

BGGN-213: FOUNDATIONS OF BIOINFORMATICS – FIND A GENE PROJECT
<http://thegrantlab.org/bgg213>

Name: Torrey Rhyne
PID: A14397504

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: TATA-box binding protein (TBP)

Species: Homo sapiens

Accession: NP_003185.1

Function: subunit of the basal transcription factor TFIID; binds promoter DNA at the TATA box; plays a major role in the initiation of transcription by RNA polymerase II.

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

Method: TBLASTN (2.14.1) search against Laupala ESTs

Database: Expressed sequence tags (est)

Organism: Laupala (taxid:109023)

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.

Translated BLAST: tblastn

blastn blastp blastx **tblastn** tblastx

TBLASTN search translated nucleotide databases using a protein query. [more...](#)

Enter Query Sequence

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

NP_003185.1

Query subrange [?](#)

From

To

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

Job Title

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

☐ Align two or more sequences [?](#)

Choose Search Set

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

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

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

☒ Show results in a new window

i Your search is limited to records that include: Laupala (taxid:109023)

Job Title **NP_003185:TATA-box-binding protein isoform...**

RID [PUJWP4A4013](#) Search expires on 12-05 10:46 am [Download All](#) [v](#)

Program TBLASTN [?](#) [Citation](#) [v](#)

Database est [See details](#) [v](#)

Query ID [NP_003185.1](#)

Description TATA-box-binding protein isoform 1 [Homo sapiens]

Molecule type amino acid

Query Length 339

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

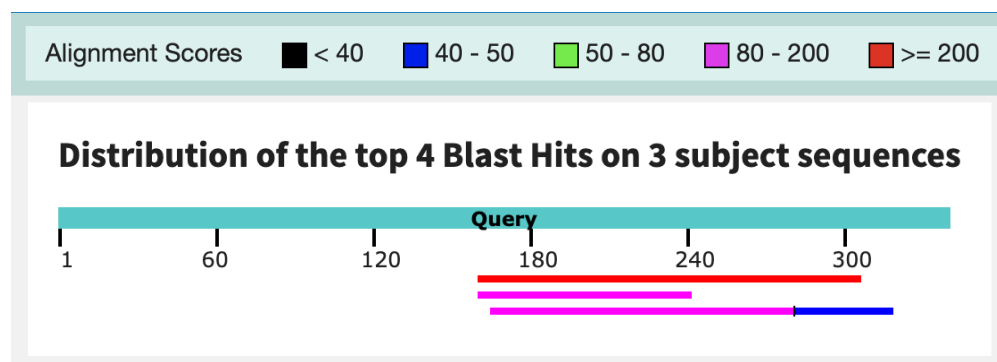
Download [v](#) Select columns [v](#) Show 100 [?](#)

☒ select all 3 sequences selected [GenBank](#) [Graphics](#)

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	EST4852 LK04 Laupala kohalensis cDNA clone 1061021814135 5' mRNA sequence	Laupala kohalensis	288	288	43%	1e-97	93.84%	771	EH633745.1
<input checked="" type="checkbox"/>	EST434 LK04 Laupala kohalensis cDNA clone 1061021306693 5' mRNA sequence	Laupala kohalensis	160	160	24%	8e-48	92.68%	790	EH629327.1
<input checked="" type="checkbox"/>	EST12419 LK04 Laupala kohalensis cDNA clone 1061021810461 5' mRNA sequence	Laupala kohalensis	105	147	45%	7e-30	40.87%	806	EH641311.1

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.

Chosen match: Accession EH633745.1 – a 771 base pair cDNA clone from *Laupala kohalensis*. See below for alignment.



