# Canadian Bioinformatics Workshops

www.bioinformatics.ca

# creative commons

## Attribution-Share Alike 2.5 Canada

### You are free:

**to Share** — to copy, distribute and transmit the work

**to Remix** — to adapt the work

*Free Cultural Works*
**APPROVED FOR**

### Under the following conditions:

**Attribution**. You must attribute the work in the manner specified by the author or licensor (but not in any way that suggests that they endorse you or your use of the work).

**Share Alike**. If you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar licence to this one.

- For any reuse or distribution, you must make clear to others the licence terms of this work.
- Any of the above conditions can be waived if you get permission from the copyright holder.
- The author's moral rights are retained in this licence.

Disclaimer

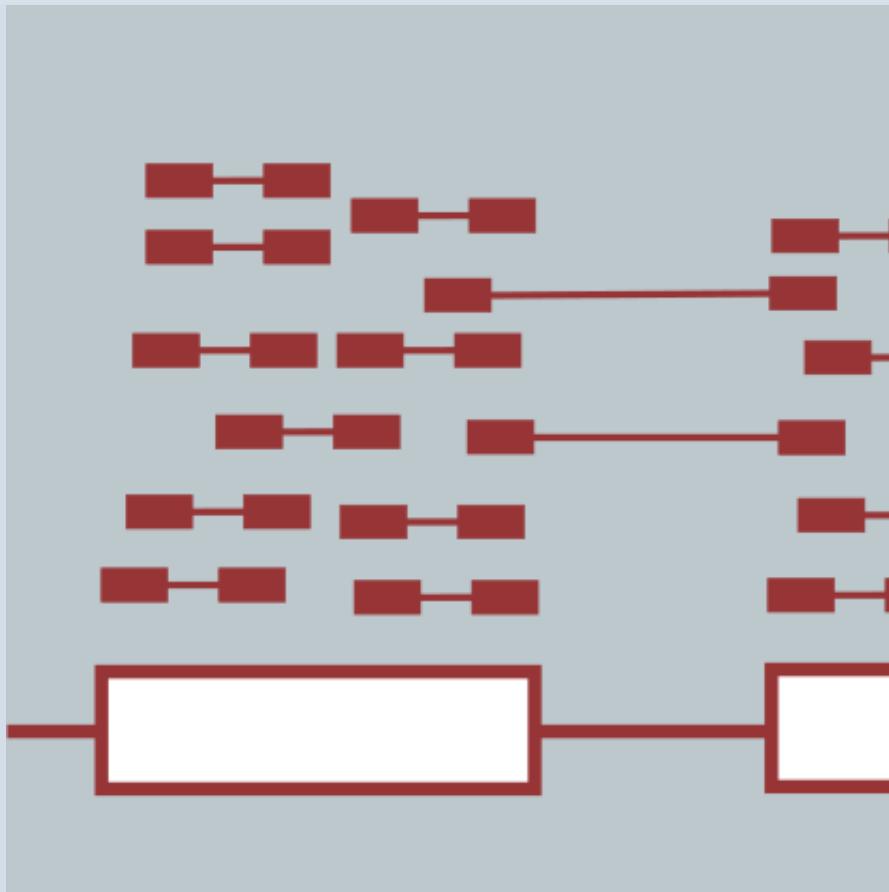Your fair dealing and other rights are in no way affected by the above.
This is a human-readable summary of the Legal Code (the full licence) available in the following languages:
English French

Learn how to distribute your work using this licence

# Functional Annotation and Analysis of Transcripts

Brian Haas

Informatics for RNA-Seq Analysis

July 10-12, 2017

bioinformatics.ca

# Learning Objectives of Module

- Explore methods to glean biological function from transcript sequences.

- Differentiate between homology-based and sequence composition-based functional inference.

