# Find a Gene Assignment

Alex Cagle

2023-04-30

## Question 1

**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.**

> I am interested in studying the huntingtin (Htt) protein, specifically huntingtin isoform 1, which is involved in Huntington's Disease. This protein is found in humans (Homo sapiens) and its accession number from NCBI RefSeq is NP_001375421.1.

> Link: https://www.ncbi.nlm.nih.gov/protein/NP_001375421.1

## Question 2

**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).**

**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 . 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.**

**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.**

**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.**

> I used NCBI BLAST, specifically tblastn, to find genes relating to the huntingtin isoform 1 protein (NP_001375421.1) with the organism set to Homo (taxid: 9605). The database I searched was the nucleotide collection (nt) database. The top hit was Homo sapiens huntingtin (HTT), transcript

variant 1, mRNA, which corresponds to Homo sapiens (humans), has a query cover of 100%, an E-value of 0.0, a percent identity of 92.62%, and the accession number NM_001388492.1. Screenshot is attached with results from the BLAST search.

**Homo sapiens huntingtin (HTT), transcript variant 1, mRNA**

Sequence ID: **NM_001388492.1**   Length: **13472**   Number of Matches: **1**

Range 1: 146 to 9571 GenBank    Graphics                  ▼ Next Match ▲ Previous Match

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 5907 bits(15323) | 0.0 | Compositional matrix adjust. | 3142/3142(100%) | 3142/3142(100%) | 0/3142(0%) | +2 |

```
Query  1     MATLEKLMKAFESLKSFqqqqqqqqqqqqqqqqqqqqqppppppppppppqlpqpppqaqp  60
             MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPPPQLPQPPPQAQP
Sbjct  146   MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPPPQLPQPPPQAQP  325

Query  61    llpqpqppppppppppGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE  120
             LLPQPQPPPPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
Sbjct  326   LLPQPQPPPPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE  505

Query  121   FQKLLGIAMELFLLCSDDAESDVRMVADECLNKVIKALMDSNLPRLQLELYKEIKKNGAP  180
             FQKLLGIAMELFLLCSDDAESDVRMVADECLNKVIKALMDSNLPRLQLELYKEIKKNGAP
Sbjct  506   FQKLLGIAMELFLLCSDDAESDVRMVADECLNKVIKALMDSNLPRLQLELYKEIKKNGAP  685

Query  181   RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG  240
             RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG
Sbjct  686   RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG  865

Query  241   NFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWllnvllgllv  300
             NFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWLLNVLLGLLV
Sbjct  866   NFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWLLNVLLGLLV  1045

Query  301   pvEDEHSTllilgvlltlrylvpllQQQVKDTSLKGSFGVTRKEMEVSPSAEQLVQVYEL  360
             PVEDEHSTLLILGVLLTLRYLVPLLQQQVKDTSLKGSFGVTRKEMEVSPSAEQLVQVYEL
Sbjct  1046  PVEDEHSTLLILGVLLTLRYLVPLLQQQVKDTSLKGSFGVTRKEMEVSPSAEQLVQVYEL  1225

Query  361   TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAAKEEsggrsrsgsI  420
             TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAAKEESGGRSRSGSI
Sbjct  1226  TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAAKEESGGRSRSGSI  1405

Query  421   VELIAGGGSSCSPVLSRKQKGKVLLGEEEALeddsesrsdvsssALTASVKDEISGELAA  480
             VELIAGGGSSCSPVLSRKQKGKVLLGEEEALEDDSESRSDVSSSALTASVKDEISGELAA
Sbjct  1406  VELIAGGGSSCSPVLSRKQKGKVLLGEEEALEDDSESRSDVSSSALTASVKDEISGELAA  1585

Query  481   SSGVSTPGSAGHDIITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILshsssqvsav  540
             SSGVSTPGSAGHDIITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSSQVSAV
Sbjct  1586  SSGVSTPGSAGHDIITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSSQVSAV  1765

Query  541   psdpaMDLNDGTQASSPISDSSQTTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD  600
             PSDPAMDLNDGTQASSPISDSSQTTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
Sbjct  1766  PSDPAMDLNDGTQASSPISDSSQTTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD  1945
```

# Question 3

**Gather information about this "novel" protein. At a minimum, show me the protein sequence of the "novel" protein as displayed in your BLAST results from [Q2] as FASTA format (you can copy and paste the aligned sequence subject lines from your BLAST result page if necessary) or translate your novel DNA sequence using a tool called EMBOSS Transeq at the EBI. Don't forget to translate all six reading frames; the ORF (open reading frame) is likely to be the longest sequence without a stop codon. It may not start with a methionine if you don't have the complete coding region. Make sure the sequence you provide includes a header/subject line and is in traditional FASTA format. Here, tell me the name of the novel protein, and the species from which it derives. It is very unlikely (but still definitely possible) that you will find a novel gene from an organism such as S. cerevisiae, human or mouse, because those genomes have already been thoroughly annotated. It is more likely that you will discover a new gene in a genome that is currently being sequenced, such as bacteria or plants or protozoa.**

The protein I found is from humans and is called huntingtin isoform 1.

