

## BIMM-143: INTRODUCTION TO BIOINFORMATICS

The find-a-gene project assignment

<http://thegrantlab.org/bimm143>

Dr. Barry Grant

### **Overview:**

The find-a-gene project is a required assignment for BIMM-143. You should prepare a written report in **PDF** format that has responses to each question labeled **[Q1] - [Q10]** below. You may wish to consult the scoring rubric at the end of this document and the example report provided online (note that the example report is from a previous quarter and the questions may differ).

The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered in class.

### **Due Date:**

Your responses to questions Q1-Q4 are due at 12pm on the **Monday of Week 5** (see the Assignments and Grading section of our website for details). Note that these first set of answers can be obtained very quickly (at best within 15 or 20 minutes), so if you don't succeed at first, just keep trying.

The complete assignment, including responses to all questions, is due at 12pm on the **Monday of Week 10**.

### **Submission instructions:**

Your report formatted as a **PDF document** should be uploaded to **GradeScope**. Please make sure to include your UCSD email and PID number on the first page.

**Be sure to include your UCSD email and PID number on the first page of your report.**

Submit your preliminary report with answers to Q1-Q4 as soon as you can so we can determine if you have found a novel gene. Submit this preliminary report as one document with screen shots of the results inserted appropriately.

See the demonstration report linked to on the course website for an example of format. I will email you my decision; proceed with subsequent questions only after we are sure you have found a novel gene (and thus be successful in the later stages of the project).

For the final report add your results for Q5-Q10 to the preliminary report and submit the final document containing your results for all questions - **Please do not send only Q5-Q10 answers as the final report.**

**Name: Yi-Hung Lee**  
**PID: A16587141**

**Questions:**

**[Q1]** Tell me the name of a protein you are interested in. Include the species and the accession number. This can be a human protein or a protein from any other species as long as its function is known.

If you do not have a favorite protein, select human RBP4 or KIF11. Do not use beta globin as this is in the worked example report that I provide you with online.

Name: albumin preproprotein [Homo sapiens]

Accession: NP\_000468

Species: Homo sapiens

**[Q2]** Perform a BLAST search against a DNA database, such as a database consisting of genomic DNA or ESTs. The BLAST server can be at NCBI or elsewhere. Include details of the BLAST method used, database searched and any limits applied (e.g. Organism).

Method: TBLASTN 2.15.0+ search against nematode ESTs

Database: Expressed Sequence Tags (est)

Organism: Amphibia (taxid:8292)

Also include the output of that BLAST search in your document. If appropriate, change the font to Courier size 10 so that the results are displayed neatly. You can also screen capture a BLAST output (e.g. alt print screen on a PC or on a MAC press ⌘-shift-4. The pointer becomes a bulls eye. Select the area you wish to capture and release. The image is saved as a file called Screen Shot [].png in your Desktop directory). It is not necessary to print out all of the blast results if there are many pages.

On the BLAST results, clearly indicate a match that represents a protein sequence, encoded from some DNA sequence, that is homologous to your query protein. I need to be able to inspect the pairwise alignment you have selected, including the E value and score. It should be labeled a "genomic clone" or "mRNA sequence", etc. - but include no functional annotation.

In general, [Q2] is the most difficult for students because it requires you to have a “feel” for how to interpret BLAST results. You need to distinguish between a perfect match to your query (i.e. a sequence that is not “novel”), a near match (something that might be “novel”, depending on the results of [Q4]), and a non-homologous result.

If you are having trouble finding a novel gene try restricting your search to an organism that is

poorly annotated.

blastn

blastp

blastx

**tblastn**

tblastx

Enter Query Sequence

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

NP\_000468

Query subrange [?](#)

From To

Or, upload file

Choose File

No file chosen [?](#)

Job Title

NP\_000468:albumin preproprotein [Homo sapiens]

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database

Expressed sequence tags (est) [?](#)

Organism

Optional

Amphibia (taxid:8292) ☐ 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

[YouTube](#) Create custom database

Enter an Entrez query to limit search [?](#)

BLAST

Search database est using Tblastn (search translated nucleotide databases using a protein query)

☐ Show results in a new window

Chosen match: Accession CN062421.1, a 935 base pair clone from Ambystoma tigrinum. See below for alignment details.

<input checked="" type="checkbox"/>	<a href="#">JGI_CAARK10485.fwd.NIH_XGC_tropLiv1_Xenopus tropicalis cDNA clone IMAGE:7742695 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	250	313	78%	7e-77	40.83%	870	<a href="#">DN034543.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC6253.fwd.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7802597 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	250	312	70%	7e-77	40.07%	883	<a href="#">DR836429.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC1289.fwd.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7799341 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	311	73%	7e-77	42.40%	852	<a href="#">DN060715.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CAAR11724.fwd.NIH_XGC_tropLiv1_Xenopus tropicalis cDNA clone IMAGE:7743774 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	312	78%	8e-77	41.20%	863	<a href="#">DN036844.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CAAR3900.fwd.NIH_XGC_tropLiv1_Xenopus tropicalis cDNA clone IMAGE:7736489 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	313	78%	8e-77	41.20%	861	<a href="#">DN023018.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC10616.fwd.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7805416 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	250	312	73%	1e-76	41.61%	919	<a href="#">DR841374.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC1217.rev.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7799323 3', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	311	78%	1e-76	40.83%	877	<a href="#">DN060582.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABD10473.fwd.NIH_XGC_tropLun1_Xenopus tropicalis cDNA clone CABD10473 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	311	73%	1e-76	40.13%	901	<a href="#">DN077644.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABD14766.rev.NIH_XGC_tropLun1_Xenopus tropicalis cDNA clone CABD14766 3', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	249	312	74%	1e-76	42.16%	871	<a href="#">DN084472.1</a>
<input checked="" type="checkbox"/>	<a href="#">Ag2_p36_M20_M13R AG Ambystoma tigrinum tigrinum cDNA, mRNA sequence</a>	<a href="#">Ambystoma tigrin...</a>	249	249	47%	4e-76	41.18%	935	<a href="#">CN062421.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABD12445.fwd.NIH_XGC_tropLun1_Xenopus tropicalis cDNA clone CABD12445 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	248	310	76%	7e-76	41.64%	898	<a href="#">DN080425.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC11038.fwd.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7805714 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	248	310	72%	7e-76	42.14%	898	<a href="#">DR842130.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CAAR4381.fwd.NIH_XGC_tropLiv1_Xenopus tropicalis cDNA clone IMAGE:7736945 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	248	333	84%	8e-76	39.38%	959	<a href="#">DN023916.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CABC10416.fwd.NIH_XGC_tropFat1_Xenopus tropicalis cDNA clone IMAGE:7805413 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	247	309	70%	1e-75	40.14%	870	<a href="#">DR841020.1</a>
<input checked="" type="checkbox"/>	<a href="#">JGI_CAAR6459.fwd.NIH_XGC_tropLiv1_Xenopus tropicalis cDNA clone IMAGE:7738762 5', mRNA sequence</a>	<a href="#">Xenopus tropicalis</a>	246	310	77%	1e-75	40.78%	858	<a href="#">DN027791.1</a>

