

Sequence analysis

DNA, RNA & proteins



ILLINOIS INSTITUTE OF TECHNOLOGY

Jean-François Pombert, Ph.D.
Office PS 296 (Lab PS 340)

Sequence analysis – Topics

- Sequence alignment
- Homology search
- Types of sequencing data
- Read mapping + variant calling
- Genome assembly
- Genome annotation

Course section:

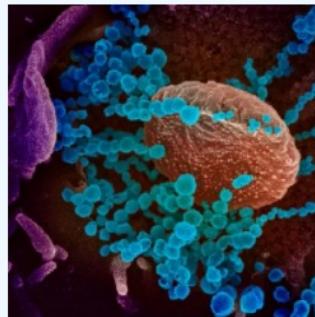
05 – Sequence analysis

06 – Sequencing data

07 – Genome Annotation

Featured communities

Need help uploading? Contact us

**Chicago COVID-19 Response**[Browse](#) [New upload](#)

This repository community collects research outputs and information objects relevant to the COVID-19 / SARS-CoV-2 efforts in Chicago. Users are encouraged to upload their research objects in this collection to facilitate sharing and discovery of information. Although Open Access articles and...

Curated by: saragon

Open source data

<https://zenodo.org/>

Data **accessibility** is becoming essential for publication in high impact journals

NCBI databases - <https://www.ncbi.nlm.nih.gov/>

National Institutes of Health (NIH)

Zenodo - <https://zenodo.org/>

European Organization for Nuclear Research (CERN)



How developers work

Support your workflow with lightweight tools and features. Then work how you work best—we'll follow your lead.

New to GitHub? See how it works



Code review



Project management



Integrations



Team management



Social coding



Documentation



Code hosting

<https://github.com/>

Open source software (OSS)

The standard in scientific publications

Many, many programs available

Free, various licenses (MIT, GNU...)

Source-code available <-> can be modified

Often required for publication in high impact journals

EMBOSS Programs

[!\[\]\(8bba887393ca45b761e5cb49e755e762_img.jpg\) Feedback](#)

Tools > EMBOSS Programs

<https://www.ebi.ac.uk/Tools/emboss/>

EMBOSS

European Molecular Biology Open Software Suite

Free

Other software packages in the 1990s were often very expensive

On Redhat/Fedora:

dnf install EMBOSS

On Debian/Ubuntu:

apt-get install emboss

EMBOSS: the European Molecular Biology Open Software Suite

Rice P, Longden I, Bleasby A

Trends Genet. 2000. 16(6):276-7

DOI: 10.1016/s0168-9525(00)02024-2

<http://emboss.sourceforge.net/>



5.30.0



That's why we love Perl

25,000 extensions on CPAN

Perl is a highly capable, feature-rich programming language with over 30 years of development.



DOWNLOAD AND GET STARTED



<https://www.perl.org/>



For loop on a listi series up to n

```
>>> numbers = [2, 4, 6, 8]
>>> product = 1
>>> for number in numbers:
...     product *= number
...
...     print('The product is:', product)
The product is: 384
0 1 1 2 3 5 8 13 21 34 55 89 144 233 377 610 987
```

All the Flow You'd Expect

Python knows the usual control flow statements that other languages speak — `if`, `for`, `while` and `range` — with some of its own twists, of course. [More control flow tools in Python 3](#)

1 2 3 4 5

<https://www.python.org/>

Perl & Python

Scripting languages

Bio-specific libraries are available

Helpful communities

Many tools are available on [GitHub](#)

Design/adapt software to your specific needs

BioPerl, BioPython



IUPAC Top Ten

Emerging Technologies in Chemistry

Enter now your nomination for the 2020 Call

Let's showcase the value of Chemistry (and chemists!) and illustrate how the chemical sciences contribute to the well-being of Society and the sustainability of Planet Earth.

#IUPACTopTen

IUPAC – Standard codes

International **U**nion of **P**ure and **A**pplied **C**hemistry

<https://iupac.org/>

<http://www.bioinformatics.org/sms2/iupac.html>

About understanding algorithms

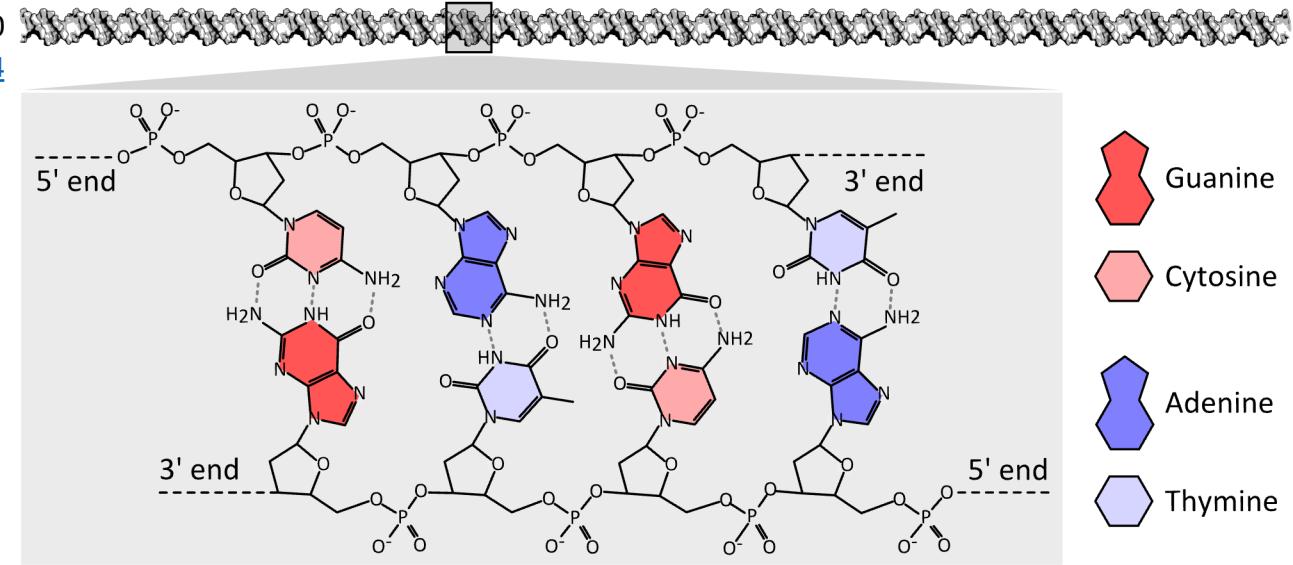
This is a hands-on class

We'll focus on the **general ideas**

Don't worry about the math

Alignments & homology

Chicken or the egg?



WARNING! – Orientation

Sequences can be **reversed** and/or **complemented**

5' – 3'; common

Most sequencers work from 5' to 3' (direction of synthesis)

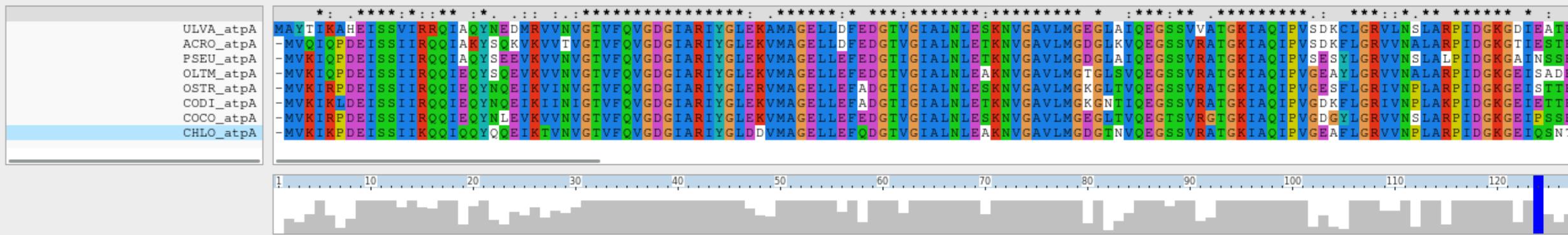
3' – 5'; very rare

Maxam-Gilbert's degradation sequencing

Wrong strand?

Common in databases, especially for RNA-encoding genes

Mode: Multiple Alignment Mode Font: 12

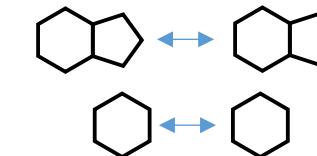


CLUSTAL-Alignment file created [/home/jpombert/BIOL550/05_Sequence_analysis/Ex_01/atpA.aln]

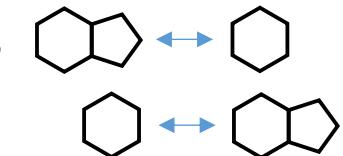
Identity vs. homology

DNA/RNA – transitions, transversions

Transitions



Transversions



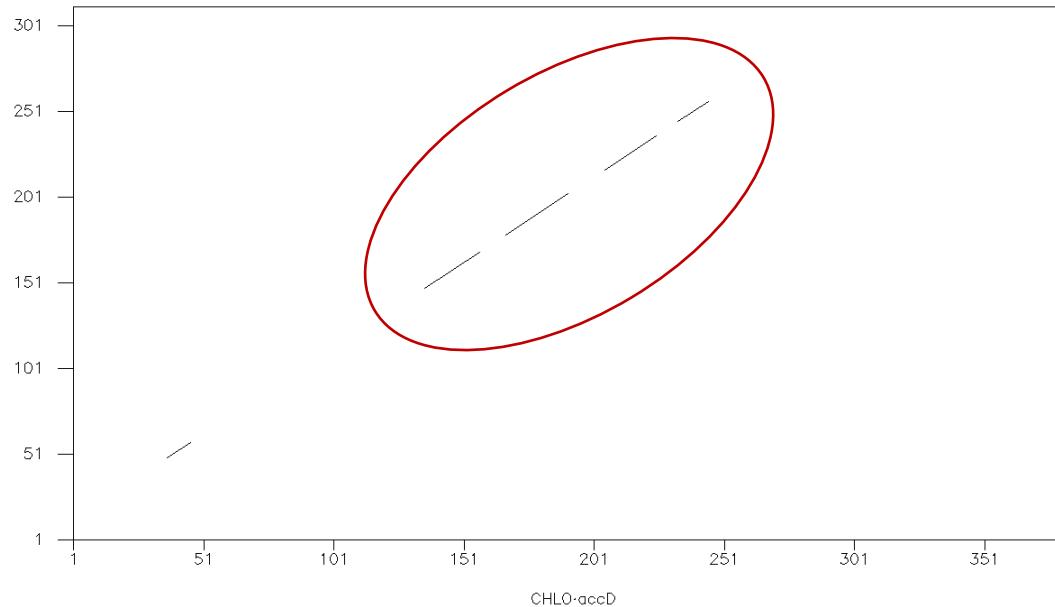
Proteins – amino acid (aa) properties:

Aromatic, aliphatic, hydrophobic, polar, acidic, basic? ## aa property [chart](#) from Sigma-Aldrich

Identity dot plots – Easy to visualize

Dottup: fasta::accD-chlo.fasta:CHLO-accD vs fasta::accD...

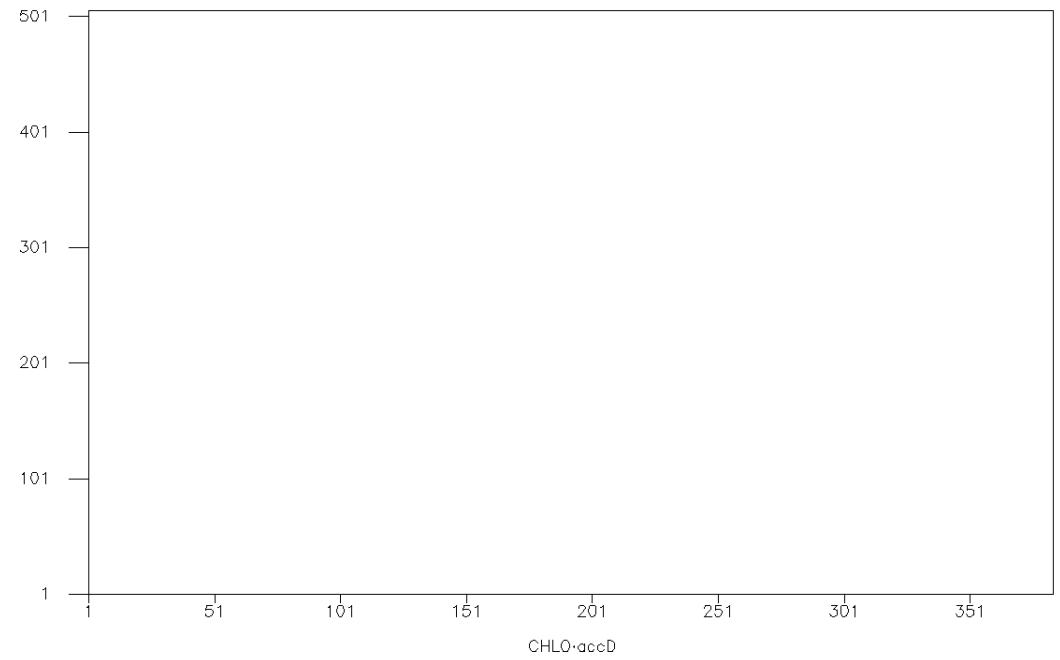
Mon 2 Sep 2013 14:25:38



Decent similarity

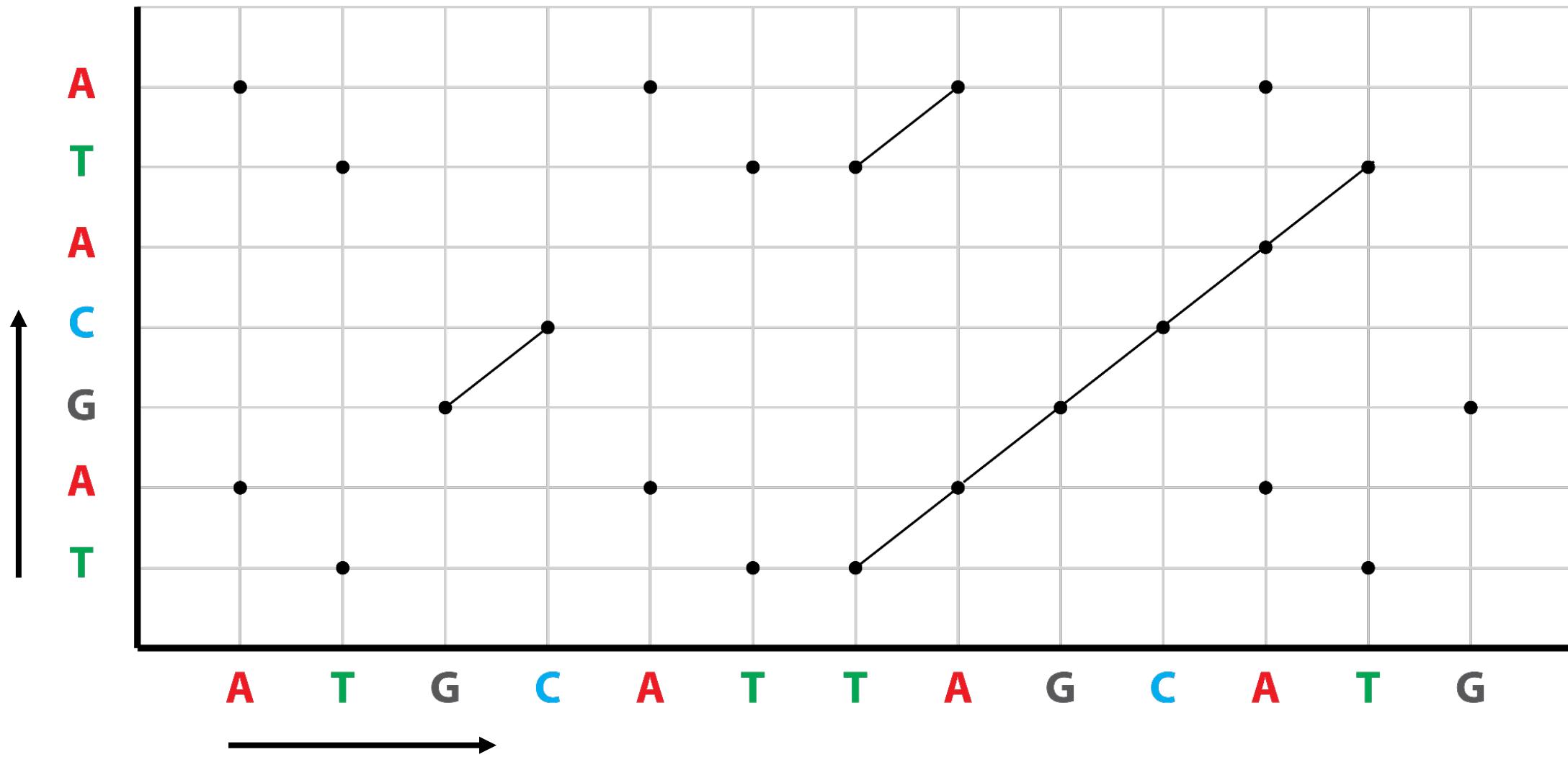
Dottup: fasta::accD-chlo.fasta:CHLO-accD vs fasta::atpA...

Mon 2 Sep 2013 14:27:03



No similarity (obviously)

Dot plots – How do they work?



