

BIMM-143: INTRODUCTION TO BIOINFORMATICS

The find-a-gene project assignment
https://bioboot.github.io/bimm143_S20/
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.

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 the beginning of class **Tuesday May 5th** (05/05/20) at 12pm San Diego time. Note that these answers can be obtained very quickly (at best within 10 or 15 minutes), so if you don't succeed at first, just keep trying.

The complete assignment, including responses to all questions, is due **Friday June 5th** (06/05/20) at 12pm San Diego time.

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.
JTT009@ucsd.edu
A17197773

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.

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.**

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 it's 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: RBP4 (Retinol-binding protein 4)

Accession: P02753

Species: Homo Sapiens

Function Known: Retinol binding protein 4, also known as RBP4, is a transporter protein for retinol.

[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: BLAST search using the BLAST method "blastp" against the "nr" (non-redundant) protein database at NCBI. I set the organism limit to "Bacteria" to increase the chances of finding a novel gene.

Database: Non-redundant protein sequences (nr)

Organism: Bacteria (taxid:2)

Chosen Match:

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/> lipocalin/fatty-acid binding family protein [Salmonella enterica]	Salmonella enterica	373	373	100%	8e-130	85.57%	201	MCQ7614318.1

Standard Protein BLAST

blastn **blastp** blastx tblastn tblastx

BLASTP programs search protein databases using a protein query. more...

Enter Query Sequence

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

Query subrange ?

From

To

Or, upload file

Choose File No file chosen ?

Job Title

P02753:RecName: Full=Retinol-binding protein...

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

Non-redundant protein sequences (nr) ?

Organism

Optional

Bacteria (taxid:2) ☐ exclude [Add organism](#)

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

Exclude

Optional

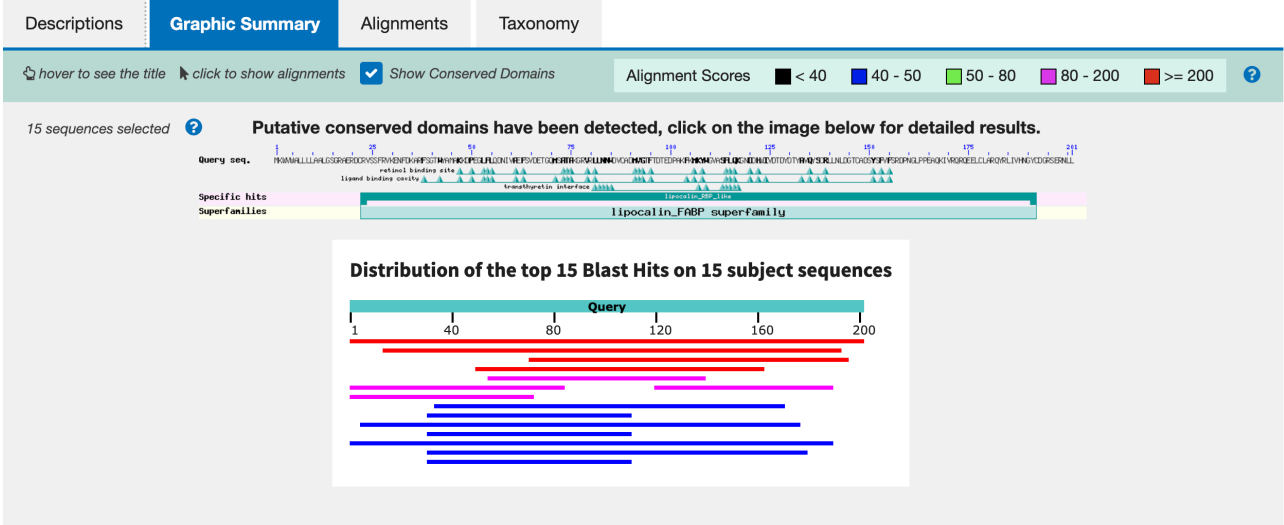
☐ Models (XM/XP) ☐ Non-redundant RefSeq proteins (WP) ☐ Uncultured/environmental sample sequences

[Try experimental clustered nr database](#) ?

