

Basic Bioinformatics, Sequence Alignment, and Homology

Biochemistry Boot Camp 2023

Session #11

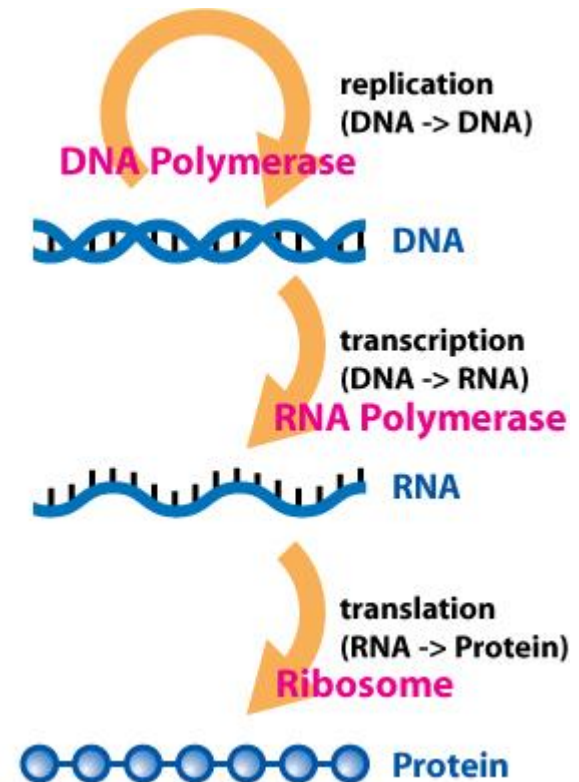
Nick Fitzkee

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



The Central Dogma of
Molecular Biology

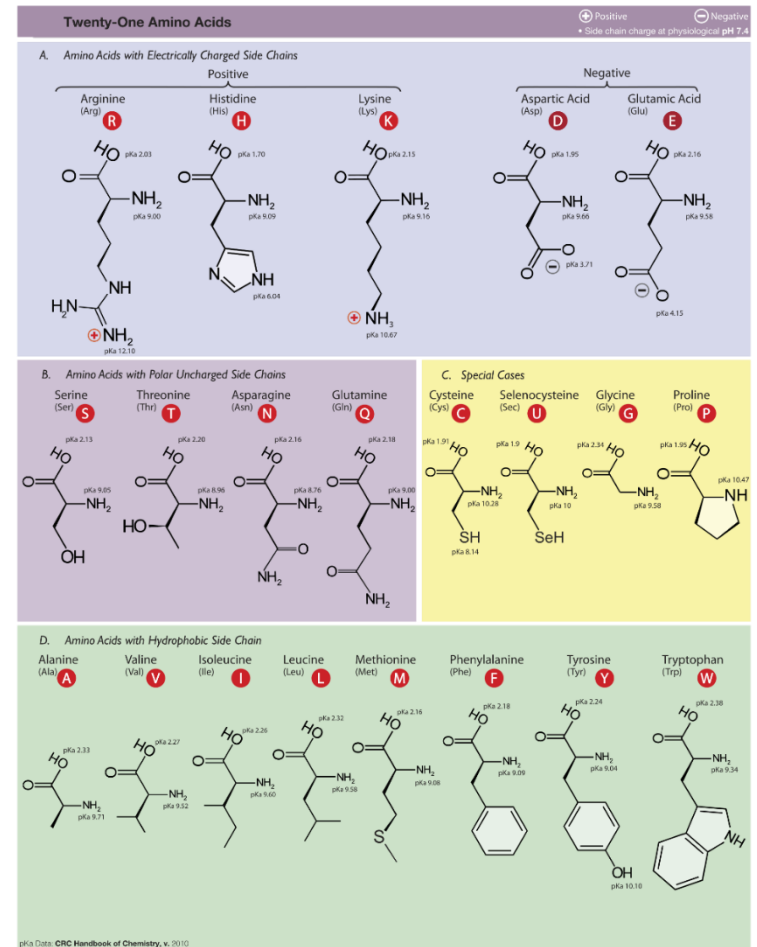
Primary Structure (Sequence)

- **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'



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



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.**

Storing Sequences

- FASTA format
- Can represent nucleotide sequences or peptide sequences using single letter codes
- FASTQ format
- Represents nucleotide sequences and their corresponding quality scores

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]  
LCLYTHIGRNIYYGSYLYSETWNTGIMLLITMATAFMGYVLPWQMSFWGATVITNLFSAIPYIGTNLV  
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYTIKDFLG  
LLILLLLLLLLLLALLSPDMLGDPDNHMPADPLNTPHLIKPEWYFLFAYAILRSVPNKLGGVLALFLSIVIL  
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFPLIAGX  
IENY
```

```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT  
+  
! '* (((**+))%%%+) (%%%) .1***-+*' ) **55CCF>>>>>>>CCCCCC65
```

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])

Sequence Alignment

Sequence 1: GARVEY

Sequence 2: AVERY

Global Alignment:

GARVE-Y

-A-VERY

Global Sequence Alignment

- 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
(download `sequences.txt` global sequence #1 and #2)

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?

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 MADQLTEEQIAEFKEAFSLFDKDGDTITTKELGTVMRSLGQNPTEAELQ      50
          |||
Nematode   1 MADQLTEEQIAEFKEAFSLFDKDGDTITTKELGTVMRSLGQNPTEAELQ      50

Human     51 DMINEVDADGNGTIDFPEFLTMMARKMKDSDSEEEIREAFRVFDKDGNGY     100
          |||
Nematode  51 DMINEVDADGNGTIDFPEFLTMMARKMKDSDSEEEIREAFRVFDKDGNGF     100

Human    101 ISAAELRHVMTNLGEKLTDEEVDEMIREADIDGDGQVNYEEFVQMMIAK     149
          |||
Nematode 101 ISAAELRHVMTNLGEKLTDEEVDEMIREADIDGDGQVNYEEFVIMMITK     149
```

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

- Pretty darn similar!

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,
<https://www.megasoftware.net/>