# Transcript Functional Annotation

```
GGAGCTGGAGGCCCCCAGGCAACTACACCGTCCACGTACCCAGAGGGGCTGGGCCCTCCC
ACCAGAGACCACGCCCTGGTGTGCCTTAGGGGCCCTGGTTTGTTAGTCTCTGAGTGTGCA
GTTGCTGCACATGGGGCCCTGGCGCTTGCTGCACCAACTTCCTGTTGGGCCCGTGGTCCT
TGGAGGCATGCAGTTCAGCAGACAGTGACTCAGCCATCCACCCAACATGCGGAACGTGTC
TCTTCTGCAGGTCCCGGTCCACAGCAGGATTCCCCTCTGTGAAAGGCACGCTGATCTG
TCTGGA                                                  TCGAC
TCTCCG                                                  TCCCA
AAAGAC                                                  CCTGG
GGCTTC                                                  CCTAA
TGACCT                                                  TGCTG
GAAAAG                                                  CAGCC
TTGTCA                                                  TTCCA
GGAAGCACATAATTGAAGGACTGAAAGCGTCCCTGGAGCGGCTGCAGCTGGAGTACGTGG
ATGTGGTTTTTGCCAACCGCCCAGACCCCAACACGCCCATGGAAGAGACCGTGCGGGCCA
TGACCCATGTCATCAACCAGGGGATGGCCATGTACTGGGGCACATCACGCTGGAGCTCCA
TGGAGATCATGGAGGCCTACTCGGTGGCTCGGCAGTTCAACCTGATCCCGCCCATCTGCG
AGCAAGCGGAATATCACATGTTCCAGAGGGAGAAGGTGGAGGTCCAGCTGCCAGAGCTGT
TCCACAAGATAGGAGTAGGTGCCATGACCTGGTCCCCTCTGGCGTGCGGCATCGTCTCAG
GGAAGTATGACAGCGGGATCCCACCCTACTCCAGAGCCTCCCTGAAGGGCTACCAGTGGT
TGAAGGACAAGATCCTGAGTGAGGAGGGTCGCCGCCAGCAGGCCAAGCTGAAGGAACTGC
```

## Can we gather hints of biological function from sequence?

# Methods used to predict function from sequence

- ## Sequence homology
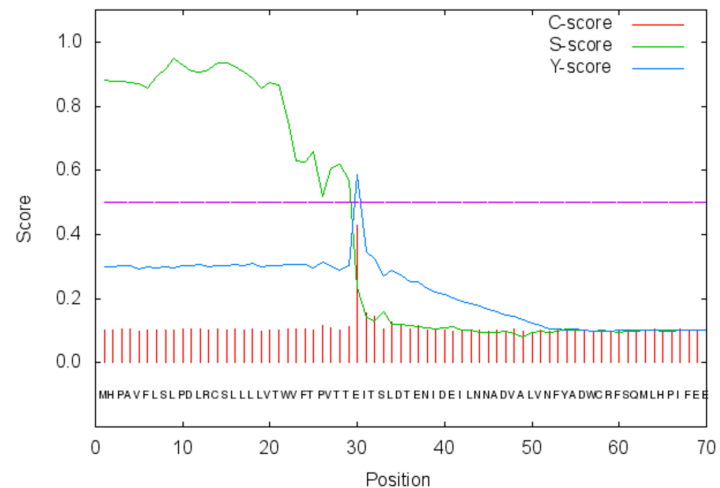
Searching protein database for sequence similarity

```
Query          THVHRPYNEHKSLSGTARYMSINTHLGREQSRRDDLESMGHVFMYFLRGSLPW--QGLKA
               T   P +  K   GT  Y S + HLG    RR DLE +G       L   LPW  Q L A
Database Match TGDFKP-DPKKMHNGTIEYTSRDAHLG-VPTRRADLEILGYNLIEWLGAELPWVTQKLLA
```

- ## Sequence composition

Predict functions of sequence using machine learning methods for pattern recognition.
- Neural Networks
- Hidden Markov Models

# Use BLAST to search for sequence similarity to known proteins

# The Swiss-Prot database is a valuable source of proteins with known functions



(as of July, 2017)

# Example of a Swiss-Prot Record



**Gene Ontology (GO)**: Structured vocabulary for defining molecular functions, biological processes, and cellular components.

