# Final 'Find a Gene' Project

## Renee Zuhars (PID A17329856)

## Question 1

Beginning my search, I knew I wanted to limit my organism to some kind of fungus, because they are very understudied and I think their diversity and versatility are fascinating!

I decided to narrow my search to those proteins that help the fungus *Ophiocordyceps Unilateralis* "zombify" a host insect by taking over the hosts' neurological systems, eventually killing the host.

**protein name**: serine/threonine-protein kinase MAK, partial

**species**: *Ophiocordyceps Unilateralis*

**accession number**: ADI72911.1

**function**: The role of MAK-like kinases in this species is to induce behavioral changes in the host by interfering with Mitogen- Activated Protein Kinase signaling pathways. (ChatGPT)



Figure 1: Ophiocordyceps Unilaterialis

# Question 2

Attempting to find a homologous protein:

**blast method:** NCBI tblastn

**database:** est

**limits/restrictions:** none

My BLAST results were as followed:



Figure 2: BLAST results, original query

**Alignment of choice:** MgA0137f MgA Library Zymoseptoria tritici cDNA clone MgA0137 5', mRNA sequence

**E Value:** 3e-46

**Percent Identity:** 59.84%

**Percent Coverage:** 91%

# Question 3

Here is some information about the homolog I am looking into:



Figure 3: Zymoseptoria Triciti

**Name**: MgA0137f MgA Library Zymoseptoria tritici cDNA clone MgA0137 5', mRNA sequence

**Species Derived from**: Zymoseptoria tritici : this is a pathogenic fungus that attacks wheat plants. It is resistant to multiple fungicides, and causes septoria leaf blotch.

**FASTA format sequence, translated using EMBOSS Transeq**:

AW180074.1_1 MgA0137f MgA Library Zymoseptoria tritici cDNA clone MgA0137 5', mRNA sequence

RQLSVNSQGNHYAEIHRQEAERALVGASALKSPTGSQRESFFSHLRKRARRLSGRNSGVI
TPSMDAMETSAGCVPWAANKQTTFDTHSIASAAADPSSDPNFAELDRALQSVRYSLDAAA
NATQQARKPTNRVVEQPSLKRHHSLPHGVRHKTNPTTVYHDEH*STPRAADTRPPTKKKN
SRRSHELSASRRTAFSX

DNYQSTHRAITTPKFTGRKLSVLWLAQALSSHRLAAKEKASSLICARGREDFPAATQVSS
HLQWMLWKPALGAFLGLLTNKPPSTPTRSRLPQPIRHQTPISLSWIVHCKVYDTAWMPPR
TRLNKLGSLRTALSNHHSVTTRFLTALDTRPTQPPYTTTSTEARHEQPIQDPRRRRRI
LDEVMNSAHLAARRSR

TIISQLTGQSLRRNSPAGSACSGWRKRSQVTDWQPKRKLLLSSAQEGEKTFRPQLRCHH
TFNGCYGNQRWVRSLGCQTNHLRHPLDRVCRSRSVIRPQFRAGSCTAKCTIQPGCRRE
RDSTSEAYEPRSATIIEASPLASSRRTQDQPNHRIPRRALKHATSSRYKTPDEEEF
STKSTQRISPHGVLX

RERRAARCAEFMTSSRILLLRRGSCIGCSWRASVLVVVYGGWVGLVSNAVRKRVVTLQ*W
LLNYAVRRLPSLLSRVRGGIQAVSYTLQCTIQLSEIGV**RIGCGRRDRVGVEGGLFVSS
PRNAPSAGFHSIHRCDDTVAAGKSSRPLAQMREEAFSLAASRLESACANQSTLSFLP
VNFGVVIALVDLS

SRTPCGEMRVHDFVENSSSSSGVLYRLLVACFSARRGIRWLGWSCVRREEASGDASMM
VAQLRGSASLVESRSRRHPGCIVHFAVHDPAQRNWGLMTDRLRQTRSSGCRRWFVCQ
PKERTQRWFPHPLKVHLSCGRKVFSPSCADERRSFLFGCQSVTERLRQPEHAQLPA
GEFRRSDCPVSLIIVX

ENAVRRDALSSLRREFFFFVGGLVSAARGVLQCSSWYTVVGLVLCLTPGSEWRFNDG
CSTTRFVGFLACVAFAAASRLYRTLCSARSSSAKLGSDDGSAAADAIEWVSKVVCLLAA
QGTHPALVSIASIEGVMTPELRPESLLALLRR*EKKLSLWLPVGDLRALAPTRARSASCR
ISALPCELTDNCR

## Question 4

To determine if this protein is novel:

I used NCBI blastp, in the nr database.



