# Phenotypic data

**Learning objectives:**
- Explore how to combine different phenotypic data
- Explore high throughput mutagenesis data
- Explore curated phenotypic data
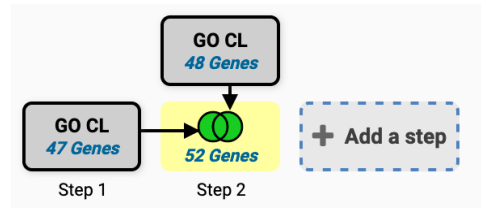- Explore high throughput subcellular localization data

1. **Identify genes that are targeted to the ciliary tip of *Trypanosoma brucei* that are also essential for parasite fitness.** Note for this exercise use http://tritrypdb.org

   TriTrypDB integrates data from the TrypTag project (http://tryptag.org). Genes from *T. brucei* were N- and C-terminally tagged with a fluorescent protein and subcellular localization determined by microscopy. The description of the localization was done using gene ontology terms.

   a. Start by finding the "Cellular Localization Imaging" search.



   b. Configure the search to identify the GO term "Ciliary Tip" – notice that when you start typing the autocomplete function offers you selectable options.  Since the experiment examined both N and C terimnal fusions proteins, you will have to run the search twice and combine the results from both searches.  Did you use a union or an intersect to combine the results?

c. Explore the results you got. Scroll down to the results section, then scroll to the right of the results window to reveal the subcellular localization images. These are very small, but you can right click on them to open a larger image in a new window. If you do not see the images, you may need to add the data column. Click Add Columns and choose the Cellular localization images column





d. Add a step to identify how many genes are essential for the fitness of the parasite. Click on Add step, then search for the phenotype searches. Click on the Phenotype Evidence option. Select the "High-throughput phenotyping using RNAi target sequencing (David Horn).

e. Configure the search to return genes that are decreased in coverage by 1.5 fold when comparing the maximum expression value of all induced samples to the uninduced sample. How many genes did you get?

2. **Finding genes based on high throughput mutagenesis and fitness analysis.** Note for this exercise use http://toxodb.org

   a. Navigate to the CRISPR phenotype search. Note that this search form is quite simple just requiring a range of fitness values. The defaults return all genes without limiting the search at all. This returns all genes that were assayed, which is nearly the entire genome. The tricky bit is deciding where to make the cutoffs. The description on the search form is very helpful in this regard (as is the link to the paper … These phenotypes were assayed under specific conditions. If a particular gene doesn't show a phenotype, it might show a phenotype in other conditions (or infecting an actual host).
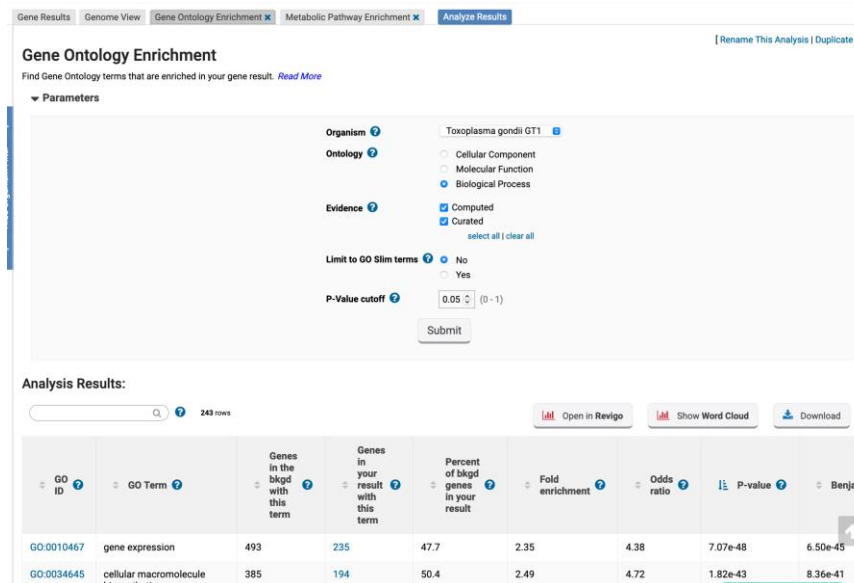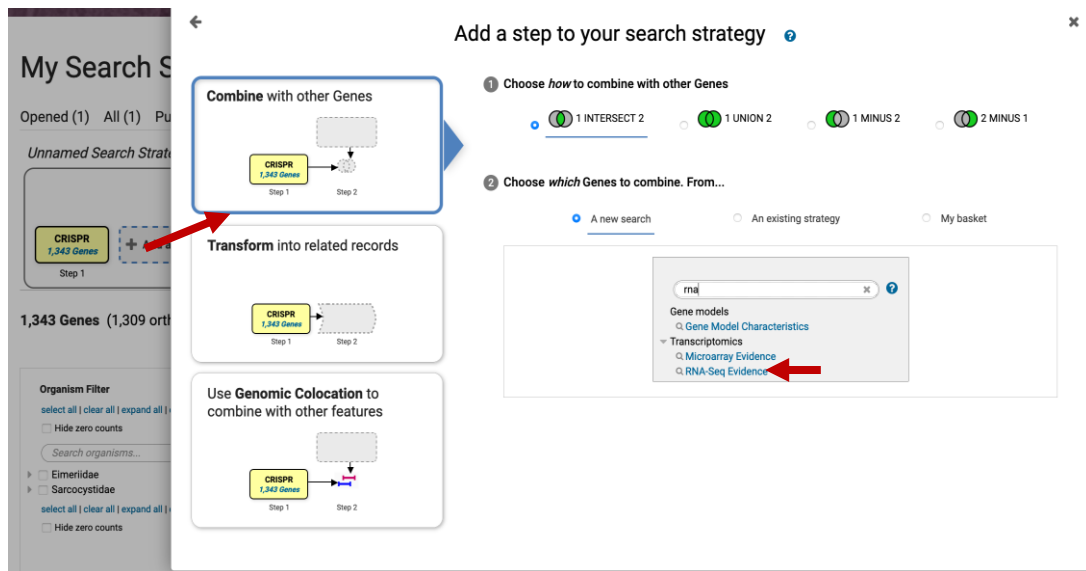