**bio**informatics.ca

# Gene Ontology: a structured relational vocabulary for describing biological functions



Gene Ontology terms are organized into a directed acyclic graph. Terms are organized from general (top) to more specific (bottom).

The GO structure enables computations such as exploring function enrichment among sets of transcripts.

# Gene ontology functional enrichment

|  | (+) Differentially Expressed | (-) Not Differentially Expressed | Totals |
|---|---|---|---|
| + Gene Ontology | 50 | 200 | 250 |
| - Gene Ontology | 1950 | 17800 | 19750 |
| Totals | 2000 | 18000 | 20000 |

|  | drawn | not drawn | total |
|---|---|---|---|
| green marbles | $k$ | $K - k$ | $K$ |
| red marbles | $n - k$ | $N + k - n - K$ | $N - K$ |
| total | $n$ | $N - n$ | $N$ |

The probability of drawing exactly $k$ green marbles can be calculated by the formula

$$P(X = k) = f(k; N, K, n) = \frac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}.$$

# No significant sequence similarity…  What else?

```
GGAGCTGGAGGCCCCCAGGCAACTACACCGTCCACGTACCCAGAGGGGCTGGGCCCTCCC
ACCAGAGACCACGCCCTGGTGTGCCTTAGGGGCCCTGGTTTGTTAGTCTCTGAGTGTGCA
GTTGCTGCACATGGGGCCCTGGCGCTTGCTGCACCAACTTCCTGTTGGGCCCGTGGTCCT
TGGAGGCATGCAGTTCAGCAGACAGTGACTCAGCCATCCACCCAACATGCGGAACGTGTC
TCTTCTGCAGGTCCCGGTCCACAGCAGGATTCCCCTCTGTGAAAAGGCACGCTGATCTG
TCTGGATAAGTGTGGCCGGCCCCATGTATCCGGAATCAACCACGGGGTCCCCAGCTCGAC
TCTCCCTGCGGCAGACAGGCTCCCCGGGATGATCTACAGTACTCGTTATGGGAGTCCCA
AAAGACAGCTCCAGTTTTACAGGAATCTGGGCAAATCTGGCCTTCGGGTCTCCTGCCTGG
GGCTTGGAACATGGGTGACCTTCGGGGGCCAGATCACGGATGAGATGGCAGAGCACCTAA
TGACCTTGGCCTACGATAATGGCATCAACCTGTTCGATACGGCGGAGGTCTACGCTGCTG
GAAAAGCTGAAGTGGTATTAGGGAACATCATTAAGAAGAAGGGATGGAGACGGTCCAGCC
TTGTCATCACCACCAAGATCTTCTGGGGTGGAAAAGCGGAGACTGAGAGAGGCCTTTCCA
GGAAGCACATAATTGAAGGACTGAAAGCGTCCCTGGAGCGGCTGCAGCTGGAGTACGTGG
ATGTGGTTTTTGCCAACCGCCCAGACCCCAACACGCCCATGGAAGAGACCGTGCGGGCCA
TGACCCATGTCATCAACCAGGGGATGGCCATGTACTGGGGCACATCACGCTGGAGCTCCA
TGGAGATCATGGAGGCCTACTCGGTGGCTCGGCAGTTCAACCTGATCCCGCCCATCTGCG
AGCAAGCGGAATATCACATGTTCCAGAGGGAGAAGGTGGAGGTCCAGCTGCCAGAGCTGT
TCCACAAGATAGGAGTAGGTGCCATGACCTGGTCCCCTCTGGCGTGCGGCATCGTCTCAG
GGAAGTATGACAGCGGGATCCCACCCTACTCCAGAGCCTCCCTGAAGGGCTACCAGTGGT
TGAAGGACAAGATCCTGAGTGAGGAGGGTCGCCGCCAGCAGGCCAAGCTGAAGGAACTGC
```

