

## BIMM-143: INTRODUCTION TO BIOINFORMATICS

The find-a-gene project assignment  
[https://bioboot.github.io/bimm143\\_S20/](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.**

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:** **neurexin-1 isoform alpha4 precursor**

**Accession:** NP\_001317007

**Organism:** 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 (BLOSUM62) against crustaceans ESTs (expect threshold = 10)

**Database:** Expressed sequence tags (est)

**Organism:** Crustaceans (Taxid: 6657)

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.

blastn   blastp   blastx   **tblastn**   tblastx

Enter Query Sequence

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

NP\_001317007.1

Query subrange

From

To

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

Job Title

NP\_001317007:neurexin-1 isoform alpha4 precursor...

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

☐ Align two or more sequences [?](#)

Choose Search Set

Database  [?](#)

Organism

Optional  ☐ 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  [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

Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign

Algorithm parameters

General Parameters

Max target sequences  [?](#)

Select the maximum number of aligned sequences to display [?](#)

Expect threshold  [?](#)

Word size  [?](#)

Max matches in a query range  [?](#)

Scoring Parameters

Matrix  [?](#)

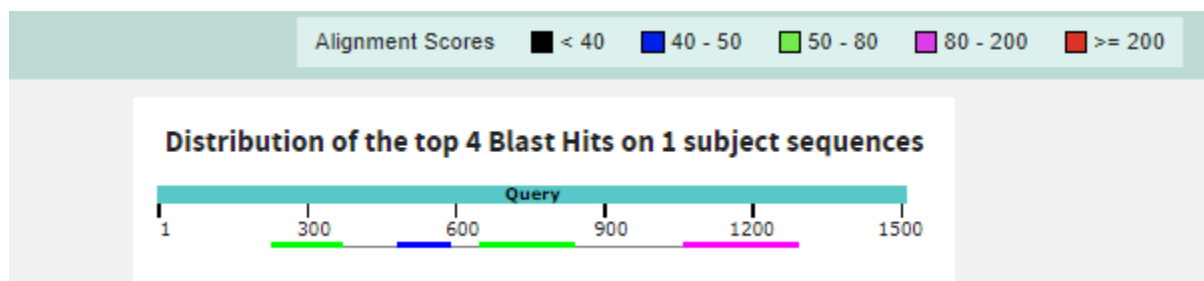
Gap Costs  [?](#)

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.

Accession number: [CI998089.1](#)

cDNA clone from immature *Penaeus japonicus*, 635 bp, mRNA sequence

Query: neurexin-1 isoform alpha4 precursor [Homo sapiens] Query ID: NP\_001317007.1 Length: 1507



>CI998089 Marsupenaeus japonicus eyestalk immature Penaeus japonicus cDNA clone  
YAQ01A02NGRM0003\_C10 5', mRNA sequence

Sequence ID: CI998089.1 Length: 635

Range 1: 28 to 633

Score:188 bits(477), Expect:6e-53,

Method:Compositional matrix adjust.,

Identities:96/233(41%), Positives:138/233(59%), Gaps:32/233(13%)