FASTA file for the protein I chose:

>NP_001375421.1 huntingtin isoform 1 [Homo sapiens]
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPPPQLPQPPPQAQPLLPQPQPPPPPP
AEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPEFQKLLGIAMELFLLCSDDAESDVRMVADECLNKVII
SNLPRLQLELYKEIKKNGAPRSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKII
NFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWLLNVLLGLLVPVEDEHSTLLILGVLLTI
LVPLLQQQVKDTSLKGSFGVTRKEMEVSPSAEQLVQVYELTLHHTQHQDHNVVTGALELLQQLFRTPPPELLQT
IGQLTAAKEESGGRSRSGSIVELIAGGGSSCSPVLSRKQKGKVLLGEEEALEDDSESRSDVSSSALTASVKDEISGEI

SSGVSTPGSAGHDIITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSSQVSAVPSDPAMDLNDGTQASS

SSQTTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQDEDEEATGILPDEASEAFRNSSMALQQAHLLKNMSH

DSSVDKFVLRDEATEPGDQENKPCRIKGDIGQSTDDDSAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALA

AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSILSRSRFHVGDWMGTIRT

TFSLADCIPLLRKTLKDESSVTCKLACTAVRNCVMSLCSSSYSELGLQLIIDVLTLRNSSYWLVRTELLETLAEIDF

SFLEAKAENLHRGAHHYTGLLKLQERVLNNVVIHLLGDEDPRVRHVAAASLIRLVPKLFYKCDQGQADPVVAVA

YLKLLMHETQPPSHFSVSTITRIYRGYNLLPSITDVTMENNLSRVIAAVSHELITSTTRALTFGCCEALCLLSTAFP

WSLGWHCGVPPLSASDESRKSCTVGMATMILTLLSSAWFPLDLSAHQDALILAGNLLAASAPKSLRSSWASEEEA

KQEEVWPALGDRALVPMVEQLFSHLLKVINICAHVLDDVAPGPAIKAALPSLTNPPSLSPIRRKGKEKEPGEQAS

KKGSEASAASRQSDTSGPVTTSKSSSLGSFYHLPSYLKLHDVLKATHANYKVTLDLQNSTEKFGGFLRSALDVLS

ATLQDIGKCVEEILGYLKSCFSREPMMATVCVQQLLKTLFGTNLASQFDGLSSNPSKSQGRAQRLGSSSVRPGLY

APYTHFTQALADASLRNMVQAEQENDTSGWFDVLQKVSTQLKTNLTSVTKNRADKNAIHNHIRLFEPLVIKALK

CVQLQKQVLDLLAQLVQLRVNYCLLDSDQVFIGFVLKQFEYIEVGQFRESEAIIPNIFFFLVLLSYERYHSKQIIGIP

IQLCDGIMASGRKAVTHAIPALQPIVHDLFVLRGTNKADAGKELETQKEVVVSMLLRLIQYHQVLEMFILVLQQC

DKWKRLSRQIADIILPMLAKQQMHIDSHEALGVLNTLFEILAPSSLRPVDMLLRSMFVTPNTMASVSTVQLWISG

RVLISQSTEDIVLSRIQELSFSPYLISCTVINRLRDGDSTSTLEEHSEGKQIKNLPEETFSRFLLQLVGILLEDIVTKQL

KVEMSEQQHTFYCQELGTLLMCLIHIFKSGMFRRITAAATRLFRSDGCGGSFYTLDSLNLRARSMITTHPALVLL

LLVNHTDYRWWAEVQQTPKRHSLSSTKLLSPQMSGEEEDSDLAAKLGMCNREIVRRGALILFCDYVCQNLHDSF

VNHIQDLISLSHEPPVQDFISAVHRNSAASGLFIQAIQSRCENLSTPTMLKKTLQCLEGIHLSQSGAVLTLYVDRLLC

FRVLARMVDILACRRVEMLLAANLQSSMAQLPMEELNRIQEYLQSSGLAQRHQRLYSLLDRFRLSTMQDSLSPSP

PLDGDGHVSLETVSPDKDWYVHLVKSQCWTRSDSALLEGAELVNRIPAEDMNAFMMNSEFNLSLLAPCLSLGM

KSALFEAAREVTLARVSGTVQQLPAVHHVFQPELPAEPAAYWSKLNDLFGDAALYQSLPTLARALAQYLVVVSK

LPPEKEKDIVKFVVATLEALSWHLIHEQIPLSLDLQAGLDCCCLALQLPGLWSVVSSTEFVTHACSLIYCVHFILEA

QPGEQLLSPERRTNTPKAISEEEEEVDPNTQNPKYITAACEMVAEMVESLQSVLALGHKRNSGVPAFLTPLLRNI

RLPLVNSYTRVPPLVWKLGWSPKPGGDFGTAFPEIPVEFLQEKEVFKEFIYRINTLGWTSRTQFEETWATLLGV

VMEQEESPPEEDTERTQINVLAVQAITSLVLSAMTVPVAGNPAVSCLEQQPRNKPLKALDTRFGRKLSIIRGIVEQ

MVSKRENIATHHLYQAWDPVPSLSPATTGALISHEKLLLQINPERELGSMSYKLGQVSIHSVWLGNSITPLREEEW

EEEADAPAPSSPPTSPVNSRKHRAGVDIHSCSQFLLELYSRWILPSSSARRTPAILISEVVRSLLVVSDLFTERNQFE

YVTLTELRRVHPSEDEILAQYLVPATCKAAAVLGMDKAVAEPVSRLLESTLRSSHLPSRVGALHGVLYVLECDLLI

QLIPVISDYLLSNLKGIAHCVNIHSQQHVLVMCATAFYLIENYPLDVGPEFSASIIQMCGVMLSGSEESTPSIIYHCA

GLERLLLSEQLSRLDAESLVKLSVDRVNVHSPHRAMAALGLMLTCMYTGKEKVSPGRTSDPNPAAPDSESVIVAM

LFDRIRKGFPCEARVVARILPQFLDDFFPPQDIMNKVIGEFLSNQQPYPQFMATVVYKVFQTLHSTGQSSMVRDW

SNFTQRAPVAMATWSLSCFFVSASTSPWVAAILPHVISRMGKLEQVDVNLFCLVATDFYRHQIEEELDRRAFQSV

APGSPYHRLLTCLRNVHKVTTC

FASTA file for the gene I found:

>NM_001388492.1 Homo sapiens huntingtin (HTT), transcript variant 1, mRNA

GCTGCCGGGACGGGTCCAAGATGGACGGCCGCTCAGGTTCTGCTTTTACCTGCGGCCCAGAGCCCCATTC

TGCTGAGCGGCGCCGCGAGTCGGCCCGAGGCCTCCGGGGACTGCCGTGCCGGGCGGGAGACCGCCATGG

AAGCTGATGAAGGCCTTCGAGTCCCTCAAGTCCTTCCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGC

GCAGCAGCAGCAGCAACAGCCGCCACCGCCGCCGCCGCCGCCGCCTCCTCAGCTTCCTCAGCCGCCG

AGCCGCTGCTGCCTCAGCCGCAGCCGCCCCCGCCGCCGCCCCCGCCGCCACCCGGCCCGGCTGTGGCTGA
CACCGACCAAAGAAAGAACTTTCAGCTACCAAGAAAGACCGTGTGAATCATTGTCTGACAATATGTGAAAA
ACAGTCTGTCAGAAATTCTCCAGAATTTCAGAAACTTCTGGGCATCGCTATGGAACTTTTTCTGCTGTGCA
CAGAGTCAGATGTCAGGATGGTGGCTGACGAATGCCTCAACAAAGTTATCAAAGCTTTGATGGATTCTAAT
TTACAGCTCGAGCTCTATAAGGAAATTAAAAAGAATGGTGCCCCTCGGAGTTTGCGTGCTGCCCTGTGGAC
GCTGGCTCACCTGGTTCGGCCTCAGAAATGCAGGCCTTACCTGGTGAACCTTCTGCCGTGCCTGACTCGA
GACCCGAAGAATCAGTCCAGGAGACCTTGGCTGCAGCTGTTCCCAAAATTATGGCTTCTTTTGGCAATTT
AATGAAATTAAGGTTTTGTTAAAGGCCTTCATAGCGAACCTGAAGTCAAGCTCCCCACCATTCGGCGGAC
ATCAGCAGTGAGCATCTGCCAGCACTCAAGAAGGACACAATATTTCTATAGTTGGCTACTAAATGTGCTCT
TCGTTCCTGTCGAGGATGAACACTCCACTCTGCTGATTCTTGGCGTGCTGCTCACCCTGAGGTATTTGGT
CAGCAGCAGGTCAAGGACACAAGCCTGAAAGGCAGCTTCGGAGTGACAAGGAAAGAAATGGAAGTCTCTC
GCAGCTTGTCCAGGTTTATGAACTGACGTTACATCATACACAGCACCAAGACCACAATGTTGTGACCGGAG
TGTTGCAGCAGCTCTTCAGAACGCCTCCACCCGAGCTTCTGCAAACCCTGACCGCAGTCGGGGGCATTGG
GCTGCTAAGGAGGAGTCTGGTGGCCGAAGCCGTAGTGGGAGTATTGTGGAACTTATAGCTGGAGGGGGTT
CCCTGTCCTTTCAAGAAAACAAAAAGGCAAAGTGCTCTTAGGAGAAGAAGAAGCCTTGGAGGATGACTCT
CGGATGTCAGCAGCTCTGCCTTAACAGCCTCAGTGAAGGATGAGATCAGTGGAGAGCTGGCTGCTTCTTC
ACTCCAGGGTCAGCAGGTCATGACATCATCACAGAACAGCCACGGTCACAGCACACACTGCAGGCGGACT
GGCCAGCTGTGACTTGACAAGCTCTGCCACTGATGGGGATGAGGAGGATATCTTGAGCCACAGCTCCAGC
