

Find a Gene Project – BGGN 213 AY2021 Fall

Jack Reddan (PID: A59010543)

[Q01] 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.

Name: glutathione reductase
Accession: NP_191026
Species: *Arabidopsis thaliana*
Function: glutathione-disulfide reductase activity

Obtained from NCBI (1).

[Q02] 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: NCBI TBLASTN (v2.12.0) searched against all ESTs
Database: Expressed Sequence Tags [est] database
Organism: None
Results:

[← Edit Search](#) [Save Search](#) [Search Summary ▾](#) [How to read this report?](#) [BLAST Help Videos](#) [Back to Traditional Results Page](#)

Job TitleNP_191026:glutathione reductase [Arabidopsis...]

RID[RONE0BNT013](#) Search expires on 10-22 01:45 am [Download All ▾](#)

ProgramTBLASTN [Citation ▾](#)

Databaseest [See details ▾](#)

Query ID[NP_191026.1](#)

Descriptionglutathione reductase [Arabidopsis thaliana]

Molecule typeamino acid

Query Length565

Other reports [?](#)

Filter Results

Organism only top 20 will appear ☐ exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

Percent Identity to E value to Query Coverage to

[Filter](#) [Reset](#)

Descriptions

Graphic Summary

Alignments

Taxonomy

Sequences producing significant alignments

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

☒ select all 100 sequences selected [GenBank](#) [Graphics](#)

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	chlorokybus_71836.3_lrc301_c Chlorokybus atmophyticus EST library Chlorokybus atmophyticus cDNA 5' mRNA...	Chlorokybus atm...	578	578	87%	0.0	58.61%	4069	HO413015.1
<input checked="" type="checkbox"/>	CLS_cLIFproots_12a4_1_j12cLibkit5LD_E06 CLS_cLIFproots_plant Festuca arundinacea cDNA clone 12j12 5'...	Lolium arundinac...	569	569	59%	0.0	81.66%	1063	GT033527.1
<input checked="" type="checkbox"/>	KP1B_113E21F06011777 KP1B Nicotiana tabacum cDNA clone KP1B_113E21 mRNA sequence	Nicotiana tabacum	560	560	54%	0.0	85.39%	926	EB682653.1
<input checked="" type="checkbox"/>	CLS_cLIFpEISpn_24a2_1_a12cLibkit5LD_A06 CLS_cLIFpEISpn_plant Festuca pratensis cDNA clone 24a12 5'...	Festuca pratensis	560	560	58%	0.0	82.07%	987	GO853961.1
<input checked="" type="checkbox"/>	GR_Ea44J17.r.GR_Ea Gossypium raimondii cDNA clone GR_Ea44J17 3' mRNA sequence	Gossypium raim...	545	545	52%	0.0	84.80%	891	CO080876.1
<input checked="" type="checkbox"/>	UFL_352_72 Cotton fiber 0-10 day post anthesis Gossypium hirsutum cDNA mRNA sequence	Gossypium hirsu...	540	540	54%	0.0	82.79%	1014	ES820783.1
<input checked="" type="checkbox"/>	CHTM8760.b1_P06.ab1 CHT(LMS) Jerusalem artichoke Helianthus tuberosus cDNA clone CHTM8760 mRNA s...	Helianthus tuber...	518	518	53%	1e-180	80.67%	901	EL456095.1
<input checked="" type="checkbox"/>	KT7C_109E11F051221T7 KT7 Nicotiana tabacum cDNA clone KT7C_109E11 mRNA sequence	Nicotiana tabacum	516	516	50%	3e-180	85.21%	856	EB450985.1
<input checked="" type="checkbox"/>	SLA_T3_196_C04_20APRIL2006_028.1 SLA (leaves from unfertilized Solanum tuberosum Shepody) Solanum t...	Solanum tuberos...	521	607	69%	7e-180	68.80%	1307	JG558174.1
<input checked="" type="checkbox"/>	RR4A214TF.RR4(PB) Raphanus raphanistrum subsp. landra cDNA 5' mRNA sequence	Raphanus rapha...	513	513	46%	3e-179	93.89%	804	EV568609.1

Chosen match: Accession H0413015.1, 4069 bp *Chlorokybus atmophyticus* mRNA sequence, highlighted in purple above (2). Alignment details are printed below:

Query: glutathione reductase [Arabidopsis thaliana] Query ID: NP_191026.1 Length: 565

>chlorokybus_71836.3_lrc301_c Chlorokybus atmophyticus EST library Chlorokybus atmophyticus cDNA 5', mRNA sequence
Sequence ID: H0413015.1 Length: 4069
Range 1: 2451 to 3941

Score:578 bits(1489), Expect:0.0,
Method:Compositional matrix adjust.,
Identities:296/505(59%), Positives:363/505(71%), Gaps:17/505(3%)

Query	52	LRPRIALLSNHRYYHSRRFS-----VCASTDNGAESDRHYDFDLFTIGAGSGGVRA	102
		LRP L HSR+ S V AS++ YD+D+ TIGAGSGGVRA	
Sbjct	3941	LRPS-GLAQQSRPHSRQQSRTVQRYGLRVIASSNGSG-----YDYDVITIGAGSGGVRA	3780
Query	103	SRFATSFSGASAAVCELPFSTISSDTAGGVGGTCVLRGCVPKKLLVYASKYSHEFEDSHGF	162
		SR A+ GA A E+PF+ ++SDT GGVGGTCVLRGCVPKKLLVY S +S+EF+DS GF	
Sbjct	3779	SRIASQLGAKVACVEMPFPNNVASDTEGGVGGTCVLRGCVPKKLLVYGSIFSNEFDDDSAGF	3600
Query	163	GWKYETEPSHDWTTLIANKNAELQRLTGIYKNILSKANVKLIEGRGKVIDPHTVDVDGKI	222
		GWK EP W TL NKN EL RL +Y+NILSKANV+L+EGR ++D HTVD+DGK	
Sbjct	3599	GWKLPGEPKFTWQTLNENKNKELTRLNNVYRNILSKANVELLEGRASLVDAHTVDIDGKQ	3420
Query	223	YTTRNIIIAVGGRPFIPIPGKEFAIDSDAALDLPSKPKKIAIVGGGYIALEFAGIFNGL	282
		T +NI++A GGR F IPG E AIDSD AL L PK+IAI GGGYIALEFA IF+G	
Sbjct	3419	LTAKNIILATGGRSFALPIPGAHAIDSALKSLDEVKRIAIYGGGYIALEFACIFSGF	3240
Query	283	NCEVHVFIHQKKVLRGFDEEDVRDFVGEQMSLRGIEFHTEESPEAIIKAGDGSFSLKTSKG	342
		+V VF R LRGFDE++R+ + E++ +GI H + + E I K +G ++LKT+ G	
Sbjct	3239	GAKVDVFYRAPLPLRGFDEEIRNALVEELGKKGINLHPKCTAEEIRKEANGEYTLKTNCG	3060
Query	343	TVEGFSHVMFATGRKPNTKNLGLENVGVKMAKNGAIEVDEYSQTSVPSIWAVGDVTDRIN	402
		+ VMFATGR PNTK L L+ VGV + GAI VDEYS+T+VP+I+A+GDVT+RIN	
Sbjct	3059	EFKA-DLVMFATGRTPNTKYLNLDAVGVDTEKGAIVVDEYSRTTVPNIFAIGDVTNRIN	2883
Query	403	LTPVALMEGGALAKTLFQNEPTKPDYRAVPCAVFSQPPIGTVGLTEEQAIEQYGDVDVYT	462
		LTPVALMEG A+AKT+ Q EPTKPD+ VP AVF+QPPIGT GLTEE+A EQ+ +VDVYT	
Sbjct	2882	LTPVALMEGTAVAKTI-QGEPTKPDHVNVPASVFTQPPIGTAGLTEEEAKEQFDEVVYT	2706
Query	463	SNFRPLKATLSGLPDRVFMKLIVCANTNKVLGVHMCGEDSPEIIQGFGVAVKAGLTKADF	522
		S+FRP+K T+SG +R MK+IV T+KVLG+HM GE SPEI+QGF VA+K G TK	
Sbjct	2705	SSFRPMKHTISGRDERSLMKIIVDVKTDKVLGIHMLGESSPEILQGFVAVALKCGATKKQL	2526
Query	523	DATVG VHPTAAEEFVTMRAPTRKFR	547
		DAT+G+HPTAAEEFVTMR TR+ R	
Sbjct	2525	DATIGIHPTAAEEFVTMRTVTRQHR	2451

[Q03] 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.

