

BGGN-213: FOUNDATIONS OF BIOINFORMATICS

The find-a-gene project assignment

<http://thegrantlab.org/bggn213>

Dr. Barry Grant

Juan Gonzalez; jhg001@ucsd.edu; PID: A69036681

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: Longifolia 1 (LNG1)

Accession: NP_197062.1

Species: Arabidopsis thaliana

Function: Regulates leaf morphology by promoting cell expansion in the leaf-length direction

[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).

LNG1 Arabidopsis thaliana DNA sequence:

ATGTCGGCAAGCTTGTATAACTTGTGGATGAGAATCAAATCTGAATAACAGATGGATGTATGAATGGGATCTTCAGGTGTTTACCGGAAACATTATCCACCGAGACGTGTCACC
GGAGATGAGCTCAAGCTCTCCCTCAGGAAAGCAAGTGCACATGGTGGATCACCAACATTCTCAGGGACAAGAGAAAAGGAGAAGAGACTGCAAAAGGAGA
AACAGAGGGAGATCTCTCATGCGCTCGAGGTGTTCTCATCACCATGCTCCAGCTTCATCTCAGATATTAGCACCCGCTCTCAGTTAACACGCCGGTTTG
AGTAATGAGTGAAGAACGACCAAGGGCTTAAATGATGCCAAGTGTATGGCTTCAAGGAGATATAAGGGAGCTTGTGAGAAGCTTATTCTATAAGGAGAACAGAAAC
AGAGATGAAGAACGCTTGTCTCAGGCTTAATCAGGAGCTAATGTTCTCTCTCAAAAGAATCATCACCATCTCGGAATTCTTAATGAATGGAGTGGAGGAGAGTAGTGAG
CTGAAAAGACAGTCTCCGTTCTCTACGATGAGAGGAGACGAGAAAGACAGGGGCCAAGTGTGAAAGAGACCCCCGAGGGTTGTCATTAGACAGTAGATCGAATTCTCTAGGAGC
AGTAGTTGAGCTTCAACAGAGACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
ATCTAGGAGCTTCAACAGAGACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
AGAGAAATCGCTCTCGCACTCCCAGGGCAACTTCCCGAGTGGAAAGTAGTACAAGATCAAAGATGAGGTTTGTATCAATCAAAGATGATGCCCTGAAAG
CTCAGTGGAAAGAGCAGTGTGGTGGCCAAGAACCAACTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
AACACCTCAGGGCTCTTAAAGCAAAACTCTGAAGCAAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
AGCAAACTTCCATCTGCAATAAACACCTGTCATGAATTAAATCATCTCTATCGTGTGTTGAAAGCAGCTACCGCTCCAGTCTCAAAGACACAGGCATTCGAGGTTCTG
CCCGCGGAATGTTGCTTACCAAAATGTCAGGTTGAAAGGAAACCTGAGGCAAGGGCAGAAAGCTTCCGGAGGAAAGCAGAGTGTGACCCGAGGGCAGGATATTCAAGGGCC
AGACAAATCCACAATGAAAACACTAGTACCCAGACCTTAACTCGAAAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
GGTTTGGAGAAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
AAATCTGTGGCTTACAGCAATCTGAAGACGGTTAACAGTGTGAAAGCAGTGTGAAAGTCTAAGATCTGACAGCAAGCTGGGCTCTAACCTTGATACTGAGGTTAACAGCA
GGTATAACTATGAGGAGAACAGGCACATTACGGAGCAGCACACCCCGAACAAACAGGAGTGGCAGACTTGGGAATGAGGTCGCTGCAAACACTCTGAAAGTTAACAGTGGAGGCCAG
CCCGTGTAGITCTGAGTGTAGCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAGGAGATCTGAG
AGACAAACAACAACTTATGAGTGTGGCTTGGCTGAGGAGTAACAGCAGCTAACAGCAGCTAACAGCAGCTAACAGCAGCTAACAGCAGCTAACAGCAGCTAACAGCAG
ACAAGTACTCTCGAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
ACAAGAACAG
CCATCAACGAGACTTAGCTCACAGATTGCTGAGAAGGGTGTACAAAGCAACCACATCATAACATTACATCAGCACGCAAGGACACAGAACAGAACAGAACAGAACAG
TGAACAAACTCTGTGAGGAGATCTGAGTCAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
TGAAGGAGAGAACACGGGGTTAGTGTAGACATTGAGGAGGCTTAATCTCAAAGAGCTTGTGAGTGGAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAGCTAACAGGAG
CGAGCTTTTCTGCTAA

LNG1 Arabidopsis thaliana protein sequence:

MSAKLLYNSDENPNLNKQIGCMNGIQFVYRQHYPPRRVTGDELKSLPSKGASDNVGDTNISADKKTEKSKKKTAKEKQRGSSESSRLSFSSPCSSSSADSINTTASQFEQPLSNSGE
NPVREPTNGSPRWGLMPSDIRELVRSSIHKETRTRDEEALSQQPKSARANVSLKESSPNSNEWSEGRVVKLDSPRSYDERETRKTGAKLKETPRLSDRSNSRSARSSCPPEQ
ELVTGHRRRTSSVAKLMLGLEVPDEPTVQNRNRFCDSPRPTSRVEVDLQRSGRFDISKMMPAKFPMKASPWAQVGDGAKNVKIPDATTLYVEIJKRQLSLEFKKSEKDRLALKQILEAME
KTQQLISKDDDNKTLCSNSFMQRNNQPSPSAINTSSMFKSSSIIVVMKAATAPVFKDTGIASSAFSPRNVALPVKGNLROAQKVIPRKQSAMDVTPRGYKKQTESTMKNTSTRPLQSKS
DMAKSGKIQKPSVSLRTPKKLGFEKQSRPTSPKPELNKNQRQQLSRQQTESASPRRKPGIKSRLQSQSEDRLSDESSLRSRSDNSVNLASNLDETEVTSRYNERNSDITEQHTPKQRSPDL
GMRSLSKPLKVTVEQSPSPVSLDAFEDDSPSVPRKIVSFKEEDDNLSSEESHWMNKNNLCSRIVWESNTSLKQPDALTEGFMEDDAEFKNGDHKYISEIMLASGLLRDIDYSMISIQLHQAH
LPINPSSLFFVLEQNKTSNVLQDNKHKGRRGFGQQQTVNLVERSKRKLIFDTEILHRAFAEGCTKQPSLTSISTQRTHEKSSRGEELLQTLCEIDRLQDNSKCILDEDDEDLIWEDLQSHGMNW
KEIEGETPGLVDIERLIFKDLIGEVTFSEFAFPRLMSGQPRQLFHC

Method: blastn (to use protein sequence to search nucleotide databases)
Database: Core nucleotide database (core_nt)
Organism: Brassica rapa (taxid:3711); Capsella rubella (taxid:81985)

Search Screenshot:

The screenshot shows the NCBI BLAST search interface. The 'tblastn' tab is active. The 'Enter Query Sequence' field contains the accession number NP_197062.1. The 'Choose Search Set' section is configured to search the 'Core nucleotide database (core_nt)' for 'Brassica rapa (taxid:3711)'. The results table below shows 34 significant alignments.

Blast Result Screenshot:

Descriptions	Graphic Summary	Alignments	Taxonomy	Download	Select columns	Show 100	?		
Sequences producing significant alignments									
<input checked="" type="checkbox"/> select all 34 sequences selected									
	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	PREDICTED: Capsella rubella protein LONGIFOLIA 1 (LOC17881062).mRNA	Capsella rub...	1472	1472	100%	0.0	79.29%	3129	XM_006286950.2
<input checked="" type="checkbox"/>	PREDICTED: Brassica rapa protein LONGIFOLIA 1 (LOC103846402).mRNA	Brassica rapa	1194	1194	100%	0.0	69.56%	2953	XM_009123315.3
<input checked="" type="checkbox"/>	PREDICTED: Brassica rapa protein LONGIFOLIA 1 (LOC103850992).mRNA	Brassica rapa	1116	1116	100%	0.0	66.60%	8739	XM_009127809.3
<input checked="" type="checkbox"/>	Brassica rapa genome_scaffold_A02	Brassica rapa	994	1395	92%	0.0	64.98%	31546239	LR031573.1
<input checked="" type="checkbox"/>	PREDICTED: Capsella rubella protein LONGIFOLIA 2 (LOC17892073).mRNA	Capsella rub...	835	835	98%	0.0	52.92%	3029	XM_006296884.2
<input checked="" type="checkbox"/>	PREDICTED: Brassica rapa protein LONGIFOLIA 2 (LOC103858943).mRNA	Brassica rapa	775	775	98%	0.0	51.86%	3239	XM_009136402.3
<input checked="" type="checkbox"/>	PREDICTED: Brassica rapa protein LONGIFOLIA 2 (LOC103870931).transcript...	Brassica rapa	710	710	98%	0.0	49.68%	2750	XM_009149128.3
<input checked="" type="checkbox"/>	PREDICTED: Brassica rapa protein LONGIFOLIA 2 (LOC103870931).transcript...	Brassica rapa	705	705	98%	0.0	49.52%	2759	XM_009149127.3
<input checked="" type="checkbox"/>	Brassica rapa genome_scaffold_A10	Brassica rapa	616	1199	93%	0.0	69.74%	17648834	LR031577.1
<input checked="" type="checkbox"/>	Brassica rapa subsp. pekinensis cultivar Inbred line "Chiifu" clone KBrB020F06_c...	Brassica rap...	313	724	92%	1e-168	46.24%	135946	AC232459.1
<input checked="" type="checkbox"/>	Brassica rapa genome_scaffold_A05	Brassica rapa	304	709	91%	2e-164	46.09%	47572232	LR031570.1