CCGTCCCATCTGACCCTGCCATGGACCTGAATGATGGGACCCAGGCCTCGTCGCCCATCAGCGACAGCTC
ACCGAAGGGCCTGATTCAGCTGTTACCCCTTCAGACAGTTCTGAAATTGTGTTAGACGGTACCGACAACCA
CCTGCAGATTGGACAGCCCCAGGATGAAGATGAGGAAGCCACAGGTATTCTTCCTGATGAAGCCTCGGAG
ACTCTTCCATGGCCCTTCAACAGGCACATTTATTGAAAAACATGAGTCACTGCAGGCAGCCTTCTGACAG
AAATTTGTGTTGAGAGATGAAGCTACTGAACCGGGTGATCAAGAAAACAAGCCTTGCCGCATCAAAGGTG
GTCCACTGATGATGACTCTGCACCTCTTGTCCATTGTGTCCGCCTTTTATCTGCTTCGTTTTTGCTAACAG
ATGTGCTGGTTCCGGACAGGGATGTGAGGGTCAGCGTGAAGGCCCTGGCCCTCAGCTGTGTGGGAGCAGC
CACCCGGAATCTTTCTTCAGCAAACTCTATAAAGTTCCTCTTGACACCACGGAATACCCTGAGGAACAGTAT
CATCTTGAACTACATCGATCATGGAGACCCACAGGTTCGAGGAGCCACTGCCATTCTCTGTGGGACCCTCA
TCCTCAGCAGGTCCCGCTTCCACGTGGGAGATTGGATGGGCACCATTAGAACCCTCACAGGAAATACATT
GATTGCATTCCTTTGCTGCGGAAAACACTGAAGGATGAGTCTTCTGTTACTTGCAAGTTAGCTTGTACAGC
CTGTGTCATGAGTCTCTGCAGCAGCAGCTACAGTGAGTTAGGACTGCAGCTGATCATCGATGTGCTGACTC
GTTCCTATTGGCTGGTGAGGACAGAGCTTCTGGAAACCCTTGCAGAGATTGACTTCAGGCTGGTGAGCTT
AAAGCAGAAAACTTACACAGAGGGGCTCATCATTATACAGGGCTTTTAAAACTGCAAGAACGAGTGCTCAA
CATCCATTTGCTTGGAGATGAAGACCCCAGGGTGCGACATGTTGCCGCAGCATCACTAATTAGGCTTGTCC
TTTATAAATGTGACCAAGGACAAGCTGATCCAGTAGTGGCCGTGGCAAGAGATCAAAGCAGTGTTTACCTC
ATGCATGAGACGCAGCCTCCATCTCATTTCTCCGTCAGCACAATAACCAGAATATATAGAGGCTATAACCTA
CATAACAGACGTCACTATGGAAAATAACCTTTCAAGAGTTATTGCAGCAGTTTCTCATGAACTAATCACATC
GAGCACTCACATTTGGATGCTGTGAAGCTTTGTGTCTTCTTTCCACTGCCTTCCCAGTTTGCATTTGGAGT
CACTGTGGAGTGCCTCCACTGAGTGCCTCAGATGAGTCTAGGAAGAGCTGTACCGTTGGGATGGCCACAA
CCTGCTCTCGTCAGCTTGGTTCCCATTGGATCTCTCAGCCCATCAAGATGCTTTGATTTTGGCCGGAAACT
CCAGTGCTCCCAAATCTCTGAGAAGTTCATGGGCCTCTGAAGAAGAAGCCAACCCAGCAGCCACCAAGCA
TGGCCAGCCCTGGGGGACCGGGCCCTGGTGCCCATGGTGGAGCAGCTCTTCTCTCACCTGCTGAAGGTGA