Download ▾ [GenBank](#) [Graphics](#)

EST4852 LK04 *Laupala kohalensis* cDNA clone 1061021814135 5', mRNA sequence

Sequence ID: [EH633745.1](#) Length: 771 Number of Matches: 1

Range 1: 304 to 741 [GenBank](#) [Graphics](#)

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

Score	Expect	Method	Identities	Positives	Gaps	Frame
288 bits(736)	1e-97	Compositional matrix adjust.	137/146(94%)	140/146(95%)	0/146(0%)	+1
Query 160	GIVPQLQNIVSTVNLGCKLDLKTIALRARNAEYNPKRFAAVIMRIREPRTTALIFSSGKM					219
	IVPQLQNIVSTVNLGCKLDLK IAL ARNAEYNPKRFAAVIMRIREPRTTALIFSSGKM					
Sbjct 304	AIVPQLQNIVSTVNLGCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKM					483
Query 220	VCTGAKSEEQSRLAARKYARVVQKLGFPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHQQ					279
	VCTGAKSEE SRLAARKYAR++QKLGFPKFLDFKIQNMVGSCDVKFPIRLEGLVLTH Q					
Sbjct 484	VCTGAKSEEDSRLAARKYARI IQKLGFPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQ					663
Query 280	FSSYEPELFPGLIYRMIKPRIVLLIF					305
	FSSYEPELFPGLIYRM+KPRIVLLI					
Sbjct 664	FSSYEPELFPGLIYRMVKPRIVLLIL					741

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.

> EH633745.1 Laupala kohalensis protein sequence (used EBI EMBOSS Transeq)
GTPLHQPEEDQQILPQAQQQQQQNQQGTPLLQMPLLATPQKIMHTYAPSGFTTPQSL
MQPQTPQNMMSPMVASSSQIIPPSSLGPATPSPMTPMTPSSADPAIVPQLQNIVSTVN
LGCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLA
ARKYARIIQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIY
RMVKPRIVLLILYLESGLTVQS

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

Name: *Laupala* TBP

Species: *Laupala kohalensis* (a cricket)

Eukaryota; Opisthokonta; Metazoa; Eumetazoa; Bilateria; Protostomia; Ecdysozoa;
Panarthropoda; Arthropoda; Mandibulata; Pancrustacea; Hexapoda; Insecta;
Dicondylia; Pterygota; Neoptera; Polyneoptera; Orthoptera; Ensifera; Gryllidea;
Grylloidea; Gryllidae; Trigonidiinae; *Laupala*

[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.

The top hits from a BLASTP search against the NR database are TBP from many other species, but not *Laupala*. The top hit is TBP from *Gryllus bimaculatus* (another cricket). See below for setup, top hits, and alignment.

TATA-box-binding protein [Gryllus bimaculatus]

Sequence ID: [GLH06850.1](#) Length: 294 Number of Matches: 1

Range 1: 16 to 261 [GenPept](#) [Graphics](#)

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

Score	Expect	Method	Identities	Positives	Gaps
453 bits(1166)	6e-159	Compositional matrix adjust.	231/247(94%)	235/247(95%)	1/247(0%)
Query 1	60	GTPLHQPEEDQQILPQAQQQQQNNQOGTPLLQMPLLATPQKIMHTYAPSGFTTPQSLMQP			
Sbjct 16	74	GTPLHQPEEDQQILPQAQQQQQ NNQOGTPL QMPLL +PQK MHTYAP+GFTTPQSLMQP			
Query 61	120	QTPQNMMSPMVASSSQIIPPSSLGPATPSPMTMPMTSSADPAIVPQLQNIIVSTVNLGCKL			
Sbjct 75	134	QTPHNMMSPMVASASQIVAPSSLGPATPGPMTMPMTSSADPGIVPQLQNIIVSTVNLCKL			
Query 121	180	DLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYA			
Sbjct 135	194	DLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYA			
Query 181	240	RIIQKLGFPAPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKP			
Sbjct 195	254	RIIQKLGFPAPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKP			
Query 241	247	RIVLLIL			
Sbjct 255	261	RIVLLIF			

[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.