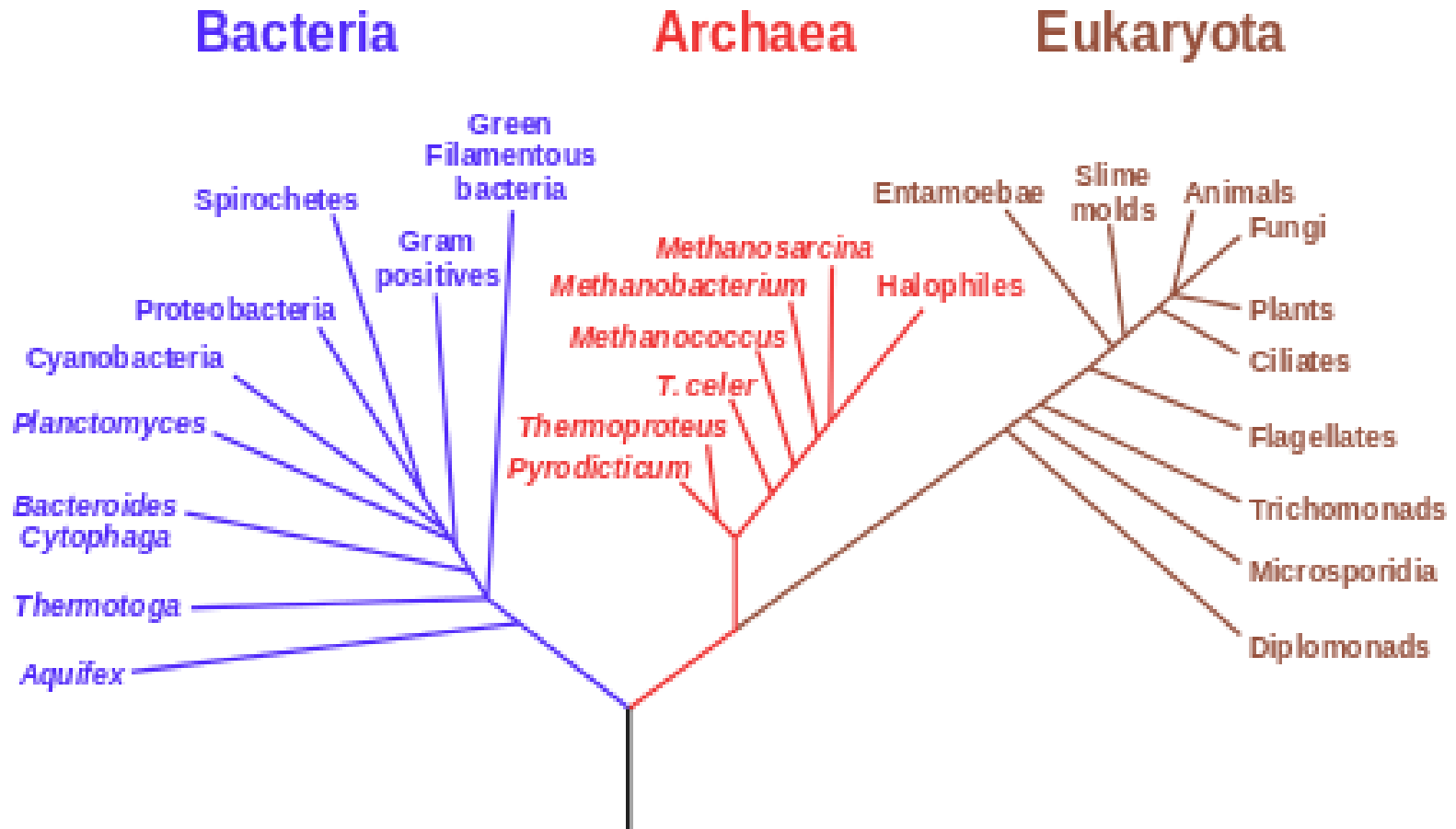


MSA Example

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--------------|--|-------------|-----------|----------|---------------|-----------|-------------|------------|-----------|------------|-----------|---------|---------|-------------|--------|---------|-------|-------|---------|---------|-------|------|------|------|------|------|-------|------|------|------|------|----|
| | | | * | . | : | . | * | : | : | : | . | | | | | | | | | | | | | | | | | | | | | |
| Q5E940_BOVIN | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_HUMAN | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_MOUSE | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_RAT | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_CHICK | ---- | MPREDRATWKS | NYFMKIIQL | LDDYPKCF | VVADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_RANSY | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQQIRMS | LRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | SALE | 76 | | | | | | | | | | | | | | | | | | | |
| Q7ZUG3_BRARE | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQTI | IRLSLRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0 ICTPU | ---- | MPREDRATWKS | NYFLKIIQL | LDDYPKCF | IVGADNVGS | KOMQTI | IRLSLRGK-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PALE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_DROME | ---- | MVRENKA | AWKAQYFIK | VVLFDEF | PKCFIVGADNVGS | KOMQNI | IRTSIRGL-AV | VL | MGKNTMMR | KAIRGHLENN | -- | PQLE | 76 | | | | | | | | | | | | | | | | | | | |
| RLA0_DICDI | ---- | MSGAG-SKR | KKLFIEKAT | KLFTTYDK | MIVAEAD | FVGS | SQLOKIRKS | IRGI-GAV | LMGKKT | MIRK | VIRDLADSK | -- | PELD | 75 | | | | | | | | | | | | | | | | | | |
| Q54LP0_DICDI | ---- | MSGAG-SKR | KNVFIEKAT | KLFTTYDK | MIVAEAD | FVGS | SQLOKIRKS | IRGI-GAV | LMGKKT | MIRK | VIRDLADSK | -- | PELD | 75 | | | | | | | | | | | | | | | | | | |
| RLA0_PLAF8 | ---- | MAKLSK | QKKQMYIEK | LSSLIQQ | YSKILLI | VHVDNVGS | NQMASV | RKSLRGK-AT | ILMGKNT | RIRTA | LKKNLQAV | -- | PQIE | 76 | | | | | | | | | | | | | | | | | | |
| RLA0_SULAC | ---- | MIGLAV | TTTKKIAK | WKVDEV | AELTEK | LKTKHT | IIIANIE | GFPADK | LHEIRK | KL | LRGK-AD | IKVT | KNL | FNIALKNAG | ---- | YDTK | 79 | | | | | | | | | | | | | | | |
| RLA0_SULTO | ---- | MRIMAV | ITQERKIA | KKIEV | KEVKE | QKLREY | HTIIANIE | GFPADK | LHDIRK | KK | MRGM-AE | IKVT | KNL | FLGIAAKNAG | ---- | LDVS | 80 | | | | | | | | | | | | | | | |
| RLA0_SULSO | ---- | MKRLAL | ALQKRVAS | WKLEEV | KELETEL | IKNSNT | LIGNLE | GFPADK | LHEIRK | KL | LRGK-AT | IKVT | KNL | FLFKIAAKNAG | ---- | IDIE | 80 | | | | | | | | | | | | | | | |
| RLA0_AERPE | MSVVS | LVGQMYK | REKPIPE | WKTLM | LELEEL | FSK | HRVVL | FADLT | GTPT | FFVQ | RVRKK | LWKK-Y | P | MVAK | KRIIL | RAMKA | AAGLE | ---- | LDDN | 86 | | | | | | | | | | | | |
| RLA0_PYRAE | MMLAI | GKRRY | VRTRQY | PARKV | KIVSE | ATELL | QKYP | YVFL | FDLH | GLSS | RILHE | YRY | RLRRY-G | V | IKI | IKP | TLFK | IAFTK | VYGG | ---- | IPAE | 85 | | | | | | | | | | |
| RLA0_METAC | ---- | MAERH | HTHEHIP | QWKKDE | IENIKEL | IQSHK | VFGM | VGIEG | ILATK | MOKIR | RD | LKDV-AV | L | KVSR | NL | TE | RALN | QLG | ---- | ETIP | 78 | | | | | | | | | | | |
| RLA0_METMA | ---- | MAERH | HTHEHIP | QWKKDE | IENIKEL | IQSHK | VFGM | VRIEG | ILATK | IOKIR | RD | LKDV-AV | L | KVSR | NL | TE | RALN | QLG | ---- | ESIP | 78 | | | | | | | | | | | |
| RLA0_ARCFU | ---- | MAAVR | GS--P | PEYK | VRAVEE | EIKRM | ISSK | PVVAI | V | FRNVP | AGOM | QKIR | REFR | GK-AE | IKV | VKN | LLE | RALD | ALG | ---- | GDYL | 75 | | | | | | | | | | |
| RLA0_METKA | MAVKAK | GQPP | SGYE | PKVAE | WKRRE | VEKEL | ELMDE | YENV | GLVD | LEGIP | APQLOE | IRAKL | RE | DTI | IRMS | RNTLM | RIALE | EK | LDER | -- | PELE | 88 | | | | | | | | | | |
| RLA0_METTH | ---- | MAHVAE | WKKKEV | QELHDL | IKGYE | VVGIAN | LADIPAR | LOKMR | QTLRDS-AL | IRMS | SKTLL | ISL | ALEK | AGREL | -- | ENVD | 74 | | | | | | | | | | | | | | | |
| RLA0_METTL | ---- | MITAE | SEHKI | APWK | IEEVN | KLKEL | LKNGQ | IVAL | VDMME | V | PARQLOE | IRDK | IR-G | T | MTL | KMSR | NL | LIE | RAI | KEVA | EETGN | PEFA | 82 | | | | | | | | | |
| RLA0_METVA | ---- | MIDAK | SEHKI | APWK | IEEVN | KLKEL | LKSNV | IAL | IDMME | V | PAVQLOE | IRDK | IR-D | Q | MTL | KMSR | NL | L | KRAVEE | VAE | EETGN | PEFA | 82 | | | | | | | | | |
| RLA0_METJA | ---- | METKV | KAHVAP | WK | IEEVK | TLKGL | IKSKP | VVAI | V | DMM | VPAPQLOE | IRDK | IR-D | K | VKL | RMSR | NL | L | IRAL | KEAAE | LNN | PKLA | 81 | | | | | | | | | |
| RLA0_PYRAB | ---- | MAHVAE | WKKKEV | EELAN | LKSYP | VIAL | VDVSS | M | PAYPLS | QMRRL | IRE | NG | LLR | VSR | NL | LIE | LAIK | KAQEL | LG | KPELE | 77 | | | | | | | | | | | |
| RLA0_PYRHO | ---- | MAHVAE | WKKKEV | EELAKL | LKSYP | VIAL | VDVSS | M | PAYPLS | QMRRL | IRE | NG | LLR | VSR | NL | LIE | LAIK | KA | AKEL | LG | KPELE | 77 | | | | | | | | | | |
| RLA0_PYRFU | ---- | MAHVAE | WKKKEV | EELAN | LKSYP | VVAL | VDVSS | M | PAYPLS | QMRRL | IRE | NG | LLR | VSR | NL | LIE | LAIK | KA | QEL | LG | KPELE | 77 | | | | | | | | | | |
| RLA0_PYRKO | ---- | MAHVAE | WKKKEV | EELANI | LKSYP | VIAL | VDVAG | V | PAYPLSK | MRDKLR-G | K | ALL | VSR | NL | LIE | LAIK | RA | QEL | LG | KPELE | 76 | | | | | | | | | | | |
| RLA0_HALMA | ---- | MSAES | ERKTET | IP | EWKQ | EEVD | AI | EMIES | Y | SVVNI | AGIP | SRQLO | DM | RRDL | HGT-AE | L | VSR | NL | L | LE | RA | DDVD | ---- | DGLE | 79 | | | | | | | |
| RLA0_HALVO | ---- | MSESE | VRQTE | VI | PQWK | RE | VD | LVDF | IES | YES | V | GV | VAG | IP | SRQLO | S | MR | REL | LHGS-AA | V | RMSR | NL | V | N | RA | L | DEVN | ---- | DGFE | 79 | | |
| RLA0_HALSA | ---- | MSAEE | QRTTE | EV | PEWK | RQ | EVAEL | VDL | LET | Y | DSV | V | V | VTG | IP | SKO | L | DM | RRGL | LHGO-AA | L | RMSR | NL | L | V | RA | LEEAG | ---- | DGLD | 79 | | |
| RLA0_THEAC | ---- | MKEV | SQKKEL | VNEIT | ORIKAS | RSVAI | VD | TAG | IRT | ROI | Q | DIR | G | KNR | GK-IN | L | KV | IK | TL | L | L | F | KA | LEN | LG | ---- | EKLS | 72 | | | | |
| RLA0_THEVO | ---- | MRKIN | PKKKE | IV | SELA | QD | IT | SKAV | AI | VD | IK | G | VRT | Q | MOD | IRAK | NR | DK-V | K | V | V | K | T | L | L | F | KA | LDS | IND | ---- | EKLT | 72 |
| RLA0_PICTO | ---- | MTEPA | QWKID | FVKN | LENE | INSR | KVA | AI | VS | IK | GLRN | NE | FQ | KIR | NS | IRDK-AR | IKV | SR | AR | LL | RLA | IEN | TGK | ---- | NNIV | 72 | | | | | | |
| ruler | 1.....10.....20.....30.....40.....50.....60.....70.....80.....90 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