Figure 4: BLAST results, novel query

**There is no match with 100% identity!**

4
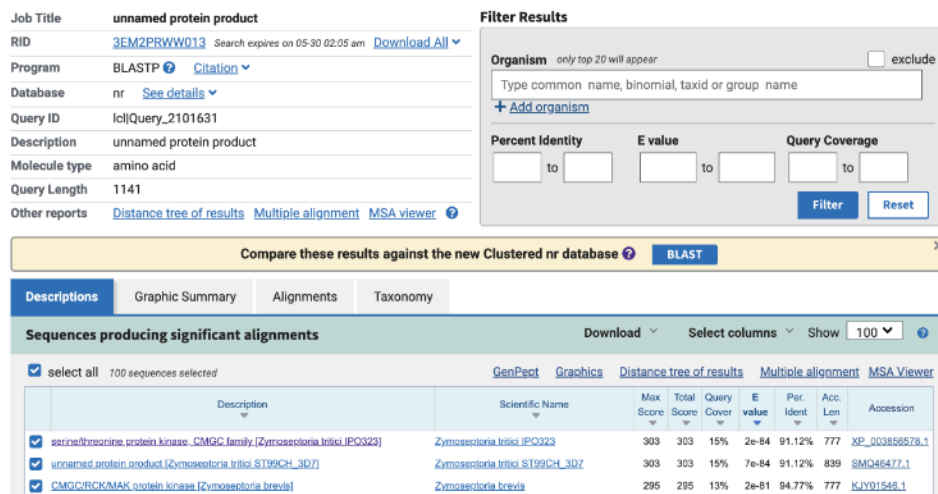
# Question 5

I will use MUSCLE at EBI to produce a multiple sequence alignment of the following proteins:

**My novel protein**

>Zymoseptoria unknown protein (novel protein from BLAST results)

**My original sequence for ophiocordyceps:**

>ADI72911.1 serine/threonine-protein kinase MAK, partial [Ophiocordyceps unilateralis]

**Other proteins of interest, based on BLAST results of original and novel proteins:**

*1.* >KAK4508075.1 hypothetical protein PRZ48_001812 [Zasmidium cellare]

*2.* >XP_047755397.1 Sporulation protein kinase pit1 [Fulvia fulva]

*3.* >KJY01546.1 CMGC/RCK/MAK protein kinase [Zymoseptoria brevis]

*4.* >KAK4981489.1 hypothetical protein LTR28_003093, partial [Elasticomyces elasticus]

*5.* >KAI5369935.1 putative serine/threonine-protein kinase, active [Septoria linicola]