Re-labeled sequences for alignment:

```
> Human_TBP | NP_003185.1 | TATA-box-binding protein isoform 1 [Homo sapiens]
MDQNNSLPPYAQGLASPPQAMTPGIPFI FSPMPYPGTGLTPQPIQNTNSLSILEEQQRQQQQQQQQQQQQQ
QQQQQQQQQQQQQQQQQQQQQQQVAAA VQQSTSQQATQGTSGQAPQLFHSQTLTTAPLP GTTPLYP
SPMTMPITPATPASESSGIVPQLQNIIVSTVNLGCKLDLKTIALRARNAEYNPKRFAAVIMRIEPRTT
ALIFSSGKMVCTGAKSEEQSRLAARKYARVVQKLGFPAPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQF
SSYEPELFPGLIYRMIPRIVLLIFVSGKVLTGAKVRAEIIYEAFENIYPILKGFRKTT
```

```
> Cricket_(Laupala) | EH633745.1 | translated protein sequence from
EBI EMBOSSTranseq [Laupala kohalensis]
GTPLHQPEEDQQILPQAQQQQQNNQOGTPLLQMPLLATPQKIMHTYAPSGFTTPQSLMQPQTPQNMMSP
VASSSQIIPPSSLGPATPSPMTMPMTSSADPAIVPQLQNIIVSTVNLGCKLDLKKIALHARNAEYNPKRFA
AVIMRIEPRTTALIFSSGKMVCTGAKSEEDSRLAARKYARIIQKLGFPAPKFLDFKIQNMVGSCDVKFPI
RLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLILYLESGLTVQS
```

```
> Cricket_(Gryllus) | GLH06850.1 | TATA-box-binding protein [Gryllus
bimaculatus]
```

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPMAQQQQQQNQGGTPLQQMPLLGSPPQKTMHTYAPAGFTTPQS
LMQPQTPHNMMSPMVASASQIVAPSSLGPATPGPMTMPMPSSADPGIVPQLQNIVSTVNLCKLDLKKIA
LHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYARI IQKLGFPKFLDFK
IQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEA
FDNIYPILKSFKKQ

> Termite | XP_023703955.1 | TATA-box-binding protein [Cryptotermes secundus]

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPHAQQQQQQQQQQQLPPPPTSLQQMPVLSTPQKTMHTYAPSG
FTPQSLMQPQTPQNLMSPMVTTPSSQMATLSNMGPATPSPMTMPMPSSADPGIIPQLQNIVSTVNLGCKLD
LKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYARI IQKLGFPK
FLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLIFVSGKVVLTGAKVRQ
EIYEAFDNIYPILKSFKKQ

> Sawfly_(Athalia) | XP_012257098.1 | TATA-box-binding protein [Athalia rosae]

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPNALQQQQQQVQQSHQIQQMHPMGSNHGMPLATPHKTMH
TYTPAFATPQSLMQPQTPQNMMSPMVQPASQIAPPSIGPSTPGPMTMPMPASADPGILPQLQNIVSTVNL
NCKLELKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYARI IQKL
GFPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLIFVSGKVVLTG
AKVRQEIYEAFDNIYPILKSFKKQ

> Sawfly_(Neodiprion) | XP_015524603.1 | TATA-box-binding protein [Neodiprion lecontei]

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPNALQQQQQQVQQNQLQQMHPMGSNVHGMPLGTPHKTMHTY
APAFATPQSLMQPQTPQNMMSPMVQPASQIAPPSIGPSTPGPMTMPMPASADPGILPQLQNIVSTVNLNC
KLNLKEIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKYARI IQKLG
PAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLIFVSGKVVLTGAK
VRQEIYEAFDNIYPILKSFKKQ

> Wasp_(Venturia) | XP_043276984.1 | TATA-box-binding protein [Venturia canescens]

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPNALQQQQQQQQQQQQQQQTSQQYQMQQLHSMSPTMQANSMM
MLSTPQKTMHTYAPTPTFATPQSLMQPQTPQNMMSPMVQPTSQIAPSSIGPSTPGPMTMPMPASADPGIL
PQLQNIVSTVNLTKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRL
AARKYARI IQKLGFPKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVL
LIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ

> Wasp_(Orussus) | XP_012281073.1 | TATA-box-binding protein [Orussus abietinus]

MDQMLPSPGFSIPSIGTPLHQPEEDQQILPNALQQQQQQQQQQQQSQQYQLQQQLHSMSPNMQGSMLMISTP
QKTMHTYAPTPTAFATPQSLMQPQTPQNMMSPIVQPSSQIAPSSIGPSTPGPMTMPMPASADPGILPQLN
IVSTVNLNCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSRLAARKY
ARI IQKLGFAAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGLIYRMVKPRIVLLIFVS
GKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ

> Ant | XP_011330203.1 | TATA-box-binding protein [Ooceraea biroi]

MDQMLPSPGGFSIPSIGTPLHQPEEDQQILPNALQQQQQQSQQYQLQQQLHSMSPNMQSGMLMITTPQKTMH
TYAPTPTFATPQSLMQPQTPQNMMSPIVQSNQIAPPSIGPATPGPMTMPMPASADPGILPQLQNIVSTV
NLNCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEADSRLAARKYARI IQ

KLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLSHGQFSSYEPELFPGLIYRMVKPRIVLLIFVSGKVVL
TGAKVRKEIYEAFDNIYPILKSFKKQ

Alignment:

CLUSTAL multiple sequence alignment by MUSCLE (3.8)

```
Human_TBP      MDQNNSLPPYAQGLASPGAMTPGIPFISPMMPYGTGLTPQPIQNTNSLSILEEQQRQQQ
Cricket_(Laupala) -----GTPLHQ-----EEDQQILP
Sawfly_(Athalia) MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
Sawfly_(Neodiprion) MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
Ant            MDQ---MLPSPGGFSIP---SIGTPLHQ-----EEDQQILP
Wasp_(Venturia) MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
Wasp_(Orussus)  MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
Cricket_(Gryllus) MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
Termite        MDQ---MLPSP-GFSIP---SIGTPLHQ-----EEDQQILP
                  *  *:. . *                      **:. .
```

```
Human_TBP      QQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQA--VAAAQVQSTSQQATQGTSGQA
Cricket_(Laupala) QAQQQQQQ-----NQQTPLLQ--MPLLATPQKIMHTYAPS-GFTT
Sawfly_(Athalia) NALQQQQQ-----QVQSSHQIQMHMPMGSNIHG--MPMLATPHKTMHTYTPA--FAT
Sawfly_(Neodiprion) NALQQQQQ-----VQQNQLQQMHMPMGSNVHG--MPMLGTPHKTMHTYAPA--FAT
Ant            NALQQQQQ-----SQYQLQLHSMSPNMQS--GMLMITTPQKTMHTYAPTPTFAT
Wasp_(Venturia)  NALQQQQQQQQQQQQQQTSQQYQMQQLHSMSPTMQANSMMMLSTPQKTMHTYAPTPTFAT
Wasp_(Orussus)   NALQQQQQ-----QQQQQSQYQLQLHSMSPNMQG--SMLMISTPQKTMHTYAPTPTFAT
Cricket_(Gryllus) MAQQQQQN-----QQTPLQ--MPLLGSPQKTMHTYAPA-GFTT
Termite        HAQQQQQQ-----QQQLPPPTSLQ--MPVLSTPQKTMHTYAPS-GF-T
                ****:                      :      :.  :  :  :
```

```
Human_TBP      PQ-LFHSQT---LTTAPLPGTTPPL-YFSPMTPMT--PITPATPASESSGIVPQLQNIYST
Cricket_(Laupala) PQSLMQPQTPQNMMSPMVASSSQIIPSSSLGPATPSPMTPMTPSADPAIVPQLQNIYST
Sawfly_(Athalia)  PQSLMQPQTPQNMMSPMVQPASQI-APPSIGPSTPGPMTMTPASADPGILPQLQNIYST
Sawfly_(Neodiprion) PQSLMQPQTPQNMMSPMVQPASQI-APPSIGPSTPGPMTMTPASADPGILPQLQNIYST
Ant            PQSLMQPQTPQNMMSPIVQSNQI-APPSIGPATPGPMTMTPASADPGILPQLQNIYST
Wasp_(Venturia)  PQSLMQPQTPQNMMSPMVQPTSQI-APSSIGPSTPGPMTMTPASADPGILPQLQNIYST
Wasp_(Orussus)   PQSLMQPQTPQNMMSPIVQSSQI-APSSIGPSTPGPMTMTPASADPGILPQLQNIYST
Cricket_(Gryllus) PQSLMQPQTPHNMMSPMVASASQIVAPSSSLGPATPGPMTMTPSADPGIVPQLQNIYST
Termite        PQSLMQPQTPQNLMSPMVTSSQMATLSNMGPATPSPMTPMTPSADPGIIPQLQNIYST
                ** *:. . **      :. . :  :. :  . :  *  *  *:. ** *:. *  . . . :. *****
```

```
Human_TBP      VNLGCKLDLKTIALRARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Cricket_(Laupala) VNLGCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Sawfly_(Athalia)  VNLNCKLELKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Sawfly_(Neodiprion) VNLNCKLNLKEIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Ant            VNLNCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Wasp_(Venturia)  VNLCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Wasp_(Orussus)   VNLNCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Cricket_(Gryllus) VNLCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
Termite        VNLGCKLDLKKIALHARNAEYNPKRFAAVIMRIREPRTTALIFSSGKMVCTGAKSEEDSR
                ***  ***: *  ***.*****:*****:*****:*****:*****:*****:*****:*****:*****
```