EMBOSS Transeq for TBLASTN result in Q2 (3) -

Input:

DNA/RNA: Chlorokybus atmophyticus cDNA (Acc HO413015.1)
FRAME: 6 (All six frames)
CODON TABLE: Standard Code

Chosen sequence:

```
>2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus atmophyticus EST
library Chlorokybus atmophyticus cDNA 5', mRNA sequence
ALLPAPRVA AVATRSSSSRRSELPGLLARPLGVSR SFRGFSGLRPSGLAQQQSRPHSRQQ
SRTVQRYGLRVIASSNGSGYD YDVITIGAGSGGVRASRIASQLGAKVACVEMPFN NVASD
TEGGVGGTCVLRGCVPKLLVYGSIFSNEFDD SAGFGWKLPGE PKFTWQTLNENKNKELT
RLNNVYRNILSKANVELLEGRASLVDAHTVDIDGKQLTAKNIILATGGRSFALPIPGA EH
AIDSDKALS LDEV PKRIAIYGGGYIALEFACIFSGFGAKVDVFYRAPLPLRGFDEEIRNA
LVEELGKKGINLHPKCTAEEIRKEANGEYTLKTNCGEFKADLVMFATGRTPNTKYLNLDA
VGVDTTEKGAIVVDEYSRTTVPNIFAIGDVTNRINLTPVALMEGTAVAKTIQGEPTKPDH
VNVPSAVFTQPPIGTAGLTEEEAKEQFDEVDVYTSSFRPMKHTISGRDERSLMKIIVDVK
TDKVLGIHMLGESSPEILQGFVAALKCGATKKQLDATIGIHPTAAEEFVTMRTVTRQHRK
EKQQQQQEEKEKVAAAK*
```

Name (Unofficial): *Chlorokybus* putative glutathione reductase.

Species (4): *Chlorokybus atmophyticus*: Cellular organisms;
 Eukaryota; Viridiplantae; Streptophyta;
 Chlorokybophyceae; Chlorokybales; Chlorokybaceae;
 Chlorokybus

[Q04] 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.

Conducted a BLASTP search against the NR database using the protein sequences listed in [Q3] (2):

Method: NCBI BLASTP (v2.12.0) searched all nr protein entries
Database: non-redundant protein sequences [nr] database
Organism: None
Results:

BLAST[®] » blastp suite » results for RID-R0R28H65016

[Home](#) [Recent Results](#) [Saved Strategies](#) [Help](#)

[Edit Search](#)

[Save Search](#)

[Search Summary](#)

[How to read this report?](#)

[BLAST Help Videos](#)

[Back to Traditional Results Page](#)

Job Title

2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus...

RID

[R0R28H65016](#) Search expires on 10-22 02:13 am [Download All](#)

Program

BLASTP [Citation](#)

Database

nr [See details](#)

Query ID

lcl|Query_53894

Description

2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus ...

Molecule type

amino acid

Query Length

557

Other reports

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

Filter Results

Organism

only top 20 will appear ☐ exclude

Type common name, binomial, taxid or group name

[Add organism](#)

Percent Identity

E value

Query Coverage

to

to

to

Filter

Reset

Descriptions

Graphic Summary

Alignments

Taxonomy

Sequences producing significant alignments

[Download](#) [Select columns](#) [Show](#) 100

☒ 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"/>	Glutathione reductase [Klebsormidium nitens]	Klebsormidium n...	616	616	86%	0.0	63.04%	580	GAQ78357.1
<input checked="" type="checkbox"/>	glutathione reductase [Pyrrus ussuriensis x Pyrrus communis]	Pyrrus ussuriensi...	603	603	95%	0.0	56.34%	559	KAB2627620.1
<input checked="" type="checkbox"/>	hypothetical protein C1H46_000718 [Malus baccata]	Malus baccata	602	602	96%	0.0	56.19%	560	TQE13711.1
<input checked="" type="checkbox"/>	hypothetical protein DVH24_010675 [Malus domestica]	Malus domestica	599	599	96%	0.0	55.82%	590	RXH98350.1
<input checked="" type="checkbox"/>	glutathione reductase_chloroplastic-like [Malus domestica]	Malus domestica	599	599	96%	0.0	55.82%	560	XP_028958518.1
<input checked="" type="checkbox"/>	PREDICTED: glutathione reductase_chloroplastic [Pyrrus x bretschnideri]	Pyrrus x bretschn...	595	595	95%	0.0	56.05%	559	XP_009367745.1
<input checked="" type="checkbox"/>	glutathione reductase [Pyrrus ussuriensis x Pyrrus communis]	Pyrrus ussuriensi...	595	595	96%	0.0	55.64%	560	KAB2636727.1
<input checked="" type="checkbox"/>	hypothetical protein DVH24_021775 [Malus domestica]	Malus domestica	595	595	95%	0.0	55.68%	600	RXH86502.1
<input checked="" type="checkbox"/>	glutathione reductase_chloroplastic-like [Malus domestica]	Malus domestica	595	595	95%	0.0	55.68%	560	XP_008393844.1
<input checked="" type="checkbox"/>	hypothetical protein H5410_047701 [Solanum commersonii]	Solanum comme...	594	594	88%	0.0	59.48%	557	KAG5587267.1
<input checked="" type="checkbox"/>	PREDICTED: glutathione reductase_chloroplastic-like [Pyrrus x bretschnideri]	Pyrrus x bretschn...	593	593	95%	0.0	55.78%	560	XP_009354878.1
<input checked="" type="checkbox"/>	hypothetical protein KY289_036763 [Solanum tuberosum]	Solanum tuberos...	592	592	89%	0.0	58.97%	557	KAH0636848.1

The top result was glutathione reductase in *Klebsormidium nitens* highlighted in purple above. Alignment is printed below:

Query: 2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus atmophyticus EST
library Chlorokybus atmophyticus cDNA 5', mRNA sequence Query ID: lcl|Query_53894
Length: 557

>Glutathione reductase [Klebsormidium nitens]
Sequence ID: GAQ78357.1 Length: 580
Range 1: 93 to 579

Score:616 bits(1588), Expect:0.0,
Method:Compositional matrix adjust.,
Identities:307/487(63%), Positives:366/487(75%), Gaps:3/487(0%)

Query	73	ASSNGSGYDYDVITIGAGSGGVRASRIASQLGAKVACVEMPFNNVASDTEGGVGGTCVLR	132
		++ +G +DYD+ TIGAGSGGVRASR ASQ GAKVA E+PF+ ASD +GGVGGTCVLR	
Sbjct	93	STEDGQQFDYDLFTIGAGSGGVRASRFASQYGAKVAVCELPFSTKASDDKGGVGGTCVLR	152
Query	133	GCVPKKLLVYGSIFSNEFDDSDAGFGWKLPG-EPKFTWQTLNENKNKELTRLNNVYRNILS	191
		GCVPKKLLVYGS F++ F+DS GFGW PG EP+ W L E KNKEL RLNN Y+ L	
Sbjct	153	GCVPKKLLVYGSFADYFEDSRGFGWSFPGGEPEVDWSHLIEKKNKELDRLNNAYKTTLK	212
Query	192	KANVELLEGRASLVDAHTVDIDGKQLTAKNIILATGGRSFALPIPGAHAIDSDKALS LD	251
		A V+L+EG+ ++VD HTVD+DGK+ KNI++ATGGR F PIPGAEH I SD AL L	
Sbjct	213	NAKVDLIEGKGTIVDRHTVDVDGKRFKVNIL IATGGRIFVPPIPGAEHVITSDDALDLT	272
Query	252	EVPKRIAIYGGGYIALEFACIFSGFGAKVDVFYRAPLPLRGFDEEIRNALVEELGKKGIN	311
		VP +IAI GGGYIALEFA IF+ GA+VD+F R LRGFD+E+R L E+L +GI	
Sbjct	273	SVPSKIAIVGGGYIALEFAGIFNSAGAEVDIFVRGDKLLRGFDDEVREFLAEQLQAQGIR	332
Query	312	LHPKCTAEEIRKEANGEYTLKTNCGE-FKADLVMFATGRTPNTKYLNLDAVGVD TTEKGA	370
		+H EI K + TLKT G+ ++ VMFATGR PN K L L+ GVD +K A	
Sbjct	333	IHFGAKPVEIEKRDEDQLTLKTEQGDTWQGSVMFATGRRPNIKGLGLEEAGVDVDDKTA	392
Query	371	IVVDEYSRTTPVNIFAIGDVTNRINLTPVALMEGTAVAKT-IQGEPTKPDHVNVP SAVFT	429
		I VDEYSRT+V NI+A+GDVT+RINLTPVALMEG A AKT Q EPTKPDH NVPSAVFT	
Sbjct	393	IKVDEYSRTSVDNIWAVGDVTDRINLTPVALMEGMAFAKTAFQDEPTKPDHTNVPSAVFT	452
Query	430	QPPIGTAGLTEEEAKEQFDEVDVYTSSFRPMKHTISGRDERSLMKIIVDVKTDKVLGIHM	489
		PPIGT GLTE EA EQ+ +VDV+TS+FRPMK TISG R+ +KI+VD TDKV+G+HM	
Sbjct	453	NPPIGTVGLTEAEAVEQYGDVDVFTSTFRPMKSTISGNPVRTFVKILVDAATDKVIGLHM	512
Query	490	LGESSPEILQGFVAVALKCGATKKQLDATIGIHPTAAEEFVTMRTVTRQHRKEKQQQQQEE	549
		GE PEI+QGFAVA++ G TKKQ+D+T+GIHPT+AEE VTMRT TRQ RKE+ Q + E	
Sbjct	513	CGEDGPEIMQGFAVAVRMGVTKKQMDSTVGIHPTSAAEELVTMRTPTRQIRKEEAQNGKGE	572
Query	550	KEKVAAA 556	
		KE AAA	
Sbjct	573	KEMAAAA 579	

This result matches almost all classification of "novel" for this class project:

- There is no 100% match found in the database for the protein sequence in the original species (*Chlorokybus atmophyticus*).
- The top match reported has less than 100% identity.
- There is no 100% match to a different species (redundant to above).

