# Exploring Transcriptomic data

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

a. 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.



   - 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. Choose up-regulated.

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

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

- **Between each gene's AVERAGE expression value:** This parameter appears once you have chosen two Reference Samples and defines the operation applied to reference samples. Fold change is calculated as the ratio of two values (upregulated ratio = expression in comparison)/(expression in reference). 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>

- **(or a Floor of 10 reads):** This parameter defines a lower limit of aligned reads for a gene to avoid unreliable fold change calculations. (Low numbers of aligned reads means low expression but the low values may be may be technically inaccurate.  Dividing by small numbers creates large numbers.  2000FPKM/10 = 200; 2000/0.1 =  20,000) If a gene has fewer than 10 aligned reads, it is assigned 10 reads before the fold change calculation is made.  <mark>Leave this as default at 10 reads.</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>

- **And its AVERAGE expression value:** This parameter appears once you have chosen two Comparison Samples and defines the operation applied to comparison samples.  See explanation above. <mark>Choose average</mark>

**b.** 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.). IDC 48 hr Marray – Expr Graph** shows normalized expression values. 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".



OPTIONAL: You can also run a fold change search using this experiment to compare results on a genome scale. Add a step to your strategy and intersect your current results (genes upregulated 12 fold in later IDC time periods) with a fold change search using the "Erythrocytic expression time series (3D7, Dd2, HB3) (Bozdech et al. and Linas et al.)" experiment (under microarray evidence). Configure it similarly to the RNA-seq experiment although you will probably need to make the fold change smaller (try 2 or 3) due to the decreased dynamic range of microarrays compared to RNA-seq.

**Add Step 2 : P.falciparum Erythrocytic expression time series (3D7, DD2, HB3) Microarray (fold change)**

For the **Experiment** iRBC HB3 (48 Hour scaled) ▼

return protein coding ▼ **Genes**

that are up-regulated ▼

with a **Fold change** >= 2

between each gene's average ▼ **expression value**

in the following **Reference Samples**

*28 selected, out of 46*

Filter list below...

▸ ☑ 1-16 Hours
▸ ☑ 17-30 Hours
▸ ☐ 31-48 Hours

select all | clear all | expand all | collapse all

and its average ▼ **expression value**

in the following **Comparison Samples**

*18 selected, out of 46*

Filter list below...

▸ ☐ 1-16 Hours
▸ ☐ 17-30 Hours
▾ ☑ 31-48 Hours
  ▸ ☑ 31-39 Hours
  ▸ ☑ 40-48 Hours

select all | clear all | expand all | collapse all

**Example showing one gene that would meet search criteria**

(Dots represent this gene's expression values for selected samples)

**Up-regulated**

Expression

2 fold

**Average Expression Level Comparison**

**Average Expression Level Reference**

Reference Samples    Comparison Samples

*A maximum of four samples are shown when more than four are selected.*

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

For each gene, the search calculates:

$$\text{fold change} = \frac{\textit{average expression level in comparison}}{\textit{average expression level in reference}}$$
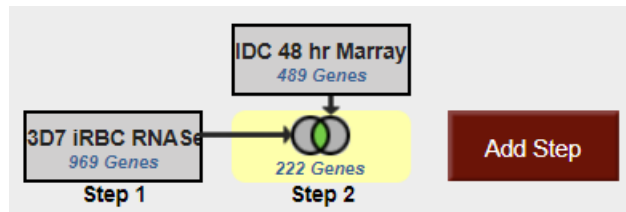
and returns genes when **fold change** >= 2.

To narrow the window, use the maximum reference value, or minimum comparison value. To broaden the window, use the minimum reference value, or maximum comparison value.