# Ag2\_p36\_M20\_M13R AG Ambystoma tigrinum tigrinum cDNA, mRNA sequence

Sequence ID: [CN062421.1](#) Length: 935 Number of Matches: 1

Range 1: 70 to 930 [GenBank](#) [Graphics](#)

▾ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps	Frame
249 bits(635)	4e-76	Compositional matrix adjust.	119/289(41%)	176/289(60%)	3/289(1%)	+1
Query 218	AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT	AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT	AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT	AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT	AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT	277
Sbjct 70	AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS	AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS	AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS	AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS	AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS	249
Query 278	ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK	ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK	ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK	ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK	ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK	337
Sbjct 250	MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP	MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP	MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP	MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP	MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP	429
Query 338	DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV	DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV	DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV	DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV	DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV	397
Sbjct 430	EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA	EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA	EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA	EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA	EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA	609
Query 398	FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV	FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV	FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV	FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV	FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV	457
Sbjct 610	EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI	EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI	EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI	EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI	EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI	789
Query 458	GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL	GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL	GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL	GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL	GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL	505
Sbjct 790	GQRCKLPEDQQMPCSEGLSLVLGQVC---QKTPFEIEKVAHCCKDSL	GQRCKLPEDQQMPCSEGLSLVLGQVC---QKTPFEIEKVAHCCKDSL	GQRCKLPEDQQMPCSEGLSLVLGQVC---QKTPFEIEKVAHCCKDSL	GQRCKLPEDQQMPCSEGLSLVLGQVC---QKTPFEIEKVAHCCKDSL	GQRCKLPEDQQMPCSEGLSLVLGQVC---QKTPFEIEKVAHCCKDSL	930

>Ag2\_p36\_M20\_M13R AG Ambystoma tigrinum tigrinum cDNA, mRNA sequence

Sequence ID: CN062421.1 Length: 935

Range 1: 70 to 930

Score:249 bits(635), Expect:4e-76,

Method:Compositional matrix adjust.,

Identities:119/289(41%), Positives:176/289(60%), Gaps:3/289(1%)

Query 218 AKQRLKCASLQKFGERAFAKAWAVARLSQRFPKAEFAEVSKLVTDLT

Sbjct 70 AVQKHNCYILESFKERALQAYKAVQTCQKFPHASFENAHS

Query 278 ADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESK

Sbjct 250 MVERMKLTTKTCEKKDELSTHLKECCEKPVRLERSACIVRLPNDEKPADLSQQVRQYIDDP

Query 338 DVCKNYAEAKDVFLGMFLYFYARRHPDYSVVLRLAKTYETTLEKCCAAADPHECYAKV

Sbjct 430 EVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNECIAKA

Query 398 FDEFKPLVEEPQNLIKQNCLEFQELGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKV

Sbjct 610 EETMKHEIEASKTLLKTTCAALEKMGYPFFQNLHITKYTPKLPCKVENLLHITKSMTTI

Query 458 GSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVS-DRVTKCCTESL 505

```

      G +CCK PE ++MPC+E  LS+VL Q+C   +KTP   ++V  CC +SL
Sbjct  790  GQRCKLPEDQQMPCSEGGLSLVLGQVC--QQKTPFEIEKVAHCKDSL  930

```

**[Q3]** Gather information about this “novel” **protein**. At a minimum, show me the protein sequence of the “novel” protein as displayed in your BLAST results from [Q2] as FASTA format (you can copy and paste the aligned sequence subject lines from your BLAST result page if necessary) or translate your novel DNA sequence using a tool called EMBOSS Transeq at the EBI. Don’t forget to translate all six reading frames; the ORF (open reading frame) is likely to be the longest sequence without a stop codon. It may not start with a methionine if you don’t have the complete coding region. Make sure the sequence you provide includes a header/subject line and is in traditional FASTA format.

Here, tell me the name of the novel protein, and the species from which it derives. It is very unlikely (but still definitely possible) that you will find a novel gene from an organism such as *S. cerevisiae*, human or mouse, because those genomes have already been thoroughly annotated. It is more likely that you will discover a new gene in a genome that is currently being sequenced, such as bacteria or plants or protozoa.

Chosen sequence:

```

>serum albumin precursor(sequence taken from BLAST result)
MMACMVERMKLTTKTCEKKDELSTHLKECCEKPVLESAIVRLPNDEKPADLSQQVRQYIDDPEVCKHFKEEGDTF
MGRFLCDYSKRHQDYSQELILRIGSGYEEVLKKCCAGEAHNECIAKAETMKHEIEASKTLLKTTCAALEKMGPIYFF
QNHLITKYTPKLPRCKVENLLHITKSMTTIGQRCKLPEDQQMPCSEGGLSLVLGQVCQQKTPFEIEKVAHCKDSL
S

```

Name:serum albumin precursor

Species: Ambystoma tigrinum

```

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Amphibia; Batrachia; Caudata; Salamandroidea; Ambystomatidae; Ambystoma

```

**[Q4]** Prove that this gene, and its corresponding protein, are novel. For the purposes of this project, “novel” is defined as follows. Take the protein sequence (your answer to [Q3]), and use it as a query in a blastp search of the nr database at NCBI.