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

MSA-Derived Phylogenetic Tree



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.

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/>

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

BLAST Main Page

BLAST: Basic Local Alignment Search Tool

https://blast.ncbi.nlm.nih.gov/Blast.cgi

NIH National Library of Medicine
National Center for Biotechnology Information

Log in

BLAST® Home Recent Results Saved Strategies Help

Take the BLAST survey today [Start survey](#)

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

NEWS
BLAST+ 2.14.0 is here!
BLASTP, BLASTX, and TBLASTN are faster than before.
Fri, 28 Apr 2023 [More BLAST news...](#)

Web BLAST

Nucleotide BLAST
nucleotide ▶ nucleotide

blastx
translated nucleotide ▶ protein

tblastn
protein ▶ translated nucleotide

Protein BLAST
protein ▶ protein

BLAST Genomes

Enter organism common name, scientific name, or tax id [Search](#)

Human Mouse Rat Microbes

Nucleotide BLAST: Search nucle X +

https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blast 80%

BLAST® » blastn suite Home Recent Results Saved Strategies Help

Take the BLAST survey today Start survey

Standard Nucleotide BLAST

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. more... Reset page Bookmark

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) Clear Query subrange

From To

Or, upload file Browse... No file selected.

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database Standard databases (nr etc.): rRNA/ITS databases Genomic + transcript databases Betacoronavirus

Experimental databases [Try experimental taxonomic nt databases](#) [Download](#)

[For more info see What are taxonomic nt databases?](#)

Nucleotide collection (nr/nt)

Organism Optional

Enter organism name or id--completions will be suggested exclude [Add organism](#)

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

Exclude Optional

Models (XM/XP) Uncultured/environmental sample sequences

Limit to Optional

Sequences from type material

Entrez Query Optional

Enter an Entrez query to limit search [YouTube](#) [Create custom database](#)

Program Selection

Optimize for Highly similar sequences (megablast) More dissimilar sequences (discontiguous megablast) Somewhat similar sequences (blastn)

Choose a BLAST algorithm

BLAST Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences) Show results in a new window

+ Algorithm parameters

Sequence Input

Databases to Search Against

Program Selection

Click to Run!

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&BLAST_SF

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

Standard Protein BLAST

blastn blastp

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastx&BLAST_PROGRAMS=blastx&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&BLAST_SPE

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

Translated BLAST: blastx

blastn blastp blastx

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=tblastn&BLAST_PROGRAMS=tblastn&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&BLAST_SF

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

Translated BLAST: tblastn

blastn blastp blastx tblastn tblastx

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=tblastx&BLAST_PROGRAMS=tblastx&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&BLAST_SF

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

Translated BLAST: tblastx

blastn blastp blastx tblastn tblastx

TBLASTX search translated nucleotide databases using a translated nucleotide query. [more...](#) [Reset page](#) [Bookmark](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#) [Query subrange](#)

From

To

Or, upload file No file chosen

Genetic code

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database

Organism Exclude

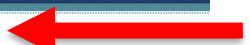
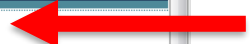
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Models (XM/XP) Uncultured/environmental sample sequences

Exclude

Entrez Query

Same Page Organization



BLAST Example

- What gene is this?

```
>unknown_sequence_1
```

```
TGATGTCAAGACCCTCTATGAGACTGAAGTCTTTTCTACCGACTTCTCCAACATTTCTGCAGCCAAGCAG  
GAGATTAACAGTCATGTGGAGATGCAAACCAAAGGGAAAGTTGTGGGTCTAATTC AAGACCTCAAGCCAA  
ACACCATCATGGTCTTAGTGA ACTATATTC ACTTTAAAGCCCAGTGGGCAAATCCTTTTGATCCATCCAA  
GACAGAAGACAGTTCCAGCTTCTTAATAGACAAGACCACCACTGTTCAAGTGCCCATGATGCACCAGATG  
GAACAATACTATCACCTAGTGGATATGGAATTGAACTGCACAGTTCTGCAAATGGACTACAGCAAGAATG  
CTCTGGCACTCTTTGTTCTTCCCAAGGAGGGACAGATGGAGTCAGTGGAAAGCTGCCATGTCATCTAAAAC  
ACTGAAGAAGTGGAACCGCTTACTACAGAAGGGATGGGTGACTTGTTTGTTCCAAAGTTTTCCATTTCT  
GCCACATATGACCTTGGAGCCACACTTTTGAAGATGGGCATTCAGCATGCCTATTCTGAAAATGCTGATT  
TTTCTGGACTCACAGAGGACAATGGTCTGAAACTTTCCAATGCTGCCCATAAGGCTGTGCTGCACATTGG  
TGAAAAGGGAACTGAAGCTGCAGCTGTCCCTGAAGTTGAACTTTCGGATCAGCCTGAAAACACTTTCCTA  
CACCTATTATCCAAATTGATAGATCTTTCATGTTGTTGATTTTGGAGAGAAGCACAAGGAGTATTCTCT  
TTCTAGGGAAAGTTGTGAACCCAACGGAAGCGTAGTTGGGAAAAGGCCATTGGCTAATTGCACGTGTGT  
ATTGCAATGGGAAATAAATAAATAATATAGCCTGGTGTGATTGATGTGAGCTTGGACTTGCATTCCTTA  
TGATGGGATGAAGATTGAACCCTGGCTGAACTTTGTTGGCTGTGGAAGAGGCCAATCCTATGGCAGAGCA  
TTCAGAATGTCAATGAGTAATTCATTATTATCCAAAGCATAGGAAGGCTCTATGTTTGTATATTTCTCTT  
TGTCAGAATACCCCTCAACTCATTTGCTCTAATAAATTTGACTGGGTGAAAAATTAAAA
```

