

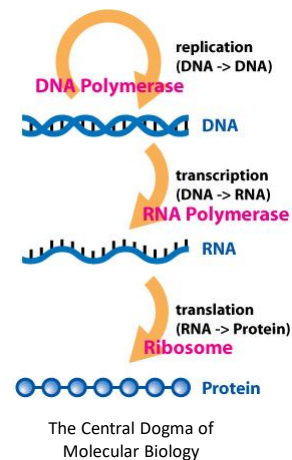
# Basic Bioinformatics, Sequence Alignment, and Homology

Biochemistry Boot Camp 2018  
Session #10  
Nick Fitzkee  
nfitzkee@chemistry.msstate.edu

\* BLAST slides have been adapted from an earlier presentation by W. Shane Sanders.

## Biology Review

- Genome is the genetic material of an organism, normally DNA but RNA possible (viruses)
- Central Dogma:
  - DNA → RNA → Protein



## Primary Structure (Sequence)

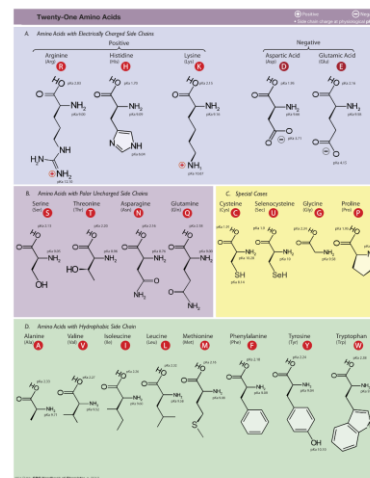
- **DNA and Proteins are chemically complex**, but their “alphabets” are rather simple.
  - 4 nucleobases (A, C, T, G)
  - 20 amino acids
- DNA sequences are represented from 5' to 3'



3

## Primary Structure (Sequence)

- **DNA and Proteins are chemically complex**, but their “alphabets” are rather simple.
  - 4 nucleobases (A, C, T, G)
  - 20 amino acids
- Protein sequences are represented from NT to CT



4

## Storing Sequences

- GenBank ( \*.gb | \*.genbank)
  - National Center for Biotechnology's (NCBI) Flat File Format (text)
  - Provides a large amount of information about a given sequence record
  - <http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>
  - We've seen this before! (Remember NCBI Protein result?)
- FASTA (\*.fasta | \*.fa)
  - Pronounced "FAST-A"
  - Simple text file format for storing nucleotide or peptide sequences
  - Each record begins with a single line description starting with ">" and is followed by one or more lines of sequence
- FASTQ (\*.fastq | \*.fq)
  - Pronounced "FAST-Q"
  - Text based file format for storing nucleotide sequences and their corresponding quality scores
  - Quality scores are generated as the nucleotide is sequenced and correspond to a probability that a given nucleotide has been correctly sequenced by the sequencer
- Text files are also okay in many cases.

5

## Storing Sequences

- |   |  |
|---|--|
| <ul style="list-style-type: none"> <li>• FASTA format</li> <li>• Can represent nucleotide sequences or peptide sequences using single letter codes</li> </ul> | <ul style="list-style-type: none"> <li>• FASTQ format</li> <li>• Represents nucleotide sequences and their corresponding quality scores</li> </ul> |
|---|--|

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus]
LCLYTHIGENNIYGSVLYSEWVQIMLLITMATAGVYLVWQMSFWGATVITWLSAIPYICINLV
EWINGQFSVDKATINRFAPHFILPFWALAGVHLFLHETGSSNNPLGLTSDGKIPFHPYTIKDFLG
LLILILLLLLLALLSDFMLGDFONHMDADFLNTEHLIKEDWYFLFAYAILRSVFNKLGVLALFLSIVIL
GLMPFLHTSKHRSNMLRPLSQALEWLTMDLLTLTWIGSQFVEYPTTIIGQMASILYFSIILAFLPIAGX
IENY
```

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAAATAGTAAATCCATTGTTCAACTCACAGTTT
+
!''*(((****+))%%++) (%%%) .1***-+*') **55CCF>>>>>CCCCCCC65
```

6

## Sequence Alignment

Sequence alignment is the procedure of comparing two (pairwise) or more (multiple) sequences and searching for a series of individual characters or character patterns that are the same in the set of sequences.

