

###### General

- Nucleotide Databases
- Protein Databases

###### Specialised

- Genomes Databases
- WGS Databases
- LGIC Protein Databases

#### STEP 2 - Enter your input sequence

Enter or paste a **PROTEIN** sequence in any supported format:

or upload a file: **Choose File** No file chosen

#### STEP 3 - Set your parameters

##### PROGRAM

SSEARCH ▲

- Tools Home
- Tools A-Z
- ID Mapping
- Literature
- Microarray Analysis
- Protein Functional Analysis
- Proteomic Services
- Sequence Analysis
- Similarity & Homology
- Structural Analysis
- Tools - Miscellaneous
- Web Services

- Databases
- Downloads

EBI &gt; Tools

<http://www.ebi.ac.uk/Tools/>

## Tools at the EBI

We provide a comprehensive range of bioinformatics tools.

This page shows a selection of those that are used most frequently. These include tools for the analysis and comparison of nucleotide and protein sequences, data from functional genomics experiments, text mining of the scientific literature and tools for determination and visualisation of macromolecular structures.

All these tools can be accessed over the web and most provide [Web Services](#) interfaces using SOAP or REST APIs.

### Nucleotide and Protein sequence searching

#### **Nucleotide sequence searches**

The sequence databases that can be searched with the tools outlined below include EMBL-Bank, Coding Sequences, immunoglobulins and High throughput cDNA:

- [ENA Search](#)
- [BLAST Nucleotide](#)
- [Fasta Genomes](#)
- [Search Genomes](#)

#### **Protein sequence searches**

The protein sequence databases available to search below include UniProtKB, sequences derived from macro molecular structures, immunoglobulins and sequences from patents:

- [BLAST Protein](#)
- [PSI-Search](#)
- [Fasta Proteomes](#)
- [Search Proteomes](#)

### Multiple Sequence Alignment

Alignment of three or more sequences to identify regions of conservation which may indicate functional constraints and infer evolutionary relationships.

- [Clustal Omega](#)

### Pairwise Sequence Alignments

Alignment of two sequences to identify regions where the sequence is conserved and conversely regions where the sequence is not conserved.

- [Needle](#)
- [J Align](#)

## Pairwise Sequence Alignment

**Pairwise Sequence Alignment** is used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences (protein or nucleic acid).

By contrast, **Multiple Sequence Alignment (MSA)** is the alignment of three or more biological sequences of similar length. From the output of MSA applications, homology can be inferred and the evolutionary relationship between the sequences studied.

### Global Alignment

Global alignment tools create an end-to-end alignment of the sequences to be aligned. There are separate forms for protein or nucleotide sequences.

**Needle**   
**(EMBOSS)**

EMBOSS Needle creates an optimal global alignment of two sequences using the Needleman-Wunsch algorithm.

**Needleman-wunsch algorithm**



**Stretcher**   
**(EMBOSS)**

EMBOSS Stretcher uses a modification of the Needleman-Wunsch algorithm that allows larger sequences to be globally aligned.



### Local Alignment

Local alignment tools find one, or more, alignments describing the most similar region(s) within the sequences to be aligned. There are separate forms for protein or nucleotide sequences.

**Water**   
**(EMBOSS)**

EMBOSS Water uses the Smith-Waterman algorithm (modified for speed enhancements) to calculate the local alignment of two sequences.



**Matcher**   
**(EMBOSS)**

EMBOSS Matcher identifies local similarities between two sequences using a rigorous algorithm based on the LALIGN application.



**LALIGN**

LALIGN finds internal duplications by calculating non-intersecting local alignments of protein or DNA sequences.



### Genomic Alignment

Help

EMBOSS website

Programmatic Access

Download

Nucleotide form

Similar Applications

CsIA: 539 aa

CesA: 1089 aa

EBI > Tools > Pairwise Sequence Alignment > EMBOSS Needle

## EMBOSS Needle - Pairwise Sequence Alignment (PROTEIN)

EMBOSS Needle reads two input sequences and writes their optimal global sequence alignment to file.

This is the form for protein sequences. Please go to the [nucleotide form](#) if you wish to align DNA or RNA sequences.

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Enter your protein sequences

Enter or paste your first **protein** sequence in any [supported](#) format:

```
>AT5G22740.1|AT5G22740.1|csIA
MDGVSPKFVLPETFDGVRMEITGQLGMIWELVKAPVIVPLLQLAVYICLLMSVMILLCERVYMGIVIVLVKLFKKPDKRY
KFEPIHDDEELGSSNFPVVVLVQIPMFNEREVYKLSIGAACGLSWPSDRLVIQVLDDSTDPTVKQMVEVECQRWASKGINIR
YQIRENRVCYKAGALKEGLKR SYVKHCEYVVFIDFADQPEPDPFLRRSIPFLMHNPNIALVQARWRFVNSDECCLTRMQE
MSLDYHFTVEQEVSSTHAFFGFNGTAGIWRIA AINEAGGWKDKRTTVEDMDLAVRASLRGWKFLYLGDLQVKSELPST
FRAFRFQQHRWSCGPANLFRKMVMEIVRNKKVRFWKKVYVIYSFFFVRKIIAHWTFCFYCVVPLTILVPEVKVPIWGSV
YIPSIITLNSVGTPRSIIHLFYWILFENVMSLHRTKATLIGLFEAGRANEVVTAKLGSQSAKGNTKGIKRFPRIFKLPDRL
```

Or, upload a file:  No file chosen

AND

Enter or paste your second **protein** sequence in any [supported](#) format:

```
>AT2G21770.1|AT2G21770.1|cesA
MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTNGEPIACNECAFPTCRPCYEYERREGNQA
CPQCGTRYKRIKGSPRVEGDEEDDDIDDLHEFYGMDPEHVTAA LYMLNTGRGTDEVSHLYSASPGSEVPLLTYCD
EDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSFAIHPMPVQKDLTVCYGSVAWKDRMEVWKKQQIEKLQVVKN
ERVNDGDGDGFIVDELDDPGLPMMDEGRQPLSRKLPIRRSSRINPYRMLIFCRLA ILGLFFHYRILHPVNDAGLWLTSVICE
IWFAVSWILDQFPKWYPIERETYLDRLSLRYEKEGKPSELAPVDVFVSTVDPLKEPPLITANTVLSILA VDYPVEKVACYVSD
DGAAMLTFEALSYTAEFARKWVPFCKKSIEPRAPEWYFSQKMDYLKHKVDPAFVMERRAMKRDYEEFKVKINALVSVS
```

Or, upload a file:  No file chosen

#### STEP 2 - Set your pairwise alignment options

```

=====
#
# Aligned_sequences: 2
# 1: cs1A
# 2: cesA
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 1162
# Identity: 125/1162 (10.8%)
# Similarity: 219/1162 (18.8%)
# Gaps: 700/1162 (60.2%)
# Score: 45.5
#
=====

cs1A      1 ----- 0
cesA      1 MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTD 50
cs1A      1 ----- 0
cesA      51 NGEPFIAACNECAFPTCRPCYEYERREGNQACPQCGTRYKRIKGSPRVEGD 100
cs1A      1 ----- 0
cesA      101 EEDDDIDDLLEHEFYGMDPEHVTEAAALYYMRLNTGRGTDEVSHLYSAPGS 150
cs1A      1 ----- 0
cesA      151 EVPLLTYCDEDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSFAIHPRM 200
cs1A      1 ----- 0
cesA      201 VPQKDLTVYGYGSVAWKDRMEVWKKQQIEKLQVVKNERVNDGDGDFIVD 250
cs1A      1 ----- 0
cesA      251 ELDDPGLPMMDGRQPLSRKLPIRSSRINPYRMLIFCRLAILGLFFHYRI 300
cs1A      1 ----- 0
cesA      301 LHPVNDAFGLWLTSVICEIWFAVSWILDQFPKWYPIERETYLDRLSLRYE 350

```

[Help](#)[EMBOSS website](#)[Programmatic Access](#)[Download](#)[Nucleotide form](#)[Similar Applications](#)

## EMBOSS Water - Pairwise Sequence Alignment (PROTEIN)

EMBOSS Water uses the Smith-Waterman algorithm (modified for speed enhancements) to calculate the local alignment of two sequences.

This is the form for protein sequences. Please go to the [nucleotide form](#) if you wish to align DNA or RNA sequences.