[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 of “novel” protein:

```
>B. rapa protein LONGIFOLIA 1 (sequence taken from BLAST result)
MSAKLLYNLSDENPNLNKQFGCMNGIFQVFYRQHYPARRVSAGDELKSLPSGKTSDNVGVTNGSTDKKETEKSKKK
KAAKEKQKVSVSESSSRLSFSSPCSSSFSSADISTTSQFEQPMNSGETPAREPTYGSPRWGGLVMSSDLRELVRS
SIHKETRTRVEEALSQQPKSARANVSLLKELSPRSSNEWSEGRRVVKLKSPRFSYDEREARKTGAKFKETPRLS
LDSRSNSFRSAKSSCPEPQELVTGHRRRTSSVIAKLMLGVDVVSDEPVTDQSRHENFCDSPRPAPRVEADLPRSRGS
DSFKKMMPAAKFPAKTAPWTQADGARNQVKAADAAAATLTVYGEIQRQLSQLEFKKSEKDLRALQQILEAMEKTQQLM
SKDDDNSSLSSTNFMQPSPKSIRSSSIVVMKAASAPVFKETGSSSTSSSPRSVALPNVKVSNQKGTRKQSAMD
VTPRPATKNTSTRPLQSKIEMAKSGKPSVSPRTQPKKLGFQSRPTSPKPEPNKNQRQQLSRQOTESPSPRKPGM
KSRLGLQQSEDRSSDESSLRSLRSDNSVSSASNFDIEVTSRHKCDLTEQHTPKQRSPELGMRSLPKPLKITVEQPSP
VSILDVAFDDDESPPVRKISIVFKDDDHIRSEESIWMKKHNNLCRSIVWPESNTSLNQPAVLTESFMEEGADLRN
GDRKYISEILSASGLLDIDYSMLSISQLHQAHLPINPSLFFVLEQNKTNSVTHRGRGFGQQTANLIGRSRRKLFDT
VNEILARKFAAEGCTKQPYITSSISPLMKTDKSSRGKELLEALCSEIDRLQDNSNCILDEDDEDLIWEDLQSQGMNW
KEIEGETPGLVLDIERLIFKDLISEVVTSEVAAPGNKLSGQPRQLFH
```

Name: protein LONGIFOLIA 1

Species: *Brassica rapa*

**Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
Spermatophyta; Magnoliopsida; eudicotyledons; Gunneridae; Pentapetalae;
rosids; malvids; Brassicales; Brassicaceae; Brassiceae; Brassica**

[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 the top match reported has less than 100% identity, then it is likely that your protein is novel, and you have succeeded.

Chosen match: Accession ID: XM_009123315.3 (see below for alignment details)

Query Coverage: 100%, E-value: 0.0: Percent Identity: 69.56%

PREDICTED: *Brassica rapa* protein LONGIFOLIA 1 (LOC103846402), mRNA

Sequence ID: XM_009123315.3 Length: 2953 Number of Matches: 1

Range 1: 86 to 2773	GenBank	Graphics	▼ Next Match	▲ Previous Match		
Score	Expect	Method	Distance	Position	GenBank	Frame
1194 bits(308B)	0.0	Conventional matrix adjust:	746/933(0%)	800/933(0%)	43/933(4%)	+2
Query 1	MSAKLLYLNIDENPNLANKIDCNGCIGFOVYRQHNPYPRVNT—GDEKLKLPGLPSASDNNNG					
Sbjct 86	MSAKLLYLNIDENPNLANKIDCNGCIGFOVYRQHNPYPRVNT—GDEKLKLPGLPSASDNNNG					
Query 59	DNTNIsadkKETEKKKKKKKAKEKQVSESSSRLsffsspcssssfsadittas0F0Q					
Sbjct 266	DNTNIsadkKETEKKKKKKKAKEKQVSESSSRLsffsspcssssfsadittas0F0Q					
Query 119	PLGSNGENPVEEPTNGSPRMGLIMPQDIRELVNSLHKEITRTR-DEEALSQDQPKSARAN					
Sbjct 444	PLGSNGENPVEEPTNGSPRMGLIMPQDIRELVNSLHKEITRTR-DEEALSQDQPKSARAN					
Query 178	VSLKKESSPNSNEMESEGRRVVKLQKOSPRFSDYDERETTKGAKLKETPRLSDLsnsf					
Sbjct 623	VSLKKESSPNSNEMESEGRRVVKLQKOSPRFSDYDERETTKGAKLKETPRLSDLsnsf					
Query 238	PSKSSCPEPEQELVTGHRTTSSVVAKLNGLLEVIPDEFUTVNRNCDSPRPTSRVE					
Sbjct 883	PSKSSCPEPEQELVTGHRTTSSVVAKLNGLLEVIPDEFUTVNRNCDSPRPTSRVE					
Query 294	VDLORSRGFDSTLQKQKAF-mb-SPWAQVQGAKNQV/PIPADT-ITLVYGEIOKRISLQL					
Sbjct 983	VDLORSRGFDSTLQKQKAF-mb-SPWAQVQGAKNQV/PIPADT-ITLVYGEIOKRISLQL					
Query 356	EFKKSEKOLRAKQILEMEKTQQLSKQDDDNNTLSSSNMORNNP0IPSAINTS9WIF					
Sbjct 1163	EFKKSEKOLRAKQILEMEKTQQLSKQDDDNNTLSSSNMORNNP0IPSAINTS9WIF					
Query 416	KSSSIVMKAAATAPVFTGAGASASTPSPRNVALPNVKVQNLRAOAKVTPKQSMADVP					
Sbjct 1313	KSSSIVMKAAATAPVFTGAGASASTPSPRNVALPNVKVQNLRAOAKVTPKQSMADVP					
Query 476	RPGYYKGQTESTMWNTSTRPLQKPSDMSAKSGC10KSVSLRTPPKLGFEKQSPRTSPKP					
Sbjct 1481	RPA-----RPTNSTRPLQKPSDMSAKSGC10KSVSLRTPPKLGFEKQSPRTSPKP					
Query 536	ELNKNQRQLSRQTEASPRRKPGKSRGLQ0edrlsdeddsrlsrls-nVSLASN					
Sbjct 1625	ELNKNQRQLSRQTEASPRRKPGKSRGLQ0edrlsdeddsrlsrls-nVSLASN					
Query 596	DTEVTSRNYEINSDTIEHTHPKRSRDLGMSLSPKLPKVITVEDPSPVSPVLSVDAFEDDS					
Sbjct 1805	DTEVTSRNYEINSDTIEHTHPKRSRDLGMSLSPKLPKVITVEDPSPVSPVLSVDAFEDDS					
Query 656	PSPVKRCIS1VPEKEDDNLSSEEHMMNNQNLCRS1VPEKEDDNLSSEHSKQKAELETEGFMEDDA					
Sbjct 1973	PSPVKRCIS1VPEKODDHIRSEESLMMKXHNLCRS1VPEKEDDNLSSEHSKQKAELETEGFMEDDA					
Query 716	EFKGNQHYTSEIILASGLRBDITYMS10HOAH1PTNPSLFLV1EONKTSVSLQDN					
Sbjct 2153	EFKGNQHYTSEIILASGLRBDITYMS10HOAH1PTNPSLFLV1EONKTSVSLQDN					
Query 776	KKHGRGFQGQOTVNLVYKRRKLFDTIMEILAMHAAEGCTKOPSTILSISTORTHEKS					
Sbjct 2321	KKHGRGFQGQOTVNLVYKRRKLFDTIMEILAMHAAEGCTKOPSTILSISTORTHEKS					
Query 836	SRGEFELLQTLSETEDRLOONSKC1dedded1lwend1OSHGNWKEITEGETPGVLVDIER					
Sbjct 2495	SRGEFELLQTLSETEDRLOONSKC1dedded1lwend1OSHGNWKEITEGETPGVLVDIER					
Query 896	LIFKQD1LWVWVFAFP-RNQ-EGCPDFC_HC 927					
Sbjct 2675	LIFKQD1LWVWVFAFP-RNQ-EGCPDFC_HC 927					

To verify if I have identified a novel gene, I will run a BLASTp search against the non-redundant (nr) database. This will allow me to determine if there are any matches and confirm the novelty. 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