Gaps – Phased identity/homology increase

ATGTCTTACTG
|| x | x x | x

ATTTACTG

Not very similar

4 identities

4 mismatches

0 gap

ATGTCTTACTG
|| x x x | | | | |

AT---TTACTG

Better

8 identities

0 mismatch

3 gaps (codon?)

Gaps – Maximise homology, minimize holes

ATGTCTTACTG
| x |
AAG

Better
2 identities
1 mismatch
0 gap

ATGTCTTACTG
| x x x x x | x x |
A-----A--G

Illogical
3 identities
0 mismatch
7 gaps

We need scoring matrices

Why?

So that we can attribute a score to alignments

Algorithms will return the highest scoring alignments

A simple identity matrix

Points are awarded for a match

Mismatches get nothing (or penalties)

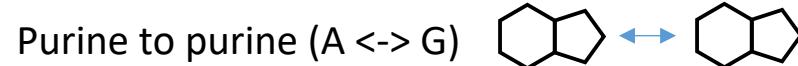
Values are arbitrary and can be modified

	A	T	G	C
A	1	0	0	0
T		1	0	0
G			1	0
C				1

A transition/transversion matrix

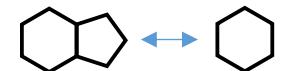
Transitions are more likely to occur than transversions

Transitions get lower penalties



Transversions get higher penalties

Purine \leftrightarrow pyrimidine (A or $G \leftrightarrow C$ or T)



Values are arbitrary and can be modified

	A	T	G	C
A	1	-1	0	-1
T		1	-1	0
G			1	-1
C				1

Amino acids tables are more complex

BLOSUM62

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	B	J	Z	X	*			
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0	-2	-1	-1	-1	-1	-4		
R		5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3	-1	-2	0	-1	-1	-4		
N			6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3	4	-3	0	-1	-1	-4		
D				6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3	4	-3	1	-1	-1	-4		
C					9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1	-3	-1	-1	-3	-1	-4		
Q						5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2	0	-2	4	-1	-4			
E							5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2	1	-3	4	-1	-4			
G								6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3	-1	-4	-2	-1	-1	-4		
H									8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3	0	-3	0	-1	-1	-4		
I										4	2	-3	1	0	-3	-2	-1	-3	-1	3	-3	3	-3	-1	-1	-4		
L											4	-2	2	0	-3	-2	-1	-2	-1	1	-4	3	-3	-1	-1	-4		
K												5	-1	-3	-1	0	-1	-3	-2	-2	0	-3	1	-1	-1	-4		
M													5	0	-2	-1	-1	-1	1	-3	2	-1	-1	-1	-4			
F														6	-4	-2	-2	1	3	-1	-3	0	-3	-1	-1	-4		
P															7	-1	-1	-4	-3	-2	-2	-3	-1	-1	-1	-4		
S																4	1	-3	-2	-2	0	-2	0	-1	-1	-1	-4	
T																	5	-2	-2	0	-1	-1	-1	-1	-1	-4		
W																		11	2	-3	-4	-2	-2	-1	-1	-1	-4	
Y																			7	-1	-3	-1	-2	-1	-1	-1	-4	
V																				4	-3	2	-2	-1	-1	-4		
B																					4	-3	0	-1	-1	-4		
J																						3	-3	-1	-1	-4		
Z																							4	-1	-1	-4		
X																								-1	-1	-4		
*																										1		

Extra non-IUPAC groups

1992. PNAS. Henikoff & Henikoff. 89:10915-10919

2004. Nature Biotechnology. Eddy. 22(8):1035-1036 <https://www.nature.com/articles/nbt0804-1035>

B Asx
Z Glx

asparagi(ne)c acid
glutami(ne)c acid

L Lxe
X Xaa

(iso)leucine
unknown

Properties of Common Amino Acids

Name	Abbr.	Molecular Weight	Molecular Formula	Residue Formula	Residue Weight (-H ₂ O)	pK _a ¹	pK _b ²	pK _x ³	pI ⁴	
Alanine	Ala	A	89.10	C ₃ H ₇ NO ₂	C ₃ H ₇ NO	71.08	2.34	9.69	—	6.00
Arginine	Arg	R	174.20	C ₆ H ₁₄ N ₄ O ₂	C ₆ H ₁₂ N ₄ O	156.19	2.17	9.04	12.48	10.76
Asparagine	Asn	N	132.12	C ₄ H ₈ N ₂ O ₃	C ₄ H ₆ N ₂ O ₂	114.11	2.02	8.80	—	5.41
Aspartic acid	Asp	D	133.11	C ₂ H ₅ NO ₄	C ₂ H ₅ NO ₃	115.09	1.88	9.60	3.65	2.77
Cysteine	Cys	C	121.16	C ₃ H ₇ NO ₂ S	C ₃ H ₇ NO ₂ S	103.15	1.96	10.28	8.18	5.07
Glutamic acid	Glu	E	147.13	C ₅ H ₉ NO ₄	C ₅ H ₇ NO ₃	129.12	2.19	9.67	4.25	3.22
Glutamine	Gln	Q	146.15	C ₅ H ₁₁ N ₂ O ₃	C ₅ H ₉ N ₂ O ₂	128.13	2.17	9.13	—	5.65
Glycine	Gly	G	75.07	C ₂ H ₅ NO ₂	C ₂ H ₅ NO	57.05	2.34	9.60	—	5.97
Histidine	His	H	155.16	C ₆ H ₁₁ N ₃ O ₂	C ₆ H ₉ N ₃ O	137.14	1.82	9.17	6.00	7.59
Hydroxyproline	Hyp	O	131.13	C ₅ H ₉ NO ₃	C ₅ H ₇ NO ₂	113.11	1.82	9.65	—	—
Isoleucine	Ile	I	131.18	C ₆ H ₁₁ N ₂ O ₂	C ₆ H ₉ NO	113.16	2.36	9.60	—	6.02
Leucine	Leu	L	131.18	C ₆ H ₁₃ N ₂ O ₂	C ₆ H ₁₁ NO	113.16	2.36	9.60	—	5.98
Lysine	Lys	K	146.19	C ₆ H ₁₄ N ₂ O ₂	C ₆ H ₁₂ NO ₂	128.18	2.18	8.95	10.53	9.74
Methionine	Met	M	149.21	C ₅ H ₁₁ NO ₂ S	C ₅ H ₉ NOS	131.20	2.28	9.21	—	5.74
Phenylalanine	Phe	F	165.19	C ₉ H ₁₁ NO ₂	C ₉ H ₉ NO	147.18	1.83	9.13	—	5.48
Proline	Pro	P	115.13	C ₅ H ₉ NO ₂	C ₅ H ₇ NO	97.12	1.99	10.60	—	6.30
Pyroglutamic	Glp	U	139.11	C ₅ H ₇ NO ₃	C ₅ H ₅ NO ₂	121.09	—	—	—	5.68
Serine	Ser	S	105.09	C ₃ H ₇ NO ₃	C ₃ H ₅ NO ₂	87.08	2.21	9.15	—	5.68
Threonine	Thr	T	119.12	C ₄ H ₉ NO ₃	C ₄ H ₇ NO ₂	101.11	2.09	9.10	—	5.60
Tryptophan	Trp	W	204.23	C ₁₁ H ₁₂ N ₂ O ₂	C ₁₁ H ₁₀ N ₂ O	186.22	2.83	9.39	—	5.89
Tyrosine	Tyr	Y	181.19	C ₉ H ₁₁ NO ₃	C ₉ H ₉ NO ₂	163.18	2.20	9.11	10.07	5.66
Valine	Val	V	117.15	C ₅ H ₁₁ NO ₂	C ₅ H ₉ NO	99.13	2.32	9.62	—	5.96

¹ pKa is the negative of the logarithm of the dissociation constant for the -COOH group.

² pKb is the negative of the logarithm of the dissociation constant for the -NH₃⁺ group.

³ pKx is the negative of the logarithm of the dissociation constant for any other group in the molecule.

⁴ pI is the pH at the isoelectric point.

Reference: D. R. Lide, *Handbook of Chemistry and Physics*, 72nd Edition, CRC Press, Boca Raton, FL, 1991.

Amino Acids Reference Charts

Hydrophobic - aliphatic
Hydrophobic - aromatic
Neutral - polar side chains
Acidic

Basic
Unique
Properties Table
Hydrophobicity Index

Technical Articles

Amino acids are the compounds or building blocks that make up peptides and proteins. Each amino acid is structured from an amino group and a carboxyl group bound to a tetrahedral carbon. This carbon is designated as the α-carbon (alpha-carbon).

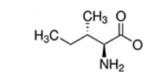
Amino acids differ from each other with respect to their side chains, which are referred to as R groups. The R group for each of the amino acids will differ in structure, electrical charge, and polarity.

Refer to the charts and structures below to explore amino acid properties, types, applications, and availability.

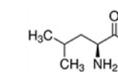
Amino Acids with Hydrophobic Side Chain – Aliphatic



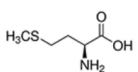
Alanine, Ala, A



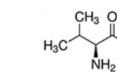
Isoleucine, Ile, I



Leucine, Leu, L

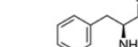


Methionine, Met, M

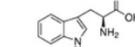


Valine, Val, V

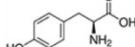
Amino Acids with Hydrophobic Side Chain – Aromatic



Phenylalanine, Phe, F

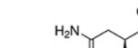


Tryptophan, Trp, W

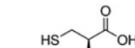


Tyrosine, Tyr, Y

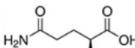
Amino Acids with Polar Neutral Side Chains



Asparagine, Asn, N

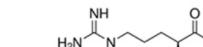


Cysteine, Cys, C

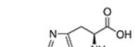


Glutamine, Gln, Q

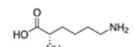
Amino Acids with Electrically Charged Side Chains – Basic



Arginine, Arg, R



Histidine, His, H

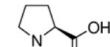


Lysine, Lys, K

Unique Amino Acids



Glycine, Gly, G



Proline, Pro, P

Hydrophobicity Index for Common Amino Acids

The hydrophobicity index is a measure of the relative hydrophobicity, or how soluble an amino acid is in water. In a protein, hydrophobic amino acids are likely to be found in the interior, whereas hydrophilic amino acids are likely to be in contact with the aqueous environment.

The values in the table below are normalized so that the most hydrophobic residue is given a value of 100 relative to glycine, which is considered neutral (0 value). The scales were extrapolated to residues which are more hydrophilic than glycine.

	At pH 2 ^A	At pH 7 ^B
	Very Hydrophobic	Hydrophilic
Leu	100	Phe
Ile	100	Ile
Phe	92	Trp
Trp	84	Leu
Val	79	Val
Met	74	Met
	Hydrophobic	Hydrophilic
Cys	52	Tyr
Tyr	49	Cys
Ala	47	Ala
	Neutral	Neutral
Thr	13	Thr
Glu	8	His
Gly	0	Gly
Ser	-7	Ser
Gln	-18	Gln
Asp	-18	Asp
	Hydrophilic	Hydrophilic
Arg	-26	Arg
Lys	-37	Lys
Asn	-41	Asn
His	-42	Glu
Pro	-46	Pro
	<small>-46 (used pH 2)</small>	<small>-55</small>

^ApH 2 values: Normalized from Sereca et al., *J. Chrom.* 676: 139-153 (1994).

^BpH 7 values: Monera et al., *J. Protein Sci.* 1: 319-329 (1996).

Gap penalties – In general

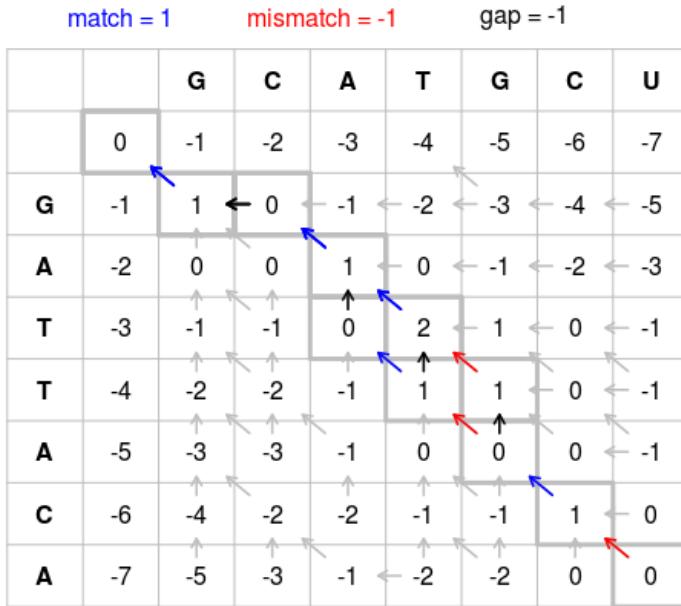
Opening high

Extension medium to high

Codons? Lower (in protein-coding sequences, triplets are less penalized)

Behavior can be modified to reflect biological properties

Needleman-Wunsch



Picture from: Kamil Slowikowski
<https://gist.github.com/slowkow/508393#file-needlemanwunsch-png>

Multiple Sequence Alignment Using a Genetic Algorithm and GLOCSA
 Arenas-Diaz ED, Ochoterena H, Rodriguez-Vazquez K. JAEA. 2009.
<https://doi.org/10.1155/2009/963150>

TABLE 1: Number of possible alignments for two sequences; **m** and **n** are the respective sizes of two given sequences.

<i>m, n</i>	No. of possible alignments
1,1	3
2,2	13
3,3	63
4,4	321
5,5	1683
6,6	8989
7,7	48639
8,8	265729
9,9	1462563
10,10	8097453

Needleman-Wunsch – 1970

Dynamic programming ## Using past knowledge to solve future problems
 ## Breaks down large problems into smaller overlapping subproblems
 ## Alignments of long sequences would be near intractable otherwise

Struggles with very long sequences ## We can't align genomes with NW

Smith–Waterman algorithm (1981) differ in scoring scheme: **local** (SW) vs. **global** (NW)

DNA sequencing wasn't invented yet!

Local *vs.* global alignments

- Most algorithms are fine-tuned for local alignments
- Local optima issues
- Global alignments sometimes make no sense ## motifs and domains

Example
LAGLIDADG
endonuclease

Query	12	LSYLAGFLDGDGCINAQIVRRSDYKLKFQIRVSITFFQKTNRHWFLIWLDKKLDCGTL-R	70
		L Y+AGF DG+G + + VR Y+ +++ + F QK L + + L G L R	
Sbjct	13	LDYIAGFFDGEGSVVVRFVRDGRYRAGYRVSTKVVFVQKERD--VLEEIHETLGMGHLYR	70
Query	71	KRPDGMEYAIIGIASVRNLLSILKPYLKLKK---RQAILLKIIIEKMPHIQNDP	122
		+ DG+ I +R + ++ +K+ + +L+++E H D	
Sbjct	71	RGSDGVWYLEIYRREDLREFVELIGNRTMVKRDALERLATVLELLEGGVHGSRDG	125

Locally? Good alignment of a known biological motif
Globally? Pretty bad alignment

Left -> Right **ne** Right -> Left

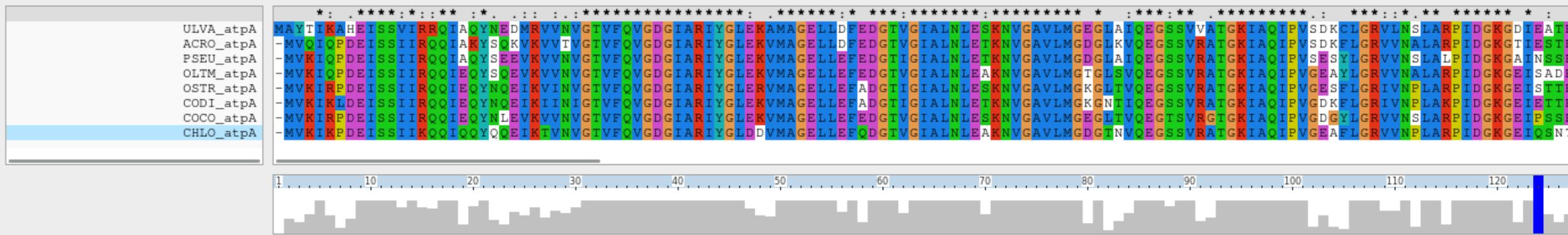
Alignments can differ if started from the **right** rather than the **left**!

An algorithm issue

Not a big problem with highly conserved sequences

Mode: Multiple Alignment Mode

Font: 12



CLUSTAL-Alignment file created [/home/jpombert/BIOL550/05_Sequence_analysis/Ex_01/atpA.aln]

MSAs – Multiple sequences align.

The most common scenario

MSA programs can be used for 2 or more sequences

Speed is key ## heuristics or bust

Common MSA tools

Some of the early ones; ClustalX is very visual and thus a good introduction

ClustalW, X, O <http://www.clustal.org/>

T-COFFEE <http://www.tcoffee.org/>

The recommended ones

MAFFT <http://mafft.cbrc.jp/alignment/software/>

MUSCLE <http://www.drive5.com/muscle/>

A special tool capable of aligning amino acid sequences from pseudogenes

MACSE <http://mbb.univ-montp2.fr/macse>

Exercise 1 – ClustalX

- 1) Connect to Mozart with -X enabled (ssh -X ID@216.47.151.148)
-X is required to display the graphical interface (or -Y with MacOS/Xquartz)
- 2) Type `/opt/clustal/clustalx` (Note: can be slow depending on bandwidth)
- 3) Open the file `atpA.fasta` (File -> Open or ctrl+o)
- 4) Look at the options in the menu
- 5) Do a complete alignment (ctrl+l)
- 6) Noticed the improvements?