**Alignment**

*Here is the alignment I obtained after running the above 7 proteins through MUSCLE at EBI (labeled by species), displayed in a code chunk because PDF formatting was giving me issues*

```
"
Zymoseptoria      LLRRGSCIGCSWRASVLVVVYGGWVGLVSNAVRKRVVTLQWLLNYAVRRLPSLLSRVRGG
Elasticomyces     ----------------VASHGNHYADAHRHEAEQALN----------------------
Zymoseptoria_br   ----------------VNSQGNHYAEIHRQEAERALV----------------------
Ophiocordyceps    --------------------------QEAERALS------------------------
Septoria          ----------------VNSQGNHYAELHRQEAERALN----------------------
Zasmidium         ----------------VNSQGNHYADIHRQEAERALT----------------------
Fulvia            ----------------VNSQGNHYADMHRQEAERALT----------------------


Zymoseptoria      IQAVSYTLQCTIQLSEIGVRIGCGRRDRVGVEGGLFVSSPRNAPSAGFHSIHRCDDTVAA
Elasticomyces     --------------------------GRNG---LASPTSSQRGSFFAHLRKRARRLS
Zymoseptoria_br   --------------------------GASA---LKSPTGSQRESFFSHLRKRARRLS
Ophiocordyceps    --------------------------GANG---RKSPTGTLLESFFSHLRKRARRLS
Septoria.         --------------------------GASG---LKSPTGSQRESFFSHLRKRARRFS
Zasmidium         --------------------------GANG---LKSPTGSQRESFFSHLRKRARRLS
```

```
Fulvia            -----------------------------GATG---LQSPTGSQRESFFSHLRKRARRLS

Zymoseptoria      GKSSRPLAQMREEAFSLAASRLESACANQSTLSFLPVNFGVVIALVDLSSRTPCGEMRVH
Elasticomyces     GRNQAPVSPSVDD--------------------IEASAG------------CAP----
Zymoseptoria_br   GRNSGVITPSMDA--------------------METNAG------------CVP----
Ophiocordyceps    GRNQGPMSPGAED--------------------LEANAG------------CAP----
Septoria          GKPSGLASPTAED--------------------MEANVG------------CAP----
Zasmidium         GRNQGPMSPGAED--------------------IEASVG------------CAP----
Fulvia            GRNSGPMSPSAED--------------------AEANVG------------CAP----

Zymoseptoria      DFVENSSSSSGVLYRLLVACFSARRGIRWLGWSCVRREEASGDASMMVAQLRGSASLVES
Elasticomyces     ----------------------------WAS-NRQSMAIE----------------
Zymoseptoria_br   ----------------------------WAA-NKQPVF-D----------------
Ophiocordyceps    ----------------------------WSS-NRGSIQ-E----------------
Septoria          ----------------------------WTT-NRQSIP-D----------------
Zasmidium         ----------------------------WSTNNRGSIQ-E----------------
Fulvia            ----------------------------WASNNRQSVQ-E----------------

Zymoseptoria      RSRRHPGCIVHFAVHDPAQRNWGLMTDRLRQTRSSGCRRWFVCQPKERTQRWFPHPLKVH
Elasticomyces     -----SLAITTHATDPSSDPNFAELDRALQNVRYSLDAGSYSNNNVQK------PVQKV-
Zymoseptoria_br   -----THSVASAAADPSSDPNFAELDRALQSVRYSLDAAANATQQARK------PTNRV-
Ophiocordyceps    -----PQPIEAVASDPSSDPNFAELDRALQNVRYSLDATANTSNNQPK------HPTKM-
Septoria          -----AQAIAPTAADPSVDPNFAELDRALQSVRYSLDATAGAMPTQPK------PPVKM-
Zasmidium         -----PQPIAPAAADPSTDPNFAELDRALQNVRYSLDAAANPANMQPK------HPTKM-
Fulvia            -----PQSIASVAVDPSSDPNFAELDRALQNVRYSLDAAAGAANPQPK------QPTKM-

Zymoseptoria      LSCGRKVFSPSCADERRSFLFGCQSVTERLRQPEHAQLPAGEFRRSDCPVSLIIVXENAV
Elasticomyces     ---------PSNPMLKR----------------HHSLPFGQDERISPVPAVNGPISSRT
Zymoseptoria_br   ---------VSNPSLKR----------------HHSLPHGVDDKTQ-------ANHRIP
Ophiocordyceps    ---------ASNPSLKR----------------HQSSHSG------------------
Septoria          ---------ASNPALKR----------------HHSLPYGKEEL--------SVVNRT
Zasmidium         ---------PSNQSLKR----------------HHSLPYGKEEIMSQT---GGSTSNRT
Fulvia            ---------ASNPTLKR----------------HHSVPCSKEEVTSNP-----SMANRT

Zymoseptoria      RRDALSSLRREFFFFVGGLVSAARGVLQCSSWYTVVGLVLCLTPGSEWRFNDGCSTTRFV
Elasticomyces     RRS--------------------------------VRQAPHPGHRYETPDEEEELL
Zymoseptoria_br   RRA--------------------------------LKHATSS--RYETPDEEEELL
Ophiocordyceps    ------------------------------------------------------
Septoria          RRS--------------------------------VKQAPSN-IRYETPCEEDELL
Zasmidium         RRS--------------------------------LRHAPSS--RYETPCEEDELL
Fulvia            RRS--------------------------------LRHAPSS--RYETPCEEDELL
```

```
Zymoseptoria       GFLACVAFAAASRLYRTLCSARSSSAKLGSDDGSAAADAIEWVSKV--VCLLAAQGTHPA
Elasticomyces      DEVLASAHRAARRLDRYIQQDNSPLPSVTSQQERARPPVQQVTSDP--GCFVPYLTPSPS
Zymoseptoria_br    DEVMTSAHLAARRL--------------DNEQLSRPPLPHVISEPVTVYTAPYLTPSHS
Ophiocordyceps     -----------------------------------------------------PAPS
Septoria           DEAIASAHQAVTRL------------DNGITQPARPHLPHVTSEP--TYNVPYLTPSPS
Zasmidium          DEALASVHAAATRL------DKGTS-NVAGTYQPSRPPIGHVTSEP--VYPAPYLTPSHS
Fulvia             DEALSNAHQAAQNL------DNAPSMNTTATYQPARPLLPQAISEP--AYAAPYLTPSPS


Zymoseptoria       LVSIASI
Elasticomyces      KDRNG--
Zymoseptoria_br    KDQMSLD
Ophiocordyceps     RKP----
Septoria           KDHMAVD
Zasmidium          KDQMNVT
Fulvia             KDQMNIS
"
```