# Is there an ORF for a potential Coding Region?

```
GGAGCTGGAGGCCCCCAGGCAACTACACCGTCCACGTACCCAGAGGGGCTGGGCCCTCCC
ACCAGAGACCACGCCCTGGTGTGCCTTAGGGGCCCTGGTTTGTTAGTCTCTGAGTGTGCA
GTTGCTGCACATGGGGCCCTGGCGCTTGCTGCACCAACTTCCTGTTGGGCCCGTGGTCCT
TGGAGGCATGCAGTTCAGCAGACAGTGACTCAGCCATCCACCCAACATGCGGAACGTGTC
TCTTCTGCAGGTCCCGGTCCACAGCAGGATTCCCCTCTGTGAAAGGCACGCTGATCTG
TCTGGATAAGTGTGGCCGGCCCCATGTATCCGGAATCAACCACGGGGTCCCCAGCTCGAC
TCTCCCTGCGGCAGACAGGCTCCCCGGGATGATCTACAGTACTCGTTATGGGAGTCCCA
AAAGACAGCTCCAGTTTTACAGGAATCTGGGCAAATCTGGCCTTCGGGTCTCCTGCCTGG
GGCTTGGAACATGGGTGACCTTCGGGGGCCAGATCACGGATGAGATGGCAGAGCACCTAA
TGACCTTGGCCTACGATAATGGCATCAACCTGTTCGATACGGCGGAGGTCTACGCTGCTG
GAAAAGCTGAAGTGGTATTAGGGAACATCATTAAGAAGAAGGGATGGAGACGGTCCAGCC
TTGTCATCACCACCAAGATCTTCTGGGGTGGAAAAGCGGAGACTGAGAGAGGCCTTTCCA
GGAAGCACATAATTGAAGGACTGAAAGCGTCCCTGGAGCGGCTGCAGCTGGAGTACGTGG
ATGTGGTTTTTGCCAACCGCCCAGACCCCAACACGCCCATGGAAGAGACCGTGCGGGCCA
TGACCCATGTCATCAACCAGGGGATGGCCATGTACTGGGGCACATCACGCTGGAGCTCCA
TGGAGATCATGGAGGCCTACTCGGTGGCTCGGCAGTTCAACCTGATCCCGCCCATCTGCG
AGCAAGCGGAATATCACATGTTCCAGAGGGAGAAGGTGGAGGTCCAGCTGCCAGAGCTGT
TCCACAAGATAGGAGTAGGTGCCATGACCTGGTCCCCTCTGGCGTGCGGCATCGTCTCAG
GGAAGTATGACAGCGGGATCCCACCCTACTCCAGAGCCTCCCTGAAGGGCTACCAGTGGT
TGAAGGACAAGATCCTGAGTGAGGAGGGTCGCCGCCAGCAGGCCAAGCTGAAGGAACTGC
```
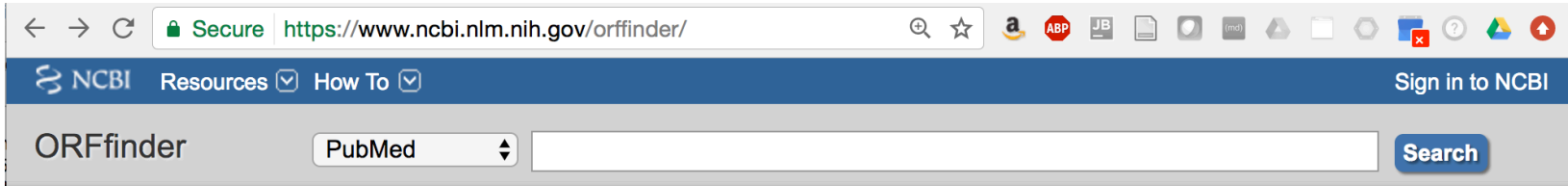
# Is there an ORF for a potential Coding Region?