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Enter your protein sequences

Enter or paste your first **protein** sequence in any [supported](#) format:

```
MDGVSPKFVLPETFDGVRMEITGQLGMIWELVKAPVIVPLLQLAVYICLLMSVMLLCERVYMGIVIVLVKLFKKPDKRY  
KFEPIHDDEELGSSNFPVVVLQIPMFNEREVYKLSIGAACGLSWPSDRLVIQVLDSTDPTVKQMVEVECQRWASKGINIR  
YQIRENRVCGYKAGALKEGLKRSYVKHCEYVVIFDADFQPEPDFLRRSIPFLMHNPNIALVQARWRFVNSDECLLTRMQE  
MSLDYHFTVEQEVGSSSTHAFFGFNGTAGIWRIA AINEAGGWKDRTTVEDMDLAVRASLRGWKFLYLGDLQVKSELPST  
FRAFRFQQHWRWSGGPANLFRKMVMEIVRNKKVRFWKKVYYVIYSFFFVRKIIAHWTFCFYCVVLPLTILVPEVKVPIWGSV  
YIPSIIILNSVGTPRSIIHLFYWILFENVMSLHRTKATLIGLFEAGRANEWWVTAKGSGQSAKGNTKGIKRFPRIFKLDPRL  
NTLELGFAAFLVCCGYDFVHGKNYYFIYLFLQTMSFFISGLGWIGTYVPS*
```

Or, upload a file:  Choose File No file chosen

**AND**

Enter or paste your second **protein** sequence in any [supported](#) format:

```
>AT2G21770.1|AT2G21770.1|cesA  
MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTNDGEFFIAACNECAFPTCRPCYEYERREGNQA  
CPQCGTRYKRIKGSPRVEGDEEDDDIDDLHEFYGMDPEHVTEAALYYMRLNTGRGTDEVSHLYSASPGSEVPLLYCD  
EDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSFASIHTRPMVPQKDLTVYGYGSVAWKDRMEVKKQQIEKLQVVKN  
ERVNDGDGGFIVDELDDPGLPMMDEGRQLSRKLPIRSSRINPYRMLIFCRLAILGLFFHYRILHPVNDAGLWLTSVICE  
IWFAVSWILDQFPKWYPIERETYLDRLSLRYKEGKPSELAPVDVFVSTVDPLKEPPLITANTVLSILAVDYPVEKVACYVSD  
DGAAMLTFEALSYTAEFARKWVPFCKKSIEPRAPEWYFSQKMDYLKHKVDPAFVMERRAMKRDYEEFKVKINALVSVS
```

Or, upload a file:  Choose File No file chosen

#### STEP 2 - Set your pairwise alignment options

```

=====
#
# Aligned_sequences: 2
# 1: cslA
# 2: cesA
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#

```

The best way to find the optimally aligned regions and calculate the similarity between two sequences

```

# Length: 207
# Identity:      38/207 (18.4%)
# Similarity:    76/207 (36.7%)
# Gaps:          48/207 (23.2%)
# Score: 73.5
#
#
=====
```

cslA	280	GWKDRTTVEDMDLAVRASLRGKFLYLGDLQV--KSELPSTFRAFRFQQH	327
		.....   :....:....   :.. ..... ... ..... ...	
cesA	779	GWIYGSVTEDILTGFKMHCHGWRSVYCMPKRAAFKGSAPINLSDRLHQVL	828
cslA	328	RWSCGPANLFRKMVMIEVRNKKVRF-----WKKVYVIYSFFFVRKIIA	370
		.. ....        :.. ....:               ... ....:..	
cesA	829	RWALGSVEIF-----LSRHCPIWGYGGGLKW-----LERFSYINSVVY	867
cslA	371	HW-----VTFCFYCVVLPLT--ILVPEVK-----VPIWGSVYIPSIIT	406
		.        .. .... ..    ..   :.       :.... .... ..	
cesA	868	PWTSLPLLVYCSLPAICLLTGKFIVPEISNYAGILFLLMFMSIAVTGILE	917
cslA	407	I-LNSVGTPRSIHLFYWILFENVMSLHRTKATLIGLF-----AGRAN	449
		: .. .... .: ..    .. ..         ..   :         ..	
cesA	918	MQWGKIGIDDDWWRNEQFWVI--GGVSSH-----LFALFQGLLKVLGVST	960
cslA	450	EWVVTAK      456	
		..:  :	
cesA	961	NFTVTSK      967	

```

=====
#-----
```

[blastn](#) [blastp](#) [blastx](#) [tblastn](#) [tblastx](#)

BLASTP programs search protein subjects using a protein query. [more...](#)

### Enter Query Sequence

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

[Clear](#)

Query subrange [?](#)

```
FNGTAGIWRIAINEAGGWKDRTTVEDMDLAVRASLRGWKFLYLGDLQVKSELPSTFRAFRFQQ
HRWSCGPANLFRKMVMEIVRNKKVRFWKVYVIYSFFFVRKIIAHWTFCFYCVVLPLTILVPEVK
VPIWGSVYIPSIITILNSVGTPRSIHLLFYWILFENVMSLHRTKATLIGLFEAGRANEVVVTAKLGSG
QSAKGNTKGIKRFPRIKLPDRLNTLELGFAAFLVCGCYDFVHGKNYYFIYLFLQTMSFFISGLG
WIGTYVPS*
```

Or, upload file

[Choose File](#)

No file chosen



Job Title

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

Align two or more sequences [?](#)

### Enter Subject Sequence

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

[Clear](#)

Subject subrange [?](#)

```
>AT2G21770.1|AT2G21770.1|cesA
MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTDNGEPEFIACNECAFPTC
RPCYEYERREGNQACPQCGCTRYKRIKGSPRVEGDEEDDDIDDLHEFYGMDPEHVTEAALYYMR
LNTGRGTDEVSHLYSASPGEVPLLTYCDEDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSASI
HTRPMVPQKDLTVYGYGSVAWKDRMEVWKKQQIEKLQVVKNERVNDGDGDGFIVDELDDPGL
```

Or, upload file

[Choose File](#)

No file chosen



### Program Selection

Algorithm

**blastp** (protein-protein BLAST)

[Choose a BLAST algorithm](#) [?](#)

[From](#)

[To](#)

**BLAST**

Search **protein sequence** using **Blastp (protein-protein BLAST)**

**blast2seq**

AT2G21770.1|AT2G21770.1|cesA

Sequence ID: Icl|34283 Length: 1089 Number of Matches: 4

# Fragmented alignments

Range 1: 779 to 838 [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
23.5 bits(49)	0.044	Compositional matrix adjust.	14/60(23%)	22/60(36%)	2/60(3%)
Query 280	GWKDRTTVEDMDILAVRASLRGKFLYLGDLQV--KSELPSTFRAFRFQQHRWSCGPANLF	337			
Sbjct 779	GW+ ED+ + GW+ +Y + K P Q RW+ G +F				
	GWIYGSVTEDILTGFKMHC HGWR SVY CMPKRAAFKG SAPI NLS DR LH QV L RWA LGS VEIF	838			

Range 2: 414 to 426 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
17.3 bits(33)	3.5	Compositional matrix adjust.	8/17(47%)	9/17(52%)	4/17(23%)
Query 360	YSFFFVRKIIIAHWVTFC	376			
Sbjct 414	Y+ F RK WV FC				
	YTAEFARK----WVPFC	426			

Range 3: 777 to 798 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
16.5 bits(31)	5.7	Compositional matrix adjust.	5/22(23%)	11/22(50%)	0/22(0%)
Query 168	RVG YKAG ALKE GLKRS YVKH C	189			
Sbjct 777	+G+ G++ E + + HC				
	EIGWIYGSVTEDILTGF KMCH	798			

Range 4: 25 to 43 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
15.8 bits(29)	9.2	Compositional matrix adjust.	7/19(37%)	10/19(52%)	0/19(0%)
Query 454	TAKLGSGQSAKGNTKG IKR	472			
Sbjct 25	TA++ S + G T I R				
	TARI RSAEELSGQTCKICR	43			

- Help
- FAQ
- Programmatic Access
- Similar Applications

**■ Database Information**

- UniProt
- UniParc

EBI > Tools > Sequence Similarity Searching > PSI-Search

## PSI-Search - Protein Similarity Search

PSI-Search combines the Smith-Waterman search algorithm with the PSI-BLAST profile construction strategy to find distantly related protein sequences.

**N.B. PSI-SEARCH does not accept sequences containing bases containing translation STOP residues represented with the '\*' (asterisk) character in the sequence.**

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Select your database

##### PROTEIN DATABASES

UniProt Knowledgebase

#### STEP 2 - Enter your input sequence

Enter or paste a PROTEIN sequence in any supported format:

Upload a file:  No file chosen

#### STEP 3 - Set your parameters

##### PSSM E-VALUE CUT-OFF

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

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

# Multiple sequence alignment tools

Foundation for many other further analyses: phylogeny, evolution, motif, protein family etc.

## Multiple Sequence Alignment

**Multiple Sequence Alignment (MSA)** is generally the alignment of three or more biological sequences (protein or nucleic acid) of similar length. From the output, homology can be inferred and the evolutionary relationships between the sequences studied.

By contrast, **Pairwise Sequence Alignment** tools are used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences.

**Clustal Omega**  New MSA tool that uses seeded guide trees and HMM profile-profile techniques to generate alignments. Suitable for medium-large alignments.

 [Launch Clustal Omega](#)

**ClustalW2**  Popular MSA tool that uses tree-based progressive alignments. Suitable for medium alignments.

 [Launch ClustalW2](#)

**DbClustal**  Create a Multiple Sequence Alignment from a protein BLAST result using the DbClustal program.

 [Launch DbClustal](#)

**Kalign**  Very fast MSA tool that concentrates on local regions. Suitable for large alignments.

 [Launch Kalign](#)

**MAFFT**  MSA tool that uses Fast Fourier Transforms. Suitable for medium-large alignments.

 [Launch MAFFT](#)

**MUSCLE**  Accurate MSA tool, especially good with proteins. Suitable for medium alignments.

 [Launch MUSCLE](#)

**MView**  Transform a Sequence Similarity Search result into a Multiple Sequence Alignment or reformat a Multiple Sequence Alignment using the MView program.