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

[Reset page](#) [Bookmark](#)

Caption



[Download](#) [GenPept](#) [Graphics](#) [Next](#) [Previous](#) [Descriptions](#)

lipocalin/fatty-acid binding family protein [Salmonella enterica]
Sequence ID: [MCQ7614318.1](#) Length: 201 Number of Matches: 1

Range 1: 1 to 201 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
373 bits(957)	8e-130	Compositional matrix adjust.	172/201(86%)	189/201(94%)	0/201(0%)
Query 1	MKWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIV				60
Sbjct 1	M+W+WAL+LLAA+GSGRAERDCRVSSFRVKENFDKARFSGTWYA+AKKDPEGLFLQDNI+				60
Query 61	AEFSVDETGQMSATAKGRVRLNNWDVCADMVGTFDTEDPAKFKMKYWGVASFLQKGND				120
Sbjct 61	AEFSVDE G MSATAKGRVRL+NW+VCADMVGTFDTEDPAKFKMKYWGVASFLQ+GND				120
Query 121	DHWIVDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEELCLA				180
Sbjct 121	DHWIIDTDYETFALQYSCRLQNLGTCADSYSFVFSRDPNGL PE +K+VRQRQEELCL				180
Query 181	RQYRLIVHNGYCDGRSERNLL				201
Sbjct 181	RQYR I HNGYC + RN+L				201

Alignment statistics for match #1

Score	Expect	Method	Identities	Positives	Gaps
373	8E-130	Compositional	172/201(189/201(0/201(
Query 1					
MKWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIV 60					
M+W+WAL+LLAA+GSGRAERDCRVSSFRVKENFDKARFSGTWYA+AKKDPEGLFLQDNI+					
Sbjct 1					
MEWMWALVLLAAVGSRAERDCRVSSFRVKENFDKARFSGTWYAI AKKDPEGLFLQDNII 60					
Query 61					
AEFSVDETGQMSATAKGRVRLNNWDVCADMVGTFDTEDPAKFKMKYWGVASFLQKGND 120					
AEFSVDE G					
MSATAKGRVRL+NW+VCADMVGTFDTEDPAKFKMKYWGVASFLQ+GND					
Sbjct 61					
AEFSVDENGHMSATAKGRVRLSNWEVCADMVGTFDTEDPAKFKMKYWGVASFLQRGND 120					
Query 121					
DHWIVDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEELCLA 180					
DHWI+DTDY+T+A+QYSCRL NLDGTCADSYSFVFSRDPNGL PE					
+K+VRQRQEELCL					
Sbjct 121					
DHWIIDTDYETFALQYSCRLQNLGTCADSYSFVFSRDPNGLTPETRKLVRQRQEELCLD 180					
Query 181					
RQYRLIVHNGYCDGRSERNLL 201					
RQYR I HNGYC + RN+L					
Sbjct 181					
RQYRWIEHNGYCQSKLSRNIL 201					

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.

[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.

Chosen Sequence:

```
>MCQ7614318.1 lipocalin/fatty-acid binding family protein
[Salmonella enterica]
MEWMWALVLLAAVGSGRAERDCRVSSFRVKENFDKARFSGTWYAIKKDPEGLFLQDNIIA
EFSVDENGH
MSATAKGRVRLLSNWEVCADMVGTFDTEDPAKFVKMYWGVASFLQRGNDHWDIDTDYET
FALQYSCRL
QNLDGTCADSYSFVFSRDPNGLTPETRKLVRQRQEELCLDRQYRWIEHNGYCQSKLSRNIL
```

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.

