

Find a Gene Project
BGGN 213
Kalodiah Toma
PID: A07606689
Email: k1toma@ucsd.edu

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: TOE1 target of EGR1
Accession: NP_079353.3 (NM_025077.4)
Species: Homo Sapiens

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

Method: TBLASTN (2.7.1) search against mouse
Database: ESTs Expressed Sequence Tags (est)
Organism: Mouse (Taxid: 10090)
Chosen match: Accession CB202140.1, a 860 base pair clone from Mus musculus.
See below for alignment details.

BLAST® = tblastn

Home Recent Results Saved Strategies Help

Check out the ClusteredNR database on BLAST+ Learn more Give us feedback

blastn blastp blastx **tblastn** tblastx

Translated BLAST: tblastn

TBLASTN search translated nucleotide databases using a protein query. more...

Reset page Bookmark

Enter Query Sequence

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

NP_079353.3

Query subrange

From To

Or, upload file Choose file no file selected

Job Title

Enter a descriptive title for your BLAST search

☐ Align two or more sequences

Choose Search Set

Database

Expressed sequence tags (est)

Organism

Optional

mouse (taxid:10090) 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

Enter an Entrez query to limit search

BLAST

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

☒ Show results in a new window

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

Algorithm parameters

Alignment view Pairwise [Restore defaults](#) Download

100 sequences selected

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

AGENCOURT_11289261 NIH_MGC_135 Mus musculus cDNA clone IMAGE:30140611 5', mRNA sequence
 Sequence ID: [CB202140.1](#) Length: 860 Number of Matches: 1

Range 1: 22 to 822 [GenBank](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps	Frame
451 bits(1159)	2e-156	Compositional matrix adjust.	227/267(85%)	246/267(92%)	1/267(0%)	+1
Query 1	MAADSDDGAVSAPAASDGGVSKSTTSGEELVvqvpvvdvqSNNFKEMWPSLLLAIKTANF					
Sbjct 22	MAADSDDG P SDGGV+K+T S EE VV+VPVVDVQS+NFKE+WPSLLLA+KTA+F					
Query 61	VAVDTELSGLGDRKSLNQIEERYKAVCHAARTRSILSLGLACFRQPKGEHSYLAQV					
Sbjct 202	VAVDTELSGLGDRKSLNQIEERYKAVCHAARTRS+LSLGLACF++QPKGE+SYLAQV					
Query 121	FNLTLCCMEYVIEPKSVQFLIQHGFNFNQYAAQGIPIYHKGNKGDQSQSVRTLFLLEL					
Sbjct 382	FNLTLCCMEYVIEPKSVQFL+QHGFFNF+QYAAQGIPIYHKGNKGDQSQSVRTLFLLEL					
Query 181	IRARRPLVLHNLIDLFLYQNFYAHLPESLGTFTADLCMFPAgiYDTKYAAEFHARFV					
Sbjct 562	IRARRPLVLHNLIDLFLYQNFYAHLP+LGTFTADLCMFPAgiYDTKYAAEFHARFV					
Query 241	ASYLEYAFRCER-ENGKQRAAGSPHL 266					
Sbjct 742	ASYLEY F + G +R + + H+ 822					



[Q3] Gather information about this “novel” **protein**. At a minimum, show me the protein sequence of the “novel” protein as displayed in your BLAST results from [Q2] as FASTA format (you can copy and paste the aligned sequence subject lines from your BLAST

result page if necessary) or translate your novel DNA sequence using a tool called EMBOSS Transeq at the EBI. Don't forget to translate all six reading frames; the ORF (open reading frame) is likely to be the longest sequence without a stop codon. It may not start with a methionine if you don't have the complete coding region. Make sure the sequence you provide includes a header/subject line and is in traditional FASTA format.

Chosen sequence:

```
>A. TOE1 like protein (sequence taken from BLAST result){1}SEP
MAADSDDGVPVPVTPSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTSVLSLGLACFRQQPDKGNSYLAQVFNLTLTCMEYVIEPKSVQFLVQHGFNFRQYAO
GIPYHKGNKGDGDESQSQSVRTLFLLELIRARRPLVLHNLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
AAEFHARFVASYLEYFPFEMSTGRMGSKRGLATHI
Species: Mus musculus
Eukaryota; Animalia; Chordata; Mammalia; Rodentia; Muridae; Mus; Mus;
M. musculus.
```

[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.

A BLASTP search against NR database yielded a top hit result is to a protein from *Rattus norvegicus* (Norway rat).

See additional screen shots below for top hits and selected alignment details:

blastn

blastp

blastx

tblastn

tblastx

BLASTP programs search protein databases using a protein query. [more...](#)

Enter Query Sequence

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

>Protein|

MAADSDDGVPPVPTPSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSSLALKTASFV

AVDTELSGLGDRKSLLNQCIEERYKAVCHAARTSVLSLGLACFRQQPKGENSYLAQVFNL

TLLCMEEYVIEPKSVQFLVQHGFNFRQYAQGIPYHKGNDKGDESQSQSVRTLFLLELRARR

PLVLHNLGLDLVFLYQNFYAHLPENLGTFTADLCEMFPAgiYDTKYAAEFHARFVASYLEYPPF

EMTCGMGCKDGLATL

Query subrange [?](#)

From

To

Or, upload file

Choose File no file selected [?](#)

Job Title

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

☐ Align two or more sequences [?](#)

Choose Search Set

Databases

☒ Standard databases (nr etc.): **New**
☐ Experimental databases

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

Compare

☐ Select to compare standard and experimental database [?](#)

Standard

Database

Non-redundant protein sequences (nr) [?](#)

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

Program Selection

Algorithm

☐ Quick BLASTP (Accelerated protein-protein BLAST)
☒ blastp (protein-protein BLAST)
☐ PSI-BLAST (Position-Specific Iterated BLAST)
☐ PHI-BLAST (Pattern Hit Initiated BLAST)
☐ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm [?](#)

BLAST

Search database nr using Blastp (protein-protein BLAST)