 [Launch MView](#)

**T-Coffee**  Consistency-based MSA tool that attempts to mitigate the pitfalls of progressive alignment methods. Suitable for small alignments.

 [Launch T-Coffee](#)

### WebPRANK

The EBI has a new phylogeny-aware multiple sequence alignment program which makes use of evolutionary information to help place insertions and deletions.

Try it out at [WebPRANK](#).

[Help](#)[FAQ](#)[Clustal website](#)[Jalview](#)[Programmatic Access](#)[Download](#)

#### Related Applications

[Pairwise Sequence Alignment](#)[Multiple Sequence Alignment](#)[Phylogeny](#)

## ClustalW2 - Multiple Sequence Alignment

ClustalW2 is a general purpose multiple sequence alignment program for DNA or proteins.

Note: **ClustalW2 is no longer being maintained.** Please consider using the new version instead: [Clustal Omega](#)

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Enter your input sequences

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

```
>AT2G21770.1|AT2G21770.1|cesA
MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTDNCEPFIACNECAFPTCRPCYEYERREGNQA
CPQCGTRYKRIKGSPRVEGDEEDDDIDDLHEFYGMDPEHVTAAALYYMRLNTGRGTDEVSHLYSASPCSEVPLTYCD
EDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSFASIHTPMVPQKDLTVYGYGSVAWKDRMEVWKQQIEKLQVVKN
ERVNDGDDGFIVDELDDPGPLMMDEGRQPLSRKLPIRSSRNPYRMLIFCRLAILGLFFHYRILHPVNDAFGLWLTSVICE
IWFAVSWILDQFPKWYPIERETYLDRLSLRYKEGKPSELAPVDVFVSTVDPLKEPPLITANTVLSILAVDYPVEKVACYVSD
DGAAMLTFEALSYTAEFARKWVPFCKFSIEPRAPEWYFSQKMDYLKHKVDPAFVMERRAMKRDYEEFKVKINALVSVS
```

Or, upload a file:  No file chosen

#### STEP 2 - Set your Pairwise Alignment Options

Alignment Type:  Slow  Fast

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

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

#### STEP 3 - Set your Multiple Sequence Alignment Options

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

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

#### STEP 4 - Submit your job

Help

FAQ

Jalview

Related Applications

Multiple Sequence Alignment

Phylogeny

## ClustalW2 Results

[Alignments](#)[Result Summary](#)[Guide Tree](#)[Submission Details](#)[Submit Another Job](#)

### Alignment

[Download Alignment File](#)[Hide Colors](#)

CLUSTAL 2.1 multiple sequence alignment

AT1G02730.1|AT1G02730.1|cs1D  
os\_42915|LOC\_Os07g36610.1|cs1F  
AT2G21770.1|AT2G21770.1|cesA  
AT2G32530.1|AT2G32530.1|cs1B  
os\_25268|LOC\_Os04g35020.1|cs1H  
AT1G55850.1|AT1G55850.1|cs1E  
AT4G23990.1|AT4G23990.1|cs1G  
AT5G22740.1|AT5G22740.1|cs1A  
AT2G24630.1|AT2G24630.1|cs1C

AT1G02730.1|AT1G02730.1|cs1D  
os\_42915|LOC\_Os07g36610.1|cs1F  
AT2G21770.1|AT2G21770.1|cesA  
AT2G32530.1|AT2G32530.1|cs1B  
os\_25268|LOC\_Os04g35020.1|cs1H  
AT1G55850.1|AT1G55850.1|cs1E  
AT4G23990.1|AT4G23990.1|cs1G  
AT5G22740.1|AT5G22740.1|cs1A  
AT2G24630.1|AT2G24630.1|cs1C

AT1G02730.1|AT1G02730.1|cs1D  
os\_42915|LOC\_Os07g36610.1|cs1F  
AT2G21770.1|AT2G21770.1|cesA  
AT2G32530.1|AT2G32530.1|cs1B  
os\_25268|LOC\_Os04g35020.1|cs1H  
AT1G55850.1|AT1G55850.1|cs1E

----- MVKSAASQSPSPVTITVTPCKGSGDRSLGLTSPIPRASVITNQ 43

MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTD 50

N ----- SPLSSRATRRTSISSGNRRSNGDEGRYCSMSVEDLTAETTNSE 87

----- MALS PAAAGRT ----- GRNNNNNDAG ----- 20

NGEPPFIA CNECAFPTCRPCYERYERREGNQACPQCG ----- 85

CVLSYT VHIPP TP DHQTVF ASQE EED EMLKGNSNQKSFLSGTIFTGGFK 137

TRY KRIK GS PR VEG D E E D D I D D L E H E F Y G M D P E 119

- Help
- MAFFT website
- Jalview
- Programmatic Access
- Download

- Related Applications
- Pairwise Sequence Alignment
- Multiple Sequence Alignment
- Phylogeny

EBI > Tools > Multiple Sequence Alignment > MAFFT

## MAFFT - Multiple Sequence Alignment

MAFFT (Multiple Alignment using Fast Fourier Transform) is a high speed multiple sequence alignment program.

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Enter your input sequences

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

```
>AT2G21770.1|AT2G21770.1|cesA
MNTGGRLIAGSHNRNEFVLINADDTARIRSAEELSGQTCKICRDEIELTDNGEPIACNECAFPTCRPCYERYERREGNQA
CPQCGTRYKRIKGSPRVEGDEEDDDIDDLEHEFYGMDFEHVTEAALYYMRLNTGRGTDEVSHLYSAPGSEVPPLTYCD
EDSDMYSDRHALIVPPSTGLGNRVHHVPFTDSFAIHTRPMVPQKDLTVCYGCSVAWKDRMEVWKKQQIEKLQVVKN
ERVNDGDGDFIVDELDDPGLPMMDEGRQPLSRKLPIRRSSRNPYRMLIFCRLAILGLFFHYRILHPVNDAFGLWLTSVICE
IWFAVSWILDQFPKWYPIERETYLDRLSLRYKEGKPSLAPVDVFVSTVDPLKEPPLITANTVLSILAVDYPVEKVACYVSD
DGAAMLTFEALSYTAEFARKWVPFCKKFSIEPRAPEWYFSQKMDYLHKVDPAFVMERRAMKRDYEEFKVKINALVSVS
```

Or upload a file:  No file chosen

#### STEP 2 - Set your Parameters

OUTPUT FORMAT

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

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

#### STEP 3 - Submit your job

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

Help  
Jalview

Related Applications  
Multiple Sequence Alignment  
Phylogeny

EBI > Tools > Multiple Sequence Alignment > MAFFT

## MAFFT Results

Alignments

Result Summary

Guide Tree

Submission Details

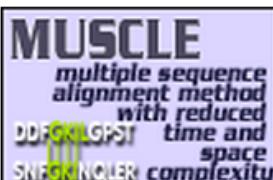
Submit Another Job

### Alignment

Download Alignment File

>AT2G21770.1|AT2G21770.1|cesA

M-----	NTGGRLIAGSHNRNEFVLI
N-----	
-----ADDTARIRSAEELSGQTC--KICRDEIELTDNGEPF	
IACNECAFPTCRPCYERREGNQACPQCGTRYKRIKGSPRVEGDEEDDDIDDLHEFYG	
MDPEHVTTEAALYYMRLNTGRGTDEVSHLYSASPGSEVPLLTYCDEDSDMYSDRHALIVPP	
STGLGNRVHHVPFTDSFA-----HTRPMVPQKDLTVYGYGSVAW-KDRMEVWKQQ	
IEKLQVVKNERVNDGDGDFIVDELDDPGPLPMMDEG-RQP-----LSRKLPIRSSRINPY	
RMLIFCRLAILGLFFHYRILHPVNDAF-----GLWLTSVICEIWFAWSWILDQFPKWYPIE	
RETYLDRL-----SLRYKEG-----KPSELAPVDV-FVSTVDPLKEPPLITANT	
VLSILAVDYP--VEKVACYVSDDGAAMLTFEALSYTAEFARKWVP-FCKKFSIEPRAPEW	
YFSQKMD-YLKHKVDPAFVMERRAMKRDYEEFKVKINALVSVSQK-----	
-----VPEDGWTMQDGTPWPG-----NNVRDHPGMIQVFLGHSGVC	
DMDGNE-----LPRLVYVSREKRPGFDHHKKAGAMNSLIRVSAVLNAPY	
LLNVDCDHYINNSKAIREAMCFMMDPQ-SGKKICYVQFPQRFDGIDRH--DRYSNRNVVF	
FDINMKGLDGIQGPI--YVGTGCVFRRQALYGFADPKKKQPPGRTNCWPWKCCLCGMR	
KKKTGKVKDNRKKPKETSKQIHALEHIEGLQVTNAENNSETAQLKLEKKFGQSPVLVA	
STLL-----LNG-----GVPSNVNPASLLRESIQVISCGYEKTEWGKE	
IGWIYGSVTEDIITGFKMHCHGWRSSVYCMPKRAAFKGSAPINLSDRLHQVLRWALGSVEI	
F-----LSRHCPWIYGYGGGLKWLERFSYINSVVYPWTSLLVYCSLPAICL-LTG	
KFIVPEI-SNYAGILFLLMFMSIAVTGILEMQWGKIGIDDWWRNEQFWVIGGVSSHFLAL	
FQGLLKVLAGVS-TNFTVT-----SKAAD---D----GEFSELY---IFKW---	
TSLL-----IPPTTLLIIINVGVIVGVSDAINNGY--DSWGPLFGRRLFFALWVIV	
HLY-----PFLKGLLGKQD--RVPTIILVW-----SILLASILTLL-----WV-	
-----RVNPFFVSKDGPVLEICGLDCLK----*	
>AT1G02730.1 AT1G02730.1 ces1D	



## MUSCLE - Multiple Sequence Alignment

MUSCLE stands for **MU**ltiple **S**equence **C**omparison by Log- **E**xpectation. MUSCLE is claimed to achieve both better average accuracy and better speed than ClustalW2 or T-Coffee, depending on the chosen options.