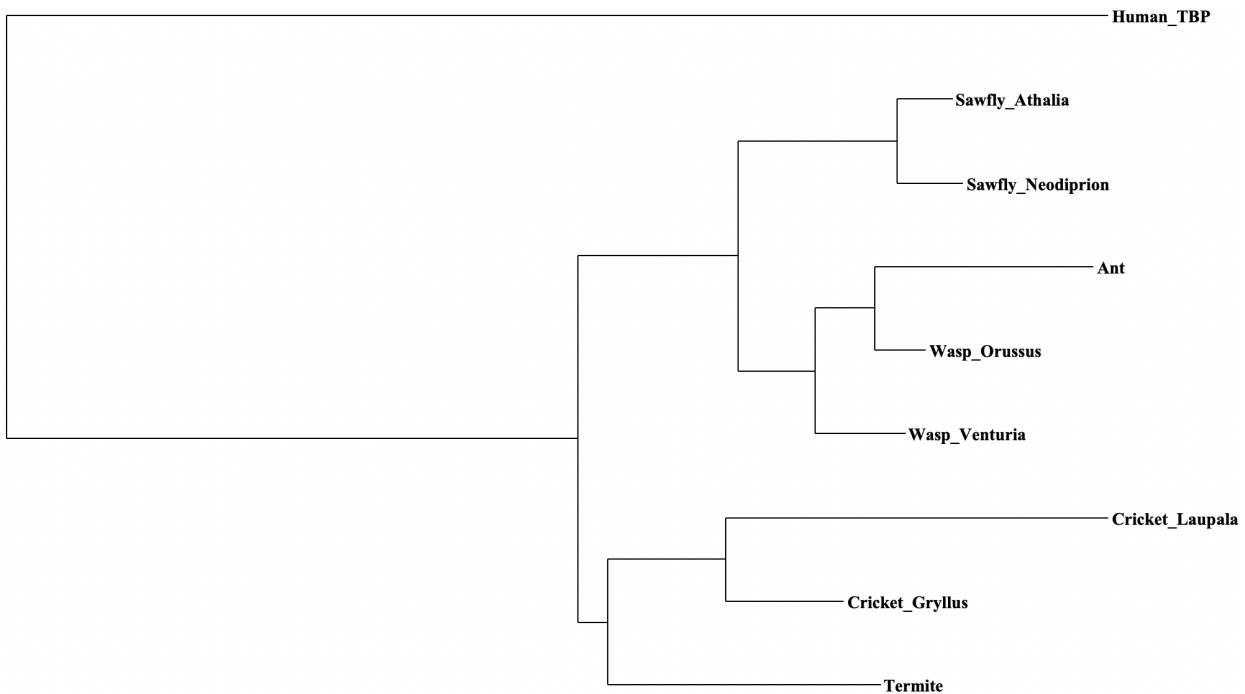
```
Human_TBP      LAARKYARVVQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Cricket_(Laupala) LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Sawfly_(Athalia)  LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Sawfly_(Neodiprion) LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Ant            LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Wasp_(Venturia)  LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Wasp_(Orussus)   LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Cricket_(Gryllus) LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
Termite        LAARKYARI IQKLGFPAKFLDFKIQNMVGSCDVKFPIRLEGLVLTHGQFSSYEPELFPGL
                *****:*****.*****:*****:*****:*****:*****:*****:*****:*****
```

```
Human_TBP      IYRMIKPRIVLLIFVSGKVLTGAKVRAEIYEAFENIYPILKGFRKT
Cricket_(Laupala) IYRMVKPRIVLLILYLES----GLTVQS-----
```