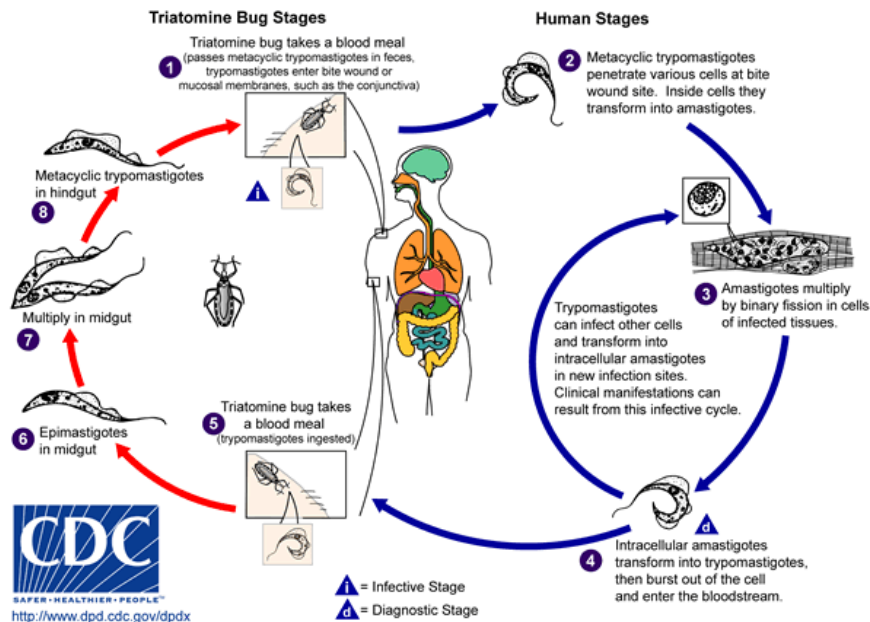See the detailed help for this search.

**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



IDC 48 hr Marray
*489 Genes*

3D7 iRBC RNASe
*969 Genes*

*222 Genes*

**Add Step**

**Step 1**     **Step 2**

2. *Optional (come back if time).* Exploring microarray data in TriTrypDB.
   Note:  For this exercise use http://www.tritrypdb.org



a. Find *T. cruzi* protein coding genes that are upregulated in amastigotes compared to trypomastigotes. Go to the transcript expression section then select **microarray**. Choose the fold change (FC) search for the data set called: **Transcriptomes of Four Life-Cycle Stages (Minning et al.)**.

- Select the direction of regulation, your reference sample and your comparison sample. For the fold change keep the default value 2.

- How many genes did you find? Do the results seem plausible?

- Are any of these genes also up-regulated in the replicative insect stage compared to the transmissive insect stage? How can you find this out? (*Hint*: add a step and run a microarray search comparing expression of epimastigotes to metacyclics).



- Do these genes have orthologs in other kinetoplastids? (*Hint*: add a step and transform your results into orthologs in all other organisms in TriTrypDB (select all for the ortholog transform).

  How many orthologs exist in *L. braziliensis* MHOM/BR/75/M2903? (*Hint*: look at the filter table between the strategy panel and your result list. Click on the number in the table under a species to view results from a specific species). Explore your results. Scan the product descriptions for this list of genes. Did you find anything interesting? Perhaps a GO enrichment analysis would support your ideas.

3. Finding genes based on RNAseq evidence and inferring function of hypothetical genes.
   Note: Use http://plasmodb.org for this exercise.

a.   Find all genes in *P. falciparum* that are up-regulated at least 50-fold in ookinetes compared to other stages: "**Transcriptomes of 7 sexual and asexual life stages (Lopez-Barragan et al.)**". For this search select "average" for the operation applied on the reference samples.



b.   The above search will give you all genes that are up-regulated by 50 fold in ookinetes compared to the average expression level of other stages. Despite the high fold change, some genes in the list may be highly expressed in the other stages. How can you remove genes from the list that are highly expressed in the other stages?



-   *Hint: Add a search for genes based on RNA Seq evidence from the same experiment, but this time select the percentile search: P.f. seven stages - RNA Seq (percentile). What*

*minimal percentile values should you choose?  40 – 100%? How does setting the any / all samples impact the result …. Which would be better in this case?*



- *Hint II: Try changing the operator from average to maximum for the set of non-ookinete stages in your initial fold change search. What does this do?  How do the resulting genes compare with the two step strategy you generated in the first hint?  Which hint do you think works better?*



c.   Which metabolic pathways are represented in this gene list? *Hint: add a step and transform results to pathways.*  How does this result compare to running a pathways enrichment on step 2?