[1] "\nZymoseptoria        LLRRGSCIGCSWRASVLVVVYGGWVGLVSNAVRKRVVTLQWLLNYAVRRLPSLLSRVRGG\nEla

# Question 6

Given the alignment above, I can now create a phylogenetic tree using the EBI's Simple Phylogeny tool.
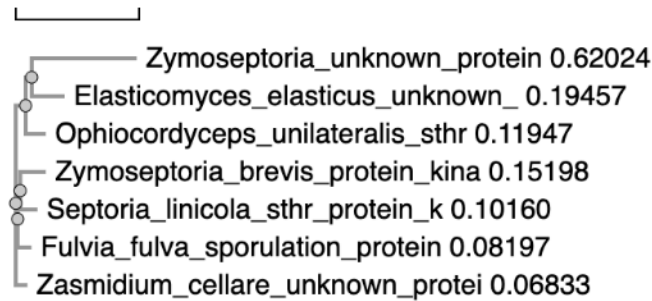


Figure 5: Phylogenetic tree for various fungal proteins

# Question 7

To generate a heatmap, the msa must be in fasta format. One of the MUSCLE outputs is 'the alignment in FASTA format converted by Seqret', so I used that.

```r
library(bio3d)

# Read multiple sequence alignment into R

msa <- read.fasta("finalseq.fasta")

# Calculate sequence identity matrix

seq_id_m <- seqidentity(msa)
```
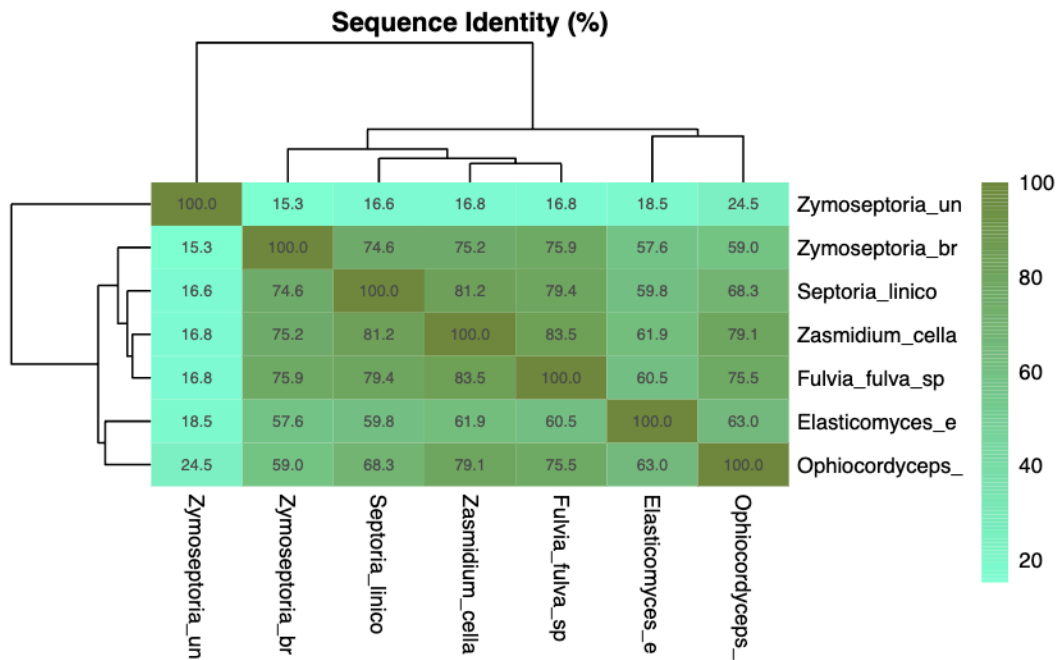
```r
library(pheatmap)

#Generate a heatmap of the msa (nicer enhancements courtesy of ChatGPT)

idm_percent <- seq_id_m * 100
rownames(idm_percent) <- substr(msa$id, 1, 15)
colnames(idm_percent) <- substr(msa$id, 1, 15) #shortening the labels

pheatmap(idm_percent,
         display_numbers = TRUE,      # show exact % in each cell
         number_format = "%.1f",      # format: 1 decimal place
         color = colorRampPalette(c("aquamarine", "darkolivegreen4"))(100),  # color gradient
         clustering_distance_rows = "euclidean",
         clustering_distance_cols = "euclidean",
         main = "Sequence Identity (%)",
         fontsize = 8,
         border_color = NA)           # cleaner look without gridlines
```