☐ Show results in a new window

Details: The top result is to a protein from *Rattus norvegicus* (Norway rat), see second screen shot below for alignment details:

[Download](#)
[GenPept](#)
[Graphics](#)

[Next](#)
[Previous](#)
[Descriptions](#)

target of EGR1 protein 1 isoform X4 [Rattus norvegicus]

Sequence ID: [XP_038965533.1](#) Length: 257 Number of Matches: 1

Range 1: 1 to 248 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
504 bits(1299)	1e-179	Compositional matrix adjust.	241/248(97%)	245/248(98%)	0/248(0%)
Query 1	MAADSDDGVPVPPTPSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSSLALKKTASF				
Sbjct 1	MAADSDDGVP VPTPSDGGVN+ TQSAEEFVVRVPVVDVQSDNFKEIWPSSLALKKTASF				
Query 61	VAVDTELSGLDGRKSLLNQCIEERYKAVCHAAARTSRVLSGLACFRQDPKGENSYLAQV				
Sbjct 61	VAVDTELSGLDGRKSLLNQCIEERYKAVCHAAARTSRVLSGLACF+QDPKG++SYLAQV				
Query 121	FNLTLTLMCEEYIEPKSVQFLVQHGFNFRQYAQGIPYHKGNKDGDSESQSVRTLFLFL				
Sbjct 121	FNLTLTLMCEEYIEPKSVQFLVQHGFNFRQYAQGIPYHKGNKDGDSESQSVRTLFLFL				
Query 181	IRARRPLVLHNLIDLVLFYQNIFYAHLPENLGTFTADLCMFAGIYDTKYAAEFHARFV				
Sbjct 181	IRARRPLVLHNLIDLVLFYQNIFYAHLPENLGTFTADLCMFAGIYDTKYAAEFHARFV				
Query 241	ASYLEYPF 248				
Sbjct 241	ASYLEYAF 248				

Related Information

[Gene](#) - associated gene details
[Genome Data Viewer](#) - aligned genomic context

[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.

```
>Homo sapiens NP_079353.3 target of EGR1 protein 1 [Homo sapiens]
MAADSDDGAVSAPAASDGGVSKSTTSGEELVVQVPVVDVQSNFKEIWPSSLLLAIKTANFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTRSILSLGLACFKRQPPDKGEHSYLAQVFNLTLTCMEEYVIEPKSVQFLVQHGFNFNRQY
GIPYHKGNKDGDESQSQSVRTLFLLELIRARRPLVLHNGLIDLVLFLYQNFYAHLPESLGTFTADLCMFPAgiYDTKY
AAEFHARFVASYLEYAFRKCERENGKQRAAGSPHLTLEFCNYPSSMRDHIDYRCCLPPATHRPHPTSICDNFSAYGW
CPLGPPQCPQSHDIDLIIIDTDEAAAEDKRRRRRRREKRKRALLNLPQTSGEAKDGPCKQVCGDSIKPEETEQEVA
ADETRNLPHSKQGNKNLDEMGIIKAARPEIADRATSEVPGSQASPFPVPGDGLHRAGFADFMTGYVMAYVEVSQGPQP
CSSGPWLPECHNKVYLSGKAVPLTVAKSQFSRSSKAHNQKMKLTWGSS
```

```
>Mus musculus TOE1 like protein (sequence taken from BLAST result)[1]SEP
MAADSDDGVPVPVPTPSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSSLLLALKTASFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTRSVLSLGLACFRQPPDKGENSYLAQVFNLTLTCMEEYVIEPKSVQFLVQHGFNFNRQY
GIPYHKGNKDGDESQSQSVRTLFLLELIRARRPLVLHNGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPAgiYDTKY
AAEFHARFVASYLEYFPPEMSTGRMGSKRGLATHI
```

```
>Rattus norvegicus XP_038965533.1 target of EGR1 protein 1 isoform X4 [Rattus
norvegicus]
MAADSDDGVPSPVPTPSDGGVNRTTQSAEEFVVRVPVVDVQSDNFKEIWPSSLLLALKTASFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTRSVLSLGLACFKQPPDKGDSSYLAQVFNLTLTCMEEYVIEPKSVQFLVQHGFNFNRQY
GIPYHKGNKDGDESQSQSVRTLFLLELIRARRPLVLHNGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPAgiYDTKY
AAEFHARFVASYLEYAFRKWPMAGAL
```

```
>Cricetulus griseus XP_027258385.1 target of EGR1 protein 1 [Cricetulus
griseus]
MAADSDDGVPSPVPTTSDGGVNKNTKSAEEFVIRVPVVDVQSDNFKEIWPSSLLLALKTASFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTRSVLSLGLACFKQPPDKGENSYLTQVFNLTLTCMEEYVIEPKSVQFLVQHGFNFNRQY
GIPYHKGNKDGDESQSQSVRTLFLLELIRARRPLVLHNGLIDLVLFLYQNFYAHLPETLGTFTADLCMFPAgiYDTKY
AAEFHARFVASYLEYAFRKCERENGKQRAAGTPHLALEFCSYPSSMRGHIDYRCMSPTVYRRSHTTGICDKFSAYG
WCPLGPPQCPQSHDIDLIIIDTDEAVAEDKRRRRRRRKDKRKRALLSQPETQTFEEAEDGPPTKQVCEDSLKTEEMEQRV
TEGETRDELGSKQAHKSGLEMEHKATSSETVDVATTELPVSQASPFPVPGDGLHRAGFADFMTGYVMAYVGLSQGPQ
LCSSRPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAWGSS
```

```
>Chionomys nivalis XP_057640769.1 target of EGR1 protein 1 [Chionomys
nivalis]
MAADSDDGVPSPVPTSDGDVNKNTKSAEEFVVRVPVVDVQSDNFKEIWPSSLLLALKTASFVAVDTELSGLGDRKSL
NQCI EERYKAVCHAARTRSVLSLGLACFKQPPDKGENSYLTQVFNLTLTCMEEYVIEPKSVQFLVQHGFNFNRQY
GIPYHKGNKDGDESQSQSVRTLFLLELIRARRPLVLHNGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPAgiYDTKY
```