- If there is a match with 100% amino acid identity to a protein in the database, from the same species, then your protein is NOT novel (even if the match is to a protein with a name such as “unknown”). Someone has already found and annotated this sequence, and assigned it an accession number.
- If the top match reported has less than 100% identity, then it is likely that your protein is novel, and you have succeeded.
- If there is a match with 100% identity, but to a different species than the one you started with, then you have likely succeeded in finding a novel gene.
- If there are no database matches to the original query from [Q1], this indicates that you have partially succeeded: yes, you may have found a new gene, but no, it is not

actually homologous to the original query. You should probably start over.

blastn

**blastp**

blastx

tblastn

tblastx

Standard Protein BLAST

Search protein databases using a translated nucleotide query

BLASTP programs search protein databases using a protein query. [more...](#)

Enter Query Sequence

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

>serum albumin precursor (sequence taken from BLAST result)  
MMACMVERMKLTTKTEKKDELSTHLKECKEPVLERSAIVRLPNDEKPADL  
SQQVRQYIDDPEVCKHF  
KEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKKCCAGEAHNECIAKA

Query subrange [?](#)  
From   
To

Or, upload file  No file chosen [?](#)

Job Title   
Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Databases ☒ Standard databases (nr etc.): New ☐ Experimental databases

Compare ☐ Select to compare standard and experimental database [?](#)

Standard

Database ?

Organism  ☐ exclude

Try experimental clustered nr database

For more info see [What is clustered nr?](#)

The top result is to a protein from *Ambystoma texanum*, see second screen shot below for alignment details:

Job Title  
serum albumin precursor (sequence taken from...

RID  
[S09KJFK3016](#) Search expires on 06-04 10:12 am [Download All](#) [?](#)

Program  
BLASTP [?](#) [Citation](#) [?](#)

Database  
nr [See details](#) [?](#)

Query ID  
Id|Query\_6781625

Description  
serum albumin precursor (sequence taken from BLAST result)

Molecule type  
amino acid

Query Length  
232

Other reports  
[Distance tree of results](#) [Multiple alignment](#) [MSA viewer](#) [?](#)

Filter Results

Organism only top 20 will appear ☐ exclude  
  
[+ Add organism](#)

Percent Identity  to  E value  to  Query Coverage  to

Compare these results against the new Clustered nr database [?](#)

Descriptions

Graphic Summary

Alignments

Taxonomy

Sequences producing significant alignments

Download [?](#) Select columns [?](#) Show  [?](#)

☒ select all 100 sequences selected

[GenPept](#) [Graphics](#) [Distance tree of results](#) [Multiple alignment](#) [MSA Viewer](#)

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/> <a href="#">serum albumin precursor [Ambystoma texanum]</a>	<a href="#">Ambystoma texa...</a>	345	397	100%	4e-112	73.39%	624	<a href="#">AAL56645.1</a>
<input checked="" type="checkbox"/> <a href="#">serum albumin precursor [Ambystoma maculatum]</a>	<a href="#">Ambystoma mac...</a>	328	328	100%	3e-105	68.67%	626	<a href="#">AAL56646.1</a>
<input checked="" type="checkbox"/> <a href="#">hypothetical protein NDU88_001408 [Pleurodeles waltl]</a>	<a href="#">Pleurodeles waltl</a>	301	301	99%	1e-94	59.31%	653	<a href="#">KAJ1205990.1</a>
<input checked="" type="checkbox"/> <a href="#">hypothetical protein NDU88_001408 [Pleurodeles waltl]</a>	<a href="#">Pleurodeles waltl</a>	266	266	99%	1e-81	51.95%	621	<a href="#">KAJ1205991.1</a>
<input checked="" type="checkbox"/> <a href="#">serum albumin precursor [Plethodon chatahooshee]</a>	<a href="#">Plethodon chatta...</a>	261	261	100%	9e-80	52.79%	602	<a href="#">AFM52271.1</a>

[Download](#) ▼
 [GenPept](#)
[Graphics](#)
 Sort by: E value ▼

## serum albumin precursor [Ambystoma texanum]

Sequence ID: [AAL56645.1](#) Length: 624 Number of Matches: 2

Range 1: 283 to 507 [GenPept](#) [Graphics](#)

▼ [Next Match](#) ▲ [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
345 bits(885)	4e-112	Compositional matrix adjust.	171/233(73%)	190/233(81%)	9/233(3%)
Query 1	MMACMVERMKLT	TTKTCEKKDELSTHLKECCEKPV	LERSACIVRLPNDEKPADLSQQVRQY	60	
Sbjct 283	MMACMERMKLTT	+TCEKK	+CCEKPV	LERS CIVRLPNDEKPADLS +VR Y	334
Query 61	IDDPEVCKHFKEEGD	TFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKCCAGEAHNEC	120		
Sbjct 335	DDPEVCK FKEEGD	FMGRFLCDY+K H ++S EL LRI SG E+ K CCAGEAHNEC	394		
Query 121	IAKAEETMKHEIEASK	TLLKTTCAALEKMGYPYFFQNH	LITKYTPKLPRCKVENLLHITKS	180	
Sbjct 395	IAK EET++HEIEASK	LKTTC ALEK+GPY FQN +I +YT LPR LL+ITK+	454		
Query 181	MTTIGQRCKLPEDQQMPCSEGG	LSLVLGQVCQ-QKTPFEIEKVAHCCKDSLS	232		
Sbjct 455	+T IGQ+CCKLPEDQQMPCSEGG	L +V Q+CQ QKTPFE EK+AHCCCKDSLS	507		

Range 2: 121 to 279 [GenPept](#) [Graphics](#)

▼ [Next Match](#) ▲ [Previous Match](#) ▲ [First Match](#)

Score	Expect	Method	Identities	Positives	Gaps
52.4 bits(124)	0.002	Compositional matrix adjust.	38/163(23%)	64/163(39%)	6/163(3%)
Query	29	CCEKPV LERSACIVRLPNDEKPADLSQQVRQYIDDPE-VCKHFKEEGDTFMGRFLCDYSK			87
Sbjct	121	CC+K ER C + + K D ++ PE +CK E D F+G ++ +			176
Query	88	RHQDYSQELILRIGSGYEEVLKCCAGEAH-NECIAKAEETMKHEIEASKTLLKTTCAAL			146
Sbjct	177	H IL ++ ++ CC EA +C+++ K E+E + K C L			236
Query	147	EKMGPYFFQNH LITKYTPKLPRCKVENLLHITKSMTTIGRCC		189	
Sbjct	237	+ + K P EN+ +T + + Q CC		279	
		QNFNERALRASKAAHACSKFPHASFENVQRLTDGIVHLHQTCC			

**[Q5]** Generate a multiple sequence alignment with your novel protein, your original query protein, and a group of other members of this family from different species. A typical number of proteins to use in a multiple sequence alignment for this assignment purpose is a minimum of 5 and a maximum of 20 - although the exact number is up to you. Include the multiple sequence alignment in your report. Use Courier font with a size appropriate to fit page width.