Query 1056 VDLNGRLPDLI-

SDALFCNGQIERGCEGPSTTCQEDSCSNQGVCLQQWDGFSCDCSMTSF 1114

+DLNG PD D + + GC+GPST C ++C+N G C+QQW+ +SC+C MTSF

Sbjct 28

LDLNGEAPDPAGKDVPLISSPVFPGCDGPSTKCHRNACANGGTCVQQWNSYSCNCD

MTSF 207

Query 1115

SGPLCNDPGTTYIFSKGGGQITYKWPPNDRPSTRADRLAIGFSTVQKEAVLVRVDSSS

GL 1174

+GP C+D T Y F G G +T+++P P TR D LA+GF T Q++AV++R+DS++

Sbjct 208

TGPTCSDESTAYEFGGGSGGLMTFQYPEGRWPDTRRDLLALGFMTSQEDAVMLRLDS

ANS- 384

Query 1175

GDYLELHIHQGKIGVKFNVTDDIAIEESNAIINDGKYHVVRFRTRSGGNATLQVDSWPVI

1234

DY+EL I G I + +N+GT+D + E A +NDG YHVVRF RSG NAT+Q+D + V

Sbjct 385

NDYMELEIVDGNIFMVYNMGTEDHPVGEVMAKVNDGIYHVVRFVRSGSNATVQIDDYE

VR 564

Query 1235

ERYPAGNNDNERLAIARQRIPYRLGRVVDEWLLDKGRQLTIFNSQATIIIGGK 1287

+ P KG QL++FN+Q+ + IGG+

Sbjct 565 SKNP-----KGHQLSVFNAQSRLQIGGR 633

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.

```
>CI998089.1 CI998089 Marsupinaeus japonicus eyestalk immature Penaeus
japonicus cDNA clone YAQ01A02NGRM0003_C10 5', mRNA sequence
GAATTCCCGGGTTCGACCCACGCGTCCGCTCGACTTGAACGGCGAAGCTCCCGACCCAGCAGGCAAAGACG
TTCCTCTGATCTCGAGCCCCGTCTTCCCAGGCTGTGACGGACCAAGTACAAAATGCCATCGGAATGCATG
TGCAAATGGTGGGACGTGCGTGCAACAGTGGAAATTCCTATTTCATGCAACTGTGACATGACCTCCTTTACC
GGACCTACGTGTTCTGATGAGTCCACGGCGTATGAGTTCGGCGGGCGGCTCTGGCCTCATGACCTTCCAGT
ACCCTGAGGGCAGATGGCCTGATACCCGCCGTGACCTGTTGGCACTCGGTTTTATGACCTCGCAGGAGGA
TGCGGTCATGCTGAGGCTGGATTTCGGCCAATTTCGAATGATTACATGGAATTGGAAATCGTCGACGGCAAT
ATCTTCATGGTCTACAACATGGGCACGGAGGACCACCCGTAGGCGAGGTCATGGCAAAGGTCAACGACG
GCATATATCACGTGGTTTCGTTTCGTTAGGTCAGGTTCCAACGCCACGGTACAGATTGACGACTATGAAGT
TCGTTCAAAGAATCCAAAAGGACACCAGCTCTCCGTCTTCAACGCACAGTCACGCCTCCAGATCGGTGGA
CGAAG
```

```
>CI998089.1_1 Marsupinaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
EFPGRPTRPLDLNGEAPDPAGKDVPLISSPVFPGCDGPSTKCHRNACANGGTCVQQWNSY
SCNCDMTSFTGPTCSDESTAYEFGGGSLMTFQYPEGRWPDTRRDLLALGFMTSQEDAVM
LRLDSANSNDYMELEIVDGNIFMVYNMGTEHPVGEVMAKVNDGIYHVVRFVRSNATV
QIDDIYEVRSKNPKGHQLSVFNAQSRLQIGGRX
>CI998089.1_2 Marsupinaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
NSRVDPRVRST*TAKLPTQQAKTFL*SRAPSSQAVTDQVQNAIGMHVQMVGRACNSGIPI
HATVT*PPLPDLRLVMSPRRMSSAAALAS*PSSTLRADGLIPAVTCWHSVL*PRRRMRSC
```

```

*GWIRPIRMITWNWKSSTAISWSTTWARRTP*ARSWQRSTTAYITWFASLGQVPTPRY
RLTTMKFVQRIQKDTSSPSSTHSHASRSVDEX
>CI998089.1_3 Marsupenaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
IPGSTHASARLERRSSRPSRQRRSSDLEPRLPRL*RTKYKMPSECMCKWWDVRATVEFLF
MQL*HDLLYRTYVF**VHGV*VRRRLWPHDLVPV*GQMA*YPP*PVGTRFYDLAGGCGHA
EAGFGQFE*LHGIGNRRRQYLHGLQHGHGGPPRRRGHGKGQRRHISRGSLR*VRFQRHGT
D*RL*SSFKESKRTPALRLQRTVTTPDRWTK
>CI998089.1_4 Marsupenaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
SSTDLEA*LCVEDGELVSFWIL*TNFIVVNLYRGVGT*PNEANHVIYAVVDLCHDLAYGV
VLRAHVVDHEDIAVDDFQFHVIIRIGRIQPQHDRILLRGHKTECQQVTAGIRPSALRVLE
GHEARAAELIRRLIRTRRSKGKGGHVTVA*IGIPLLHARPTICTCIPMAFCTWSVTAW
DGARDQRNVFACWVGSFAVQVERTRGSTREF
>CI998089.1_5 Marsupenaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
FVHRSGGVTVR*RRRAGVLLDSLNELHSRQSVPRWNLT*RSEPRDICRR*PLP*PRLRG
GPPCPCCRP*RYCRRRFPIPCNHSNWPNPASA*PHPPARS*NRVPTGHGGYQAICPQGTG
RS*GQSRRTHTPTWTHQNT*VR*RRSCHSCMNRNSTVARTSHHLMHSDGILYLVRHSLG
RRGSRSEERLCLLGRELRRSSRADAWVDPGIX
>CI998089.1_6 Marsupenaeus japonicus eyestalk immature Penaeus japonicus cDNA
clone YAQ01A02NGRM0003_C10 5', mRNA sequence
LRPPIWRRDCALKTESWC PF GFFERTS*SSICTVALEPDLTKRTT*YMPSLTFAMTSPTG
WSSVPML*TMKILPSTISNSM*SFELAESSLSMTASSCEVIKPSANRSRRVSGHLPSGYW
KVMRPEPPPNYSYAVDSSEHVGPVKEVMSQLHE*EFHCCTHVPPFAHAFRWHFVLGPSQPG
KTGLEIRGTS L P A G S G A S P F K S S G R V G R P G N S