BLAST Results

NCBI Blast:Nucleotide Sequence X

https://blast.ncbi.nlm.nih.gov/Blast.cgi

National Library of Medicine
National Center for Biotechnology Information

BLAST® » blastn suite » results for RID-7YY6WY0C016

Take the BLAST survey today [Start survey](#)

[← Edit Search](#)
[Save Search](#)
[Search Summary ▾](#)
[How to read this report?](#)
[BLAST Help Videos](#)
[Back to Traditional Results Page](#)

Job Title: Nucleotide Sequence
 RID: [7YY6WY0C016](#) Search expires on 06-08 02:20 am [Download All ▾](#)
 Program: BLASTN [Citation ▾](#)
 Database: nt [See details ▾](#)
 Query ID: lcl|Query_27249
 Description: None
 Molecule type: dna
 Query Length: 1110
 Other reports: [Distance tree of results](#) [MSA viewer](#) [?](#)

Filter Results

Organism only top 20 will appear exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

Percent Identity: to E value: to Query Coverage: to

[Filter](#) [Reset](#)

[Descriptions](#)
[Graphic Summary](#)
[Alignments](#)
[Taxonomy](#)

Sequences producing significant alignments Download ▾ Select columns ▾ Show 100 ▾ [?](#)

select all 100 sequences selected

[GenBank](#)
[Graphics](#)
[Distance tree of results](#)
[MSA Viewer](#)

| Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|--|---------------------------------|-----------|-------------|-------------|---------|------------|----------|--------------------------------|
| <input checked="" type="checkbox"/> Homo sapiens serpin family A member 7 (SERPINA7), mRNA | Homo sapiens | 2043 | 2043 | 99% | 0.0 | 100.00% | 2360 | NM_000354.6 |
| <input checked="" type="checkbox"/> Human thyroxine-binding globulin mRNA, complete cds | Homo sapiens | 2019 | 2019 | 99% | 0.0 | 99.82% | 1872 | M14091.1 |
| <input checked="" type="checkbox"/> Pan troglodytes serpin family A member 7 (SERPINA7), mRNA | Pan troglodytes | 2012 | 2012 | 99% | 0.0 | 99.64% | 1589 | NM_001009109.1 |
| <input checked="" type="checkbox"/> PREDICTED: Pongo pygmaeus serpin family A member 7 (SERPINA7), mRNA | Pongo pygmaeus | 1977 | 1977 | 99% | 0.0 | 98.92% | 2371 | XM_054471959.1 |
| <input checked="" type="checkbox"/> PREDICTED: Pongo abelii serpin family A member 7 (SERPINA7), mRNA | Pongo abelii | 1965 | 1965 | 99% | 0.0 | 98.73% | 2371 | XM_002831954.5 |
| <input checked="" type="checkbox"/> PREDICTED: Symphalangus syndactylus serpin family A member 7 (SERPINA7), mRNA | Symphalangus... | 1943 | 1943 | 99% | 0.0 | 98.37% | 2365 | XM_055269099.1 |
| <input checked="" type="checkbox"/> PREDICTED: Papio anubis serpin family A member 7 (SERPINA7), transcript variant X2, mRNA | Papio anubis | 1908 | 1908 | 99% | 0.0 | 97.83% | 1645 | XM_003918082.5 |
| <input checked="" type="checkbox"/> PREDICTED: Macaca mulatta serpin family A member 7 (SERPINA7), transcript variant X1, mRNA | Macaca mulatta | 1897 | 1897 | 99% | 0.0 | 97.65% | 1656 | XM_001088790.4 |

[Feedback](#)

BLAST Results – Graphical Summary and Alignments

NCBI Blast:Nucleotide Sequenc... X

https://blast.ncbi.nlm.nih.gov/ 80%

Import bookmarks... Most Visited 16 Tips to Make Your ... Generate a CSR (certifi... Other Bookmarks

Descriptions **Graphic Summary** Alignments Taxonomy

hover to see the title click to show alignments

Alignment Scores < 40 40 - 50 50 - 80 80 - 200 >= 200

100 sequences selected

Distribution of the top 113 Blast Hits on 100 subject sequences

Query

1 200 400 600 800 1000

Feedback

NCBI Blast:Nucleotide Sequenc... X

https://blast.ncbi.nlm.nih.gov/Blast.cgi 80%

Import bookmarks... Most Visited 16 Tips to Make Your ... Generate a CSR (certifi... Other Bookmarks

Graphic Summary **Alignments** Taxonomy

Pairwise CDS feature Restore defaults Download

selected

GenBank Graphics

serpins serpin family A member 7 (SERPINA7), mRNA

ID: [NM_000354.6](#) Length: 2360 Number of Matches: 1

491 to 1596 GenBank Graphics

| | Expect | Identities | Gaps | Strand |
|--------|--------|-----------------|------------|-----------|
| (1106) | 0.0 | 1106/1106(100%) | 0/1106(0%) | Plus/Plus |

```
Query 61 AGCCAAGCAGGAGATTAAACAGTCATGTGGAGATGCAAACCAAAGGGAAAAGTTGTGGGCTCT
Sbjct 551 AGCCAAGCAGGAGATTAAACAGTCATGTGGAGATGCAAACCAAAGGGAAAAGTTGTGGGCTCT
Query 121 AATTCAAGACCTCAAGCCAAACACCATCATGGTCTTAGTGAACATATTCACCTTTAAAGC
Sbjct 611 AATTCAAGACCTCAAGCCAAACACCATCATGGTCTTAGTGAACATATTCACCTTTAAAGC
Query 181 CCAGTGGGCAAATCCTTTTGATCCATCCAAGACAGAAGACAGTCCAGCTTCTTAATAGA
Sbjct 671 CCAGTGGGCAAATCCTTTTGATCCATCCAAGACAGAAGACAGTCCAGCTTCTTAATAGA
Query 241 CAAGACCACCACTGTTCAAGTGCCCATGATGCACCAAGTGAACAATACTACTCACCTAGT
Sbjct 731 CAAGACCACCACTGTTCAAGTGCCCATGATGCACCAAGTGAACAATACTACTCACCTAGT
```

Feedback

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

Installing BLAST Locally

Executables and documentation available at:

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

Documentation:

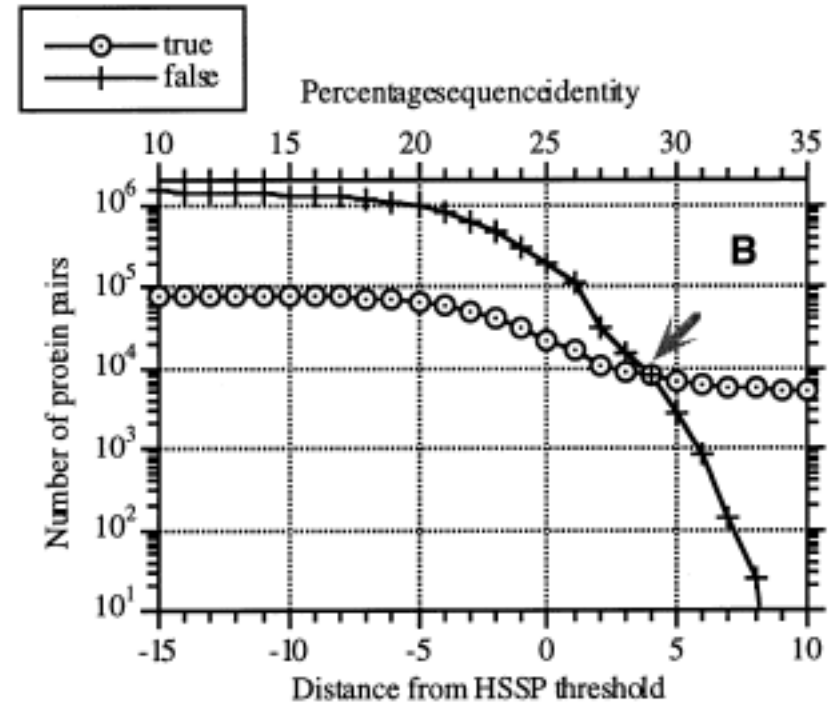
<https://www.ncbi.nlm.nih.gov/books/NBK1762/>

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>

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 Lauren Perskie-Porter, Phil Bryan and “fold switching”
- Can we predict structures?