5

TGCCCACGTCCTGGATGACGTGGCTCCTGGACCCGCAATAAAGGCAGCCTTGCCTTCTCTAACAAACCCC
GTCCCATCCGACGAAAGGGGAAGGAGAAAGAACCAGGAGAACAAGCATCTGTACCGTTGAGTCCCAAGAA
GCCAGTGCAGCTTCTAGACAATCTGATACCTCAGGTCCTGTTACAACAAGTAAATCCTCATCACTGGGGAG
TCTTCCTTCATACCTCAAACTGCATGATGTCCTGAAAGCTACACACGCTAACTACAAGGTCACGCTGGATC
GCACGGAAAAGTTTGGAGGGTTTCTCCGCTCAGCCTTGGATGTTCTTTCTCAGATACTAGAGCTGGCCAC
ATTGGGAAGTGTGTTGAAGAGATCCTAGGATACCTGAAATCCTGCTTTAGTCGAGAACCAATGATGGCAA
TCAACAATTGTTGAAGACTCTCTTTGGCACAAACTTGGCCTCCCAGTTTGATGGCTTATCTTCCAACCCCA
AAGGCCGAGCACAGCGCCTTGGCTCCTCCAGTGTGAGGCCAGGCTTGTACCACTACTGCTTCATGGCCCC
TTCACCCAGGCCCTCGCTGACGCCAGCCTGAGGAACATGGTGCAGGCGGAGCAGGAGAACGACACCTCGG
TGTCCTCCAGAAAGTGTCTACCCAGTTGAAGACAAACCTCACGAGTGTCACAAAGAACCGTGCAGATAAGA
ATAATCACATTCGTTTGTTTGAACCTCTTGTTATAAAAGCTTTAAAACAGTACACGACTACAACATGTGTG
AAGCAGGTTTTAGATTTGCTGGCGCAGCTGGTTCAGTTACGGGTTAATTACTGTCTTCTGGATTCAGATCA
TGGCTTTGTATTGAAACAGTTTGAATACATTGAAGTGGGCCAGTTCAGGGAATCAGAGGCAATCATTCCAA
TCTTCTTGGTATTACTATCTTATGAACGCTATCATTCAAAACAGATCATTGGAATTCCTAAAATCATTCAGC
GGCATCATGGCCAGTGGAAGGAAGGCTGTGACACATGCCATACCGGCTCTGCAGCCCATAGTCCACGACCT
AAGAGGAACAAATAAAGCTGATGCAGGAAAAGAGCTTGAAACCCAAAAAGAGGTGGTGGTGTCAATGTTA
TCCAGTACCATCAGGTGTTGGAGATGTTCATTCTTGTCCTGCAGCAGTGCCACAAGGAGAATGAAGACAAG
CTGTCTCGACAGATAGCTGACATCATCCTCCCAATGTTAGCCAAACAGCAGATGCACATTGACTCTCATGA
AGTGTTAAATACATTATTTGAGATTTTGGCCCCTTCCTCCCTCCGTCCGGTAGACATGCTTTTACGGAGTA
CTCCAAACACAATGGCGTCCGTGAGCACTGTTCAACTGTGGATATCGGGAATTCTGGCCATTTTGAGGGT
CAGTCAACTGAAGATATTGTTCTTTCTCGTATTCAGGAGCTCTCCTTCTCTCCGTATTTAATCTCCTGTACA
TAGGTTAAGAGATGGGGACAGTACTTCAACGCTAGAAGAACACAGTGAAGGGAAACAAATAAAGAATTTG
CATTTTCAAGGTTTCTATTACAACTGGTTGGTATTCTTTTAGAAGACATTGTTACAAAACAGCTGAAGGTG
GAGCAGCAACATACTTTCTATTGCCAGGAACTAGGCACACTGCTAATGTGTCTGATCCACATCTTCAAGTC
CCGGAGAATCACAGCAGCTGCCACTAGGCTGTTCCGCAGTGATGGCTGTGGCGGCAGTTTCTACACCCTG
ACTTGCGGGCTCGTTCCATGATCACCACCCACCCGGCCCTGGTGCTGCTCTGGTGTCAGATACTGCTGCT
ACCGACTACCGCTGGTGGGCAGAAGTGCAGCAGACCCCGAAAAGACACAGTCTGTCCAGCACAAAGTTAC
GATGTCTGGAGAAGAGGAGGATTCTGACTTGGCAGCCAAACTTGGAATGTGCAATAGAGAAATAGTACGA
TCATTCTCTTCTGTGATTATGTCTGTCAGAACCTCCATGACTCCGAGCACTTAACGTGGCTCATTGTAAAT
GATCTGATCAGCCTTTCCCACGAGCCTCCAGTACAGGACTTCATCAGTGCCGTTCATCGGAACTCTGCTGC
GTTCATCCAGGCAATTCAGTCTCGTTGTGAAAACCTTTCAACTCCAACCATGCTGAAGAAAACTCTTCAGT
GGATCCATCTCAGCCAGTCGGGAGCTGTGCTCACGCTGTATGTGGACAGGCTTCTGTGCACCCCTTTCCGT
CGCATGGTCGACATCCTTGCTTGTCGCCGGGTAGAAATGCTTCTGGCTGCAAATTTACAGAGCAGCATGGC
AATGGAAGAACTCAACAGAATCCAGGAATACCTTCAGAGCAGCGGGCTCGCTCAGAGACACCAAAGGCTC
TGGACAGGTTTCGTCTCTCCACCATGCAAGACTCACTTAGTCCCTCTCCTCCAGTCTCTTCCCACCCGCTG
GGGCACGTGTCACTGGAAACAGTGAGTCCGGACAAAGACTGGTACGTTCATCTTGTCAAATCCCAGTGTT
AGATTCTGCACTGCTGGAAGGTGCAGAGCTGGTGAATCGGATTCCTGCTGAAGATATGAATGCCTTCATGA
AGTTCAACCTAAGCCTGCTAGCTCCATGCTTAAGCCTAGGGATGAGTGAAATTTCTGGTGGCCAGAAGAGT
GAAGCAGCCCGTGAGGTGACTCTGGCCCGTGTGAGCGGCACCGTGCAGCAGCTCCCTGCTGTCCATCATG
CGAGCTGCCTGCAGAGCCGGCGGCCTACTGGAGCAAGTTGAATGATCTGTTTGGGGATGCTGCACTGTAT
CCACTCTGGCCCGGGCCCTGGCACAGTACCTGGTGGTGGTCTCCAAACTGCCCAGTCATTTGCACCTTCC