**Internet Explorer users: If button presses (including copy/paste operations) don't appear to work please try enabling Compatibility View.**

### Use this tool

#### STEP 1 - Enter your input sequences

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

```
DDWTFMKDEYDKLVRRIKNTDERSLLRHGGGEFFAEFLNVERRNHPTIVKTRVSAMTNAPIMLNMDCDMFVNPNQAVLHAMCLLLGFDDAESGFGVQAPQRFYDALKDDPFGNQMECFKRFISGVQGVQGAFYAGTCFHRKAVYGVPPNFNGAEREDTIGSSSYKELHTRFCNSEELNESARNIIWDLSSKPMVDISSLRIEVAKAVSACNYDIGTCWGQEVGWVYGSLTEILDITGQRIHAMGWRCSVLMVTEPPAFMGSAPICGGPACLTQFKRWATGQSEIIISRNNPILATMFKRLKFRQCLAYLIVLGWRRAPFELCYGLLGPYCILTNQSFLPKASEDGFSVPLALFISYNTYNFMEYMACGLSARAWWNHNRMQRIISVSAWTLAFLT VLLKSLGLSETVFETGKDMSDDDDNTDGADPGRTFDSPVFIPTALAMLNIVAVTVGACRVAFTGAEGVPCAPGIGEFMCCGWLVLCFFPFVRGIVWGKGSYGIPWSVVLKASLLVAMFVTFCRNL
```

Or upload a file:  No file chosen

#### STEP 2 - Set your Parameters

OUTPUT FORMAT:

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

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

#### STEP 3 - Submit your job

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

[Help](#)[Jalview](#)**Related Applications**[Multiple Sequence Alignment](#)[Phylogeny](#)

EBI &gt; Tools &gt; Multiple Sequence Alignment &gt; MUSCLE

**MUSCLE Results**[Alignments](#) [Result Summary](#) [Submission Details](#) [Submit Another Job](#)**Alignment**[Download Alignment File](#) [Show Colors](#)

MUSCLE (3.8) multiple sequence alignment

AT5G22740.1|AT5G22740.1|cs1A  
 AT2G24630.1|AT2G24630.1|cs1C  
 os\_25268|LOC\_Os04g35020.1|cs1H  
 AT2G32530.1|AT2G32530.1|cs1B  
 AT2G21770.1|AT2G21770.1|cesA  
 AT1G02730.1|AT1G02730.1|cs1D  
 os\_42915|LOC\_Os07g36610.1|cs1F  
 AT4G23990.1|AT4G23990.1|cs1G  
 AT1G55850.1|AT1G55850.1|cs1E

-----  
 MAPRFDSDLWAKETRRG-----  
 MAVVAAAAATGST-----  
 MADSSSSL-----  
 MNTGGRLIAGSHNRNEFVLINADDTARI-----  
 MVKSAASQSPSPVTTTPCKGSGDRSLGLTSPIPRASVITQNQNSPLSSRATRRRTSISSG-----  
 MALSPAAGRTG-----  
 MYQVSLKQFVFLLKIKSTM-----  
 MVNKDDRI-----

AT5G22740.1|AT5G22740.1|cs1A  
 AT2G24630.1|AT2G24630.1|cs1C  
 os\_25268|LOC\_Os04g35020.1|cs1H  
 AT2G32530.1|AT2G32530.1|cs1B  
 AT2G21770.1|AT2G21770.1|cesA  
 AT1G02730.1|AT1G02730.1|cs1D  
 os\_42915|LOC\_Os07g36610.1|cs1F  
 AT4G23990.1|AT4G23990.1|cs1G  
 AT1G55850.1|AT1G55850.1|cs1E

-----  
 -----  
 -----  
 RSAEEELSGQTCKICRDEIELTDNGEPIACNECAFPTC---RPCYEYERREGNQACPQC  
 NRRSNGDEGRYCSMSVEDLTAETTNSECVLSYTWHIPPTPDHQTVFASQESEEDEMLKGN  
 -RNNNNNDAG-----  
 -----  
 -----

AT5G22740.1|AT5G22740.1|cs1A  
 AT2G24630.1|AT2G24630.1|cs1C  
 os\_25268|LOC\_Os04g35020.1|cs1H  
 AT2G32530.1|AT2G32530.1|cs1B  
 AT2G21770.1|AT2G21770.1|cesA  
 AT1G02730.1|AT1G02730.1|cs1D

-----  
 -----  
 -----  
 -----  
 GTRYKRIKGSPRVEGDEEDDDIDDLHEFYGMDPEHVTEAALYY-----M  
 SNOKSFTSCTTETCCFKSVTCHVTDOSMDRADDEKKSGOTCWTLKGCDPEKUVVHCRCECCF

## MAFFT > Clustal Omega > MUSCLE >> ClustalW

accuracy

speed

**Table I** BALiBASE results

Aligner	Av score (218 families)	BB11 (38 families)	BB12 (44 families)	BB2 (41 families)	BB3 (30 families)	BB4 (49 families)	BB5 (16 families)	Tot time (s)	Consistency
MSAprobs	0.607	0.441	0.865	0.464	0.607	0.622	0.608	12 382.00	Yes
Probalign	0.589	0.453	0.862	0.439	0.566	0.603	0.549	10 095.20	Yes
MAFFT (auto)	0.588	0.439	0.831	0.450	0.581	0.605	0.591	1475.40	Mostly (203/218)
Probcons	0.558	0.417	0.855	0.406	0.544	0.532	0.573	13 086.30	Yes
Clustal Ω	0.554	0.358	0.789	0.450	0.575	0.579	0.533	539.91	No
T-Coffee	0.551	0.410	0.848	0.402	0.491	0.545	0.587	81 041.50	Yes
Kalign	0.501	0.365	0.790	0.360	0.476	0.504	0.435	21.88	No
MUSCLE	0.475	0.318	0.804	0.350	0.409	0.450	0.460	789.57	No
MAFFT (default)	0.458	0.258	0.749	0.316	0.425	0.480	0.496	68.24	No
FSA	0.419	0.270	0.818	0.187	0.259	0.474	0.398	53 648.10	No
Dialign	0.415	0.265	0.696	0.292	0.312	0.441	0.425	3977.44	No
PRANK	0.376	0.223	0.680	0.257	0.321	0.360	0.356	128 355.00	No
ClustalW	0.374	0.227	0.712	0.220	0.272	0.396	0.308	766.47	No

The figures are total column scores produced using bali score on core columns only. The average score over all families is given in the second column. The results for BALiBASE subgroupings are in columns 3–8. The total run time for all 218 families is given in the second last column. The last column indicates whether the method is consistency based.

*Molecular Systems Biology* 7:539, 2011

<http://mafft.cbrc.jp/alignment/software/about.html>



Query all databases

 help**Visual Guidance****Categories****proteomics**

protein sequences and identification  
mass spectrometry and 2-DE data  
protein characterisation and function  
families, patterns and profiles  
post-translational modification  
protein structure  
protein-protein interaction  
similarity search/alignment

**genomics**

structural bioinformatics  
systems biology  
phylogeny/evolution  
population genetics  
transcriptomics  
biophysics  
imaging  
IT infrastructure  
drug design

**Resources A-Z****Links/Documentation**

## SIB resources

External resources - (No support from the ExPASy Team)

**Databases**

- [UniProtKB](#) • functional information on proteins • [more]
- [UniProtKB/Swiss-Prot](#) • protein sequence database • [more]
- [STRING](#) • protein-protein interactions • [more]
- [SWISS-MODEL Repository](#) • protein structure homology models • [more]
- [PROSITE](#) • protein domains and families • [more]
- [ViralZone](#) • portal to viral UniProtKB entries • [more]
- [neXtProt](#) • human proteins • [more]

[EMBnet services](#) • bioinformatics tools, databases and courses • [more]

