# FoldSeek: Fast search and comparison of protein structures

## INTRODUCTION

For today's workshop, we're going to be running through an example of how to create a database of protein structures and query the database for structural similarity to a desired protein. We'll also be going over how to directly compare the structure of two proteins directly against each other, and display the results in tab-delimited and HTML format.

## SET UP  [## Estimated runtime: < 2 minutes ]

1. Login to Atlas Open OnDemand using your web browser (https://atlas-ood.hpc.msstate.edu/) and navigate to your working directory for this workshop. Mine, for example, is in the shared directory under my username ( Files > /90daydata > Change directory > /90daydata/shared/olivia.haley )

2. Once in your working directory, select **Open in Terminal**. A new window should open.

3. Copy the shared directory containing the scripts and structures for this demo to your working directory, then activate the conda environment for this workshop.

```
# Copy the shared directory from the shared folder
cp -r /90daydata/shared/protein_structure_workshop/FoldSeek/ .

# Activate the workshop conda environment
module load miniconda3
source activate /90daydata/shared/protein_structure_conda/foldseek_env

#Navigate into the directory and view its contents
cd FoldSeek
ls -ltr
```

The directory's contents should look like the following:

```
drwxr-s---  2 olivia.haley proj-maizegdb   4096 Nov  7 14:56 log
drwxr-s--- 16 olivia.haley proj-maizegdb   4096 Nov  7 15:01 tmp
drwxr-s---  2 olivia.haley proj-maizegdb   4096 Nov  7 15:01 databases
drwxr-s---  2 olivia.haley proj-maizegdb 139264 Nov  7 15:02 viridiplantae_PDB_structures
drwxr-s---  2 olivia.haley proj-maizegdb   4096 Nov  7 15:02 scripts
drwxr-s---  2 olivia.haley proj-maizegdb   4096 Nov  7 15:02 query_structures
drwxr-s---  2 olivia.haley proj-maizegdb   4096 Nov  7 15:02 examples
-rw-r-----  1 olivia.haley proj-maizegdb 721547 Nov  7 15:02 AF-Q6XFQ4-F1-model_v4.pdb
```

# TUTORIAL

The first step is to create the target database, which will contain the structures that you will search against using your query protein(s). Target databases can be established using a directory of protein structure files or fasta files. Some pre-compiled databases are available for downloading directly. In this tutorial, we'll be looking for structural homologs of proteins with potential applications in agriculture, such as:

- KWL1, a defense protein in maize with antifungal properties [6FPG]
- GST-I, an enzyme contributing to herbicide detoxification in maize [1AXD]
- PPO, an enzyme which leads to browning in apple [6ELS]

For simplicity, our target database will be a subset of experimentally-determined structures from the *Viridiplantae* (or 'green plants') clade.

## Step 1. Create the target database [## Estimated runtime: < 1 minute ]

```
#Run the script to create the target database of Viridiplantae structures
sbatch scripts/s0_create_database.sh

#View the files generated while creating the database
ls -ltr databases
```

```
-rw-r----- 1 olivia.haley proj-maizegdb    36002 Oct 28 14:19 viridiplantae_PDBdb_ss.index
-rw-r----- 1 olivia.haley proj-maizegdb        4 Oct 28 14:19 viridiplantae_PDBdb_ss.dbtype
-rw-r----- 1 olivia.haley proj-maizegdb   677295 Oct 28 14:19 viridiplantae_PDBdb_ss
-rw-r----- 1 olivia.haley proj-maizegdb    27330 Oct 28 14:19 viridiplantae_PDBdb.source
-rw-r----- 1 olivia.haley proj-maizegdb    38070 Oct 28 14:19 viridiplantae_PDBdb.lookup
-rw-r----- 1 olivia.haley proj-maizegdb    36002 Oct 28 14:19 viridiplantae_PDBdb.index
-rw-r----- 1 olivia.haley proj-maizegdb    28310 Oct 28 14:19 viridiplantae_PDBdb_h.index
-rw-r----- 1 olivia.haley proj-maizegdb        4 Oct 28 14:19 viridiplantae_PDBdb_h.dbtype
-rw-r----- 1 olivia.haley proj-maizegdb    18960 Oct 28 14:19 viridiplantae_PDBdb_h
-rw-r----- 1 olivia.haley proj-maizegdb        4 Oct 28 14:19 viridiplantae_PDBdb.dbtype
-rw-r----- 1 olivia.haley proj-maizegdb    40066 Oct 28 14:19 viridiplantae_PDBdb_ca.index
-rw-r----- 1 olivia.haley proj-maizegdb        4 Oct 28 14:19 viridiplantae_PDBdb_ca.dbtype
-rw-r----- 1 olivia.haley proj-maizegdb  4146149 Oct 28 14:19 viridiplantae_PDBdb_ca
-rw-r----- 1 olivia.haley proj-maizegdb   677295 Oct 28 14:19 viridiplantae_PDBdb
```