Descriptions	Graphic Summary	Alignments	Taxonomy	Download	Select columns	Show	100	?
Sequences producing significant alignments								
<input checked="" type="checkbox"/> 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
protein LONGIFOLIA 1 (Brassica rapa)	Brassica rapa	1832	1832	100%	0.0	100.0%	896	XP_009121563.2
hypothetical protein YC2023_123701 (Brassica napus)	Brassica napus	1809	1809	100%	0.0	99.11%	896	WZ222317.1
unamed protein product (Brassica napus)	Brassica napus	1806	1806	100%	0.0	99.88%	896	CAF2348203.1
protein LONGIFOLIA 1-like (Brassica napus)	Brassica napus	1771	1771	100%	0.0	97.10%	897	XP_013667780.2
hypothetical protein Rca101_033550 (Brassica carinata)	Brassica cari...	1551	1551	100%	0.0	92.25%	902	KAL0731552.1
PREDICTED: protein LONGIFOLIA 1-like (Brassica oleracea var. oleracea)	Brassica oler...	1548	1548	100%	0.0	92.25%	902	XP_01311590.1
unamed protein product (Brassica oleracea)	Brassica oler...	1547	1547	100%	0.0	92.54%	896	VD033669.1
Protein LONGIFOLIA 1 (Hirschfeldia incana)	Hirschfeldia incana	1536	1536	100%	0.0	89.06%	921	KAJ0243573.1
hypothetical protein GI01_041384 (Brassica rapa subsp. trilocularis)	Brassica rapa	1526	1526	86%	0.0	99.65%	776	KAO5378788.1
hypothetical protein HD58_039127 (Brassica napus)	Brassica napus	1515	1673	93%	0.0	96.63%	842	KAH0907300.1
BraC05g422800 (Brassica napus)	Brassica napus	1515	1515	100%	0.0	88.16%	893	CDY26440.1
hypothetical protein Rca101_011831 (Brassica carinata)	Brassica cari...	1511	1511	100%	0.0	88.67%	925	KAL0887848.1
hypothetical protein N655_014560007 (Sinapis alba)	Sinapis alba	1510	1510	100%	0.0	87.92%	924	KAF8108188.1
hypothetical protein BraC012_046135 (Brassica campestris)	Brassica cam...	1502	1502	100%	0.0	90.37%	884	KAU047844.1
protein LONGIFOLIA 1 (Raphanus sativus)	Raphanus sativus	1494	1494	100%	0.0	86.62%	913	XP_018445962.1
unamed protein product (Eruca vesicaria subsp. sativa)	Eruca vesicaria	1456	1456	100%	0.0	87.19%	905	CAH832769.1
hypothetical protein F2Q68_00022715 (Brassica cretica)	Brassica creti...	1441	1441	94%	0.0	91.72%	845	KAF235731.1
unamed protein product (Thlaspi arvense)	Thlaspi arven...	1418	1418	100%	0.0	81.12%	928	CAH2072673.1
protein LONGIFOLIA 1 (Capsella rubella)	Capsella rubell...	1403	1403	100%	0.0	81.50%	931	XP_006287012.1
unamed protein product (Arabidopsis thaliana)	Arabidopsis th...	1365	1365	100%	0.0	80.60%	927	CAA0402799.1
protein LONGIFOLIA 1 isoform X2 (Eutrema sativum)	Eutrema sativ...	1364	1364	99%	0.0	79.96%	930	XP_006400082.1
hypothetical protein AALP_AAO158400 (Arabis alpina)	Arabis alpina	1383	1383	100%	0.0	79.70%	932	KPF25775.1
unamed protein product (Arabidopsis thaliana)	Arabidopsis th...	1383	1383	100%	0.0	79.94%	926	VY366927.1
longifolia1 (Arabidopsis thaliana)	Arabidopsis th...	1382	1382	100%	0.0	80.49%	927	NP_197062.1

[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 alignment for building a phylogenetic tree that illustrates species divergence.

Alignment: /Users/juangonzalez/Desktop/UCSD_Classes/Bioinformatics/Find_Gene_Project/LNG1_Pro
Seaview [blocks=10 fontsize=10 A4] on Mon Nov 25 20:07:59 2024

	1	MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TGDELKSIPS	GKASDNVGDT
Thaliana		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	SVAGDELKSL	PSGKTSDNVG
Brassica_rapa		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	SVAGDELKSL	PSGKTSDNVG
Shortpod_mustard		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	AVTGDELKSL	PSGKTSDNVG
False_Flax		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	VVAGDEVKSM	PSGKTTEENVG
Radish		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TVAGDELKSL	PSGKTSDNVG
Suecica		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TVAGDELKSL	PSGKTSDNVG
Pennycress		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	AGAGDELKSL	PSCKAASDSVG
Saltwater_Cress		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	VVAGDELKSR	PSEKTTEENVG
White_Mustard		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TVAGDELKSL	PSGKTSDNVG
Arenosa		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TVAGDELKSL	PSGKTSDNVG
Halleri		MSAKLLYNLS	DENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	TVAGDELKSL	PSGKTSDNVG
Arabis		MNGIFHVFYR	QHYPARRVAV	AGDELKSLPS	GKTGDNVGDT	NISTDKKEFE	KSRKNKAKE
Grape_Vine		MSAKLLHTLS	DENPDLOQKOI	GCMNGIFOLF	DRHHFLGGRR	INGHTHKELP	PGQGMEPNNA
Bittercress		MSAKLLYNLS	EENPNLNKOI	GCMNGIFQVF	YRQHYPARRV	AVVGYEKLSL	PSGKTSDNVG
	61	NISADKKETE	KSKKKKTAKE	KORGVSSESS	SRLSFSSSPC	SSSFSSADIS	TTASQFEQPG
Thaliana		VTNGSTDKKE	TEKSKKKKAA	KEKQKVVSSE	SSRLSFSSS	PCSSSFSSAD	ISTTTSQFEQ
Brassica_rapa		VTNVLTDNKE	KEKSKKKKAA	KEKQKVVSSE	SSRLSFSSS	PCSSSFSSAD	ISTTTSQFEQ
Shortpod_mustard		DINISTDNKE	TEKSKKKKAA	AKEKHKGVSS	ESSSSRLSFS	SPCSSSFSSA	DISTTASQFE
False_Flax		DTNVTTDNKE	TEKSKKKKKK	VAKEQORVVS	ESSSSRLSFS	SSOCSSSFSS	ADISTTTSQF
Radish		DTNISTDKKE	TEKSKKKKKT	KEKQRGVSSE	ESSSSRLSFS	PCSSSFSSAD	ISTTTSQFEQ
Suecica		DTNISTDNKE	TEKSKKKKKK	KEKQRVSSS	ESSSSRLSFS	SPCSSSFSSA	DISTTTSQFE
Pennycress		DTNISTDSKE	TEKSKKKKKK	AKEKQRGASS	ESSSSRLSFS	SPCSSSFSSA	DISTTTSQFE
Saltwater_Cress		DTNIVATDNKE	TEKSKKRKKK	AAKEQORVVS	ESSSSRLSFS	SSPCSSSFSS	ADISTTTSQF
White_Mustard		DTNIVATDNKE	TEKSKKKKKT	KEKQRGVSSE	ESSSSRLSFS	PCSSSFSSAD	ISTTTASQFEQ
Arenosa		DTNIVATDNKE	TEKSKKKKKT	KEKQRGVSSE	ESSSSRLSFS	PCSSSFSSAD	ISTTTASQFEQ
Halleri		DTNISRDKKE	TEKSKKKKAA	KEKQRVSSS	SSRLSFSSS	PCSSSFSSAD	ISTTTASQFEQ
Arabis		KQRGVSSSESS	SRLSFSSSPC	SSSSFSSADIS	TTTSQFEOPG	LVQTNNGENP	VREPINGSPR
Grape_Vine		PHAKDKNPK	KFVKKEQORIS	JESSRTSFSS	SSCSSTFSSV	DCNRTAQTES	FSHSQTGFPN
Bittercress		DTNISTVKSK	KKAAKEKQR	GVSSSESSARL	SFSSPCSSSS	FSSADISTTA	SQFEQSGLIQ
	121	LSNGENPVRE	PTNGSPRWGG	LMMPSDIREL	VRSSIHKETR	TRDEEALSQQ	PKSARANVSL
Thaliana		PMSNGETPAR	EPTYGSPRWG	GLVMSSDLRE	LVRSSIHKET	RTRVEEEALS	QQPKSARANV
Brassica_rapa		PGLVQTSNGE	TPVRELTPNG	SPRWGGGLMM	SSDLRELVRSS	SIHKETRTRN	EEEALSOQPK
Shortpod_mustard		QPCLIQTNNG	ENSAREPTYG	SPRWGGGLMP	SDLRELVRSS	IHKETRTRNDE	EALSOQPKSA
False_Flax		EQPGLMOTSN	GENPIRLLSSG	LRELVRNSIH	KETRDEEEDE	AVSQOOPKSA	KSNVSSLNES
Radish		SGLIQTNSNE	NPVREPTNGS	PRWGGLMMP	DLRELVKSSI	HKETRTRDEE	ALSOQPKSAR
Suecica		QPSNGENPVR	EPTNGSPRWG	GLMMSTDLRE	LVRSSIHKET	RTRDDEETLT	QOPKSARANV
Pennycress		QPGLQITNNG	ENPVREPTYG	SPRLGGLMMP	SDLRELVRSS	IHKETRTRDE	EENLSSQPKS
Saltwater_Cress		EQPGLMOTSN	VENPVRRGGLM	MSSDLRELVR	NSIHKESRAR	DDEEEDLSQQ	PKSARAKVSL
White_Mustard		SGLIQTNSNE	NPVREPTNGS	PRWGGLMMP	DLRELVKSSI	HKETRTRDEE	ALSOQPKSAR
Arenosa		SGLIQTNSNE	NPVREPTNGS	PRWGGLTMPS	DLRELVKSSI	HKETRTRDEE	ALSOQPKSAR
Halleri		SGLIQTNSNE	NPVREPTNGS	PRWGGLTMPS	DLRELVKSSI	HKETRTRDEE	ALSOQPKSAR
Arabis		WGGLMMPSDL	REFVRSSIHK	EITRTRDEEGL	SOQQPKSTR	NVSLLKESPP	SRNSNEWNEG
Grape_Vine		TPSRDILPMTO	PDASPLRLGRO	SLDLIRDIVK	SIYREACGPM	RLSKEPIKVP	VLDESLRTFG
Bittercress		TSNGENSVRE	STNVSPRWGG	KMMPSDMREL	VKSSIHKETG	TRLDALDSQQ	PKPARANVSL
	181	LKESSPSRNS	NEWSEGRRVV	KLKDSPRFSY	DEREKTRKTGA	KLKETPRLSL	DSRSNSFRSA
Thaliana		SLLKELPSR	SSNEWSEGRR	VVKLKDSPRF	SYDEREARKT	GAKFKETPRL	SLDSRSNSFR
Brassica_rapa		SARANVSLLR	EPSPSRSSNE	WSEGRRIVKL	KDSPRFSYDE	RESRKTGAKF	KETPRLSLDSR
Shortpod_mustard		RANVSLLKES	SPSRNSNEWS	EGRRVVVLKD	SPRFSYDERE	TRKTVPKLKE	TPRLSLDSSR
False_Flax		SPSRNSNEWS	EGRRVVKLKD	SPRFSYDERA	ERKIGAKFKE	TPRLSLDSSR	NSFRSARSSS
Radish		ANVSLLKESS	PSRNSNEWSE	GRRVVVLKDS	PRFSYDERET	RKTGAKLKE	PRLSLDSRSN
Suecica		ANVSLLKESS	PSRNSNEWSE	VVKLKDSPRF	SYDERETRKT	GAKLKETPRL	SLDSRSSTSFR
Pennycress		PLPKESSPSR	NSNEWSQGR	SEGRRVVKL	DSPRFSYDER	EARKTGAKLK	ETPRLSLDSSR
Saltwater_Cress		ARANVSLLLKE	SSPSRNSNEW	KLKDSPRFSY	DERETRKTGA	KFKETPRLSL	DSRSNSFRSA
White_Mustard		LKESSPSRNS	NEWSEGRRVV	KLKDSPRFSY	DERETRKTGA	RKTGAKLKE	PRLSLDSRSN
Arenosa		ANVSLLKESS	PSRNSNEWNE	GRRVVVLKDS	PRFSYDERET	TRKTGAKLKE	TPRLSLDSSR
Halleri		RANVSLLKES	SPSRNSNEWS	EGRRVVVLKDS	SPRFSYDERE	FRSARSTCSP	EPQELVTGHR
Arabis		RRVVMLKDSP	RFSYDERETR	KKGTKLKEP	RLSLDSRSNS	AIKLKDIPRL	SLDSRESSMR
Grape_Vine		KLRGPPRNSN	ERKDGSVLVLT	PRDAPRFSYD	GRESDRTFKS	KLKETPRLSL	DSRSNSFRSA
Bittercress		LKESSPSRNS	NEWSEGRRVV	KLKDSPRFSY	DEREMRKTGA		