**Sequence Identity (%)**

|  | Zymoseptoria_un | Zymoseptoria_br | Septoria_linico | Zasmidium_cella | Fulvia_fulva_sp | Elasticomyces_e | Ophiocordyceps_ |
|---|---|---|---|---|---|---|---|
| Zymoseptoria_un | 100.0 | 15.3 | 16.6 | 16.8 | 16.8 | 18.5 | 24.5 |
| Zymoseptoria_br | 15.3 | 100.0 | 74.6 | 75.2 | 75.9 | 57.6 | 59.0 |
| Septoria_linico | 16.6 | 74.6 | 100.0 | 81.2 | 79.4 | 59.8 | 68.3 |
| Zasmidium_cella | 16.8 | 75.2 | 81.2 | 100.0 | 83.5 | 61.9 | 79.1 |
| Fulvia_fulva_sp | 16.8 | 75.9 | 79.4 | 83.5 | 100.0 | 60.5 | 75.5 |
| Elasticomyces_e | 18.5 | 57.6 | 59.8 | 61.9 | 60.5 | 100.0 | 63.0 |
| Ophiocordyceps_ | 24.5 | 59.0 | 68.3 | 79.1 | 75.5 | 63.0 | 100.0 |

**Quick note:** because the sequence of my novel protein is quite long (over 2000 characters), some shortening is needed before it can be processed by many of the tools that I'll be using going forward.

I was unsure of the best way to shorten my sequence, so I asked ChatGPT to do it for me. It identified a central 1000-amino acid segment that avoids the N-terminal signal peptide and C-terminal tail, which are less structured and may not be essential for the core function.

The resulting novel sequence, printed below, will be saved as "zymo_novel".

```
zymo_novel <- "VITPSMDAMETSAGCVPWAANKQTTFDTHSIASAAADPSSDPNFAELDRALQSVRYSLDAAANATQQARKPTNRVVEG
```

# Question 8

I am going to try to create a consensus for all my aligned sequences (msa), using the Bio3d package `consensus()`.

```
cons <- consensus(msa)
```

Because this sequence appears to have a lot of gaps in it, I will use the protein sequences that I aligned instead. I have saved them as text in some of the code chunks above.

I will use the `blast.pdb()` function to search the seven fungal protein sequences against the pdb database.

```
# blast.pdb(zymo_novel) yielded no results

# blast.pdb(ophiocordyceps) yielded no results

# blast.pdb(zasmidium)
```
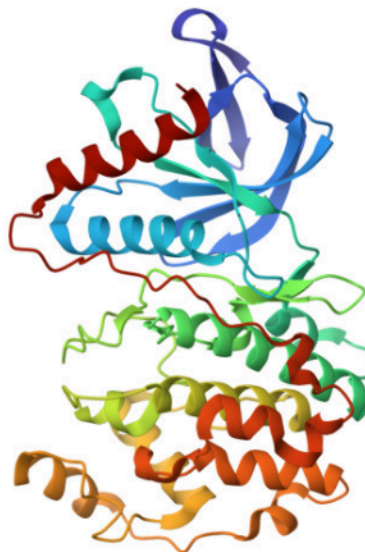


