

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>

- a. 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 search results page in PlasmoDB. At the top, there is a 'My Strategies:' section with buttons for 'New', 'Opened (1)', 'All (258)', 'Basket', and 'Public Strategies'. Below this, there is a 'Text' strategy with '223 Genes' and 'Step 1'. A red arrow points to the 'Text' strategy. Below the strategy, there is a table showing the distribution of results across different organisms. The table is titled '223 Genes from Step 1' and 'Strategy: Text(?)'. The table has columns for 'All Results', 'Ortholog Groups', and several Plasmodium species: 'Pberghel ANKA', 'Pchabaudi chabaudi', 'Pcynomol strain', 'Pfalciarum (nr genes: 222)', 'IT', and 'Pgallinaceum 8A'. The 'Pfalciarum (nr genes: 222)' column is circled in red, and the value '223' in the 'All Results' row is also circled in red.

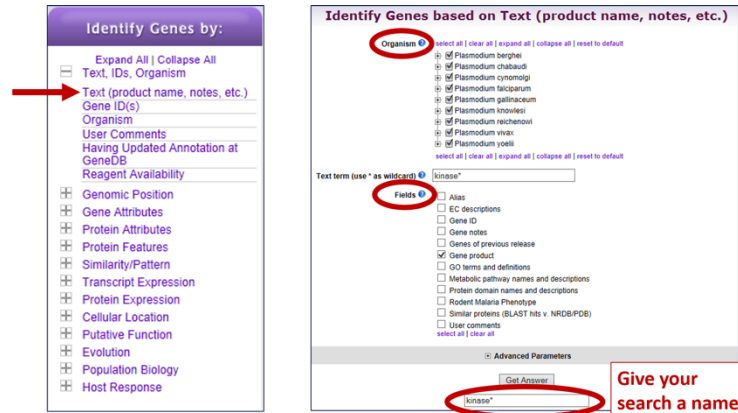
All Results	Ortholog Groups	Plasmodium					
		Pberghel ANKA	Pchabaudi chabaudi	Pcynomol strain	Pfalciarum (nr genes: 222)	Pgallinaceum 8A	
2037	243	173	174	171	223	196	0

- What happens if you search using the term **kinases** in the Gene Text Search box? How many results are returned?
- b. Find only the kinases that specifically have the word “kinase” in the gene product name.

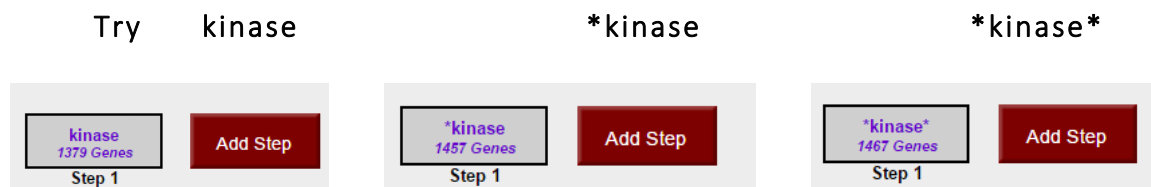
The search you ran in step 1a using the Gene Text Search box initiates a preconfigured search. Initiating the search from the full text search form - **Identify Genes based on**

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 (wildcard = asterisk and means any character) in your search term will broaden your search. Use the full text search, the specific page where you can define the fields to be searched (Fields = Gene Product).



- Give each new search a name to help you keep track of the searches.
- 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?

c. 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. Notice that there is an extra asterisk at the end of an unsaved strategy name. The list of available searches will have an * at the end of the name.

The image shows a search strategy builder interface. On the left, there are two strategy boxes: 'Step 1' with '*kinase*' (1467 Genes) and 'Step 2' with 'Copy of kinase' (1379 Genes). A yellow highlight is under 'Step 1'. A red 'Add Step' button is next to it. Below this, a diagram shows 'Step 1' and 'Step 2' with a Venn diagram and '88 Genes' in a yellow box. A red 'Add Step' button is next to it. On the right, a large 'Add Step' dialog box is open. It has a menu on the left with options like 'Run a new Search for', 'Transform by Orthology', etc. The main area shows a list of strategies: 'Gene', 'Genomic Segment', 'SNP', 'Opened', 'Saved', 'kinase*', and 'kinase*'. A red box points to 'kinase*' and 'kinase*'. Below this, there's a section 'Add Step 2 from existing strategy:' with 'kinase' listed. Underneath, it says 'Combine Genes in Step 1 with Genes in Step 2:' and lists five options: '1 Intersect 2', '1 Union 2', '1 Relative to 2, using genomic colocation', '1 Minus 2', and '2 Minus 1'. A red box points to these options with the text 'Which operation will return genes from step 1 (*kinase*) but not step 2 (kinase)?'. At the bottom of the dialog is a 'Run Step' button.