241

Thaliana	RSSCSPEPQE	LVTGHRRRTS	SVVAKLMGLE	VIPDEPVTIQ	NRENRCDS	RPTSRVEVDL
Brassica_rapa	AKSSCSPEP	OELVTGHRR	TSSVIAKLMG	LDVVSDEPVT	DOSRENHFC	SPRPAPRVEA
Shortpod_mustard	RSNSFRSAKS	SCSPEPODFV	TTGHRRTTSS	VVAKLMGLDV	ILDEPVADQS	RENHFCDSPR
False_Flax	NSFRSAKSC	SPEPOEVVIG	HRRRTSSVVA	KLMGLEVIPD	EPVVAVODRE	NRFCDSPRPT
Radish	SPEPOELVTG	HRRRTSSVVA	KLMGLEFIPD	ESVTDVREN	RFCDSPRPTS	FVADLQRSR
Suecica	SFRSARSSYS	TEPQELVIGH	RRRTSSVVAK	LMGLEVIPDE	PVTVHNREN	FCDSPRPASR
Pennycress	GARSSCSPEP	RELVTGHRR	TSSVVAKLMG	LEGIPDEPVT	DQNREIQNRF	CDSPRPTFRP
Saltwater_Cress	SNSFRSARSS	CSPEPOEFVT	GHRRRTSSVV	AKLMGLEVIP	DESVDQNRG	NONRFCDSPR
White_Mustard	RSSCSPEPQE	LVTGHRRRTS	SVVAKLMGLE	VIPDESDTQ	NRENRCDS	RPTSRVQADL
Arenosa	SFRSARSSCS	PEPQELVIGH	RRRTSSVVAK	LMGLEVIPDE	PVTVHNREN	FCDSPRPASR
Halleri	NSFRSARSSY	STEPQELVTG	HRRRTSSVVA	KLMGLEVIPD	EPVTVHNREN	RPCDSPRAS
Arabis	RTTSSVVAKL	MGLEVIPDES	VTDOSRENR	CDSRPPFRV	EADLQRSSR	DSIKKFPTKA
Grape_Vine	GSASELKSNY	LPOQPSGSNK	RPSGVVAKML	GLDAFPDSSI	SSKAAGESKO	HRISGSPRNS
Bittercress	RYSCSPEPQE	LVAGHRRSTS	SVIAKLMGLE	VIPDESDATDQ	NRENRCDS	RPTPRVEADL

301

Thaliana	ORSRGFDTSIK	KMMPAKFPMK	ASPWAQVDG	KNOVKIPDAT	TLTVYGEIOK	RLSQLEFKKS
Brassica_rapa	DLPRSRGSDS	FKKMMPAAKF	PAKTAWPTOA	DGARNOVKAA	DAAATLTVY	EIQKRLSQLE
Shortpod_mustard	PPRVEADLQR	SRSSDSFKKM	MSAKFPMKNA	PWTQVDGAKN	QVKAADAAAT	LTVYGEIOKR
False_Flax	SRVDIDLQR	RSSDSIKTMM	PAKPFMKVAP	WTQVDGAKNO	VKAADATT	VVGEIOKRLS
Radish	SSDPKKMMP	VKFPMAVADA	ATTVTVYGEI	QKRLSOLEFK	KSENDLSALK	QILEAMETTO
Suecica	VEEDLQRSG	SDSFVKMMTA	TFPMKASPR	AQVDCAKNOV	KAADATT	YGEIOKRLSQ
Pennycress	EAELORSRSS	DSAKRMMFFP	MKAAPWSQVD	GAKNQVKAAD	AAATLTVY	IQKRLSQLEF
Saltwater_Cress	PTSRAEADLQ	RSRNSDSVKK	MMPAKFPMK	APWSQVDGAK	TQGIAADAAT	TLTVYGEIOK
White_Mustard	ORSSDSSSIK	KMMPAKFVFK	AAPWTQVDGA	KNQFAIDAA	TTTVYGEIO	KRLTHLEFKK
Arenosa	VEEDLQRSS	SDSFVKMMTA	KFPMKASPAW	QVDCAKNOV	ADATT	GEIOKRLSQL
Halleri	RVEEDLQRSSR	GSDSFVKMMT	AKFPMKASPW	AQVDCAKNOV	KAADATT	YGEIOKRLSQ
Arabis	APWTQVDVFK	NOVKAADAAAT	TLTVYGEIOK	RLSQLEFKKS	EKDLRALKQI	LEAMEKTHOL
Grape_Vine	HKDPVSPRIR	NAGSVMKPTS	TSRFPIEPAP	WKQLDGSQGP	QKPTFKHREA	ATKTLNSTPS
Bittercress	ORSRSSDSIK	KMMPAKFPM	KVAPWTQVDG	AKNQAKAADA	TTTVYGEIO	KRLSQLEFKK

361

Thaliana	EKDRLRALQOI	LEAMEKTOQL	ISKDDDDNKT	LCSSNFMQRN	NQPIPSAINT	SSMNFKSSSI
Brassica_rapa	FKKSEKDRLA	LQQILEAMEK	TOQLMSKDDD	NSSSLSTTNFM	OPSPSSKSIR	SSSIVVMKAA
Shortpod_mustard	LSQLEFKKSE	KDLRALQOQL	EAMEKTOQLI	SKDDDNSSLS	TSFMPQSPSP	KSISSSSIVV
False_Flax	OLEFKKSEKD	LRALQOILEA	MEKTOQLHISK	DDDNKTNFMQ	GIDQPVPSAT	SPSSSKNFAS
Radish	KLINKDDNN	TLSSTNFMQP	VPSAATTSP	SKSSRSSTV	MNAATASVFK	ETGNYGSASY
Suecica	LEFKKSEKDL	RALQOILEAM	EKTQOLISKD	DDNKNLICSTN	FMORTDOPIP	SATNPPSKNF
Pennycress	NKSEKDLSAL	KOILEAMEKT	QQLVSKDDDD	NNTRLSTSFM	KOIPPATTSP	SSNFRSSSI
Saltwater_Cress	RLSOLFEFKK	EKDRLRALQOI	LEAMEKTOOL	ISKDDDDNNT	SSTNFMORAD	QIPSATATTSP
White_Mustard	SEKDLRALKQ	ILEAMEKTK	LISKDDDNNT	LRSTNCLOPI	PSAATASGFK	ETGSYGSASF
Arenosa	EFKKSEKDLR	ALQOILEAME	KTOQLISKDD	DNKNLICSTNF	MORNDQPIPS	ATNPPSKNF
Halleri	LEFKKSEKDL	RALKQIILKAM	EKTQOLISKD	DDNKNLICSTN	FMQPTDQPIP	SATNPPSKNF
Arabis	TSKDDDDNNL	SSTTFMORAD	QPIPFATTSP	SSKNFRSSSI	VVMKAATAPV	FNETGNSSA
Grape_Vine	IYGEIEKRIT	ELEFKKSKGD	LRALKRILEA	MOKTKEITEA	KKDHNNSNSVS	QTSNRTSSP
Bittercress	SEKDLRALKQ	ILEAMEKTHQ	LINRDDDDNK	TLSSTNFMQ	TDNPIPSATS	PSKKNFRSSS