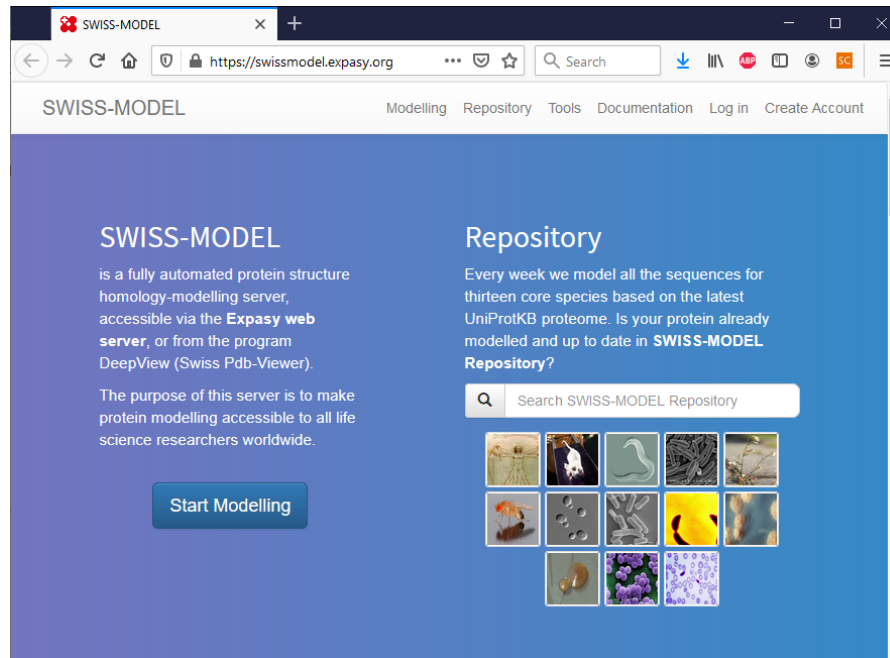


Below ~28% sequence identity, the number of structurally dissimilar aligned pairs explodes.

What is Homology Modeling?

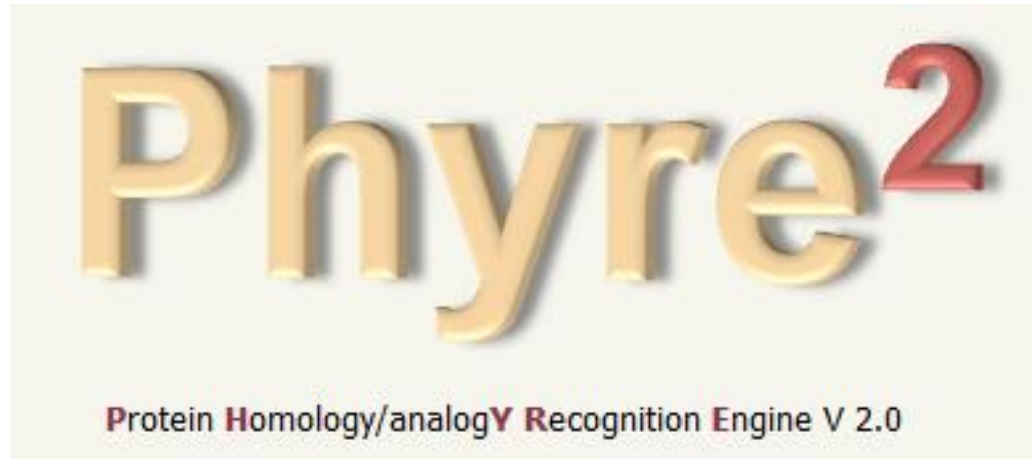
- **Consider:** Protein with known sequence, but unknown structure
- Use sequence alignment (protein BLAST) to identify similar sequences with known structures
 - These are termed “template structures”
- “Map” unknown sequence onto known backbone
 - Side chains may be more ill-defined: it’s a model!

Homology Modeling Servers: **SWISS-MODEL**



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

Homology Modeling Servers: **Phyre²**



- 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

Homology Modeling Servers: I-TASSER



I-TASSER Protein Structure & Function Predictions

(The server completed predictions for 739548 proteins submitted by 182114 users from 160 countries)
(The template library was updated on 2023/05/01)

- Web page: <https://zhanggroup.org/I-TASSER/>
- Slowest option by far; can take a day or more
- Uses multiple templates and performs sophisticated refinement