## Step 2. Run FoldSeek (TSV output) to get initial query matches [## Estimated runtime: < 1 minute ]

```
#Run the script
sbatch scripts/s1_foldseek_run_initial_query.sh

#View the output tab-delimited file
head initial_query_results.tsv
```

```
6ELS_A  6ELS_A  1.000  459  0    0  1   459  1   459  0.000E+00  4154
6ELS_A  4Z13_A  0.461  466  246  0  1   458  40  505  8.126E-57  1958
6ELS_A  6HQJ_A  0.454  461  250  0  1   459  11  471  7.796E-58  1926
6FPG_D  6FPG_D  1.000  153  0    0  1   153  1   153  3.613E-33  1415
6FPG_D  6TI2_E  0.593  160  62   0  1   153  1   160  2.563E-23  942
6FPG_D  4PMK_A  0.569  149  60   0  3   151  18  158  2.299E-20  794
6FPG_D  4X9U_A  0.572  150  60   0  2   151  41  182  4.488E-20  786
6FPG_D  1N10_A  0.158  121  101  0  31  151  6   126  5.000E-04  109
6FPG_D  7KSN_A  0.134  121  101  0  35  152  1   121  3.146E-03  96
6FPG_D  4JP7_A  0.128  116  99   0  38  152  3   118  1.133E-02  92
```

## Step 3. Run FoldSeek (HTML output) to explore the initial query-target matches
## [## Estimated runtime: < 1 minute ]

We have some matches from each of our proteins, but let's explore one of the matches for maize kiwellin defense protein (6FPG_D). In particular, let's look at the alignment between the kiwellin and 4PMK_A. This script will output the alignment's results in HTML form as a file called **6FPGD_4PMKA.html**. To move the file to your local machine using the Atlas Open On Demand interface, click on the open tab, then Refresh > FoldSeek). Next to the file name there should be a drop-down box, click the drop-down then select Download.

```
#Run the script
sbatch scripts/s2_view_query_target_alignment.sh
```



As it turns out, this protein 4PMK_A is also a kiwellin! It comes from *Actinidia chinensis* var. *chinensis* (the Chinese soft-hair kiwi) where it was first identified as an allergen (Tamburrini et al., 2005). Its homolog in maize was later found to have antifungal properties (Han et al., 2019).

***Bonus Question: What insights can you make from the other proteins?***

## Step 4. Expand the search to other clades [## Estimated time: < 2 minutes ]

Until this point, we've been using a database of experimentally-determined protein structures. But less than 10% of the structures in the Protein Data Bank come from plants! For this exercise, we'll expand our structural homology search by using a database of computational protein structures in FoldSeek. FoldSeek has pre-compiled databases such as AlphaFoldDB, PDB, and ESMAtlas. In this example we'll use the AlphaFold database of Swiss-Prot proteins. This script will create an HTML file that you'll need to transfer to your local computer with Atlas Open On Demand.

```
#Generates the database files for the AlphaFold/Swiss-Prot database
#Runs the query protein against the AlphaFold/Swiss-Prot database
sbatch scripts/s3_create_and_query_AF2_database.sh
```