Name: lipocalin/fatty-acid binding family protein

Species: *Salmonella enterica*

Bacteria; Pseudomonadota; Gammaproteobacteria; Enterobacterales;
Enterobacteriaceae; *Salmonella*

[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.

Details:

A BLASTP search against NR database (see setup in first screen-shot below) yielded a top hit result is to a novel protein from *Jaculus Jaculus* (Lesser Egyptian jerboa). See additional screen shots below for top hits and selected alignment details:

The screenshot shows the NCBI BLASTP search interface. At the top, there are tabs for different BLAST programs: blastn, **blastp**, blastx, tblastn, and tblastx. Below the tabs, there's a header bar with "BLASTP programs search protein databases using a protein query, more..." and buttons for "Reset page" and "Bookmark".

The main section is titled "Enter Query Sequence". It contains a text area for the query sequence, which is populated with a protein sequence from *Salmonella enterica*. There are fields for "Query subrange" (From and To) and a section for "Or, upload file" with a "Choose File" button. Below this is a "Job Title" field and a checkbox for "Align two or more sequences".

The next section is "Choose Search Set". It has a "Databases" section with radio buttons for "Standard databases (nr etc.)" (selected) and "Experimental databases". There's a link to "Try experimental clustered nr database". The "Compare" section has a checkbox for "Select to compare standard and experimental database". The "Standard" section has a "Database" dropdown set to "Non-redundant protein sequences (nr)". The "Organism" section has a text field and an "Add organism" button. The "Exclude" section has checkboxes for "Models (XM/XP)", "Non-redundant RefSeq proteins (WP)", and "Uncultured/environmental sample sequences".

The bottom section is "Program Selection". It has a "Algorithm" section with radio buttons for "Quick BLASTP (Accelerated protein-protein BLAST)", **blastp (protein-protein BLAST)** (selected), "PSI-BLAST (Position-Specific Iterated BLAST)", "PHI-BLAST (Pattern Hit Initiated BLAST)", and "DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)". There's a "Choose a BLAST algorithm" link.

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
✓	lipocalin/fatty-acid binding family protein [Salmonella enterica]	Salmonella ente...	422	422	100%	4e-149	100.00%	201	MCQ7614318.1
✓	retinol-binding protein 4 [Jaculus jaculus]	Jaculus jaculus	408	408	100%	2e-143	96.02%	201	XP_044986072.1
✓	PREDICTED: retinol-binding protein 4 [Dipodomys ordii]	Dipodomys ordii	402	402	100%	3e-141	95.02%	201	XP_012889769.1
✓	retinol-binding protein 4 isoform X1 [Mastomys coucha]	Mastomys coucha	400	400	100%	1e-140	94.53%	201	XP_031246073.1
✓	retinol-binding protein 4 precursor [Rattus norvegicus]	Rattus norvegicus	400	400	100%	3e-140	94.03%	201	NP_037294.1
✓	retinol-binding protein 4 [Arvicanthis niloticus]	Arvicanthis niloti...	400	400	100%	3e-140	94.53%	201	XP_034365668.1
✓	plasma retinol binding protein 4 precursor [Castor fiber]	Castor fiber	399	399	100%	4e-140	94.53%	201	APD32941.1
✓	retinol-binding protein 4 [Grammomys surdaster]	Grammomys sur...	399	399	100%	4e-140	94.03%	201	XP_028629352.1
✓	retinol-binding protein 4 [Sciurus carolinensis]	Sciurus caroline...	399	399	100%	7e-140	93.53%	201	XP_047410179.1
✓	retinol-binding protein 4 isoform X2 [Mastomys coucha]	Mastomys coucha	400	400	100%	9e-140	94.53%	246	XP_031246074.1
✓	retinol-binding protein 4 isoform 2 precursor [Mus musculus]	Mus musculus	398	398	100%	1e-139	93.53%	201	NP_035385.1

The top result is to a protein from Jaculus Jaculus (Lesser Egyptian jerboa), see second screen shot below for alignment details:

Range 1: 1 to 201

[GenPept](#)

[Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps	
422 bits(1085)	4e-149	Compositional matrix adjust.	201/201(100%)	201/201(100%)	0/201(0%)	
Query 1	MEWMWALVLLAAVGS	GRAERDCRVSSFRV	KENFDKARFSGTW	YAI	AKKDPEGLFLQDNII	60
Sbjct 1	MEWMWALVLLAAVGS	GRAERDCRVSSFRV	KENFDKARFSGTW	YAI	AKKDPEGLFLQDNII	60
Query 61	AEFSVDENGHMSATA	KGRVRLLSNWEVC	ADMVGTF	TDTEDPAKFKMKY	WGVASFLQRGND	120
Sbjct 61	AEFSVDENGHMSATA	KGRVRLLSNWEVC	ADMVGTF	TDTEDPAKFKMKY	WGVASFLQRGND	120
Query 121	DHWIIDTDYETFALQ	YSCLQNL	DGTCADSYSFV	SRDPNGLTPETR	KLVRQRQEELCLD	180
Sbjct 121	DHWIIDTDYETFALQ	YSCLQNL	DGTCADSYSFV	SRDPNGLTPETR	KLVRQRQEELCLD	180
Query 181	RQYRWIEHNGYCQSK	LSRNIL				201
Sbjct 181	RQYRWIEHNGYCQSK	LSRNIL				201

Download ▼

[GenPept](#)

[Graphics](#)

retinol-binding protein 4 [Jaculus jaculus]

Sequence ID: [XP_044986072.1](#) Length: 201 Number of Matches: 1

Range 1: 1 to 201

[GenPept](#)

[Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps	
408 bits(1048)	2e-143	Compositional matrix adjust.	193/201(96%)	198/201(98%)	0/201(0%)	
Query 1	MEWMWALVLLAAVGS	GRAERDCRVSSFRV	KENFDKARFSGTW	YAI	AKKDPEGLFLQDNII	60
Sbjct 1	MEWMWALVLLAA+G	SGRAERDCRVSSFRV	KENFDKARFSGTW	YAI	AKKDPEGLFLQDNII	60
Query 61	AEFSVDENGHMSATA	KGRVRLLSNWEVC	ADMVGTF	TDTEDPAKFKMKY	WGVASFLQRGND	120
Sbjct 61	AEF+VDENGHMSATA	KGRVRLLSNWEVC	ADMVGTF	TDTEDPAKFKMKY	WGVASFLQ+GND	120
Query 121	DHWIIDTDYETFALQ	YSCLQNL	DGTCADSYSFV	SRDPNGLTPETR	KLVRQRQEELCLD	180
Sbjct 121	DHWIIDTDY+T+ALQ	YSCL NLDGTCADSY	SFVSRDPNGL PETR	KLVRQRQEELCLD		180
Query 181	RQYRWIEHNGYCQSK	LSRNIL				201
Sbjct 181	RQYRWIEHNGYCQSK	L SGNIL				201