   b. The plot showing the phenotype score (fitness) is particularly useful. Red points along the plot are genes known to be essential under these conditions, while yellow are known to be expendable. This graph will help you determine where to set the values. The scores range from 2.96 (least "essential") to -6.89 (most "essential").
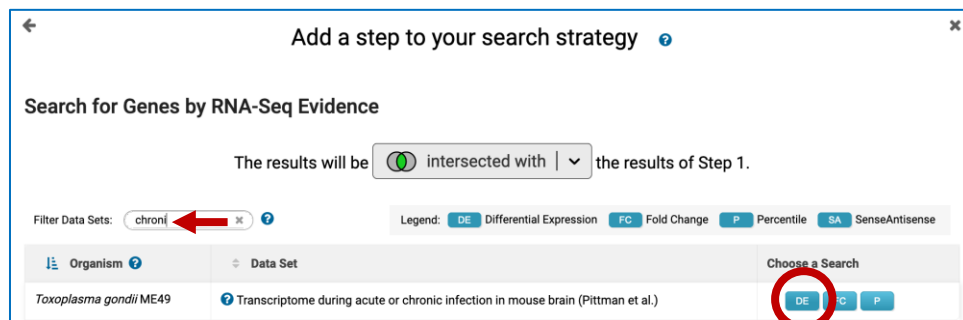
c. Try running this search by limiting the range from -6.89 to -4. Do you get the expected results based on the above graph and the number of genes returned in your search results?

- What kinds of genes are in your results? What kinds of genes would you expect to be essential? One way to explore the data is to run a GO enrichment analysis to determine if any biological processes are enriched in your results. Give this a try. What do you results look like and do they make sense?



d. How many of these genes are upregulated in *in vivo* chronic stages of *T. gondii*? Click on Add Step and elect the RNAseq searches under the Transcriptomics category.

e. Find the experiment with chronic stages and run a search based on differentially expressed genes (DE). Intersect genes that are 2-fold upregulated in chronic stages compared to acute stages

- Did you get zero results? This is to be expected since the CRISPR data was analyzed using the GT1 strain of *Toxoplasma* and the RNA-Seq data is from the ME49 strain. How can you fix this?
- Hint: transform the results in step 2 from *T. gondii* ME49 to *T. gondii GT1*. Click on the step edit button (move your mouse over the step and select edit).



- Select **orthologs** from the menu items at the top of the pop window

- Select *T. gondii* GT1 from the list of organisms and click on Run Step.



- Now what do your results look like?

3. **Identify essential *Plasmodium falciparum* genes that are highly expressed in schizont stages of the parasite.** Note for this exercise use https://plasmodb.org

   a. You can start by exploring the phenotype data in PlasmoDB. Select and run the search associated with the dataset: piggyBac insertion mutagenesis (John Adams).



   b. Configure the search to identify genes with a *mutant fitness score* of less that -3. (This example shows -4 to -3.089) Note that you can select the range by either clicking and dragging you mouse over the histogram or by typing the values in the selection boxes.

c. How many genes did you identify? Which gene has the lowest fitness score? Note that you might need to add the fitness score column, by clicking on add columns then filtering the options with the word "fitness".



d. Click on Add Step and find the RNA-Seq searches.

e. Find the search called "Intraerythrocytic development cycle transcriptome (2019) (Wichers et al. 2019)" and select the percentile search.



f. Configure the search to identify all genes that are in the 80-100 percentile in all three available schizont samples. Remember to change the parameter to **require matching all samples**.

g. How many genes did you get? Are any of these genes interesting? How many are predicted to be secreted?



h. How did you identify the secreted genes? Hint, add a step and search for genes that have a predicted secretory signal peptide.



4. **Identify Neurospora crassa genes that affect conidia formation. Note for the exercise use https://fungidb.org**
   - Start by locating the phenotype searches.

- This search provides you the option to filter based on categories on the left. Notice how when you select a different category on the left the filtering options in the middle change. Select the **Conidia number** category. Next select the "Reduced" value.
- Notice that this search allows you to explore your results even before you click on the "Get Answer" button! Click around on the other categories on the left and see if the genes that are involved in a reduced number of conidia may also be involved in other phenotypes. For example, click on the **Ascospore Number** category, how maybe of your genes also have a phenotype with no ascospore formation?

- Click on get answer. What kinds of genes are in your results? Try analysing the results to see if there are any biological processes enriched in your results.