421

Thaliana	VVMKAAATAPV	FKDTGIACSA	SFSPRVNALP	NVKGVLNQRA	OKVIPRKOSA	MDVTPRPGYY
Brassica_rapa	SAPVFKEIGS	SSSTSSSPRS	VALPNVKS	QKGITRKOSA	MDVTPRPATK	NTSTRPLQSK
Shortpod_mustard	MKAATAPVFK	ETGNSSSTSS	SPRTVALPNV	KVSNLQNOK	VIPRKOSAMD	VTPRPGFYKG
False_Flax	SIVVMKVA	PVFKETIGISG	SASFAPRNV	LANVKGVLN	QTONVIPRK	SAMDVTPRPG
Radish	SPRSVTLP	KVSNQDRQSR	VNPKQKSAMD	VTPRPGVYKG	QTDSPTKNTG	SRQLLPKNE
Suecica	KSSSIVVMKA	AAAPVFKETG	ISGSASFS	NVALPNVKG	NLQTOKVIP	RKQSAMDETS
Pennycress	VVMKAAAAPP	VFKETGSSGS	ASFSPRSVAL	PNVKVLNQRO	PQXVAPRKOS	AMDVTPRGF
Saltwater_Cress	SSSSKNLNSSS	IVVMKAAEEA	VFKETIGNC	TSFSPRSVAL	PNVKVALNLR	TOKVTORKQS
White_Mustard	SPRSVTLPND	KVTTNLQSO	KVTPRKOSAM	DVTPRPGVYK	GQTDSPTKNT	GSROLLSKNE
Arenosa	SSSIVVMKA	AAAPVFKETG	SGSASFS	VALPNVKG	LRQTOKVIP	KOSAMDETS
Halleri	KSSSIVVMKA	AATPVFKETG	ISGSASFS	NVALPNVKG	NLQTOKVIP	RKQSAMDETS
Arabis	FSPRNVALPN	VKVGVLNLRQK	KVTPRKOSAM	DVTPRPGVYK	GQIDSATKNT	ITRQLOQKSD
Grape_Vine	SFKSPVIMK	PAKLEKSHN	LASSAIPIDG	LSGIPRQLOTG	DLVGSRKDSV	DKQTAKDLTP
Bittercress	IVVMKAAATAP	VFKETGISDS	ASFSRPNVAL	PNVKVLNQRO	TORAIPRKOS	AMDVTPRPGF

481

Thaliana	KQGTESTMKN	TSTRPLQSKS	DMAKSGKIQK	PSVSLRTPPK	KLGFEKQSRP	TSPKPPELNK
Brassica_rapa	LEMAKSGKPS	VSPRTOPKKL	GFEKOSRPTS	PKPBPKNQR	QQLSROOTES	PSPRRKPGMK
Shortpod_mustard	OTDSPTKNTT	TRPMQSIITEM	AKSGKOSKPS	VSPRTOPKKL	GFEKOSRPTS	PKPEPNKQR
False_Flax	FYKGQTDSTM	KSASTRPLQS	KNDMAKSGKI	QKPSVSPRTP	PKLIGFEKQOS	RPTSPKPEPN
Radish	AKSGKPSVSP	RTQSKKIGLE	KQSRPTSPKP	EONRTOQQL	SROOTESASP	RRKPRSVQOS
Suecica	RPGFYKGQID	SKMKNTSTR	QSKSDMARS	GKIQKPSVSP	RTPKKLGFE	KOSRPTSPKPK
Pennycress	YKGQADSAFK	NTSTRLLQSK	SDMAKSGKSO	KPSVSPRTPP	KLGFEKQSR	PTSPKPEPNK
Saltwater_Cress	GMDVTPRPGF	YKGQTDSATK	NTTTRPLQSK	SDMRSGKSO	KPSVSPRTQP	KKLAFEKQSR
White_Mustard	MAKSGKQDCK	SVSPRTOPKK	LGFEKOSRPT	SPKREPENKIQ	QQLSROQTE	SASPRRKPGI
Arenosa	PGFYKGQID	TMKNTSTRLL	QSKSDMARS	KIQKPSVSPR	TPPKKLGF	QSRPTSPKPE
Halleri	RPGFYKGQID	STMKNTSTR	QSKSDMARS	KIQKPSVSPR	TPPKKLGF	QSRPTSPKPE
Arabis	MAKSGKQDCK	SVSPRTOPKK	LGFEKOSRPT	SPKREPENKIQ	QQLSROQTE	SAPPRKKTGS
Grape_Vine	RNKHKLKEPPS	QPSRLLDKSS	ADRSSLRLTKT	SKVHOKINEN	ETSSGRNG	AVSPRLQQKK
Bittercress	YKGQTDNTT	KNTSTRPLQS	KSDMAKSGKS	QKPSVSPRTP	PKFGFKEQOS	RPTSPKSEPN

541

Thaliana	QROOQLSROOQ	ESASPRRKPG	IKSRGQLQOSE	DRLSDDESSDL	RSLRSDSNV	LASNLDTEVTI
Brassica_rapa	SRGLQQSEDR	SSDESSDLRS	LRSDSNVSSA	SNFDIEVTSR	HKCDLITEHQT	PKQRSPELGM
Shortpod_mustard	QQLSROOTES	ASPRRKPPMK	SRGMOQSEDR	LSDEGSDLRS	LRSDSNLSSA	SNLDTEVTSR
False_Flax	KTQRQQLSROQ	QIESASPRK	OQIKSRGLHQ	SEDRLSDESS	DLRSQRDSN	VSLASNLDE
Radish	EDRLSDESSD	WRSLRSDSNV	SSASNLDSEV	TSRYKYERNNS	DFTEQHTPKQ	RSPELGMRSL
Suecica	EPNKIQORQLS	ROQTESASPR	RKPGIKSRGL	QSEDRLSDESS	SSDLRSLRSN	SVNLASNLDD
Pennycress	IQRQQLSROQ	TESASPRRK	GTKTRGLQOS	EDRLSDESSD	LRSLRSDSNV	SLASNVMDV
Saltwater_Cress	PTSPKPENPK	IQRQQLSROQ	TESASPRRK	GIKSRVQQQS	EDHFSDETS	TEQHTPKQRS
White_Mustard	TSRSMQQSVD	RLSDESSLR	SLRSDSNVSS	ASNLDSEATS	RYRYERNNSD	TEQHTPKQRS
Arenosa	PNKIQORQLSR	QOTESASPR	KPGIKSRGLQ	QSEDRLSDES	SDLRSLRSN	NVSLASNLDT
Halleri	PNKIQORQLSR	QOTESASPR	KPGIKSRGLQ	QSEDRLSSEES	SDLRSLRSN	NVSLASNLDT
Arabis	KSRGLQOSED	RLSDESSLR	SLRSDSNVSL	ASNLDIEVTS	RYRSENRNDI	TEQYTPKQRS
Grape_Vine	LELDKQSRST	TPSPESSRVR	RQSSRQLTEP	SSPARKLQR	APNLLQSDDO	LSEISGDSRN
Bittercress	KIQRQQLSKQ	OTESASPRK	OQIKSRALQQ	SEDRLSDESS	DLRSLRSDSN	VS LTSNLDT

601

Thaliana	SRYNYERNNSD	ITEQHTPKQR	SPDLGMRSL	KPLKVTVEQP	SPVSVIDVAF	DEDDSPSPV
Brassica_rapa	RSLPKPLKIT	VEQPSPVSL	DVAFDDDESP	SPVRKISIVE	KDDDHIRSEE	SLWMKKHNNL
Shortpod_mustard	YKCDLITEHQT	PKQRSPSELGM	RSLPKPLKIT	VEQPSPVSL	DVTFDDDESP	SPVRKINIVF
False_Flax	VTSRYNFERN	SDIIIEQHTPK	QKSPELGMRS	LPKPLKVTV	QSPSPVSLDV	AFDEDDSPSP
Radish	PKPLKVTVEQ	PSPVSVLDVA	FDEDESPSPV	RKISVVFKDD	DHLRSEESEW	MNHRTLRRS
Suecica	TEVTSRYNYE	RNSDITEQHT	PKQSPSPDLM	RSLPKPLKV	VEQPSPVSVL	DVAFDEDESP
Pennycress	TSRYKYERN	DIAEQHTPKQ	RSPDPGLM	RSLPKPLKV	PSPVSVLDVA	FDEDESPSPV
Saltwater_Cress	SLTSNLIDIEG	TSRYKYERNN	DITEQHTPKQ	RSPELGMRS	PKPLKITVEQ	PSPVSVLDVA
White_Mustard	PEFGMRSLEPK	PLKVTVEOPS	PVSVIDVAFD	EDESPSPV	ISVVFKDDDH	LRSEESQWMN
Arenosa	EVTTSRYNYER	NSDITEQHTP	KQRSPDLM	SLPKPLKV	EOPSPVSVL	VAFDEDESPS
Halleri	EVTTSRYNYER	NSDITEQHTP	KQRSPDLM	SLPKPLKV	EOPSPVSVL	VAFDEDESPS
Arabis	PDLGMKLLPK	PLKVTVEOPS	PVSVIDVAFD	EDESPSPV	ISIVFKDDNL	SSSEEPQWMN
Grape_Vine	LSYQVTSIDR	SGGINSIISFO	HGGQKHNGD	GTMTKFATAT	QEOPSPVSVL	DAAFYKDDLP
Bittercress	VTSRYKYERN	SDIMEQHTPK	ORSELGMRS	PKPLKLTVEQ	PSPVSVLDVA	FDEDESPSPV