Side-note: Indicate your sequence in the alignment by choosing an appropriate name for each sequence in the input unaligned sequence file (i.e. edit the sequence file so that the species, or short common, names (rather than accession numbers) display in the output alignment and in the subsequent answers below). The goal in this step is to



create an interesting an alignment for building a phylogenetic tree that illustrates species divergence.

```
>NP_000468.1 albumin preproprotein [Homo sapiens]
MKWVTFISLLFLFSSAYSRGVFRDAHKSEVAHRFKDLGEENFKALVLIAFAQYLQQCPFEDHVKLVNEV
TEFAKTCVADESAENCDKSLHTLFGDKLCTVATLRETYGEMADCCAKQEPERNECFHQKDDNPPLRLV
RPEVDVMCTAFHDNEETFLKKLYEIIARRHPYFYAPELLFFAKRYKAAFTTECCQAADKAAACLLPKLDEL
RDEGKASSAKQRLKQKASLQKFGERAFAKAWAVARLSQRFPAEFAEVSKLVTDLTQVHTECCHGDLLECADD
RADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESKDVCKNYAEAKDVF
LGMFLYEYARRHPDYSVVLRLAKTYETTLKCCAAADPHECYAKVFDEFKPLVEEPQNLKQNCSELF
QLGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKVGSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEK
TPVSDRVTKCTESLVNRRPCFSALEVDETYVPKEFNAETFTFHADICTLSEKERQIKKQTALVELVKHK
PKATKEQLKAVMDDFAAFVEKCKADDKETCFAEEGKKLVAASQAALGL
```

```
>serum albumin precursor [Ambystoma texanum] (sequence taken from BLAST result)
MMACMVERMKLTTKTCEKKDELSTHLKECCEKPVLSACIVRLPNDEKPADLSQQVRQYIDDPEVCKHFKEEGDTFM
GRFLCDYSKRHQDYSQELILRIGSGYEEVLKKCCAGEAHNECIAKEETMKHEIEASKTLLKTTCAALEKMGPFYFFQN
HLITKYTPKLPRCKVENLLHITKSMTTIGQRCKLPEDQQMPCSEGGLSLVLGQVCQKTPFEIEKVAHCCKDSLS
```

```
>AAL56645.1:283-507 serum albumin precursor [Ambystoma texanum]
MMACMAERMKLTTQTCEKKKCEKPVLERSECIVRLPNDEKPADLSPEVRYFDDPEVCKRFKEEGDAFMGRFLCDYA
KIHPEHSAELNLRISGLEKAYKTCCAGEAHNECIAKEEETLRHEIEASKTKLKTTCGALEKLGPHYFQNMIVRYTG
ILPRSSDAFLLYITKTLTNIGQKCKLPEDQQMPCSEGGLGMVFAQICQNKTPFENEKLAHCCKDSLS
```

```
>KAJ1205990.1:309-539 hypothetical protein NDU88_001408 [Pleurodeles waltl]
MMGCMIERLHLTTRTCEKKDRISKHLKDCCDKDVIERSACIVKMNDEKPADLSPQVREYLEGPDVCKHYADEKDLYL
AKFSCDYAKRHPEFSLQLLLRVSKGYQDLLTKCCEENSHDCLIKGEEALKKEIESSTLLKTTCAAFEKLGPFYFQN
ELLVKYTRNIPQLTDESLLHITSGMTRIGQKCKIPEEKQMPCEGSLSLVIGEMCEKMPANFPNEKVTHCCSDT
```

```
>AFM52292.1:266-498 serum albumin precursor, partial [Plethodon yonahlossee]
MIACMEDRLALTTKTCAKKDELSSKLAACCEKPVVERSACIVKMDNDDRPADLSPQVREYIDVSVCKRFEDDKNELL
NHFLYDYSRRHPMSTEMLLKIVIGYDGVLVKCCCHKEDKLACLGKAEGEMKKEVQSSVELLKTNCAALEKVGSYHFEV
MLLGKYTLTMPQVTTPTLIHLIDDMTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCQKQEGHFVNNQVAHCCSDSYA
```

```
>AFM52295.1:266-498 serum albumin precursor, partial [Plethodon ouachitae]
MMACMEDRLALTTKTCAKKDELSSKLAACCEKPIVERSACIVKMDNDDRPADLSPQVREYIDVSVCKRFEDDKNELL
SHFLYDYSRRHPMSTEMLLKIVIGYDGLLVKCCHEEDKLACLGKAEGEMKKEVQSSVELLKTNCAALEKVGSYHFEV
MLLGKYTTMPQVTTPTLIHLIDDMTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCEKQEGHFVNNQVAHCCSDSYA
```

```
>AFM52300.1:266-498 serum albumin precursor, partial [Plethodon kentucki]
MMACMEDRLALTTKTCAKKDELSSKLAACCXKPVVERSACIVKMDNDDRPADLSPQVREYIDVSVCKRFEDDKNELL
SHFLYDYSRRHPMSTEMLLKIVIGYDGLLVKCCHEEDKLACLGKAEGEMKKEVQSSVELLKTNCAALEKVGSYHFEV
MLLGKYTTMPQVTTPTLIHLIDDMTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCEKQEGHFVNNQVAHCCSDSYA
```

## Alignment:

Obtained from CLUSTAL 1.7 multiple sequence alignment from Seaview

Query\_10001: albumin preproprotein [Homo sapiens]