d.      What happens if you revise the first step and modify the fold difference to a lower value -
        10 for example?  Compare results when you also modify the "between each genes"
        parameter.  What happens if you set this to maximum?  Which value do you think is most
        stringent for ensuring a 10 fold up regulation compared to the other samples?

e.      PlasmoDB also has an experiment examining gene expression during sexual development
        in *Plasmodium berghei* (rodent malaria).  Can you determine if there are genes that are up-
        regulated in both human and rodent ookinetes (compared to all other stages)? *Hint:* start
        *by deleting the last step you added in this exercise (transform to pathways). To do this click*
        *on edit then delete in the popup. Next, add steps for the P. berghei experiments "P berghei*
        *ANKA 5 asexual and sexual stage transcriptomes RNASeq". Why did you get 0 results? Hint:*
        *we are comparing results from different organisms …. Click the edit link in either step and*
        *choose orthologs to transform to appropriate organism.  Try it both ways … do you get the*
        *same number of genes?  Why does the strategy make a nested strategy when you*
        *transform the last step and not when  you transform the second to last step?*



## High-throughput phenotyping searches.

4.      Find genes that are essential in procyclics but not in blood form *T. brucei*.  Note that this
        search uses the same search form as RNASequencing searches but it is NOT
        RNASequencing.  Read the search description at the bottom of the High Throughput
        Sequencing search page for information about this assay.
        For this exercise use http://TriTrypDB.org.

-       Find the query for High Throughput Phenotyping. Think about how to set up this query
        (*Hint*: you will have to set up a two-step strategy). Remember you can play around with
        the parameters but there is no one correct way of setting them up –

## Identify Genes based on High-Throughput Phenotyping

Tutorial

For the **Experiment** Quantitated from the CDS Sequence
return protein coding ▼ **Genes**
that are Decrease in coverage ▼
with a **Fold change** >= 1.5

between each gene's **expression value**
in the following **Reference Samples**

- ⦿ Uninduced sample

and its **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**

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

For each gene, the search calculates:

$$fold\ change = \frac{reference\ expression\ level}{comparison\ expression\ level}$$

and returns genes when **fold change** >= **1.5**.
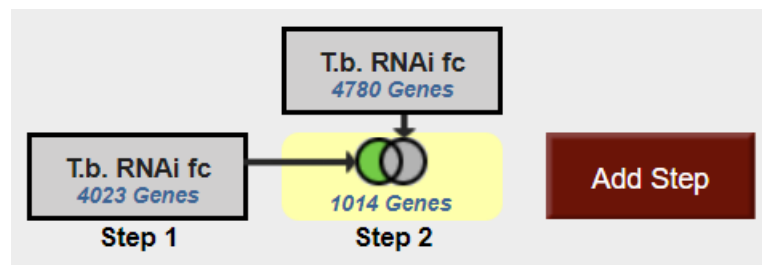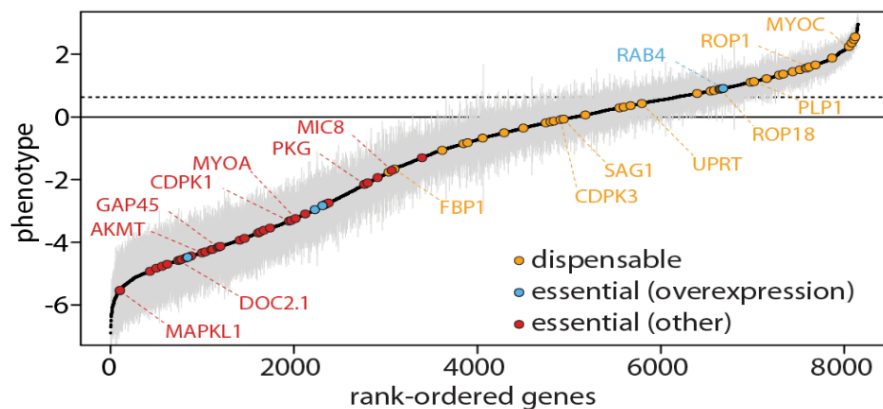
See the detailed help for this search.

**Get Answer**

- Next add a step and run the same search except this time select the "induced bloodstream form" samples.