Exercise 2 – Alignment formats

Using the same alignment from exercise 1

- 1) File -> Save Sequences As (ctrl+s)
- 2) Save as **CLUSTAL**, **NEXUS** and **PHYLIP** formats
- 3) Look at the differences with less in the shell and compare them with the **FASTA** one
- 4) All three are interleaved formats, **FASTA** is sequential. Note that **PHYLIP** can also be written in a sequential format

Exercise 3 – MAFFT

- 1) Type **mafft -h**, look at the options
- 2) Perform a mafft alignment of **atpB.fasta**
- 3) Look at the mafft output
- 4) Write a Perl script that will run mafft on all fasta files using 2 threads

Exercise 4 – MACSE: codon based alignments

- 1) MACSE is a java program. To launch type:

```
java -jar /opt/MACSE/macse_v2.05.jar
```

New version with GUI

- 2) Look at the CMD line options with -help

Not very helpful right?

- 3) Re-launch with the appropriate parameters for the standard genetic code:

```
-gc_def 1. Use EC_0360.fasta as input.
```

- 4) Look at the outputs. AA => amino acids; NT => nucleotides. When present, putative frameshifts are indicated by !.

MACSE takes longer to compute but returns both amino acid and DNA alignments.

MACSE can align pseudogenes at the amino acid level <- Powerful tool

For a Perl script, see:

https://github.com/PombertLab/Publication_scripts/blob/master/2020_GBE_Hamilton_Oordospore_nonadaptive_processes/dNdS/run_macse.pl

Homology search

Sequence homology

- | | |
|--------------------|--|
| Primary sequences | - BLAST (alignment-based); MASH (alignment-free) |
| Domains and motifs | - Pfam, CDD |

Structural homology

- 3D structures

Sequence homology – primary sequences

BLAST

Basic Local Alignment Search Tool

Fast queries against databases (online/local)

<http://blast.ncbi.nlm.nih.gov/Blast.cgi> (searches against GenBank)

DIAMOND

Recent (2015); 500x to 20,000x faster than BLAST

Basic Local Alignment Search Tool

Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ

J. Mol. Biol. 1990. 215:403–410

Fast and sensitive protein alignment using DIAMOND

Buchfink B, Xie C, Huson DH

Nature Methods. 2015. 12:59–60

Gapped BLAST and PSI-BLAST: a new generation of protein database search programs

Altschul SF, et al.

Nucl. Acids Res. 1997. 25(17):3389–3402.

BLAST: improvements for better sequence analysis

Ye J, McGinnis S, Madden TL

Nucleic Acids Res. 2006. 34:W6–W9

Basic Local Alignment Search Tool

Stephen F. Altschul¹, Warren Gish¹, Webb Miller²
Eugene W. Myers³ and David J. Lipman¹

¹National Center for Biotechnology Information
National Library of Medicine, National Institutes of Health
Bethesda, MD 20894, U.S.A.

²Department of Computer Science
The Pennsylvania State University, University Park, PA 16802, U.S.A.

³Department of Computer Science
University of Arizona, Tucson, AZ 85721, U.S.A.

(Received 26 February 1990; accepted 15 May 1990)

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST), directly approximates alignments that optimize a measure of local similarity, the maximal segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP scores allow an analysis of the performance of this method as well as the statistical significance of alignments it generates. The basic algorithm is simple and robust; it can be implemented in a number of ways and applied in a variety of contexts including straightforward DNA and protein sequence database searches, motif searches, gene identification searches, and in the analysis of multiple regions of similarity in long DNA sequences. In addition to its flexibility and tractability to mathematical analysis, BLAST is an order of magnitude faster than existing sequence comparison tools of comparable sensitivity.

1. Introduction

The discovery of sequence homology to a known protein or family of proteins often provides the first clues about the function of a newly sequenced gene. As the DNA and amino acid sequence databases continue to grow in size they become increasingly useful in the analysis of newly sequenced genes and proteins because of the greater chance of finding such homologies. There are a number of software tools for searching sequence databases but all use some measure of similarity between sequences to distinguish biologically significant relationships from chance similarities. Perhaps the best studied measures are those used in conjunction with variations of the dynamic programming algorithm (Needleman & Wunsch, 1970; Sellers, 1974; Sankoff & Kruskal, 1983; Waterman, 1984). These methods assign scores to insertions, deletions and replacements, and compute an alignment of two sequences that corresponds to the least costly set of such mutations. Such an alignment may be thought of as minimizing the evolutionary distance or maximizing the similarity between the two sequences compared. In either case, the cost of this alignment is a measure of similarity; the algorithm guarantees it is

optimal, based on the given scores. Because of their computational requirements, dynamic programming algorithms are impractical for searching large databases without the use of a supercomputer (Gotoh & Tagashira, 1986) or other special purpose hardware (Coulson *et al.*, 1987).

Rapid heuristic algorithms that attempt to approximate the above methods have been developed (Waterman, 1984), allowing large databases to be searched on commonly available computers. In many heuristic methods the measure of similarity is not explicitly defined as a minimal cost set of mutations, but instead is implicit in the algorithm itself. For example, the FASTP program (Lipman & Pearson, 1985; Pearson & Lipman, 1988) first finds locally similar regions between two sequences based on identities but not gaps, and then rescores these regions using a measure of similarity between residues, such as a PAM matrix (Dayhoff *et al.*, 1978) which allows conservative replacements as well as identities to increment the similarity score. Despite their rather indirect approximation of minimal evolution measures, heuristic tools such as FASTP have been quite popular and have identified many distant but biologically significant relationships.

Sequence homology – motifs and domains

Not every portion of a protein has functional significance

Functional domains, or **motifs**, are conserved throughout evolution

Hidden Markov Models:

- Search for motifs
- Search for genes

nhmmmer: DNA Homology Search With Profile HMMs
Wheeler TJ and Eddy SR
Bioinformatics. 2013. 29: 2487–2489

Sequence analysis

Advance Access publication July 9, 2013

nhmmmer: DNA homology search with profile HMMs

Travis J. Wheeler* and Sean R. Eddy

HHMI Janelia Farm Research Campus, Ashburn, VA 20147, USA

Associate Editor: Alfonso Valencia

ABSTRACT

Summary: Sequence database searches are an essential part of molecular biology, providing information about the function and evolutionary history of proteins, RNA molecules and DNA sequence elements. We present a tool for DNA/DNA sequence comparison that is built on the HMMER framework, which applies probabilistic inference methods based on hidden Markov models to the problem of homology search. This tool, called nhmmmer, enables improved detection of remote DNA homologs, and has been used in combination with Dfam and RepeatMasker to improve annotation of transposable elements in the human genome.

Availability: nhmmmer is a part of the new HMMER3.1 release. Source code and documentation can be downloaded from <http://hmmer.org>. HMMER3.1 is freely licensed under the GNU GPLv3 and should be portable to any POSIX-compliant operating system, including Linux and Mac OS/X.

Contact: wheeler@janelia.hhmi.org

Received and revised on June 17, 2013; accepted on July 5, 2013

1 INTRODUCTION

A widely used general purpose tool for DNA/DNA sequence comparison is blastn (Altschul *et al.*, 1990; Camacho *et al.*, 2009), which heuristically approximates the Smith–Waterman algorithm (Smith and Waterman, 1981) for recognizing local regions of similarity between two sequences. In recent years, most advances in DNA/DNA comparison have related to accelerating search for near-exact matches (Kent, 2002; Langmead *et al.*, 2009; Li and Durbin, 2009), and to improving whole-genome alignment (Kurtz *et al.*, 2004; Schwartz *et al.*, 2003). Another area that deserves attention is the development of methods that maximize the power of computational sequence comparison tools to detect remote homologs.

Profile hidden Markov models (profile HMMs) (Durbin *et al.*, 1998; Krogh *et al.*, 1994) represent an important advance in terms of sensitivity of sequence searches for remote homology. They provide a formal probabilistic framework for sequence comparison and improve detection of remote homologs by (i) enabling position-specific residue and gap scoring based on a query profile, and (ii) calculating the signal of homology based on the more powerful ‘Forward/Backward’ HMM algorithm that computes not just one best-scoring alignment, but a sum of support over all possible alignments. In the past, this improved sensitivity came at a significant computational cost, but recent advances in HMMER3 have increased speed for

protein search by ~100-fold, reaching blastp-like speed through a combination of filtering heuristics (Eddy, 2008) and computer engineering (Eddy, 2011; Farrar, 2007). Tools based on profile HMMs (Eddy, 2009; Karplus *et al.*, 1998) have historically focused on protein search, with little concentration on the challenges presented by (i) chromosome-length target sequences, and (ii) the extreme composition bias often seen in genomic DNA. With attention to the details of DNA search, nhmmmer builds upon the speed advances of HMMER3, bringing the power of profile HMMs to DNA homology search, at speeds nearly as fast as blastn with sensitive settings.

An example of a biological problem requiring sensitive detection of remote DNA homologs is the annotation of genomic sequence derived from ancient transposable element (TE) expansions. A prerelease version of the nhmmmer tools has recently been shown to provide increased sensitivity over blastn and other single-sequence search methods, with reduced false discovery rate and reasonable runtime, in searching for TEs (Wheeler *et al.*, 2013). For example, when nhmmmer was used within the recently released RepeatMasker 4.0 (Smit and Hubley, 2013), an additional 150 Mb (5%) of the human genome was reliably annotated as derived from TEs.

2 USAGE AND PERFORMANCE

Usage. The program nhmmmer is used to search one or more nucleotide queries against a nucleotide sequence database. For each query, nhmmmer searches the target database and outputs a ranked list of the hits with the most significant matches to the query. A query may consist of a single sequence, a multiple sequence alignment, or a profile HMM built using the HMMER program hmmbuild. Each hit represents a region of local similarity between a portion of the query and a subsequence of the full target database sequence, and is assigned a similarity score S in bits, along with an E -value (Eddy, 2008) indicating the expected number of false positives at a threshold of score S . Each hit is also accompanied by an alignment of the matched sequence to the model, with values indicating the confidence with which each position is aligned.

The final score, boundaries and alignment of a hit are computed based on filling in a Forward/Backward dynamic programming matrix, but the computational burden of doing this for the full target database is prohibitive. Therefore, nhmmmer uses a series of acceleration filters that depend on simpler approximations of the final Forward score of a hit. These filters are based on those used in the HMMER3 protein search tools (Eddy, 2011), but have been modified to work in the context of long (potentially chromosome length) target sequences. The initial filter, called ‘single segment ungapped Viterbi’, scans along

*To whom correspondence should be addressed.

Structural homology – 3D structures

Sometimes the sequence itself is irrelevant

Structure is often key

If it looks like a wheel, it probably is...

The structural alignment between two proteins: is there a unique answer?

Adam Godzik, 1996

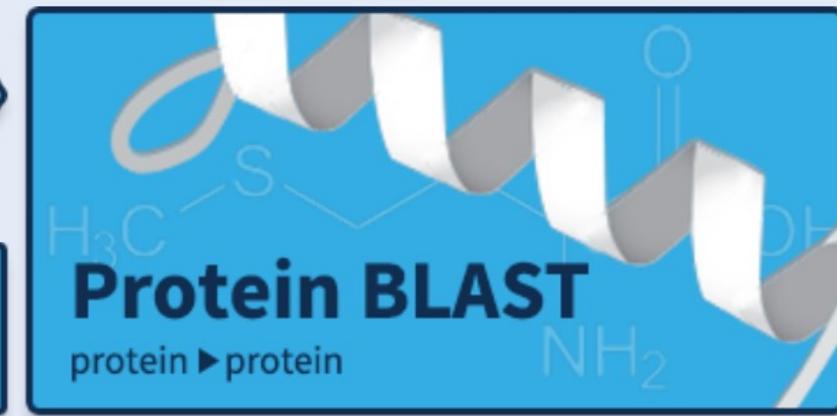
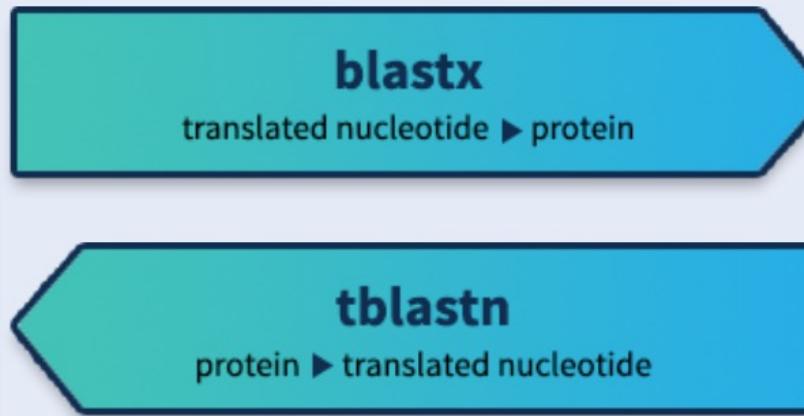
Abstract Finding common molecular substructures in complex 3D protein structures is still challenging. This is especially visible when scanning entire databases containing tens or even hundreds of thousands protein structures. Graphics processing units (GPUs) and general purpose graphics processing units (GPGPUs) promise to give a high speedup of many time-consuming and computationally demanding processes over their original implementations on CPUs. In this chapter, we will see that a massive parallelization of the 3D structure similarity searching on many core CUDA-enabled GPU devices leads to reduction of the execution time of the process and allows to perform it in real time.

Keywords Proteins · 3D protein structure · Tertiary structure · Similarity searching · Structure matching · Structure comparison · Structure alignment · Parallel computing · GPU · CUDA

3.1 Introduction

Protein 3D structure similarity searching is a process in which a given protein structure is compared to another protein structure or a set of protein structures collected in a database. The aim of the process is to find matching fragments of compared protein structures. On the basis of the similarities found during this process, scientists can draw useful conclusions about the common ancestry of the proteins, and thus the organisms (that the proteins came from), their evolutionary relationships, functional similarities, existence of common functional regions, and many other things [6]. This process is especially important in situations, where sequence similarity searches fail or deliver too few clues [15]. There are also other processes in which protein structure

Web BLAST



NCBI BLAST

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

Basic Local Alignment Search Tool

Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ
J. Mol. Biol. 1990. 215:403–410

BLAST: improvements for better sequence analysis

Ye J, McGinnis S, Madden TL
Nucleic Acids Res. 2006. 34:W6–W9

Gapped BLAST and PSI-BLAST: a new generation of protein database search programs

Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, & Lipman DJ
Nucl. Acids Res. 1997. 25(17):3389–3402.

Types of BLAST

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

<u>Type</u>	<u>Query</u>	<u>Database</u>
blastn	nucleotide	nucleotide
blastp	protein	protein
blastx	translated nucleotide	protein
tblastn	protein	translated nucleotide
tblastx	translated nucleotide	translated nucleotide

Note: To perform protein <-> nucleotide searches, the nucleotide sequences are translated to amino acids, then searches are done at the amino acid level