Looking at the results, many of the top hits are kiwellins in plant species like rice (*Oryzae sativa*) and kiwi. What's interesting is that we do see structural homology for our maize kiwellin and a protein from *Streptomyces mobaraensis*, a spore-forming bacterium that is known to produce antimicrobial compounds (Zindel et al., 2013; P86242).
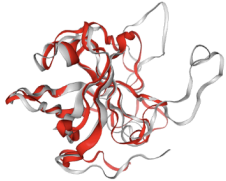
**Results**

| Target | Sequence Id. | Score | E-Value | Query Pos. | Target Pos. | |
|---|---|---|---|---|---|---|
| AF-A0A1D6GNR3-F1-model_v4 | 0.955 | 1140 | 1.45e-27 | 1-151 (153) | 41-198 (198) | ≡ |
| AF-Q9FWT5-F1-model_v4 | 0.683 | 948 | 6.41e-23 | 1-151 (153) | 53-213 (213) | ≡ |
| AF-P84527-F1-model_v4 | 0.578 | 817 | 3.89e-20 | 2-151 (153) | 72-213 (213) | ≡ |
| AF-Q7XVA8-F1-model_v4 | 0.557 | 807 | 1.66e-19 | 1-151 (153) | 29-183 (183) | ≡ |
| AF-Q9FWU1-F1-model_v4 | 0.546 | 806 | 5.44e-20 | 2-151 (153) | 56-216 (216) | ≡ |
| AF-Q9M4H4-F1-model_v4 | 0.568 | 803 | 3.68e-20 | 1-151 (153) | 77-220 (220) | ≡ |
| AF-Q6H5X0-F1-model_v4 | 0.577 | 798 | 1.19e-19 | 2-151 (153) | 40-192 (192) | ≡ |
| AF-P85261-F1-model_v4 | 0.565 | 786 | 1.19e-19 | 2-151 (153) | 72-213 (213) | ≡ |
| AF-Q7XD66-F1-model_v4 | 0.492 | 578 | 2.96e-14 | 30-151 (153) | 51-167 (167) | ≡ |
| AF-Q8LN49-F1-model_v4 | 0.462 | 465 | 2.37e-14 | 1-153 (153) | 26-170 (276) | ≡ |
| AF-Q7XD65-F1-model_v4 | 0.429 | 409 | 3.51e-11 | 35-151 (153) | 45-162 (162) | ≡ |
| AF-O42799-F1-model_v4 | 0.227 | 192 | 0.0000038 | 1-151 (153) | 140-270 (270) | ≡ |
| AF-P86242-F1-model_v4 | 0.227 | 190 | 0.0000145 | 29-152 (153) | 32-143 (143) | ≡ |

TM-Score: 0.71608

```
Q  29 GCSPPVTGSTRAVLTLNSFAEGGGGAAACTGKFYDDSK-KVVALSTGW---YNGGS--RCR-KHIMIHAGNGNSVSALVV
         P+     + +T  +      G +AC G   D S   +VA+  +W     N  +   CR   + + NG ++    V
T  32 SADIPIGQKMTGKMTYYTDK----GYGAC-GTPIDASSQDLVAIPAAWWTTPNPNNDPLCRGVSVEV-SYNGRTIRVPVR
```

**Step 4. When might FoldSeek not perform as expected? [## Estimated time: < 2 minute ]**

FoldSeek performs a rigid structural alignment, meaning that it doesn't account for the flexibility of protein backbones during the structural similarity search. There are cases (particularly when using computational protein structures) where this can lead to inaccurate conclusions. For example, let's compare the experimental structure of a maize photosystem I protein, with its AlphaFold2 structure. Run the script, and then download the HTML output file (maize-phytochrome-comparison.html)

```
#Generates the database files for the AlphaFold/Swiss-Prot database
#Runs the query protein against the AlphaFold/Swiss-Prot database
sbatch scripts/s4_compare_protein_structures.sh
```

Although these two are the same protein, this is not a great alignment. The TM-score is < 0.50 (indicating they're not assuming the same general fold). In this case, it looks like we have a protein domain in our computational structure (red) that is not in the same orientation as the domain in the experimental structure (gray).