Sawfly_(Athalia)	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Sawfly_(Neodiprion)	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Ant	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Wasp_(Venturia)	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Wasp_(Orussus)	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Cricket_(Gryllus)	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
Termite	IYRMVKPRIVLLIFVSGKVVLTGAKVRQEIYEAFDNIYPILKSFKKQ
	****:*****: . * .*

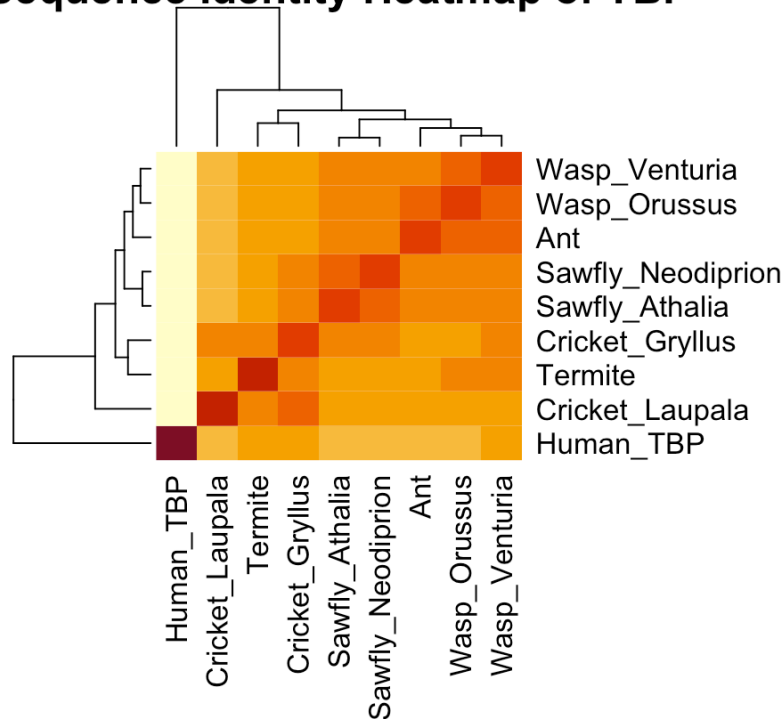
[Q6] Create a phylogenetic tree, using either a parsimony or distance-based approach. Bootstrapping and tree rooting are optional. Use “simple phylogeny” online from the EBI or any respected phylogeny program (such as MEGA, PAUP, or Phylip). Paste an image of your Cladogram or tree output in your report.

Aligned with MUSCLE and created a distance-based phylogenetic tree:



[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.

Sequence Identity Heatmap of TBP



[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 (experimental Technique), resolution (resolution), and source organism (source).

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

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

The consensus sequence has lots of gaps (mostly caused by human TBP), so I chose the sequence most similar to all the others in the alignment (Wasp_Venturia). Below are the top 3 unique hits from different species:

ID	Technique	Resolution	Source	E-value	Identity
5FUR_A	electron microscopy	8.5	Homo sapiens	2.14 e-125	86.34
1RM1_A	X-ray diffraction	2.5	Saccharomyces cerevisiae	9.24 e-111	80.77
7Z7N_D	electron microscopy	5.1	Thermochaetoides thermophila	1.78e-110	74.52

