# 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

   a. 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.

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

   ### Identify Genes based on Cellular Localization Imaging

   

   - 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?

   

- 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.
b. 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)".

- 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.



For the **Experiment**
☑ Quantitated from the CDS Sequence
☐ Quantitated from gene model (5 prime UTR + CDS)
select all | clear all

❓

return   protein coding  ⊟ ❓ **Genes**
that are  Decrease in coverage  ⊟ ❓
with a **Fold change** >=  1.5  ❓ ⬅
    between each gene's   maximum  ⊟ ❓ **expression value**
    in the following  **Reference Samples**  ❓
    ☑ Uninduced sample  ⬅

    select all | clear all

and if   maximum  ⊟ ❓  **expression value**
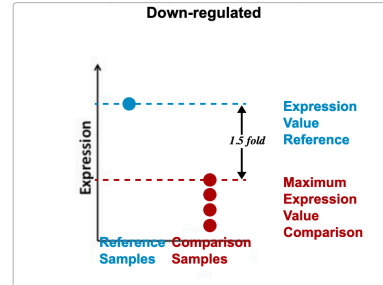in the following  **Comparison Samples**  ❓
☑ Induced in bloodstream (BS) forms, 3 days (10 doublings)
☑ Induced in bloodstream (BS) forms, 6 days (20 doublings)  ⬅
☑ Induced in procyclic forms (PS) forms, 9 days (9 doublings)
☑ Induced throughout differentiation (DIF = 7 BS doublings + 6 PS doublings)

select all | clear all

**Example showing one gene that would meet search criteria**
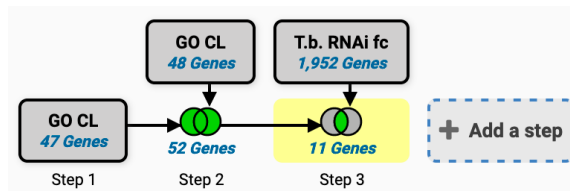(Dots represent this gene's expression values for selected samples)

**Down-regulated**

Expression

*1.5 fold*   Expression Value Reference

Maximum Expression Value Comparison

Reference Samples    Comparison Samples

For each gene, the search calculates:

$$fold\ change = \frac{reference\ \text{expression value}}{maximum\ \text{expression value in } comparison}$$

and returns genes when **fold change >= 1.5**.

You are searching for genes that are **down-regulated** between one **reference sample** and at least two **comparison samples**.

This calculation creates the **narrowest** window of expression values in which to look for genes that meet your fold change cutoff. To broaden the window, use the average or minimum comparison value.
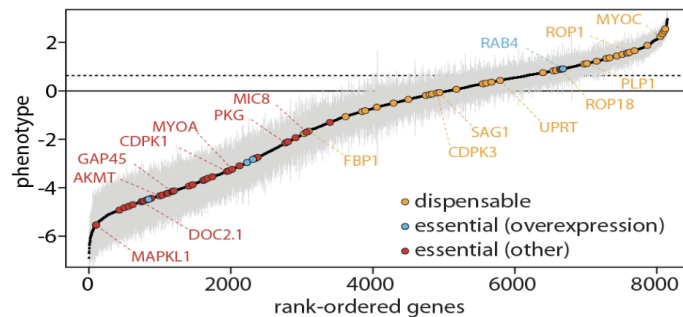
- How many genes did you get?



GO CL
48 Genes

T.b. RNAi fc
1,952 Genes

GO CL
47 Genes

52 Genes

11 Genes

➕ Add a step

Step 1          Step 2          Step 3

2. **Finding genes based on high throughput mutagenesis and fitness analysis.**
   Note for this exercise use http://toxodb.org
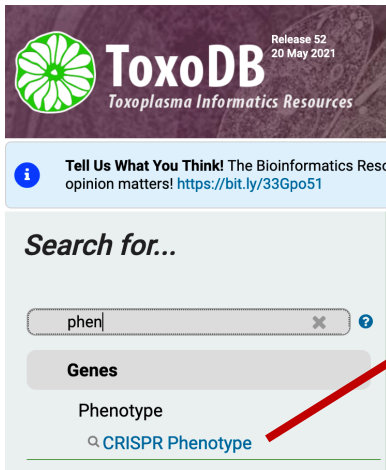
   - Navigate to the CRISPR phenotype search.  Note that this search form is quite simple just requiring a range of fitness values.  The defaults return all genes not limiting the search at all.  This is only useful in as much as it tells you which genes were assayed which is nearly the entire genome.  The tricky bit is deciding where to make the cutoffs.  Again, the description on the search form is very helpful in this regard (as is the link to the paper … remember these phenotypes were assayed under specific conditions so just because a particular gene doesn't show a phenotype doesn't mean it wouldn't in other conditions (or infecting an

actual host).  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 will help you determine where to set the values.  The scores range from 2.96 (least "essential) to -6.89 (most "essential). Try it 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?

- 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



- 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.