- How did you combine the results? Remember you want to find genes that are essential in procyclics and not in blood form.
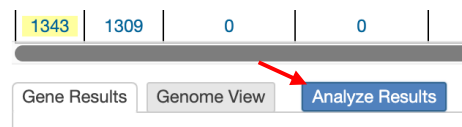
5. **Finding genes based on high throughput mutagenesis and fitness analysis.**

In EuPathDB we have a variety of studies where genome scale phenotypic analyses were carried out.   In this exercise we'll use ToxoDB.org and look at fitness following CRISPR mutagenesis.  You could also explore phenotyping studies in PlasmoDB or FungiDB if you prefer, the principles are the same.

- 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 last essential gene has a fitness score just > than -2 so setting the phenotype score <= -4 would provide a pretty stringent search but still return more than 1000 genes.  Try it.  Do you get the expected results based on the number of genes returned?
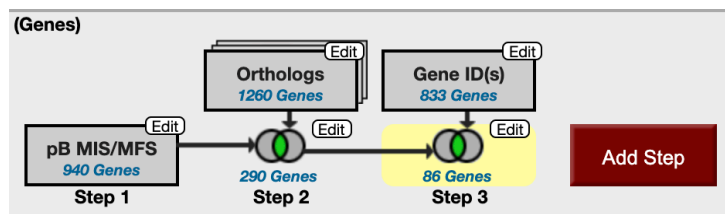


- Can you find additional evidence that these genes are essential?  One way is to use the analysis tools to assess biological process and go function.  Are the results what you would expect?



- Try adding columns to show additional data or intersecting these results with other queries, perhaps expression queries, to further assess this list.  NOTE: this experiment was done in GT1 while all *T. gondii* functional data in ToxoDB is mapped to ME49 so an ortholog transform to ME49 is required before adding any additional functional studies.

- Optional, try intersecting your results with the results from the previous exercise or one of the experiments in PlasmoDB. NOTE that we don't make this easy ☺. Due to technical reasons, most of these searches are not available in EuPathDB and we can't currently generate orthologs between component websites. I accomplished this by re-running the CRISPR search in EuPathDB, doing an ortholog transform to *P. falciparum* 3D7 and downloading the results as a list of IDs. I then went to PlasmoDB and ran a strategy intersecting the results of the two high throughput assays (incorrectly called curated phenotype searches ☹). Then added an id search using the list I downloaded from the CRISPR results. The following strategy shows this effort (https://tinyurl.com/ThreeOrgPheno).



- click the edit link by the Boolean operator in step 2 and try excluding the first or second step and compare your results.



Try ignoring each of the first two steps. Which gives you the most results when the other is ignored? Can you think of plausible reasons for this difference?

- You've now identified a core set of genes conserved across three Apicomplexan species that are likely to be essential. One criterion you might want in a parasite drug target would be that it not be conserved in humans as you wouldn't want the drug to impact the host. Are any of these genes also not conserved in humans (or more stringently mammals)? *Hint: add an Orthology Phylogenetic Profile search.*

6. *Optional (come back if time).* Finding oocyst expressed genes in *T. gondii* based on microarray evidence.

   Note: For this exercise use http://toxodb.org



a. Find genes that are expressed at 10 fold higher levels in one of the oocyst stages than in any other stage in the "**Oocyst, tachyzoite, and bradyzoite developmental expression profiles (M4) (John Boothroyd)**" microarray experiment.
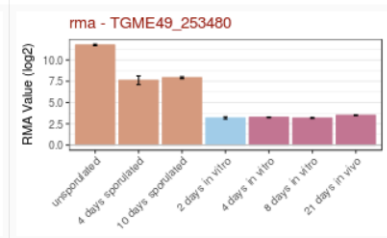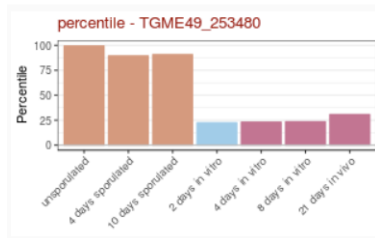
b.  <u>Add a step</u> to limit this set of genes to only those for which all the non-oocyst stages are expressed below 50$^{th}$ percentile … ie likely not expressed at those stages. *(Hint*: after you click on add step find the same experiment under microarray expression and chose the percentile search).