Figure 6: pdb match 1, 7W5C

**Sequence details**

PDB ID: 7W5C

E Value: 6.985e-33

Experimental technique: X-RAY DIFFRACTION

Resolution: 2.20 Å

Source organism: Arabidopsis thaliana
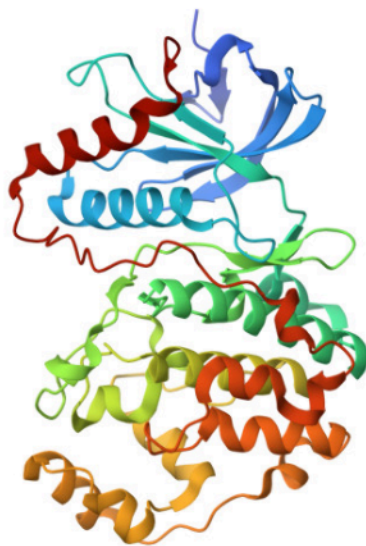
% identity: 34%

```
# blast.pdv(fulvia)
```

Figure 7: pdb match 2, 7E75

**Sequence details**

PDB ID: 7E75

E Value: 2.141e-28

Experimental technique: X-RAY DIFFRACTION

Resolution: 2.48 Å

Source organism: Homo Sapiens

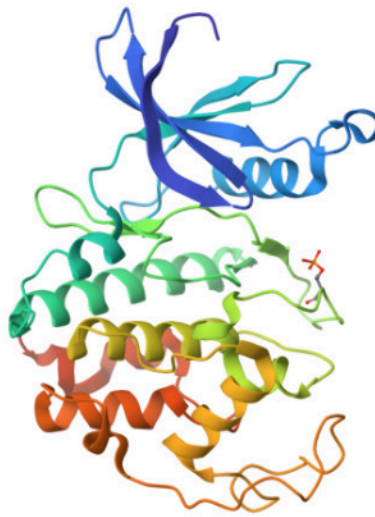% identity: 34%

```
# blast.pdb(z_brevis)
```

Figure 8: pdb match 3, 7NJ0

**Sequence details**

PDB ID: 7NJ0

E Value: 1.682e-40

Experimental technique: ELECTRON MICROSCOPY

Resolution: 3.60 Å

Source organism: Homo Sapiens

% identity: 34%

```
# blast.pdb(elasticomyces) yielded no results

# blast.pdb(septoria) yielded no results with a higher identity score than those above
```

# Question 9

I will run "zymo_novel" through AlphaFold.

It took a couple hours, but my output included several files that were saved to my project directory.