✗ There are no database matches to the original query when no organism is specified, but when you refine the search to the original query species (*Arabidopsis thaliana*) there is a database match (see below):

BLAST® » blastp suite » results for RID-RORYHSJJ013 [Home](#) [Recent Results](#) [Saved Strategies](#) [Help](#)

[← Edit Search](#) [Save Search](#) [Search Summary ▼](#) [? How to read this report?](#) [▶ BLAST Help Videos](#) [↩ Back to Traditional Results Page](#)

i Your search is limited to records that include: *Arabidopsis thaliana* (taxid:3702)

Job Title	2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus...
RID	RORYHSJJ013 Search expires on 10-22 02:28 am Download All ▼
Program	BLASTP ? Citation ▼
Database	nr See details ▼
Query ID	Icl Query_77954
Description	2451-3941_6 chlorokybus_71836.3_lrc301_c Chlorokybus ...
Molecule type	amino acid
Query Length	557
Other reports	Distance tree of results Multiple alignment MSA viewer ?

Filter Results

Organism only top 20 will appear ☐ exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

Percent Identity to E value to Query Coverage to

[Filter](#) [Reset](#)

[Descriptions](#) [Graphic Summary](#) [Alignments](#) [Taxonomy](#)

Sequences producing significant alignments [Download ▼](#) [New Select columns ▼](#) Show [?](#)

☒ select all 46 sequences selected [GenPept](#) [Graphics](#) [Distance tree of results](#) [Multiple alignment](#) [New MSA Viewer](#)

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	Glutathione reductase, chloroplast precursor [Arabidopsis thaliana]	Arabidopsis thaliana	588	588	93%	0.0	56.90%	565	AAK96868.1
<input checked="" type="checkbox"/>	glutathione reductase [Arabidopsis thaliana]	Arabidopsis thaliana	586	586	93%	0.0	56.71%	565	NP_191026.1
<input checked="" type="checkbox"/>	GR [Arabidopsis thaliana]	Arabidopsis thaliana	586	586	93%	0.0	56.71%	565	OAP01378.1
<input checked="" type="checkbox"/>	glutathione-disulfide reductase [Arabidopsis thaliana]	Arabidopsis thaliana	519	519	83%	2e-180	53.83%	499	NP_001030756.2
<input checked="" type="checkbox"/>	putative glutathione reductase [Arabidopsis thaliana]	Arabidopsis thaliana	519	519	83%	2e-180	53.83%	499	AAK25938.1
<input checked="" type="checkbox"/>	unnamed protein product [Arabidopsis thaliana]	Arabidopsis thaliana	493	493	83%	2e-159	52.77%	1445	CAD5323997.1
<input checked="" type="checkbox"/>	glutathione reductase, cytosolic [Arabidopsis thaliana]	Arabidopsis thaliana	275	275	41%	8e-89	58.01%	242	BAD95212.1
<input checked="" type="checkbox"/>	lipamide dehydrogenase 1 [Arabidopsis thaliana]	Arabidopsis thaliana	186	186	92%	4e-51	30.36%	570	NP_566562.1
<input checked="" type="checkbox"/>	unnamed protein product [Arabidopsis thaliana]	Arabidopsis thaliana	186	186	92%	8e-51	30.36%	623	CAA0382694.1
<input checked="" type="checkbox"/>	lipamide dehydrogenase 1 [Arabidopsis thaliana]	Arabidopsis thaliana	186	186	92%	1e-50	30.36%	623	NP_001078165.1

[Q05] 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.

MAFFT MSA using the GUIDANCE2 server for the novel protein [*Chlorokybus atmopyticus*], the original query [*Arabidopsis thaliana*], and the top 18 sequences (e-value) from a BLASTx using the novel protein cDNA (5, 2).

GUIDANCE2 parameters:

Sequences: All accession numbers in MSA
Sequence Type: Amino Acids
MSA Algorithm: MAFFT
No. of bootstrap repeats: 100
Output order: Aligned
MAFFT Options
{
Max-Iterate: 0
Pairwise alignment method: 6mer
}

BLASTx parameters:

Method: NCBI BLASTX (v2.12.0) searched against all
reference proteins using HO413015.1
Database: Reference proteins [refseq_protein] database
Organism: None

MSA:

Malus domestica [XP_028958518.1]	-----
Pyrus x bretschneideri [XP_009367745.1]	-----P--
Pistacia vera [XP_031258822.1]	-----
Gossypium hirsutum [XP_016709857.2]	-----
Rosa chinensis [XP_024166425.1]	-----
Fragaria vesca subsp. vesca [XP_004289892.1]	-----
Arachis hypogaea [XP_025604791.1]	-----
Arachis duranensis [XP_015969716.1]	-----
Arachis ipaensis [XP_016204709.1]	-----
Prosopis alba [XP_028758428.1]	-----
Prunus persica [XP_007199772.1]	-----P--
Solanum tuberosum [XP_006349748.1]	-----
Solanum lycopersicum [NP_001234243.2]	-----
Solanum pennellii [XP_015087986.1]	-----
Capsicum annuum [XP_016565503.1]	-----
Sesamum indicum [XP_011079991.1]	-----
Arabidopsis lyrata subsp. lyrata [XP_020880886.1]	-----
Camelina sativa [XP_010427174.1]	-----
Arabidopsis thaliana [NP_191026.1]	-MASTPKLTSTISSSSPSLQFLCKKLPIAIHLPSS
Chlorokybus atmophyticus [tHO_413015.1]	ALLPAPRVA AVATRSSSSRR---SELPGLLARP--
XP_028958518.1	-----LSLPKTLTSLSHLR---RTSTSHPHHH-----LNSRRHFSI--RASDSGNGA-DSTRHYDFDLF
XP_009367745.1	-----LSLPKTLSPLSHLR---RTSTSHPHHH-----LHSRRRFSV--RAVDSGNGA-DSTRHYDFDLF
XP_031258822.1	-----YDFDLF
XP_016709857.2	-----SRFHLHHHHTPRFHPRRLFV--RA-ESENGA-EPLRHYDFDLF
XP_024166425.1	-----YDFDLF
XP_004289892.1	-----SDNGAGDPSRHYDFDLF
XP_025604791.1	-----HH---TRTRRSFTL--SA--SANPH-----NYDFDLF
XP_015969716.1	-----HH---TRTRRSFTL--SA--SANPH-----NYDFDLF
XP_016204709.1	-----HH---TRTRRSFTL--SA--SANPH-----NYDFDLF
XP_028758428.1	-----YDFDLF

XP_007199772.1 -----LSLRRTLTSLSHLH---RTSISPLHHH-----SRRRSFSV--RA-DSGNGA-DSGSHYDYDLF
XP_006349748.1 -----HPR--TSSLSYGRRFTTT-RA-ESSNGA-ETPRHYDFDLF
NP_001234243.2 -----HTR--TSSLSYGRRFTTP-RA-ESSNGA-ETPRHYDFDLF
XP_015087986.1 -----HTR--TSSLSYGRRFTTP-RA-ESSNGV-ETPRHYDFDLF
XP_016565503.1 -----YDFDLF
XP_011079991.1 -----SSPHPR-----RAFTTSIRA-DSTNGS-EPPRNYDFDLF
XP_020880886.1 -----LSLPKTLTSLYSLRP---RIAVLSNHRYY-----HSRRFSV--RA-STDNGA-DSEHYDFDLF
XP_010427174.1 -----LSLPKTLTSLYSLRP---RVAVLSNHRYY-----HHSRRFSV--SA-SSDNGT-DSEHYDFDLF
NP_191026.1 SSSSFLSLPKTLTSLYSLRP---RIALLSNHRYY-----HSRRFSV--CA-STDNGA-ESDRHYDFDLF
tHO_413015.1 -----LGVSRFSRFGFSGLRPSGLAQGQSRPHSRQ---QSRTVQRYGL--RVIASSNGS-----GYDYDVI