- Select the 4 **non-oocyst** samples.

- We want all to have less than 50$^{th}$ percentile so set *minimum percentile* to 0 and *maximum percentile* to 50.

- Since we want all of them to be in this range, choose ALL in the "*Matches Any or All Selected Samples*".



- To view the graphs in the final result table, turn on the columns called "TgM4 OoTachyBrady Marray - Expr Graph" and "TgM4 OoTachyBrady Marray - %ile Graph" (inside the "T. gondii ME49 Oocyst, tachyzoite, and bradyzoite developmental expression profiles (M4) (Fritz and Buchholz et al.)" Microarray).

7.  Comparing RNA abundance and Protein abundance data.
    Note: for this exercise use http://TriTrypDB.org.

    In this exercise we will compare genes that show differential RNA abundance levels between procyclic and blood form stages in *T. brucei* with genes that show differential protein abundance in these same stages.

a.  Find genes that are down-regulated 2-fold in procyclic form cells.  Go to the search page for Genes by Microarray Evidence and select the fold change search for the "Expression profiling of five life cycle stages (Marilyn Parsons)" experiment and configure the search to return protein-coding genes that are down-regulated 2 fold in procyclic form (PCF) relative to the Blood Form reference sample. Since there are two PCF samples, it is reasonable to choose both and average them.



b.  Add a step to compare with quantitative protein expression.  Select protein expression then "Quantitative Mass Spec Evidence" and the "Quantitative phosphoproteomes of bloodstream and procyclic forms (Tb427) (Urbaniak et al.)" experiment. Configure this search to return genes that are down-regulated in procyclic form relative to blood form.

**Add Step**

**Add Step 2 : T. brucei brucei TREU927 Quantitative phosphoproteomes of bloodstream and procyclic forms (Tb427) Proteomics (direct comparison)**

❓ Experiment

[ Quantitative phosphoproteomes of bloodstream and procyclic forms (Tb427) ▾ ]

❓ Direction

[ down-regulated ▾ ]

❓ Comparison

◉ Pcf-Bsf ratio

❓ Fold difference >=

[ 2 ]

**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 ]

c.    How many genes are in the intersection?  Does this make sense? Make certain that you set the directions correctly.



d.    Try changing directions and compare up-regulated genes/proteins. (*Hint:* revise the existing strategy … you might want to duplicate it so you can keep both).  When you change one of the steps but not the other do you have any genes in the intersection?  Why might this be?

e.    Can you think of ways to provide more confidence (or cast a broader net) in the microarray step? (*Hint:* you could insert steps to restrict based on percentile or add a RNA Sequencing step that has the same samples).

8.  Find genes with evidence of protein phosphorylation in intracellular *Toxoplasma* tachyzoites. For this exercise use http://www.toxodb.org

Phosphorylated peptides can be identified by searching the appropriate experiments in the Mass Spec Evidence search page.

**8a.** Find all genes with evidence of protein phosphorylation in intracellular tachyzoites. Navigate to the Post-Translational Modification search. Select the "**Infected host cell, phosphopeptide-enriched** (peptide discovery against **TgME49**)" sample under the experiment called "**Tachyzoite phosphoproteome from purified parasite or infected host cell (RH) (Treeck et al.)**"



**8b.** Remove all genes with phosphorylation evidence from purified tachyzoites and the phosphopeptide depleted fractions.

Hint: Use the Mass Spec Evidence search to access the tachyzoite and depleted fractions. Subtract (1 minus 2) these results from your first search.

**8d.** Explore your results. What kinds of genes did you find? *Hint: use the Product description word column or perform a GO enrichment analysis of your results.*

**8e.** Are any of these genes likely to be secreted? Hint: add a step searching for genes with secretory signal peptides.



**8f.** Pick one or two of the hypothetical genes in your results and visit their gene pages. Can you infer anything about their function? Hint: explore the protein and expression sections.

**8g.** What about polymorphism data? Go back to your strategy and add columns for SNP data found under the population biology section. Explore the gene page for the gene that has the highet number of non-synonymous SNPs. Hint: you can sort the columns by clicking on the up/down arrows next to the column names.