- 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. Combining text search results with results from other searches

a. Find kinase genes that are likely secreted.

In exercise 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. Navigate the Add Step popup to the Predicted Signal Peptide search. Choose intersect and click Run Step.

- How did you combine the search results?
- How many kinases are predicted to have a signal peptide?

The screenshot shows a workflow in a bioinformatics tool. Step 1 is a search for '*kinase*' resulting in 1467 Genes. An 'Add Step' button is clicked, opening a dialog with various search categories. 'Protein Features' is expanded, and 'Predicted Signal Peptide' is selected. A second 'Add Step' dialog is shown, titled 'Add Step 2 : Predicted Signal Peptide', with a list of organisms and a 'Run Step' button. Below this, a 'Combine Genes in Step 1 with Genes in Step 2' section shows several options: '1 INTERSECT 2' (selected), '1 UNION 2', '1 MINUS 2', '2 MINUS 1', and '1 Relative to 2, using genomic colocation'. A red box with the text 'Which operation will return genes that are in both search result sets?' points to the '1 INTERSECT 2' option. A summary diagram at the bottom left shows Step 1 (*kinase*, 1467 Genes) and Step 2 (Signal Pep, 10714 Genes) with an arrow pointing to a Venn diagram showing the intersection of the two sets, resulting in 91 Genes.

Operator	:	Combined Result will contain:
<input type="radio"/> 1 INTERSECT 2	:	IDs in common between the two lists
<input checked="" type="radio"/> 1 UNION 2	:	IDs from list 1 and list 2
<input type="radio"/> 1 MINUS 2	:	IDs unique to 1
<input type="radio"/> 2 MINUS 1	:	IDs unique to 2
<input type="radio"/> 1 Relative to 2	:	IDs whose features are near each other (colocated) in the genome

b. 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?