Sequences producing significant alignments

Download ▾

Manage Columns ▾

Show

100 ▾



select all 94 sequences selected

GenPept

Graphics

Distance tree of results

Multiple alignment

	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input checked="" type="checkbox"/>	DNA-directed RNA polymerase [Mitosporidium daphniae]	1373	1373	99%	0.0	62.68%	XP_013238900.1
<input checked="" type="checkbox"/>	DNA-directed RNA polymerase II second largest subunit [Encephalitozoon cuniculi]	1281	1281	99%	0.0	55.63%	AGE96471.1
<input checked="" type="checkbox"/>	DNA-DIRECTED RNA POLYMERASE II SECOND LARGEST SUBUNIT [Encephalitozoon cuniculi]	1280	1280	99%	0.0	55.63%	NP_586140.1
<input checked="" type="checkbox"/>	DNA-directed RNA polymerase subunit B [Encephalitozoon hellem ATCC 50504]	1272	1272	99%	0.0	55.54%	XP_003888100.1
<input checked="" type="checkbox"/>	subunit beta of DNA-directed RNA polymerase [Hamiltosporidium tvaermannensis]	488	665	98%	9e-154	32.19%	TBU01157.1
<input checked="" type="checkbox"/>	subunit beta of DNA-directed RNA polymerase [Hamiltosporidium tvaermannensis]	488	665	98%	1e-153	32.19%	TBU10981.1
<input checked="" type="checkbox"/>	DNA-directed RNA polymerase II subunit RPB2 [Vittaforma cornea ATCC 50505]	469	469	44%	1e-153	45.67%	XP_007605419.1
<input checked="" type="checkbox"/>	subunit beta of DNA-directed RNA polymerase [Hamiltosporidium magnivora]	488	665	98%	2e-153	32.19%	TBU08960.1

bit scores E-values

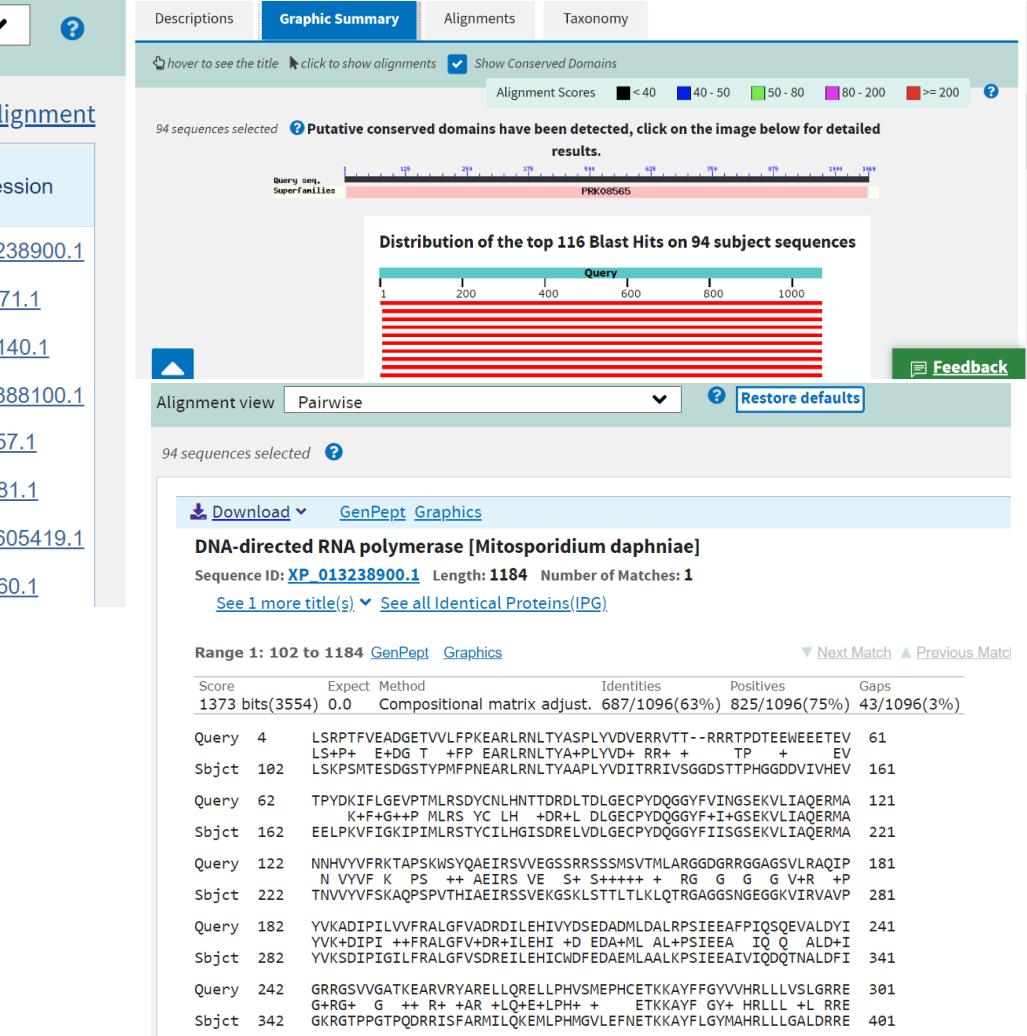
BLAST – E-values

BLAST analyses will return many hits

Not all of them are good (quite a few are garbage actually)

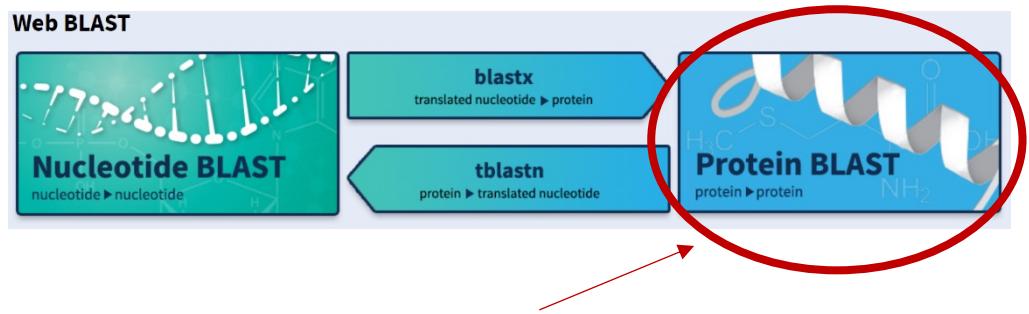
E-values are probability estimates

bit scores are the sum of all the aligned bits



Lower is better

Higher is better



Exercise 5 – Online BLAST (BLASTP)

- 1) Goto <http://blast.ncbi.nlm.nih.gov/Blast.cgi>
- 2) Open **BLAST_1.fasta**, perform a **BLASTP** on it
- 3) Find anything good?
- 4) Open **BLAST_2.fasta**, perform a **BLASTP** on it
- 5) What about this time? Try with **BLAST_3.fasta**
- 6) Explore the various options by yourself, feel free to ask questions

Enter accession number(s), gi(s), or FASTA sequence(s)

[Clear](#)

Query subrange

From
To

Or, upload file

 Choose File No file chosen

Job Title

 Enter a descriptive title for your BLAST search Align two or more sequences

Choose Search Set

Database

Organism
OptionalExclude
Optional

Non-redundant protein sequences (nr)

Microsporidia (taxid:6029)

 exclude

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

 Models (AM/XP) Non-redundant RefSeq proteins (WP) Uncultured/environmental sample sequences

Assembly txid6029[ORGN]

Create alert Advanced Browse by organism

Organism group Summary 20 per page Sort by Significance

Status clear

Latest (61)
Latest GenBank (61)
Latest RefSeq (9)

Assembly level

Complete genome (0)
Chromosome (5)
Scaffold (21)
Contig (35)

RefSeq category

Reference (0)
Representative (31)

Exclude clear

Exclude partial (0)
Exclude derived from

Search results
Items: 1 to 20 of 61

Filters activated: Latest, Exclude derived from surveillance project, Exclude anomalous. [Clear all](#) to show 67 items.

Quoted phrase not found.

ASM9122v2

1. Organism: Encephalitozoon cuniculi GB-M1 (microsporidians)
Submitter: Genoscope
Date: 2017/03/09
Assembly level: Chromosome
Genome representation: full

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC

Search for as complete name lock Go Clear

Display 3 levels using filter: none

Microsporidia

L taxonomy ID: 6029 (for references in articles please use NCBI:txid6029)
current name

Microsporidia

NCBI BLAST name: microsporidians
Rank: phylum
Genetic code: [Translation table 1 \(Standard\)](#)
Mitochondrial genetic code: [Translation table 1 \(Standard\)](#)
Other names:

<https://www.ncbi.nlm.nih.gov/taxonomy>

txidXXXX[ORGN]

We can restrict our searches with a taxonomic ID

The txid is a unique ID given by NCBI to an organism or group

It can be retrieved from the scrolling menu on the BLAST query page

The txid can also be found from the NCBI taxonomy database

BLAST – Command line

No web latency

Can run batch searches

Can be automated

Runs privately (you can make your own databases) ## Useful for data privacy

Makeblastdb – Private databases

```
makeblastdb \
  -in file.fasta \
  -dbtype nucl \
  -out nickname \
  -title 'short description'
```

FASTA input file (single/multifasta; DNA, RNA or protein)
Database type: nucl (DNA/RNA) or prot (proteins)
Name of the database to be created
Description (optional)

```
(t)blast(n,p,x) \
-num_threads 16 \
-query input.fasta \
-db database_name \
-evalue 1e-40 \
-culling_limit 10 \
-outfmt 6 \
-out test.blastp.6
```

BLAST program to use
Number of threads to use
FASTA file to query
Database to query against
Desired minimum *E*-value
Desired number of top hits
Desired output format; 0 = pairwise alignments, 6 = tab-delimited
Desired output name

Running local BLAST searches

Works against NCBI and/or custom databases

Great when running several queries

No CPU usage limit, unlike the web-based tool

Exercise 6 – Creating your own BLAST db

- 1) Concatenate all fasta files from the Alignments folder into a single file (`my.fasta`)
- 2) Create a folder `BLAST` and a subfolder `db`
- 3) Move (or copy) the `my.fasta` file into the `BLAST/db` folder
- 4) In `db`, create your own database called `my`, for title use `PlastidProteins`
- 5) Keep it, we will use it for the next exercise

Exercise 7 – Running BLAST on my db

- 1) Use the **my** database from exercise 7
- 2) Copy **ARABthali.fasta** into folder BLAST
- 3) From the BLAST folder, run a **BLASTP** with **ARABthali.fasta** against **my** with:
2 threads, evalue of **1e-10**, and outfmt **6** (tabular), save it as **ARABthali.blastp**
- 4) Browse the output with less
- 5) Let's parse it with Perl

Exercise 8 – A simple BLAST parsing

While simple, the script may be beyond your current abilities. If so, read the ## answer and try to understand what every little piece of code does.

Using the `ARABthali.blastp` output from exercise 8, create a script that:

- 1) Takes the best hit for each protein query ## the top hit is the best for each query
Hints: `%protein = ()`, `if (exists $protein{$query}){ do something; }`
- 2) Assigns `$query` (e.g. `NP_085480.1`) and `$gene` values (e.g. `CHLO_atpA`)
- 3) Prints only the best hit to the output file e.g.:

`NP_085480.1 rpl16`

`NP_085492.1 petB`

`NP_085496.1 rbcL`

```
taxdb.tar.gz          ## Taxonomy information for the databases
nr.*tar.gz           ## Non-redundant protein sequences
nt.*tar.gz           ## Partially non-redundant nucleotide sequences
16SMicrobial.tar.gz ## Bacterial + Archaeal 16S rRNA sequences (BioProjects 33175 and 33117)
Representative_Genomes.*tar.gz ## Representative bacterial/archaeal genomes
human_genomic.*tar.gz    ## Human RefSeq (NC_#####) chromosome record
other_genomic.*tar.gz   ## RefSeq chromosome records (NC_#####) for non-human organisms
cdd_delta.tar.gz       ## Conserved Domain Database sequences for use with standalone deltablast
env_nr.*tar.gz         ## Protein sequences for metagenomes
env_nt.*tar.gz         ## Nucleotide sequences for metagenomes
refseq_genomic.*tar.gz ## NCBI genomic reference sequences
refseq_protein.*tar.gz ## NCBI protein reference sequences
refseq_rna.*tar.gz     ## NCBI Transcript reference sequences
```

Downloading NCBI databases

<ftp://ftp.ncbi.nlm.nih.gov/blast/db/>

The files are compressed to save space

Must be decompressed (tar -zxvf *.tar.gz)

See NCBI_scripts for automatic downloading

Large files, NR > 28Gb and growing

```
blastdb_aliastool \
-dbtype nucl \
-gilist file.gi \
-db nr \
-out nickname \
-title 'short description' \
## Database type (nucl or prot)
## List of GIs (sequence identifiers)
## Input database
## Name of the desired subset
## Description (optional)
```

Creating local db subsets

Creates smaller subsets of NCBI databases for to reduce computation time

Works with NR and NT

Get taxonomic ID (e.g. txid2759[ORGN]) ## NCBI taxonomy or BLAST (Organism field)

Get GI list from Entrez: <http://www.ncbi.nlm.nih.gov/protein/> - NR, [/nucleotide/](http://www.ncbi.nlm.nih.gov/nucleotide/) - NT

How to get GI list from NCBI?

1. Go to <http://www.ncbi.nlm.nih.gov/protein/>
or <https://www.ncbi.nlm.nih.gov/nucleotide>

The screenshot shows the NCBI Protein search results page. At the top, there is a search bar with the query "txid2759[ORGN]". To the right of the search bar is a "Search" button. A red circle highlights both the search bar and the "Search" button, with the text "2. Enter desired taxonomic ID, then click search" above them.

A large red arrow points from the text "3. Select GI list from the Send to menu" down to a tooltip that appears over a "Send to" dropdown menu. The tooltip lists various options: Summary, GenPept, GenPept (full), FASTA, ASN.1, XML, INSDSeq XML, TinySeq XML, Feature Table, FASTA CDS, Accession List, GI List, GFF3, and Summary. The "GI List" option is highlighted with a red circle.

The main content area displays search results for items 1 to 20 of 87471920. One result is shown: "protachykinin [Syngnathus acus]" with an accession number XP_037131029.1 and GI number 1929476210. Navigation links for "GenPept", "Identical Proteins", "FASTA", and "Graphics" are at the bottom of the result card.

On the left sidebar, there are categories for "Species" (e.g., Animals, Plants, Fungi, Protists, Bacteria, Viruses) and "Source databases" (e.g., PDB, RefSeq, UniProtKB / Swiss-Prot). On the right sidebar, there are filters and a "Tree" view.

Alternatives to db subsets

Since **BLAST 2.4.0+**, you can run blast searches with the following options:

- | | |
|---|--|
| -gilist Filename | # Restrict search to GI's listed in this file. Local searches only |
| -negative_gilist Filename | # Opposite. Restrict search to everything but the GI's listed |
| -seqidlist/-negative_seqidlist Filename | # Restrict search to list of accession numbers |
| -query X.fasta -subject Y.fasta | # Directly search query against subject |

Kingdom

Phylum

Class

Family

Genus

Species

Subspecies

Taxonomy: Classification of living organisms according to their evolutionary relationships

Taxonomized BLASTs

Requires local NCBI Taxonomy database (<ftp://ftp.ncbi.nlm.nih.gov/blast/db/taxdb.tar.gz>)

Databases must be in the path (`export BLASTDB=/path/to/NCBI/TaxDB:/path/to/NCBI/NR:/path/to/NCBI/NT`)

Uses a modified outfmt

`## e.g. -outfmt '6 qseqid sseqid qstart bitscore eval staxids sscinames sskingdoms sblastnames'`

`## Fields can be defined; http://www.metagenomics.wiki/tools/blast/blastn-output-format-6`

Exercise 9 – Taxonomized BLAST

- 1) Add Mozart's local NCBI NR, NT and TaxDB databases located in /media/Data_1/NCBI/ to your .bash_profile:

```
BLASTDB=$BLASTDB:/media/Data_1/NCBI/TaxDB  
BLASTDB=$BLASTDB:/media/Data_1/NCBI/NR  
BLASTDB=$BLASTDB:/media/Data_1/NCBI/NT  
export BLASTDB
```

- 2) Run a blastp analysis using TaxBlast.fasta as query against the NCBI nr database using 4 threads the following outfmt format:
-outfmt '6 qseqid sseqid bitscore eval staxids sskingdoms sscinames'



The DIAMOND protein aligner - <http://www.diamondsearch.org>

Diamond – Faster BLASTP/BLASTX

<https://github.com/bbuchfink/diamond>

Up to 20,000X times faster than BLASTP/BLASTX

Works from FASTA or FASTQ files ## Fast enough to work with illumina data!

Easy to use

Diamond – Running searches

```
diamond makedb --in input.fasta -d custom
```

Creating a custom database

```
diamond blastx -d custom -q reads.fna -o matches.m8
```

Running blastx equivalent

```
diamond blastp -d custom -q reads.fq -o matches.m8
```

Running blastp equivalent

```
### Downloading NCBI taxonomy accession numbers to taxon ids file
nice -n +15 wget ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/accession2taxid/prot.accession2taxid.gz
### Downloading NCBI taxonomy dmp files (nodes.dmp + names.dmp)
nice -n +15 wget ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/taxdmp.zip
unzip -o taxdmp.zip

### Downloading NCBI nr FASTA file => huge file!!!
nice -n +15 wget ftp://ftp.ncbi.nlm.nih.gov/blast/db/FASTA/nr.gz

### Creating diamond database from NCBI NR fasta files; this will take time!
diamond makedb \
--threads 64 --in nr.gz \
--db diamond_nr \
--taxonnames names.dmp \
--taxonmap prot.accession2taxid.gz \
--taxonnodes nodes.dmp
```

Diamond – Setting up NCBI NR database

We can use the **NCBI NR database** with Diamond
Searches can be performed with taxonomy
Requires the NCBI NR fasta file + taxonomy files

Fast and sensitive protein alignment using DIAMOND
Buchfink B et al.
Nat Methods. 2015. 12(1):59-60
doi: 10.1038/nmeth.3176

```
### Running a taxonomized search with Diamond
diamond blastx \
-d /media/Data_1/NCBI/Diamond/diamond_nr \
-q test.fasta \
-o test.blastx.diamond \
-f 6 qseqid sseqid pident length mismatch gapopen qstart qend sstart send evalue bitscore \
    staxids skingdoms sphylums sscinames sblastnames saltitles
                                                ## Blastx or blastp
                                                ## Database
                                                ## Query
                                                ## Output
                                                ## Desired fields
                                                ## Desired fields
```

Diamond – Taxonomized searches

We can use the **NCBI NR database** with Diamond

Searches can be performed with taxonomy

Requires the NCBI NR fasta file + taxonomy files

Annotations from homology searches

Don't trust them without cross-validation! Why?