***Bonus Question: What can we use to perform the structural alignment and 'fix' the domain orientation?***

**Step 5 (Optional). Flexible structure alignment [## Estimated time: < 5 minute ]**
Often, these cases can be corrected by using a flexible structural alignment program. We'll use one such program, called FATCAT to perform the flexible alignment. For a small set of alignments, it'll likely be easier to use their web interface https://fatcat.godziklab.org/fatcat/fatcat_pair.html. However, FATCAT does have a local implementation (https://github.com/GodzikLab/FATCAT-dist) if you have a larger set of structural alignments to perform. On the web platform, you can use PDB codes to upload structure files directly from the PDB, or input your own .pdb files. Note that this may not perform as intended for mmCIF files, or for structures in the PDB which only have mmCIF files available.

Under ' *Enter the 1st structure* ' provide the PDB code **8ISK** and input **A** under Chain. Alternatively, you can download the structure from our workshop directory and upload it. Make sure to select ' *Upload PDB file:* '

Under ' *Enter the 2nd structure* ', select ' *Upload PDB file:* ' and upload the AlphaFold2 structure file of Q6XFQ4 (AF-Q6XFQ4-F1-v4.pdb)

Enter the 1st structure

Enter a name for your structure: [Experimental] (optional)
○ Upload PDB file:
   [Choose File] no file selected      Chain: [    ]
● Provide PDB code:
   [8isk]      Chain: [A]
○ Provide SCOP domain code:
   [            ]

Enter the 2nd structure

Enter a name for your structure: [Computational] (optional)
● Upload PDB file:
   [Choose File] AF-Q6XFQ4…el_v4.pdb   Chain: [A]
○ Provide PDB code:
   [            ]      Chain: [    ]
○ Provide SCOP domain code:
   [            ]

FATCAT will provide a couple of outputs. Of note is the *P*-value which indicates the statistical similarity of the structural alignment (testing the hypothesis if the alignment score occurred randomly). It will also give a breakdown of the number of residues included in the alignment, the RMSD, and how many 'twists' were needed to align the structures. To view the alignment, download the superimposed structures (.pdb file), and then drag the file into the Mol* viewer (https://molstar.org/viewer/).

The alignment of these two proteins is much better once we allow for flexibility in the protein backbone!

These two structures are significantly similar with P-value ⓘ of 0.00e+00 (raw FATCAT score ⓘ is 2074.01)
They have 819 equivalent positions with an RMSD ⓘ of 3.83Å and 3 twists ⓘ.

