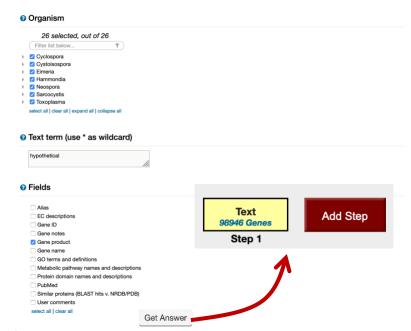# Sequence Exercises: Motifs, Domains and Colocation

1. **Using InterPro domain searches to identify unannotated kinesin motor proteins.**
   **Note: For this exercise use http://toxodb.org**

   a. Identify all genes annotated as hypothetical in organisms in ToxoDB (select the gene product field). Use the full text search and look for genes with the word *hypothetical* in their Gene products.
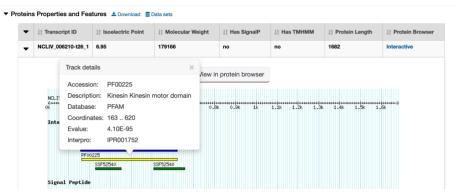
   ## Identify Genes based on Text (product name, notes, etc.)

   

   b. How many of these hypothetical genes have a kinesin-motor protein PFAM domain?
   - Add a step to the strategy. Go to the "Interpro Domain" search under 'Protein features and properties' similarity/pattern, start typing the work kinesin and it should autocomplete.

   

c. Go to the gene page for NCLIV_006210 and look at the protein feature section. Does this look like a possible motor protein?
   - Click on the ID for NCLIV_006210 in the result table to go to the gene page. Scroll down to the Protein Properties and Features section and mouse over the glyphs in the InterPro domain section.
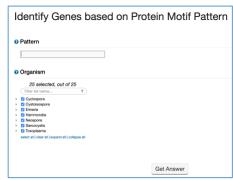


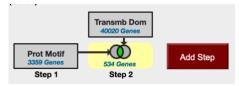   - What other evidence on the gene page supports your conclusion?

2. **Find Cryptosporidium genes with the YXXΦ receptor signal motif.**

   The YXXΦ (Y=tyrosine, X=any amino acid, Φ=bulky hydrophobic [phenylalanine, tyrosine, threonine]) motif is conserved in many eukaryotic membrane proteins that are recognized by adaptor proteins for sorting in the endosomal/lysosomal pathway. This motif is typically located in the c-terminal end of the protein. ***Note: do not look for the Φ symbol on your keyboard – this will not work. Rather you should use the amino acid symbols.*

a. Use the "protein motif pattern" search to find all proteins in ToxoDB that contain this motif anywhere in the terminal 10 amino acids of proteins. (hint: for your regular expression, remember that you want the first amino acid to be a tyrosine, followed any two amino acids, followed by any bulky hydrophobic amino acid (phenylalanine, tyrosine, threonine). Refer to regular expression tutorial if you need to).



b. How many of these proteins also contain at least one transmembrane domain.



c. What would happen if you revise the first step (the motif pattern step) to include genes with the sorting motif in the C-terminal 20 amino acids? (hint: edit the first step and modify your regular expression) https://toxodb.org/toxo/im.do?s=e036d471ad1ccc7c