

## GENE PAGE EXERCISES

### FINDING GENES, BUILDING SEARCH STRATEGIES AND VISITING A GENE PAGE

#### 1. Finding a gene using text search.

For this exercise use <http://www.plasmodb.org>

##### 1a. Find all possible kinases in *Plasmodium*.

Hint: use the keyword “kinase” (without quotations) in the **Gene Text Search** box.



- How many genes did you get?
- Look closely at the sections of the result page. How many of those are in *P. falciparum*? How did you find this out?

Hint – the filter table is located between the strategy panel and the result table and shows the distribution of results across the organisms that you searched. Click on a number to display on that species’ portion of the results.

The screenshot shows the 'My Strategies' panel with a strategy named 'Text' containing 210 genes. Below this is a filter table for 'Plasmodium' species. The table has columns for various species and their respective gene counts. The 'P. falciparum' column is circled in red, showing 210 results. The table also includes a 'All Results' column and a 'Click on a number in this table to limit/filter your results' instruction.

All Results	Ortholog Groups	<i>P. berghei</i>	<i>P. chabaudi</i>	<i>P. cynomolgi</i>	<i>P. falciparum</i>	<i>P. gallinaceum</i>	<i>P. knowlesi</i>	<i>P. reichenowi</i>	<i>P. vivax</i>	<i>P. yoelii</i> (nr Genes: 174)
1734	235	161	162	166	210	187	0	166	0	177
		ANKA	chabaudi	strain B	3D7	IT	8A	strain H	Dennis	Sal-1
										yoelli 17XNL
										yoelli 17X
										yoelli YM

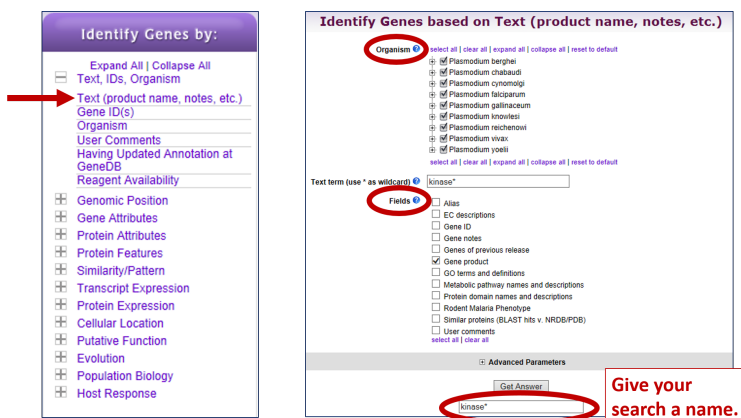
- What happens if you search using the term **kinases** in the Gene Text Search box? How many results are returned?

##### 1b. Find only the kinases that specifically have the word “kinase” in the gene product name.

The search you ran in step 1.1a using the Gene Text Search box initiates a preconfigured search. Initiating the search from the full text search form - **Identify Genes based on Text**,

allows you to configure the search yourself, choosing parameters that best meet your needs. Use the search form to search for genes that have the word kinase in their **gene product** name/description.

- There are several ways to navigate to the **Identify Genes based on Text** page. Notice the sections of the search page. At the top are parameters and the Get Answer button followed by a search description and a list of datasets used by the search.



- How can you make sure to find your text term in plural form or in compound words like “kinases” or “6-phosphofruktokinase”. Adding a wild card, in your search term will broaden your search. A wildcard is an asterisk (\*) and indicates that any character can be in this position. Use the full text search, the specific page where you can define the fields to be searched (Fields = Gene Product).

Try      kinase      \*kinase      \*kinase\*

- **Give each new search a name** to help you keep track of the searches. Notice that in the strategy panel, search names are appended with an asterisk. So the search you save as \*kinase\* will appear as \*kinase\*\*.
- How did you get to the Text Search page?
- How does limiting the number of fields searched affect your results?
- Did you remember to use the wild card?
- How many genes have the word kinase in their product names?