Detailed results:
○ FATCAT alignment file 🔗
○ Graph of FATCAT chaining result 🔗 (postscript version ⬇)
○ Superimposed structures ⬇ (a pdb file with structure 8ISKA and modified structure pdb2 stored as chains A and B)
○ Transformation matrices for alignment blocks ⬇
○ Differential Distance Matrix 🔗
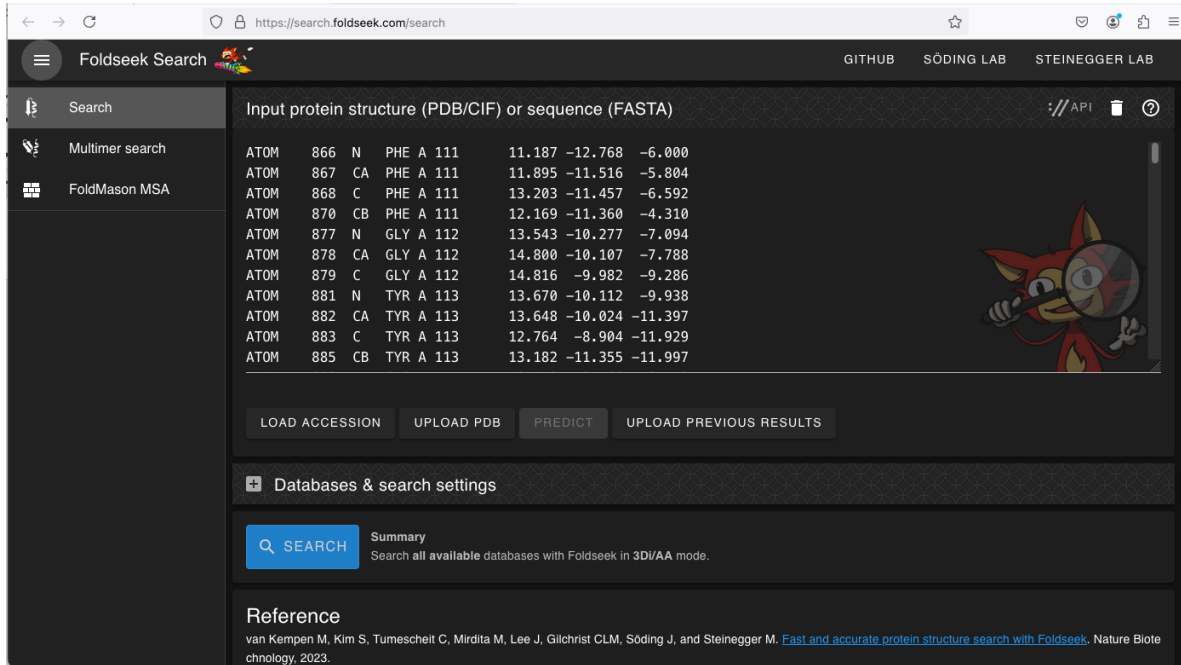○ Interactive viewer 🔗 (structures, alignment, contact map)