- [ENZYME](#) • enzyme nomenclature • [more]
- [GlycoSuiteDB](#) • glycan database • [more]
- [GPSDB](#) • gene and protein synonyms • [more]
- [HAMAP](#) • UniProtKB family classification and annotation • [more]
- [MIAPEGelDB](#) • MIAPE document edition • [more]
- [MyHits](#) • protein domains database and tools • [more]
- [PANDITplus](#) • protein families and domains resources • [more]
- [PaxDb](#) • protein abundance database • [more]
- [Prolune](#) • Popular science articles (in French) • [more]
- [Protein Model Portal](#) • structural information for a protein • [more]
- [Protein Spotlight](#) • Informally written reviews on proteins • [more]
- [SugarBind](#) • pathogen sugar-binding • [more]
- [SWISS-2DPAGE](#) • proteins on 2-D and SDS PAGE maps • [more]
- [SwissSidechain](#) • non-natural amino-acid sidechains • [more]
- [SwissVar](#) • variants in UniProtKB entries • [more]
- [TCS](#) • interaction specificity in two-component systems • [more]
- [UniMES \(UniProt metagenomic samples\)](#) • UniProt Metagenomic and Environmental Sequences • [more]

**Tools**

- [SWISS-MODEL Workspace](#) • structure homology-modeling • [more]
- [SwissDock](#) • protein ligand docking server • [more]
- [2ZIP](#) • Prediction of leucine zipper domains • [more]
- [3of5](#) • find user-defined patterns in protein sequences • [more]
- [AACompIdent](#) • protein identification by aa composition • [more]
- [AACompSim](#) • amino acid composition comparison • [more]
- [Agadir](#) • Prediction of the helical content of peptides • [more]
- [ALF](#) • simulation of genome evolution • [more]
- [Alignment tools](#) • Four tools for multiple alignments • [more]
- [AllAll](#) • protein sequences comparisons • [more]
- [APSSP](#) • Advanced Protein Secondary Structure Prediction • [more]
- [Ascalaph](#) • Molecular modeling software • [more]
- [big-PI](#) • predict GPI modification sites • [more]
- [Biochemical Pathways](#) • Biochemical Pathways • [more]
- [BLAST](#) • sequence similarity search • [more]
- [BLAST - NCBI](#) • Biological sequence similarity search • [more]
- [BLAST - PBIL](#) • BLAST search on protein sequence databases • [more]
- [Blast2Fasta](#) • Blast to Fasta conversion • [more]
- [boxshade](#) • MSA pretty printer • [more]
- [CFSSP](#) • Protein secondary structure prediction • [more]
- [ChloroP](#) • chloroplast transit peptides & cleavage sites • [more]
- [Click2Drug](#) • Directory of computational drug design tools • [more]
- [ClustalO \(UniProt\)](#) • Align two or more protein sequences • [more]
- [ClustalW](#) • Multiple sequence alignment • [more]
- [ClustalW - PBIL](#) • Multiple sequence alignment program • [more]

# boxshade

## Visual Guidance

## Categories

### proteomics

protein sequences and identification  
mass spectrometry and 2-DE data  
protein characterisation and function  
families, patterns and profiles  
post-translational modification  
protein structure

protein-protein interaction

### similarity search/alignment

genomics

structural bioinformatics

systems biology

phylogeny/evolution

population genetics

transcriptomics

biophysics

imaging

IT infrastructure

drug design

## Resources A..Z

## Links/Documentation

### SIB resources

External resources - (No support from the ExPASy Team)

## Databases



UniProtKB • functional information on proteins • [\[more\]](#)



MyHits • protein domains database and tools • [\[more\]](#)

## Tools

Alignment tools • Four tools for multiple alignments • [\[more\]](#)

BLAST • sequence similarity search • [\[more\]](#)

BLAST - NCBI • Biological sequence similarity search • [\[more\]](#)

BLAST - PBIL • BLAST search on protein sequence databases • [\[more\]](#)

Blast2Fasta • Blast to Fasta conversion • [\[more\]](#)

boxshade • MSA pretty printer • [\[more\]](#)

ClustalO (Only 10!) • Align two or more protein sequences • [\[more\]](#)

ClustalW • Multiple sequence alignment • [\[more\]](#)

ClustalW - PBIL • Multiple sequence alignment program • [\[more\]](#)

ClustalW2 • Multiple sequence alignment program • [\[more\]](#)

Decrease redundancy • Sequence redundancy reduction • [\[more\]](#)

DIALIGN • Local multiple sequence alignment • [\[more\]](#)

Dotlet • sequence similarity plots • [\[more\]](#)

FASTA/SSEARCH/GGSEARCH/GLSEARCH • Sequence similarity searching of protein db • [\[more\]](#)

GENIO/logo • RNA/DNA & Amino Acid Sequence Logos • [\[more\]](#)

Kalign - EBI • Fast and accurate multiple sequence alignment • [\[more\]](#)

Kalign - SBC • Fast and accurate multiple sequence alignment • [\[more\]](#)

LALIGN • Pairwise alignment • [\[more\]](#)

MAFFT - CBRC • Multiple sequence alignment • [\[more\]](#)

MAFFT - EBI • Multiple sequence alignment • [\[more\]](#)

MaxAlign • Gap removal from alignments • [\[more\]](#)

Multalin • Multiple sequence alignment • [\[more\]](#)

MUSCLE • Multiple alignment server • [\[more\]](#)

MyHits • protein domains database and tools • [\[more\]](#)

psort • Generalized profile build and search tools • [\[more\]](#)

**Note:** Starting July 17, version 3.21 is running on this server. The only changes are some improvements in the RTF output routine that enables RTF with shaded background for users of MS-Word 7.0

I have also compiled a list of [List of frequently asked questions](#).

This server uses version 3.21 of BOX SHADE, written by K. Hofmann and M. Baron.

[BOX SHADE](#) is in the public domain and available from Source Forge

<http://sourceforge.net/projects/boxshade/> The available version runs on PCs, VMS- and OSF1-machines and includes much more options than are implemented on this server.

This **server** takes a multiple-alignment file in either GCG's MSF-format or Clustal ALN-format.

Output can be created in the following formats:

- Postscript/EPS (using shaded background)
- RTF old (using colors)
- RTF new (using shaded background)
- XFIG-files (using shaded background)
- ASCII (showing similarities)
- ASCII (showing differences)
- HPGL (using colors)
- PICT (for later editing on MACs and PCs)

If you have problems using this server (like getting no result), [read this](#) and see the [FAQ list](#).

Output format	ASCII_similarities
Font Size	10
Consensus Line	consensus line with symbols
Fraction of sequences:	0.5 (that must agree for shading)
Enter sequence number:	9 only if 'consensus to a single sequence' is required
Query title (optional)	csl
<b>When pasting MSF or ClustalW files, please make sure that the pasted text starts with the header line of the alignment and contains no extra blank lines at the bottom.</b>	
Input sequence format	ALN
Paste your multiple-alignment file (see above for valid formats)	CLUSTAL format alignment by MAFFT L-INS-1 (v6.850b) AT2G21770.1 AT2 M----- NTGGRLIAGSHNRNEFVLI
<input type="button" value="Run BOX SHADE..."/> <input type="button" value="Clear Input"/>	

-  [Seq2Logo](#) • Protein sequence logo method • [more]
-  [SeqLogos](#) • Sequence logo from amino acid alignment • [more]
-  [Sequence Variability Server](#) • Protein sequence variability in MSA • [more]
-  [Sequerome](#) • BLAST similarity search and sequence profiling • [more]
-  [SIM](#) • binary sequence alignment • [more]
-  [T-Coffee](#) • sequence and structure multiple alignments • [more]
-  [T-Coffee - EBI](#) • Multiple sequence alignment program • [more]
-  [T-Coffee - WUR](#) • Multiple sequence alignment program • [more]
-  [WebLogo](#) • Sequence logos • [more]
-  [WU BLAST](#) • Sequence similarity search in protein databases • [more]

**Multiple Sequence Alignment**

CLUSTAL format alignment by MAFFT L-INS-1 (v6.850b)

```
AT2G21770.1|AT2 M-----NTGGRLIAGSHNRNEFVLI
AT1C02730.1|AT1 MVKSAASQSPSPVTTVTPCKGSGDRSLGLTSPIPRASVITNQNPLSSRATRTSISSG
os_42915|LOC_Os_MA-----LSPAAAGRTG-----
AT1G55850.1|AT1 MVN-----
AT4G23990.1|AT4 MYQ-----
AT2G32530.1|AT2 MA-----
os_25268|LOC_Os_MA-----
```

**Upload Sequence Data:** No file chosen**Image Format & Size****Image Format:****Logo Size per Line:**

18 X 5 cm

**Advanced Logo Options****Sequence Type:**  amino acid  DNA / RNA  Automatic Detection**First Position Number:** **Logo Range:**  - **Small Sample Correction:** **Frequency Plot:** **Multiline Logo (Symbols per Line):**  ( 32 )**Multiline Logo****Advanced Image Options****Bitmap Resolution:**  pixels/inch (dpi)**Antialias Bitmaps:** **Title:** **Y-Axis Height:**  (bits)**Show Y-Axis:** **Y-Axis Label:**  bits**Show X-Axis:** **X-Axis Label:** **Show Error Bars:** **Label Sequence Ends:** **Boxed / Boxed Shrink Factor:**  / **Outline Symbols:** **Show fine print:** **Y-Axis Tic Spacing:**  (bits)**Colors****Color Scheme:**  Default  Black & White  Custom (See Below.)**Symbols****Color****RGB****Symbols****Color****RGB**

KRH

green

purple

## Categories

proteomics

genomics

structural bioinformatics

systems biology

phylogeny/evolution

population genetics

transcriptomics

biophysics

imaging

IT infrastructure

drug design

## Resources A..Z

## Links/Documentation

algorithm • [more]