GGAGCTGGAGGCCCCCAGGCAACTACACCGTCCACGTACCCAGAGGGGCTGGGCCCTCCC
ACCAGAGACCACGCCCTGGTGTGCCTTAGGGGCCCTGGTTTGTTAGTCTCTGAGTGTGCA
GTTGCTGCAC**ATGGGGCCCTGGCGCTTGCTGCACCAACTTCCTGTTGGGCCCGTGGTCCT**
**TGGAGGCATGCAGTTCAGCAGACAGTGA**CTCAGCCATCCACCCAACATGCGGAACGTGTC
TCTTCTGCAGGTCCCGGTCCACAGCAGGATTCCCCCTCTGTGAAAAGGCACGCTGATCTG
TCTGGATAAGTGTGGCCGGCCCCATGTATCCGGAATCAACCACGGGGTCCCCAGCTCGAC
TCTCCCTGCGGCAGACAGGCTCCCCGGGATGATCTACAGTACTCGTTATGGGAGTCCCA
AAAGACAGCTCCAGTTTTACAGGAATCTGGGCAAATCTGGCCTTCGGGTCTCCTGCCTGG
GGCTTGGAACATGGGTGACCTTCGGGGGCCAGATCACGGATGAGATGGCAGAGCACCTAA
TGACCTTGGCCTACGATAATGGCATCAACCTGTTCGATACGGCGGAGGTCTACGCTGCTG
GAAAAGCTGAAGTGGTATTAGGGAACATCATTAAGAAGAAGGGATGGAGACGGTCCAGCC
TTGTCATCACCACCAAGATCTTCTGGGGTGGAAAGCGGAGACTGAGAGAGGCCTTTCCA
GGAAGCACATAATTGAAGGACTGAAAGCGTCCCTGGAGCGGCTGCAGCTGGAGTACGTGG
ATGTGGTTTTTGCCAACCGCCCAGACCCCAACACGCCCATGGAAGAGACCGTGCGGGCCA
TGACCCATGTCATCAACCAGGGGATGGCCATGTACTGGGGCACATCACGCTGGAGCTCCA
TGGAGATCATGGAGGCCTACTCGGTGGCTCGGCAGTTCAACCTGATCCCGCCCATCTGCG
AGCAAGCGGAATATCACATGTTCCAGAGGGAGAAGGTGGAGGTCCAGCTGCCAGAGCTGT
TCCACAAGATAGGAGTAGGTGCCATGACCTGGTCCCCTCTGGCGTGCGGCATCGTCTCAG
GGAAGTATGACAGCGGGATCCCACCCTACTCCAGAGCCTCCCTGAAGGGCTACCAGTGGT
TGAAGGACAAGATCCTGAGTGAGGAGGGTCGCCGCCAGCAGGCCAAGCTGAAGGAACTGC

# Find all ORFs using ORFfinder

# ORFfinder finds all open reading frames and provides translations



**Open Reading Frame Viewer**

**Sequence**

ORFs can appear in random sequence – so further analysis is required

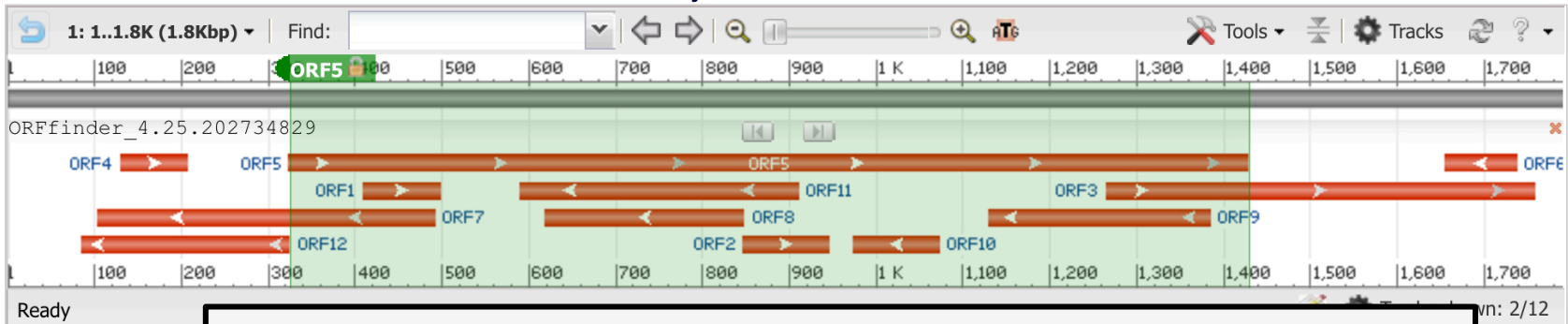ORFs found: **12**    Genetic code: **1**    Start codon: **'ATG' only**