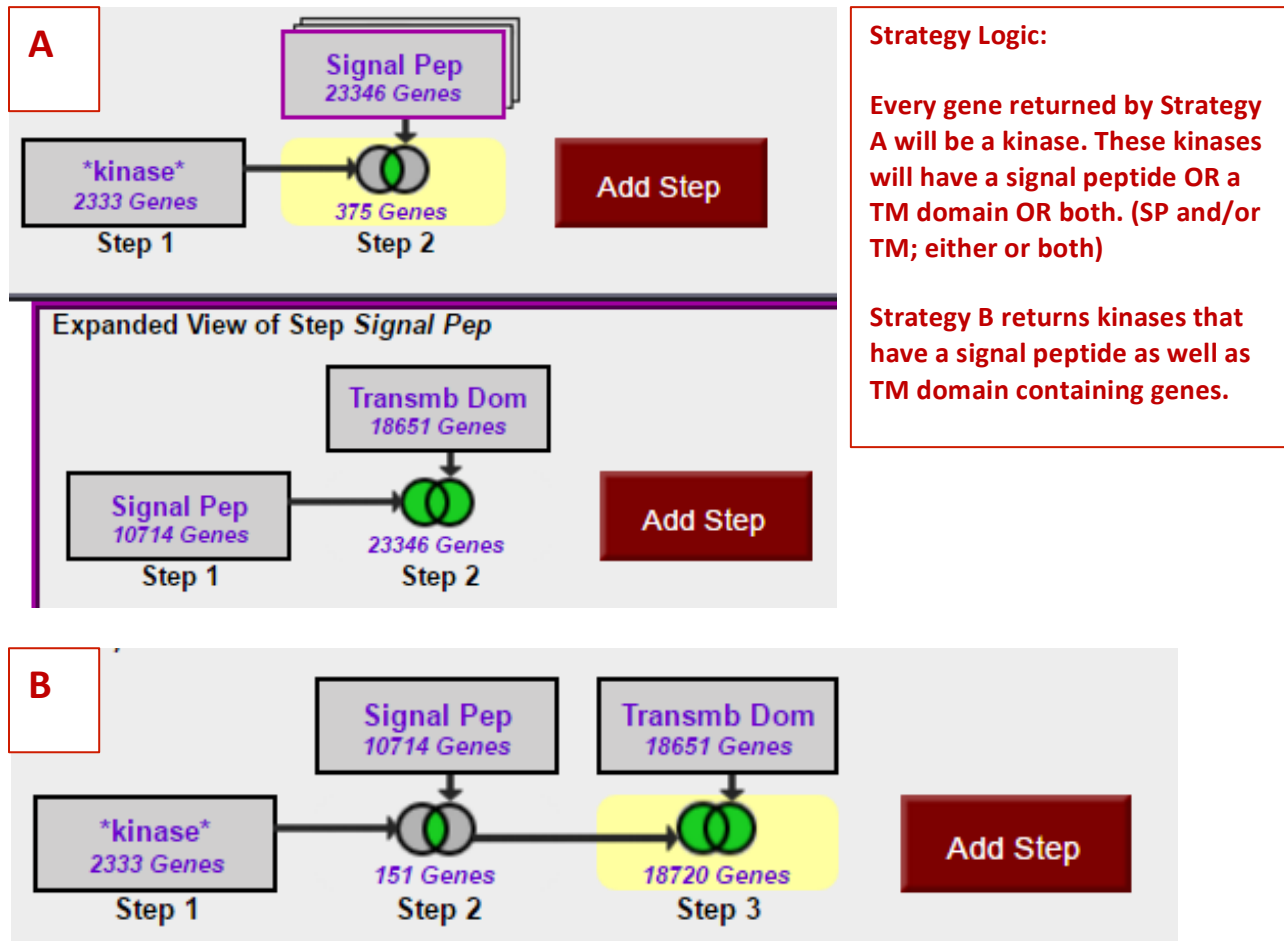
c. In the above example, how can you define kinases that have either a secretory signal peptide AND/OR a transmembrane domain(s)?

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

The image shows a bioinformatics strategy builder interface. It features two main steps: Step 1, labeled '*kinase*' with 1467 Genes, and Step 2, labeled 'Signal Pep' with 10714 Genes. An arrow points from Step 1 to Step 2. Below Step 2 is a red 'Add Step' button. An 'Expanded View of Step Signal Pep' section is visible at the bottom, showing a 'Signal Pep' box with 10714 Genes and another 'Add Step' button. A red circle highlights the 'Make Nested Strategy' option in the top menu of the 'Signal Pep' step's configuration window. The configuration window also displays a list of organisms, minimum scores for SignalP-NN Conclusion, SignalP-NN D-Score, and SignalP-HMM Signal Probability, and a 'Matches any or all advanced parameters' field set to 'any'. The results section shows 'Results: 9366 Genes' and a 'Give this search a weight' dropdown menu.



3. Finding a gene by BLAST Similarity.

Note: For this exercise start with <http://www.toxodb.org>

Imagine that you generated an insertion mutant in *Toxoplasma* that is providing you with some of the most interesting results in your career! You sequence the flanking region and you are only able to get sequence from one side of the insertion (the sequence shown below). You immediately go to ToxoDB to find any information about this sequence. What do you do?

- aaaggagagaaagataaaaatatacaaaaggtccccagagacacgatagtgttactgacaa
catacagaatcaggtcgagcaatggaagaaccaagcaccggcgccagagattgaactcgc
ttggattgccgtagcgttttatgagttgatagcttggctctaaaaaacaaggctgaaaa
atggaaaaaatgtctccaat
- Sequence is also available from this URL:

<http://tinyurl.com/ex1blast>

- Try using the BLAST search with this sequence

The screenshot shows the ToxoDB website interface. At the top, there is a search bar with 'Gene ID: TGME49_239250' and a 'Gene Text Search:' field. Below the search bar is a navigation menu with options like 'Home', 'New Search', 'My Strategies', 'My Basket (0)', 'Tools', 'Data Summary', 'Downloads', and 'Community'. The 'Tools' menu is open, showing a list of options including 'BLAST', 'Results Analysis', 'Sequence Retrieval', 'Pathogen Portal', 'PubMed and Entrez', 'Genome Browser', 'Ancillary Genome Browser', and 'Searches via Web Services'. The 'BLAST' option is highlighted with a red arrow. In the background, there is a sidebar with 'Identify Gen' and a list of categories like 'Text, IDs, Organism', 'Genomic Position', 'Gene Attributes', 'Protein Attributes', 'Protein Features', 'Similarity/Pattern', 'Transcript Expression', 'Protein Expression', 'Cellular Location', 'Putative Function', 'Evolution', and 'Population Biology'. The 'BLAST' option is also highlighted in the sidebar with a red arrow. On the right side, there is a 'Tools:' section with a list of tools including 'BLAST', 'Identify Sequence Similarity', 'Results Analysis', 'Sequence Retrieval', 'Pathogen Portal', 'PubMed and Entrez', 'Genome Browser', and 'Ancillary Genome Browser'. The 'BLAST' option is highlighted in this section with a red arrow. At the bottom right, there is a note: 'For additional tools, use the...'

- Which blast program should you use? (hint: try different Blast programs, just keep in mind that you have a nucleotide sequence so you have to use an appropriate BLAST program).

1. Choose your target data type. What type of sequence in the database do you want to match your sequence to?
2. Choose the BLAST program to use.
3. Choose the target organism. What genome do you want to match your sequence to?

Target Data Type Transcripts
 Proteins
 Genome
 EST
 ORF
 Isolates

BLAST Program blastn
 blastp
 blastx
 tblastn
 tblastx

Target Organism

- Elmeria
- Hammondia
- Neospora
- Sarcocystis
- Toxoplasma
 - Toxoplasma gondii GT1
 - Toxoplasma gondii ME49
 - Toxoplasma gondii VEG

Input Sequence

Note: only one input sequence allowed.
 maximum allowed sequence length is 31K bases.

Expectation value

Maximum descriptions/alignments (V=B)

Low complexity filter

Note on BLAST programs:

- blastn compares a nucleotide sequence against a nucleotide sequence database;
- blastp compares an amino acid sequence against a protein sequence database;
- blastx compares the six-frame conceptual translation products of a nucleotide sequence (both strands) against a protein sequence database;
- tblastn compares a protein sequence against a nucleotide sequence database dynamically translated in all six reading frames (both strands);
- tblastx compares the six-frame translations of a nucleotide sequence against the six-frame translations of a nucleotide sequence database.

- Are you getting any results from blastx? tblastn? What about blastn?
- What is your gene? (hint: after running a blastn against *Toxoplasma* ME49 (Target organism) genomic sequence (Target Data Type), click on the “link to the genome browser”. In the genome browser zoom out to see what gene is in the area).

4. Viewing data on a gene page.

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

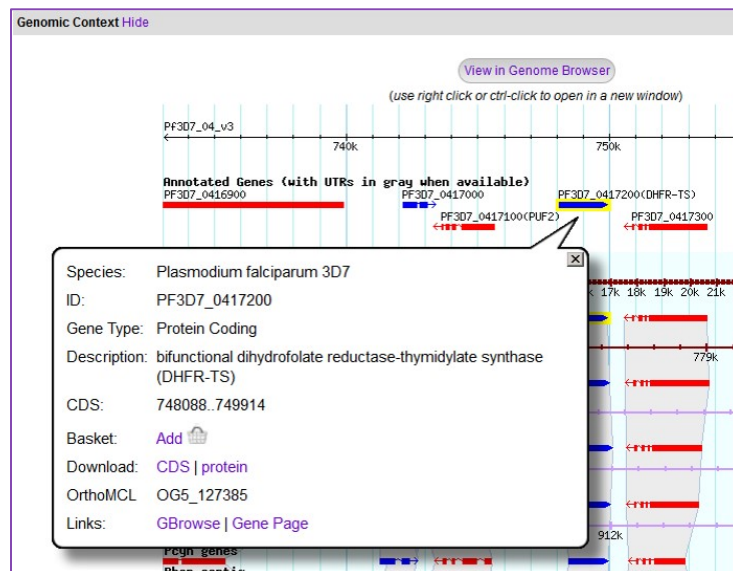
a. Find the gene page for one of the following *P. falciparum* genes and explore the information there to answer these questions.