```

Name: *Penaeus neurexin*

CI998089

Species: [Penaeus japonicus](#)

Eukaryota; Metazoa; Ecdysozoa; Arthropoda; Crustacea;  
 Multicrustacea; Malacostraca; Eumalacostraca; Eucarida; Decapoda;  
 Dendrobranchiata; Penaeoidea; Penaeidae; Penaeus.

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.

**[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.

BLAST® » blastn suite

Standard Nucleotide

blastnblasttblastntblastx

BLASTN programs search nucleotide data

Enter Query Sequence

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

dbj|C1998089

Query subrange [?](#)

From

To

Or, upload file

Choose FileNo file chosen [?](#)

Job Title

C1998089:C1998089 Marsupenaeus japonicus eyestalk...

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

☐ Align two or more sequences [?](#)

Choose Search Set

Database

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

Nucleotide collection (nr/nt) [?](#)

Organism

Optional

Enter organism name or id--completions will be suggested

☐ exclude [Add organism](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

Exclude

Optional

☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences

Limit to

Optional

☐ Sequences from type material

Entrez Query

Optional

[YouTube](#) [Create custom database](#)

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

Program Selection

Optimize for

☒ Highly similar sequences (megablast)
 ☐ More dissimilar sequences (discontiguous megablast)
 ☐ Somewhat similar sequences (blastn)

Choose a BLAST algorithm [?](#)

BLAST

Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)

☐ Show results in a new window

+ Algorithm parameters

The top result is to a protein from *Penaeus Japonicus* (caridean shrimp), see second



screen shot below for alignment details:

Alignment view Pairwise ☐ CDS feature [Restore defaults](#) Download

15 sequences selected

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

**PREDICTED: *Penaeus japonicus* neurexin-1-like (LOC122254534), transcript variant X3, mRNA**  
Sequence ID: [XM\\_043018271.1](#) Length: 7590 Number of Matches: 1