ORFfinder_4.25.202734829

Predict coding vs. non-coding ORFs:  http://TransDecoder.github.io
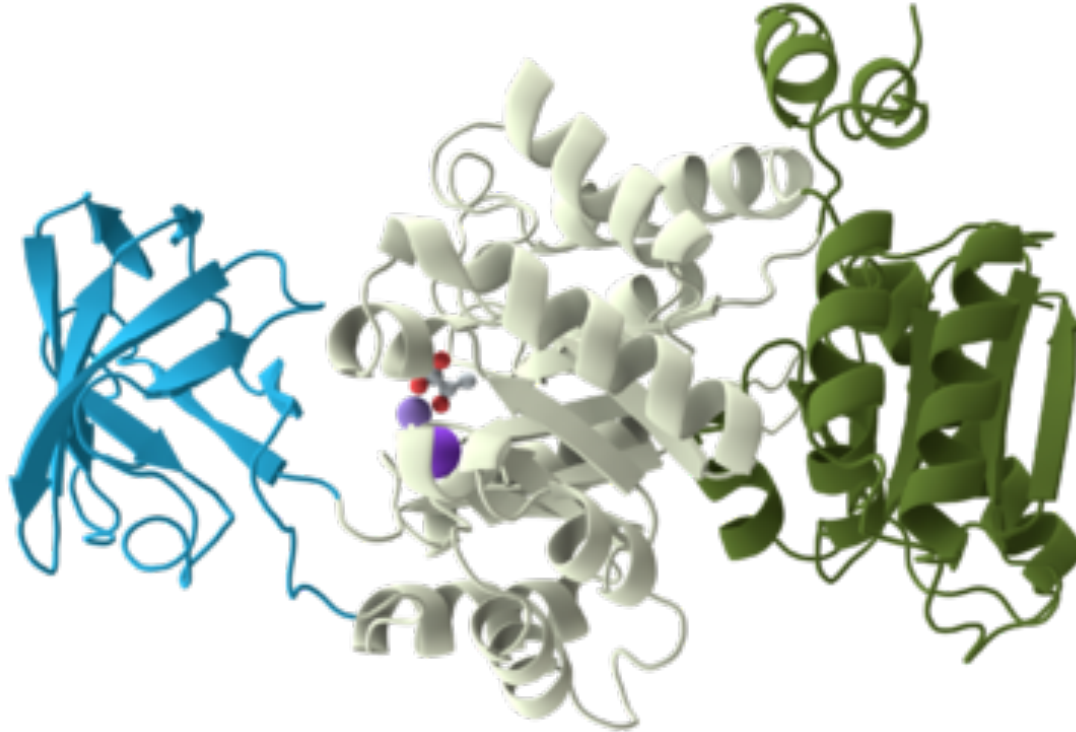
Add six-frame translation track