Query\_10002: serum albumin precursor [Ambystoma texanum](sequence taken from BLAST result)  
Query\_10003: serum albumin precursor [Ambystoma texanum]  
Query\_10004: hypothetical protein NDU88\_001408 [Pleurodeles waltl]  
Query\_10005: serum albumin precursor [Plethodon yonahlossee]  
Query\_10006: serum albumin precursor [Plethodon ouachitae]  
Query\_10007: serum albumin precursor [Plethodon kentucki]

CLUSTAL W (1.7) multiple sequence alignment

```
Query_10001  LLECADDRADLAKYICENQDSISSKLKECCEKPLLEKSHCIAEVENDEMPADLPSLAADF
Query_10002  MMACMVERMKLTTKTCEKKDELSTHLKECCEKPVLEERSACIVRLPNDEKPADLSQQVRQY
Query_10003  MMACMAERMKLTTQTCEKK-----KCCEKPVLERSECIVRLPNDEKPADLSPEVRY
Query_10004  MMGCMIERLHLTTRTCEKKDRISKHLKDCCDKDVIERSACIVKMDENDEKPADLSPQVREY
Query_10005  MIACMEDRLALTTKTCAKKDELSSKLAACCEKPVVERSACIVKMDNDDRPADLSPQVREY
Query_10006  MMACMEDRLALTTKTCAKKDELSSKLAACCEKPIVERSACIVKMDNDDRPADLSPQVREY
Query_10007  MMACMEDRLALTTKTCAKKDELSSKLAACCXKPVVERSACIVKMDNDDRPADLSPQVREY
```

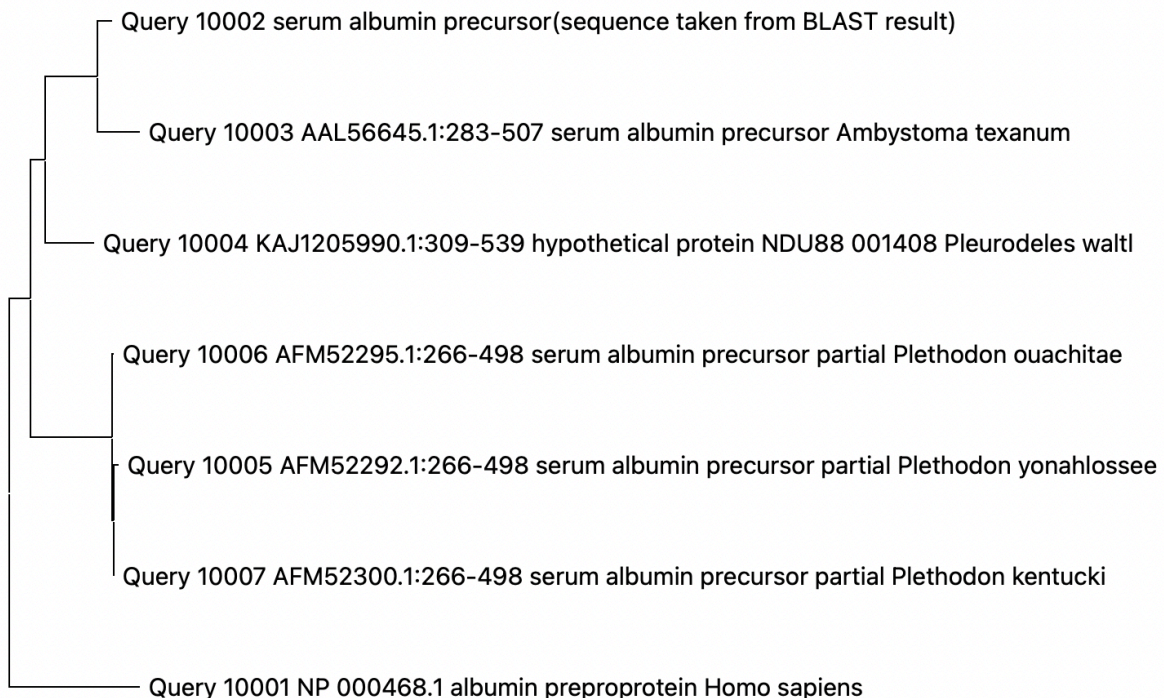
```
Query_10001  VESKDVCKNYAEAKDVFLGMFLYHEYARRHPDYSVVLRLAKTYETTLEKCCAAADPHEC
Query_10002  IDDPEVCKHFKEEGDTFMGRFLCDYSKRHQDYSQELILRIGSGYEEVLKKCCAGEAHNEC
Query_10003  FDDPEVCKRFKEEGDAFMGRFLCDYAKIHPEHSAELNLRASGLEKAYKTCCAGEAHNEC
Query_10004  LEGPDVCKHYADEKDLYLAKFSCDYAKRHPEFSLQLLLRVSKGYQDLLTKCCAEENSHDC
Query_10005  IDDVSVCKRFEDDKNELLNHFLYDYSRRHPEMSTEMLLKIVIGYDGLLVKCCHEEDKLAC
Query_10006  IDDVSVCKRFEDDKNELLSHFLYDYSRRHPEMSTEMLLKIVIGYDGLLVKCCHEEDKLAC
Query_10007  IDDVSVCKRFEDDKNELLSHFLYDYSRRHPEMSTEMLLKIVIGYDGLLVKCCHEEDKLAC
```

```
Query_10001  YAKVFDEFKPLVEEPQNLIKQNCSELFQLG EYKFQNALLVRYTKKVPQVSTPTLVEVSRN
Query_10002  IAKAEETMKHEIEASKTLLKTTCAALEKMGPFYFQNHLLITKYTPKLPRCKVENLLHITKS
Query_10003  IAKEEETLRHEIEASKTKLKTTCGALEKLGPHYFQNMIVRYTGILPRSSDAFLLYITKT
Query_10004  LIKGEEALKKEIESSTLLKTTCAAFEKLG PYSFQNELLVKYTRNIPQLTDESLHITSG
Query_10005  LGKAEGEMKKEVQSSVELLKTNCAALEKVG SYHFEVMLLGKYTLTMPQVTTPTLIHLIDD
Query_10006  LGKAEGEMKKEVQSSVELLKTNCAALEKVG SYHFEVMLLGKYTTMPQVTTPTLIHLIDD
Query_10007  LGKAEGEMKKEVQSSVELLKTNCAALEKVG SYHFEVMLLGKYTPTMPQVTTPTLIHLIDD
```