AAEFHARFVASYLEYAFRK CERENGKRQATGSPHLALEFCSYPSSMRGHIDYRCCMSPVTSRRSHTAGICTKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQPGTGT FEEAEDGPPTKQVCEDSLKTETE QKVT
EGESRDQAEDGLPTKQVCEGSLKTEIEQKVTEGESRDQLGSKQCHNSDLEVEHKATSSEIADVAASELPASQASPNP
VPGDGLHRAGFDAFMTGYVMAYVGLRHGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAWG
SS

>Phodopus roborovskii XP_051044502.1 target of EGR1 protein 1 [Phodopus roborovskii]

MAADSDDGVPSPVPTTSDGGVNKNTKSAAEFVIRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTSVLSLGLACFKHQADKGENSYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKRRAAGTPHLVLEFCSYPSSMRGHIDYRCC LSPVTYRRSHTTSICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQ LGTQNFEEADDPPTKQVCEDNLKTQELEQ RV
TERETRDELDSKQGHKSDLEMEHKATGSETADVAISELPVSQASPNMPG DGLHRAGFDAFMTGYVMAYVGLS QGPQ
LCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAWSSS

>Peromyscus californicus insignis XP_052572255.1 target of EGR1 protein 1 [Peromyscus californicus insignis]

MAADSV DGVSPVPTTSDGGVNKNTKSAD E FIVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTSVLSLGLACFKQ QPDKGENSYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCC MSPVTYRRFHTSGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQ PGAQT FEEAEDGPPTKQVCEDSLKTEEIEQKV
TEGETRDQLGSQQGHKSGLAVERKATSSETAEVATSEL PVSQANPNPGDGLHRAGFDAFMTGYVMAYVGLS QGPQ
LCSSGPWLPECHNKVYLSGKTVPLTVARSQFSRSSKAHNQKMKLAWGSS

>Mesocricetus auratus XP_005072112.3 target of EGR1 protein 1 [Mesocricetus auratus]

MAADSDDGVPSPVPTTSDGGVNKNTKSAEKFVIRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTSVLSLGLACFKQ QPDKGENCYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHPETLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGTPHLALEFCSYPSSMRSHIDYRCCISPVTYRRSHTTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQ PGTQT FEEAEDGPPTKQVCEDNLKTEEIEQ RV
TEGEIRDELGSKQAHKSGSEMEHKATSSETVDVATSEL PVSQASPNVP G DGLHRAGFDAFMTGYVMAYVGLS QGPQ
LCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAWGSS

>Myodes glareolus XP_048286556.1 target of EGR1 protein 1 [Myodes glareolus]

MAADSDDGVPSPV PATSDGDVNKTTKSAAEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTSVLSLGLACFKQ QPDKGENSYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKRQAAGSPHLALEFCSYPSSMRGHIDYRCC MSPITSRRSHTTGICTKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLNQ PGTET FEEAEDGPPTKQVCEDSLKTEIEQKV
EGETRDQLGSNQDHKSDLEVEHKATSSEIADVAASELPVSQASPNVP G DGLHRAGFDAFMTGYVMAYVGLRHGPQL
CSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAWGSS

>Mastomys coucha XP_031234176.1 target of EGR1 protein 1 isoform X1 [Mastomys coucha]

MAADSA DGMPSVPSSSDGGVNKNTQSAKEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTSVLSLGLACFKQ QPEKGENSYLAQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFEARFVASYLEYAFRK CERENG RQQAVGSPHLALEFCSYPSSMRGHIDYRCC MSPVTCRRSYTTTGICDKFSAYG
WCPLGPECPQSHDIDLII DTDEAVAEDKRRRRWQRRKDRKRALQSQ PGTQTL EEAEGGPPTKQVCEDSLKAEKMEQ
KVAEGDAGDQLGSRQGTGSLEIAHRRTSAETADVAPSEL PVSQASTNPLPGDGLHRAGFDAFMTGYVMAYVGLS QG
LQLCSSEPWL PKCHNKVYLSGKTVPLTVAKSQFSHP SKAHKQKMKLAWGSS

>Microtus ochrogaster XP_026633303.1 target of EGR1 protein 1 isoform X1 [Microtus ochrogaster]

MAADSDDGVPVSATS DGDVNKNTKPAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFKQQPDKGENSYLTQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCMSPVTSRRSHTTGICTRF SAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQPGTETFE EAEDGPPIKQVCEDSLKTEIEQKVT
EGESRDQLGSKQDQNSDLEVEHKATSSEIADVAASELPVSQASPNPVPDGLHRAGFDAFMTGYVMAYVGLRHGPQL
CSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAWGSS

>Acomys russatus XP_051027553.1 target of EGR1 protein 1 [Acomys russatus]
MAADSGDGVPLVPKTS DGGVNKSTSAAEFVARVPVVDVQTDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFRQQPDKGENSYLAQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKYA
AEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCMSPVTDRRSYATGICDKFSAYGW
CPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLRQPGTQTFEEGEDGPPTKQVREDSLNTENTEQKVA
EGETSDQLGSKQGHKGGLEMKHEATGSETADVATSELQVNQASPNPVPDGLHRAGFDAFMTGYVMAYVGLSQGLQL
CSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGNS

>Psammomys obesus XP_055474418.1 target of EGR1 protein 1 [Psammomys obesus]
MAADSDDGGLSVPPTS DGVVNKNTKSEEEFVIRVPVVDVQTDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFKQQPDKGENSYLAQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCMSPVTYHRSHTGGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALLSQPGMQTFEEAE EGPPTKQVCEDSLKTENTEQKV
AEGETRDQPGSKQGHKGGLEMEHEAVSSEIADVATSELVQNQTSNPVPDGLHRAGFDAFMTGYVMAYVGLSQGTQ
LCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Neotoma lepida OBS83234.1 hypothetical protein A6R68_22771 [Neotoma lepida]
MAADSDDGVPSPVPTS DDDGVNKNKTSAD E FVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFKQQPDKGENSYLTQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQNSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
AAEFHARFFCSYPSSMRDHIDYRCCMSPVTYRRSHTTGICNKFSAYGWCPLGPQCPQSHDIDLII DTDEAVAEDKRR
RRRRKDRKRALLSQPETQTFEEAEDGPPTKQVCEDSLKTEEIEQKVTEGETRDQLGSQQGHKSLEIEYKATSSKI
ADVATSELVPSQASPNMPDGLHRAGFDAFMTGYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQ
FSRSSKAHNQKMKLAWGSS