Range 1: 3407 to 4021 [GenBank](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Identities	Gaps	Strand
1125 bits(609)	0.0	614/616(99%)	1/616(0%)	Plus/Plus
Query 20	CGCGTCCGCTCGACTTGAACGGCGAAGCTCCCGACCCAGCAGGCAAGACGTTCTCTGA	79		
Sbjct 3407	CGCGT~CGCTCGACTTGAACGGCGAAGCTCCCGACCCAGCAGGCAAGACGTTCTCTGA	3465		
Query 80	TCTCGAGCCCCGTCTTCCCAGGCTGTGACGGACCAAGTACAAAATGCCATCGGAATGCAT	139		
Sbjct 3466	TCTCGAGCCCCGTCTTCCCAGGCTGTGACGGACCAAGTACAAAATGCCATCGGAATGCAT	3525		
Query 140	GTGCAAAATGGTGGGACGTGCGTGCAACAGTGGAAATTCCTATTGCAACTGTGACATGA	199		
Sbjct 3526	GTGCAAAATGGTGGGACGTGCGTGCAACAGTGGAAATTCCTATTGCAACTGTGACATGA	3585		
Query 200	CCTCCCTTACCGACCTACGTGTTCTGATGAGTCCACGGCGTATGAGTTCGGCGGCGCT	259		
Sbjct 3586	CCTCCCTTACCGACCTACGTGTTCTGATGAGTCCACGGCGTATGAGTTCGGCGGCGCT	3645		
Query 260	CTGGCCTCATGACCTTCCAGTACCTGAGGGCAGATGGCTGATACCCCGGTGACCTGT	319		
Sbjct 3646	CTGGCCTCATGACCTTCCAGTACCTGAGGGCAGATGGCTGATACCCCGGTGACCTGT	3705		
Query 320	TGGCACTCGGTTTTATGACCTCGCAGGAGGATGCGGTATGCTGAGGCTGGATTGCGCCA	379		
Sbjct 3706	TGGCACTCGGTTTTATGACCTCGCAGGAGGATGCGGTATGCTGAGGCTGGATTGCGCCA	3765		
Query 380	ATTGCAATGATTACATGGAATTGGAAATCGTCGACGGCAATATCTTCATGGTCTACACA	439		
Sbjct 3766	ATTGCAATGATTACATGGAATTGGAAATCGTCGACGGCAATATCTTCATGGTCTACACA	3825		
Query 440	TGGGCACGGAGGACACCCCGTAGGCGAGGTATGGCAAGGTCAACGACGGCATATATC	499		
Sbjct 3826	TGGGCACGGAGGACACCCCGTAGGCGAGGTATGGCAAGGTCAACGACGGCATATATC	3885		
Query 500	ACGTGGTTCGCTTCGTTAGGTGAGGTTCAACGCCACGGTACAGATTGACGACTATGAAG	559		
Sbjct 3886	ACGTGGTTCGCTTCGTTAGGTGAGGTTCAACGCCACGGTACAGATTGACGACTATGAAG	3945		
Query 560	TTGTTTCAAAGAAATCCAAAAGGACACCAAGCTCTCCGCTTCAACGCACAGTCACGCTCC	619		
Sbjct 3946	TTGTTTCAAAGAAATCCAAAAGGACACCAAGCTCTCCGCTTCAACGCACAGTCACGCTCC	4005		
Query 620	AGATCGGTGGACGAAG	635		
Sbjct 4006	AGATCGGTGGACGAAG	4021		

**[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.

**[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.

**[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.

**[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.

**[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?

**[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?

**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
<b>Q4</b> (3 point)	
Blastp output list with identities & Evalue	1
Top alignment shown with alignment statistics	1
Results indicates a “novel” gene found	1
<b>Q5</b> (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
<b>Q6</b> (1 point)	
Figure illustrates sequence clustering pattern	1
<b>Q7</b> (10 points)	
Heatmap figure included in report	5
Heatmap is legible (i.e. no labels obscured)	5
<b>Q8</b> (10 points)	
PDB identifiers from multiple species reported	5
Annotation of PDB source, resolution and technique	4
Annotation of Evalue and Sequence Identity	1
<b>Q9</b> (4 points)	
Structure figure provided	2
Uses white background for molecular figure	1
Figure of high resolution (i.e. not just snapshot)	1
<b>Q10</b> (1 point)	
Evidence of ChEMBEL searches	1