Query\_10001 LGKVGSKCKHPEAKRMPCAEDYLSVVLNQLCVLHEKTPVSDRVTKCCTESLV  
 Query\_10002 MTTIGQRCKLPEDQQMPCSEGLSLVLGQVCQ-QKTPFEIEKVAHCCCKDSLS  
 Query\_10003 LTNIGQKCKLPEDQQMPCSEGLGMVFAQICQNQKTPFENEKLAHCCCKDSLS  
 Query\_10004 MTRIGQKCKIPEEKQMPCEGSLSLVIGEMCEKMPANFPNEKVTHCCSDT--  
 Query\_10005 MTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCEKQEGHFVNNQVAHCCSDSYA  
 Query\_10006 MTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCEKQEGHFVNNQVAHCCSDSYA  
 Query\_10007 MTHVGEYCCKVPAEKQLPCSEGALGLIIGSMCEKQEGHFVNNQVAHCCSDSYA

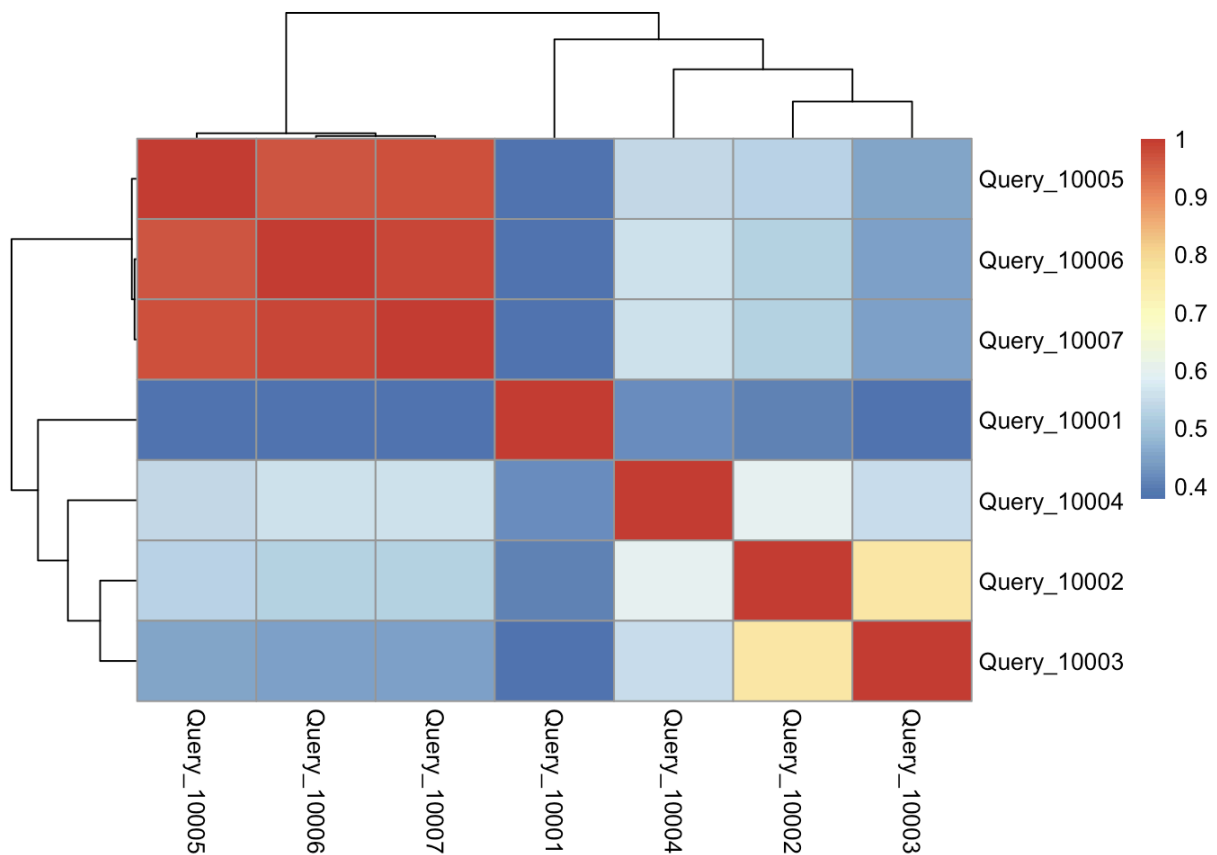
**[Q6]** Create a phylogenetic tree, using either a parsimony or distance-based approach. Bootstrapping and tree rooting are optional. Use “simple phylogeny” online from the EBI or any respected phylogeny program (such as MEGA, PAUP, or Phylip). Paste an image of your Cladogram or tree output in your report.

Import the sequences into MEGA, align with CLUSTAL, and create a neighbor-joining tree:



0.10

**[Q7]** Generate a sequence identity based **heatmap** of your aligned sequences using R. If necessary convert your sequence alignment to the ubiquitous FASTA format (Seaview can read in clustal format and “Save as” FASTA format for example). Read this FASTA format alignment into R with the help of functions in the **Bio3D package**. Calculate a sequence identity matrix (again using a function within the Bio3D package). Then generate a heatmap plot and add to your report. Do make sure your labels are visible and not cut at the figure margins.



Query\_10001: albumin preproprotein [Homo sapiens]

Query\_10002: serum albumin precursor [Ambystoma texanum](sequence taken from BLAST result)

Query\_10003: serum albumin precursor [Ambystoma texanum]

Query\_10004: hypothetical protein NDU88\_001408 [Pleurodeles waltl]

Query\_10005: serum albumin precursor [Plathodon yonahlossee]

Query\_10006: serum albumin precursor [Plathodon ouachitae]

Query\_10007: serum albumin precursor [Plathodon kentucki]

**[Q8]** Using R/Bio3D (or an online blast server if you prefer), search the main protein structure database for the most similar atomic resolution structures to your aligned sequences.

List the top 3 *unique* hits (i.e. not hits representing different chains from the same structure) along with their Evalue and sequence identity to your query. Please also add annotation details of these structures. For example include the annotation terms PDB identifier (structureId), Method used to solve the structure (experimentalTechnique), resolution (resolution), and source organism (source).