**Step 6. Deactivate the environment**
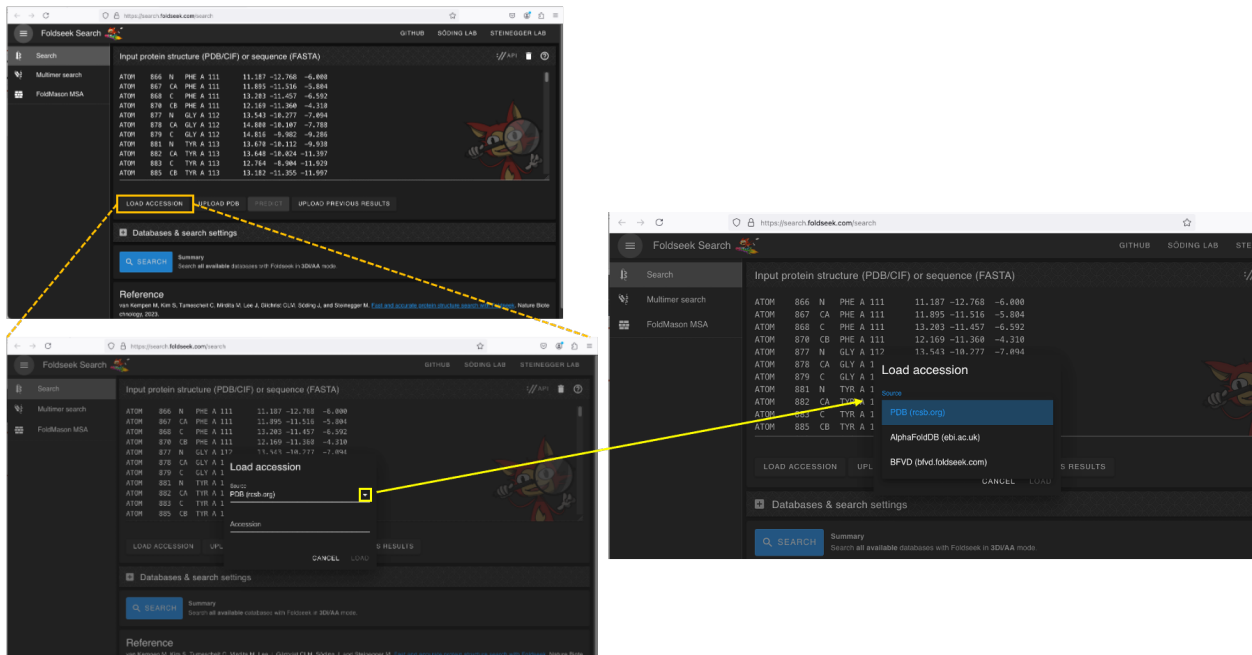
```
conda deactivate
```

# Foldseek (Online Server)

**Foldseek Search online server can be accessed here: [https://search.foldseek.com/search](https://search.foldseek.com/search)**

Following the link opens to the search portal.



Data can be uploaded from a local machine in either PDB or CIF formats using the "UPLOAD PDB" button. The "LOAD ACCESSION" button allows you to import a PDB from a repository (e.g, RCSB Protein Data Bank).

Results from a previous analysis can be uploaded using the "UPLOAD PREVIOUS RESULTS" button.



If desired, the databases and search settings can be modified (e.g., alignment and taxonomic filtering).