- **Global alignment** – find matches along the entire sequence (use for sequences that are quite similar)
- **Local alignment** – finds regions or islands of strong similarity (use for comparing less similar regions [finding conserved regions])

7

## Sequence Alignment

Sequence 1: GARVEY

Sequence 2: AVERY

### Global Alignment:

GARVE-Y

-A-VERY

8

## Global Sequence Alignment

- EMBOSS Needle  
[http://www.ebi.ac.uk/Tools/psa/emboss\\_needle/](http://www.ebi.ac.uk/Tools/psa/emboss_needle/)  
– Command line version also available
- Alternative: Biopython (library for the python programming language)
- **Example:** Human vs. Nematode Calmodulin (global sequence #1 and #2)

10

## Global Sequence Alignment

- EMBOSS Needle Options:

How to compare residues?

How much penalty to open a gap in the sequence?

STEP 2 - Set your pairwise alignment options

MATRIX	GAP OPEN	GAP EXTEND	OUTPUT FORMAT
BLOSUM62	10	0.5	pair
END GAP PENALTY	END GAP OPEN	END GAP EXTEND	
false	10	0.5	

Worry about the ends?

How much penalty to have overhang at each end?

11

## Global Sequence Alignment

```
# Length: 149
# Identity:   146/149 (98.0%)
# Similarity: 147/149 (98.7%)
# Gaps:      0/149 ( 0.0%)
# Score: 745.0
```

Percent Identity and Similarity quantify alignment.

Human	1	MADQLTEEQIAEFKEAFSLFDKDGDTITTTKELGTVMRSLGQNPTEAELQ	50
Nematode	1	MADQLTEEQIAEFKEAFSLFDKDGDTITTTKELGTVMRSLGQNPTEAELQ	50
Human	51	DMINEVDADGNGTIDFPEFLTMMARKMKDIDSEEEIREAFRVFDKDGNGY	100
Nematode	51	DMINEVDADGNGTIDFPEFLTMMARKMKDIDSEEEIREAFRVFDKDGNGF	100
Human	101	ISAAELRHVMTNLGEKLTDEEVDEMIREADIDGGQVNYEEFVQMMTAK	149
Nematode	101	ISAAELRHVMTNLGEKLTDEEVDEMIREADIDGGQVNYEEFVIMMITK	149

- Pretty darn similar!

Identical residues shown with |,  
similar residues with : and ., and  
blanks represent dissimilar  
residues.

12

## Multiple Sequence Alignment

- Align many sequences simultaneously, normally from multiple organisms
- Mathematically much more challenging, and requires assumptions about data analysis
- Results can be used to generate phylogenetic tree
  - <https://www.ebi.ac.uk/Tools/msa/clustalo/>
- Example software: MEGA, ClustalX
  - <http://www.megasoftware.net/>
  - <http://www.clustal.org/>



