# Functional Genomics
## Transcriptomics and Proteomics

1. **Exploring RNA sequence data in *Plasmodium falciparum.***
   **Note: For this exercise use http://www.plasmodb.org**

**1a.** Find all genes in *P. falciparum* that are up-regulated during the later stages of the intraerythrocytic cycle.
   - Hint: Use the fold change search for the data set "Transcriptome during intraerythrocytic development (Bartfai *et al.*)". For this data set, synchronized Pf3D7 parasites were assayed by RNA-seq at 8 time-points during the iRBC cycle. We want to find genes that are up-regulated in the later time points (30, 35, 40 hours) using the early time points (5, 10, 15, 20, 25 hours) as reference.

- Hint: there are a number of parameters to manipulate in this search. As you modify parameters on the left side note the dynamic help on the right side. See screenshots.

- **Direction**: the direction of change in expression. <mark>Choose up-regulated</mark>.

- **Fold Change>=:** the intensity of difference in expression needed before a gene is returned by the search. <mark>Choose 12</mark> but feel free to modify this.

- **Between each gene's AVERAGE expression value:** This parameter sets the operation applied to reference samples. Fold change is calculated as the ratio of two values (expression in reference)/(expression in comparison). When you choose multiple samples to serve as reference, we generate one number for the fold change calculation by using the minimum, maximum, or average. <mark>Choose average</mark>

- **Reference Sample**: the samples that will serve as the reference when comparing expression between samples. <mark>choose 5, 10, 15, 20, 25</mark>

- **And it's AVERAGE expression value:** This is the operation applied to comparison samples. see explanation above. <mark>Choose average</mark>

- **Comparison Sample**: the sample that you are comparing to the reference. In this case you are interested in genes that are up-regulated in later time points <mark>choose 30, 35, 40</mark>

**1b.** For the genes returned by the search, how does the RNA-sequence data compare to microarray data?

- Hint: PlasmoDB contains data from a similar experiment that was analyzed by microarray instead of RNA sequencing. This experiment is called: Erythrocytic expression time series (3D7, DD2, HB3) (Bozdech et al. and Linas et al.). To directly compare the data for genes returned by the RNA seq search that you just ran, add the column called "Pf-iRBC 48hr - Graph".



**1c.** How many genes in this result have 16 exons?

- Hint: add a column for number of exons. To help you find the genes with 16 exons, you can sort the columns using the arrows that precede the column heading. Also, clicking the histogram icon in the column heading will provides options for viewing the column data as a table or histogram.

- There are three gene IDs with 16 exons each. Two have similar gene IDs. What does this mean?

**1d.** Click on one of the two similar gene IDs from above. Look at the gene page. Take note of the Gene ID. Mouse over the gene models in the genomic context view and explore the popup. What information does it contain? Note that the CDS section includes exon coordinates. Compare the coordinates for the two alternative splice variants of this gene - can you identify the difference (it is very subtle)?



**1e.** View this gene in the genome browser and load the RNA-seq tracks for this experiment. The track is named: "Transcriptome during intraerythrocytic development mRNAseq Coverage aligned to P falciparum 3D7 (Bartfai et al.) (log plot).

- Do these tracks match the differential expression results you got above? Is this gene differentially regulated between the early time points and the late ones?

- Do you agree with the alternative splice call? Are there other possible splice variants? (*Hint*: turn on the track called "Splice Site Junctions (Union of All Experiments)").

**1f.** What other data type can you load to help in looking at gene structure? (*Hint*: Look in the transcript expression section of the gbrowse tracks... how about ESTs?).

**1g.** You decide that you would like to present this data at your lab meeting and possibly use the data as supporting evidence in a grant proposal. The data you want to present is much like the graphs on individual gene pages but you would like to make a composite graph showing the profiles of several genes from your list.