>Mus caroli XP_021015031.1 target of EGR1 protein 1 [Mus caroli]
MAADSDDGVPVPVPTS DSGGVNKNNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKNLL
NQCI EERYKAVCHAARTSVLSLGLACFRQQPDKGENSYLAQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRMLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
ASEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCTSPGTCRRSRPTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRSLQSQPGTQALAEAEDGPPTKQVCEDSLKTEKMEQKV
AEGEAGDQPGSREGHTSSLEMAHRRTSAETADVATSELLVNQASTNPVPDGLHRAGFDAFMTGYVMAYVGLSQGLQ
LCSSEPWLPECHNKVYLSGKTVPLTVTKSQFSRPSKAHNQKMKLAWGSS

>Mus pahari XP_021055456.1 target of EGR1 protein 1 [Mus pahari]
MAADSDDGVPVPVATSS DGGVNKNNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFRQQPDKGENSYLAQVFNLTLTLCV E EYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCMSPGTCRRSHTTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDRKRALQSQPGTQTLAEAEDGPPTKQVCEDSLKTEKMEQKV
AEGEAGDQPGSRQDHTGSLEMEHRRTRAETA EVATSELLVSQARTDPVPDGLHRAGFDAFMTGYVMAYVGLSQGLQ
LCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Rattus rattus XP_032755624.1 target of EGR1 protein 1 [Rattus rattus]
MAADSDDGVPSPVPTS DSGGVNRTTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLN
NQCI EERYKAVCHAARTSVLSLGLACFKQQPDKGDSSYLAQVFNLTLTLCMEEYVIEPKSVQFLVQHGFNFNRQYAQ
GIPYHKGNDKGDESQSQSVRTLFLLELIRARRPLVLHNLGLIDLVLFLYQNFYAHLPENLGTFTADLCMFPA GIYDTKY

AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHVALEFCSYPSSMRGHIDYRCCMSPVSCRRSHTTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDKRKRALQSQPGTQNLEEAEDGPPTKQVCEDSLKTEKIEQKV
AEGDQLGSTQGHKDSLEMACKRTADVPTSELLVNQASPNVPVPGDGLHRAGFDAFMTGYVMAYVGLSKGLQLCSSEPW
LPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Apodemus sylvaticus XP_052033054.1 target of EGR1 protein 1 [Apodemus
sylvaticus]
MAADSDDGVPSPVPTSSDGGVNKNTQTAEFFVVRVPVVDVQNDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTRSVLSLGLACFKQQPDKGENSYLAQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHVALEFCSYPSSMRGHIDYRCCMPPVTCRPPHTTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDKRKRALQSQPGTQ TLEEAEDGPPTKQVCEDNVKTEKVEQKV
AEGEAGDELGSRQGHTGSP EAHRTSADTADVATSEL PVNQANANPVPVPGDGLHRAGFDAFMTGYVMAYVGLSQGLQL
CSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Meriones unguiculatus XP_021489081.1 target of EGR1 protein 1 [Meriones
unguiculatus]
MAADSDVGGLSVPPTSDG VVNKNTKSEEEFVIRVPVVDVQTDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTRSVLSLGLACFKQQPDKGENSYLVQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPHVALEFCSYPSSMRGHIDYRCCMSPVTYHRSHTSGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDKRRRALLSQPGMQTFEEAE EGPPTKQVCEDSLKTENPEQKV
AEGETRDQVGSKQGHEG GLEMEHEAPSSEIADVATSEL PVNQASPNVPVPGDGLHRAGFDAFMTGYVMAYVGLSQGTQ
LCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Arvicanthis niloticus XP_034357978.1 target of EGR1 protein 1 [Arvicanthis
niloticus]
MAADSDDGVPSPVPTSSDGGVNKNTQSAAEFFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTRSVLSLGLACFKQQPDKGENSYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQNQSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYVFRK CERENGKQRAAGSPHVALEFCSYPSSMRGHIDYRCCMSPVTCRRSHTTGICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDKRKRALQSQPGTQ TLEEAEDGPPTKQVCEDSLQTEKIEQIM
AEGEAKDQLGSKQGHTGSLVMAHKRTSSETADMATSDLLVNQGSTNPVPGDGLHRAGFDAFMTGYVMAYVGLSQGLQ
LCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAWGSS

>Onychomys torridus XP_036033233.1 target of EGR1 protein 1 [Onychomys
torridus]
MAADSDDGVPSPVPTTSDGGVNKNTKSADEFVVRVPVVDVQSDNFKEIWPSLLLALKTASFVAVDTELSGLGDRKSLL
NQCI EERYKAVCHAARTRSVLSLGLACFKQQPAKGENSYLTQVFNLTL CME EYVIEPKSVQFLVQHGFNFNRQY AQ
GIPYHKGN DKGDESQS QSVRTL FLELIRARRPLVLHNGLIDL VFLYQNFY AHLPENLGTFTADLC EMFPAGIYDTKY
AAEFHARFVASYLEYAFRK CERENGKQRAAGSPH LTFECSYPSSMGGHIDYRCCMSPVTHRRSHTSSICDKFSAYG
WCPLGPQCPQSHDIDLII DTDEAVAEDKRRRRRRKDKRKRALLSQPGAQT FEEAE EGPPTKQICEDSLKTEETE QKV
TEGETKDQLG SQQDHKSGLAIRRKATSSETADVATSEL PVSQANPNVPVPGDGLHRAGFDAFMTGYVMAYVGLSQGPQ
LCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAWGSS

Alignment:

Homo-sapiens	MAADSDDGAVSAPAASDGGVSKSTTSGEELVVQVPVVDVQSNNFKEMWPSLLLAIKATNF
Psammomys-obesus	MAADSDDGGLSVPPTSDG VVNKNTKSEEEFVIRVPVVDVQTDNFKEIWPSLLLALKTASF
Meriones-unguiculatus	MAADSDVGGLSVPPTSDG VVNKNTKSEEEFVIRVPVVDVQTDNFKEIWPSLLLALKTASF
Acomys-russatus	MAADSGDGVLPVKTS DGGVNKST-SAEFFVARVPVVDVQTDNFKEIWPSLLLALKTASF
Mastomys-coucha	MAADSADGMPSPVSSSDGGVNKNTQSAKEFVVRVPVVDVQSDNFKEIWPSLLLALKTASF
Phodopus-roborovskii	MAADSDDGVPSPVPTTSDGGVNKNTKSAEEFVIRVPVVDVQSDNFKEIWPSLLLALKTASF
Cricetulus-griseus	MAADSDDGVPSPVPTTSDGGVNKNTKSAEEFVIRVPVVDVQSDNFKEIWPSLLLALKTASF
Mesocricetus-auratus	MAADSDDGVPSPVPTTSDGGVNKNTKSAKEFVIRVPVVDVQSDNFKEIWPSLLLALKTASF
Neotoma-lepida	MAADSDDGVPSPVPTTSDGGVNKNTKSADEFVVRVPVVDVQSDNFKEIWPSLLLALKTASF
Onychomys-torridus	MAADSDDGVPSPVPTTSDGGVNKNTKSADEFVVRVPVVDVQSDNFKEIWPSLLLALKTASF
Peromyscus-californicus-insignis	MAADSDDGVPSPVPTTSDGGVNKNTKSADEFVIRVPVVDVQSDNFKEIWPSLLLALKTASF
Mus-caroli	MAADSDDGVPSPVPTTSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASF
Mus-pahari	MAADSDDGVPVATSSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLLALKTASF

MAADSDDGVPSPVPTSSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPVATSDGGVNKTKSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPATSDGGVNKTKSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPVATSDGGVNKNTKPAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPVPTSSDGGVNKNTQTAEEFVVRVPVVDVQNDNFKEIWPSLLALKTASF
MAADSDDGVPSPVPTSDGGVNKNTQSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPVPTSDGGVNRTTQSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF
MAADSDDGVPSPVPTSDGGVNRTTQSAEEFVVRVPVVDVQSDNFKEIWPSLLALKTASF

[illegible][illegible][illegible]

Rattus-rattus	IRARRPLVLHNGLIDLVLFLYQNFYAHLPENLGTFTADLCEMFPAGIYDTKYAAEFHARFV *****.*****;** **
Homo-sapiens	ASYLEYAFRK CERENGKQRAAGSPHLTLEFCNYPSSMRDHIDYRCCLPPAT-HRPHTSI
Psammomys-obesus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTYHRSHTGGI
Meriones-unguiculatus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTYHRSHTSGI
Acomys-russatus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTDRRSYATGI
Mastomys-coucha	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTCRRSYTTGI
Phodopus-roborovskii	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCLSPVTYHRSHTTGI
Cricetulus-griseus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTYHRSHTTGI
Mesocricetus-auratus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRSHIDYRCISPVTYHRSHTTGI
Neotoma-lepida	-----FCSYPSSMRDHIDYRCMSPVTYHRSHTTGI
Onychomys-torridus	ASYLEYAFRK CERENGKQRAAGSPHLTLEFCNYPSSMRGHIDYRCMSPVTHRRSHTSSI
Peromyscus-californicus-insignis	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTYHRSHTSGI
Mus-caroli	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCCTSPGTTCRRSPTGI
Mus-pahari	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPGTTCRRSHTTGI
Arvicanthis-niloticus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTCRRSHTTGI
Myodes-glareolus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPITSRRSHTTGI
Chionomys-nivalis	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTSRRSHTAGI
Microtus-ochrogaster	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTSRRSHTTGI
Apodemus-sylvaticus	ASYLEYAFRK CERENGKQRAAGSPHLALEFCSYPSSMRGHIDYRCMSPVTCRPHTTGI
Mus-musculus	ASLEYPFPEM-----
Rattus-norvegicus	ASYLEYAFRK-----
Rattus-rattus	ASYLEYAFRK CERENGKQRAAGSPHVALEFCSYPSSMRGHIDYRCMSPVSCRRSHTTGI
Homo-sapiens	CDNFSAYGWCPLGPQCPQSHDIDLIDTDEAAAEDKRRRR--RRREKRKRALLNLPGTQT
Psammomys-obesus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPMQT
Meriones-unguiculatus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPMQT
Acomys-russatus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Mastomys-coucha	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Phodopus-roborovskii	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Cricetulus-griseus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Mesocricetus-auratus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Neotoma-lepida	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Onychomys-torridus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Peromyscus-californicus-insignis	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Mus-caroli	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Mus-pahari	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Arvicanthis-niloticus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Myodes-glareolus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Chionomys-nivalis	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Microtus-ochrogaster	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Apodemus-sylvaticus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Mus-musculus	-----
Rattus-norvegicus	-----
Rattus-rattus	CDKFSAYGWCPLGPQCPQSHDIDLIDTDEAVAEDKRRRR--RRDKRRRALLSQPGTQT
Homo-sapiens	SGEAKDGPPTKQVCGDSIKPE-----ETEQVAAD
Psammomys-obesus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Meriones-unguiculatus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Acomys-russatus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Mastomys-coucha	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Phodopus-roborovskii	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Cricetulus-griseus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Mesocricetus-auratus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Neotoma-lepida	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Onychomys-torridus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Peromyscus-californicus-insignis	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Mus-caroli	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Mus-pahari	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Arvicanthis-niloticus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Myodes-glareolus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Chionomys-nivalis	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Microtus-ochrogaster	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Apodemus-sylvaticus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Mus-musculus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Rattus-norvegicus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Rattus-rattus	FEAEEDGPPTKQVCEDSLKTE-----NTEQKVAEG
Homo-sapiens	ETRNLPHSKQGNKNDLEMGIAARPEIADRATSEVPGSQASPFPVPGDGLHRAGFADFMT
Psammomys-obesus	ETRDQVGSKQGHKGGLMEHEAVSSEIADVATSELPVNQTSNPNVPGDGLHRAGFADFMT
Meriones-unguiculatus	ETRDQVGSKQGHKGGLMEHEAVSSEIADVATSELPVNQTSNPNVPGDGLHRAGFADFMT
Acomys-russatus	ETSDQLGSKQGHKGGLMEHEAVSSEIADVATSELPVNQTSNPNVPGDGLHRAGFADFMT