I then ran "/test_62382/test_62382_unrelaxed_rank_001_alphafold2_ptm_model_2_seed_000.pdb" through the online Mol* viewer (the conserved regions are rendered in spacefill):
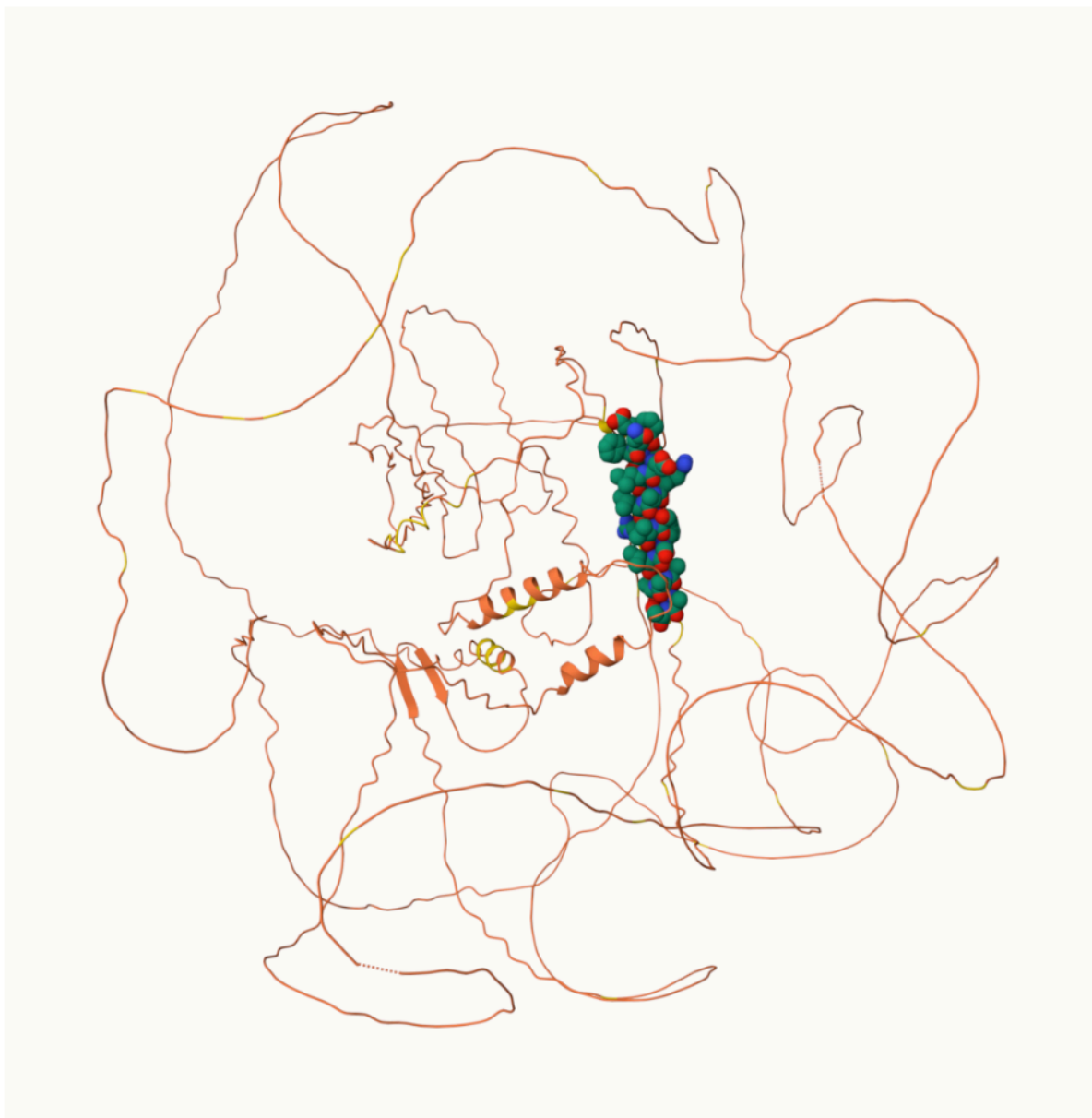


Figure 9: Final Unknown Protein Structure- colored by PLDDT confidence

Here it is without the spacefill, so the PLDDT coloring can be seen more clearly:
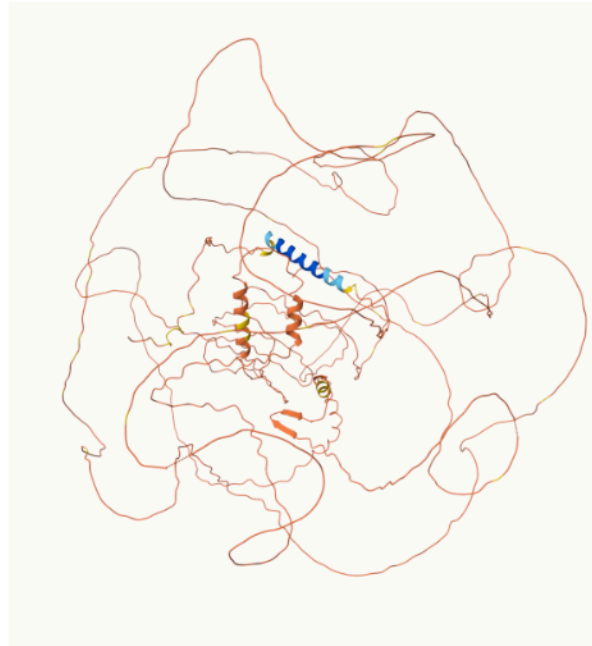
Figure 10: Final Protein Structure (no spacefill)