← **Add a step to your search strategy** ⊘

**❷ Experiment**

  ⦿ Acute and chronic T.gondii infection of mouse. unstranded

**❷ Reference Sample**

  ⦿ acute infection 10 days p.i.
  ○ chronic infection 28 days p.i.

**❷ Comparator Sample**

  ○ acute infection 10 days p.i.
  ⦿ chronic infection 28 days p.i.

**❷ Direction**

  [ up-regulated    ⬍ ]

**❷ fold difference >=**

  [ 2                    ]

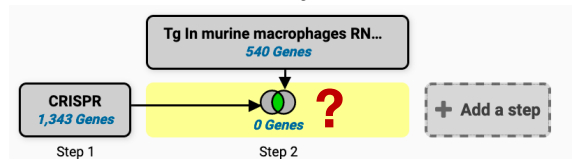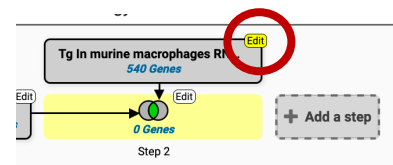**❷ adjusted P value less than or equal to**

  [ 0.1                  ]

- 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.

View | Analyze | **Revise** | **Make nested strategy** | **Insert step before** | **Orthologs** | Delete ✖

**Details for step** *Tg In murine macrophages RNA-Seq (de)* ✏
540 Genes

| | |
|---|---|
| **Experiment** | Acute and chronic T.gondii infection of mouse. unstranded |
| **Reference Sample** | acute infection 10 days p.i. |
| **Comparator Sample** | chronic infection 28 days p.i. |
| **Direction** | up-regulated |
| **fold difference >=** | 2 |
| **adjusted P value less than or equal to** | 0.1 |

▶ **Give this search a weight**

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

**❷ Organism**

*1 selected, out of 31*

add these | clear these | select only these
select all | clear all
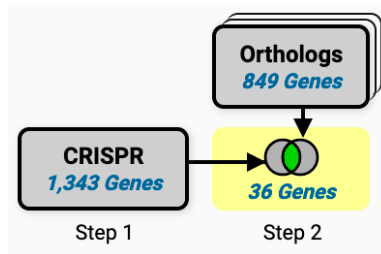
gt1 ← ✖ ❷

⊟ Sarcocystidae
⊟ Toxoplasma
← ☑ Toxoplasma gondii GT1
add these | clear these | select only these
select all | clear all

**❷ Syntenic Orthologs Only?**

no ◉

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



   - Configure the search to identify genes with a *mutant fitness score* of less that -3. 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.

- 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".



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

- Find the search called "Intraerythrocytic development cycle transcriptome (2019)" and select the percentile search.

**Search for Genes by RNA-Seq Evidence**

The results will be [⬤ intersected with ▾] the results of Step 2.

Filter Data Sets: [intraer] →    ❓    Legend: [DE] Differential Expression [FC] Fold Change [P] Percentile [SA] SenseAntisense

| ↓≡ Organism ❓ | ⇕ Data Set | Choose a Search |
|---|---|---|
| *Plasmodium falciparum* 3D7 | ❓ Intraerythrocytic development cycle transcriptome (2019) (Wichers et al. 2019) | [DE] [FC] [P] |
| *Plasmodium falciparum* 3D7 | ❓ Intraerythrocytic development cycle transcriptome (2018) (Toenhake et al.) | [FC] [P] [SA] |
| *Plasmodium falciparum* 3D7 | ❓ Transcriptome during intraerythrocytic development (Bartfai et al.) | [FC] [P] |
| *Plasmodium falciparum* 3D7 | ❓ Blood stage transcriptome (3D7) (Otto et al.) | [FC] [P] |
| *Plasmodium falciparum* 3D7 | ❓ Intraerythrocytic cycle transcriptome (3D7) (Hoeijmakers et al.) | [FC] [P] [SA] |
| *Plasmodium falciparum* 3D7 | ❓ Strand specific transcriptome of the intraerythrocytic developmental cycle (Siegel et al.) | [FC] [P] [SA] |
| *Plasmodium vivax* P01 | ❓ Transcription profile of intraerythrocytic cycle (Zhu et al.) | [FC] [P] |

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

**❓ Samples**

☐ young ring 8 hpi
☐ late ring_early trophozoite 16 hpi
☐ mid trophozoite 24 hpi
☐ late trophozoite 32 hpi
☑ early schizont 40 hpi
☑ schizont 44 hpi  ←
☑ late schizont 48 hpi
☐ purified merozoites 0 hpi
select all | clear all

**❓ Minimum expression percentile**

[ 80 ]  ←

**❓ Maximum expression percentile**

[ 100 ]