661

Thaliana	KISIVFKEDD	NLSSEESHW	NKNNNLCRS	VWPESNTSLK	QPDAELTEGF	MEDDAEFKNG
Brassica_rapa	CRSIVWPESN	TSLNQPAVL	TESFMEEGAD	LRNGDRKYIS	EILSASGLLK	DIDYSMLSIO
Shortpod_mustard	KDDDHHLKSEE	SQWMSKHN	CRSIVWPESN	TSLNQPAEL	TESLMEEGA	LRNGDHKVIS
False_Flax	VRKISIFKED	DNLSSQDAEL	MNKHSNLCRS	IVWPESNASL	KOPDAELIEG	FMEEDAEFKN
Radish	IVWPETE	SLNQDAEL	EGIREEGAEL	NNGDHKYISE	ILLASGLL	IDYSMLSIO
Suecica	SPVRKISIVF	KEDANLSYEE	SQWMNKHN	CRSIVWPESN	ATLQPD	EGFMEEEAEF
Pennycress	RKISIAFKED	HLSYEEQWM	NKHSNLCRS	WPESNASPK	QPDDELIEGS	MEEDAKFRNG
Saltwater_Cress	FDEDESPSPV	RKISIVFKED	DNLGEES	MNKORNICRS	IVWPENNASK	QPDAEPIMGF
White_Mustard	KHRTLRRNSIV	WPESGNSLQ	SDAELNEGIM	EEGAELNNGD	HKYISEILL	SGLLRIDYD
Arenosa	PVRKISIVFK	EDANLSSEES	QWMNKHNSLC	RSIVWPESNA	TLKOPDAELM	EGFMEEEAEF
Halleri	PVRKISIVFK	EDANLSSEES	QWMNKHNSLC	RSIVWPESNA	TLKOPDAEV	DGFMEEDADEF
Arabis	KHSNLCRSIV	WPESNVSLQ	PDHNEGYMEE	GVELKNGDHK	YISEILLASG	LLRDIDYSMM
Grape_Vine	SPVKKISIFN	KDDETLYNDE	MEWATKLEN	ENLVQRIREL	NSTHNEFSDV	LIA
Bittercress	RKISIVFKEE	DNLSSSEES	MNKQRNLCRS	IVWPESNMSV	KLPDT	MEEGAELKN

721

Thaliana	DHKYISEIML	ASGLLRLIDY	SMISIOLHOA	HPINPSLFF	VLEQNKTSNV	SLQDNKHKGR
Brassica_rapa	LHQAHLPINP	SLFFVLEEQNK	TSNVTHRGRG	FGQQTANLIG	RSRRKLVFD	VNEILARKFA
Shortpod_mustard	EILSASGLLK	DIDYNMISTQ	LHQAHLPINP	SLFFVLEQSK	TGNVRHDNKH	RGRGFGOOOA
False_Flax	GDHKYAAEIL	SAAGLRLID	YSMISIQLHQ	ALHPLNPNLF	FVLEQNKTSN	VSLQDNKHKG
Radish	HOAHLPINP	LFFFVLEQSKT	SNVTOQHDSK	RIGIFGQOOT	ANLIERSR	LVDFTINEIL
Suecica	KNGDHKYISE	ILLASGLL	IDYNMISIQL	HOAHLPINP	LFFFVLEQSKT	SNVSPQDNK
Pennycress	DHKYISEIML	ASGLLRLIDY	SMMSIQLHQ	HPINPTLFF	VLEQSKTSNV	SQHDNKHKG
Saltwater_Cress	MDEGVFEFKNG	DDKYISEILL	ASGLLRLIDY	SMISIQLHQ	HPINPSLFF	VLEQNKTSNL
White_Mustard	MLSIQLYQAH	LPINPSLFFV	LEQNQNTSNAT	QHDKHHKGIG	FGQQTNVNLIE	RSRRKLVFD
Arenosa	KNGDHKYISE	ILLASGLL	IDYNMISIQL	HOAHLPINP	LFFFVLEQSKT	SNVSPQDNK
Halleri	KNGDHKYVSE	ILLASGLL	DIDYSMISTQ	LHQAHLPINP	SLFFVLEQNK	TSNVSPQDNK
Arabis	SIOLHQAHLP	INPSLFFVLE	ONQTSNTVQQ	DNKHSNRGFG	QOQTTNLLIDR	IRRKLVFDII
Grape_Vine	PDHYRIL	LASGLLRLDCS	GLMITKLHQ	SHPINPKFL	VLEQNRDVAN	ILNDKYSSQ
Bittercress	GDHEYISEIL	LTSGLLRLID	YSMISIQLHQ	AHLPINPSLF	FVLEQNKTSN	VSLQDNKHKG

781

Thaliana	GFGQQQTVNL	VERS SKRKLIF	DTINEILAH	FAAE GCTKQ	SITLSISTOR	THEKSSRGEE
Brassica_rapa	AEGCTKOPYI	TSSI SPMLKT	DKSS RKGELL	EALC SEIDRL	QDNS NCILDE	DDED LIWEDL
Shortpod_mustard	ANLIERSR	LVFD NVNEIL	ARKFAAE GCT	KOPY ITSS	QLMT TEKSS	GKELLES ICS
False_Flax	RGFGQQQTAN	LIERS KRLV	FDTINE ILAR	RFAADG CTQ	PSIIS SISPL	RITE KNSRGK
Radish	ARKFTAEGCT	KQCP YITTS	SQRL TTNKS	KG KELL KTMC	SEID RL DNS	KCIL DDEDE
Suecica	KGRGFGQQQD	VNL IERS KRR	LVFD TINEIL	AHRFAAE GCT	KOPS II SIS	TQ RIPVKSSR
Pennycress	GFGQQQTANL	MERS RRKLV	DTINEIL A	FAAE GCTKQ	SIV SSIS RLR	TT DKSS RGEE
Saltwater_Cress	SLQDDKHKAR	GFGQQQTANL	IDRS RRKLV	DTINEIL A	FAA EGCTKQ	SMISS I RPLN
White_Mustard	MNEIL A	AEGCTKOPYI	TPSIS QLRTT	KSS SKG KELL	ETMC SEID RL	QDN SKCIL DE
Arenosa	KGRGFGQQQD	VNL IERS KRR	LVFD TINEIL	AHRFAAE GCT	KOPS II SIS	TQ RIPVKSSR
Halleri	HKGRGFGQQQ	TVNL IERS KRR	KL VFDT NEI	LAHR FAA EGC	TKRPS I SIS	STQ RIPVKSS
Arabis	NEIL ARN F	EGCTK QT SII	SSIS PLR TT	KSS SKG KELL	ILC SEID RL	DDS NCIL DED
Grape_Vine	TAOS KLQ	IF DVV NEIL	OKLA FT GSSE	PC FL PN KIVR	RS QNG QELL R	ELC SEID QLQ
Bittercress	RGFG M	RS RR KLV P	INE IARR FA	AEG CT KQPSI	ISSI SPL RTT	EKSS RGKELL