GAGAAGGACATTGTGAAATTCGTGGTGGCAACCCTTGAGGCCCTGTCCTGGCATTTGATCCATGAGCAGA
TCTGGATCTCCAGGCAGGGCTGGACTGCTGCTGCCTGGCCCTGCAGCTGCCTGGCCTCTGGAGCGTGGTG
AGTTTGTGACCCACGCCTGCTCCCTCATCTACTGTGTGCACTTCATCCTGGAGGCCGTTGCAGTGCAGCC
CTTCTTAGTCCAGAAAGAAGGACAAATACCCCAAAAGCCATCAGCGAGGAGGAGGAGGAAGTAGATCCAA
TCCTAAGTATATCACTGCAGCCTGTGAGATGGTGGCAGAAATGGTGGAGTCTCTGCAGTCGGTGTTGGCC
AAAGGAATAGCGGCGTGCCGGCGTTTCTCACGCCATTGCTAAGGAACATCATCATCAGCCTGGCCCGCCTG
AACAGCTACACACGTGTGCCCCCACTGGTGTGGAAGCTTGGATGGTCACCCAAACCGGGAGGGGATTTTG
CCCTGAGATCCCCGTGGAGTTCCTCCAGGAAAAGGAAGTCTTTAAGGAGTTCATCTACCGCATCAACACAC
CCAGTCGTACTCAGTTTGAAGAAACTTGGGCCACCCTCCTTGGTGTCCTGGTGACGCAGCCCCTCGTGATC
GAGAGCCCACCAGAAGAAGACACAGAGAGGACCCAGATCAACGTCCTGGCCGTGCAGGCCATCACCTCAC
TGCAATGACTGTGCCTGTGGCCGGCAACCCAGCTGTAAGCTGCTTGGAGCAGCAGCCCCGGAACAAGCCT
TCGACACCAGGTTTGGGAGGAAGCTGAGCATTATCAGAGGGATTGTGGAGCAAGAGATTCAAGCAATGGT
GAGAATATTGCCACCCATCATTTATATCAGGCATGGGATCCTGTCCCTTCTCTGTCTCCGGCTACTACAGGT
CAGCCACGAGAAGCTGCTGCTACAGATCAACCCCGAGCGGGAGCTGGGGAGCATGAGCTACAAACTCGGC
TACACTCCGTGTGGCTGGGGAACAGCATCACACCCCTGAGGGAGGAGGAATGGGACGAGGAAGAGGAGGA
GCCCCTGCACCTTCGTCACCACCCACGTCTCCAGTCAACTCCAGGAAACACCGGGCTGGAGTTGACATCC
GCAGTTTTTGCTTGAGTTGTACAGCCGCTGGATCCTGCCGTCCAGCTCAGCCAGGAGGACCCCGGCCATC
AGGTGGTCAGATCCCTTCTAGTGGTCTCAGACTTGTTCACCGAGCGCAACCAGTTTGAGCTGATGTATGTC
GAACTGCGAAGGGTGCACCCTTCAGAAGACGAGATCCTCGCTCAGTACCTGGTGCCTGCCACCTGCAAGG
CCTTGGGATGGACAAGGCCGTGGCGGAGCCTGTCAGCCGCCTGCTGGAGAGCACGCTCAGGAGCAGCCA
GGGGTTGGAGCCCTGCACGGCGTCCTCTATGTGCTGGAGTGCGACCTGCTGGACGACACTGCCAAGCAGCT
ATCAGCGACTATCTCCTCTCCAACCTGAAAGGGATCGCCCACTGCGTGAACATTCACAGCCAGCAGCACGT
GTGTGCCACTGCGTTTTACCTCATTGAGAACTATCCTCTGGACGTAGGGCCGGAATTTTCAGCATCAATAA
GTGGGGTGATGCTGTCTGGAAGTGAGGAGTCCACCCCCTCCATCATTTACCACTGTGCCCTCAGAGGCCTC
CTGCTCTCTGAGCAGCTCTCCCGCCTGGATGCAGAATCGCTGGTCAAGCTGAGTGTGGACAGAGTGAACG
GCACCGGGCCATGGCGGCTCTGGGCCTGATGCTCACCTGCATGTACACAGGAAAGGAGAAAGTCAGTCCG
CAGACCCTAATCCTGCAGCCCCCGACAGCGAGTCAGTGATTGTTGCTATGGAGCGGGTATCTGTTCTTTTT
AGGAAAGGCTTTCCTTGTGAAGCCAGAGTGGTGGCCAGGATCCTGCCCCAGTTTCTAGACGACTTCTTCC
CATCATGAACAAAGTCATCGGAGAGTTTCTGTCCAACCAGCAGCCATACCCCCAGTTCATGGCCACCGTGG
TGTTTCAGACTCTGCACAGCACCGGGCAGTCGTCCATGGTCCGGGACTGGGTCATGCTGTCCCTCTCCAA
AGGGCCCCGGTCGCCATGGCCACGTGGAGCCTCTCCTGCTTCTTTGTCAGCGCGTCCACCAGCCCGTGGG
CCTCCCACATGTCATCAGCAGGATGGGCAAGCTGGAGCAGGTGGACGTGAACCTTTTCTGCCTGGTCGCC
ACAGACACCAGATAGAGGAGGAGCTCGACCGCAGGGCCTTCCAGTCTGTGCTTGAGGTGGTTGCAGCCCC
TATCACCGGCTGCTGACTTGTTTACGAAATGTCCACAAGGTCACCACCTGCTGAGCGCCATGGTGGGAGA
CGGCAGCTGGGGCCGGAGCCTTTGGAAGTCTGCGCCCTTGTGCCCTGCCTCCACCGAGCCAGCTTGGTCC
CCGCACATGCCGCGGGCGGCCAGGCAACGTGCGTGTCTCTGCCATGTGGCAGAAGTGCTCTTTGTGGCAC
GGGAGTGTCTGCAGTCCTGGTGGGGCTGAGCCTGAGGCCTTCCAGAAAGCAGGAGCAGCTGTGCTGCACC
GACCAGGTCCTTTCTCCTGATAGTCACCTGCTGGTTGTTGCCAGGTTGCAGCTGCTCTTGCATCTGGGCCA
CCTCCTGCAGGCTGGCTGTTGGCCCCTCTGCTGTCCTGCAGTAGAAGGTGCCGTGAGCAGGCTTTGGGAA
GGTCTCCCTGGTGGGGTGTGCATGCCACGCCCCGTGTCTGGATGCACAGATGCCATGGCCTGTGCTGGGC
GGGTGCTAGACACCCGGCACCATTCTCCCTTCTCTCTTTTCTTCTCAGGATTTAAAATTTAATTATATCAGT

