

Exploring Transcriptomic and Proteomic data

1. Find all *P. falciparum* genes that are up-regulated during the later stages of the intra-erythrocytic cycle.

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

- a. Use the fold change search for the data set “Transcriptome during intra-erythrocytic 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.

Identify Genes by:

- Expand All | Collapse All
- Text, IDs, Organism
- Genomic Position
- Gene Attributes
- Protein Attributes
- Protein Features
- Similarity/Pattern
- Transcript Expression
 - EST Evidence
 - SAGE Tag Evidence
 - Microarray Evidence
 - RNA Seq Evidence
 - ChIP on Chip Evidence
 - TF Binding Site Evidence
- Protein Expression
- Cellular Location
- Putative Function
- Evolution
- Population Biology

Identify Genes based on RNA Seq Evidence

Filter Data Sets: Type keyword(s) to filter

Legend: ☒ FC Fold Change ☐ FqV Fold Change... ☐ P Percentile

Organism	Data Set	Choose a search
<i>P. falciparum</i> 3D7	Transcriptome during intraerythrocytic development (Bartfai et al.)	<input checked="" type="button"/> FC <input type="button"/> FqV <input type="button"/> P
<i>P. falciparum</i> 3D7	Transcriptomes of 7 sexual and asexual life stages (Lopez-Barragan et al.)	<input type="button"/> FC <input type="button"/> FqV <input type="button"/> P
<i>P. falciparum</i> 3D7	Strand specific transcriptomes of 4 life cycle stages (Lopez-Barragan et al.)	<input type="button"/> FC <input type="button"/> P
<i>P. falciparum</i> 3D7	NSR-seq Transcript Profiling of malaria-infected pregnant women and children (Vignali et al.)	<input type="button"/> FC <input type="button"/> FqV <input type="button"/> P

Identify Genes based on P.f. post infection (RBC) RNA-seq time series (fold change)

For the Experiment Post-Infection (RBC) RNA-Seq time Series

return protein coding Genes

that are up or down regulated

with a Fold change >= 2

between each gene's expression value

in the following **Reference Samples**

☐ Hour 5
☐ Hour 10
☐ Hour 15
☐ Hour 20
☐ Hour 25
☐ Hour 30
 select all | clear all

and its expression value

in the following **Comparison Samples**

☐ Hour 5
☐ Hour 10
☐ Hour 15
☐ Hour 20
☐ Hour 25
☐ Hour 30
 select all | clear all

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

Up or down regulated

Expression

This graphic will help you visualize the parameter choices you make at the left. It will begin to display when you choose a Reference Sample or a Comparison Sample.

See the [detailed help for this search.](#)

Advanced Parameters

Get Answer

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** \geq the intensity of difference in expression needed before a gene is returned by the search. **Choose 4** but feel free to modify this.
- **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 (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. **Choose average**
- **Reference Sample:** the samples that will serve as the reference when comparing expression between samples. **choose 5, 10, 15, 20, 25**
- **And its AVERAGE expression value:** This is the operation applied to comparison samples. see explanation above. **Choose average**
- **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 **choose 30, 35, 40**

Fold Change

Fold Change with pValue

Percentile

Identify Genes based on P.f. post infection (RBC) RNA-seq time series (fold change)

Tutorial

You Tube

For the Experiment

Post-Infection (RBC) RNA-Seq time Series

return

protein coding

Genes

that are

up-regulated

with a Fold change \geq

4

between each gene's

average

expression value

in the following

Reference Samples

☒ Hour 5
 ☒ Hour 10
 ☒ Hour 15
 ☒ Hour 20
 ☒ Hour 25

select all | clear all

and its

average

expression value

in the following

Comparison Samples

☐ Hour 20
 ☐ Hour 25
 ☒ Hour 30
 ☒ Hour 35
 ☒ Hour 40

select all | clear all

Example showing one gene that would meet search criteria

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

Up-regulated

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:

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

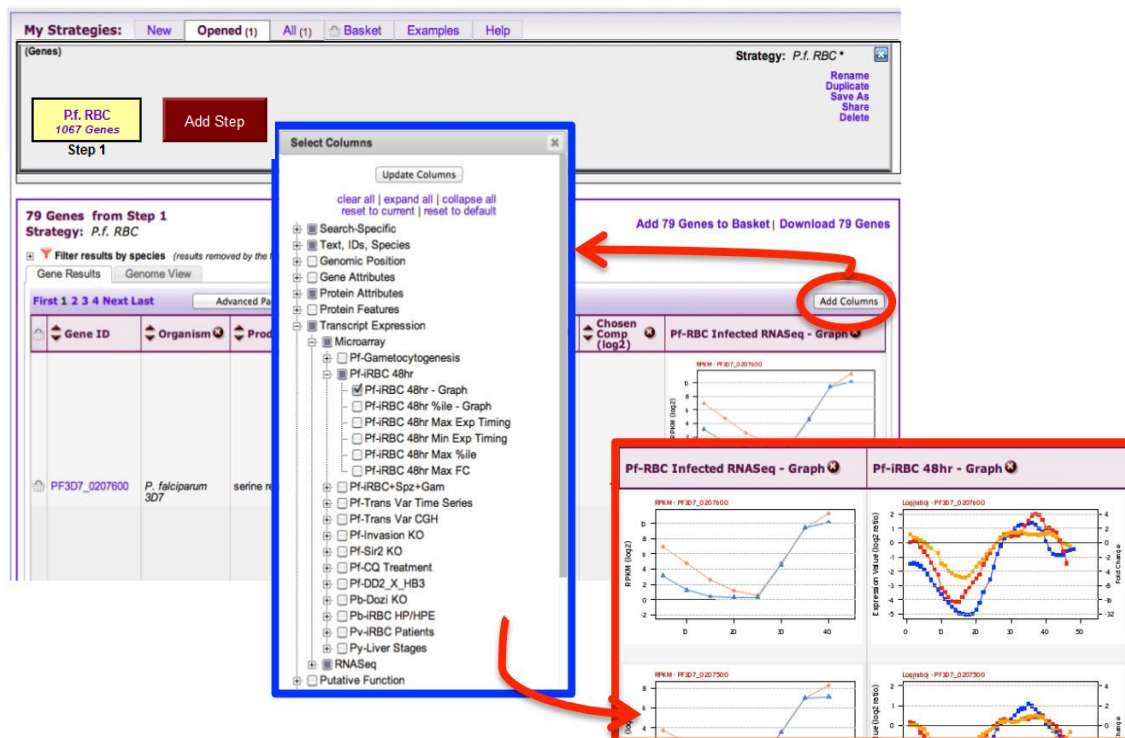
and returns genes when fold change ≥ 4 . To narrow the window, use the maximum reference value, or minimum comparison value. To broaden the window, use the minimum reference value, or maximum comparison value.

See the [detailed help for this search](#).

Advanced Parameters

Get Answer

- b. Compare the RNA Sequencing data to similar microarray data on a gene-by-gene basis.
- 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.)** or **Pf-iRBC 48hr** for shorter column headings. 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”.



- c. Compare the RNA-sequencing data to microarray data on a genome scale.

You can also run a fold-change search on the microarray data to compare results on a genome scale. Add a step to your strategy and intersect the results of 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 keeping the fold change ≥ 2 due to the decreased dynamic range of microarrays compared to RNA-seq.

- How many genes are upregulated in the later stages of the erythrocytic cycle based on microarray and RNA sequencing evidence?

P.f. RBC
1067 Genes
Step 1

Add Step

Run a new Search for
Transform by Orthology
Add contents of Basket
Add existing Strategy
Filter by assigned Weight
Transform to Pathways
Transform to Compounds

Genes
Genomic Segments
SNPs
SNPs (from Chips)
ORFs

Text, IDs, Organism
Genomic Position
Gene Attributes
Protein Attributes
Protein Features
Similarity/Pattern
Transcript Expression
Protein Expression
Cellular Location
Putative Function

EST Evidence
SAGE Tag Evidence
Microarray Evidence
RNA Seq Evidence
ChIP on Chip Evidence
TF Binding Site Evidence

Add Step 2 : Microarray Evidence

Filter Data Sets: Type keyword(s) to filter

Legend: AGS Associati... DC Direct Co... FC Fold Cha... FD Fold Diff... P Percentile S Similarity SA Similarity...

Organism	Data Set	Choose a search
<i>P. falciparum</i> 3D7	Asexual blood stage transcriptomes of clonal strains (Rovira-Graells et al.)	FC FD
<i>P. falciparum</i> 3D7	Asexual life cycle time course with biosynthetic pyrimidine labeling (Manuel Linas)	FC P
<i>P. berghei</i> ANKA	DOZI Mutant Transcript Profile (Mair et al.)	DC P
<i>P. falciparum</i> 3D7	eQTL for HB3, Dd2 and 34 progeny (Gonzales et al.)	AGS S SA
<i>P. falciparum</i> 3D7	Erythrocytic expression time series (3D7, Dd2, HB3) (Bozdech et al. and Linas et al.)	FC P S

iRBC 3D7
(48 Hour scaled)

Upregulated
2 fold

Average

1-16 hours,
17-30 hours

Average

31-48 hours

Add Step 2 : P.f. Intraerythrocytic Infection Cycle (fold change)