Mastomys-coucha	DAGDQLGSRQGGHTGSLEIAHRRRTSAETADVAPSELVPSQASTNPLPGDGLHRAGFDAMT
Phodopus-roborovskii	ETRDELDSKQGHKSLEMEHKATGSETADVAISELPVSQASPNPMPGDGLHRAGFDAMT
Cricetulus-griseus	ETRDELGSKQAHKSGLEMEHKATSSSETVDVATELPVSQASPNPVPDGLHRAGFDAMT
Mesocricetus-auratus	EIRDELGSKQAHKSGSEMEHKATSSSETVDVATSELPVSQASPNPVPDGLHRAGFDAMT
Neotoma-lepida	ETRDQLGSGQGHKSLEIEYKATSSKIADVATSELPVSQASPNPMPGDGLHRAGFDAMT
Onychomys-torridus	ETKDQLGSGQDDHKSGLAIRRKATSSSETADVATSELPVSQANPNPVPDGLHRAGFDAMT
Peromyscus-californicus-insignis	ETRDQLGSGQGHKSGLAVERKATSSSETAEVATSELPVSQANPNPVPDGLHRAGFDAMT
Mus-caroli	EAGDQPGSREGHTSSLEMAHRRRTSAETADVATSELLVNQASTNPVPGDGLHRAGFDAMT
Mus-pahari	EAGDQPGSRQDHTGSLEMEHRRTRAETAEVATSELLVQARTDPVPGDGLHRAGFDAMT
Arvicanthis-niloticus	EAKDQLGSKQGGHTGSLVMAHKRTSSETADMATSDLLVNQGSTNPVPGDGLHRAGFDAMT
Myodes-glareolus	ETRDQLGSGNDHKSLEVEHKATSSSEIADVAASELPVSQASPNPVPDGLHRAGFDAMT
Chionomys-nivalis	ESRDQLGSKQCHNSDLEVEHKATSSSEIADVAASELPASQASPNPVPDGLHRAGFDAMT
Microtus-ochrogaster	ESRDQLGSKQDQNSDLEVEHKATSSSEIADVAASELPVSQASPNPVPDGLHRAGFDAMT
Apodemus-sylvaticus	EAGDELGSRQGGHTGSPMAHR-TSADTADVATSELPVNQANANPVPDGLHRAGFDAMT
Mus-musculus	-STGRMGSKRG-----
Rattus-norvegicus	-----
Rattus-rattus	---DQLGSTQGHKDSLEMACKRT---ADVPTSELLVNQASPNPVPDGLHRAGFDAMT

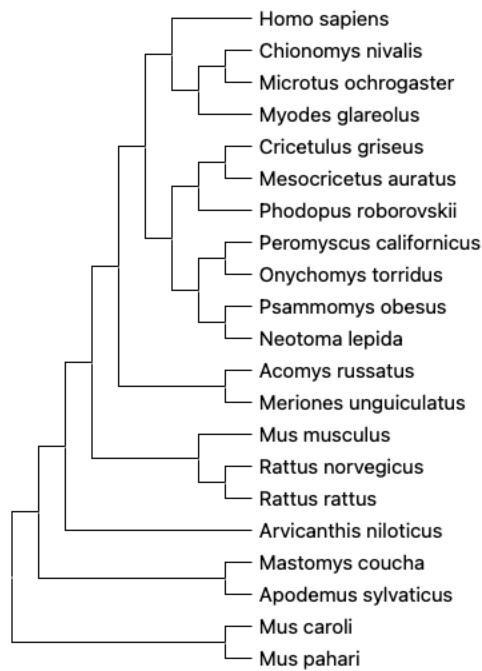
Homo-sapiens	GYVMAYVEVSQGPQPCSSGPWLPECHNKVYLSGKAVPLTVAKSQFSRSSKAHNQKMKLTW
Psammomys-obesus	GYVMAYVGLSQGTQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Meriones-unguiculatus	GYVMAYVGLSQGTQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Acomys-russatus	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Mastomys-coucha	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSHPSKAHKQKMKLAW
Phodopus-roborovskii	GYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAW
Cricetulus-griseus	GYVMAYVGLSQGPQLCSSRPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAW
Mesocricetus-auratus	GYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSHSSKAHNQKMKLAW
Neotoma-lepida	GYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAW
Onychomys-torridus	GYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAW
Peromyscus-californicus-insignis	GYVMAYVGLSQGPQLCSSGPWLPECHNKVYLSGKTVPLTVARSQFSRSSKAHNQKMKLAW
Mus-caroli	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVTKSQFSRPSKAHNQKMKLAW
Mus-pahari	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Arvicanthis-niloticus	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Myodes-glareolus	GYVMAYVGLRHGFPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAW
Chionomys-nivalis	GYVMAYVGLRHGFPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAW
Microtus-ochrogaster	GYVMAYVGLRHGFPQLCSSGPWLPECHNKVYLSGKTVPLTVAKSQFSRSSKAHNQKMKLAW
Apodemus-sylvaticus	GYVMAYVGLSQGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW
Mus-musculus	-----SLATHI-----
Rattus-norvegicus	-----W-----PMA-----
Rattus-rattus	GYVMAYVGLSKGLQLCSSEPWLPECHNKVYLSGKTVPLTVAKSQFSRPSKAHNQKMKLAW

. . .