XP_028958518.1 TIGAGSGGVRASRFAANFGASVAVCELPFSTISSD TTGGVGGTCVLRGCVPKKLLVYASKFAHEFEESNG
XP_009367745.1 TIGAGSGGVRASRFAANFGASVAVCELPFSTISSA TAGGVGGTCVLRGCVPKKLLVYASKFAHEFEESNG
XP_031258822.1 TIGAGSGGVRASRFASNFGASVAVCELPFSTISSE TTGGVGGTCVLRGCVPKKLMVYASKFSHEFDESNG
XP_016709857.2 TIGAGSGGVRASRFAANFGASVAVCELPFSTISSE TTGGVGGTCVLRGCVPKKLMVYASKYSHEFDESNG
XP_024166425.1 TIGAGSGGVRASRFASNFGAKVAVCELPFAAISSD TTGGVGGTCVLRGCVPKKLLVYASKYTHEFDDSIG
XP_004289892.1 TIGAGSGGVRASRFAANFGAKVAVCELPFATISSE TAGGVGGTCVLRGCVPKKLMVYASKYSHEFEDSIG
XP_025604791.1 TIGAGSGGVRASRFAANYGASVAICELPFSTISSD VTGGVGGTCVLRGCVPKKLLVYSSKYSHEFEESNG
XP_015969716.1 TIGAGSGGVRASRFAANYGASVAICELPFSTISSD VTGGVGGTCVLRGCVPKKLLVYSSKYSHEFEESNG
XP_016204709.1 TIGAGSGGVRASRFAANYGASVAICELPFSTISSD VTGGVGGTCVLRGCVPKKLLVYSSKYSHEFEESNG
XP_028758428.1 TIGAGSGGVRASRFAANFGASVAICELPFSTISSD TTGGVGGTCVLRGCVPKKLLVYSSKYAHEFEESNG
XP_007199772.1 TIGAGSGGVRASRFAANFGASVAICELPFATIASD TAGGVGGTCVLRGCVPKKLLVYASQFAHEFEESNG
XP_006349748.1 TIGAGSGGVRASRFASNFGASVAVCELPFSTISSD STGGVGGTCVLRGCVPKKLLVYASKYSHEFEESCG
NP_001234243.2 TIGAGSGGVRASRFASNFGASVAVCELPFSTISSD STGGVGGTCVLRGCVPKKLLVYASKYSHEFEESCG
XP_015087986.1 TIGAGSGGVRASRFASNFGASVAVCELPFSTISSD STGGVGGTCVLRGCVPKKLLVYASKYSHEFEESCG
XP_016565503.1 TIGAGSGGVRASRFASNFGASVAVCELPFSTISSD STGGVGGTCVLRGCVPKKLLVYASKYSHEFEESCG
XP_011079991.1 TIGAGSGGVRASRFAANFGAKVAVCELPFATISSD TTGGVGGTCVLRGCVPKKLLVYASKFSHEFQESCG
XP_020880886.1 TIGAGSGGVRASRISTSFGASAAVCELPFSTISSD TAGGVGGTCVLRGCVPKKLLVYASKYSHEFEDSHG
XP_010427174.1 TIGAGSGGVRASRFATSFGASAAVCELPFSTISSD TAGGVGGTCVLRGCVPKKLLVYASKYSHEFEDSHG
NP_191026.1 TIGAGSGGVRASRFATSFGASAAVCELPFSTISSD TAGGVGGTCVLRGCVPKKLLVYASKYSHEFEDSHG
tHO_413015.1 TIGAGSGGVRASRIASQLGAKVACVEMPFN NVASDTEGGVGGTCVLRGCVPKKLLVYGSIFSNEFDD SAG

XP_028958518.1 FGWRYETEPKHDWSTLLANKNAELQRLTGIYKNVLKNANVT LIEGRGKIVDPHTVDVDGKLYSARHILVS
XP_009367745.1 FGWRYETEPKHDWSTLIANKNAELKRLTGIYKNVLSNANVT LIEGRGKIVDPHTVDVEGKLYSARHILVS
XP_031258822.1 FGWKYEVEPQHDWSTLIANKNAELQRLTGIYKNVLKNANVT LLEGRGKIVDPHTVDVDGKLYTARHILIS
XP_016709857.2 FGWKYDTEPKHDWSTLMANKNAELQRLTGIYKNILNKAGVT LIEGRGKVDPHTVDVDGKLYTARHILIS
XP_024166425.1 YGWKYETEPKHDWSTLMANKNAELQRLTGIYKNVLKNANVT LVQGRGKVDPHTVDVDGKLYSARHILIS
XP_004289892.1 YGWKYETEPKHDWSTLMANKNAELQRLTGIYKNVLKNANVS LLEGRGKIVDPHTVDVDGQLYSARHILIS
XP_025604791.1 FGWKYESEPKHDWSTLIANKNAELQRLTGIYKNILKNSNVQ LIEGRGKIVDPHTVDVDGKLYTARHILVS
XP_015969716.1 FGWKYESEPKHDWSTLIANKNAELQRLTGIYKNILKNSNVQ LIEGRGKIVDPHTVDVDGKLYTARHILVS
XP_016204709.1 FGWKYESEPKHDWSTLIANKNAELQRLTGIYKNILKNSNVQ LIEGRGKIVDPHTVDVDGKLYTARHILVS
XP_028758428.1 FGWKYDSEPKHDWSTLMANKNAELQRLTGIYKNILKNSDV K LIEGRGKIVDPHTVDVDGKLYSARHIIIS
XP_007199772.1 FGWRYETEPKHDWSTLMANKNAELQRLTGIYKNVLKNAGVA LIEGRGKIVDPHTVDVDGKLYSARHIIVS
XP_006349748.1 FGWNYEAEPKHDWSILANKNAELQRLTGIYKNILKNADV T LIEGRGKVDPHTVDVDGKLYSAKNILIS
NP_001234243.2 FGWNYEAEPKHDWSTLIANKNAELQRLTGIYKNILKNADV T LIEGRGKVDPHTVDVDGKLYSAKNILIS
XP_015087986.1 FGWNYEAEPKHDWSTLIANKNAELQRLTGIYKNILKNADV T LIEGRGKVDPHTVDVDGKLYSAKNILIS
XP_016565503.1 FGWNYEAEPKHDWNTLIANKNAELQRLMGIYKNILKNANVT LIEGRGKVDPHTVDVDGKLYSAKNILIS
XP_011079991.1 FGWNYEAEPKHDWSTLIANKNAELQRLTGIYKNILKNAGVA LIEGRGKIVDPHTVEVNGKLYSAKNILIS
XP_020880886.1 FGWKYETEPSHDWTTLIANKNAELQRLTGIYKNILSKANVK LIEGRGKVIDPHTVDVDGKIYTRNIIIA
XP_010427174.1 FGWKYDTEPTHWSTLIANKNAELQRLTGIYKNILSKANVK LIEGRGKVIDPHTVDVDGKIYTRNIIIA
NP_191026.1 FGWKYETEPSHDWTTLIANKNAELQRLTGIYKNILSKANVK LIEGRGKVIDPHTVDVDGKIYTRNIIIA
tHO_413015.1 FGWKLPGEPKFTWQTLNENKNKELTRLNNVYRNILSKANV ELLEGRASLVDAHTVDIDGKQLTAKNIIIA

XP_028958518.1 VGGRPFIPFIPGSEYAI DSDAALDLP TKPEKIAIVGGGYIAVEFAGIFNGLTSDVHVFI RQKKVLRGFDE
XP_009367745.1 VGGRPFIPFIPGSEYAI DSDAALDLP TKPEKIAIVGGGYIAVEFAGIFNGLTSDVHVFI RQKKVLRGFDE
XP_031258822.1 VGGRPFIPDIPGSEHAIDSDAALDLP SKPKKIAIIGGGYIALEFAGIFNGLTSDVHVFI RQKKVLRGFDE
XP_016709857.2 VGGRPFIPDIPGSEYAI DSDAALDLP SKPEKVAIVGGGYIALEFAGIFNGLTSEVHVFI RQKKVLRGFDE
XP_024166425.1 VGGRPFIPDIPGSEYAI DSDAALDLP DPKGKIAIVGGGYIALEFAGIFNGLRSDVHVFI RQKQVLRGFDE
XP_004289892.1 VGGRPFIPDIPGSKYAI DSDAALDLP ERPGKIAIVGGGYIALEFAGIFNGLKSDVHVFI RQKQILRGFDD