7

TAATTTTAACGTAACTCTTTCTATGCCCGTGTAAAGTATGTGAATCGCAAGGCCTGTGCTGCATGCGACAG
TGGTGGACAGGGCCCCCGGCCACGCTCCCTCTCCTGTAGCCACTGGCATAGCCCTCCTGAGCACCCGCTGA
TGTACATGTTCCTGTTTATGCATTCACAAGGTGACTGGGATGTAGAGAGGCGTTAGTGGGCAGGTGGCCAG
GAGGACAGGCCCCCATTATCCTAGGGGTGCGCTCACCTGCAGCCCCTCCTCCTCGGGCACAGACGACTGTG
CCACCAGTCAGGGACAGCAGCCTCCCTGTCACTCAGCTGAGAAGGCCAGCCCTCCCTGGCTGTGAGCAGC
TCCAGAGACATGGGCCTCCCACTCCTGTTCCTTGCTAGCCCTGGGGTGGCGTCTGCCTAGGAGCTGGCTG
GGACCTGCTGCTCCATGGATGCATGCCCTAAGAGTGTCACTGAGCTGTGTTTTGTCTGAGCCTCTCTCGGT
AGCTTGGTGTCTTGGCACTGTTAGTGACAGAGCCCAGCATCCCTTCTGCCCCCGTTCCAGCTGACATCTTO
CCCTTTTAGTCAGGAGAGTGCAGATCTGTGCTCATCGGAGACTGCCCCACGGCCCTGTCAGAGCCGCCACT
AGGCCAGGTCCCTGGACCAGCCTCCTGTTTGCAGGCCCAGAGGAGCCAAGTCATTAAAATGGAAGTGGAT
CGGGCTGCTGCTGATGTAGGAGCTGGATTTGGGAGCTCTGCTTGCCGACTGGCTGTGAGACGAGGCAGGC
CTCAGCCCTAGAGGCGAGCCAGGCAAGGTTGGCGACTGTCATGTGGCTTGGTTTGGTCATGCCCGTCGAT
TTGAATGTGGTAAGTGGAGGAAATGTTGGAACTCTGTGCAGGTGCTGCCTTGAGACCCCCAAGCTTCCAC
CCTATGTGGCAGCTGGGGAGCAGCTGAGATGTGGACTTGTATGCTGCCCACATACGTGAGGGGGAGCTGA
CTCCTCTGAGCAGCCTCTGCCAGGCCTGTATGAGGCTTTTCCCACCAGCTCCCAACAGAGGCCTCCCCAC
CCTCGTCCTCGTGGCGGGGCAGCAGGAGCGGTAGAAAGGGGTCCGATGTTTGAGGAGGCCCTTAAGGGA
TATAACACGTAAGAAAATCACCATTCCGTATTGGTTGGGGGCTCCTGTTTCTCATCCTAGCTTTTTCCTGG
TAGAAGGTTTGGGAACGAGGGGAAAGTTCTCAGAACTGTTGGCTGCTCCCCACCCGCCTCCCGCCTCCCC
GTCAGCAGCTCTGAGACAGCAGTATCACAGGCCAGATGTTGTTCCTGGCTAGATGTTTACATTTGTAAGAA
TGAATGTAAAACAGAGCCATTCCCTTGGAATGCATATCGCTGGGCTCAACATAGAGTTTGTCTTCCTCTTG
TGATCTAAACCAGTCCTTAGCAAGGGGCTCAGAACACCCCGCTCTGGCAGTAGGTGTCCCCCACCCCCAAA
GTGTGCTCCGGAGATGAATATGAGCTCATTAGTAAAAATGACTTCACCCACGCATATACATAAAGTATCCAT
ATATAGACACATCTATAATTTTACACACACACCTCTCAAGACGGAGATGCATGGCCTCTAAGAGTGCCCGTG