**ORF5**  (367 aa)    **Display ORF as...**    Mark

```
>lcl|ORF5
MYPESTTGSPARLSLRQTGSPGMIYSTRYGSPKRQLQFYR
NLGKSGLRVSCLGLGTWVTFGGQITDEMAEHLMTLAYDNG
INLFDTAEVYAAGKAEVVLGNIIKKKGWRRSSLVITTKIF
WGGKAETERGLSRKHIIEGLKASLERLQLEYVDVVFANRP
DPNTPMEETVRAMTHVINQGMAMYWGTSRWSSMEIMEAYS
VARQFNLIPPICEQAEYHMFQREKVEVQLPELFHKIGVGA
MTWSPLACGIVSGKYDSGIPPYSRASLKGYQWLKDKILSE
EGRRQQAKLKELQAIAERLGCTLPQLAIAWCLRNEGVSSV
LLGASNAEQLMENIGAIQVLPKLSSSIVHEIDSILGNKPY
SKKDYRS
```

**Mark subset...**    Marked: 0    Download marked set    as    Protein FA

| Label | Strand | Frame | Start | Stop | Length (nt \| a |
|-------|--------|-------|-------|------|------------------|
| **ORF5** | **+** | **3** | **324** | **1427** | **1104 \| 36** |
| ORF3 | + | 1 | 1264 | 1758 | 495 \| 16 |
| ORF7 | - | 1 | 492 | 103 | 390 \| 12 |
| ORF11 | - | 3 | 910 | 590 | 321 \| 10 |
| ORF9 | - | 3 | 1384 | 1130 | 255 \| 8 |
| ORF12 | - | 3 | 325 | 86 | 240 \| 7 |

# Can we recognize functional domains in putative coding regions?



Hints at <u>substrate binding</u> or <u>catalytic activity</u>

| DNA, RNA, calcium, phoshate, etc. |
|---|

| Glycoslase, methylase, kinase, nuclease, lipase, protease, etc. |
|---|

**bio**informatics.ca

# Search the Pfam library of HMMs to identify potential functional domains

# Example Pfam report illustrating modular domain architecture



**Sequence search results**

Show the detailed description of this results page.

We found **9** Pfam-A matches to your search sequence (**all** significant)

Show the search options and sequence that you submitted.

Return to the search form to look for Pfam domains on a new sequence.

## Significant Pfam-A Matches

Show or hide all alignments.

| Family | Description | Entry type | Clan | Envelope | | Alignment | | HMM | | HMM length | Bit score | E-value | Predicted active sites | Show/hide alignment |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Start | End | Start | End | From | To | | | | | |
| CUB | CUB domain | Domain | CL0164 | 93 | 206 | 93 | 206 | 1 | 110 | 110 | 42.2 | 7.7e-11 | n/a | Show |
| EGF_2 | EGF-like domain | Domain | CL0001 | 249 | 280 | 249 | 280 | 1 | 32 | 32 | 22.5 | 0.0001 | n/a | Show |
| Kelch_5 | Kelch motif | Repeat | CL0186 | 351 | 393 | 352 | 392 | 2 | 41 | 42 | 33.7 | 2.2e-08 | n/a | Show |
| Kelch_4 | Galactose oxidase, central domain | Repeat | CL0186 | 466 | 518 | 468 | 514 | 3 | 44 | 49 | 20.6 | 0.0003 | n/a | Show |
| Kelch_1 | Kelch motif | Repeat | CL0186 | 520 | 574 | 520 | 573 | 1 | 45 | 46 | 20.0 | 0.00033 | n/a | Show |
| Kelch_5 | Kelch motif | Repeat | CL0186 | 579 | 614 | 581 | 613 | 5 | 40 | 42 | 25.3 | 9.7e-06 | n/a | Show |
| Lectin_C | Lectin C-type domain | Domain | CL0056 | 765 | 874 | 766 | 874 | 2 | 108 | 108 | 70.2 | 2e-19 | n/a | Show |
| PSI | Plexin repeat | Family | CL0630 | 889 | 939 | 890 | 938 | 2 | 50 | 51 | 27.8 | 2.5e-06 | n/a | Show |
| PSI | Plexin repeat | Family | CL0630 | 942 | 1012 | 942 | 1012 | 1 | 51 | 51 | 50.0 | 2.9e-13 | n/a | Show |