XP_025604791.1 VGGRRPFIPDIPGSEHAIDSDAALDLPSPKEKIAIVGGGYIALEFAGIFNGLKSDVHVFIROKKVLRGFDE
XP_015969716.1 VGGRRPFIPDIPGSEHAIDSDAALDLPSPKEKIAIVGGGYIALEFAGIFNGLKSDVHVFIROKKVLRGFDE
XP_016204709.1 VGGRRPFIPDIPGSEHAIDSDAALDLPSPKEKIAIVGGGYIALEFAGIFNGLKSDVHVFIROKKVLRGFDE
XP_028758428.1 VGGRRPFIPDIPGREYAIDSDAALDLPSPKEKIAIVGGGYIALEFAGIFNGLASEVHVFIROKKVLRGFDE
XP_007199772.1 VGGRRPFIPDIPGSEYAIDSDAALDLPSPKPKIAIVGGGYIAVEFAGIFNGLSSDVHVFIROKKVLRGFDE
XP_006349748.1 VGGRRPFIPDIPGSEYAIDSDAALDLPKPKIAIVGGGYIALEFAGIFNGLKSEVHVFIROKKVLRGFDE
NP_001234243.2 VGGRRPFIPDIPGSEYAIDSDAALDLPKPKIAIVGGGYIALEFAGIFNGLKSEVHVFIROKKVLRGFDE
XP_015087986.1 VGGRRPFIPDIPGSEYAIDSDAALDLPKPKIAIVGGGYIALEFAGIFNGLKSEVHVFIROKKVLRGFDE
XP_016565503.1 VGGRRPFIPNIPGSEYAIDSDAALDLPKPKIAIVGGGYIALEFAGIFNGLTSEVHVFIROKKVLRGFDE
XP_011079991.1 VGGRRPFIPDIPGREYVIDSDAALDLPSPKTKIAIVGGGYIALEFAGIFNGLTSSVHVFIROKKVLRGFDE
XP_020880886.1 VGGRRPFIPDIPGKEFAIDSDAALDLPSPKPKIAIVGGGYIALEFAGIFNGLNSEVHVFIROKKVLRGFDE
XP_010427174.1 VGGRRPFIPDIPGKEFAIDSDAALDLPSPKPKIAIVGGGYIALEFAGIFNGLNSEVHVFIROKKVLRGFDE
NP_191026.1 VGGRRPFIPDIPGKEFAIDSDAALDLPSPKPKIAIVGGGYIALEFAGIFNGLNCEVHVFIROKKVLRGFDE
tHO_413015.1 TGGRSFALP IPGAEHAIDSDKALSDEVKRIAIYGGGYIALEFACIFSGFGAKVDVFYRAPLPLRGFDE

XP_028958518.1 EVRDFVQEQMALRGIEFHTEESPQAIKADGSLSLKTNKGTIEGF SHIMFATGRRPNTKNLGLAIGVK
XP_009367745.1 EVRDFVQEQMALRGIEFHTEESPQAIKADGSLSLKTNKGTIEGF SHIMFATGRRPNTKDLGLEAVGVK
XP_031258822.1 EIRDFVAEQMSVRGIEFHTEESPEAILKSADGSLSLKTNKGTVEGF SHIMFATGRRPNTKNLGLKVGK
XP_016709857.2 EIRDFVGEQMALRGIQFHTEESPQAIKADGSLSLKTNKGTIEGF SHIMFATGRRPNTKNLGLSVGVK
XP_024166425.1 EIRDFVSEQMSVRGIEFHTEESPQAILKSADGSLSLKTNKGTVEGF SHVMFATGRRPNTKNLGLLEVGVK
XP_004289892.1 EIRDFLAEQMSLRGIEFHTEESPQAILKSSDGSFSLKTNKGTVEGF SHVMFATGRRPNTKNLGLAIGVK
XP_025604791.1 EIRDFVGEQMALRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SYIMFATGRRPNTKNIGLESVGK
XP_015969716.1 EIRDFVGEQMALRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRRPNTKNIGLESVGK
XP_016204709.1 EIRDFVGEQMALRGIEFHTEESPQAVKADGSLSLKTNKGTVEGF SHIMFATGRRPNTKNIGLESVGK
XP_028758428.1 EVRDFVSEQMAIRGIEFHVEETPQAIKADGSLSLKTNKGTVEGF SHIMFATGRPTNTKNLGLSVGVK
XP_007199772.1 EVRDFVQEHMSLRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRRPNTKNLGLLEVGVK
XP_006349748.1 EIRDFVGEQMSLRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRSPNTKNLGLDVTGVK
NP_001234243.2 EIRDFVGEQMSLRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRSPNTKNLGLDVTGVK
XP_015087986.1 EIRDFVGEQMSLRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRSPNTKNLGLDVTGVR
XP_016565503.1 EIRDFVGEQMSLRGIEFHTEESPQAIKADGSLSLKTNKGTVEGF SHIMFATGRRPNTKNLGLDVTGVK
XP_011079991.1 EIRDFVGEQMSLRGIEFHTEETPQAIKADGSLSLKTNKGTVDGF SHVMFATGRRPNTKNLGLAIGVK
XP_020880886.1 DVRDFVGEQMSLRGIEFHTEESPEAIKADGSLSLKTSKGTVEGF SHVMFATGRKPNTKNLGLNVGVK
XP_010427174.1 DVRDFVGEQMSLRGIEFHTEESPEAIKADGSLSLKTSKGTVEGF SHVMFATGRKPNTKNLGLNVGVK
NP_191026.1 DVRDFVGEQMSLRGIEFHTEESPEAIKADGSLSLKTSKGTVEGF SHVMFATGRKPNTKNLGLNVGVK
tHO_413015.1 EIRNALVEELGKKGINLHPKCTAEEIRKEANGEYTLKTNCGEFKA-DLVMFATGRPTNTKYNLDAVGVD

XP_028958518.1 LSKNGAIEVDEFSTRTEVPSIWAIGDVTDRVNLTPALMEGGAIKTLFLNEPTKPDYRAVPSAVFSQPP I
XP_009367745.1 LSKNGAIEVDKFSRTAVPSIWAIGDVTDRVNLTPALMEGGAIKTLFLNEPTMPDYRAVPSAVFSQPP I
XP_031258822.1 MSKNGAIEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTLFQDEPTKPDYRAVPSAVFCQPP I
XP_016709857.2 INKNGAIEVDEYSRTTVP SIWAIGDVTDRINLTPALMEGGAALAKTLFQNEPTKPDYRAVPSAVFSQPP I
XP_024166425.1 IAKNGAIEVDEFSTRTSVPSIWAIGDVTDRVNLTPALMEGGALAKTLFLNEPTKPDYRAIPSAVFSQPP I
XP_004289892.1 MANSNGAIEVDEFSTRTSVPSIWAIGDVTDRVNLTPALMEGGALAKTLFLNEPTKPDYRAIPSAVFSQPP I
XP_025604791.1 IDKKGAIEVNEYSQSSVPSIWAIGDVTDRINLTPALMEGVALAKTLFLNEPTKPDYSYVPSAVFSQPP I
XP_015969716.1 IDKKGAIEVNEYSQSSVPSIWAIGDVTDRINLTPALMEGVALAKTLFLNEPTKPDYSYVPSAVFSQPP I
XP_016204709.1 IDKKGAIEVNEYSQSSVPSIWAIGDVTDRINLTPALMEGVALAKTLFLNEPTKPEYSYVPSAVFSQPP I
XP_028758428.1 TAKNGAIEVDEYSQTSVPSIWAIGDVTDRMNLTPALMEGMALAKTLFQNNPTKPDYRAVPSAVFSQPP I
XP_007199772.1 LSKTGAIEVDEFSTRTSVPSIWAIGDVTDRVNLTPALMEGGALAKTLFLNEPTKPDYRAVPSAVFSQPP I
XP_006349748.1 MTKNGAIEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTIFAGEPTKPDYRNVPCAVFSQPP I
NP_001234243.2 MTKNGAIEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTIFAGEPTKPDYRNVPCAVFSQPP I
XP_015087986.1 MTKNGAIEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTIFAGEPTKPDYRNVPCAVFSQPP I
XP_016565503.1 MARNGAIEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTIFAGEPTKPDYRNVPCAVFSQPP I
XP_011079991.1 LSKNGAVEVDEYSRTSVPSIWAIGDVTDRINLTPALMEGGALAKTLFANPTKPDFSNVPSAVFSQPP I
XP_020880886.1 MAKNGAIEVDEYSQTSVPSIWAIGDVTDRINLTPALMEGGALAKTLFQNEPTKPDYRAVPCAVFSQPP I
XP_010427174.1 MAKNGAIEVDEYSQTSVPSIWAIGDVTDRINLTPALMEGGALAKTLFQNEPTKPDYRAVPCAVFSQPP I
NP_191026.1 MAKNGAIEVDEYSQTSVPSIWAIGDVTDRINLTPALMEGGALAKTLFQNEPTKPDYRAVPCAVFSQPP I
tHO_413015.1 TTEKGAIIVDEYSRTTVPNIFAIGDVTNRINLTPALMEGTAVAKTI-QGEPTKPDHVNVP SAVFTQPP I