13

## MSA Example

```

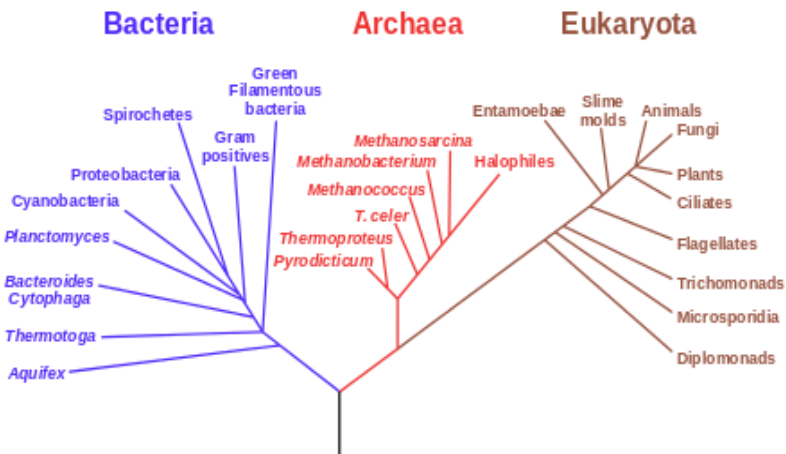
Q5E940 BOVIN -----M*PREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 HUMAN -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 MOUSE -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 RAT -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 CHICK -----MPREDRATWKSNYFMKTIIO*LLDDY*PKCFVVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 RANSY -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
Q7ZUG3 BRARE -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 ICTPU -----MPREDRATWKSNYFLKTIIO*LLDDY*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 DROME -----MVRENKRAAKQAYEIKVVELEDEF*PKCFIVGADNVGSKMQQIRMSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 DICDI -----MSGAG--SKRKKLFIEKATKLFITT*DKMIVAEADFVGS*SLQKIRKSIRGI--GAVLMGKNTMMRKAIRGHLENN--PALE 75
Q54LP0 DICDI -----MSGAG--SKRKNVFIEKATKLFITT*DKMIVAEADFVGS*SLQKIRKSIRGI--GAVLMGKNTMMRKAIRGHLENN--PALE 75
RLA0 PLAF8 -----MAKLSQKKQKQMYIEKLS*LIQQ*SKILIVHYD*VGS*NMASVRSKSLRGK--AVVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0 SULAC -----HICLAVITTTKKIAKKVDEVAELTE*KLKT*KTILLIANIEGFPADKLHEIRKKLRGK--ADIKVTKNLFNIALKNAG----VDIK 79
RLA0 SULFO -----MKIMAVITQERKIAKKVDEVAELTE*KLKT*KTILLIANIEGFPADKLHEIRKKLRGK--ADIKVTKNLFNIALKNAG----VDIK 80
RLA0 SULSO -----MKRLALAKQKQVSSKLEVEKELTE*KLKT*KTILLIANIEGFPADKLHEIRKKLRGK--ADIKVTKNLFNIALKNAG----VDIK 80
RLA0 AERPE *MSVSVSEVQMYKREK*P*PE*KTLLMIRELEELFSKRVVLFADITG*P*FVY*RVYKKLWKK--VPMVAKKRLILEAMKAGGLE----LDDN 86
RLA0 PYRAE *MMLAIGKRRYVTRQ*P*PARKVIVSEATELLQK*P*PYVFLFDLHGLSSRILHEVRYRLRY--GVIKIIP*PLFKIAFTKVVGG----IPAE 85
RLA0 METAC -----MAEERHHTTEHIPQWKDEIENIKELIQSKVFGMV*GIEGILATKMKIRRDLDV--AVLKVRNTLTERALNQLG----ETIP 78
RLA0 METMA -----MAEERHHTTEHIPQWKDEIENIKELIQSKVFGMV*GIEGILATKMKIRRDLDV--AVLKVRNTLTERALNQLG----ETIP 78
RLA0 ARCFU -----MAAVRGS-----PPEYKVRAVEEIKRMISSEPVVAVISFRNV*PAGCMKIRREFRGK--AEIKVTKNLFNIALKNAG----VDIK 75
RLA0 METKA *MAVKAKGQPS*G*GKVAEKKRREVKEL*E*EMDE*EVNGLVD*ESLPAPQLAEIRAKLRER*GLIRMSRNTLMRIALEEKLDER--PELE 80
RLA0 METH -----MAHVAEWKKEVEELANLIKSPVIALVDYSSMPAYPLSQMRRLIRENGCLLRVSRNTLIELAKKAAKEIGKPELE 77
RLA0 METTL -----MITAESEHKIAPWKIEEVNALKELLLKMQIVALVDMMEV*PARLQEI*RD*IR*ET*MLKMSRNTLIERAKVEAEETGNPEFA 82
RLA0 METVA -----MIDAKSEHKIAPWKIEEVNALKELLLKMQIVALVDMMEV*PARLQEI*RD*IR*ET*MLKMSRNTLIERAKVEAEETGNPEFA 82
RLA0 METJA -----METKYKAHVAPWKIEEVNALKELLLKMQIVALVDMMEV*PARLQEI*RD*IR*ET*MLKMSRNTLIERAKVEAEETGNPEFA 81
RLA0 PYRAB -----MAHVAEWKKEVEELANLIKSPVIALVDYSSMPAYPLSQMRRLIRENGCLLRVSRNTLIELAKKAAKEIGKPELE 77
RLA0 PYRHO -----MAHVAEWKKEVEELANLIKSPVIALVDYSSMPAYPLSQMRRLIRENGCLLRVSRNTLIELAKKAAKEIGKPELE 77
RLA0 PYRFO -----MAHVAEWKKEVEELANLIKSPVIALVDYSSMPAYPLSQMRRLIRENGCLLRVSRNTLIELAKKAAKEIGKPELE 77
RLA0 PYRKO -----MAHVAEWKKEVEELANLIKSPVIALVDYSSMPAYPLSQMRRLIRENGCLLRVSRNTLIELAKKAAKEIGKPELE 77
RLA0 HALMA -----MSESEERKTETIP*EKQKEVEEVAIV*E*IES*ESVGVVNTA*IP*RL*LD*MR*RD*H*ET--AFLRVSRNTLIERALDDVD----DELE 79
RLA0 HALVO -----MSESEVRQTEV*IPQWKREVEEVDLVD*IES*ESVGVVNTA*IP*RL*LD*MR*RD*H*ET--AFLRVSRNTLIERALDDVD----DELE 79
RLA0 HALSA -----MSESEVRQTEV*IPQWKREVEEVDLVD*IES*ESVGVVNTA*IP*RL*LD*MR*RD*H*ET--AFLRVSRNTLIERALDDVD----DELE 79
RLA0 THEAC -----MKEVVSQKKELVNEIT*RIKASRSVAIVD*ACIR*RI*Q*ID*IR*GKNRQK--INLKVIKTL*LLFKALENLGD----EKIS 72
RLA0 THEVO -----MRKINPKKKEIVSELA*Q*ITKSKAVALVDIKGV*IR*Q*MODIRAKNRQK--VKIKVVKTL*LLFKALENLGD----EKIT 72
RLA0 PICTO -----MTE*PKKK*ID*F*YKNLENE*INSKVA*VAIVS*IKGLRNNE*FQKIRNS*IRDK--ARIKVEARLLRLALENLGK----NNIV 72
ruler 1. .... 10. .... 20. .... 30. .... 40. .... 50. .... 60. .... 70. .... 80. .... 90