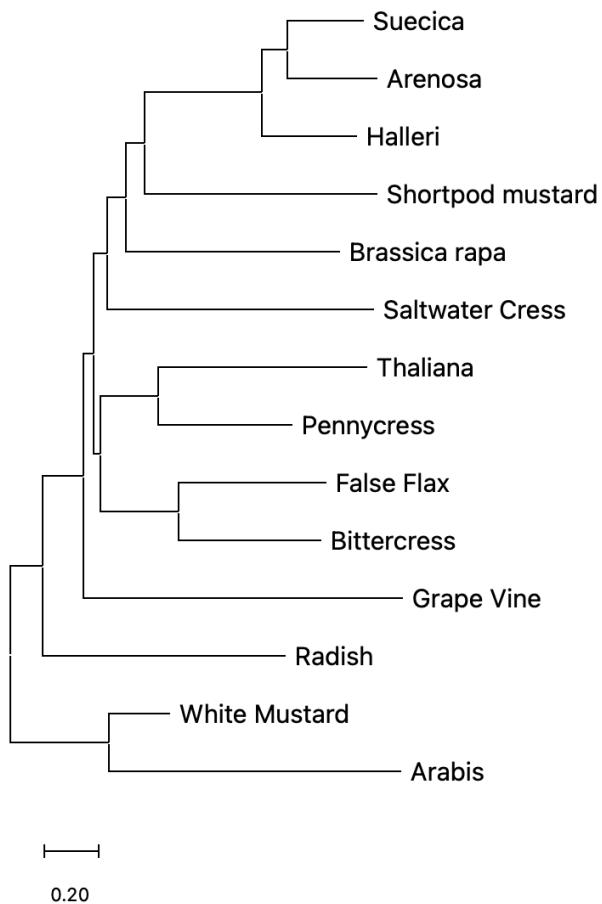
841

Thaliana	LLQTLCSEID	RQDQNSKCIL	DEDDEDLIWE	DQSHGMNWK	EIEGETPGLV	LDIERLIFKD
Brassica_rapa	QSOGMNWKEI	EGETPGLVLQ	IERLIFKDLI	SEVVTSEVAA	FPGNKLSGQP	ROLFHC
Shortpod_mustard	EIDRLQDNSN	CILDDEDDED	LIWEDMQSQG	MNWKEIDGET	PGLVLDIERL	IFKDLISEVV
False_Flax	ELLQTICSEI	DRLQNSKCI	LDDDEDDELL	WEDLQSHGMN	WKEIEGETPQ	LVLDIERLIF
Radish	DLFWEDLQSQ	GMMWKEIEGK	TPGLVLDIER	LIFKDLISEV	VTSEVAASPG	MLSGQPIRVF
Suecica	GKELLOTLCS	EIDRLQDNSK	CILDDEDDED	LIWEDLQSHGM	MNWKEIEGEI	PGLVLDIERL
Pennycress	LLETLCSEID	RQDQNSKCIL	DEDDEDLIW	EDLQSQGMNW	KEIEGETPGL	VLDIERLIFK
Saltwater_Cress	TTEKSSRGKE	LLQTLCSEID	QLODKAKCIL	DEDDEDLIWE	DLQSQGMNWK	EIEGETPGLV
White_Mustard	DDEDLIWEGL	QOGGMNWKEI	EGETPGLVLD	IERLIFKDLI	SEVVTSEVAA	SPGMLSGKPR
Arenosa	GKELLOTLCS	EIDRLQDNSK	CILDDEDDED	IWEDLQSHGM	NWKEIEGEIP	GLVLDIERLI
Halleri	RGKELLOTLCS	SEIDRLQDNS	KCILDDEDDED	LIWEDLQSHGM	MNWKEIEGEI	PGLVLDIERL
Arabis	DEDLIWEGLQ	SQGMNWKDI	GEPTGLVLDI	ERLIFKDLSI	EVVISELAAI	MLPSLGQPRO
Grape_Vine	GNNSDCSLEN	EVSWEDIMHR	SANRADFHGE	VSGIALDVER	LIFKDLIGEV	LNGEAALSLRA
Bittercress	QTLCSEIDRL	QDNSKCILDD	EDDEDLIWE	LQSHGMNWKE	IEGETPGLVL	DIERLIFKD

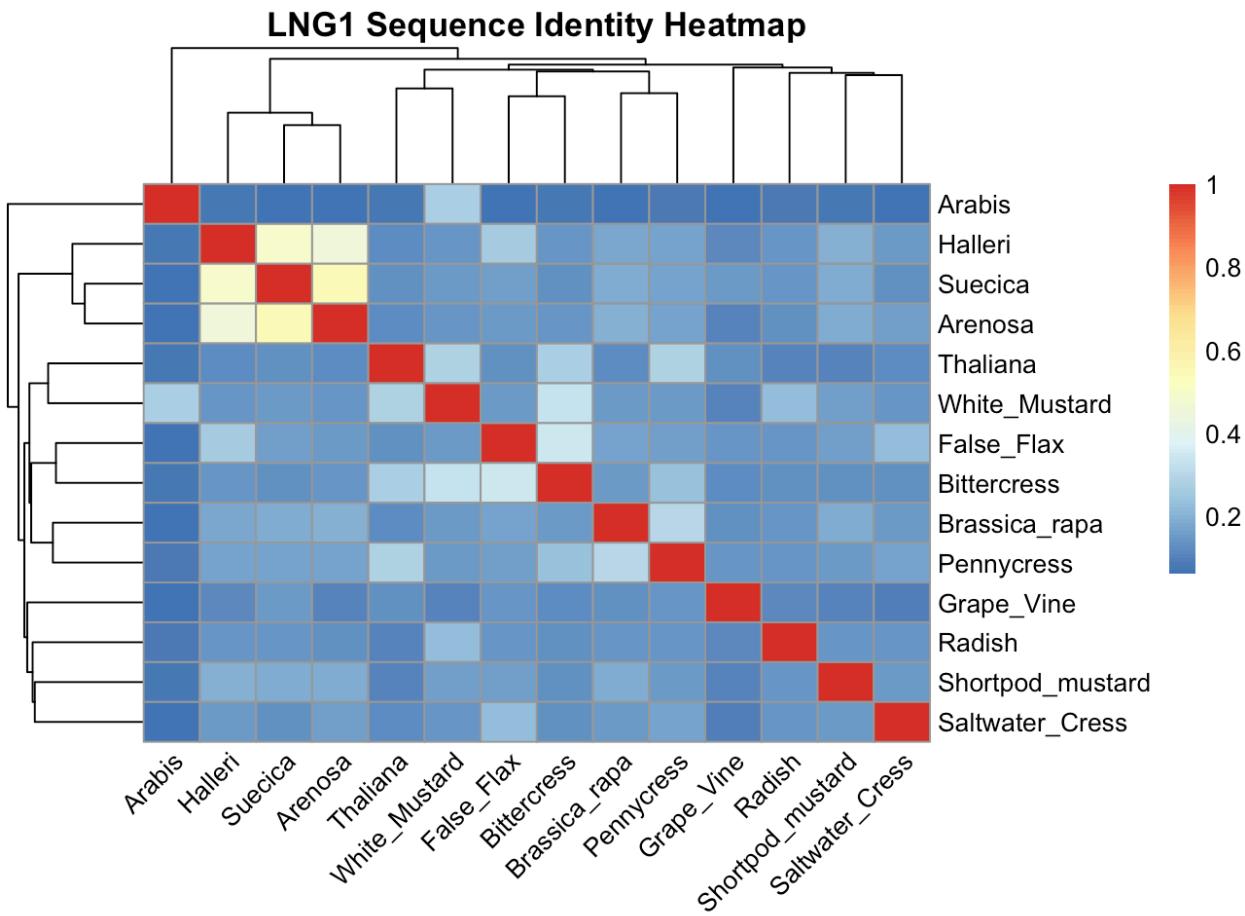
901

Thaliana	LIGEVVTSEF	AAFPRLMLSGQ	PRQLFH C
Brassica_rapa	TSEVAAPFGM	LSGQPRQLFH	C
Shortpod_mustard	KDLISeVVTS	EVAAPFGMLSGQ	GQPRQLFH C
False_Flax	NC		
Radish	IFKDLIIGEVV	TSEVAAFPGT	LRGQPRQLFH C
Suecica	DLISeVVTS	AAGMLSLGKPR	QLFH C
Pennycress	LDIERTLIFKD	LISEVVTS	AAFPKLMLSGQ PGOLFHS
Saltwater_Cress	QLFNF		
White_Mustard	FKDLIIGEVTT	SEVAAFPGTL	RGQPRQLFH C
Arenosa	IFKDLIIGEVV	TSEVAAFPGT	LRGQPRQLFH C
Halleri	LFHC		
Arabis	RPRGHYYHQRL	FPK	
Grape_Vine	ICEVVTS	AFPGMLSGQ	RQLFH C
Bittercress			

[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.

I went to PDB, used one sequence (A. thaliana) and expanded the e-value. This is what I received:

Advanced Search Query Builder [Help](#)

Full Text [Help](#)

Structure Attributes [Help](#)
 -- Type to filter and/or select an attribute --
 Add Attribute Add Subquery Remove Subquery

Add Subquery

Chemical Attributes [Help](#)
 AND
 -- Type to filter and/or select an attribute --
 Add Attribute Add Subquery Remove Subquery

Add Subquery

Sequence Similarity [Help](#)
 AND MSAKLLYNLSDENPNLNKQIGCMNGIFQVFYRQHYPPLRVGTDELKSLPSKGASDNVGDTNISADKKETEKSKKKTAKEKQRGVSSSESSRLSFSSSPCSSSSADISTTASQFEQPGLSNGENPVREPTNGSPRW
 GGLMMPSDIRELVRSSIHKETRDDEALSQQPKSARANVSLKESSPRNSNEWSEGRRV/KLKDSPRFSYDERETRKIGAKLKETPRLSDLRSNSFRSARSSCPEPQELVTGHRRTTSSWAKLMGLEVIPDEM
 Entry ID 1MBN Sequence Type Protein E-Value Cutoff 10000000 Identity Cutoff 0 % (Integer only) Count Clear