XP_028958518.1 GQVGLSEEQAVEQYGDVDIYTSNFRPLKATVSGLPDRTFMKLIVCAKTNKVLGLHMCGEDSPEIVQGF AV
XP_009367745.1 GQVGLSEEQAVEQYGDVDIYTSNFRPLKATLSGLPDRTFMKLIVCAKTNKVLGLHMCGEDSPEIVQGF AV

```

XP_031258822.1 GQVGLSEEQAIQEYGDIDVFTANFRPLKATLSGLPDRVFMKLIVCAKTNKVLGLHMCGEDAPEIVQGFAV
XP_016709857.2 GQVGLTEEQARKEYGDIDVYTANFRPLKATLSGLPDRVFMKLIVCAKTNKVLGLHMCGEDSAEIAQGFAV
XP_024166425.1 GQVGLSEEQATEQYGDVDIYTSNFRPMKATLSGLPDRVFMKLIVCAKTNKLLGLHMCGEDSPEIVQGFAV
XP_004289892.1 GQVGLSEEQATEQYGDVDIYTSNFKPMKATLSGLPDRVFMKLIVCAKTNKILGLHMCGDDSP EIVQGFAV
XP_025604791.1 GQVGLTEEQAVEQYGDVDIFTSNFRPLKATLSGLPDRTFMKLIVCAKTNKVLGLHMCGEDSPEITQGFAV
XP_015969716.1 GQVGLTEEQAVEQYGDVDIFTSNFRPLKATLSGLPDRTFMKLIVCAKTNKVLGLHMCGEDSPEITQGFAV
XP_016204709.1 GQVGLTEEQAVEQYGDVDIFTSNFRPLKATLSGLPDRTFMKLIVCAKTNKVLGLHMCGEDSPEITQGFAV
XP_028758428.1 GQVGLTEEQAVQQYGNVDIFTANFKPLKATLSGLPDRAFMKLIVCAKTNKVLGLHMCGEDSPEIVQGFAV
XP_007199772.1 GQVGLTEEQAIQEYGDVDIYTSNFRPLKATLSGLPDRVFMKLLVCAKTNKVLGLHMCGEDSAEIVQGFAV
XP_006349748.1 GLVGLTEEEAIKEYGDVDVYTANFRPLKATLSGLPDRAFMKLIVCSKTSKVLGLHMCGEDAPEIVQGFAV
NP_001234243.2 GLVGLTEEEAIKEYGDVDVYTANFRPLKATLSGLPDRVFMKLIVCAKSSKVLGLHMCGDDAPEIVQGFAV
XP_015087986.1 GLVGLTEEEAIKEYGDVDVYTANFRPLKATLSGLPDRVFMKLIVCAKSSKVLGLHMCGDDAPEIVQGFAV
XP_016565503.1 GIVGLTEEQAIN EYGDIDVYTTNFRPLKATLSGLPDRVFMKLIVCAKSSKVLGLHMCGEDAPEIVQGFAV
XP_011079991.1 GQVGLTEEQAIKEYGDIDVYTANFRPMKATLSGLPDRVFMKLIVCAKTNKVLGVHMCGEDSPEIIQGFAV
XP_020880886.1 GTVGLTEEQAIQEYGDVDVFTSNFRPLKATLSGLPDRVFMKLIVCANTNKVLGVHMCGEDSPEIIQGFV
XP_010427174.1 GTVGLTEEQAIQEYGDVDVYTSNFRPLKATLSGLPDRVFMKLIVCANTNKVGVHMCGEDSPEIIQGFV
NP_191026.1 GTVGLTEEQAIQEYGDVDVYTSNFRPLKATLSGLPDRVFMKLIVCANTNKVLGVHMCGEDSPEIIQGFV
tHO_413015.1 GTAGLTEEEAKEQFDEVVDVYTSNFRPMKHTISGRDERSLMKIIIVDVKTDKVLGIHMLGESSPEILQGFAV

XP_028958518.1 AVKAGLTKADLDSTIGIHPTAAEEFVTMRTPTRKIR-----
XP_009367745.1 AVKAGLTKADLDSTIGIHPTAAEEFVTMRTPTRKIR-----
XP_031258822.1 AVKAGLTKADFDTTVG IHP TAAEEFVTMRTPTRKIR-----
XP_016709857.2 AVKAGLTKADFDATVGIHPTSAEEFVTMRTPTRKIR-----
XP_024166425.1 AVKAGLTKADLDATIGIHPTAAEEFVTMRTPTRKIR-----
XP_004289892.1 AVKAGLTKADLDATIGIHPTAAEELVTMRTPTRKIR-----
XP_025604791.1 AIKAGLTKGDFDATVGIHPTAAEEFVTMRTPTRKIR-----
XP_015969716.1 AIKAGLTKGDFDATVGIHPTAAEEFVTMRTPTRKIR-----
XP_016204709.1 AIKAGLTKADFDATVGIHPTAAEEFVTMRTPTRKIR-----
XP_028758428.1 AVKAGLTKADFDATVG VHP TAAEEFVTMR TTRKIR-----
XP_007199772.1 VVKAGLTKADLDATIGIHPTAAEEFVTMRTPTRKIR-----
XP_006349748.1 AVKAGLTKADFDATVGIHPTAAEEFVTMRTPTRKVR-----
NP_001234243.2 AVKAGLTKADFDTTVG IHP TAAEEFVTMRTPTRKIR-----
XP_015087986.1 AVKAGLTKADFDTTVG IHP TAAEEFVTMRTPTRKIR-----
XP_016565503.1 AVKAGLTKADFDATVGIHPTAAEEFVTMRTPTRKVR-----
XP_011079991.1 AVKAGLTKADFDATVGIHPTAAEELVTMRTPTRKIR-----
XP_020880886.1 AVKAGLTKADFDATVG VHP TAAEEFVTMRTPTRKIR-----
XP_010427174.1 AVKAGLTKADFDATVG VHP TASEEFVTMRTPTRKIR-----
NP_191026.1 AVKAGLTKADFDATVG VHP TAAEEFVTMRAPTRKFRKDSSEGKASPEAKTAAGV
tHO_413015.1 ALKCGATKKQLDATIGIHPTAAEEFVTMRVTRQHRKEKQQQQQEEKEKVAAAK

```

[Q06] 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.

Phylogenetic tree parameters for 'Simple Phylogeny' on the EBI (3):

```

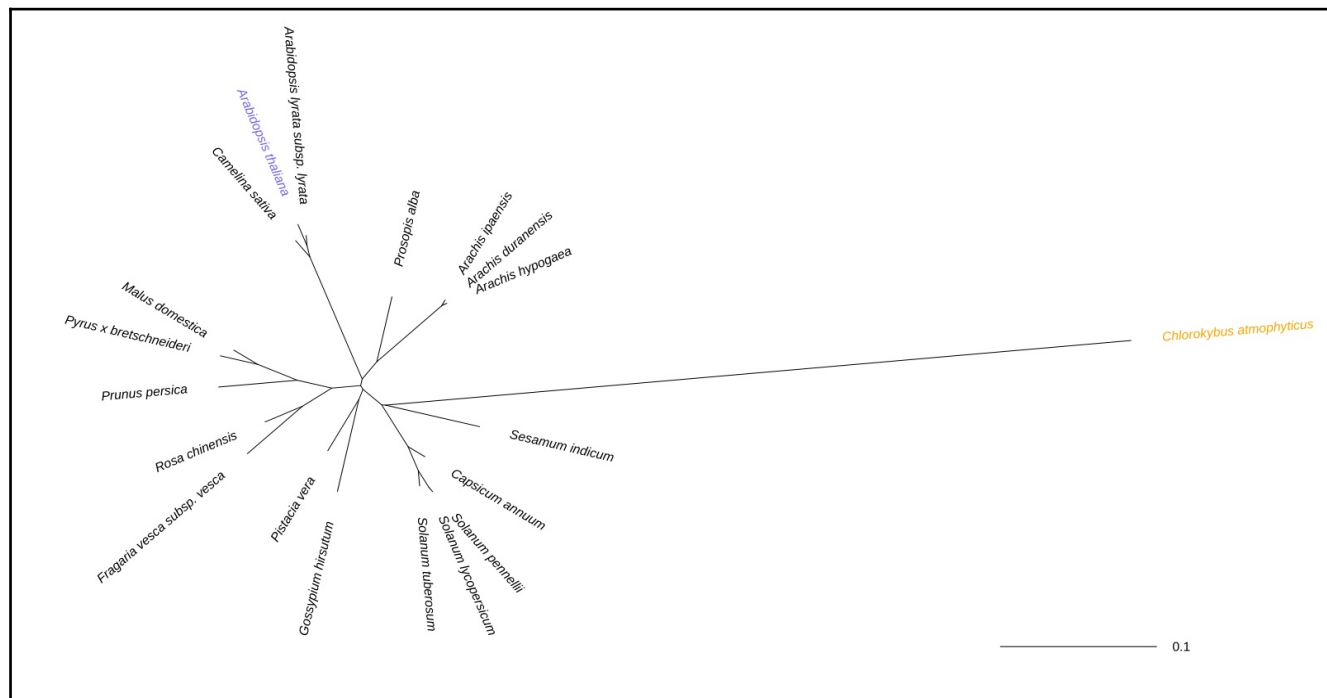
MSA: See Q05 above
TREE FORMAT: Default
DISTANCE CORRECTION: on
EXCLUDE GAPS: off
CLUSTERING METHOD: Neighbour-joining
P.I.M.: off

```

Neighbour-joining unrooted tree for alignment.