Comments or questions on the site? Send a mail to **pfam-help@ebi.ac.uk**.
**E u r o p e a n   M o l e c u l a r   B i o l o g y   L a b o r a t o r y**

Module

**bio**informatics.ca

# Transmembrane Proteins

Membrane

1

2

3

Single transmembrane α-helix
(bitopic membrane protein)

Polytopic
transmembrane
α-helical protein

Polytopic
transmembrane β-
sheet protein

# Using TMHMM to identify putative transmembrane proteins

# Trans-membrane Domains via TmHMM



TMHMM posterior probabilities for WEBSEQUENCE

Topology=i36-55o59-81i93-110o125-147i174-196o206-228i241-260o280-302i309-328o338-360i373-395o448-467i

http://www.cbs.dtu.dk/services/TMHMM/

**bio**informatics.ca

# Predicting Secreted Proteins



(from: Vaccine 23(15):1770-8)

(from: https://courses.washington.edu/conj/cell/secretion.htm)

# SignalP: Prediction of N-terminal signal peptides
## (predict secreted proteins)

CENTERFOR RBIOLOGI CALSEQU ENCEANA LYSIS CBS

**EVENTS**

**NEWS**

**RESEARCH GROUPS**

**CBS PREDICTION SERVERS**

**CBS DATA SETS**

**PUBLICATIONS**

**EDUCATION**

**STAFF**

**CONTACT**

**ABOUT CBS**

**INTERNAL**

**CBS BIOINFORMATICS TOOLS**

**CBS COURSES**

**OTHER BIOINFORMATICS LINKS**

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS ∷ TECHNICAL UNIVERSITY OF DENMARK DTU

CBS >> **CBS Prediction Servers** >> **SignalP**

## 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 prokaryotes, Gram-negative prokaryotes, and eukaryotes. The method incorporates a 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 on line, for comparison and reference.

**NEW:** The portable version of SignalP 4.1, previously only available for Mac (Darwin), Linux, and IRIX, is now also available for Windows systems. Academic users: select the "CYGWIN" option at the download page. Cygwin or MobaXterm is required to install SignalP under Windows. For details, read the installation instructions.

| FAQ | Article abstracts | Instructions | Output format | Performance | Data |
|---|---|---|---|---|---|

### SUBMISSION

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

MHPAVFLSLPDLRCSLLLLVTWVFTPVTTEITSLDTENIDEILNNADVALVNFYADWCRFSQMLHPIFEEASDVIKEEFPNENQVVFARVDCDQHSDIAQRYRISKYPTLKLFRNGMMM
KREYRGQRSVKALADYIRQQKSDPIQEIRDLAEITTLDRSKRNIIGYFEQKDSDNYRVFERVANILHDDCAFLSAFGDVSKPERYSGDNIIYKPPGHSAPDMVYLGAMTNFDVTYNWIQ
DKCVPLVREITFENGEELTEEGLPFLILFHMKEDTESLEIFQNEVARQLISEKGTINFLHADCDKFRHPLLHIQKTPADCPVIAIDSFRHMYVFGDFKDVLIPGKLKQFVFDLHSGKLHREF
HHGPDPTDTAPGEQAQDVASSPPESSFQKLAPSEYRYTLLRDRDEL

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

Choose File | 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:
  - `0.4` *D-cutoff for SignalP-noTM networks*
  - `0.5` *D-cutoff for SignalP-TM networks*

**Graphics output** (explain)
- ○ No graphics
- ● **PNG** (inline)
- ○ **PNG** (inline) and **EPS** (as links)

**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

**Positional limits** (explain)
- ☐ Minimal predicted signal peptide length. *Default: 10*
- ☐ N-terminal truncation of input sequence (0 means no truncation). *Default: Truncate sequence to a length of 70 aa*

**Module**

**bioinformatics**.ca

# Example SignalP predicted signal peptide

# Transcriptome-scale functional annotation using Trinotate

**bio**informatics.ca

# There's no substitute for experimentally validating protein functions

# We are on a Coffee Break & Networking Session

**bio**informatics.ca