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

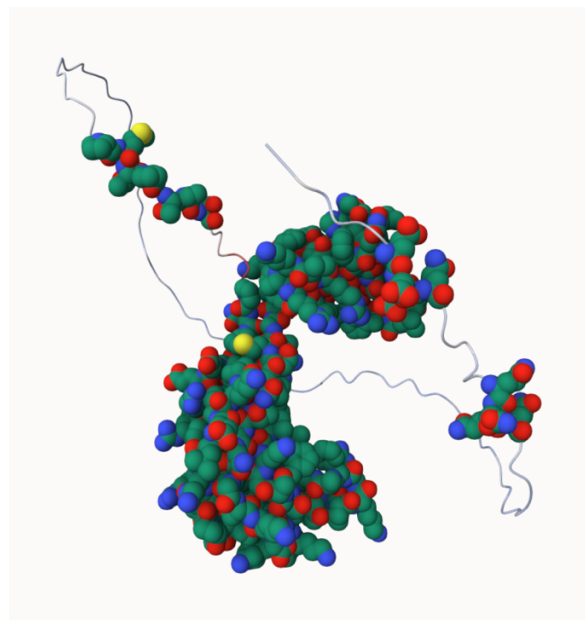
Note that this can take some time depending upon your sequence length. If your model is taking many hours to generate or your input sequence yields a “too many amino acids” (i.e. length) error you can focus on a single domain from your sequence - identify region by searching for [PFAM](#) domain matches.

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.).

Predicted structure of *Laupala kohalensis* TBP, colored by AlphaFold pLDDT quality score.

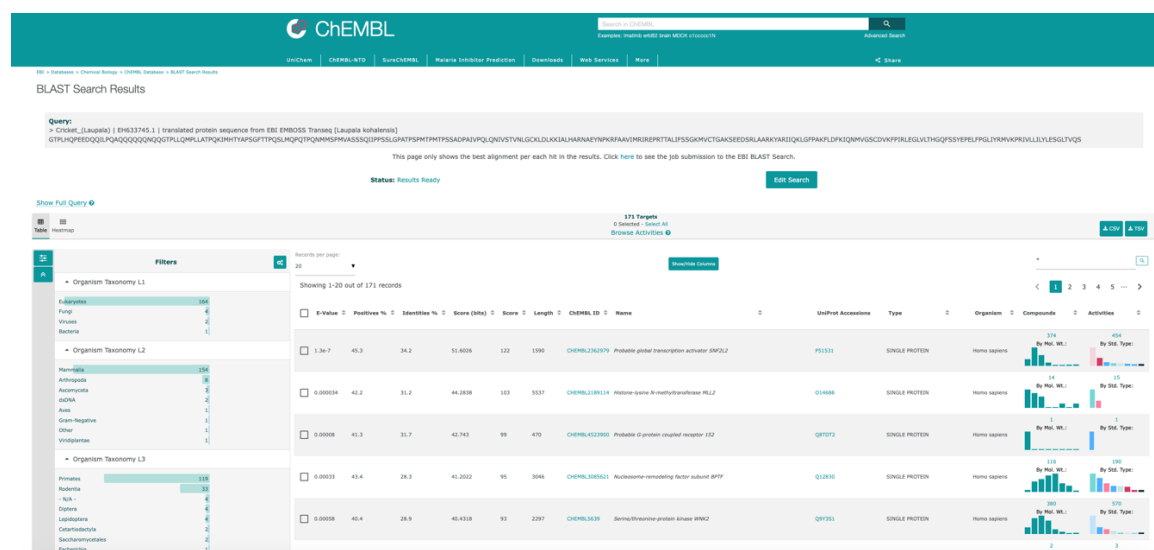


Conserved residues (based on my MSA) highlighted with spacefill. There are a few more conserved residues in the disordered region, but it was taking a long time to add them all as “components” in Mol* viewer (TBP is very conserved!).



[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]”.

My search (novel Laupala TBP sequence, default parameters, 12/6/23) generated 171 hits.



However, none of the hits are targeting TBP, even when I search the human sequence or just “TBP”.

The top hit is targeting SNF2L2, a chromatin remodeling complex subunit (ChEMBL2362979).

There are 80 associated assays, 79 binding and 1 functional.

Description of the top binding assay (ChEMBL4276258): Binding affinity to N-terminal MBP-fused human BRM RecA domain E852Q mutant expressed in Escherichia coli BL21 Star (DE3) by ITC method.

Ligand efficiency data:

Ligand Efficiencies