```

MSA of Ribosomal Protein P0 from Wikipedia, "Multiple Sequence Alignment"

14

## MSA-Derived Phylogenetic Tree



Phylogenetic Tree derived from ribosomal proteins, Wikipedia "Phylogenetic Tree"

15

## Why Sequence Alignment?

1. To determine possible functional similarity.
2. For 2 sequences:
  - a. If they're the same length, are they almost the same sequence? (global alignment)
3. For 2 sequences:
  - a. Is the prefix of one string the suffix of another? (contig assembly)
4. Given a sequence, has anyone else found a similar sequence?
5. To identify the evolutionary history of a gene or protein.
6. To identify genes or proteins.

16

## BLAST:

### Basic Local Alignment Search Tool

- A tool for determining sequence similarity
- Originated at the National Center for Biotechnology Information (NCBI)
- Sequence similarity is a powerful tool for identifying unknown sequences
- BLAST is fast and reliable
- BLAST is flexible

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

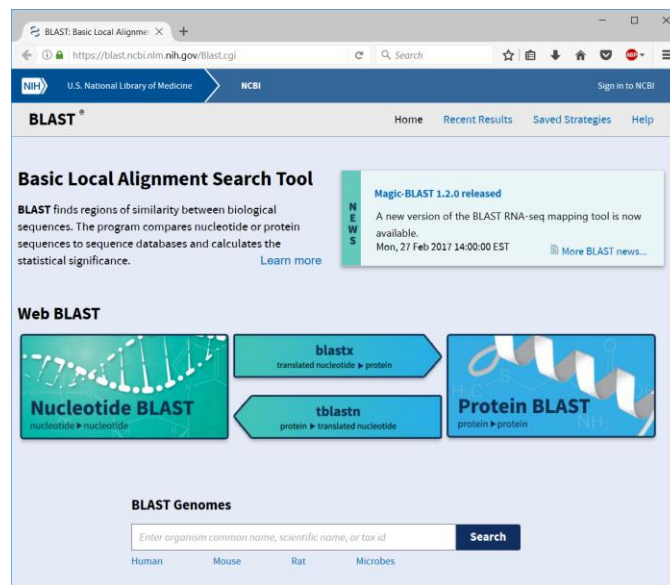
17

# Flavors of BLAST

- **blastn** – searches a nucleotide database using a nucleotide query  
*DNA/RNA sequence searched against DNA/RNA database*
- **blastp** – searches a protein database using a protein query  
*Protein sequence searched against a Protein database*
- **blastx** – search a protein database using a translated nucleotide query  
*DNA/RNA sequence -> Protein sequence searched against a Protein database*
- **tblastn** – search a translated nucleotide database using a protein query  
*Protein sequence searched against a DNA/RNA sequence database -> Protein sequence database*
- **tblastx** – search a translated nucleotide database using a translated nucleotide query  
*DNA/RNA sequence -> Protein sequence searched against a DNA/RNA sequence database -> Protein sequence database*

18

## BLAST Main Page



19

Sequence Input

Databases to Search Against

Program Selection

Click to Run!