- **Online databases can contain many errors** (contributors are not always experts)
- The homology may be too weak to be meaningful
- The homology may be with a different domain that differs from the annotated function(s)

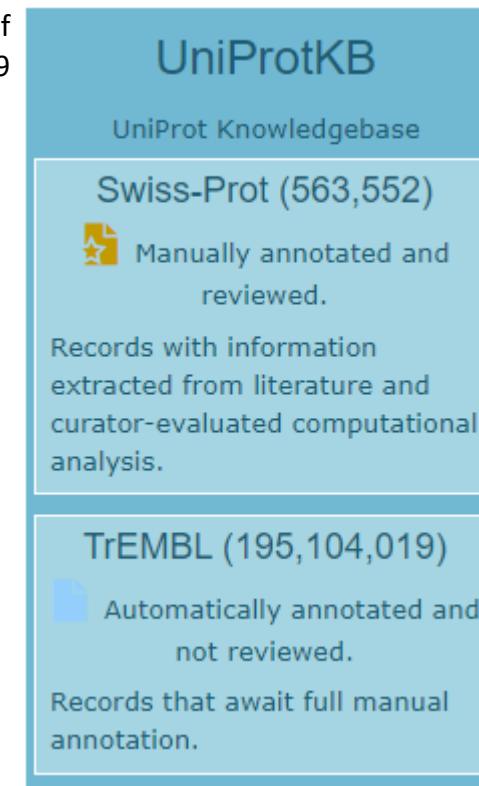
Bottom line? Use critical thinking (aka common sense)

Inferring homology using different tools is a good idea

- We must corroborate our predictions
- Using **different databases** and/or **tools** may yield different results
- Congruency across independent analyses increases confidence in our predictions

We can run BLAST searches against UniProtKB

Metrics as of
2020-10-19



UniProtKB

<http://www.uniprot.org/>

Swiss-Prot – Manually annotated and curated

TrEMBL – Automatically annotated, not curated

UniProt: a hub for protein information

UniProt Consortium

Nucleic Acids Res. 2015. 43:D204-D212 Database issue

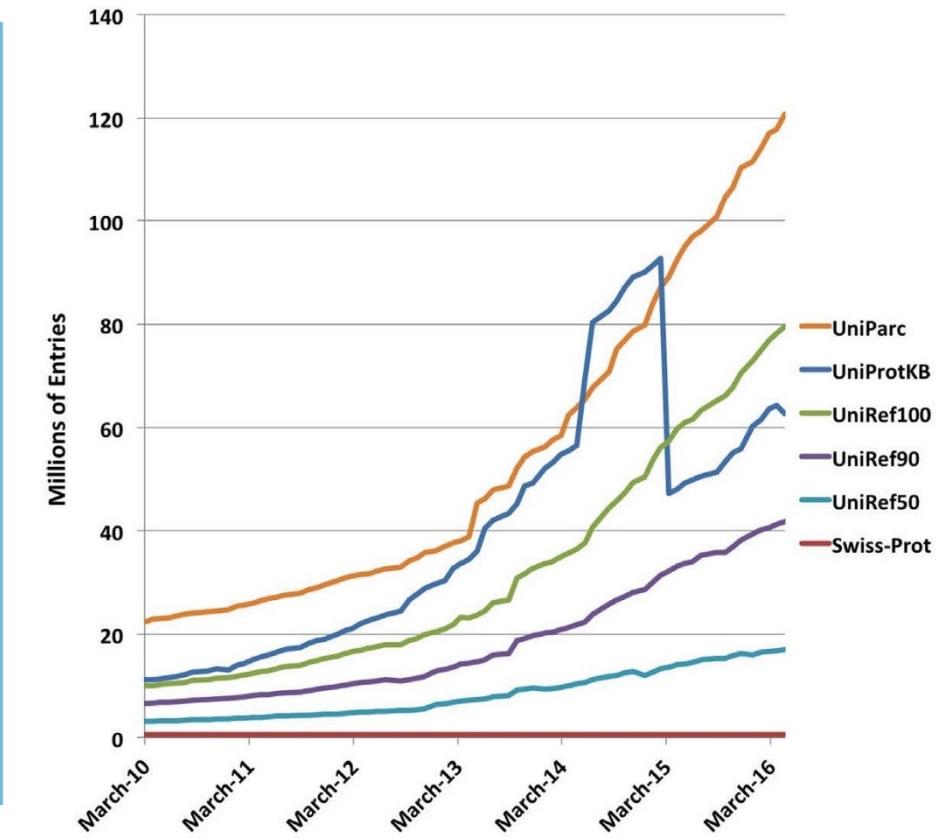


Figure 1. Growth of the number of sequences in UniProt databases. The blue line shows the growth in UniProtKB/TrEMBL entries from January 2010 to date. The sharp drop in UniProtKB entries corresponds to the proteome redundancy minimization (PRM) procedure implemented in March 2015. Note that the post-PRM growth in UniProtKB is no longer exponential.

UniProt: the universal protein knowledgebase

UniProt Consortium

Nucleic Acids Res. 2017. 45:D158-D169 Database issue

Downloads release 2020_05

<https://www.uniprot.org/downloads>

Basket ▾

UniProt is updated every eight weeks (see FAQ on [how to be notified automatically of updates](#)). You can download small data sets and subsets directly from this website by following the download link on any search result page.

We can download and search UniProtKB

For a quick how-to, see:

<https://github.com/PombertLab/A2GB#Performing-homology-searches-against-UniProt-databases>

Downloading the SwissProt and TrEMBL databases

We can use [get_UniProt.pl](#) to download the SwissProt and/or TrEMBL databases from UniProt:

```
get_UniProt.pl -s -t -f $ANNOT/UNIPROT/ -n 20 -l download.log
```

Options for [get_UniProt.pl](#) are:

-s (--swiss)	Download Swiss-Prot
-t (--trembl)	Download TrEMBL
-f (--folder)	Download folder [Default: ./]
-n (--nice)	Linux Process Priority [Default: 20] ## Runs downloads in the background
-l (--log)	Print download information to log file
-d (--decompress)	Decompresss downloaded files with gunzip ## trEMBL files will be huge, off by default

Creating tab-delimited product lists from UniProt databases

Homology searches against the UniProt databases will return positive matches against the corresponding accession numbers. However, these matches will not include product names. To facilitate downstream analyses, we can create tab-delimited lists of accession numbers and their products with [get_uniprot_products.pl](#):

```
get_uniprot_products.pl $ANNOT/UNIPROT/uniprot_*.fasta.gz
```

Running DIAMOND or BLAST searches against UniProt databases

We can use [DIAMOND](#) to perform homology searches against the [UniProt](#) databases. Documentation on how to use DIAMOND can be found [here](#).

First, let's create DIAMOND-formatted databases:

```
mkdir $ANNOT/DIAMOND/; mkdir $ANNOT/DIAMOND/DB/;

diamond makedb \
  --in $ANNOT/UNIPROT/uniprot_sprot.fasta.gz \
  -d $ANNOT/DIAMOND/DB/sprot

diamond makedb \
  --in $ANNOT/UNIPROT/uniprot_trembl.fasta.gz \
  -d $ANNOT/DIAMOND/DB/trembl
```

Second, let's perform protein-protein homology searches against the UniProt databases with a tabular output format (same as [NCBI BLAST+](#)'s -outfmt 6 format). Note that these searches will likely take a while depending the total number of proteins queried and the size of these databases:

```
diamond blastp \
  -d $ANNOT/DIAMOND/DB/sprot \
  -q $ANNOT/proteins.fasta \
  -o $ANNOT/DIAMOND/diamond.sprot.6 \
  -f 6

diamond blastp \
  -d $ANNOT/DIAMOND/DB/trembl \
  -q $ANNOT/proteins.fasta \
  -o $ANNOT/DIAMOND/diamond.trembl.6 \
  -f 6
```

We can also search custom reference datasets

Useful if you have access to:

- Well-curated references
- Your own datasets
- Datasets that must be kept private

With **DIAMOND**, we must pipe the data from the **STDIN**

Performing homology searches against reference datasets

If desired, reference datasets (custom or downloaded from NCBI) can also be used as databases in homology searches to help with annotations. NCBI datasets can be accessed from the [NCBI genome database](#) or directly from their new [dataset repository](#).

For example, using two datasets downloaded from NCBI:

```
## Downloading data from NCBI
mkdir $ANNOT/REFERENCES/;
wget -O $ANNOT/REFERENCES/ref1.faa.gz https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/007/859/695/GCA_007859695.1
wget -O $ANNOT/REFERENCES/ref2.faa.gz https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/214/015/GCF_000214015.3

## Creating DIAMOND database; using zcat to concatenate the output to STDOUT, then feed it to diamond as input
zcat $ANNOT/REFERENCES/*.gz | diamond makedb -d $ANNOT/DIAMOND/DB/reference

## Running DIAMOND
diamond blastp \
-d $ANNOT/DIAMOND/DB/reference \
-q $ANNOT/proteins.fasta \
-o $ANNOT/DIAMOND/diamond.reference.6 \
-f 6
```

To create a tab-delimited list of accession numbers and their associated proteins from the downloaded NCBI .faa.gz files, we can use [get_reference_products.pl](#).

```
get_reference_products.pl -f $ANNOT/REFERENCES/*.gz -l $ANNOT/REFERENCES/reference.list
```

For a quick how-to, see:

<https://github.com/PombertLab/A2GB#Performing-homology-searches-against-reference-datasets>

Domains and motifs-based searches

Most tools use Hidden Markov Model (HMM) approaches

Major tools are:

- Pfam
- CDD
- InterProScan 5

Pfam 33.1 (May 2020, 18259 entries)

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



The Pfam website will be decommissioned in January 2023.

You will be redirected in 3 seconds to the corresponding data page on the InterPro website.

Now part of
InterPro/InterProScan



<http://pfam.xfam.org/>

A database of curated protein families

Probabilistic, based on Hidden Markov model (HMM)

The Pfam protein families database

Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al.

Nucleic Acids Res. 2014. 42:D222-D230 Database Issue

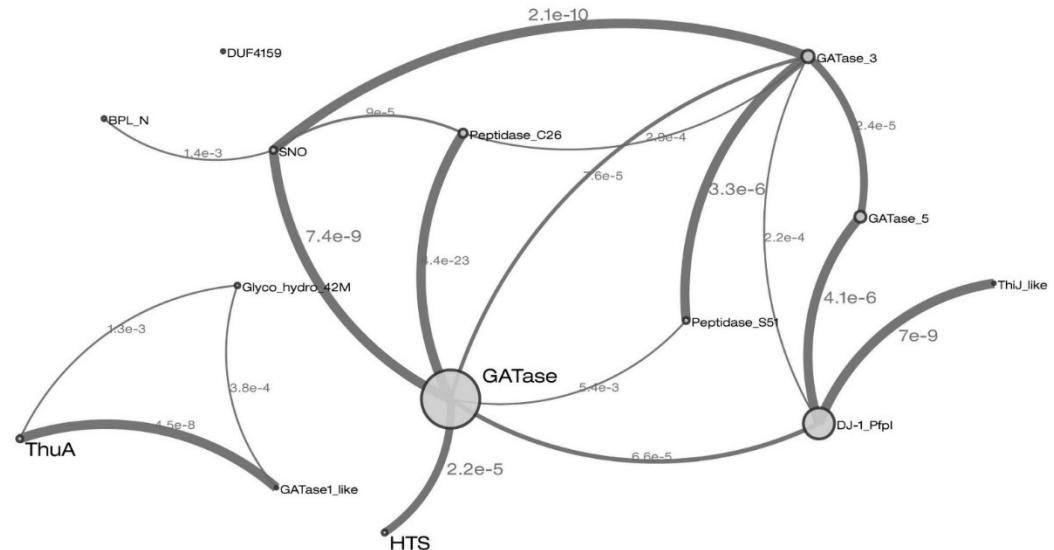


Figure 1. Example of the improved representation of relationships graph, indicating the similarity between the Pfam entries within a clan. This particular entry shows the relationship between the entries in the Glutaminase I clan (accession:CL0014). Each entry in the clan is a node in the graph and is represented as circle, with the diameter of the circle being proportional to the number of sequences in the *full* alignment. Nodes are connected (edges) based on the HHsearch results between the clan members, with the width of edges proportional to the E-value of the HHsearch similarity (E-values ≤ 0.01 are deemed significant). The clanviewer component has been included in the BioJS registry (<https://bijs.org/d/clanviewer>) and its code is freely available in github (<https://github.com/ProteinsWebTeam/clanviewer>). In this particular clan, there are three entries (ThuA (PF06283), GATase1-like (PF07090) and Glyco_hydro_42M (PF08532)) that from a disconnected sub-cluster. DUF4159 (PF13709) is also unconnected to any other entry. However, these entries are included as part of this clan based on the structural similarities to other entries in the clan.

Clans = groups of related families

The Pfam protein families database: towards a more sustainable future

Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al.

Nucleic Acids Res. 2016. 44:D279-D285 Database Issue

Rfam 14.3 (September 2020, 3446 families)

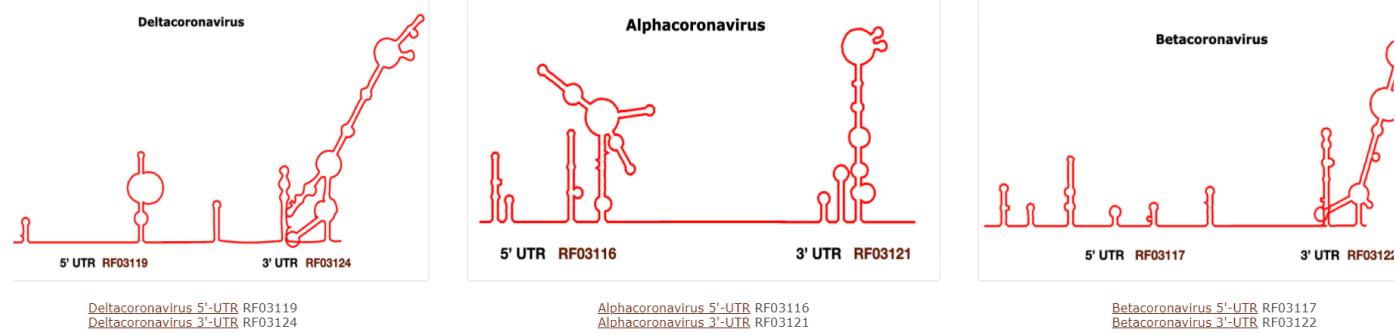
The Rfam database is a collection of RNA families, each represented by multiple sequence alignments, consensus secondary structures and covariance models (CMs). [More...](#)

Rfam COVID-19 Resources

In response to the SARS-CoV-2 outbreak, Rfam produced a special **release 14.2** that includes new and updated Coronavirus families. Read more in a [blog post](#) or [download the data](#)

Coronavirus UTR families

Untranslated regions (UTR) are important functional elements that have conserved secondary structure and are responsible for multiple functions, including replication and packaging. The following new families represent the 5'- and 3'-UTRs of the Coronavirus genomes:



Rfam

<http://rfam.xfam.org/>

A database of curated RNA families [<http://rfam.xfam.org/families>]

BLAST queries, then Infernal HMM alignments [<http://infernal.janelia.org/>]

Rfam 12.0: updates to the RNA families database

Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, et al.
Nucleic Acids Res. 2015. 43:D130-D137 Database issue

Infernal 1.1: 100-fold faster RNA homology searches

Nawrocki EP, Eddy SR
Bioinformatics 2013. 29:2933-2935



Conserved Domains

Limits Advanced search

Search

Help

Structure Group

3D Macromolecular Structures

Conserved Domains

Conserved Domains and Protein Classification

[OVERVIEW](#) **SEARCH** [HOW TO](#) [HELP](#) [NEWS](#) [FTP](#) [PUBLICATIONS](#) [DISCOVER](#)

Resources

Conserved Domain Database (CDD)

CDD is a protein annotation resource that consists of a collection of well-annotated multiple sequence alignment models for ancient domains and full-length proteins. These are available as position-specific score matrices (PSSMs) for fast identification of conserved domains in protein sequences via RPS-BLAST. CDD content includes NCBI-curated domains, which use 3D-structure information to explicitly define domain boundaries and provide insights into **sequence/structure/function relationships**, as well as domain models imported from a number of external source databases (Pfam, SMART, COG, PRK, TIGRFAMs).



Database Statistics

CDD v3.18, as of 30 April 2020:

59,951 total models from all [Source Databases](#)
 16,212 models from [NCBI CDD curation effort](#)
 183 models from [NCBIfams](#)
 1,012 models from [SMART](#) v6.0
 17,919 models from [PFAM](#) v32
 4,871 models from [COGs](#) v1.0
 10,885 models from [NCBI Protein Clusters](#)
 4,488 models from [TIGRFAM](#) v15
 organized into 4,381 multi-model [Superfamilies](#)

Click on the numbers above to retrieve the domain records from CDD; click on the source database names for additional details.