HINT: You can use a single sequence from your alignment or generate a consensus sequence from your alignment using the Bio3D function `consensus()`. The Bio3D functions `blast.pdb()`, `plot.blast()` and `pdb.annotate()` are likely to be of most relevance for completing this task. Note that the results of `blast.pdb()` contain the hits PDB identifier (or `pdb.id`) as well as Evalue and identity. The results of `pdb.annotate()` contain the other annotation terms noted above.

Note that if your consensus sequence has lots of gap positions then it will be better to use an original sequence from the alignment for your search of the PDB. In this case you could chose the sequence with the highest identity to all others in your alignment by calculating the row-wise maximum from your sequence identity matrix.

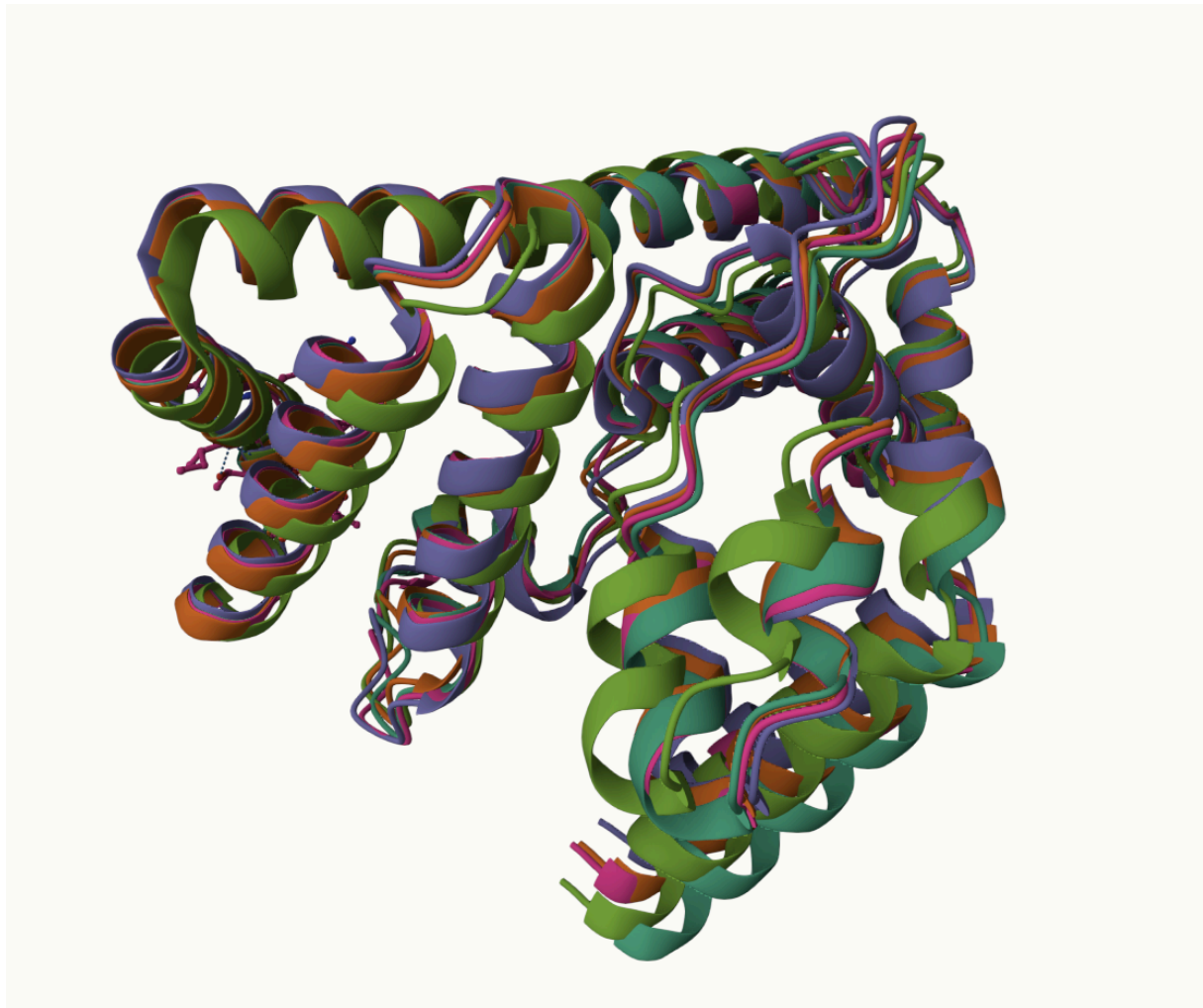
ID	Technique	Resolution	Source	E-value	Identity
4F5V	X-RAY DIFFRACTION	2.27	Oryctolagus cuniculus	1.02e-62	42.489
3V09	X-RAY DIFFRACTION	2.27	Oryctolagus cuniculus	1.16e-62	42.489
5Z0B	X-RAY DIFFRACTION	2.17	Homo sapiens	7.60e-62	41.631

**[Q9]** Using [AlphaFold notebook](#) generate a structural model using the default parameters for your novel protein sequence.