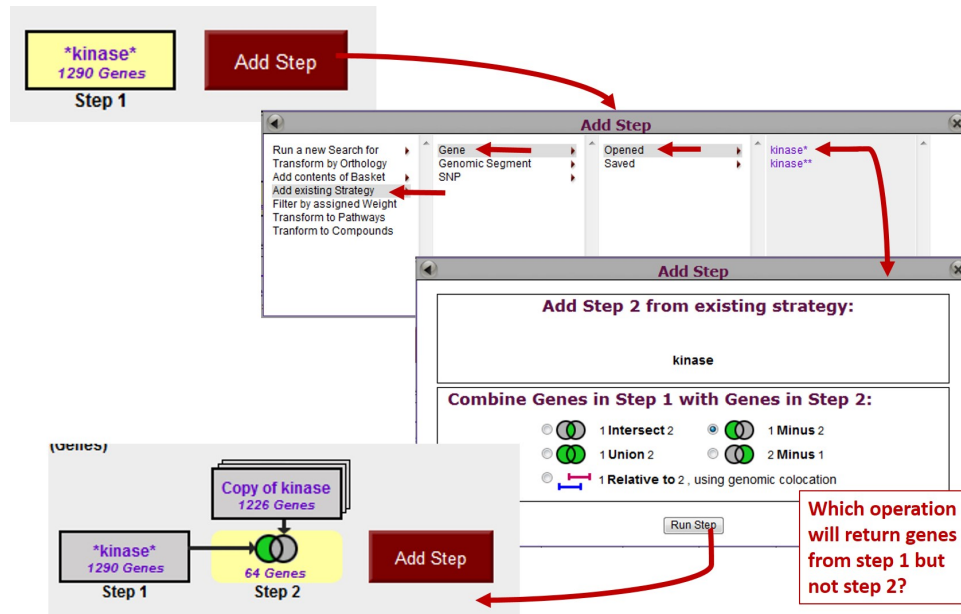
**1c. Combine the results of two text searches.**

**Find genes that were identified using the key word \*kinase\* but not the word kinase?**

- Here we will build a search strategy that combines 2 of your searches. If you are not displaying the results of the **\*kinase\*** search (the strategy box will be highlighted in

yellow), return to it by clicking on that step box in the strategy panel. To add your **kinase** search to this strategy, click on “Add Step” and select “existing strategy”:

- Select the right strategy from your list of Gene Strategies and combine the strategies with the correct operation.
- Do the results make sense? Do all the product names contain the word kinase? From the result page look at the table of gene IDs returned by the search. The Product Description column contains the gene product name.



## 2. Combing text search results with results from other searches

### 2a. Find kinase genes that are likely secreted.

In exercise 1.1b you identified genes that have the word **kinase** somewhere in their gene product name (searching \*kinase\* in gene product field). Grow your search strategy by adding a step that returns genes whose protein products are predicted to have a signal peptide. In this search you are querying the results of our genome-wide analysis that used the SignalP program to predict the presence and location of signal peptide cleavage sites in amino acid sequences.

<http://www.cbs.dtu.dk/services/SignalP/>

Focus your Strategies section on the **\*kinase\*** search and click Add Step. For the second search choose **Identify Genes based on Protein Features, Predicted Signal Peptide**

How did you combine the search results?

How many kinases are predicted to have a signal peptide?

**\*kinase\***  
1290 Genes  
Step 1

Add Step

Add Step

Run a new Search for  
Transform by Orthology  
Add contents of Basket  
Add existing Strategy  
Filter by assigned Weight  
Transform to Pathways  
Transform to Compounds

Genes  
Genomic Segments  
SNPs  
ORFs

Text, IDs, Organism  
Genomic Position  
Gene Attributes  
Protein Attributes  
Protein Features  
Similarity/Pattern  
Transcript Expression  
Protein Expression  
Cellular Location  
Putative Function

Predicted Signal Peptide  
Transmembrane Domain  
Count  
Epitope Presence

Add Step 2 : Predicted Signal Peptide

Organism  
select all | clear all | expand all | collapse all | reset to default

- Plasmodium berghei
- Plasmodium chabaudi
- Plasmodium cynomolgi
- Plasmodium falciparum
- Plasmodium gallinaceum
- Plasmodium knowlesi
- Plasmodium reichenowi
- Plasmodium vivax
- Plasmodium yoelii

Advanced Parameters

Combine Genes in Step 1 with Genes in Step 2:

- 1 Intersect 2
- 1 Minus 2
- 1 Union 2
- 2 Minus 1
- 1 Relative to 2, using genomic colocation

Run Step

Give this search a name

Which operation will return genes that are in both search result sets?

## 2b. Now that you have a list of possible secreted kinases, expand this strategy even further.

There is no wrong answer here!!