NCBI CDD & CD-Search

Conserved Domain Database

<http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>

<https://www.ncbi.nlm.nih.gov/books/NBK279685/>

deltablast -query `query.fsa` -db `cdd_delta`

We can also run it from the command line with NCBI BLAST+

DB available at: https://ftp.ncbi.nlm.nih.gov/blast/db/cdd_delta.tar.gz

CDD: NCBI's conserved domain database
Marchler-Bauer A et al.
Nucl. Acids Res. 2015. 43:D222–D226
Database issue



Classification of protein families

Home Search Browse Results

Search InterPro

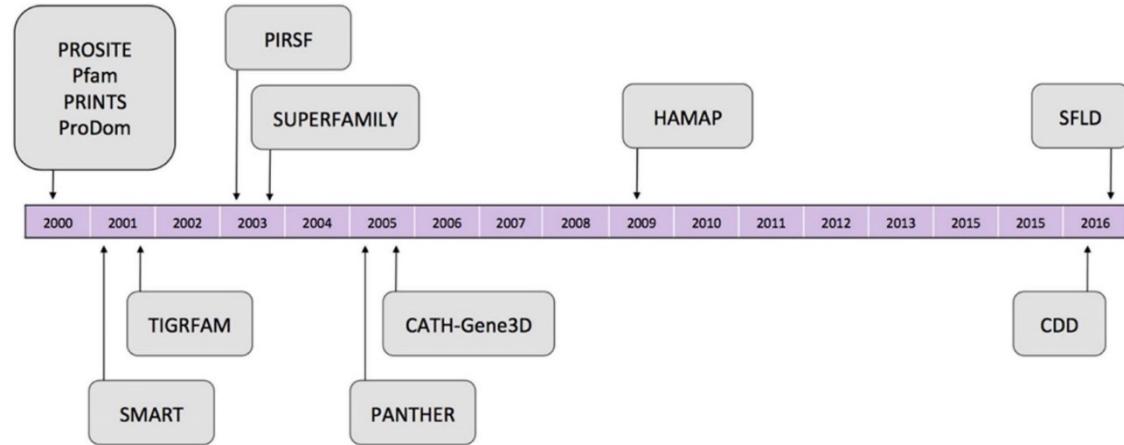
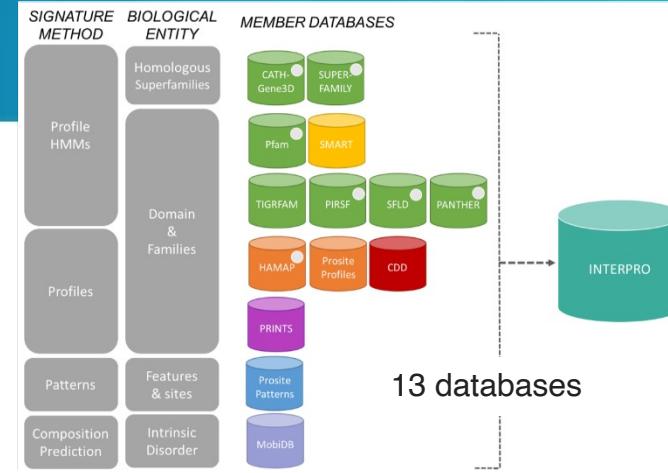


Figure 2. Timeline showing the member databases that have joined InterPro since version 1.0, released in 2000.

InterPro in 2017—beyond protein family and domain annotations

Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, et al.
Nucleic Acids Res. 2017. 45:D190-D199 Database issue

InterProScan

<http://www.ebi.ac.uk/interpro/search/sequence-search>

Searches for multiple evidences (acts as an aggregator)

Queries Pfam, CDD, Gene ontologies, KEGG pathways, and more

The InterPro protein families database: the classification resource after 15 years
Mitchell A, Chang HY, Daugherty L, Fraser M, Hunter S, Lopez R, et al.
Nucleic Acids Res. 2015. 43:D213-D221 Database issue

InterProScan 5: genome-scale protein function classification
Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al.
Bioinformatics 2014. 30(9):1236-1240



Cellular Component

- What is it part of?
- Organelle, macromolecular structure, etc.

Biological Process

- A series of events that accomplishes something
- Not a pathway, a specific result

Molecular Function

- Specific functional activity
- Catalytic activity, transporter activity, etc.

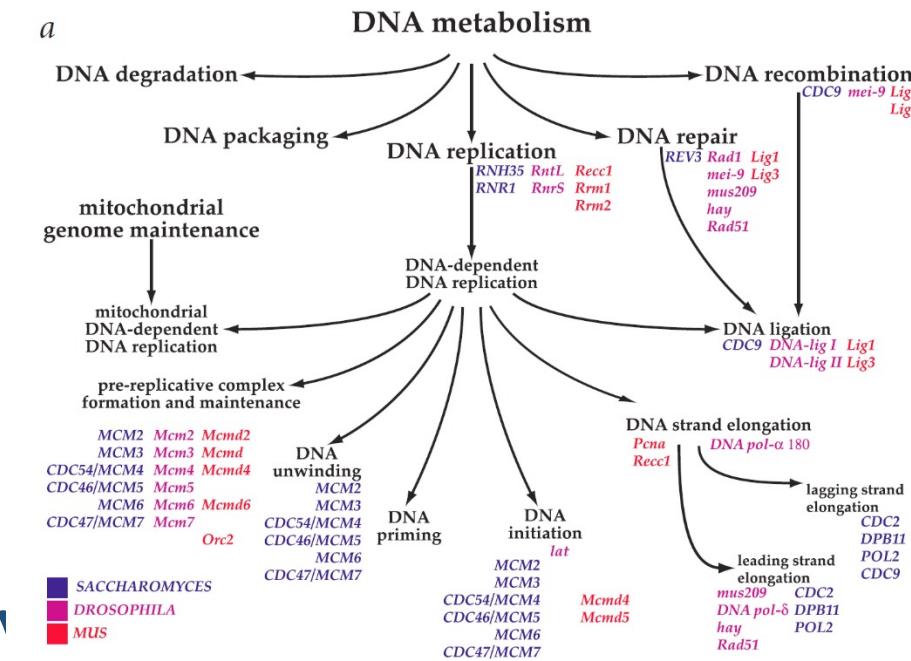
Gene ontologies (GOs)

<http://geneontology.org/>

Structured product descriptions

Consistent across databases

Directed Acyclic Graphs



Gene ontology: tool for the unification of biology
Ashburner M, Ball CA, Blake JA, Botstein D, Butler H,
Cherry JM, Davis AP, et al.

Nat Genet. 2000. 25(1):25-29



KEGG ▾

Example of a KEGG metabolic pathway

Search Help

» Japanese

[KEGG Home](#)
[Release notes](#)
[Current statistics](#)

- KEGG Database
- KEGG overview
- Searching KEGG
- KEGG mapping
- Color codes

KEGG: Kyoto Encyclopedia of Genes and Genomes

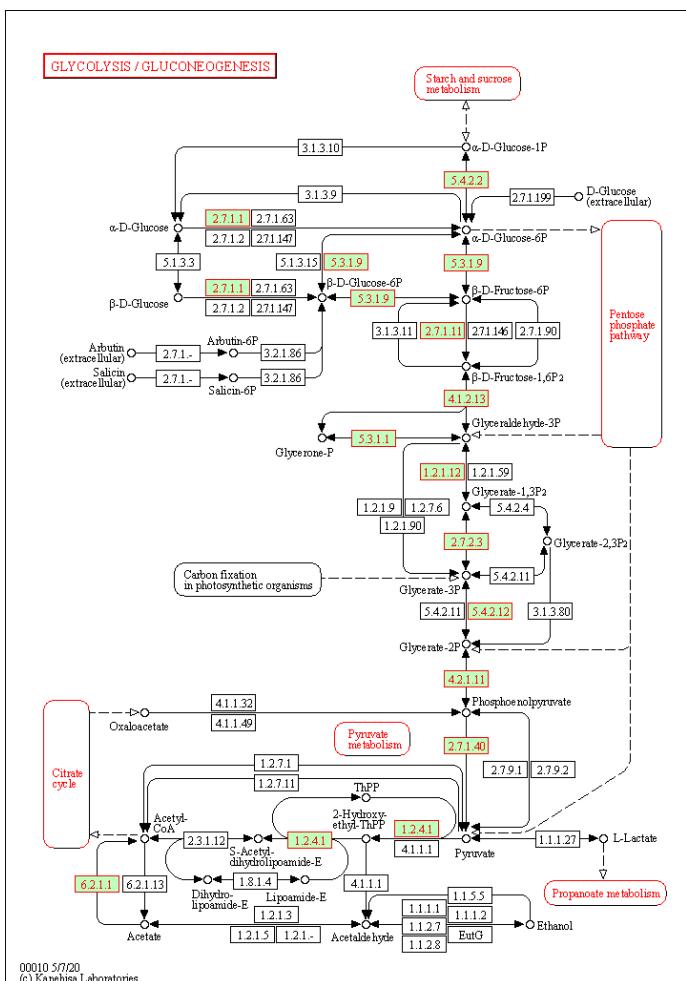
KEGG is a database resource for understanding high-level functions and utilities of the biological system, such as the cell, the organism and the ecosystem, from molecular-level information, especially large-scale molecular datasets generated by genome sequencing and other high-throughput experimental technologies. See [Release notes](#) (October 1, 2020) for new and updated features.

The KEGG databases

<http://www.genome.jp/kegg/kegg2.html>

Kyoto Encyclopedia of Genes and Genomes

The reference for metabolic pathways



KOs instead of GOs; the idea of orthology is very important
Pathways - <http://www.genome.jp/kegg/pathway.html> for pathways

Orthology - <http://www.genome.jp/kegg/ko.html> ## Stores molecular-level functions

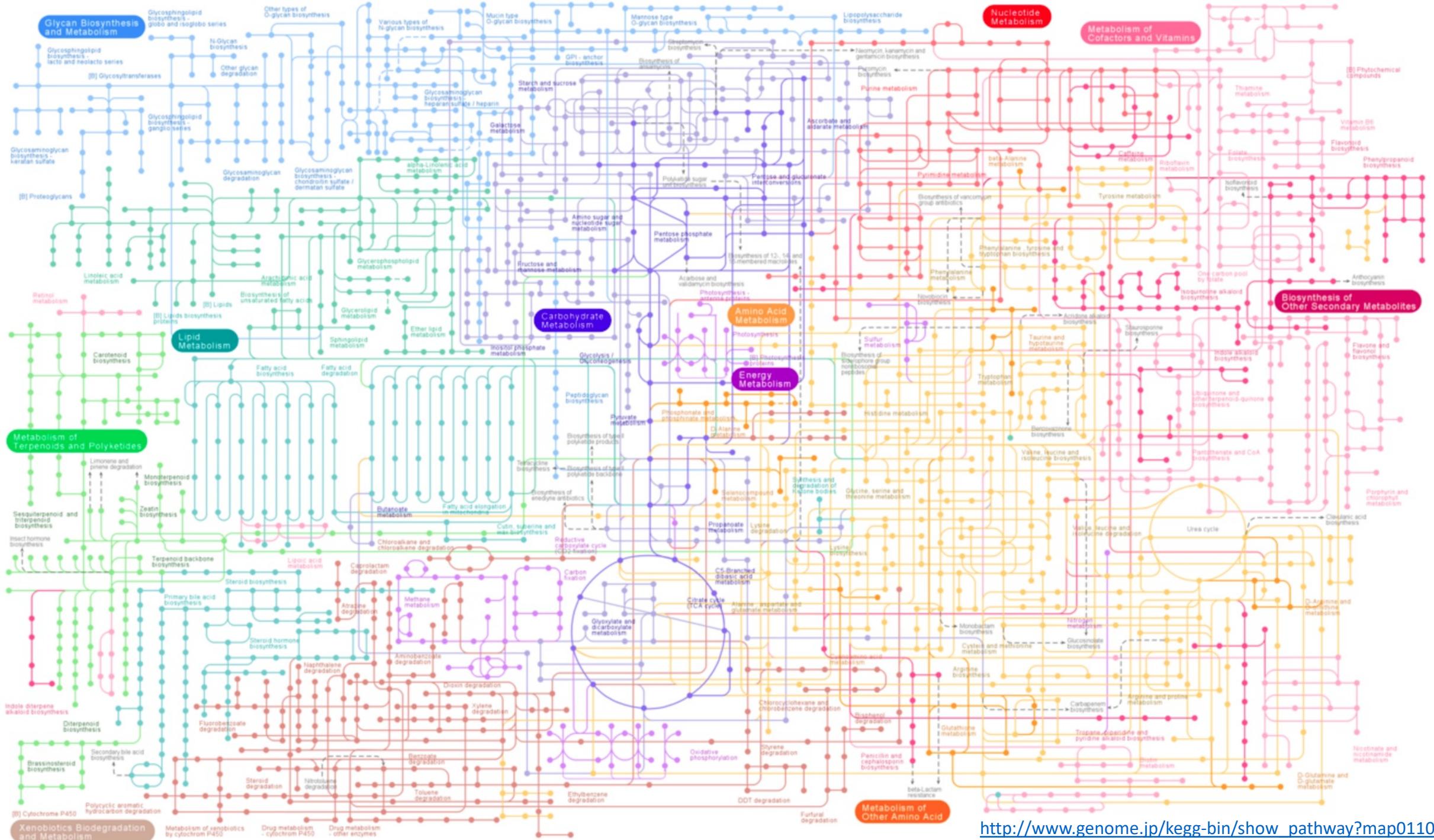
Organisms: http://www.genome.jp/kegg/catalog/org_list.html

Organisms - http://www.genome.jp/kegg/catalog/org_list.html

KEGG: new perspectives on genomes, pathways, diseases and drugs

Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K

Nucleic Acids Res. 2017, 45:D353-D361 Database issue



**ENZYME****Enzyme nomenclature database**

ENZYME is a repository of information relative to the nomenclature of enzymes. It is primarily based on the recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology (IUBMB) and it describes each type of characterized enzyme for which an EC (Enzyme Commission) number has been provided [[More details / References](#)]. ENZYME now includes entries with preliminary EC numbers. Preliminary EC numbers include an 'n' as part of the fourth (serial) digit (e.g. EC 3.5.1.n3).

Release of 07-Oct-20 (6538 active entries)

EC (Enzyme Commission) numbers

<http://enzyme.expasy.org/>

EC numbers are like IP addresses (4 digits separated by periods), e.g. EC 1.14.17.3

Unified **enzyme** nomenclature

CBS >> CBS Prediction Servers >> TargetP-2.0

TargetP-2.0 Server

TargetP-1.1

Click here to run TargetP-1.1.

Predict

Instructions/Help

Data

Abstract/Cite

Portable version

Signal motifs

Signal peptide – secretion

Target peptide – location

Locating proteins in the cell using TargetP, SignalP, and related tools
Emanuelsson O, Brunak S, von Heijne G, Nielsen H
Nat Protoc. 2007. 2:953-971

SignalP 4.0: discriminating signal peptides from transmembrane regions
Petersen TN, Brunak S, von Heijne G, Nielsen H
Nature Methods. 2011 8:785-786

Prediction of Transmembrane Regions and Orientation

The TMpred program makes a prediction of membrane-spanning regions and their orientation. The algorithm is based on the statistical analysis of TMbase, a database of naturally occurring transmembrane proteins. The prediction is made using a combination of several weight-matrices for scoring.

- **K. Hofmann & W. Stoffel (1993)**
TMbase - A database of membrane spanning proteins segments
Biol. Chem. Hoppe-Seyler 374,166

Transmembrane domains

TMHMM	https://services.healthtech.dtu.dk/service.php?TMHMM-2.0
TopPred	http://mobyle.pasteur.fr/cgi-bin/portal.py?#forms::toppred
Tmpred	http://www.ch.embnet.org/software/TMPRED_form.html

Predicting transmembrane protein topology with a hidden Markov model:
application to complete genomes

Krogh A, Larsson B, von Heijne G, Sonnhammer EL

J Mol Biol. 2001. 305:567–580

Membrane protein structure prediction. Hydrophobicity analysis
and the positive-inside rule

von Heijne G

J Mol Biol. 1992. 225:487–494

Exercise 10 – CDD

- 1) Go to: <http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>
- 2) Paste the content of **CDD.fasta** in the window
- 3) Click on Browse results for a detailed view
- 4) Select the desired query in the left panel. Click **Show selected query**.
- 5) In the main window, click on one of the domains display to inspect MSA of your query against the curated database sequences.

MSA = Multiple sequence alignment

Exercise 11 – Pfam

- 1) Go to: <http://pfam.xfam.org/>
- 2) You cannot use multifasta files in the sequence search box
- 3) Instead, click on the search link at the top of the page
- 4) Then, select the batch search option on the left
- 5) Upload **Pfam.fasta**. Enter your email address (some searches may take a long time). Submit.
- 6) Look up your emails. Results will be sent to your inbox.
- 7) You can also do the analyses one by one. Try with the 1st one. (The web display is easier to look at)

Exercise 12 – Web InterProScan

- 1) Go to: <http://www.ebi.ac.uk/interpro/>
- 2) The search only works for single fasta files.
- 3) Paste the content of **interpro.fasta**, then submit the job.
- 4) Look at the detailed signatures.
- 5) You can export your search results at the top right of the screen.
- 6) The GO at the bottom of the page are gene ontologies. You can learn more about GOs from <http://geneontology.org/>.

Exercise 13 – Local InterProScan 5

Running InterProScan locally is great for large protein datasets.