- Download to your computer the list of genes returned by the search. Include the following information in your downloaded file: Gene ID, Product Description, Genomic Location, Annotated 5' UTR length, Annotated 3' UTR length.
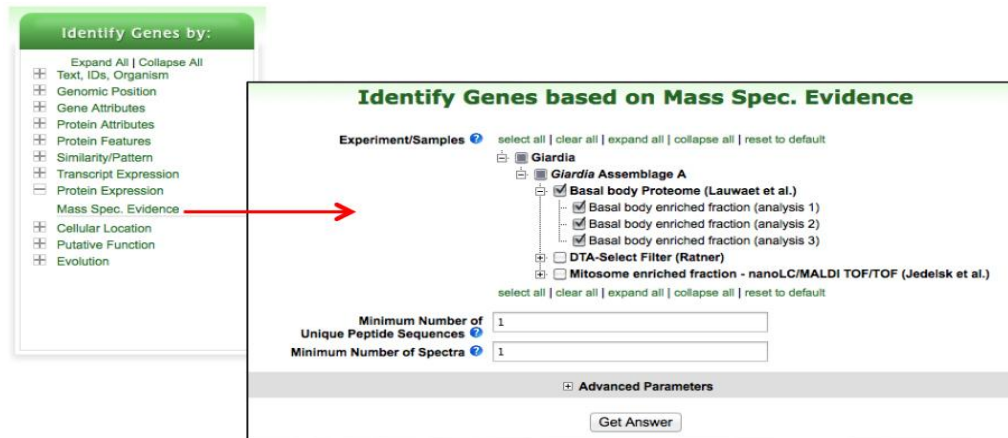


- Download to your computer the data file containing expression values (FPKM) for the RNA the sense strand uniquely aligned reads.

2. **Find all Giardia genes with evidence of basal body expression based on mass spec/proteomics data.**
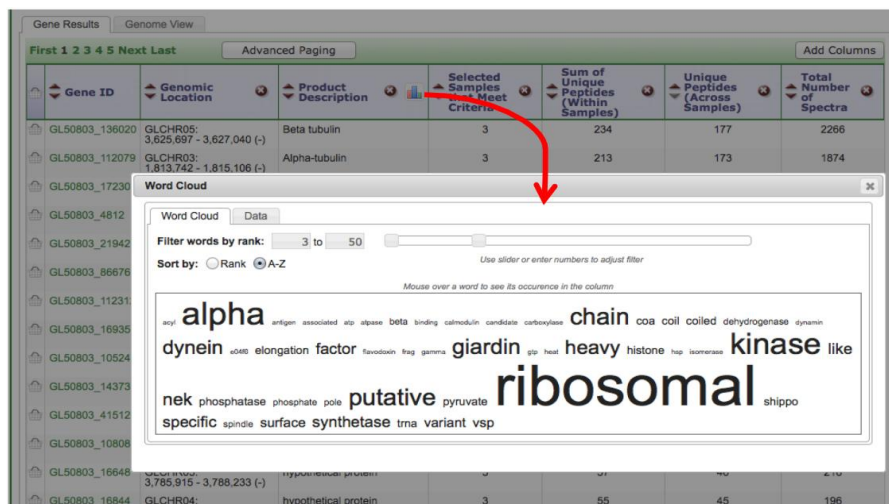   **For this exercise use http://giardiadb.org**

   The proteomics search is available under the heading "Protein Expression" in the "Identify gene by" section.



2a.    How many genes did you identify?

2b.    What kinds of genes are in your result list?
       Hint 1: analyze the Product Description column. The **word cloud tool** counts the number of times a word appears in the column and then draws a word cloud in which the size of the word reflects how many times the word appears in the product description column. Click on the little graphic icon next to the column called "Product description".

Hint 2: apply a statistical analysis to the Product Description column. The Word Enrichment Analysis Tool (Click the Blue Analyze Results button) considers the words in the Product Description column and applies the Fischer's Exact test to compare your gene result to the product descriptions of the entire genome.



Hint 3: analyze the Gene Ontology terms assigned to the genes in your result list. **The Result Analysis/Enrichment tool** applies the Fischer's Exact test to compare your gene result to the entire genome. Use the Gene Ontology Enrichment to find Biological Process ontology terms that are enriched in your gene result.

**2c.** Examine the filter table below the strategy section. Why is the number in the ortholog groups column less than the number in the Assemblage A column? Why do the other cells have a zero?



**2d.** How many genes with mass spec data from the basal proteome experiment also have mass spec evidence from the "Mitosome enriched proteome (WB) (Jedelsk et al.)"?

**2e.** The default parameters of the mass spec search are set to identify any gene with at least 1 mapped peptide identified. How will your results change if you revise the two steps in your search strategy to only return genes with at least 5 peptides identified?