**❓ Matches Any or All Selected Samples?**

[ all ⬍ ]  ←

**Pf3D7 IDC RNASeq (%ile)**
**789 Genes**

**pB MIS/MFS**
**856 Genes**
→ **122 Genes**   ➕ **Add a step**

Step 1        Step 2

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

**Pf3D7 IDC RNASeq (%ile)**
**789 Genes**

**Signal Pep**
**603 Genes**

**pB MIS/MFS**
**856 Genes**
→ **122 Genes** → **17 Genes**

Step 1        Step 2        Step 3

**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?

**1,283 Genes Total**

expand all | collapse all

[ Find a variable ] 🔍 ❓

| | |
|---|---|
| 📊 Aerial Hyphae Height | |
| ☰ Ascospore Morphology | |
| ☰ **Ascospore Number** | ⬅ |
| 📊 Basal Hyphae Growth Rate | |
| ☰ Conidia Morphology | |
| ☰ Conidia Number | |
| ☰ Perithecia Morphology | |
| ☰ Perithecia Number | |
| ☰ Protoperithecia Number | |
| ☰ Protoperithecial Morphology | |

**99 of 1,283 Genes selected**  [ Conidia Number ✕ ]

**Ascospore Number**

Check items below to apply this filter      **1,283 (100%) of 1,283** Genes have data for this variable

| ☐ | 📊 Ascospore Number | ⬍ | Remaining Genes ❓ | | ⬍ | Genes ❓ | | Distribution ❓ | % ❓ |
|---|---|---|---|---|---|---|---|---|---|
| | | | 99 | (100%) | | 1,283 | (100%) | | |
| ☐ | Normal | | 32 | (32%) | | 1,043 | (81%) | ▬▬ | (3%) |
| ☐ | Not formed ⬅ | | 56 | (57%) | | 169 | (13%) | ▬ | (33%) |
| ☐ | Reduced | | 11 | (11%) | | 65 | (5%) | ▬ | (17%) |
| ☐ | Increased | | 0 | (0%) | | 2 | (< 1%) | ▏ | (0%) |
| ☐ | Severely Reduced | | 0 | (0%) | | 5 | (< 1%) | ▏ | (0%) |
| ☐ | Severely reduced | | 0 | (0%) | | 1 | (< 1%) | ▏ | (0%) |

- 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.

**KO Mut**
**99 Genes**
Step 1    ➕ Add a step

**99 Genes** (98 ortholog groups)  [ Revise this search ]

Gene Results | Genome View | Gene Ontology Enrichment ✕ | **Analyze Results**

**Organism Filter**
select all | clear all | expand all | collapse all
☐ Hide zero counts
[ Search organisms... ] 🔍 ❓

| | |
|---|---|
| ▸ ☐ Fungi | 99 |
| ▸ ☐ Oomycota | 0 |

select all | clear all | expand all | collapse all
☐ Hide zero counts

**Gene Ontology Enrichment**                    [ Rename This Analysis | Duplicate ]

Find Gene Ontology terms that are enriched in your gene result. *Read More*

▾ **Parameters**

| | |
|---|---|
| Organism ❓ | Neurospora crassa OR74A 🔽 |
| Ontology ❓ | ⦿ Biological Process  ⚪ Cellular Component  ⚪ Molecular Function |
| Evidence ❓ | ☑ Computed  ☑ Curated  select all \| clear all |
| Limit to GO Slim terms ❓ | ⚪ No  ⚪ Yes |
| P-Value cutoff ❓ | [ 0.05 ⬍ ] (0 - 1) |

[ Submit ]

**Analysis Results:**

[ 🔍 ] ❓  **361** rows      [ 📊 Open in **Revigo** ]  [ 📊 Show **Word Cloud** ]  [ ⬇ Download ]

| ⬍ GO ID ❓ | ⬍ GO Term ❓ | ⬍ Genes in the bkgd with this term ❓ | ⬍ Genes in your result with this term ❓ | ⬍ Percent of bkgd genes in your result ❓ | ⬍ Fold enrichment ❓ | ⬍ Odds ratio ❓ | ⬍ P-value ❓ | ⬍ B |
|---|---|---|---|---|---|---|---|---|
| GO:0070787 | conidiophore development | 84 | 26 | 31.0 | 22.87 | 44.43 | 1.32e-29 | 1.28e-2 |
| GO:0032501 | multicellular organismal process | 194 | 33 | 17.0 | 12.57 | 22.24 | 2.22e-28 | 1.08e-2 |
| GO:0061458 | reproductive system development | 184 | 32 | 17.4 | 12.85 | 22.51 | 8.32e-28 | 1.61e-2 |
| GO:0048608 | reproductive structure development | 184 | 32 | 17.4 | 12.85 | 22.51 | 8.32e-28 | 1.61e-2 |
| GO:0075259 | spore-bearing structure development | 184 | 32 | 17.4 | 12.85 | 22.51 | 8.32e-28 | 1.61e-2 |
| GO:0048731 | system development | 185 | 32 | 17.3 | 12.78 | 22.36 | 9.97e-28 | 1.61e-2 |
| GO:0007275 | multicellular organism development | 187 | 32 | 17.1 | 12.64 | 22.07 | 1.43e-27 | 1.98e-2 |