- 1) From the [Ex_13](#) folder, run:

```
interproscan.sh \
-i inter.fasta \
-iprlookup \
-goterms \
-pa \
-b output_interpro \
-f GFF3 HTML TSV
## FASTA input file
## InterPro annotations
## Gene ontologies
## Pathways (KEGG)
## Output name
## Produces HTML output
```

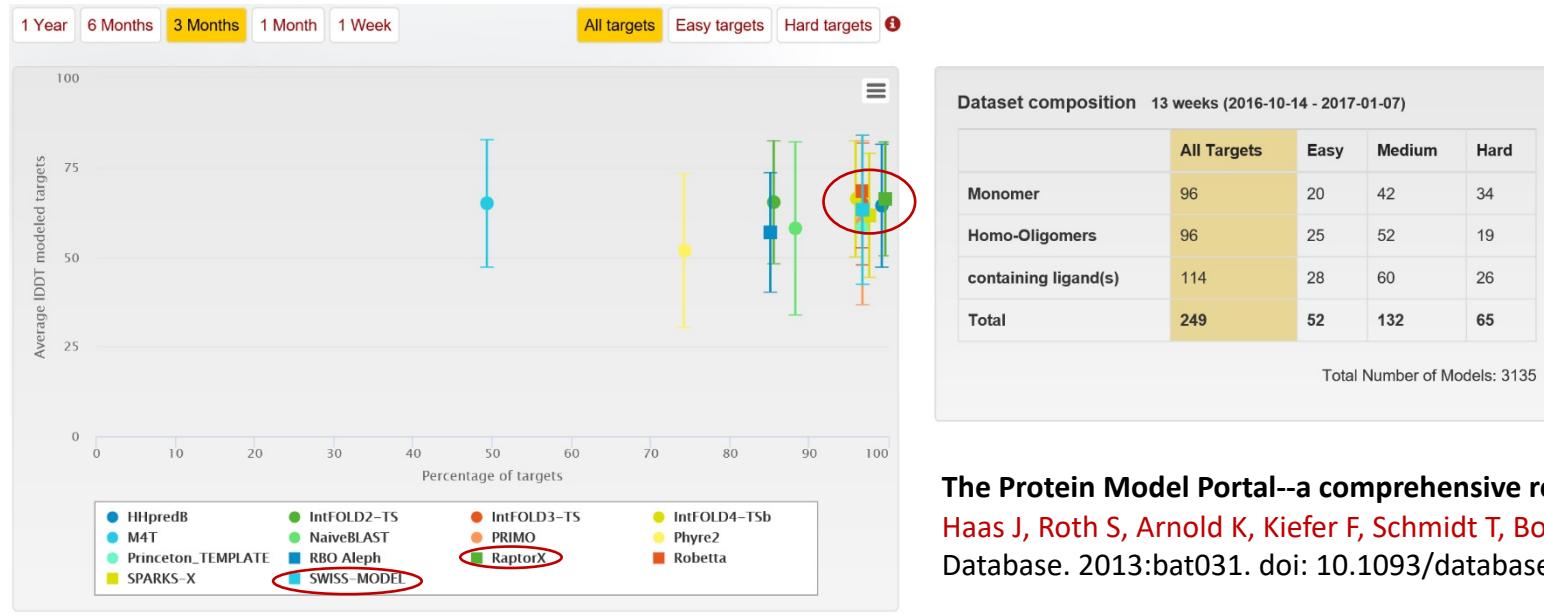
- 2) You will obtain GFF3, HTML and TSV files. The [TSV](#) (tab-separated values) file is the one you really want. [This tab-delimited text file can be parsed easily with Perl](#). It can also be opened with MS Excel or similar spreadsheet tools. Look at it.

Structural homology

- 3D structures often confer functions in biology
- Structural homology can be used as a complement to other analyses

Mandatory car analogy! :

Rubber wheels on cars and wooden wheels on ox carts serve the same function



The Protein Model Portal--a comprehensive resource for protein structure and model information
 Haas J, Roth S, Arnold K, Kiefer F, Schmidt T, Bordoli L, Schwede T
 Database. 2013;bat031. doi: 10.1093/database/bat031.

Predicting 3D structures

Many protein predictors exist

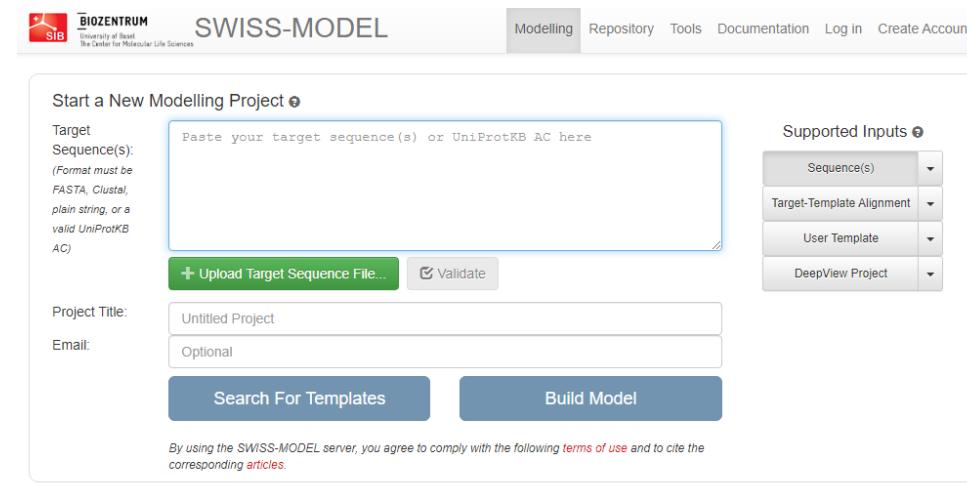
Not every predictor is reliable

No predictor is the best at everything

CAMEO C_{ontinuous} A_{utomated} M_{odel} E_{valuatiOn} – <http://cameo3d.org/>

Welcome to SWISS-MODEL

SWISS-MODEL is a fully automated protein structure homology-modelling server, accessible via the ExPASy web server, or from the program DeepView (Swiss Pdb-Viewer). The purpose of this server is to make protein modelling accessible to all life science researchers worldwide.

[Start Modelling](#)

The screenshot shows the SWISS-MODEL web interface. At the top, there's a header with the BIOZENTRUM SIB logo, the text "SWISS-MODEL", and navigation links for "Modelling", "Repository", "Tools", "Documentation", "Log in", and "Create Account". Below the header is a main form titled "Start a New Modelling Project". It has a text area for "Target Sequence(s)" with instructions: "(Format must be FASTA, Clustal, plain string, or a valid UniProtKB AC)". There's also a "Sequence(s)" dropdown under "Supported Inputs". Below the sequence input are fields for "Project Title" (set to "Untitled Project") and "Email" (set to "Optional"). At the bottom of the form are two buttons: "Search For Templates" and "Build Model". A note at the bottom states: "By using the SWISS-MODEL server, you agree to comply with the following [terms of use](#) and to cite the corresponding [articles](#)".

SWISS-MODEL

<https://swissmodel.expasy.org/>

Predicts 3D structures of proteins from FASTA files

Quick and easy to use

Web-server is hosted by the Swiss Institute of Bioinformatics

All Projects

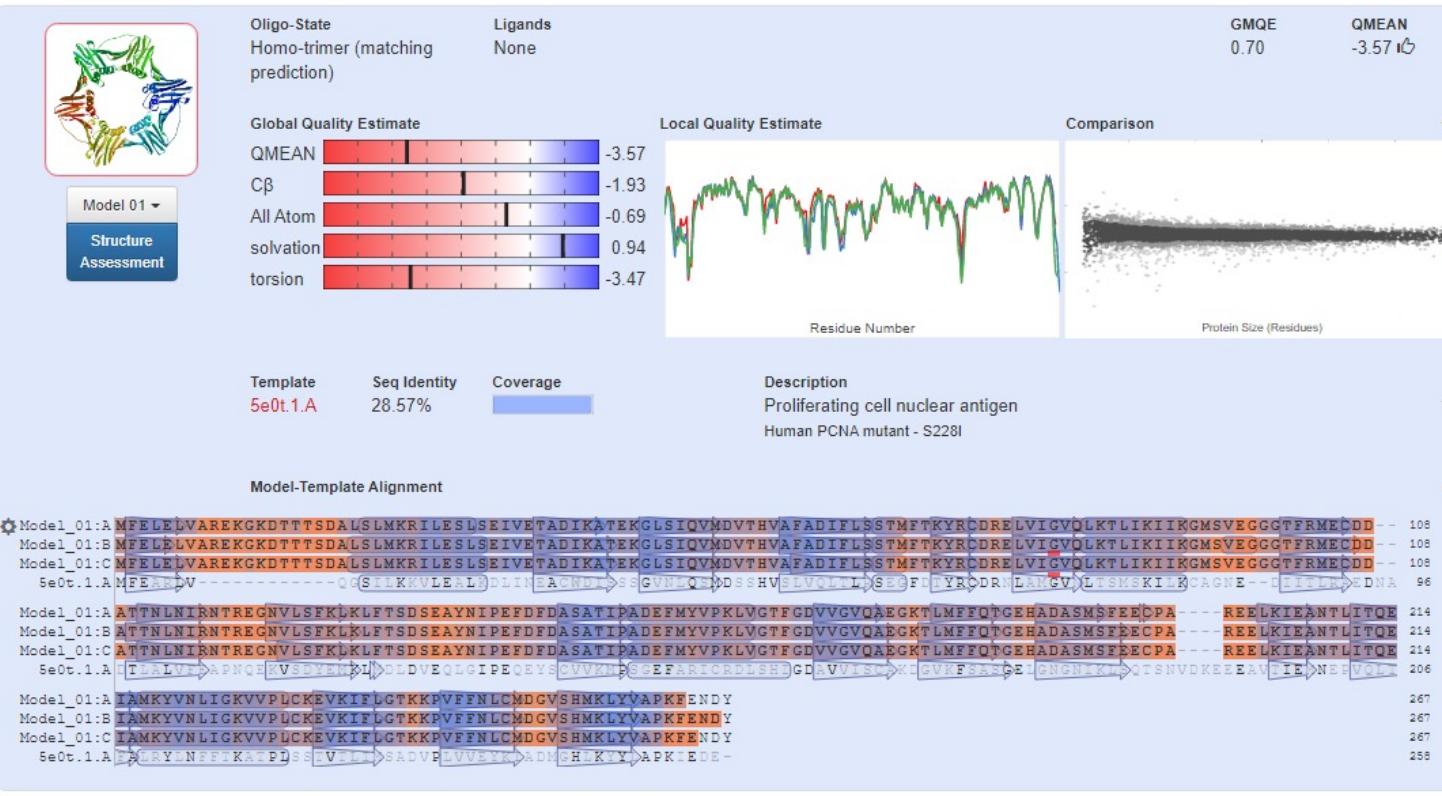
Untitled Project Created: today at 16:50

To download results

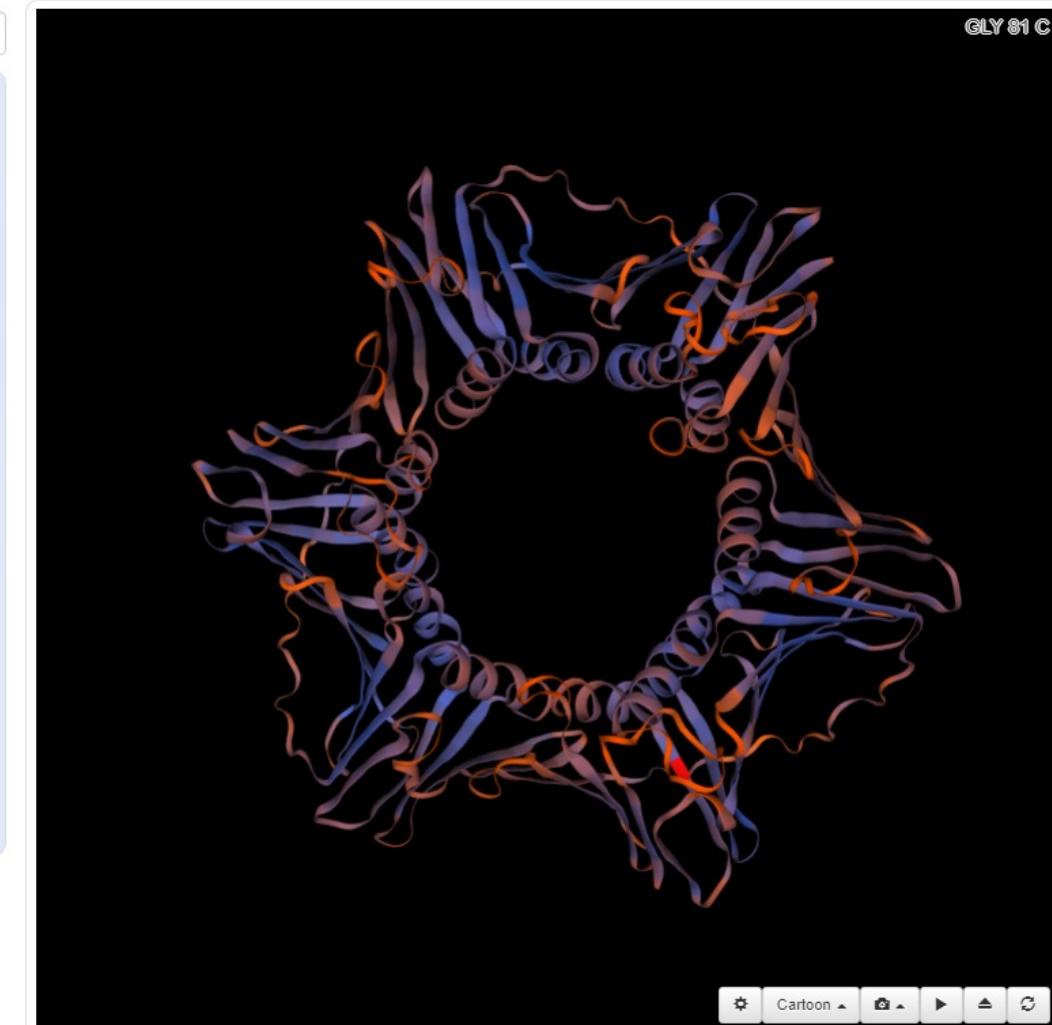
Summary Templates 50 Models 1



Model Results



Example of results with SWISS-MODEL



We can play with structure here



My Jobs | Docs | Download | Inquiry & Bug Report | About | Xu Group

RaptorX: a Web Portal for Protein Structure and Function Prediction

This web portal for protein structure and function prediction is developed by Xu group, excelling at secondary, tertiary and contact prediction for protein sequences without close homologs in the Protein Data Bank (PDB). Given a protein sequence, RaptorX predicts its secondary and tertiary structures as well as contact map, solvent accessibility, disordered regions and binding sites. RaptorX assigns the following confidence scores to indicate the quality of a predicted 3D model: P-value for the relative global quality, GDT (global distance test) and uGDT (un-normalized GDT) for the absolute global quality, and RMSD for the absolute local quality of each residue in the model. RaptorX-Binding predicts the binding sites of a protein sequence, based upon the predicted 3D model by RaptorX. More details can be found [HERE](#).

Template-based protein structure modeling using the RaptorX web server

Källberg M, Wang H, Wang S, Peng J, Wang Z, Lu H, Xu J

Nat Protoc. 2012. 7(8):1511-22. doi: 10.1038/nprot.2012.085

Protein structure alignment beyond spatial proximity

Wang S, Ma J, Peng J, Xu J

Sci Rep. 2013. 3:1448. doi: 10.1038/srep01448

Alignment of distantly related protein structures: algorithm, bound and implications to homology modeling

Wang S, Peng J, Xu J

Bioinformatics. 2011. 27(18):2537-45

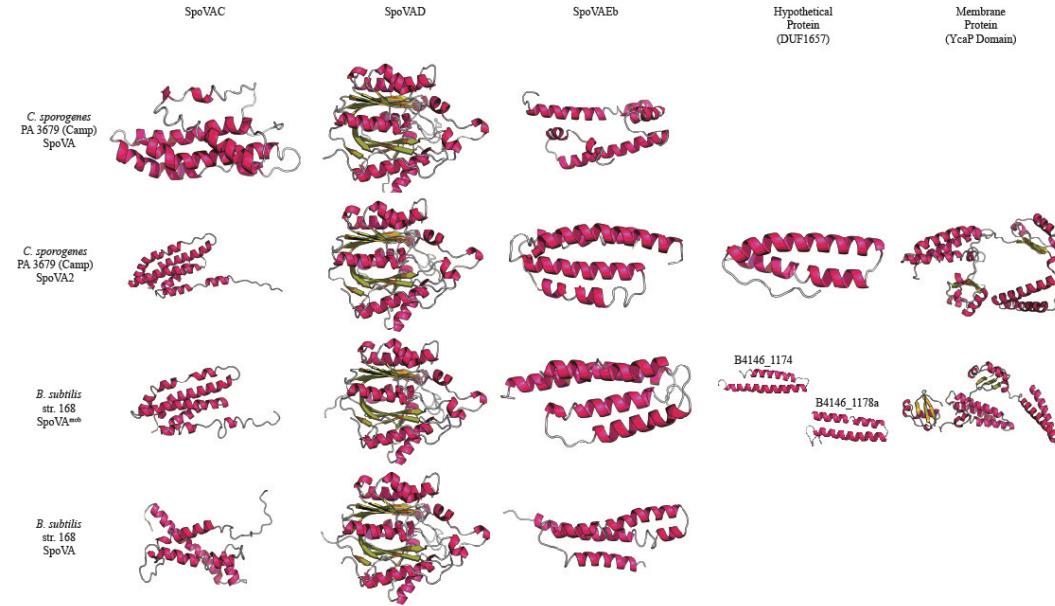
RaptorX

Developed at University of Chicago – <http://raptorg.uchicago.edu/>

One of the top predictors

Easy to use web interface

Can align protein to references



Predicted protein structures for stage V spore formation proteins. Structures predicted using sequences from Supplementary Table T1 using the RaptorX web server (Källberg et al. 2012).