Original query protein: *Arabidopsis thaliana*

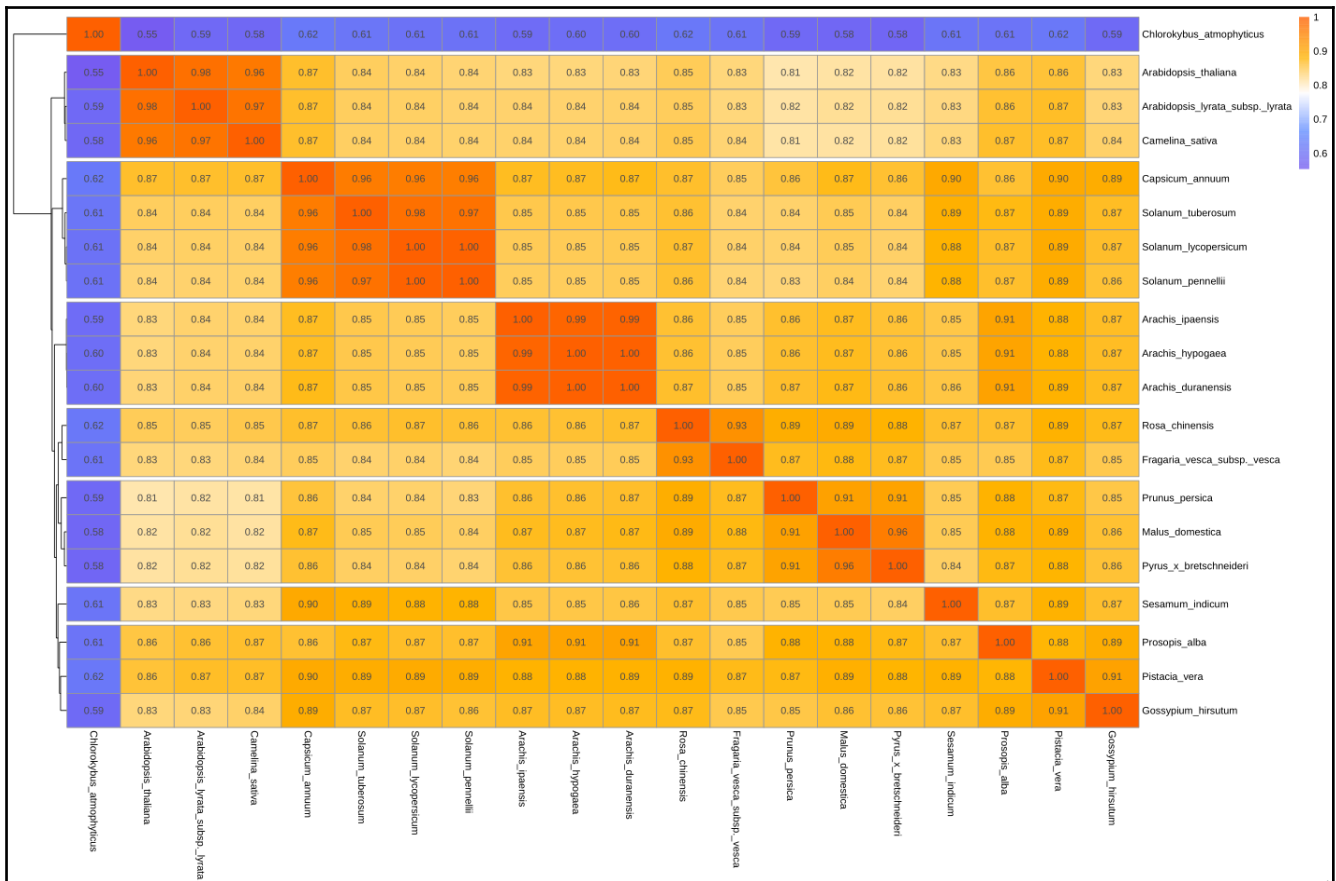
Novel protein: *Chlorokybus atmophyticus*



Made using APE in R [Appendix I] (6, 7).

[Q07] Generate a sequence identity based heatmap of your aligned sequences using R.

Percent sequence identity matrix for the MSA (Q05).



Sequence identities obtained using Bio3D in R [Appendix II] (8, 7).
Heatmap generated using pheatmap in R [Appendix III] (9, 7).

[Q08] 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).

Queried the PDB using the Bio3D package in R [Appendix IV] (10, 8, 7).

```

=====Hits for novel protein=====
*****BEGIN*****
Hit 01 - 300H_A (glutathione reductase) (11):
    E-Value [1.61e-124]
    Sequence ID [42.763%]
    Annotations
    {
        PDB identifier [300H]
        Method [X-RAY DIFFRACTION]
        Resolution [1.90 Å]
        Source Organism [Bartonella henselae str. Houston-1]
    }

Hit 02 - 4DNA_A (putative glutathione reductase) (12):
    E-Value [3.31e-121]
    Sequence ID [45.633%]
    Annotations
    {
        PDB identifier [4DNA]
        Method [X-RAY DIFFRACTION]
        Resolution [2.80 Å]
        Source Organism [Escherichia coli BL21(DE3)]
    }

Hit 03 - 6ER5_A (trypanothione reductase) (13):
    E-Value [7.59e-117]
    Sequence ID [42.766%]
    Annotations
    {
        PDB identifier [6ER5]
        Method [X-RAY DIFFRACTION]
        Resolution [3.37 Å]
        Source Organism [Leishmania infantum]
    }
*****END*****
=Hits for highest intersequence similarity protein (Capsicum annuum)=
*****BEGIN*****
Hit 01 - 300H_A (glutathione reductase) (11):
    E-Value [5.61e-131]
    Sequence ID [44.805%]
    Annotations
    {
        PDB identifier [300H]
        Method [X-RAY DIFFRACTION]
        Resolution [1.90 Å]
        Source Organism [Bartonella henselae str. Houston-1]
    }

```

Hit 02 - 4DNA_A (putative glutathione reductase) (12):

E-Value [3.09e-124]

Sequence ID [44.516%]

Annotations

{

PDB identifier [4DNA]

Method [X-RAY DIFFRACTION]

Resolution [2.80 Å]

Source Organism [*Escherichia coli* BL21(DE3)]

}

Hit 03 - 5VDN_A (glutathione reductase) (14):

E-Value [3.17e-120]

Sequence ID [45.652%]

Annotations

{

PDB identifier [5VDN]

Method [X-RAY DIFFRACTION]

Resolution [1.55 Å]

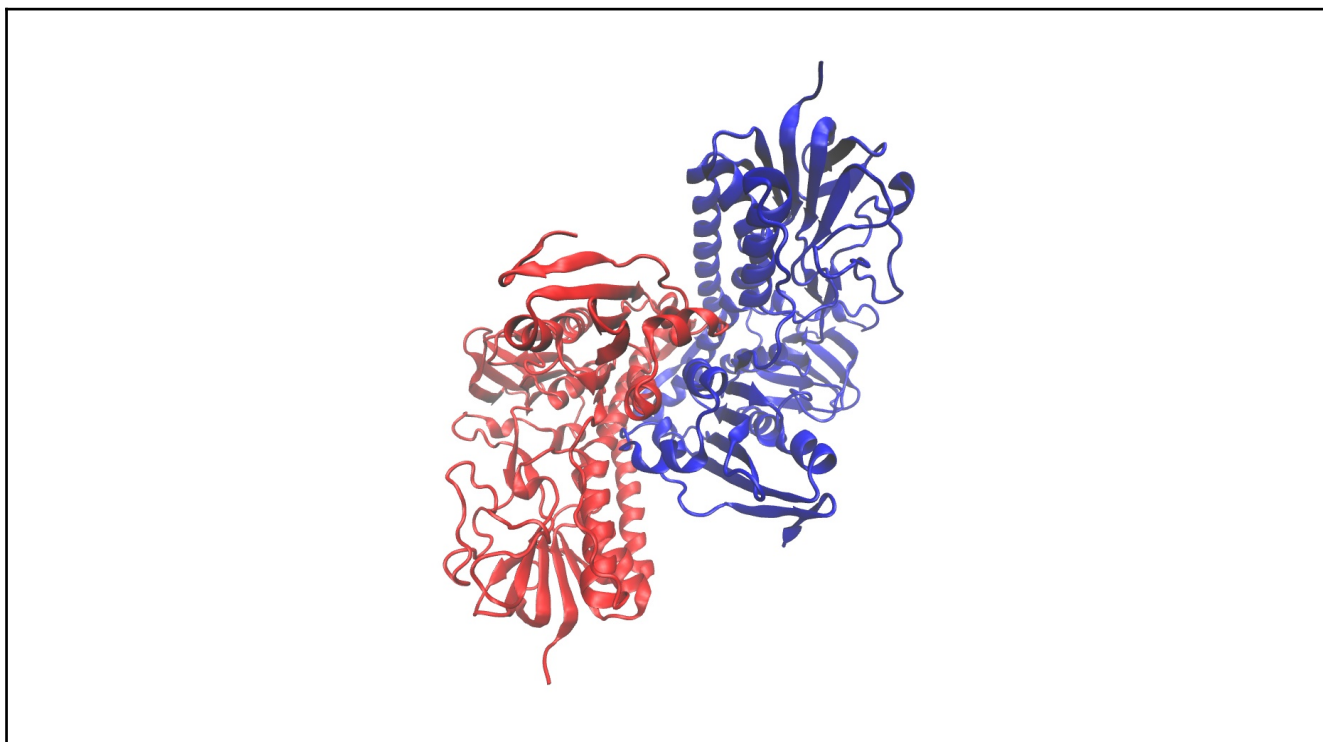
Source Organism [*Yersinia pestis* KIM10+]

}

*****END*****

[Q09] Generate a molecular figure of one of your identified PDB structures using VMD. 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?