Homo-sapiens	GSS
Psammomys-obesus	GSS
Meriones-unguiculatus	GSS
Acomys-russatus	GNS
Mastomys-coucha	GSS
Phodopus-roborovskii	SSS
Cricetulus-griseus	GSS
Mesocricetus-auratus	GSS
Neotoma-lepida	GSS
Onychomys-torridus	GSS
Peromyscus-californicus-insignis	GSS
Mus-caroli	GSS
Mus-pahari	GSS
Arvicanthis-niloticus	GSS
Myodes-glareolus	GSS
Chionomys-nivalis	GSS
Microtus-ochrogaster	GSS
Apodemus-sylvaticus	GSS
Mus-musculus	---
Rattus-norvegicus	GAL
Rattus-rattus	GSS

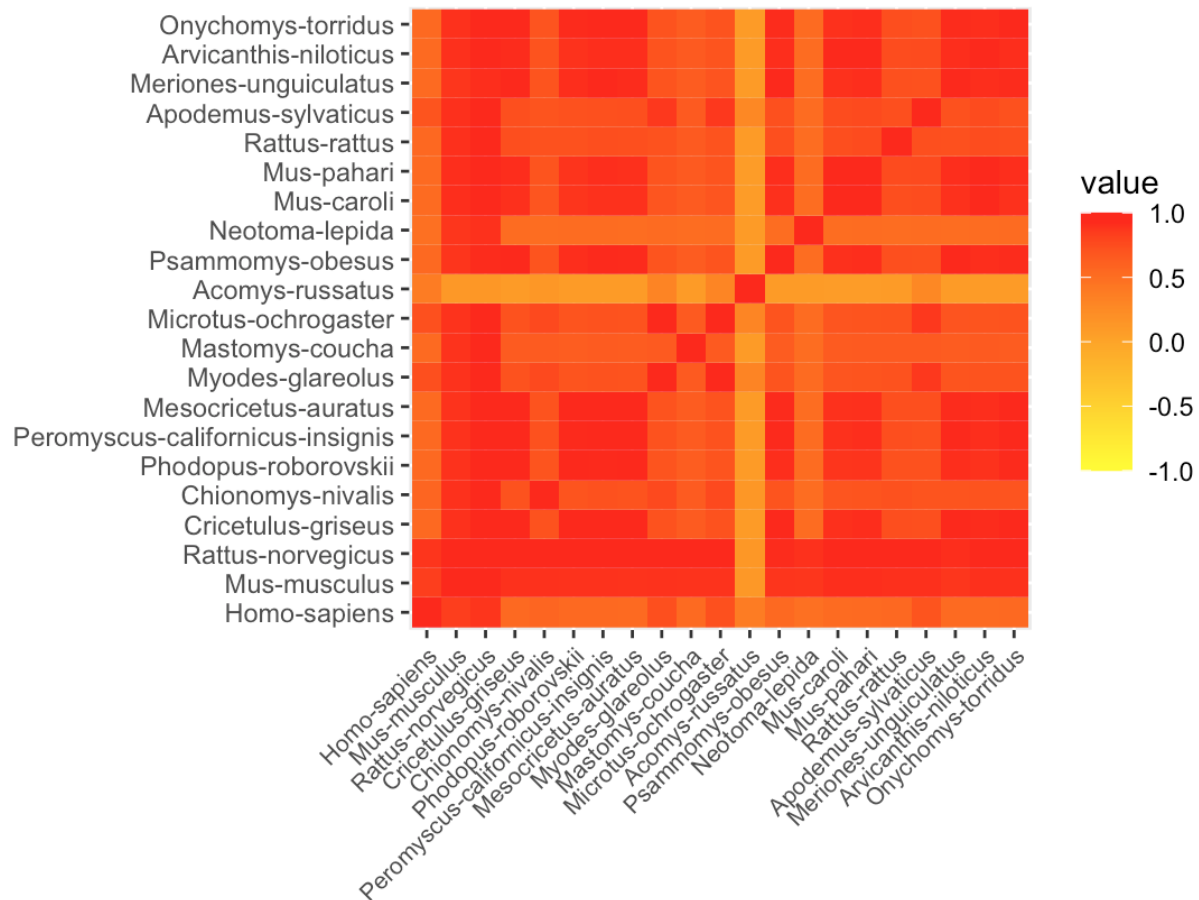
[Q6] Create a phylogenetic tree, using either a parsimony or distance-based approach. Bootstrapping and tree rooting are optional. Use “simple phylogeny” online from the EBI or any respected phylogeny program (such as MEGA, PAUP, or Phylip). Paste an image of your Cladogram or tree output in your report.

Tree created through MEGA.



[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 (structured), Method used to solve the structure (experimental Technique), resolution (resolution), and source organism (source).

Similar atomic resolution structures:

I used two methods to validate this date, PDB website and the NCBI Blast website. Both methods only provided these three structures.

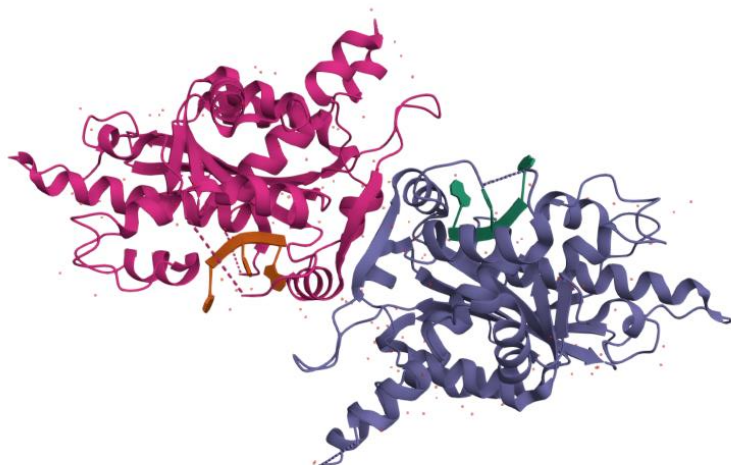
ID	Technique	Resolution	Source	Evalue	Identity
2FC6	SOLUTION NMR		Homo sapiens	1.99e-12	76

3D45	X-RAY DIFFRACTION	3 Å	Mus musculus	7.51e-7	26
2A1S	X-RAY DIFFRACTION	2.6 Å	Homo sapiens	1.068e-7	29

[Q9] 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?