☐ Hide search strategy panel

(Genes)

Strategy: *Post-Translational Mod* *

Mass Spec
4651 Genes

Signal Pep
51001 Genes

Post-Translation
2332 Genes
Step 1

67 Genes
Step 2

18 Genes
Step 3

Add Step

Rename
Duplicate
Save As
Share
Delete

18 Genes from Step 3  Revise
Strategy: *Post-Translational Mod*

☐ 🔻 Click on a number in this table to limit/filter your results

| | | Cyclospora | Cystoisospora | C.suis | | | Eimeria | | | | | | Hammondia | Neospora | Sarcocystis | | | | | | Toxoplasma | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| All Results | Ortholog Groups | C.cayetanensis | | | E.acervulina | E.brunetti | E.falciformis | E.maxima | E.mitis | E.necatrix | E.praecox | E.tenella | H.hammondi | N.caninum | S.neurona ( 0 ) | | | | | | | T.gondii ( 18 ) | | | | | | | |
| | | strain CHN_HEN01 | strain Wien I | Houghton | Houghton | Bayer Haberkorn 1970 | Weybridge | Houghton | Houghton | Houghton | strain Houghton | strain H.H.34 | Liverpool | SN3 | SO | SN1 | ARI | FOU | GAB2-2007-GAL-DOM2 | GT1 | MAS | ME49 | RH | RUB | TgCatPRC2 | VAND | VEG | p89 |
| 18 | 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 18 | 0 | 0 | 0 | 0 | 0 | 0 |

Gene Results | Genome View | Analyze Results

Rows per page: 20 ▼

Download  Add to Basket  Add Columns

| 🛒 | ⬍ Gene ID | ⬍ Transcript ID | ⬍ Product Description ❓ ✖ 📊 | ⬍ # Transcripts ❓ ✖ 📊 | Non-Coding SNPs All Strains ❓ ✖ | NonSyn/Syn SNP Ratio All Strains ❓ ✖ | NonSynonymous SNPs All Strains ❓ ✖ | SNPs with Stop Codons All Strains ❓ ✖ | Synonymous SNPs All Strains ❓ ✖ | Total SNPs All Strains ❓ ✖ |
|---|---|---|---|---|---|---|---|---|---|---|
| 🛒 | TGME49_288370 | TGME49_288370-t26_1 | hypothetical protein | 1 | 83 | 2.34 | 75 | 0 | 32 | 190 |
| 🛒 | TGME49_288880 | TGME49_288880-t26_1 | hypothetical protein | 1 | 158 | 3.29 | 56 | 0 | 17 | 231 |
| 🛒 | TGME49_243290 | TGME49_243290-t26_1 | hypothetical protein | 1 | 216 | 1.08 | 43 | 0 | 40 | 299 |
| 🛒 | TGME49_205625 | TGME49_205625-t26_1 | hypothetical protein | 1 | 207 | 1.62 | 55 | 0 | 34 | 296 |
| 🛒 | TGME49_259830 | TGME49_259830-t26_1 | diacylglycerol kinase catalytic domain-containing protein | 1 | 139 | 0.61 | 14 | 0 | 23 | 176 |
| 🛒 | TGME49_257595 | TGME49_257595-t26_1 | hypothetical protein | 1 | 131 | 2.32 | 130 | 0 | 56 | 317 |
| 🛒 | TGME49_229680 | TGME49_229680-t26_1 | hypothetical protein | 1 | 28 | 0 | 0 | 0 | 5 | 33 |

9. Find and explore the metabolic pathway for glycolysis.
For this exercise use http://plasmodb.org

Navigate to the search page for Identify Metabolic Pathways based on Pathway Name/ID.

– Metabolic pathway and compound searches are available under the "Identify Other Data Types" head on the home page. You can find metabolic pathways based on the pathway name, genes involved in the pathway, or compounds involved in the pathway. Search for the **glycolysis** pathway using the Pathway Name/ID option.

– This search is equipped with a type-ahead function for choosing the metabolic pathway name. Begin typing glycolysis and then choose the pathway name from the list that appears.