Evolutionary Trace Server (TraceSuite II) • Maps evolutionary traces to structures • [more]

fastsimcoal • coalescent simulation of genomic data • [more]

Linear Classification • simple linear classification • [more]

MLtree • maximum likelihood optimization • [more]

MLTreeMap • phylogenetics and functionalities of metagenomes • [more]

Newick Utilities • high-throughput phylogenetic tree processing • [more]

PHYLIP • Package of programs for phylogenetic analysis • [more]

Phylogenetic Tree • phylogenetic tree construction and printing • [more]

Phylogeny.fr • Simple phylogenetic analysis • [more]

Phylogeny programs • Links to phylogeny programs • [more]

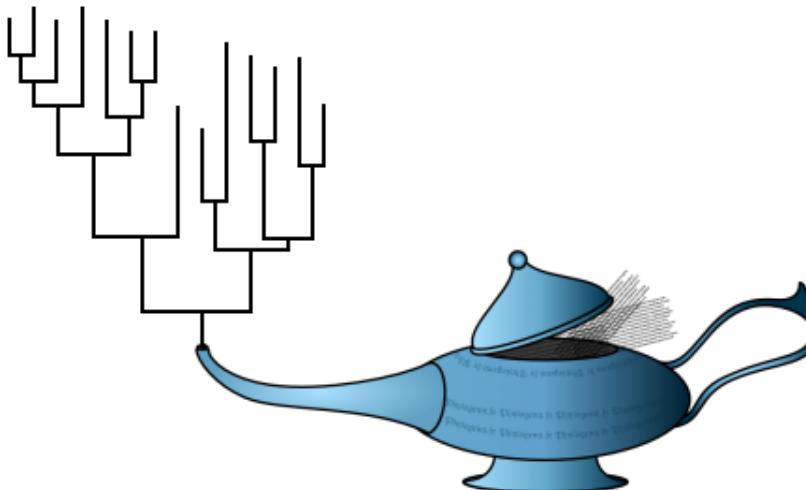
RAxML • ML inference of large phylogenetic trees • [more]

SuperTree • assemble phylogenetic trees • [more]

Home **Phylogeny Analysis** Blast Explorer Online Programs Your Workspace Documentation Downloads Contacts

"One Click"  
"Advanced"  
"A la Carte"

## Phylogeny.fr Robust Phylogenetic Analysis For The Non-Specialist



**Phylogeny.fr** is a free, simple to use web service dedicated to reconstructing and analysing phylogenetic relationships between molecular sequences.

**Phylogeny.fr** runs and connects various bioinformatics programs to reconstruct a **robust phylogenetic tree** from a set of sequences.

If you use this site, please cite:

"One Click" Mode

Alignment  
MUSCLE

Curation  
Gblocks

Phylogeny  
PhyML

Tree Rendering  
TreeDyn

1. Overview

2. Data & Settings

Name of the analysis (optional):

Upload your set of sequences in FASTA, EMBL or NEXUS format from a file:

No file chosen

Or paste it here ([load example of sequences](#))

```
LLDWWRNEQFYMIGATGVYLAALHIVLKRLGLKGVRFKLTAKQLAGGARERFAELYDVHWSPLLAPTV
VVMAVNVTAAIGAAAAGKAVVGGWTPAQVAGASAGLVFNWVVLVLLYPFALGIMGRWSKRPCALFALLVAAC
AAVAAGFVAVHAVLAAGSAAPSWLGWSRGATAILPSSWRLKRGF
>os_25268 | LOC_Os04g35020.1 | cslH
MAVVAAAAATGSTTRSGGGGEGTRSGRKKPPPLQERVPLGRRAAWWRLAGLAVLLLLLALLALRLL
RHHGGAGGDGGVWRVALVCEAWFAALCALNVSAKWSPVRFVTRPENLVAEGRTPSTTAAEYGELPAVDML
VTTADPALEPPLVTVNTVLSLLALDYPRAGERLACYVSDDGCSPLTCHALREAAGAAAWPFCRRYGVA
VRAPFRYFSSSSSESGGPADRKFLLDTFMKDEYDKLVRRIKNTDERSLLRHGGGEFFAEFLNVERRNH
PTIVKTRVSAVMTNAPIMLNMDCDMFVNNPQAVLHAMCLLGFDDAESGGFVQAPQRFYDAIKDDFFGNQ
MECFFKRFISGVQGVQGAFYAGTGCFHRRKAVYGVPPNFNGAEREDTIGSSSYKELHTRFGNSEELNES
RNIIWDLSSKPMVDISSLRIEVAKAVSACNYDIGTCWGQEVGWVYGSITEDILTGQRIHAMGWRSVLMVTE
PPAFMGSAPIGGPACLTQFKRWTGQSEIIISRNNPILATMFKRLKFRQCLAYLIVLGWPIRAPFELCYG
LLGPYCILTNQSFLPKASEDGFSPVPLALFISYNTYNFMEYMACGLSARAWNNHRMQRIISVSAWTLAFL
TVLLKSLGLSETVFEVTGDKDSMSDDDDNTDGADPGRFTFDSPVFIPTALAMLNIVAVTVGACRVAFG
TAEGVPCAPGIGEFMCCGWLVLCCFPFVRGIVWGKGSYGI PWSVKLKA SLLVAMFVTFCRKN
```

Maximum number of sequences is 200 for proteins and 200 for nucleic acids.

Maximum length of sequences is 2000 for proteins and 6000 for nucleic acids.

Use the Gblocks program to eliminate poorly aligned positions and divergent regions

To receive the results by e-mail, enter your address(es):

csl

Alignment  
MUSCLE

Curation  
Gblocks

Phylogeny  
PhyML

Tree Rendering  
TreeDyn

1. Overview

2. Data & Settings

3. Alignment

4. Curation

5. Phylogeny

6. Tree Rendering

## Tree Rendering results

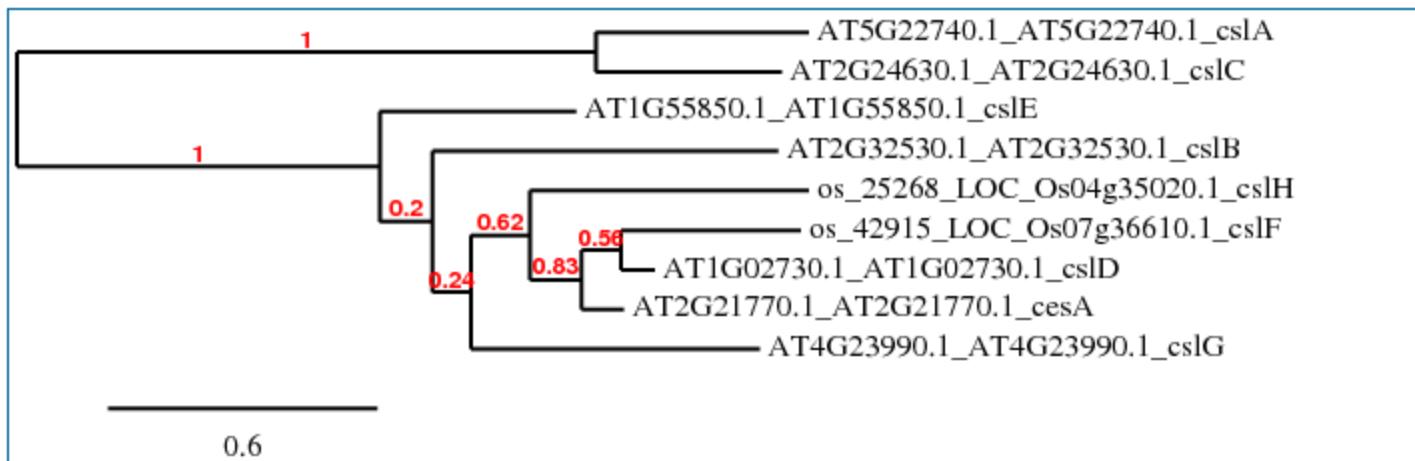


Figure 1: Phylogenetic tree.