Same Page Organization

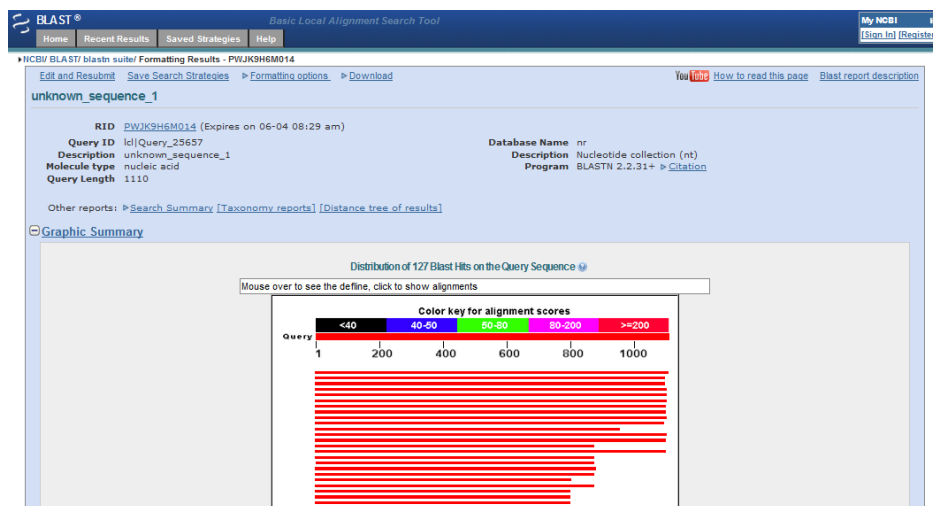
# BLAST Example

- What gene is this?

```
>unknown_sequence_1
TGATGTCAAGACCCCTCTATGAGACTGAAGTCTTTTCTACCGACTTCTCCAACATTTCTGCAGCCAAGCAG
GAGATTAACAGTCATGTGGAGATGCAAACCAAGGAAAGTTGTGGGTCTAATTCAGACCTCAAGCCAA
ACACCATCATGGTCTTAGTGAACATATTTCACTTTAAAGCCCAGTGGGCAAAATCCTTTTGATCCATCCAA
GACAGAAGACAGTTCAGCTTCTTAATAGACAAGACCACCACTGTTCAGTGCCCATGATGCACCAGATG
GAACAATACTATCACCTAGTGGATATGGAATTGAAGTGCACAGTTCTGCAAAATGGACTACAGCAAGAATG
CTCTGGCACTCTTTGTTCTTCCCAAGGAGGGACAGATGGAGTCAGTGGAGCTGCCATGTCATCTAAAC
ACTGAAGAAGTGAACCGCTTACTACAGAAGGGATGGGTTGACTTGTGTTGTTCCAAAGTTTCCATTTCT
GCCACATATGACCTTGAGCCACACTTTTGAAGATGGGCATTCAGCATGCCTATTCTGAAAATGCTGATT
TTTCTGGACTCACAGAGGACAATGGTCTGAAACTTTCCAATGCTGCCCATAGGCTGTGCTGCACATTGG
TGAAAAGGGAAGTGAAGCTGCAGCTGTCCCTGAAGTTGAAGTTTCGGATCAGCCTGAAAACACTTCCCTA
CACCTATTATCCAAATGATAGATCTTTCATGTTGTTGATTTGGAGAGAAGCACAAGGAGTATTTCTCT
TTCTAGGGAAAGTTGTGAACCCAACGGAAGCGTAGTTGGGAAAAAGGCCATTGGCTAATTGCACGTGTGT
ATTGCAATGGGAAATAAATAAATATAGCCTGGTGTGATGTGATGTGAGCTTGGACTTGCATTCCCTTA
TGATGGGATGAAGATTGAACCTGGCTGAAGTTGTTGGCTGTGGAAGAGGCCAATCCTATGGCAGAGCA
TTCAGATGTCAATGAGTAATTCATTATTATCCAAAGCATAGGAAGGCTCTATGTTGTATATTTCTCTT
TGTCAGAATACCCCTCAACTCATTGCTCTAATAAATTTGACTGGGTTGAAAAATTAAAA
```

22

# BLAST Results



23

## Interpreting BLAST Results

- **Max Score** – how well the sequences match
- **Total Score** – includes scores from non-contiguous portions of the subject sequence that match the query
- **Bit Score** – A log-scaled version of a score
  - Ex. If the bit-score is 30, you would have to score on average, about  $2^{30} = 1$  billion independent segment pairs to find a score matching this score by chance. Each additional bit doubles the size of the search space.
- **Query Coverage** – fraction of the query sequence that matches a subject sequence
- **E value** – how likely an alignment can arise by chance
- **Max ident** – the match to a subject sequence with the highest percentage of identical bases