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 on the gene page.

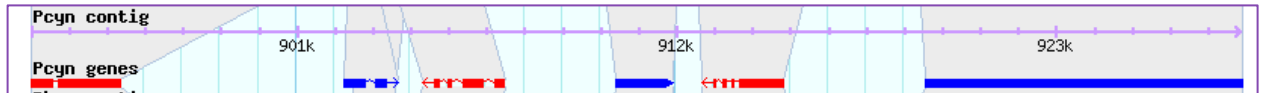
- What chromosome is this gene on?
- How many exons does this gene have? Hint: look at the graphic in the Genomic Context data track and mouse over the glyph representing the gene.
- What direction is the gene relative to the chromosome?
- How many nucleotides of coding sequence?
- Do you see a way to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



b. 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. Does the contig track contain gaps or truncations? What does this imply about the genome assembly?



- What does synteny look like across the entire chromosome? To do this:
 - Click on the “**View in Genome Browser**” 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 in the Overview panel and 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 support by looking at synteny over large genomic regions?
- c. Does the *P. falciparum* DHFR-TS (or AMA1) gene contain Single Nucleotide Polymorphisms (SNPs)?

SNPs are represented in a table called “SNP Overview” and using the “Isolate Alignments in this Gene Region” track you can view an alignment showing SNPs between specific strains/isolates.

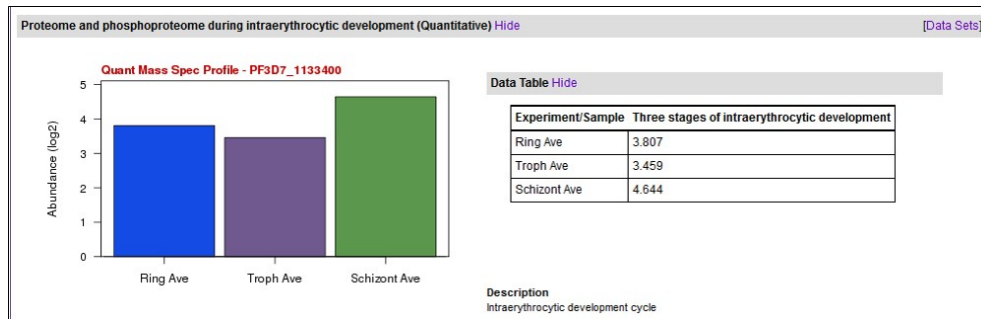
- Examine the SNP Overview table.
- What is the total number of SNPs in the gene?
- How many SNPs impact the predicted protein sequence?
- Is this likely to define the full spectrum of sequence variation in these particular strains?
- Compare the SNP characteristics of this gene to upstream and downstream genes. How do these results compare with SNP distribution in other genes?
- Open the Isolate Alignments in this Gene Region data track and run an alignment between several isolates: 303.1, 383.1, 7G8_2, GB4, N011-A, O222-A, PS097, PS206_E11, RV_3635, RV_3675

d. Is the DHFR-TS (or AMA1) gene expressed at the protein and/or transcript level? 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 protein expression? Hint: open the Protein Features graphic, which is the first data track in the Protein section.
- 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.
- Can you quickly link to the data set record for proteomics experiments?

Mass Spec.-based Expression Evidence [Data Sets]			
Experiment	Sample	Sequences	Spectra View
Blood stage phospho- and total proteome (3D7)	schizont phosphopeptide-depleted	6	16 View
Cytoplasmic and nuclear fractions from rings, trophozoites and schizonts (3D7)	Ring stage nuclear fraction 1	5	5 View
Cytoplasmic and nuclear fractions from rings, trophozoites and schizonts (3D7)	Schizont nuclear fraction 1	3	3 View

- How abundant is DHFR-TS (AMA1) protein? How confident are you of this analysis? Abundance can be estimated by counting the number of spectra supporting a peptide spectra that maps to the protein. Where do you find information about the number of spectra?
- Is the protein more abundant in the ring or schizont life cycle stage? Hint: open the quantitative proteomics track called **Proteome and phosphoproteome during intraerythrocytic development (Quantitative)**.



(PF3D7_1133400 (AMA1))

- Look at the Expression data track labeled **Life cycle expression data (3D7)**. Based on this data, 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 did you use?
- 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.**