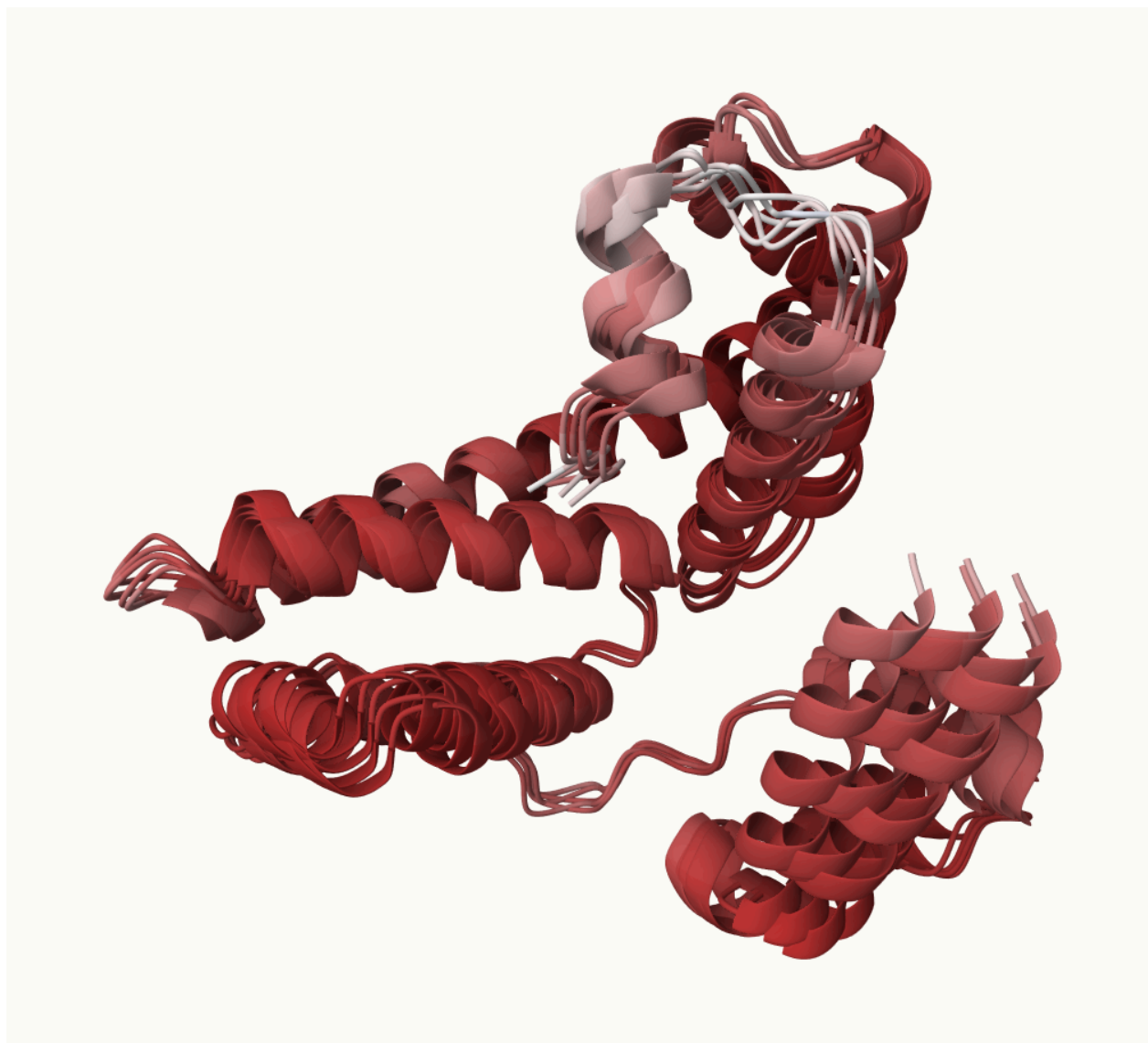
Note that this can take some time depending upon your sequence length. If your model is taking many hours to generate or your input sequence yields a “too many amino acids” (i.e. length) error you can focus on a single domain from your

sequence - identify region by searching for [PFAM](#) domain matches.

Once complete save the resulting PDB format file for your records. Finally, generate a molecular figure of your generated PDB structure using the **Mol\* viewer** online (or VMD/PyMol/Chimera if you prefer). To complete your analysis you can optionally highlight *conserved residues* that are likely to be functional as **spacefill** and the protein as **cartoon** colored by local alpha fold *pLDDT quality score*. This score is contained in the B-factor column of your PDB downloaded file. Please use a white or transparent background for your figure (i.e. not the default black in PyMol/VMD/Chimera etc.).



Annotated by pLDDT quality score:



**[Q10]** Perform a “Target” search of ChEMBL ( <https://www.ebi.ac.uk/chembl/> ) with your novel sequence. Are there any **Target Associated Assays** and **ligand efficiency data** reported that may be useful starting points for exploring potential inhibition of your novel protein? If there are no assays listed here simply list “non available as of [date]”.

CHEMBL details 1 Binding Assay (CHEMBL1055715) and 3 Functional Assays; No ligand efficiency data.

[https://www.ebi.ac.uk/chembl/assay\\_report\\_card/CHEMBL1055715/](https://www.ebi.ac.uk/chembl/assay_report_card/CHEMBL1055715/)

Binding assay linked manuscript tested a potent 3,4-disubstituted benzofuran P1' MMP-13 inhibitors. By replacing a backbone benzene with a pyridine and valine with threonine, compounds (e.g., 44) with greatly reduced plasma protein binding were also obtained.

W. Li, et al., 3,4-Disubstituted benzofuran P1' MMP-13 inhibitors: Optimization of selectivity and reduction of protein binding. Bioorganic & Medicinal Chemistry Letters 19, 4546–4550 (2009).

**Scoring Rubric:** [50 total points available]

**Q1** (4 points)

Protein name 1

Species 1

Accession number 1

Function known 1

**Q2** (6 points)

Blast method 1

Database searched 1

Limits applied 1

Search output list (top hits) 1

Alignment of choice 1

Evalue and other alignment stats 1

**Q3** (3 points)



Protein sequence of choice matches Subject above 1  
Name in header 1 Species 1

**Q4** (3 point)

Blastp output list with identities & Evalue 1 Top  
alignment shown with alignment statistics 1 Results  
indicates a “novel” gene found 1

**Q5** (3 points)

MSA labeled with useful names 1 MSA trimmed  
appropriately (i.e. no gap overhangs) 1 Pasted MSA  
fits report page width (i.e. font, format) 1

**Q6** (1 point)

Figure illustrates sequence clustering pattern 1

**Q7** (10 points)

Heatmap figure included in report 5 Heatmap is legible  
(i.e. no labels obscured) 5

**Q8** (9 points)

PDB identifiers from multiple species reported 5  
Annotation of PDB source, resolution and technique 4  
Annotation of Evalue and Sequence Identity 1

**Q9** (10 points)

Structure figure provided 2 Uses white background for  
molecular figure 1 Figure of high resolution (i.e. not just  
snapshot) 1 Conserved residues as spacefill 3 Protein  
cartoon colored by pLDDT quality score 3

**Q10** (1 point)

Evidence of ChEMBL searches 1