# Question 10

(i) I attempted to use the CASTPfold server to predict binding sites in my protein, but got the following result:
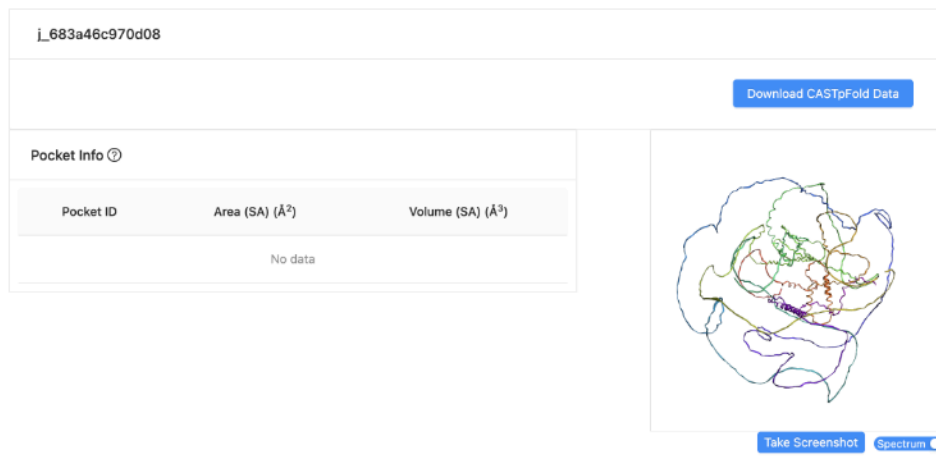


Figure 11: CASTPfold for my novel protein

So, unfortunately my protein appears to have no small molecule binding sites.

(ii) I then performed a CHEMBEL "target" search of my novel sequence. The resulting list consisted of 7 targets. I decided to focus on targets 4 and 5 (shown below) because they have the highest percent identity.

| # | E-Value | Positives % | Identities % | Score (bits) | Score | Length | ChEMBL ID | Name | UniProt Accessions | Type | Organism | Compounds | Activities |
|---|---------|-------------|--------------|--------------|-------|--------|-----------|------|--------------------|------|----------|-----------|------------|
| 1. | 0.03 | 36.8 | 25.7 | 38.1 | 87 | 1321 | CHEMBL1926492 | Tau-tubulin kinase 1 | Q5TCY1 | SINGLE PROTEIN | Homo sapiens | 142 | 185 |
| 2. | 0.21 | 39.7 | 24.8 | 34.7 | 78 | 329 | CHEMBL1255150 | G-protein coupled bile acid receptor 1 | Q80SS6 | SINGLE PROTEIN | Mus musculus | 482 | 786 |
| 3. | 2.6 | 41.9 | 27.4 | 31.6 | 70 | 1288 | CHEMBL1163123 | Mitogen-activated protein kinase kinase kinase 6 | O95382 | SINGLE PROTEIN | Homo sapiens | 447 | 537 |
| 4. | 4.2 | 52.9 | 47.1 | 30.8 | 68 | 759 | CHEMBL3317339 | Polymerase basic protein 2 | P03428 | SINGLE PROTEIN | Influenza A virus (strain A/Puerto Rico/8/1934 H1N1) | 100 | 131 |
| 5. | 4.2 | 52.9 | 47.1 | 30.8 | 68 | 759 | CHEMBL4523676 | RNA-directed RNA polymerase | P03433, P03431, P03428 | PROTEIN COMPLEX | Influenza A virus (strain A/Puerto Rico/8/1934 H1N1) | 3 | 4 |
| 6. | 5.3 | 47.5 | 27.9 | 30.4 | 67 | 571 | CHEMBL2069163 | Dual specificity testis-specific protein kinase 2 | Q96S53 | SINGLE PROTEIN | Homo sapiens | 173 | 173 |
| 7. | 9.8 | 47.3 | 29.1 | 29.6 | 65 | 910 | CHEMBL1250401 | Potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel 1 | Q88704 | SINGLE PROTEIN | Mus musculus | 14 | 22 |

Figure 12: My CHEMBL search

Target 4 appears to have the most useful data, as it has a defined compound that is effective against it, called Ribavirin.



Figure 13: Details about Ribavirin

15

This compound functions as an RNA polymerase inhibitor.

(iii) Overall, I do not believe my protein is druggable. My sequence is very long, and it was cut down so alphafold could process it. The resulting predicted structure did not have many conserved regions that I could use as starting points for further research. My CASTPfold search also did not yield any pockets. My CHEMBL searches yielded interesting results, but the highest percent identity listed was 47.1%, which does not place much confidence in the results. Perhaps if the whole sequence was run though a server with the capability to process it, a more accurate structure could be predicted.