[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.

```
>RET4_HUMAN
MKWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDETGQMSATAKGRVR
LLNNWDVCADMVGTF'TD'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYD'TYAVQYSCRLNLDGTCADSYSFVFSRDPN
GLPPEAQKIVRQRQEELCLARQYRLIVHNGYCDGRSERNLL
>Salmonella_enterica
MEWMWALVLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAIKKDPEGLFLQDNIIAEFSVDENGHMSATAKGRVR
LLSNWEVCADMVGTF'TD'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYET'FALQYSCRLQNLGTCADSYSFVFSRDPN
GLTPETRKLVRQRQEELCLDRQYRWIEHNGYCQSKLSRNIL
>Hylobates_moloch
MEASLPQGGFLGKMKWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDE
TGQMSATAKGRVRLLNNWDVCADMVGTF'TD'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYD'TYAVQYSCRLNLDGTC
ADSYSFVFSRDPNGLPPEAQKIVRQRQEELCLARQYRLIVHNGYCDGRSERNLL
>Callithrix_jacchus
MQVSPAPPPRSF'TPRGYESATPSPRRYKAAGRPQARLPSSSTRARTLQPGLLAALLLVGVLLGKMKWVWALLLLAVLGIS
RAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIIAEFSVDETGQMSATAKGRVRLLNNWDVCADMVGTF'T
D'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYD'TYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQRIIRQRQEE
LCLARQYRLIVHNGYCDGKSERNLL
>Ursus_arctos_horribilis
MAWVWALVLLAVLGSARANRDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDENGQMSATAKGRVR
LLNNWDVCADMVGTF'TD'TEDSAKFKMKYWGVASFLQKGNDHWHIIDTDYD'TYAVQYSCRLNLDGTCADSYSFVFSRDPN
GLPPEAQKIVRQRQEELCLSRQYRLIVHNGYCDGRAEHSIL
>Phyllostomus_discolor
MEWVWALVLLAALGSARAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQDNIIAEFSVDENGQMSATAKGRVR
LLNNWEVCADMVGTF'TD'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYD'TYAVQYSCRLNLD
>Ovis_aries
MEWVWALVLLAALGSARAERDCRVSSFRVKENFDKARFAGTWYAMAKKDPEGLFLQDNIVAEFSVDENGH
MSATAKGRVRLLNNWDVCADMVGTF'TD'TEDPAKFKMKYWGVASFLQKGNDHWHIIDTDYET'YAVQYSCRL
LNLDGTCADSYSFVFARDPSGFSPEVQKIVRQRQEELCLARQYRLIPHNGYCDGKSERNIL
```


CLUSTAL multiple sequence alignment by MUSCLE (3.8)

```

Salmonella_enterica  -----
Callithrix_jacchus  MQVSPAPPPRSFTPRGYESATPSPRRYKAAGRPQARARLPSSSTRARTLQPGLLAALLLVGV
RET4_HUMAN          -----
Hylobates_moloch    -----MEASLPQGG
Ursus_arctos_horribilis -----
Phyllostomus_discolor -----
Ovis                -----

```

```

Salmonella_enterica  ----MEWMWALVLLAAVSGSRAERDCRVSSFRVKENFDKARFSGTWYAI AKKDPEGLFLQ
Callithrix_jacchus  LLGKMKWVWALLLLAVLGISRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ
RET4_HUMAN          ----MKVWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ
Hylobates_moloch    FLGKMKWVWALLLLAALGSGRAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ
Ursus_arctos_horribilis ----MAVWVWALVLLAVLGSARANRDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ
Phyllostomus_discolor ----MEWVWALVLLAALGSARAERDCRVSSFRVKENFDKARFSGTWYAMAKKDPEGLFLQ
Ovis                ----MEWVWALVLLAALGSARAERDCRVSSFRVKENFDKARFAGTWYAMAKKDPEGLFLQ
                      * *:***:***.:* .*:*****:*****:*****:*****

```

```

Salmonella_enterica  DNIIAEFSVDENGHMSATAKGRVRLLSNWEVCADMVGTFTDTPAKFKMKYWGVASFLQ
Callithrix_jacchus  DNIIAEFSVDETGQMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
RET4_HUMAN          DNIIAEFSVDETGQMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
Hylobates_moloch    DNIIAEFSVDETGQMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
Ursus_arctos_horribilis DNIIAEFSVDENGQMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
Phyllostomus_discolor DNIIAEFSVDENGQMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
Ovis                DNIIAEFSVDENGHMSATAKGRVRLLNNDVCADMVGTFTDTPAKFKMKYWGVASFLQ
                      ***:*****.:*:*****:*****.:***:*****:*****:*****