24

## Installing BLAST Locally

Executables and documentation available at:

<ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/>

Documentation:

<http://www.ncbi.nlm.nih.gov/books/NBK1762/>

25

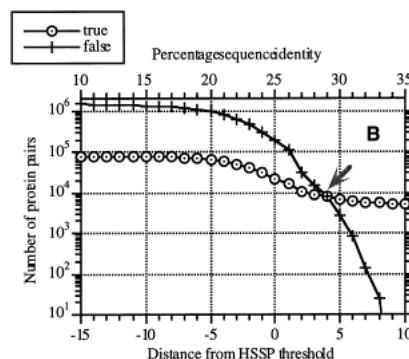
## Aligning via Structure

- So far we've focused on sequence alignment: looking at the primary (DNA or protein) sequence
- What about structural alignment? (Think shape or similar domains)
- VAST (Vector Alignment Search Tool) at NCBI: <https://structure.ncbi.nlm.nih.gov/Structure/VAST/vast.shtml>

26

## Homology Modeling

- Proteins with similar sequences tend to have similar structures.
- When sequence identity is greater than ~25%, this rule is almost guaranteed
  - Exception: See Philip Bryan and “fold switching”
- Can we use this to predict structures?



Below ~28% sequence identity,  
the number of structurally  
dissimilar aligned pairs explodes.

Rost, *Prot. Eng.* 12(2): 85-94

27

## What is Homology Modeling?

- **Consider:** Protein with known sequence, but unknown structure
- Use sequence alignment (protein BLAST) to identify similar sequences with known structures
  - These are termed “template structures”
- “Map” unknown sequence onto known backbone
  - Side chains may be more ill-defined: it’s a model!

28

## Homology Modeling Servers: **SWISS-MODEL**

SWISS-MODEL



### Welcome to SWISS-MODEL

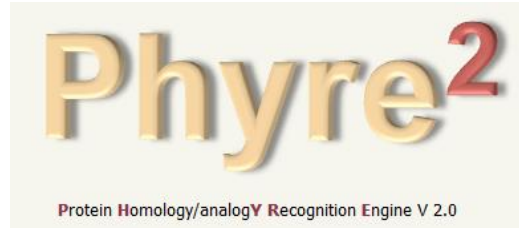
SWISS-MODEL is a fully automated protein structure homology-modelling server, accessible via the ExPASy web server, or from the program DeepView (Swiss Pdb-Viewer). The purpose of this server is to make Protein Modelling accessible to all biochemists and molecular biologists worldwide.

[Start Modelling](#)

- Web page: <http://swissmodel.expasy.org/>
- Fastest option, can take less than 5 minutes
- Final model typically based on a single template (users can upload their own)

29

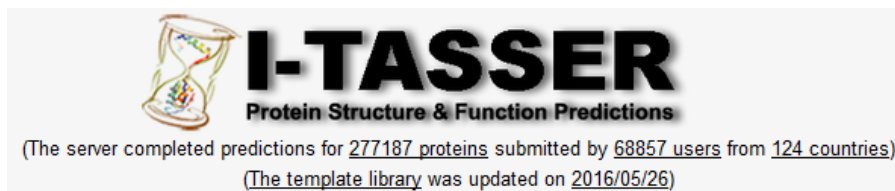
## Homology Modeling Servers: Phyre<sup>2</sup>



- Web page: <http://www.sbg.bio.ic.ac.uk/phyre2/>
- Trade off: can take 1-2 hours depending on server demand, but better structures
- Uses multiple templates, users can exclude files

30

## Homology Modeling Servers: I-TASSER



- Web page: <http://zhanglab.ccmb.med.umich.edu/I-TASSER/>
- Slowest option by far; can take a day or more
- Uses multiple templates and performs sophisticated refinement

31

## Homology Modeling Example

- Sequence for Pin1 protein:

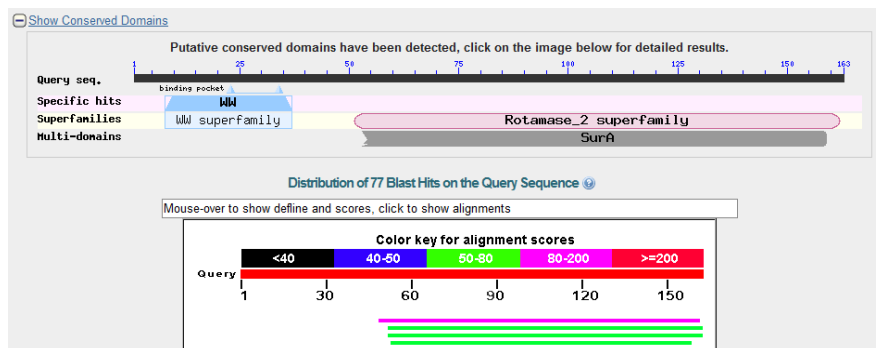
```
MADEEKLPPG WEKRMSRSSG RVYYFNHITN ASQWERPSGN SSSGGKNGQG
EPARVRCSHL LVKHSQSRRP SSWRQEKITR TKEEALELIN GYIQKIKSGE
EDFESLASQF SDCSSAKARG DLGAFSRGQM QKPFEDASFA LRTGEMSGPV
FTDSGIHIIL RTE
```

- Use BLAST to identify a homologous cis-trans prolyl isomerase in *Methanocorpusculum labreanum*

32

## Homology Modeling Example

- Initial BLASTp result:

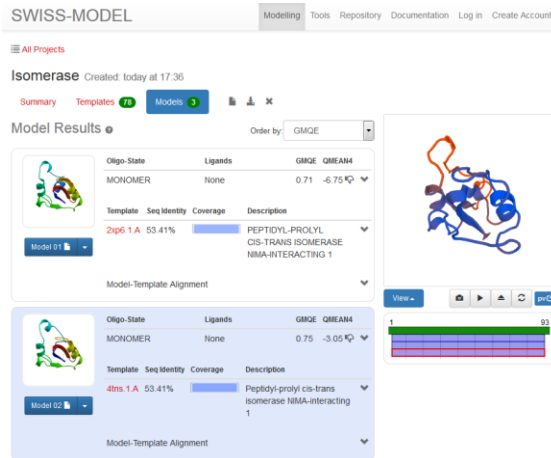


- Sequence (only second domain found):

```
MVRVKASHIL VKTEAQAKEI MQKISAGDDF AKLAKMYSQC PSGNAGGDLG
YFGKGQMVKP FEDACFKAKA GDVVGPKVTQ FGWHIIKVTD IKN
```

33

# Result: SWISS-MODEL



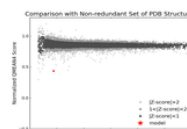
- We'll do this model in class

34

# Result: SWISS-MODEL

Model #01	File	Built with	Oligo-State	Ligands	GMQE	QMEAN4
	PDB	ProMod Version 3.70.	MONOMER	None	0.71	-6.75

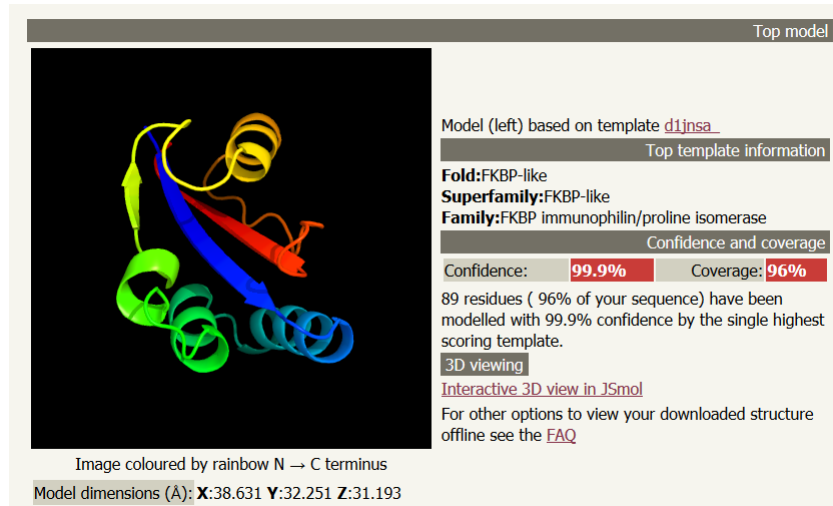
QMEAN4	-6.75
CP	-2.41
All Atom	-2.34
Solvation	-7.58
Torsion	-2.76