Search for Other Data Types

expand all | collapse all

Find a search... 🔍 ❓

▸ Popset Isolate Sequences
▸ Genomic Sequences
▸ Genomic Segments
▸ SNPs
▸ SNPs (from Array)
▸ ESTs
▸ ORFs
▾ Metabolic Pathways
  • Compounds
  • Genes
  • Identifier (pathway, gene, compound, etc.)
  • Pathway Name/ID
▸ Compounds

expand all | collapse all

a. Examine the Glycolysis / Gluconeogenesis pathway.



Identify Metabolic Pathways based on Pathway Name/ID

- The search takes you straight to the record page for the Glycolysis / Gluconeogenesis (ec00010) metabolic pathway from KEGG. The overview section of the record page contains an interactive graphical representation of the pathway. The pathway map and the legend can be repositioned.

    A. Initial pathway view is zoomed out.
    B. Zoom in to see more details including EC numbers and metabolite structures.
    C. Click on a metabolite structure to get additional information.
    D. Click on the EC number to get more info about the enzyme including links to retrieve all genes in the database assigned to this EC number.



    E. The drop-down menu under the heading "Paint Enzymes" allows you paint the pathway based on experiments or based on phyletic pattern.
    F. Painting pathway by experiment provides a graphical representation of experimental results. Click on the graph to see more details.

G. Painting pathway based on phyletic pattern provides a graphical representation of phyletic distribution. Clicking on the phyletic pattern graphic provides additional information.

- Use the Tool Box to move (drag) the map and individual nodes. Zoom in and out to help explore the map.
- What do the rectangles with numbers like 2.7.1.11 represent?
- What is the difference between the rectangular nodes that are orange and those that are not?
- Why are some enzymes grouped?
- Find the node representing 6-phosphofructokinase (EC number = 2.7.1.11). You may need to zoom and reposition the map to find the node.

- Click on the 2.7.1.11 node to open a popup with information about this enzyme.



- How many genes in the database matched this EC number?

- Try the link 'Search for Gene(s) by EC Number'. Where did you end up? What do the 90 genes in the result list represent? Is 6-phosphofructokinase unique to *P. falciparum*? Notice the two columns called "EC numbers" and "EC numbers from OrthoMCL". What do these columns represent?

‒ Use your Browser's back button to return to the Glycolysis pathway record page and open the Paint Experiment menu. Choose the experiment "Salivary gland sporozoite transcriptomes: WT vs Pfu2-KO". Be patient while the graphs appear in place of the EC numbers.

- Does 6-phosphofructokinase appear to be expressed in salivary gland sporozoites? What enzymes in this pathway are affected in knockouts of Pfu2?



- Use the Paint Genera option to determine whether 6-phosphofructokinase has orthologs across Apicomplexa and Chromerida.

- What about the enzyme that catalyzes the reverse reaction (Fructose-bisphosphatase)?



10. **Find and explore the compound record page for phosphoenolpyruvate (phosphoenolpyruvic acid or PEP).**

Compound records are accessed by running a compound search available under the "Identify Other Data Types" heading on the home page. For example, compounds may be retrieved by ID, text, metabolic pathway, molecular formula, molecular weight and metabolite levels. Compound records can also be accessed from the metabolic pathway legend after clicking on a compound (blue circle) in the map.



- Choose one of these searches and retrieve the PEP record page.

- Alternatively, you can reach the PEP record page via a metabolic pathway where it is present as a substrate or a product of an enzymatic reaction (ie. glycolysis). Click on the node representing a compound

- Which method did you use to get to the PEP record page? What compound name worked the best?

- Examine the PEP record page.

- What data sections do you see?

- Under which conditions is PEP present at highest concentrations? (Hint: navigate to the Metabolomics section)



11. Identify metabolites (compounds) that are 20-fold enriched at pH7.4 in saponin lysed infected red blood cell (iRBCs) pellets compared the pH7.4 percoll pellet.

This requires running a metabolite levels search (2-fold enriched in saponin pellet compared to the percoll pellet as the reference).

## Identify Compounds based on Metabolite levels



- How many compounds did you get?
- How many of these compounds (metabolites) are NOT enriched by 2-fold in the pH7.4 saponin media fraction compared to the percoll media as reference?

To which metabolic pathways do these compounds belong?  Click Add Step and transform
the results to metabolic pathways.