Sequence Motif

Structure Similarity

Structure Motif

Chemical Similarity

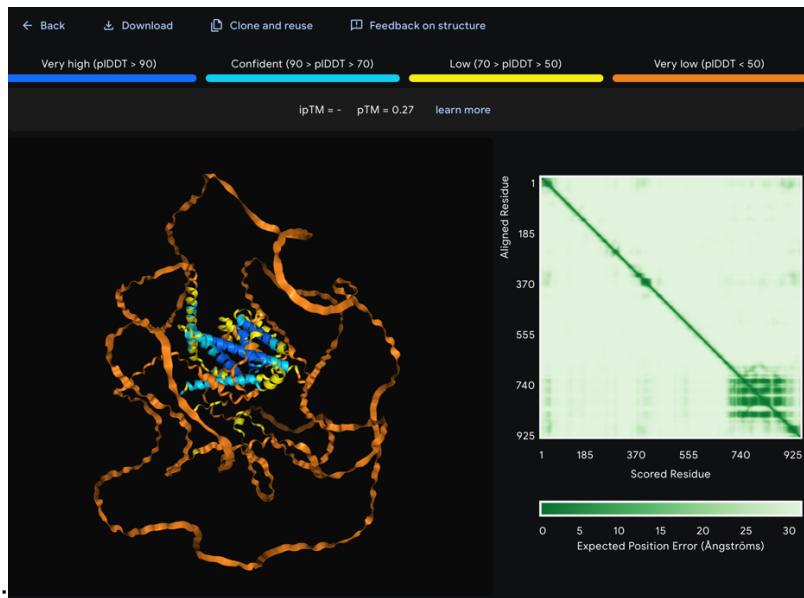
Return Structures grouped by No Grouping [Include Computed Structure Models \(CSM\)](#) Count Clear Search

PDB Identifier	Organism	Macromolecule	Technique	Resolution	E-value	Sequence Identity (%)
AF-K7MW93F1	Glycine max	DUF4378 domain-containing protein	AlphaFold Predicted Structure	N/A	48.5	80.2
AF-Q7XIU5F1	Oryza sativa Japonica Group	Os07g0109400 protein	AlphaFold Predicted Structure	N/A	48.48	78.1
AF-A0A1D6HZT8F1	Zea mays	Protein LONGIFOLIA 2	AlphaFold Predicted Structure	N/A	47.93	77.6

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

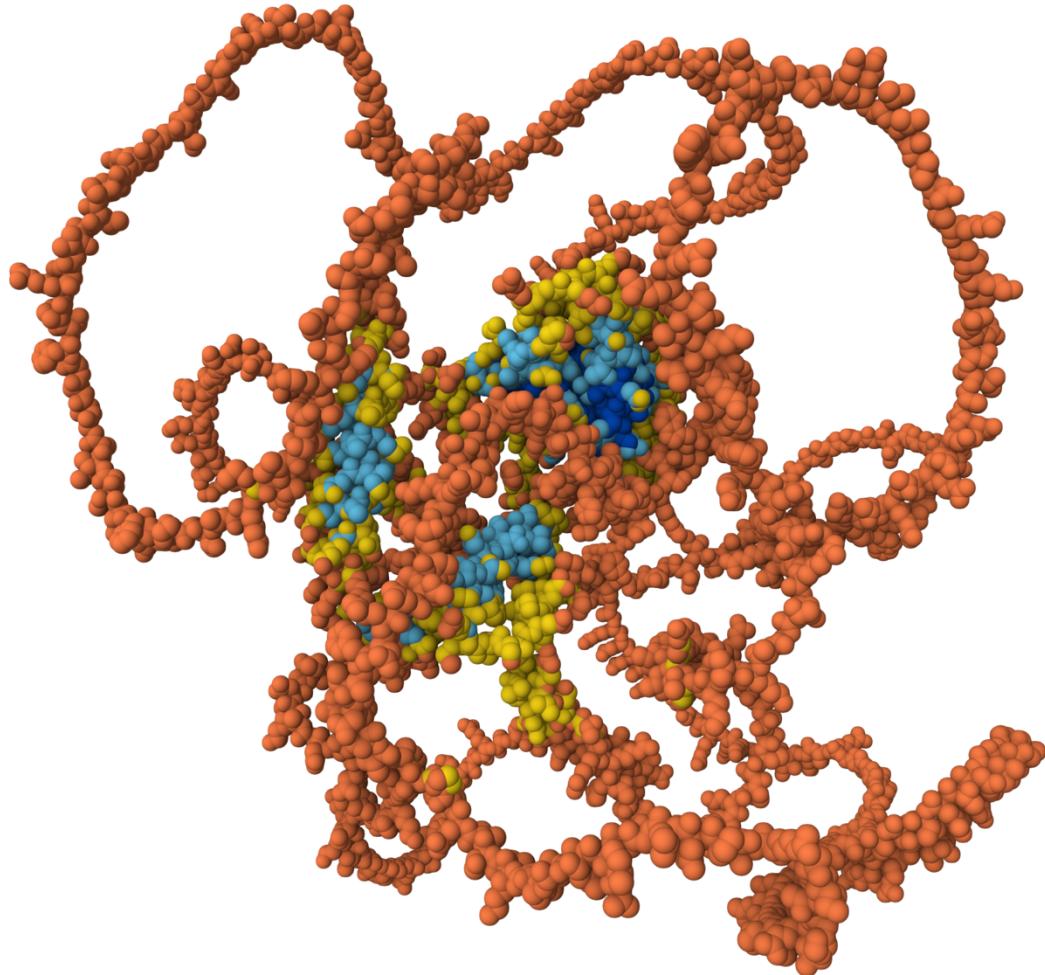
Note that this can take some time depending upon your sequence length. If your model is taking many hours to generate or your input sequence yields a “too many amino acids” (i.e. length) error you can focus on a single domain from your sequence - identify region by searching for [PFAM](#) domain matches.

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.).



AlphaFold Server Results:

Molviewer Image (UltraHD Downloaded Figure) with space fill and pLDDT validation coloring.



[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]”.

Custom Filtering

Showing 1-20 out of 119 records

	#	E-Value	Positives %	Identities %	Score (bits)	Score	Length	CHEMBL ID	Name +	UniProt Accessions	Type	Organism	Compounds	Activities
1.	0.00057	35.2	23.3	43.1	100	1007	CHEMBL1908382	Serine/threonine-protein kinase PR4 homolog	Q13523	SINGLE PROTEIN	Homo sapiens	[251]	[270]	
2.	0.0025	35.9	20.3	40.8	94	568	CHEMBL4295761	Protein AF-9	P42568	SINGLE PROTEIN	Homo sapiens	[49]	[111]	
3.	0.004	35.3	18.4	40.4	93	1334	CHEMBL4879481	WASH complex component	Q9PQLZ	SINGLE PROTEIN	Mus musculus	[2]	[2]	
4.	0.011	36	19.3	38.9	89	754	CHEMBL3707467	Peptidyl-prolyl cis-trans isomerase G	Q13422	SINGLE PROTEIN	Homo sapiens	[10]	[10]	
5.	0.013	39.2	20.5	38.5	88	602	CHEMBL4296025	Cytoskeleton-associated protein 4	Q07065	SINGLE PROTEIN	Homo sapiens	[5]	[7]	
6.	0.014	40.8	28.2	38.5	88	838	CHEMBL5291560	Transforming acidic coiled-coil containing protein 3	Q9Y6AS	SINGLE PROTEIN	Homo sapiens	[1]	[2]	
7.	0.021	36	19.1	38.1	87	2461	CHEMBL3217383	Microtubule-associated protein 1B	P15205	SINGLE PROTEIN	Rattus norvegicus	[1]	[1]	
8.	0.027	37.8	25.6	37.7	86	1309	CHEMBL2417352	Microtubule-associated serine/threonine-protein kinase 3	Q60307	SINGLE PROTEIN	Homo sapiens	[30]	[49]	
9.	0.039	40.7	22.1	37	84	565	CHEMBL2366583	Ecdysone receptor	Q4W6D0; Q4W6CB	PROTEIN COMPLEX	Leptinotarsa decemlineata	[3]	[3]	
10.	0.045	34.1	17.1	37	84	1220	CHEMBL4105852	Eukaryotic translation initiation factor 5B	Q60841	SINGLE PROTEIN	Homo sapiens	[170]	[170]	
11.	0.061	37.6	22.7	36.6	83	1616	CHEMBL1993	5'-methyltransferase 1	P26358	SINGLE PROTEIN	Homo sapiens	[767]	[978]	
12.	0.082	38	23.3	36.2	82	2564	CHEMBL3108647	Histone-lysine N-methyltransferase SETD2	Q9BYW2	SINGLE PROTEIN	Homo sapiens	[90]	[119]	
13.	0.1	42.2	24.4	35.8	81	1096	CHEMBL4295672	Lysine-specific demethylase PHF2	Q75151	SINGLE PROTEIN	Homo sapiens	[16]	[16]	
14.	0.11	37.8	23.6	35.8	81	2446	CHEMBL4523214	Transcription factor HIVEP2	P31629	SINGLE PROTEIN	Homo sapiens	[1]	[1]	
15.	0.11	43.8	26.5	35.8	81	3433	CHEMBL4523230	Urophorin	P46939	SINGLE PROTEIN	Homo sapiens	[81]	[89]	
16.	0.15	40.9	22.1	35	79	589	CHEMBL2176776	Sentrin-specific protease 2	Q9HC62	SINGLE PROTEIN	Bromodomain and WD repeat-containing	[74]	[79]	
17.	0.18	43.7	20.9	35	79	2320	CHEMBL3351192	Sentrin-specific protease 2	Q9NS16	SINGLE PROTEIN	Homo sapiens	[52]	[55]	

No target-associated assays or ligand efficiency data for plant-specific proteins are available as of 11/30/24.

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)	