```

```

Salmonella_enterica  RGNDHWHIIDTDYETFALQYSCRLNLDGTCADSYSFVFSRDPNGLTPETRKLVRQRQEE
Callithrix_jacchus  KGNDHWHIIDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQRIIRQRQEE
RET4_HUMAN          KGNDHWHIVDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEE
Hylobates_moloch    KGNDHWHIIDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEE
Ursus_arctos_horribilis KGNDHWHIIDTDYDTYAVQYSCRLNLDGTCADSYSFVFSRDPNGLPPEAQKIVRQRQEE
Phyllostomus_discolor RGNDHWHIIDTDYDTYAVQYSCRLNLD-----
Ovis                KGNDHWHIIDTDYETYAVQYSCRLNLDGTCADSYSFVFARDPSGFSPEVQKIVRQRQEE
                      .*****:*****.:*:***** ***

```

```

Salmonella_enterica  LCLDRQYRWIEHNGYCSKLSRNIL
Callithrix_jacchus  LCLARQYRLIVHNGYCDGKSERNLL
RET4_HUMAN          LCLARQYRLIVHNGYCDGRSERNLL
Hylobates_moloch    LCLARQYRLIVHNGYCDGRSERNLL
Ursus_arctos_horribilis LCLSRQYRLIVHNGYCDGRAEHSIL
Phyllostomus_discolor -----
Ovis                LCLARQYRLIPHNGYCDGKSERNLL

```

[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.



Salmonella_enterica 0.09097
Phyllostomus_discolor -0.01796
Ovis 0.03645
Ursus_arctos_horribilis 0.03865
Callithrix_jacchus 0.03587
RET4_HUMAN -0.00017
Hylobates_moloch 0.00516

[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.

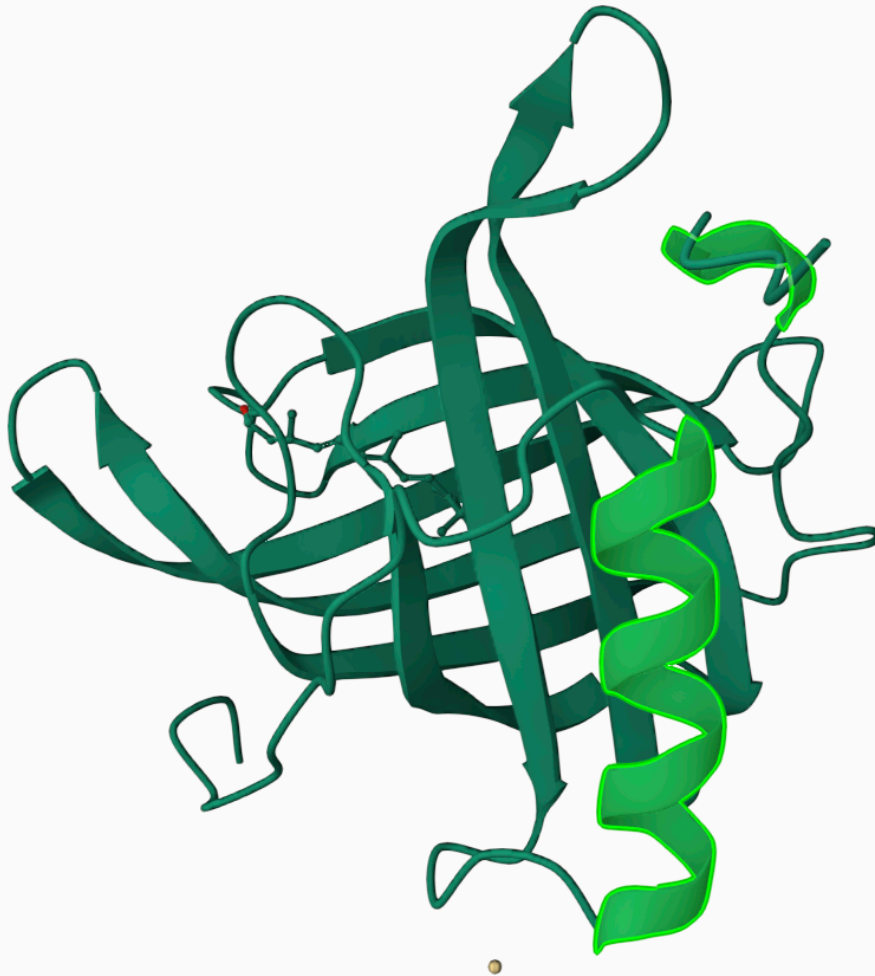
List the top 3 *unique* hits (i.e. not hits representing different chains from the same structure) along with their Evalue and sequence identity to your query. Please also add annotation details of these structures. For example include the annotation terms PDB identifier (structureId), Method used to 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 choose 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	EVALUE	Identity
1ERB	X-RAY	1.9	Bos taurus	5.661E-105	91%
1AQB	X-RAY	1.65	Sus scrofa domesticus	5.894E-118	93%
1IIU	X-RAY	2.5	Gallus gallus	6.068E-105	87%