2FC6 is very likely to be similar in structure to Mus musculus TOE1 like protein given the high sequence similarity (>76%).



[Q10] Perform a “Target” search of ChEMBEL (<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, there are five different targets. Within those targets, there are several assays and ligand efficiency data that could be useful starting points for exploring potential inhibition of the novel protein. One strong contender is purine-2,6-dione which inhibited PARN. PARN has a similar structure to the novel protein.

ChEMBL

Search in ChEMBL

Show Full Query

5 Targets
0 Selected - Select All
Browse Activities

Table Heatmap CSV TSV

Filters

- Organism Taxonomy L1
 - Eukaryotes 4
 - Bacteria 1
- Organism Taxonomy L2
 - Mammalia 4
 - Grain-Positive 1
- Organism Taxonomy L3
 - Primates 2
 - Mycobacterium 1
 - Rodentia 1
- Organism
 - Homo sapiens 3
 - Mus musculus 1
 - Mycobacterium tuberculosis 1
- Protein Classification L1
 - Enzyme 5
- Protein Classification L2
 - Hydrolase 2
 - Kinase 2
 - Ligase 1

Records per page: 20

Show/Hide Columns

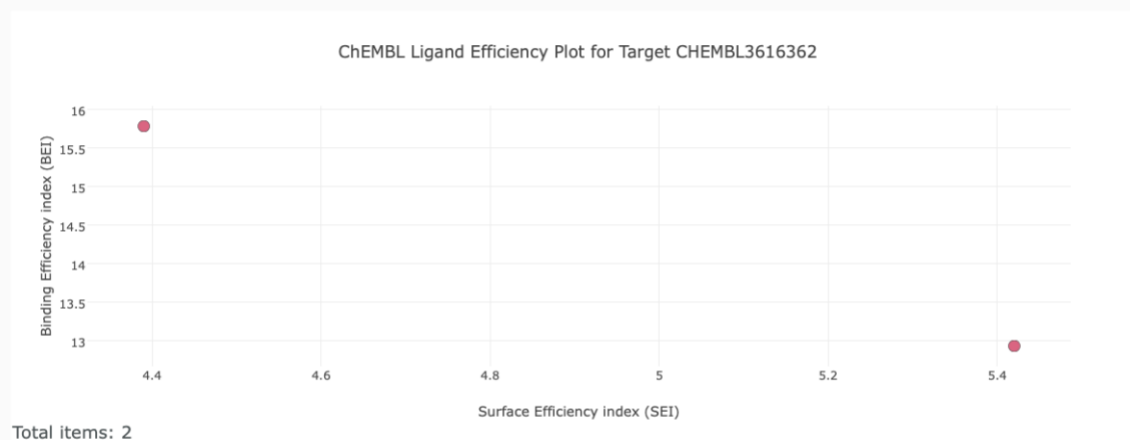
Showing 1-5 out of 5 records

<input type="checkbox"/>	E-Value	Positives %	Identities %	Score (bits)	Score	Length	ChEMBL ID	Name	UniProt Accessions	Type	Organism	Compounds	Activities
<input type="checkbox"/>	9.2e-9	46.5	27.1	55.0694	131	639	CHEMBL3616362	Poly(A)-specific ribonuclease PARN	O95453	SINGLE PROTEIN	Homo sapiens	22	22
<input type="checkbox"/>	0.21	40.9	23.6	31.9574	71	285	CHEMBL3616361	CCR4-NOT transcription complex subunit 7	Q9UIV1	SINGLE PROTEIN	Homo sapiens	32	32
<input type="checkbox"/>	3.2	51.6	30.6	28.4906	62	452	CHEMBL169231	Pup--protein ligase	P9WNU7	SINGLE PROTEIN	Mycobacterium tuberculosis	6	18
<input type="checkbox"/>	3.7	50	35.2	28.4906	62	1378	CHEMBL1795170	Macrophage-stimulating protein receptor	Q62190	SINGLE PROTEIN	Mus musculus	7	7
<input type="checkbox"/>	8.8	41.2	30.9	27.335	59	3661	CHEMBL1795195	Serine/threonine-protein kinase SMG1	Q96Q15	SINGLE PROTEIN	Homo sapiens	47	61

Showing 1-5 out of 5 records

Ligand Efficiencies

See all bioactivities for target CHEMBL3616362 used in this visualisation



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

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

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

Target Report Card

Name And Classification

ID:	CHEMBL3616362
Type:	SINGLE PROTEIN
Preferred Name:	Poly(A)-specific ribonuclease PARN
Synonyms:	DAN Deadenylating nuclease Deadenylation nuclease PARN Polyadenylate-specific ribonuclease Poly(A)-specific ribonuclease PARN
Organism:	Homo sapiens
Species Group:	No
Protein Target Classification:	- Enzyme > Hydrolase

Assay Report Card

Basic Information

Assay ID:	CHEMBL3619398
Type:	Binding
Description:	Inhibition of PARN (unknown origin)
Format:	BAO_0000357
Journal:	Bioorg Med Chem Lett (2015) 25:4219-4224
Organism:	Homo sapiens
Strain:	---
Tissue:	---
Cell Type:	---
Subcellular Fraction:	---
Target:	CHEMBL3616362
Document:	CHEMBL3616430
Cell:	
Tissue:	

Discovery, synthesis and biochemical profiling of purine-2,6-dione derivatives as inhibitors of the human poly(A)-selective ribonuclease Caf1

<https://doi.org/10.1016/j.bmcl.2015.07.095>.