Here we are going to open our FASTA file and past our protein sequence from a FASTA file directly into the input field, click predict, and use ESMFold.
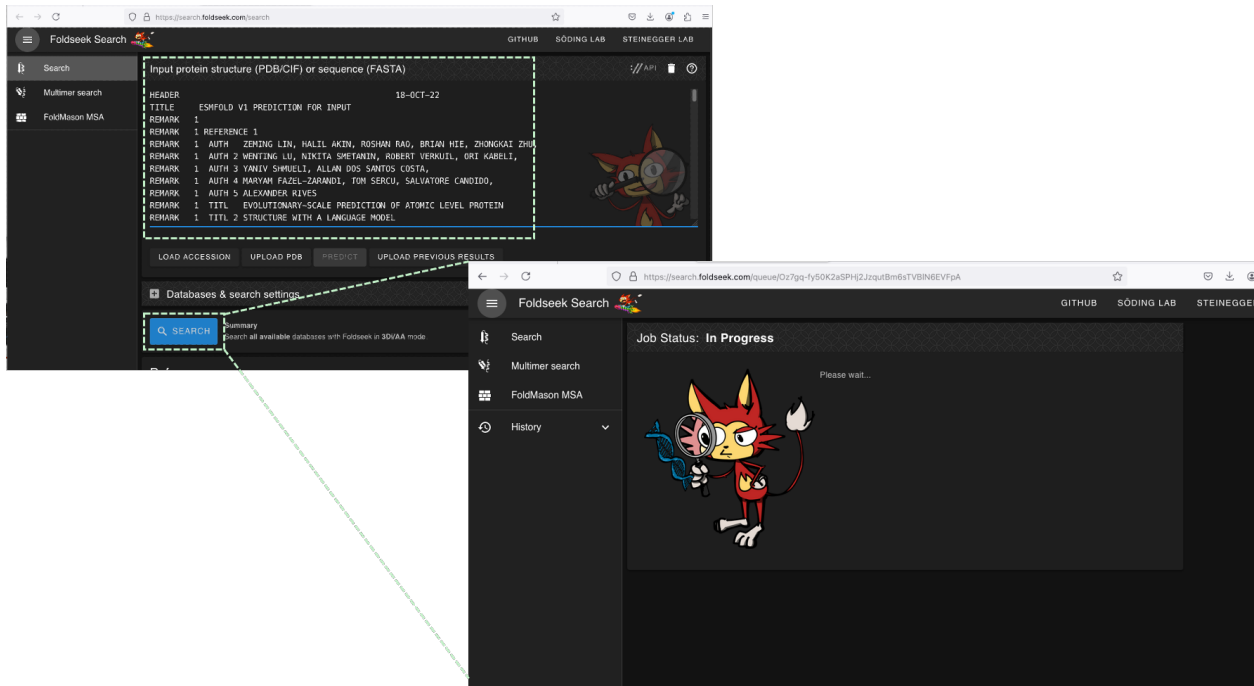
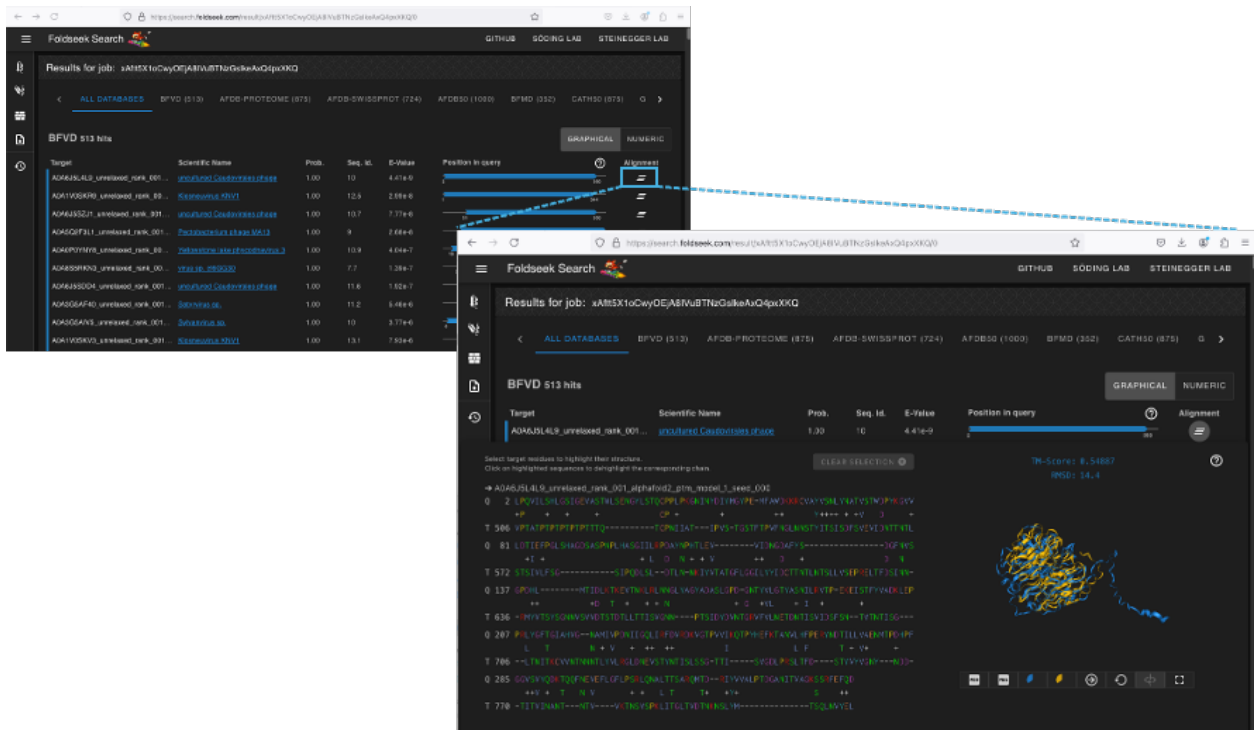**NOTE:** FASTA input cannot be uploaded and must be pasted in



After selecting the method, generate the prediction by clicking "predict".

Once the prediction completes, we can then search:



The hits from each database are made available, can now be explored, and saved.

RESULTS CAN BE SAVED IN EITHER PDB OR IMAGE FORMAT (PNG) BY CLICKING THE BUTTONS BELOW THE SUPERPOSITION.