Homology Modeling Example

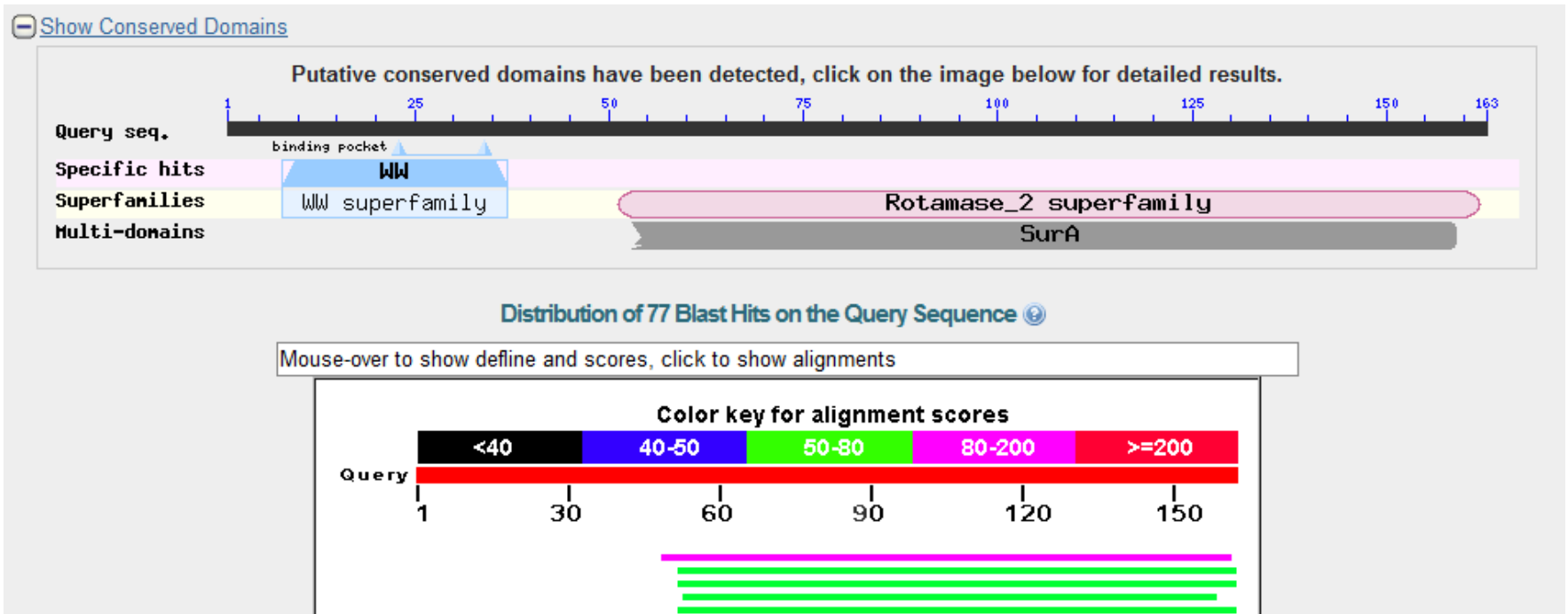
- Sequence for Pin1 protein:

```
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*

Homology Modeling Example

- Initial BLASTp result:



- Sequence (only second domain found):

MVRVKASHIL VKTEAQAKEI MQKISAGDDF AKLAKMYSQC PSGNAGGDLG
YFGKGQMVKP FEDACFKAKA GDVVGPKVTQ FGWHIIKVTD IKN

Result: SWISS-MODEL

“Searching for templates” lets you select which structure(s) are used to build your homology model.

If you aren't comfortable using AlphaFold structures, you can deselect them!

Alphafold have higher sequence similarity, but you're be building a model derived from a model

Template Results

Templates Quaternary Structure Sequence Similarity Alignment

More ▾

| Sort | Coverage | GMQE | QSQE | Identity | Method | Oligo State | Ligands |
|-------------------------------------|----------|------|------|----------|--------------|------------------------------------|---------|
| <input type="checkbox"/> | | | | | AlphaFold v2 | monomer ✓ | None |
| <input checked="" type="checkbox"/> | | | | | X-ray, 2.6Å | homo-dimer Δ 2 x 2NV ^{CS} | |
| <input type="checkbox"/> | | | | | X-ray, 1.5Å | monomer ✓ | None |
| <input type="checkbox"/> | | | | | X-ray, 2.6Å | homo-dimer Δ 2 x 2NV ^{CS} | |
| <input type="checkbox"/> | | | | | X-ray, 1.5Å | monomer ✓ | None |
| <input type="checkbox"/> | | | | | X-ray, 1.70Å | monomer ✓ | None |

Build Models 1

Clear Selection

6vj6.1.A

Result: SWISS-MODEL

The screenshot displays the SWISS-MODEL web interface for an 'Untitled Project'. The main content area shows 'Model Results' for a protein of 93 residues. A 'Structure Assessment' button is highlighted with a red box and an arrow pointing to it from the text on the left. The assessment includes a 'GMQE' score of 0.80 and a 'QMEANDisCo Global' score of 0.82 ± 0.09 . A 'QMEANDisCo Local' plot shows the local quality estimate for Chain A. Below this, 'QMEAN Z-Scores' are shown for OMEAN (0.66), C β (0.93), All Atom (0.37), solvation (-0.09), and torsion (0.45). The 'Template' section identifies the structure as 6vj6.1.A Peptidylprolyl isomerase (PrsA) from Bacillus cereus, with a sequence identity of 55.56%. The 'Model-Template Alignment' section shows the alignment of the model sequence (Model_01) with the template sequence (6vj6.1.A).

Model Results

Order by: GMQE

1 93

Model 01

Structure Assessment

Oligo-State: Monomer

GMQE: 0.80

QMEANDisCo Global: 0.82 ± 0.09

QMEANDisCo Local

Local Quality Estimate - Chain A

Predicted Local Similarity to target

Residue Number

QMEAN Z-Scores

OMEAN: 0.66

C β : 0.93

All Atom: 0.37

solvation: -0.09

torsion: 0.45

Template

6vj6.1.A Peptidylprolyl isomerase (PrsA)

2.55 Angstrom Resolution Crystal Structure of Peptidylprolyl Isomerase (PrsA) from Bacillus cereus

Seq Identity: 55.56%

Coverage

Model-Template Alignment

Model_01: MVRVVKASHLLVYKTEAQAKEIMQKISAGDDFAKLAK 35

6vj6.1.A: MVRVVKASHLLVYKTEAQAKEIMQKISAGDDFAKLAK 144

Click here to view Ramachandran plots, structure quality by residue, etc.

Click structure
to download
PDB file

Result: Phyre²

Top model

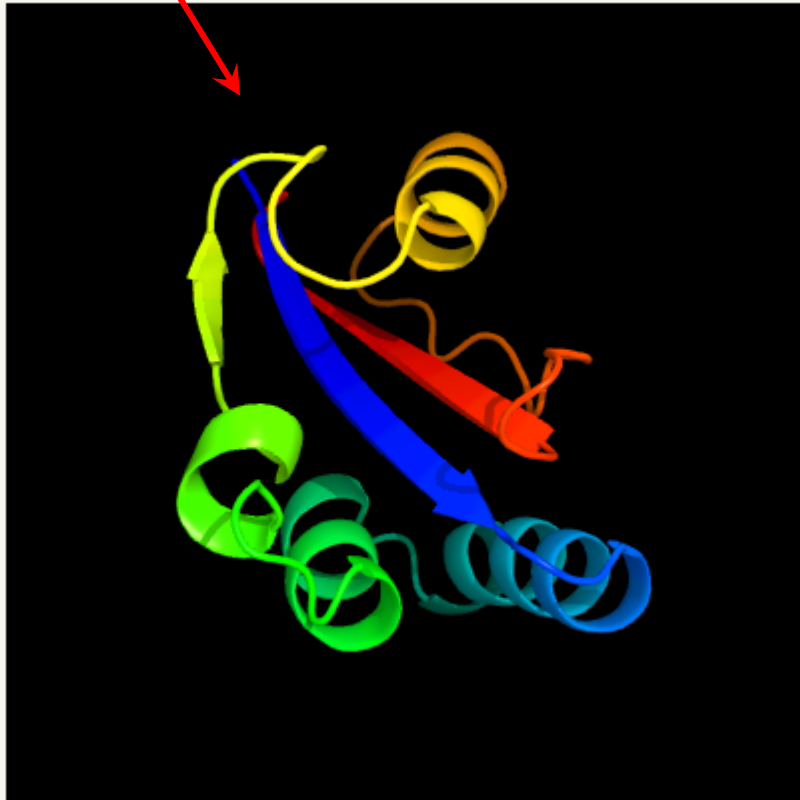


Image coloured by rainbow N → C terminus

Model dimensions (Å): **X**:38.631 **Y**:32.251 **Z**:31.193

Model (left) based on template [d1jnsa](#)

Top template information

Fold:FKBP-like

Superfamily:FKBP-like

Family:FKBP immunophilin/proline isomerase

Confidence and coverage

Confidence: **99.9%**

Coverage: **96%**

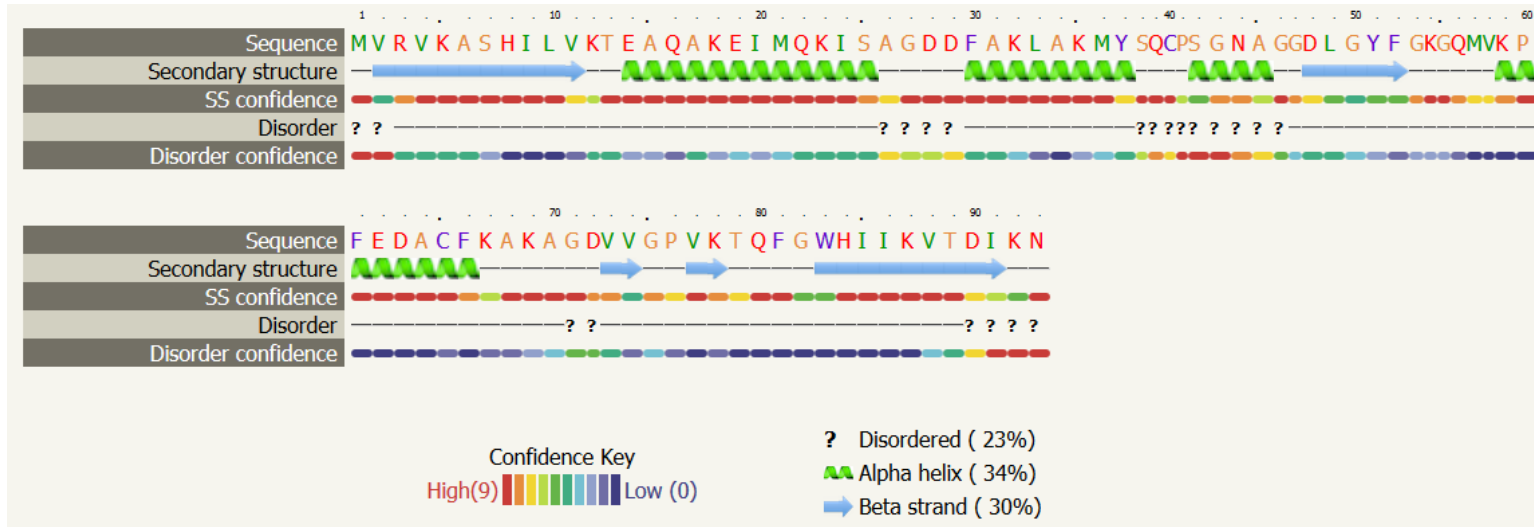
89 residues (96% of your sequence) have been modelled with 99.9% confidence by the single highest scoring template.

3D viewing

[Interactive 3D view in JSmol](#)

For other options to view your downloaded structure offline see the [FAQ](#)

Result: Phyre²



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

Result: I-TASSER

Predicted Secondary Structure

| | 20 | 40 | 60 | 80 |
|-------------|--|----|----|----|
| Sequence | MVRVKASHILVKTEAQAQAKEIMQKISAGDDFAKLAKMYSQCPSGNAGGDLGYFGKGMVKPFEDACFKAKAGDVVGPVKTFGWIIKVTDIKN | | | |
| Prediction | CCSSSSSSSSSCCHHHHHHHHHHHHCCCCHHHHHHHHHCCCCCCCCCCCCCCCCCCCCHHHHHHHHHCCCCCCCCSSCCCCSSSSSSSSSSSCC | | | |
| Conf. Score | 967998899988999999999999998879989999998688965244864553379973569999998389999788777698379999967659 | | | |