==> Download the tree: [PNG](#) - [PDF](#) - [SVG](#) - [TGF \(Treedyn format\)](#) - [Newick](#) - [Text](#)

### Select an action:



Reset (cancel all changes)



Mid-point rooting



Use Genbank information to automatically rename leaves by:  species and gi  species only  colorize

## Visual Guidance

### Categories

proteomics

protein sequences and identification

mass spectrometry and 2-DE data

protein characterisation and function

families, patterns and profiles

post-translational modification

protein structure

protein-protein interaction

similarity search/alignment

genomics

structural bioinformatics

systems biology

phylogeny/evolution

population genetics

SIB resources

External resources - (No support from the ExPASy Team)

### Databases

 UniProtKB • functional information on proteins • [\[more\]](#)

 UniProtKB/Swiss-Prot • protein sequence database • [\[more\]](#)

 ViralZone • portal to viral UniProtKB entries • [\[more\]](#)

 neXtProt • human proteins • [\[more\]](#)

 HAMAP • UniProtKB family classification and annotation • [\[more\]](#)

 SwissVar • variants in UniProtKB entries • [\[more\]](#)

 [UniMES \(UniProt metagenomic samples\)](#) • UniProt Metagenomic and Environmental Sequences • [\[more\]](#)

 [UniRef \(UniProt sequence clusters\)](#) • UniProtKB sequence clusters • [\[more\]](#)

### Tools

 AACompIdent • protein identification by aa composition • [\[more\]](#)

 Decrease redundancy • Sequence redundancy reduction • [\[more\]](#)

 EasyProt • graphical platform for proteomics analysis • [\[more\]](#)

 FindPept • peptide identification from unspecific cleavage • [\[more\]](#)

 GlycoMod • oligosaccharide structure prediction • [\[more\]](#)

 Graphical Codon Usage Analyser • Codon bias • [\[more\]](#)

 HAMAP • UniProtKB family classification and annotation • [\[more\]](#)

 LALIGN • Pairwise alignment • [\[more\]](#)

 PeptideCutter • protein cleavage sites prediction • [\[more\]](#)

 PeptideMass • peptides from protein cleavage • [\[more\]](#)

 Reverse Translate • Reverse translation • [\[more\]](#)

 TaIdent • protein identification with PI, Mw and tag • [\[more\]](#)

 Translate • nucleotide sequence translation • [\[more\]](#)

## Translate tool

**Translate** is a tool which allows the translation of a nucleotide (DNA/RNA) sequence to a protein sequence.

Please enter a DNA or RNA sequence in the box below (numbers and blanks are ignored).

```
>contig00008 gene=isogroup00001 length=862
tAAGATACCTCGAAATTTTATTGATTTATTTCATTGAGAATGGTAACATACATC
ACGCTGGCTTCACTGAGCACATTCTCAGGGATAAAAATGGACCGAATCCAGCCTTT
TTACAAAATGCCAGCCCCAAAATTAGCAAGAGGACCAGTTGTCAGGATTCAAGGCCACC
TACACAGAAGCTATGACTTGATAAGGTCCACCACCGTTGCTTAGCCCCACTCGTTA
TCATACCGCGACAACAACTTATGAAATTTCGTTAACCGCAATGCCGGCTTGCCTCG
AAGATACTCGACCTGCTATACCCAATCAGGCTGAAGAGAGCCAGGTCTTCTGTGAT
CCAAGAACATACCTTCAGTGCCCCCTGATTCTTATAACAGATTTACTTCTCG
TAAGAAGTTGTTCTGAGCTTACAGTCAGATCGACAACGTATACATCAGCAGTAGGC
ACACGGAATGCCATTCTGTTAATTCCGTTAATGCCGGAACCCCTCCAACAGCC
TTTGCAGCTCAGTTGAGCTTGAATGATGTTGAATCCAGCACCTCGTCCACCCCTCAG
TCTTGGCAGAAGGTCCATCAACTGTTCTCGCTCGCTGTTACAGCGTAAACGGTAGTC
ATAAGGCCCTAGTAATCCAAAATTCTCATGTATAACCTTAGCAAGGGGGCTAGACAA
TTAGTAGTGCAGCTTGCATTGAGATTATGTCAGTCTGGAGTATAGTCCTTTCTGGCT
ACACCTATAACGAACATAGGTGCATCCTGCTGGAGCGCTGATCACCACTTCTGGCT
```



Output format:

Genetic code:

or

## Better tool for open reading frame (ORF) prediction from EST data



The OrfPredictor (ORF-Predictor) server is designed for ORF prediction and translation of a batch of EST or cDNA sequences. [See more details](#)

Paste sequences below in **FASTA** format (Support multiple sequences in one FASTA file)

```
>contig00008 gene=isogroup00001 length=862
tAAGATACCTCGAAATTTATTGCATTTATTTGATTGAGAATGGTAACATC
ACGCTGGCTTCACTGAGCACATTCTCAGGGATGAAAAATGGACCGAATTCCAGCCTTT
TTACAAAATGCCAGCCCCAAAATTAGCAAGAGGACCAGTTGTCAAGGATTCAAGGCCACC
TACACAGAAGCTATGTACTTGATAAGGTCCACCACCGTTGCTGTAGCCCCACTCGTTA
```

Or load from disk  No file chosen

[Optional] Paste BLASTX output below in **BLASTX** output format

Or load from disk  No file chosen

Select the strand for prediction

(Sequenced from 3' end, select '-'; from 5' end, select '+'; mixed or unknown, select 'both')

E-value in BLASTX

E-mail results to: yanbin.yin@gmail.com

[http://proteomics.ysu.edu/tools/user\\_results/usr\\_WedFeb61910112013/](http://proteomics.ysu.edu/tools/user_results/usr_WedFeb61910112013/)

## Visual Guidance

### Categories

#### proteomics

protein sequences and identification

mass spectrometry and 2-DE data

protein characterisation and function

#### families, patterns and profiles

post-translational modification

protein structure

protein-protein interaction

similarity search/alignment

#### genomics

structural bioinformatics

systems biology

phylogeny/evolution

population genetics



SIB resources



External resources - (*No support from the ExPASy Team*)

## Databases



**PROSITE** • protein domains and families • [more]



**HAMAP** • UniProtKB family classification and annotation • [more]



**MyHits** • protein domains database and tools • [more]



**PANDITplus** • protein families and domains resources • [more]

## Database of protein domains, families and functional sites

PROSITE consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them [More... / References /]. PROSITE is complemented by **ProRule**, a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information on functionally and/or structurally critical amino acids [More...].

Release 20.89, of 31-Jan-2013 (1659 documentation entries, 1308 patterns, 1050 profiles and 1053 ProRule)

### PROSITE access

e.g. PDOC00022, PS50089, SH3, zinc finger  
  add wildcard <sup>\*\*</sup>

#### Browse:

- by documentation entry
- by ProRule description
- by taxonomic scope
- by number of positive hits

### PROSITE tools

#### Scan a sequence against PROSITE patterns and profiles - quick scan

(Output includes graphical view and feature detection)



Enter your sequence(s) or (a) UniProtKB (Swiss-Prot or TrEMBL) AC or ID [ Help ]:

```
>AT5G22740.1|AT5G22740.1|csIA
MDGVSPKFVLPETFDGVRMEITGQLGMWIWELVKAPVI
VPLLQLAVYICLCLMSVMILLCERYVMGIVIVLVKLFWKK
PDKRYKFEPIHDEELGSNNFPVVLVQIPMFNEREVYK
LSIGAACGLSWPSDRLVIQVLDSTDPTVKQMVEVE
CQRWASKGINIRYQIRENRVGYKAGALKEGLKRSYVK
```

exclude patterns with a high probability of occurrence

- **ScanProsite** - advanced scan
- **PRATT** - allows to interactively generate conserved patterns of aligned proteins.
- **MyDomains - Image Creator** - allows to generate custom



## NUCLEOTIDE SEQUENCES

### Whole genome visualization and analysis

« [GenomeAtlas](#)

DNA structural atlases for complete microbial Genomes

### Gene finding and splice sites

« [EasyGene](#)

Genes in prokaryotes

« [EasyGene](#)

Genes in prokaryotes

« [HMMgene](#)

Genes in eukaryotes

[MetaRanker](#)

Identification of risk genes in complex phenotypes

[NetAspGene](#)

Intron splice sites in Aspergillus DNA

« [NetGene2](#)

Intron splice sites in human, C. elegans and A. thaliana DNA

[NetPlantGene](#)

Intron splice sites in *Arabidopsis thaliana* DNA

« [NetStart](#)

Translation start in vertebrate and *A. thaliana* DNA

[NetUTR](#)

Splice sites in 5' UTR regions of human genes

« [Promoter](#)

Transcription start sites in vertebrate DNA

« [RNAmmer](#)

Ribosomal RNA sub units

« [RNAmmer](#)

Ribosomal RNA sub units

### Analysis of DNA microarray data

[GenePublisher](#)

Analysis of DNA microarray data

« [OligoWiz](#)

Design of oligonucleotides for DNA microarrays

## SMALL MOLECULES

« [ChemProt](#)

Chemical-protein interactions

## AMINO ACID SEQUENCES

### Protein sorting

ChloroP »

Chloroplast transit peptides and their cleavage sites in plant proteins

LipoP »

Signal peptidase I & II cleavage sites in gram- bacteria

NetNES »

Leucine-rich nuclear export signals (NES) in eukaryotic proteins

SecretomeP »

Non-classical and leaderless secretion of proteins

SignalP »

Signal peptide and cleavage sites in gram+, gram- and eukaryotic amino acid sequences

TargetP »

Subcellular location of proteins: mitochondrial, chloroplastic, secretory pathway, or other

TatP »

Twin-arginine signal peptides

### Post-translational modifications of proteins

DictyOGlyc

O-(alpha)-GlcNAc glycosylation sites (trained on *Dictyostelium discoideum* proteins)

NetAcet

N-terminal acetylation in eukaryotic proteins

NetCGlyc »

C-mannosylation sites in mammalian proteins

NetCorona

Coronavirus 3C-like proteinase cleavage sites in proteins

NetGlycate »

Glycation of ε amino groups of lysines in mammalian proteins

NetNGlyc »

N-linked glycosylation sites in human proteins

NetNGlyc »

N-linked glycosylation sites in human proteins

NetOGlyc »

O-GalNAc (mucin type) glycosylation sites in mammalian proteins

NetOGlyc »

O-GalNAc (mucin type) glycosylation sites in mammalian proteins

NetPhorest

Linear motif atlas for phosphorylation-dependent signaling

NetPhos »

Generic phosphorylation sites in eukaryotic proteins

NetPhosBac

Generic phosphorylation sites in bacterial proteins

NetPhosK

Kinase specific phosphorylation sites in eukaryotic proteins

NetPhosYeast

Serine and threonine phosphorylation sites in yeast proteins

Google: cbs dtu

## Instructions

## Output format

### SUBMISSION

Paste a single sequence or several sequences in **FASTA** format into the field below:

```
>contig00008 gene=isogroup00001 length=862  
tAAGATACCTCGAAATATTTATTTGCATTTCATTGAGAATGGTAAC  
TACATC
```

Submit a file in **FASTA** format directly from your local disk:

No file chosen

Vertebrate  A. thaliana

**Restrictions:** At most 50 sequences and 1,000,000 nucleotides per submission; each sequence not more than 500,000 nucleotides.