Template	Seq Identity	Oligo-state	Found by	Method	Resolution	Seq Similarity	Range	Coverage	Description
2xp6.1.A	53.41	monomer	BLAST	X-ray	1.90Å	0.45	3 - 90	0.95	PEPTIDYL-PROLYL CIS-TRANS ISOMERASE NIMA-INTERACTING 1
Ligand	Added to Model			Description					
12P	X - Binding site not conserved.			DODECAETHYLENE GLYCOL					
4G2	X - Binding site not conserved.			2-(3-CHLORO-PHENYL)-5-METHYL-1H-IMIDAZOLE-4-CARBOXYLIC ACID					

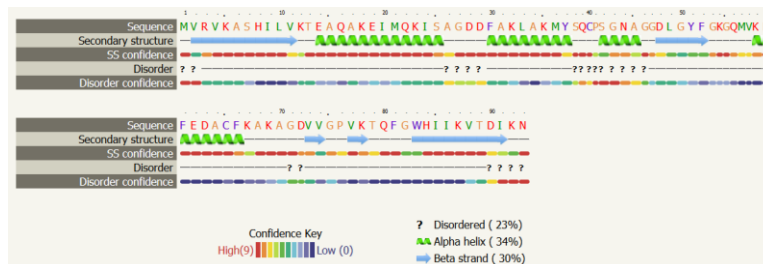
35

## Result: Phyre<sup>2</sup>



36

## Result: Phyre<sup>2</sup>



- Download entire result, which is a duplicate of the website, can be viewed here:  
<http://folding.chemistry.msstate.edu/files/bootcamp/phyre2/summary.html>
- Final result is called `final.casp.pdb`

37

# Result: I-TASSER

### Predicted Secondary Structure

Sequence MVRVKASHILVKTQAEIKIMKISAGDDFAKLAMYSQCPSGNAGDGLYFGRGQMVFPEFCADAKAKAGDVVGFKVTFGWHIIKVTDIKN  
Prediction CCSSSSSSSSSCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCS CCCCCSSSSSSSSSSCC  
Conf. Score 96799889988999999999999987998999998689652448645533799735699999838999978777698379999967659

Helix; S Strand; Coil

### Predicted Solvent Accessibility

	20	40	60	80
<b>Sequence</b>	MVVRKASHILVKTAEAKREIMQKISAGDGLAKLAFMYSCQPSGCPNSAGGDLGYPGKGQVVPFPEDACFKARAGDVVGPVPTKPGWHIIKVITDKN			
<b>Prediction</b>	V64304311116357405502630673640351056317344376323233045662243025003716645336234163100003046458			

- Results available at:  
<http://folding.chemistry.msstate.edu/files/bootcamp/itasser/>
- Final result is called `final.casp.pdb`

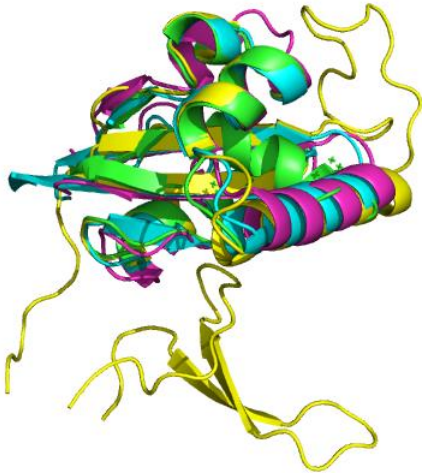
38

## Comparison of Results

- **Download the following PDBs from the Boot Camp Website:**
  - 1pin.pdb      – Original Pin1 Structure
  - swiss.pdb      – SWISS-MODEL Result
  - phyre2.pdb    – Phyre<sup>2</sup> Result
  - itasser.pdb    – I-TASSER Result
- PyMOL can help us here using the “align” command

39

## Comparison of Results



- Colors:
  - Original Pin1
  - SWISS-MODEL
  - Phyre<sup>2</sup>
  - I-TASSER
- **Important:** How much side chain accuracy do I need?

40

## Other Resources:

- EMBL-EBI (European Bioinformatics Institute) - <http://www.ebi.ac.uk/>
- DDBJ (DNA Data Bank of Japan) - <http://www.ddbj.nig.ac.jp/>
- NCBI's Sequence Read Archive (SRA) - <http://www.ncbi.nlm.nih.gov/sra>
- UCSC Genome Browser: <http://genome.ucsc.edu/>
- IGBB's Useful Links Page - <http://www.igbb.msstate.edu/links.php>

Many, many more available online, just search.

## Summary

- Sequence alignment is an important tool for searching and understanding how proteins are related
- BLAST can be used to search for similar sequences in large protein/DNA databases (and also works in tools like the PDB)
- Homology modeling can be helpful way to understand structures of unknown proteins

42