Genetic characterization of the exceptionally high heat resistance of the non-toxic surrogate *Clostridium sporogenes* PA 3679

Butler RR 3rd, Schill KM, Wang Y, Pombert JF

Front Microbiol. 2017 Apr 3;8:545. doi: 10.3389/fmicb.2017.00545. eCollection 2017

Hypothetical Protein
(DUF1657)

Membrane Protein
(YcaP Domain)

AlphaFold Protein Structure Database

Developed by DeepMind and EMBL-EBI

Search for protein, gene, UniProt accession or organism

BETA

Search

Examples: [Free fatty acid receptor 2](#) [At1g58602](#) [Q5VSL9](#) [E. coli](#) Help: [AlphaFold DB search help](#)

<https://alphafold.ebi.ac.uk/>

AlphaFold2

Developed by DeepMind – <https://deepmind.com/>

Available as a docker image

Based on deep learning

The best tool in CASP14

Leverages TensorFlow

Highly accurate protein structure prediction with AlphaFold
Jumper J. et al.

Nature. 2021. PMID: 34265844

<https://github.com/deepmind/alphafold>



- PDBeFold links
 - [FAQ](#)
 - [Visualisation](#)
 - [Performance](#)
 - [Privacy](#)
 - [Version log](#)
 - [PDBeFold Links](#)
 - [Comparisons](#)
 - [Publications](#)
 - [PDBeFOLD tutorial](#)

PDBeFold. Structure Similarity.

PDBeFold functionality:

- pairwise comparison and 3D alignment of protein structures
- multiple comparison and 3D alignment of protein structures
- examination of a protein structure for similarity with the whole [PDB archive](#) or [SCOP archive](#)
- best Ca-alignment of compared structures
- download and visualisation of best-superposed structures using [Rasmol](#) (Unix/Linux platforms), [Rastop](#) (Windows machines) and [Jmol](#)(platform-independent server-side java viewer)
- linking the results to other services - [PDBeMotif](#), [SCOP](#), [GeneCensus](#), [FSSP](#), [CATH](#), [PDBSum](#), [UniProt](#)

[Launch PDBeFold](#)

PDBeFold

<https://www.ebi.ac.uk/msd-srv/ssm/>

Searches for 3D structural homology

Requires a protein structure in PDB format

Predicted or determined experimentally

Query	Target
Source: Coordinate file <input type="button" value="Choose file"/> No file chosen	Source: Whole PDB archive <input type="button"/>
<input type="button" value="Select chains"/> <input type="button" value="Find chains"/>	
Chains: *(all) <input type="button"/>	
Lowest acceptable match (%) <input type="text" value="70"/>	Lowest acceptable match (%) <input type="text" value="70"/>
<input checked="" type="checkbox"/> match individual chains	<input checked="" type="checkbox"/> best matches only
<input checked="" type="checkbox"/> match connectivity	<input checked="" type="checkbox"/> unique matches only
<input checked="" type="checkbox"/> if no matches within limits of acceptability are found, show close ones	
Precision: <input type="button" value="normal"/>	Sort by: <input type="button" value="Q-score"/>
Viewer: <input type="button" value="Jmol"/>	

[Home](#) [Submit your query](#)

Structure Alignment Results.

Query: ECU05_1030_1sxjF.raptorx.pdb, 267 residues

Examined 163976 entries, (442651 chains). Displaying Matches 1-20 of 379.

[Back to query](#) [next](#) [last page](#) Sort by [Q-score](#) [arrange by SCOP family](#) match jump

[Example of results with PDBeFold](#)

##	Scoring 		RMSD	N_align	N_g	%seq	Query				Target (PDB entry)				Title
	Q	P	Z				%sse	Match	%sse	N_res	x				
1	0.90	41.8	19.3	0.42	252	2	22	91	1sxj:F	91	258	<input type="checkbox"/>	CRYSTAL STRUCTURE OF THE EUKARYOTIC CLAMP LOADER (REPLICATION FACTOR C, RFC) BOUND TO THE DNA SLIDING CLAMP (PROLIFERATING CELL NUCLEAR ANTIGEN, PCNA)		
2	0.84	41.8	19.4	0.91	250	3	22	100	6e49:B	96	255	<input type="checkbox"/>	PIF1 PEPTIDE BOUND TO PCNA TRIMER		
3	0.83	41.4	19.3	0.90	249	3	22	100	6b8i:B	96	255	<input type="checkbox"/>	ROLE OF THE PIF1-PCNA COMPLEX IN POL DELTA DEPENDENT STRAND DISPLACEMENT DNA SYNTHESIS AND BREAK INDUCED REPLICATION		
4	0.83	40.7	19.1	0.97	250	3	22	100	6e49:A	96	255	<input type="checkbox"/>	PIF1 PEPTIDE BOUND TO PCNA TRIMER		
5	0.83	41.3	19.3	0.84	246	4	22	100	6e49:C	96	254	<input type="checkbox"/>	PIF1 PEPTIDE BOUND TO PCNA TRIMER		
6	0.83	40.4	19.1	0.92	248	4	22	100	6b8i:A	96	255	<input type="checkbox"/>	ROLE OF THE PIF1-PCNA COMPLEX IN POL DELTA DEPENDENT STRAND DISPLACEMENT DNA SYNTHESIS AND BREAK INDUCED REPLICATION		

The maximum possible Q-score is 1.0 ## Anything above 0.7 is very similar structurally



170172 Biological
Macromolecular Structures
Enabling Breakthroughs in
Research and Education

PDB – A database of experimentally determined structures

[Advanced Search](#) | [Browse Annotations](#)



NUCLEIC ACID
DATABASE



Foundation



Welcome

Deposit

Search

Visualize

Analyze

Download

Learn

A Structural View of Biology

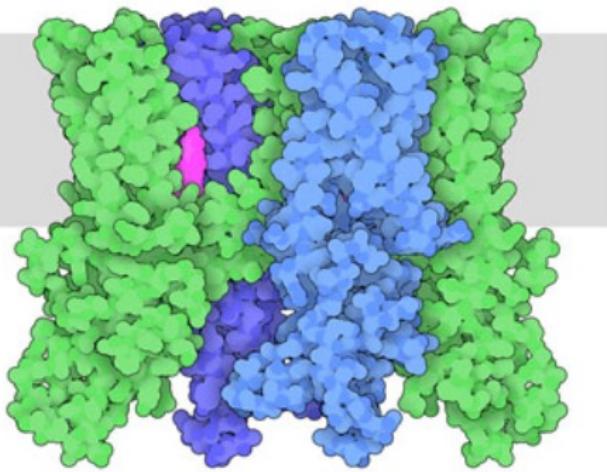
This resource is powered by the Protein Data Bank archive-information about the 3D shapes of proteins, nucleic acids, and complex assemblies that helps students and researchers understand all aspects of biomedicine and agriculture, from protein synthesis to health and disease.

As a member of the wwPDB, the RCSB PDB curates and annotates PDB data.

The RCSB PDB builds upon the data by creating tools and resources for research and education in molecular biology, structural biology, computational biology, and beyond.



October Molecule of the Month



Capsaicin Receptor TRPV1

SWISS-MODEL and the Swiss-PdbViewer: An environment for comparative protein modeling

Guex N, Peitsch MC

Electrophoresis. 1997. 18:2714-2723

Defining and searching for structural motifs using DeepView/Swiss-PdbViewer

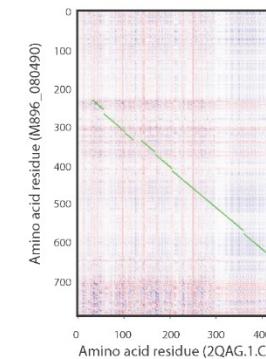
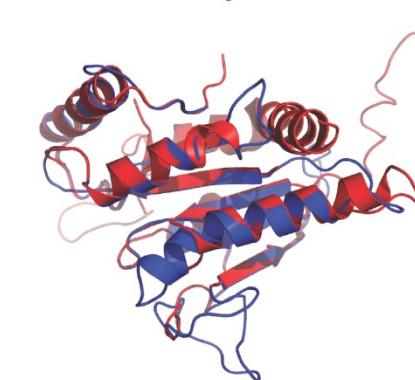
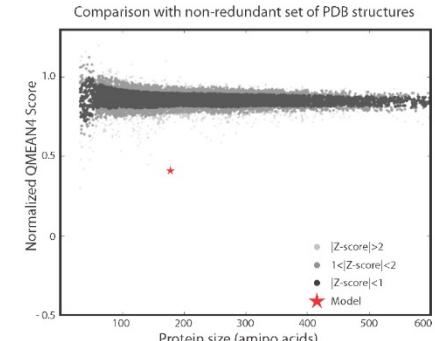
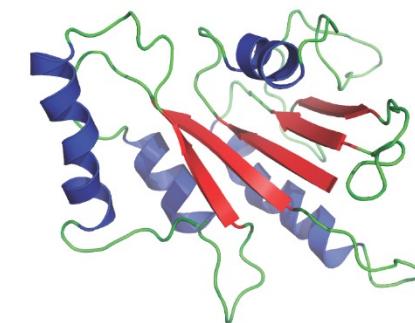
Johansson MU, Zoete V, Michielin O, Guex N

BMC Bioinformatics. 2012. 13:173

UCSF Chimera—a visualization system for exploratory research and analysis

Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE

J Comput Chem. 2004. 25(13):1605-12



The *Ordospora colligata* genome: evolution of extreme reduction in microsporidia and host-to-parasite horizontal gene transfer

Pombert JF, Haag KL, Beidas S, Ebert D, Keeling PJ

mBio. 2015. 6(1):e02400-14.

Viewing structures

Swiss-PdbViewer

<http://spdbv.vital-it.ch/>

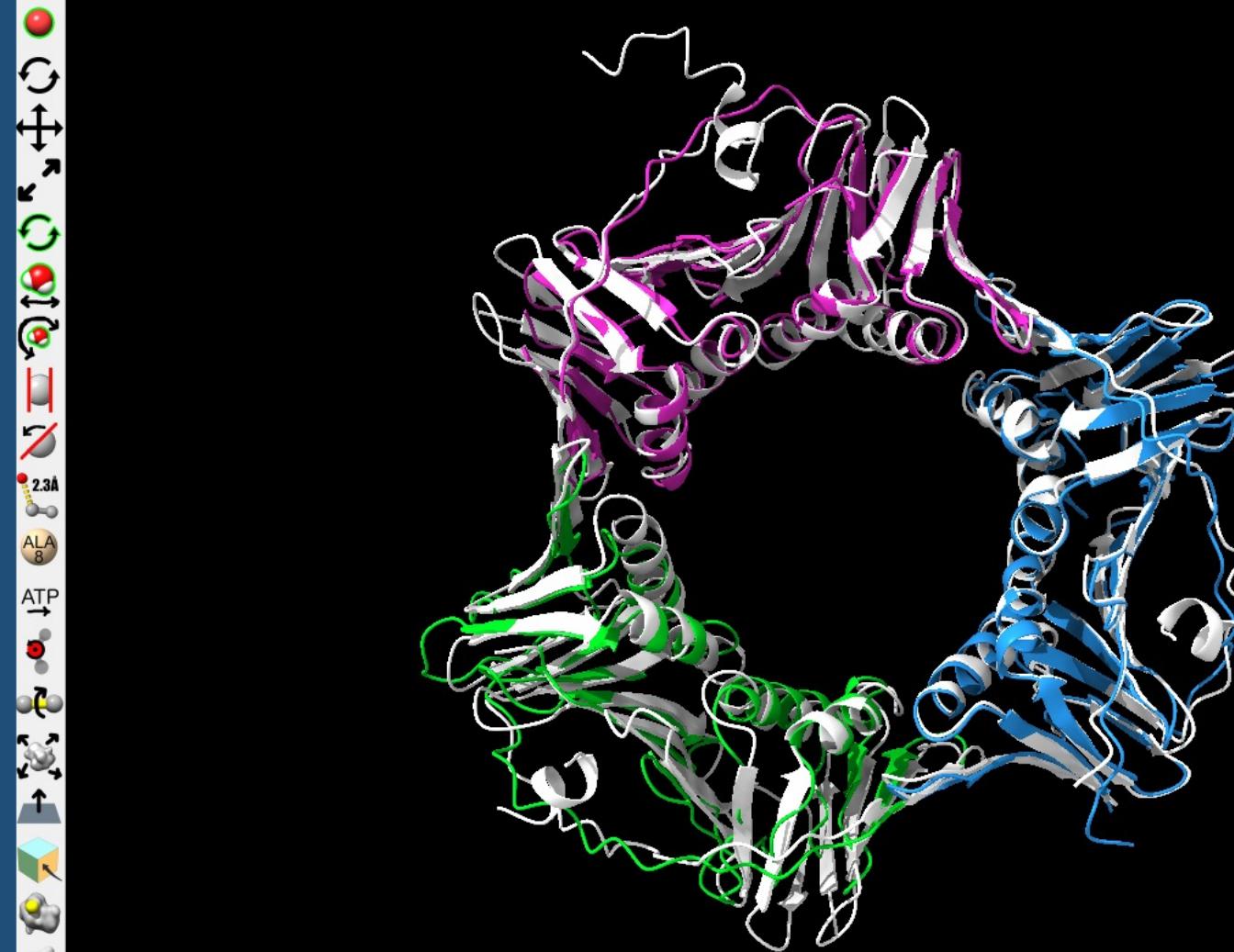
ChimeraX

<https://www.cgl.ucsf.edu/chimerax/>

PyMOL

<https://www.pymol.org/>

Some features are limited to
the non-free version



Log

SS matrix	H 0 -9 -0 S 6 -6 O 4
Iteration cutoff	2

Matchmaker 6e49.pdb, chain A (#1) with ECU08_0130_3lx2A.raptorx.pdb, chain (blank) (#4), sequence alignment score = 92
RMSD between 104 pruned atom pairs is 1.111 angstroms; (across all 203 pairs: 6.291)

matchmaker #4 to #1/c

Parameters		
Chain pairing	bb	
Alignment algorithm	Needleman-Wunsch	
Similarity matrix	BLOSUM-62	
SS fraction	0.3	
Gap open (HH/SS/other)	18/18/6	
Gap extend	1	
	H S O	
SS matrix	H 6 -9 -6 S 6 -6 O 4	
Iteration cutoff	2	

Matchmaker 6e49.pdb, chain C (#1) with ECU08_0130_3lx2A.raptorx.pdb, chain (blank) (#4), sequence alignment score = 92
RMSD between 102 pruned atom pairs is 1.131 angstroms; (across all 203 pairs: 6.297)

Models

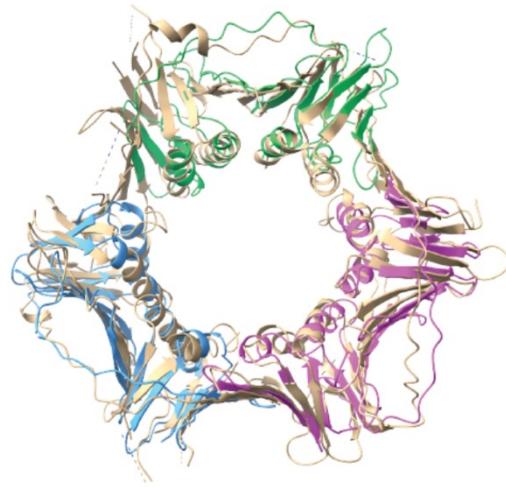
Name	ID	Eye	Close
6e49.pdb	1	<input checked="" type="checkbox"/>	<input type="checkbox"/>
ECU05_1030_1sxjF.raptorx.pdb	2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
ECU07_1290_3k4xA.raptorx.pdb	3	<input checked="" type="checkbox"/>	<input type="checkbox"/>
ECU08_0130_3lx2A.raptorx.pdb	4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Close
Hide
Show
View

Command: match #4 to #1/c



We can match structures against others with the match function; We can rotate, edit, make videos... ## Powerful visualization of 3D structural homology



3DFI

Three-dimensional functional inference using structural homology

Enhanced fold recognition using efficient short fragment clustering
Krissinel E
Mol Biochem. 2012;1(2):76-85. PMID: 27882309 PMCID: PMC5117261

3DFI

<https://github.com/PombertLab/3DFI>

Automates structural homology searches at the 3D level

A Perl pipeline

Leverages **RaptorX/AlphaFold2**, **GESAMT**, and **ChimeraX**

Can be used to query against both predicted and/or experimental structures

Optional

Additional slides and/or exercises for the curious

Exercise 14 – SignalP

- 1) Go to <http://www.cbs.dtu.dk/services/SignalP/>
- 2) Upload **signal_1.fasta**. Use the Eukaryotes group (default)
- 3) Look at the output
- 4) Now upload **signal_2.fasta** and look at the results
- 5) Notice the difference?
 - Signal and target peptides when present are located in the N-terminal portions of the protein
 - A signal peptide motif suggests that the protein is likely secreted

Exercise 15 – TargetP

- 1) Go to <http://www.cbs.dtu.dk/services/TargetP/>
- 2) Upload targetP_1.fasta. Use the Non-plant group (default)
- 3) Look at the output
- 4) Now re-upload targetP_1.fasta using the Plant group instead.
- 5) Notice a difference?
 - Target peptides are not always conserved and differ between various groups
 - Keep in mind that the localization results you obtain are *in silico* predictions only.

Exercise 16 – TMHMM

- 1) Go to <http://www.cbs.dtu.dk/services/TMHMM/>
- 2) Upload **tmhmm.fasta**
- 3) Look at the output
- 4) Here, outside means outside the membrane
- 5) In this protein, only the C-terminal domain is anchored to the membrane