**Confidentiality:** The sequences are kept confidential and will be deleted after processing.

### CITATIONS

For publication of results, please cite:

**Neural network prediction of translation initiation sites in eukaryotes: perspectives for EST and genome analysis.**

A. G. Pedersen and H. Nielsen, ISMB: 5, 226-233, 1997.

### PORTABLE VERSION

Would you prefer to run NetStart at your own site? NetStart 1.0 is available as a stand-alone software package, with the same functionality as the service above.

[download page](#) for academic users; other users are requested to contact CBS Software Package Manager at [software@cbs.dtu.dk](mailto:software@cbs.dtu.dk).

## NUCLEOTIDE SEQUENCES

### Whole genome visualization and analysis

« [GenomeAtlas](#) »

DNA structural atlases for complete microbial Genomes

### Gene finding and splice sites

« [EasyGene](#) »

Genes in prokaryotes

« [EasyGene](#) »

Genes in prokaryotes

« [HMMgene](#) »

Genes in eukaryotes

[MetaRanker](#)

Identification of risk genes in complex phenotypes

[NetAspGene](#)

Intron splice sites in Aspergillus DNA

« [NetGene2](#) »

Intron splice sites in human, C. elegans and A. thaliana DNA

[NetPlantGene](#)

Intron splice sites in Arabidopsis thaliana DNA

« [NetStart](#) »

Translation start in vertebrate and A. thaliana DNA

[NetUTR](#)

Splice sites in 5' UTR regions of human genes

« [Promoter](#) »

Transcription start sites in vertebrate DNA

« [RNAmmer](#) »

Ribosomal RNA sub units

« [RNAmmer](#) »

Ribosomal RNA sub units

### Analysis of DNA microarray data

[GenePublisher](#)

Analysis of DNA microarray data

« [OligoWiz](#) »

Design of oligonucleotides for DNA microarrays

## SMALL MOLECULES

« [ChemProt](#) »

Chemical-protein interactions

## AMINO ACID SEQUENCES

### Protein sorting

[ChloroP](#) »

Chloroplast transit peptides and their cleavage sites in plant proteins

[LipoP](#) »

Signal peptidase I & II cleavage sites in gram- bacteria

[NetNES](#) »

Leucine-rich nuclear export signals (NES) in eukaryotic proteins

[SecretomeP](#) »

Non-classical and leaderless secretion of proteins

[SignalP](#) »

Signal peptide and cleavage sites in gram+, gram- and eukaryotic amino acid sequences

[TargetP](#) »

Subcellular location of proteins: mitochondrial, chloroplastic, secretory pathway, or other

[TatP](#) »

Twin-arginine signal peptides

### Post-translational modifications of proteins

[DictyOGlyc](#)

O-(alpha)-GlcNAc glycosylation sites (trained on *Dictyostelium discoideum* proteins)

[NetAcet](#)

N-terminal acetylation in eukaryotic proteins

[NetCGlyc](#)

C-mannosylation sites in mammalian proteins

[NetCorona](#)

Coronavirus 3C-like proteinase cleavage sites in proteins

[NetGlycate](#) »

Glycation of ε amino groups of lysines in mammalian proteins

[NetNGlyc](#)

N-linked glycosylation sites in human proteins

[NetNGlyc](#) »

N-linked glycosylation sites in human proteins

[NetOGlyc](#) »

O-GalNAc (mucin type) glycosylation sites in mammalian proteins

[NetOGlyc](#) »

O-GalNAc (mucin type) glycosylation sites in mammalian proteins

[NetPhorest](#)

Linear motif atlas for phosphorylation-dependent signaling

[NetPhos](#) »

Generic phosphorylation sites in eukaryotic proteins

[NetPhosBac](#)

Generic phosphorylation sites in bacterial proteins

[NetPhosK](#)

Kinase specific phosphorylation sites in eukaryotic proteins

[NetPhosYeast](#)

Serine and threonine phosphorylation sites in yeast proteins

# SignalP 4.1 Server

SignalP 4.1 server predicts the presence and location of signal peptide cleavage sites in amino acid sequences from different organisms: Gram-positive prok prediction of cleavage sites and a signal peptide/non-signal peptide prediction based on a combination of several artificial neural networks.

View the [version history](#) of this server. All the previous versions are available online, for comparison and reference.

**New:** SignalP has been updated to version 4.1 with two new features:

- an option to choose a D-score cutoff that reproduces the sensitivity of SignalP 3.0 (this will make the false positive rate slightly higher, but still better than
- a customizable minimum length of the predicted signal peptide (default 10).

Additionally, the documentation has been rewritten. The [Instructions](#) page is expanded, the [Output format](#) page has been clarified, and there are new [Performance](#)

[FAQ](#)

[Article abstracts](#)

[Instructions](#)

[Output format](#)

## SUBMISSION

Paste a single amino acid sequence or several sequences in [FASTA](#) format into the field below:

```
GIDDWWRNEQFWVIGGVSSHLFALFQGLLKVLAGVSTNFTVTSKAADDGEFSELYIFK  
WTSLLIPPTTLINIVGVIVGVSDAINNGYDSWGPLFGRLFFALWVIVHLYPFLKGLLGK  
QDRVPTIILVWSILLASILTLLWVRVNPFVSKDGPVLEICGLDCLK|
```

Submit a file in [FASTA](#) format directly from your local disk:

No file chosen

**Organism group** ([explain](#))

- Eukaryotes
- Gram-negative bacteria
- Gram-positive bacteria

**D-cutoff values** ([explain](#))

- Default (optimized for correlation)
- Sensitive (reproduce SignalP 3.0's sensitivity)
- User defined:

D-cutoff for SignalP-noTM networks  
 D-cutoff for SignalP-TM networks

**Output format** ([explain](#))

- Standard
- Short (no graphics)
- Long
- All - SignalP-noTM and SignalP-TM output (no graphics)

**Method** ([explain](#))

- Input sequences may include TM regions
- Input sequences do not include TM regions

Pan-specific binding of peptides to MHC class I alleles of known sequence

### [NNAlign](#)

Identifying sequence motifs in quantitative peptide data

### [VDJsolver](#) »

Analysis of human immunoglobulin VDJ recombination

## **Protein function and structure**

### [ArchaeaFun](#)

Enzyme/non-enzyme and enzyme class (Archaea)

### [CPHmodels](#)

Protein structure from sequence: distance constraints

### [distanceP](#)

Protein distance constraints

### [EPipe](#) »

Functional differences of protein variants

### [InterMap3D](#)

Co-evolving amino acids in proteins

### [NetSurfP](#) »

Protein secondary structure and relative solvent accessibility

### [NetTurnP](#)

$\beta$ -turns and  $\beta$ -turn types in proteins

### [ProtFun](#) »

Protein functional category and enzyme class (Eukarya)

### [RedHom](#)

~~Reduction of sequence similarity in a data set~~

### [TMHMM](#) »

Transmembrane helices in proteins

### [VarDom](#)

Domains in the malaria antigen family PfEMP1

# TMHMM Server v. 2.0

## Prediction of transmembrane helices in proteins

NOTE: You can submit many proteins at once in one fasta file. Please limit each submission to at most 4000 proteins. Please tick the 'C

[Instructions](#)

### SUBMISSION

Submission of a local file in **FASTA** format (HTML 3.0 or higher)

No file chosen

OR by pasting sequence(s) in **FASTA** format:

```
INLSDRLHQVLRWALGSVEIFLSRHCPIWYGYGGGLKWL  
LERFSYINSVVYPWTSLLL  
VYCSLPAICLLTGKFIVPEISNYAGILFLLMFM  
SIAVTGILEMQWCKIGIDDWW  
WRNEQFWVIGGVSSH  
LFALFQQLKVLAGVSTNFTV  
TSKAADDGEFSELYIFKWT  
SLLIPPTTLLIINIVCV  
VIVCVSDAINNGYDS  
WGPLFGR  
LFFALWVIVHLYPFLK  
GLLGKQDRVPTIIL  
VWSILLASILTLL  
WVRVNP  
FVKDGPVLEI  
CGLDCLK|
```

Output format:

- Extensive, with graphics
- Extensive, no graphics
- One line per protein

Other options:

- Use old model (version 1)

---

**PORTABLE VERSION**

# Next class: JGI resource