[Q9] Generate a molecular figure of one of your identified PDB structures using the **NGL viewer** online (or **VMD/PyMol**). You can optionally highlight conserved residues that are likely to be functional. Please use a white or transparent background for your figure (i.e. not the default black).



Based on sequence similarity. How likely is this structure to be similar to your “novel” protein?

It is very likely that 1AQB is to be similar in structure to *Salmonella_enterica* given the high sequence similarity (>90%). In the figure below the alpha helix colored green is highly conserved and corresponds to the same structure seen in *Salmonella_enterica*.

[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?

CHEMBL details 20 Binding (CHEMBL3100) and 3 Functional Assays;

The Ligand Efficiency chart plots Binding Efficiency Index (BEI) against Surface Efficiency Index (SEI), where:

- **SEI** = $(-\log_{10}(\text{Standard Value} \times 10^{-9})) \times 100 / \text{PSA}$
- **BEI** = $(-\log_{10}(\text{Standard Value} \times 10^{-9})) \times 1000 / \text{MWT}$

https://www.ebi.ac.uk/chembl/target_report_card/CHEMBL3100/

Antagonists of retinol-binding protein 4 (RBP4) stop the eyes from taking in serum all-trans retinol (1) and have been shown to stop the formation of cytotoxic bisretinoid in the retinal pigment epithelium (RPE), which is linked to both dry age-related macular degeneration (AMD) and Stargardt disease. The bicyclic [3.3.0]-octahydrocyclopenta[c]pyrrolo analogue 4 is an example of a new line of nonretinoid RBP4 inhibitors that was found. Pyrimidine-4-carboxylic acid fragment was used as a good isostere for the anthranilic acid part of molecule 4, which led to the finding of the star blocker molecule 33.

Cioffi CL, Racz B, Freeman EE, Conlon MP, Chen P, Stafford DG, Schwarz DM, Zhu L, Kitchen DB, Barnes KD, Dobri N, Michelotti E, Cywin CL, Martin WH, Pearson PG, Johnson G, Petrukhin K. Bicyclic [3.3.0]-Octahydrocyclopenta[c]pyrrolo Antagonists of Retinol Binding Protein 4: Potential Treatment of Atrophic Age-Related Macular Degeneration and Stargardt Disease. *J Med Chem* (2015) 58:5863-5888 [10.1021/acs.jmedchem.5b00423](https://doi.org/10.1021/acs.jmedchem.5b00423)

Scoring Rubric:

[45 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 (10 points)

PDB identifiers from multiple species reported 5

Annotation of PDB source, resolution and technique 4

Annotation of Evalue and Sequence Identity 1

Q9 (4 points)

Structure figure provided 2

Uses white background for molecular figure 1

Figure of high resolution (i.e. not just snapshot) 1

Q10 (1 point)

Evidence of ChEMBL searches 1