- From a biological standpoint what else would be interesting to know about these kinases?  
Add more searches to grow this strategy. Open the categories under Identify Genes By: on the home page and explore the types of searches that are available. You can reduce (or expand) your result set by adding searches that are based on many types of data.
- For example, how many of the secreted kinases also have transmembrane domains?

Signal Pep  
9366 Genes

\*kinase\*  
1290 Genes  
Step 1

75 Genes  
Step 2

Add Step

Expanded View of Step Signal Pep

Signal Pep  
9366 Genes  
Step 1

Add Step

Rename | View | Refresh | Make Nested Strategy | Insert Step Before | Orthologs | Delete

Step 2 : Signal Pep

Organism : Plasmodium berghei, Plasmodium berghei ANKA, Plasmodium chabaudi, Plasmodium chabaudi chabaudi, Plasmodium cynomolgi, Plasmodium cynomolgi strain B, Plasmodium falciparum, Plasmodium falciparum 3D7, Plasmodium falciparum IT, Plasmodium gallinaceum, Plasmodium gallinaceum A, Plasmodium knowlesi, Plasmodium knowlesi strain H, Plasmodium reichenowi, Plasmodium reichenowi Dennis, Plasmodium vivax, Plasmodium vivax Sal-1, Plasmodium yoelii, Plasmodium yoelii yoelii 17XNL, Plasmodium yoelii yoelii YM

Minimum SignalP-NN Conclusion Score : 0.5

Minimum SignalP-NN D-Score : 0.5

Minimum SignalP-HMM Signal Probability : 0.5

or all advanced parameters : any

Results: 9366 Genes

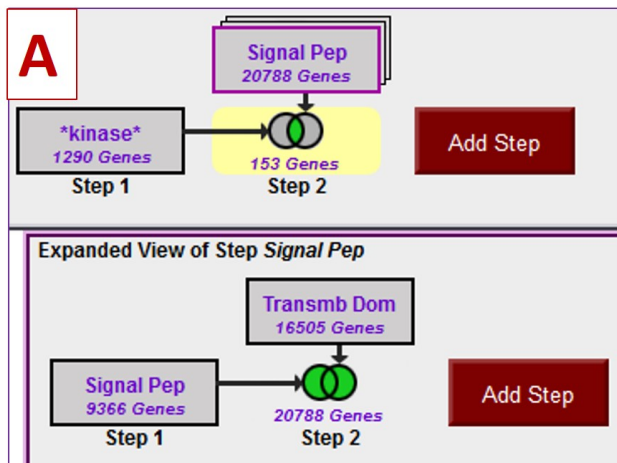
Give this search a weight

2c. In the above example, how can you define kinases that have either a secretory signal peptide AND/OR a transmembrane domain(s)?

Hint: to do this properly you will have to employ the “Nested Strategy” feature. Nesting a strategy allows you to control the order in which your result sets are combined. Think about the difference between two mathematical equations.

Equation without nesting:  $2 \times 3 + 5 = 11$

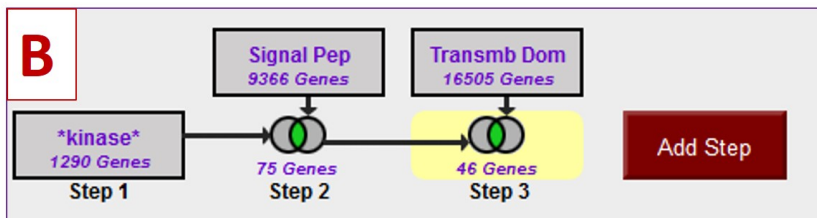
Equation with nesting:  $2 \times (3 + 5) = 16$



**Strategy Logic:**

**Strategy A returns kinases that have a signal peptide OR a TM domain OR both. (SP and/or TM)**

**Strategy B returns kinases that have a signal peptide AND a TM domain**



### 3. Visiting a specific gene record page.

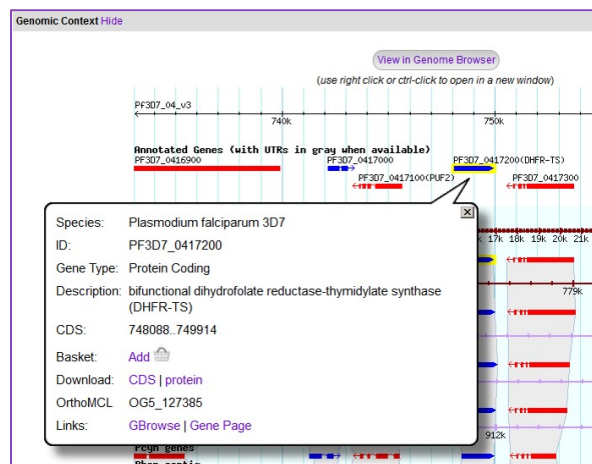
**Note:** For this exercise use <http://www.plasmodb.org>

**3a.** Find the gene page for one of the following *P. falciparum* genes and explore the information there to answer these questions.