| | H:Helix; S:Strand; C:Coil | | | |

Predicted Solvent Accessibility

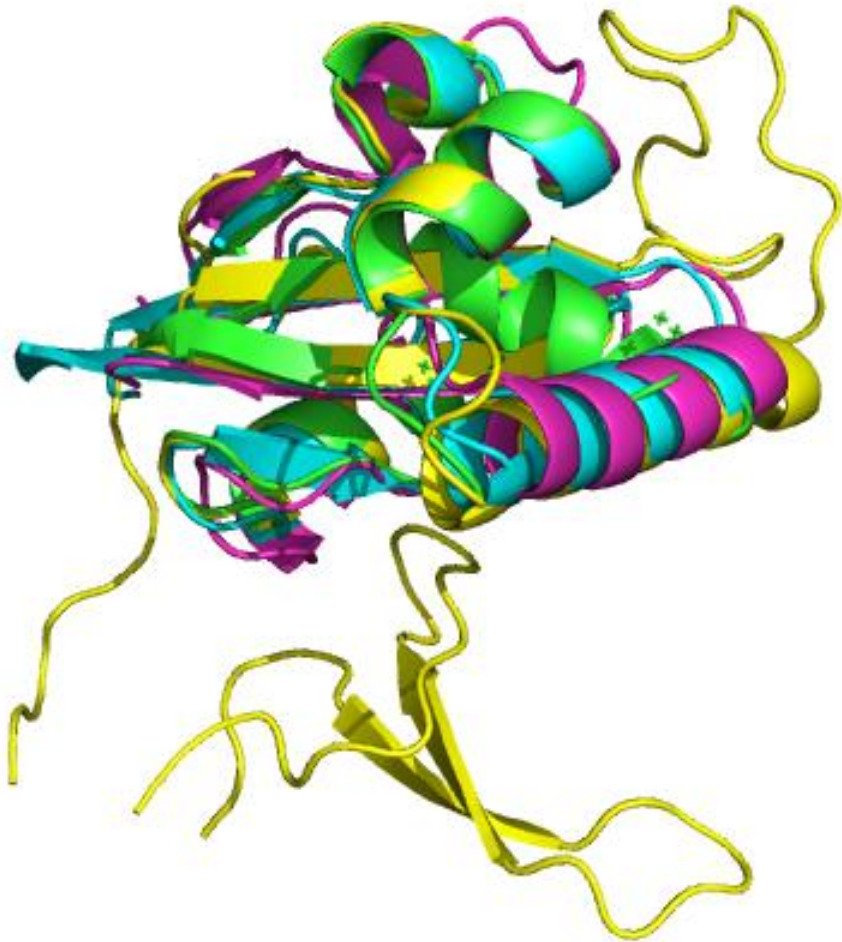
| | 20 | 40 | 60 | 80 |
|------------|---|----|----|----|
| Sequence | MVRVKASHILVKTEAQAQAKEIMQKISAGDDFAKLAKMYSQCPSGNAGGDLGYFGKGMVKPFEDACFKAKAGDVVGPVKTFGWIIKVTDIKN | | | |
| Prediction | 764340311116357405502630673640351056317344376323233045662243025003716645336234163100003046458 | | | |
| | Values range from 0 (buried residue) to 9 (highly exposed residue) | | | |

- Results available at:
<https://fitzkee.chemistry.msstate.edu/sites/default/files/bootcamp/itasser/>
- Final result is called `model1.pdb`

Comparison of Results

- **Download the following PDBs from the Boot Camp Website:**
 - 1pin.pdb – Original Pin1 Structure
 - swiss.pdb – SWISS-MODEL Result
 - phyre2.pdb – Phyre² Result
 - itasser.pdb – I-TASSER Result
- PyMOL can help us here using the “align” command (align.pse)

Comparison of Results



- Colors:
 - Original Pin1
 - SWISS-MODEL
 - Phyre²
 - I-TASSER
- **Important:** How much side chain accuracy do I need?

AlphaFold2: Neural Networks

- Google Deepmind Project: Exhaustively predict protein structure based on known structure patterns
- Not really homology modeling, not really “ab initio” or physics-based
- Extremely successful!

Article

Highly accurate protein structure prediction with AlphaFold


<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

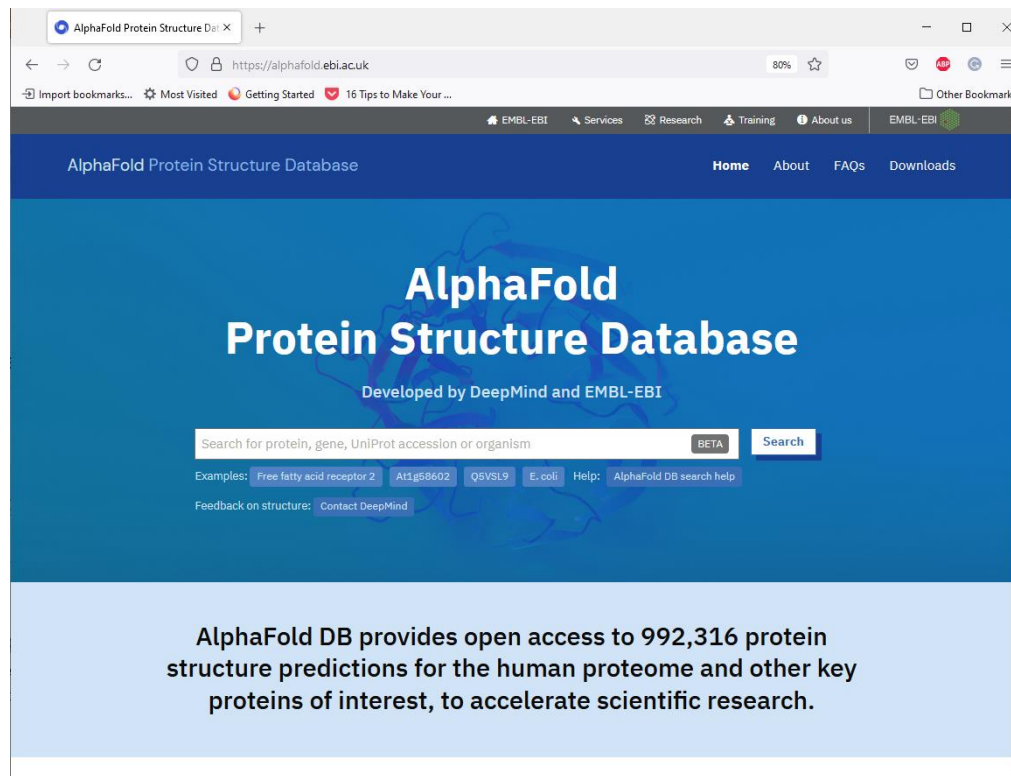
 Check for updates

John Jumper^{1,2,3,4}, Richard Evans^{1,4}, Alexander Pritzel^{1,4}, Tim Green^{1,4}, Michael Figurnov^{1,4}, Olaf Ronneberger^{1,4}, Kathryn Tunyasuvunakool^{1,4}, Russ Bates^{1,4}, Augustin Zidek^{1,4}, Anna Potapenko^{1,4}, Alex Bridgland^{1,4}, Clemens Meyer^{1,4}, Simon A. A. Kohl^{1,4}, Andrew J. Ballard^{1,4}, Andrew Cowie^{1,4}, Bernardino Romera-Paredes^{1,4}, Stanislav Nikolov^{1,4}, Rishub Jain^{1,4}, Jonas Adler¹, Trevor Back¹, Stig Petersen¹, David Reiman¹, Ellen Clancy¹, Michal Zielinski¹, Martin Steinegger^{2,3}, Michalina Pacholska¹, Tamas Berghammer¹, Sebastian Bodenstein¹, David Silver¹, Oriol Vinyals¹, Andrew W. Senior¹, Koray Kavukcuoglu¹, Pushmeet Kohli¹ & Demis Hassabis^{1,2,3,4}

Proteins are essential to life, and understanding their structure can facilitate a mechanistic understanding of their function. Through an enormous experimental effort^{1–4}, the structures of around 100,000 unique proteins have been determined⁵, but

AlphaFold2 Website

- Prediction Database: <https://alphafold.ebi.ac.uk/>



- Entry: P12104 (Human Intestinal Fatty Acid Binding Protein)

FABP Entry – P12104

- Many entries exist, but not so easy to run this yourself on a new structure
- For more information check out the DeepMind website
- <https://www.deepmind.com/research/highlighted-research/alphafold>

AlphaFold Protein Structure Database

https://alphafold.ebi.ac.uk/entry/P12104

Examples: Free fatty acid receptor 2 | A1g58602 | Q5VSL9 | E. coli | Help: AlphaFold DB search help

Fatty acid-binding protein, intestinal

AlphaFold structure prediction

Download [PDB file](#) [mmCIF file](#) [Predicted aligned error](#)

NEW Feedback on structure [Looks great](#) [Could be improved](#)

Information

| | |
|-------------------------|--|
| Protein | Fatty acid-binding protein, intestinal |
| Gene | FABP2 |
| Source organism | Homo sapiens (Human) go to search |
| UniProt | P12104 go to UniProt |
| Experimental structures | 7 structures in PDB for P12104 go to PDBe-KB |
| Biological function | FABP are thought to play a role in the intracellular transport of long-chain fatty acids and their acyl-CoA esters. FABP2 is probably involved in triglyceride-rich lipoprotein synthesis. Binds saturated long-chain fatty acids with a high affinity, but binds with a lower affinity to unsaturated long-chain fatty acids. FABP2 may also help maintain energy homeostasis by functioning as a lipid sensor. go to UniProt |

3D viewer

Model Confidence:

- Very high (pLDDT > 90)
- Confident (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

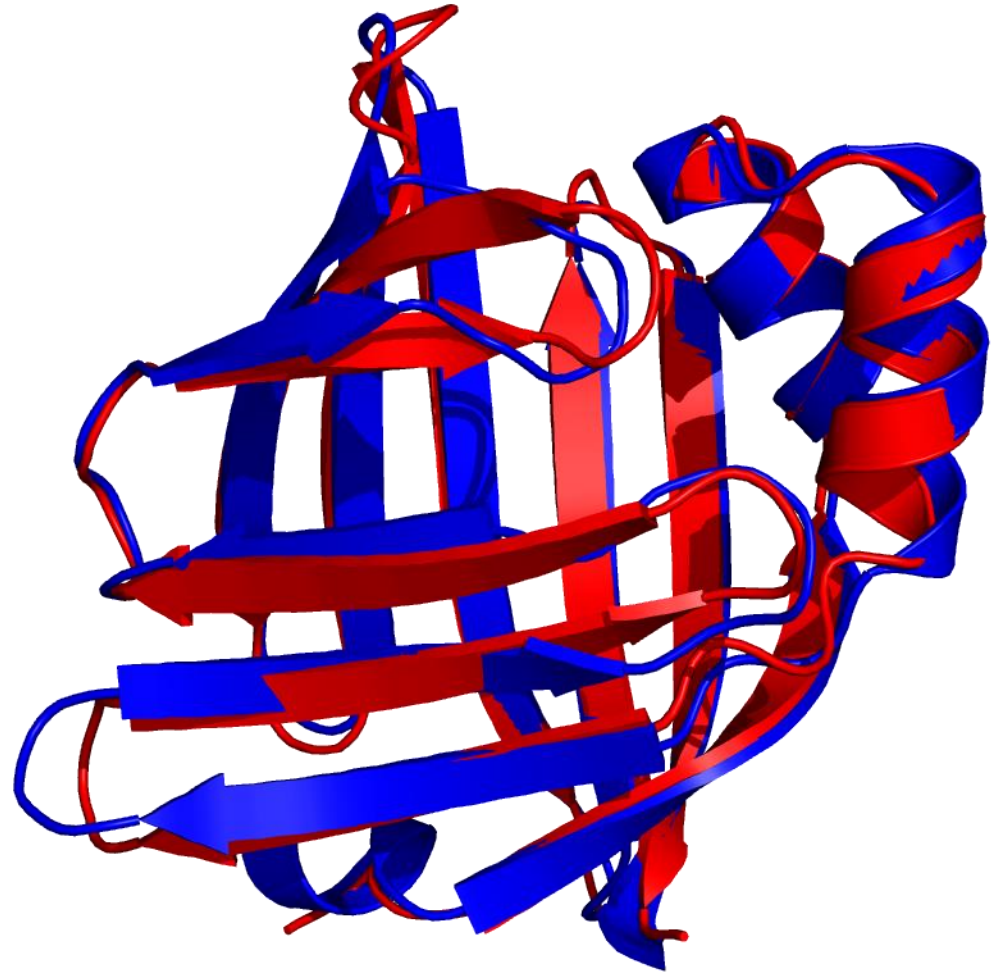
AlphaFold produces a per-residue confidence score (pLDDT) between 0 and 100. Some regions below 50 pLDDT may be unstructured in isolation.

Sequence of AF-P12104-F1 Chain 1: Fatty acid... A