TCCTGGAAGTTGACTTTCCTTAGACCCGCCAGGTCAAGTTAGCCGCGTGACGGACATCCAGGCGTGGGAC
CAGGGCTCATTCATTGCCCACTAGGATCCCACTGGCGAAGATGGTCTCCATATCAGCTCTCTGCAGAAGGG
TTATCATGTTCCTAAAAATCTGTGGCAAGCACCCATCGTATTATCCAAATTTTGTTGCAAATGTGATTAATT
AAGTTTTGGGGGTGGGCTGTGGGGAGATTGCTTTTGTTTTCCTGCTGGTAATATCGGGAAAGATTTTAAT
TAGAATTGTTTGGCAATGCACTGAAGCGTGTTTCTTTCCCAAAATGTGCCTCCCTTCCGCTGCGGGCCCAC
TGTAGGTGATGTTTCCAGCTGCCAAGTGCTCTTTGTTACTGTCCACCCTCATTTCTGCCAGCGCATGTGTC
GGAAAATGTGAAGCTGAACCCCCTCCAGACACCCAGAATGTAGCATCTGAGAAGGCCCTGTGCCCTAAAGO
GCCCCCATCTTCATGGAGGGGGTCATTTCAGAGCCCTCGGAGCCAATGAACAGCTCCTCCTCTTGGAGCTO
CCACGTGGAGCTCGGGACGGATAGTAGACAGCAATAACTCGGTGTGTGGCCGCCTGGCAGGTGGAACTTO
GGGGTGGAGTGAGGTTAGTTCTGTGTGTCTGGTGGGTGGAGTCAGGCTTCTCTTGCTACCTGTGAGCATO
GACATCCTCATCGGGCTTTGTCCCTCCCCGCTTCCTCCCTCTGCGGGGAGGACCCGGGACCACAGCTGC
AGACTTGGAGCTGTCCTCCAGAGGGGTCACGTGTAGGAGTGAGAAGAAGGAAGATCTTGAGAGCTGCTGA
AGAGCTCAGGATGGCTCAGACGAGGACACTCGCTTGCCGGGCCTGGGCCTCCTGGGAAGGAGGGAGCTGO
GCATGACAACTGAAGGCAACCTGGAAGGTTCAGGGGCCGCTCTTCCCCCATGTGCCTGTCACGCTCTGGT
GAACGCCTTCCCCTCAGTTGTTTCTAAGAGCAGAGTCTCCCGCTGCAATCTGGGTGGTAACTGCCAGCCT
TGGCCAACGTGGACCTGCCTACGGAGGGTGGGCTCTGACCCAAGTGGGGCCTCCTTGTCCAGGTCTCACT
GTGGTCAGAGGGACTGTCAGCTGAGCTTGAGCTCCCCTGGAGCCAGCAGGGCTGTGATGGGCGAGTCCC
CAGACCTGAATGCTTCTGAGAGCAAAGGGAAGGACTGACGAGAGATGTATATTTAATTTTTTAACTGCTGO
ACATCCAAATTAAAGGAAAAAAATGGAAACCA

# Question 4

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.

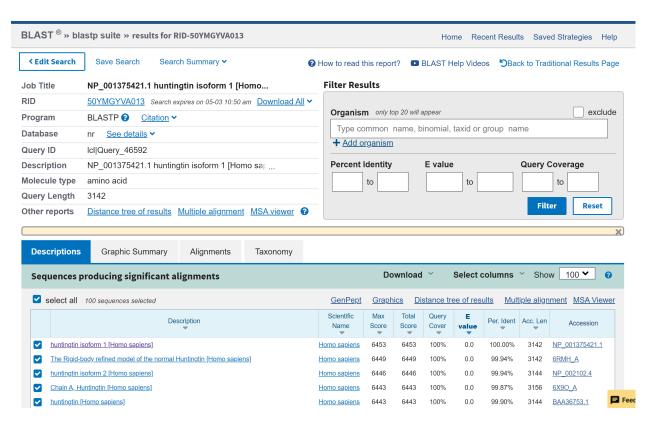> There are other proteins with percent identities less than 100%, but the top hit is 100%, since I used that exact sequence as the blastp query. Therefore, my protein is not novel?



Figure 1: Q4 search results