1. bifunctional dihydrofolate reductase-thymidylate synthase (DHFR-TS, PF3D7\_0417200)
  2. apical membrane antigen 1 gene (AMA1, PF3D7\_1133400)
- How did you navigate to this gene? What other ways could you get there? I can think of 4 ways to reach the gene page)

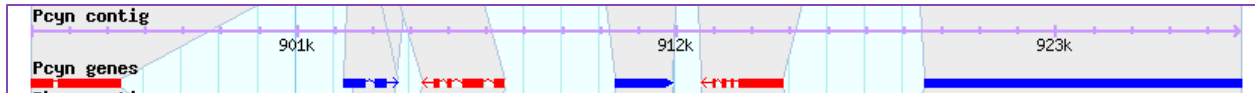
Look at the information in the Genomic Context Graphic.

- What chromosome is this gene on?
- How many exons does this gene have? Hint: look at the graphic in the Genomic Context data track and mouse over the glyph representing the gene.
- What direction is the gene relative to the chromosome?
- How many nucleotides of coding sequence does the gene have?
- Does this gene have a user comment?



**3b.** What genes are located upstream & downstream of DHFR-TS (AMA1) in *P. falciparum*?

- Is synteny (chromosome organization) in this region maintained in other species? Hint: look in the genomic context section of the gene page – what does the shading mean?
- How complete is the genome assembly for other species? Each genome is displayed as two tracks – the genomic sequence (chromosome or contig) on top and the gene models underneath. Do the contigs contain gaps or truncations?



- What does synteny look like across the entire chromosome? To do this:
  - Click on the **“View in GBrowse”** button in the genomic context section.
  - Zoom out to the entire chromosome. There are a few ways to do this – for example, drag your cursor across the entire chromosome then select “zoom” from the popup menu.
  - Click on the tab called “Select tracks”. Select the track called “Syntenic Sequences and Genes (Shaded by Orthology)”. Go back to the Browser tab (this may take a minute to load).
  - Which genome is composed of the most fragments? Are there any other interesting observations you can back by looking at synteny over large genomic regions?

**3c.** Does the *P. falciparum* DHFR-TS (or AMA1) gene contain any single Nucleotide Polymorphisms (SNPs)?

SNPs are represented graphically in the genomic context section and also in a table called “SNP Overview”. Using the SNP Alignment track you can view an alignment showing SNPs between specific strains/isolates.

- Examining the SNP track in the Genomic Context graphic. What do the different color diamonds represent? Mouse over the diamonds to get more information.
- What is the total number of SNPs in the gene?
- How many impact the predicted protein sequence?
- Is this likely to define the full spectrum of sequence variation in these particular strains?
- Compare the SNP characteristics of this gene to upstream and downstream genes. How do these results compare with SNP distribution in other genes?

**3d.** Is the DHFR-TS (or AMA1) gene expressed?

Look at the gene page sections entitled “Protein” and “Expression”. You may have to click on the **show** link to reveal the data associated with that data track.

- What kinds of data in PlasmoDB provide evidence for expression?
- Is this gene expressed at the protein level in salivary gland sporozoites? – in the blood stage phosphoproteome? Look at the Protein context graphic and the table of Mass Spec.-based Expression Evidence.
- How abundant is DHFR-TS (AMA1) protein? How confident are you of this analysis? Abundance can be estimated by counting the number of peptide spectra that map to a protein, or by using the RPKM value from RNA sequencing data.

- Look at the Expression data track labeled Life cycle expression data (3D7). At what life cycle stage is DHFR-TS (AMA1) most abundant? Does this make sense?
- Do the life cycle microarray expression profiles from different data tracks (and thus different experiments/data sets) give the same results? What tracks?
- What about RNA-sequence data, does it agree with microarray data? See these two data tracks – Strand specific transcriptomes of 4 life cycle stages; Transcriptomes of 7 sexual and asexual life stages.