Structure for 4DNA, Q08:Hit 02 for novel protein (12).



Generated using VMD (15).

Based on the sequence similarity, percent identity of 45.633%, this structure is likely to be similar to the novel protein. This further implies that the function is likely to be similar between the novel protein and that of this hit, glutathione reductase activity.

[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?

Yes, hits were found for both recombinant trypanothione reductase inhibitors and glutathione reductase inhibitors (16, 17). Both of which would be helpful in exploring potential inhibition of the novel protein. Additionally, this could help differentiate between whether the novel protein within *Chlorokybus atmophyticus* is closer in

function to glutathione reductase or trypanothione reductase, which is of particular importance for trypanosomal infections (18).

Input:

Sequence: Novel protein sequence (Q03)
BLASTp Parameters: Default

Results:

Recombinant trypanothione reductase (CHEMBL1944501) -
Inhibitors of oxidoreductase activity (top 5 pCHEMBL shown)

CHEMBL ID	IC50	pCHEMBL
CHEMBL4169040	[IC50 = 820.0 nM]	(6.09)
CHEMBL4105040	[IC50 = 900.0 nM]	(6.05)
CHEMBL4205871	[IC50 = 1200.0 nM]	(5.92)
CHEMBL4076896	[IC50 = 1200.0 nM]	(5.92)
CHEMBL4172675	[IC50 = 1200.0 nM]	(5.92)

Glutathione reductase (CHEMBL2755) -
Inhibitors of glutathione activity (top 5 pCHEMBL shown)

CHEMBL ID	IC50	pCHEMBL
CHEMBL2068507	[IC50 = 1.0 nM]	(9.00)
CHEMBL120147	[IC50 = 4.1 nM]	(8.39)
CHEMBL1824793	[IC50 = 344.0 nM]	(6.46)
CHEMBL135536	[IC50 = 750.0 nM]	(6.12)
CHEMBL39225	[IC50 = 1000.0 nM]	(6.00)

References:

1. Protein [Internet]. Bethesda (MD): National Library of Medicine (US), National Center for Biotechnology Information; 2004 - [cited 2021 Dec 03]. Available from: <https://www.ncbi.nlm.nih.gov/protein/>
2. BLAST [Internet]. Bethesda (MD): National Library of Medicine (US), National Center for Biotechnology Information; 2004 - [cited 2021 Dec 03]. Available from: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>
3. Madeira, F., Park, Y.M., Lee, J., Buso, N., Gur, T., Madhusoodanan, N., Basutkar, P., Tivey, A.R., Potter, S.C., Finn, R.D. and Lopez, R., 2019. The EMBL-EBI search and sequence analysis tools APIs in 2019. *Nucleic acids research*, 47(W1), pp.W636-W641.
4. Taxonomy [Internet]. Bethesda (MD): National Library of Medicine (US), National Center for Biotechnology Information; 2004 - [cited 2021 Dec 03]. Available from: <https://www.ncbi.nlm.nih.gov/taxonomy/>
5. Sela, I., Ashkenazy, H., Katoh, K. and Pupko, T., 2015. GUIDANCE2: accurate detection of unreliable alignment regions accounting for the uncertainty of multiple parameters. *Nucleic acids research*, 43(W1), pp.W7-W14.
6. Paradis, E. and Schliep, K., 2019. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics*, 35(3), pp.526-528.
7. R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
8. Grant, B.J., Rodrigues, A.P., ElSawy, K.M., McCammon, J.A. and Caves, L.S., 2006. Bio3d: an R package for the comparative analysis of protein structures. *Bioinformatics*, 22(21), pp.2695-2696.
9. Raivo Kolde (2019). pheatmap: Pretty Heatmaps. R package version 1.0.12. <https://CRAN.R-project.org/package=pheatmap>
10. Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E., 2000. The protein data bank. *Nucleic acids research*, 28(1), pp.235-242.
11. Seattle Structural Genomics Center for Infectious Disease (SSGCID) (2010) Crystal structure of glutathione reductase from *Bartonella henselae* doi: 10.221/pdb3o0h/pdb

12. Malashkevich, V.N., Bhosle, R., Toro, R., Seidel, R., Almo, S.C., New York Structural Genomics Research Consortium (NYSGRG) (2012) Crystal structure of putative glutathione reductase from *Sinorhizobium meliloti* 1021 doi: 10.2210/pdb4dna/pdb
13. Ilari, A., Fiorillo, A. (2018) X-ray structure of Trypanothione Reductase from *Leishmania infantum* in complex with 2-(diethylamino)ethyl 4-((3-(4-nitrophenyl)-3-oxopropyl)amino)benzoate doi: 10.2210/pdb6er5/pdb
14. Minasov, G., Shuvalova, L., Dubrovskaya, I., Cardona-Correa, A., Grimshaw, S., Kwon, K., Anderson, W.F., Center for Structural Genomics of Infectious Diseases (CSGID) (2017) 1.55 Angstrom Resolution Crystal Structure of Glutathione Reductase from *Yersinia pestis* in Complex with FAD doi: 10.2210/pdb5vbn/pdb
15. Humphrey, W., Dalke, A. and Schulten, K., 1996. VMD: visual molecular dynamics. *Journal of molecular graphics*, 14(1), pp.33-38.
16. Mendez, D., Gaulton, A., Bento, A.P., Chambers, J., De Veij, M., Félix, E., Magariños, M.P., Mosquera, J.F., Mutowo, P., Nowotka, M. and Gordillo-Marañón, M., 2019. ChEMBL: towards direct deposition of bioassay data. *Nucleic acids research*, 47(D1), pp.D930-D940.
17. Davies, M., Nowotka, M., Papadatos, G., Dedman, N., Gaulton, A., Atkinson, F., Bellis, L. and Overington, J.P., 2015. ChEMBL web services: streamlining access to drug discovery data and utilities. *Nucleic acids research*, 43(W1), pp.W612-W620.
18. Krauth-Siegel, R.L. and Inhoff, O., 2003. Parasite-specific trypanothione reductase as a drug target molecule. *Parasitology research*, 90(2), pp.S77-S85.

APPENDIX

I:

```
library(ape)

# gor-MAFFT-20-NJ_tree.tre available upon request
NJTree <- read.tree(file = "gor-MAFFT-20-NJ_tree.tre")
col_vect <- rep("black", length(NJTree$tip.label))
col_vect[grep("Chlorokybus", NJTree$tip.label)] <- "orange"
col_vect[grep("thaliana", NJTree$tip.label)] <- "slateblue2"

plot(NJTree,
      edge.width = 1,
      label.offset = 0.02,
      no.margin = TRUE,
      underscore = FALSE,
      lab4ut = "axial",
      align.tip.label = TRUE,
      type = "u",
      cex = 1,
      lwd = 2,
      rotate.tree = -40,
      font = 3,
      tip.color = col_vect)
add.scale.bar(x = 0.5, y = -0.1, lwd = 1)
```

II:

```
library(bio3d)

# gor-MAFFT_aln-20_short.faa available upon request
aln <- read.fasta("gor-MAFFT_aln-20_short.faa")
seq_id <- seqidentity(aln)
```

III:

```
library(pheatmap)

IBMColorBlindSafe = c("#785ef0", "#648fff", "#ffffff", "#ffb000",
                      "#fe6100")
colBlindScale <- colorRampPalette(IBMColorBlindSafe)
hmColors <- colBlindScale(512)

pheatmap(seq_id,
          color = hmColors,
          display_numbers = TRUE,
          fontsize_number = 10,
          cutree_rows = 8,
          treeheight_col = 0)
```

IV:

```
library(bio3d)
```

```
# chlorokybus_gor_seq.faa available upon request
```

```
hits <- blast.pdb(read.fasta("chlorokybus_gor_seq.faa"))
```

```
blast_hits <- plot.blast(hits)
```

```
pdb <- get.pdb(blast_hits$pdb.id, path="pdb", split=TRUE, gzip=TRUE)
```

```
pdb <- pdbaln(pdb[1:3], fit = TRUE)
```

```
top_3 <- hits$hit.tbl[1:3,]
```