```
1  MAFPSITWVVDRENYDFRNGRGGVNIIVKRLAAMENLRLITPQGGKFTVRESAFAFNIEVYFELGUTFDYSLADOTELAQWSLEGNKLGRFFRFDGQHELVTVREIIGELVQ  
120  TVVYEGVARRIFRKEG
```

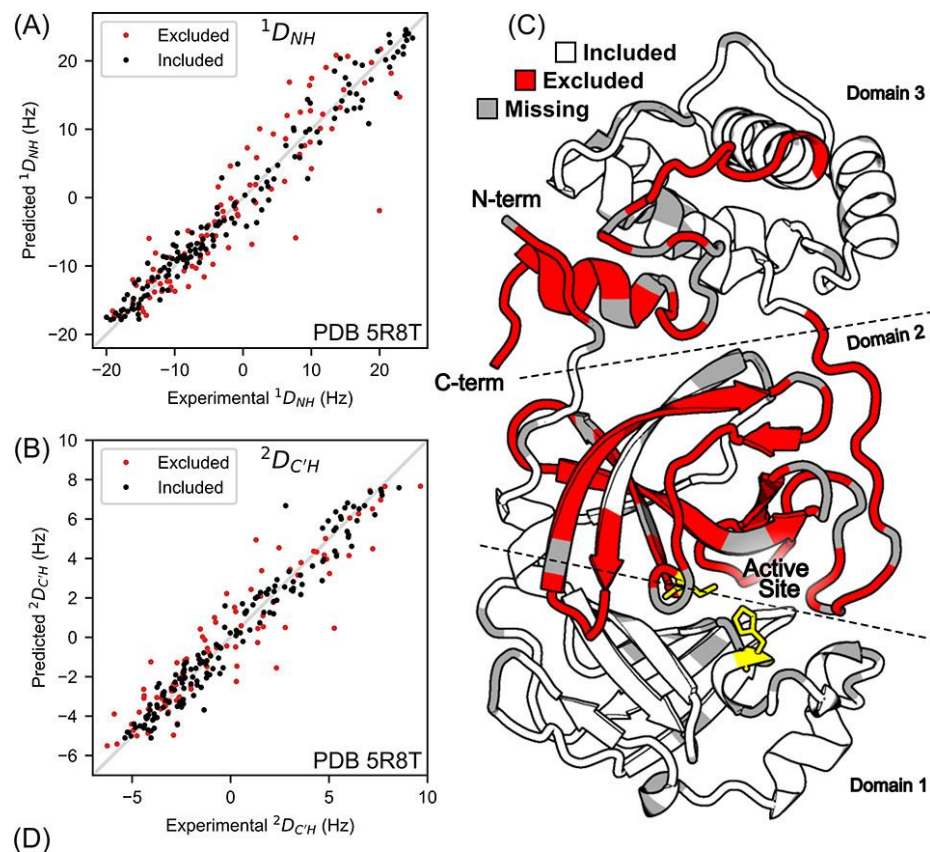
Comparison of AlphaFold2 vs 6L90

- **Red:** AlphaFold2
- **Blue:** Experimental crystal structure
- Aligned using PyMOL (align command)



AlphaFold2 Limitations

- Performs well for folded, compact regions
- Less good on loops, dynamic regions (SARS-CoV2 MPro, right)
- Very bad on disordered proteins (IDPs) → makes sense!
- **Verdict:** It's a great starting point, like many other models



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
- AlphaFold2 is probably the future, but not good for disordered proteins; it's still a model!