Experiment: iRBC 3D7 (48 Hour scaled)

return: protein coding

that are: up-regulated

Fold change >= 2

between each gene's average expression value

In the following Reference Samples

select all | clear all | expand all | collapse all | reset to default

☒ 1-16 Hours
☒ 17-30 Hours
☐ 31-48 Hours

select all | clear all | expand all | collapse all | reset to default

in the following Comparison Samples

select all | clear all | expand all | collapse all | reset to default

☐ 1-16 Hours
☐ 17-30 Hours
☒ 31-48 Hours

select all | clear all | expand all | collapse all | reset to default

Advanced Parameters

Combine Genes in Step 1 with Genes in Step 2:

Intersect

1 Intersect 2
1 Union 2
1 Relative to 2, using genomic colocation

1 Minus 2
2 Minus 1

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

Up-regulated

Expression

Average Comparison
Average Reference

2 fold

Reference Comparison Samples 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:

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

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.

Run Step

P.f. RBC
1067 Genes
Step 1

PfRBC 48HR FC
196 Genes
Step 2

Add Step

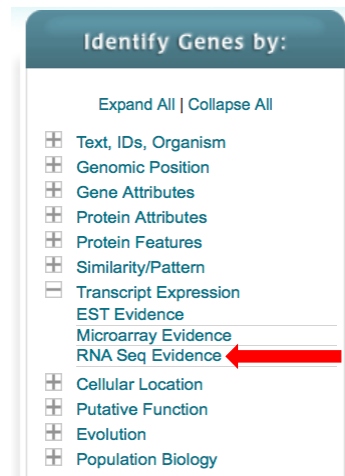
2. Exploring RNA-seq data in fungi

For this exercise use fungidb.org

C. albicans is a commensal fungal organism that inhabits skin, mucosal surfaces, and gut of humans. It can cause superficial and subcutaneous infections that may become life threatening in immunocompromised individuals. *C. albicans* has a yeast like non-invasive form, and a hyphal form associated with increased virulence and systemic infections. Let's look at the genes that are up-regulated in response to serum and high oxidative stress, conditions similar to those that are encountered by *C. albicans* during infection.

a. Identify genes that are up-regulated in YPD media supplied with serum.

- Navigate to the transcript expression data using the “Transcript Expression” menu in the “Identify Genes by” panel.



- Select the “Comprehensive Annotation of Transcriptome” by Michael Snyder.
- Find protein coding genes that are up-regulated by 2 fold in response to serum. Compare “YPD media + Serum” (comparison sample) to YPD Media (Reference control sample).

Revise Step 1 : C.alb. Comprehensive Annotation of Transcriptome RNASeq (fold change)

For the **Experiment** Comprehensive Annotation of Transcriptome ?

return protein coding ? **Genes**

that are up-regulated ?

with a **Fold change** ≥ 2 ?

between each gene's **expression value** ?

in the following **Reference Samples** ?

☐ No Nitr Stress
☐ No Oxi Stress
☐ YPD
☒ YPD Media ←
☐ YPD Media +Serum
[select all](#) | [clear all](#)

and its **expression value** ?

in the following **Comparison Samples** ?

☐ No Nitr Stress
☐ No Oxi Stress
☐ YPD
☐ YPD Media
☒ YPD Media +Serum ←
[select all](#) | [clear all](#)

[Run Step](#)

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

Up-regulated

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

For each gene, the search calculates:

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

and returns genes when **fold change** ≥ 2 .

[See the detailed help for this search.](#)

b. Identify those genes that are also up-regulated in response to high oxidative stress conditions.

- Click “Add step” and navigate to the same RNA-seq data set.

(Genes)

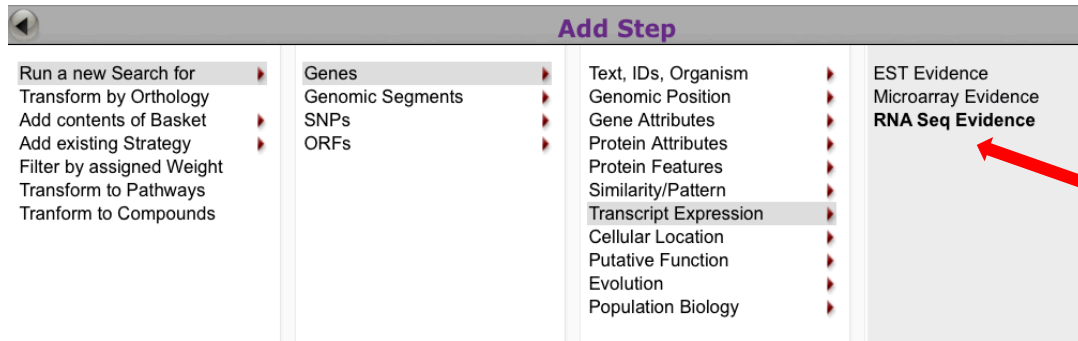
YPD+Serum
 411 Genes
 Step 1

Edit

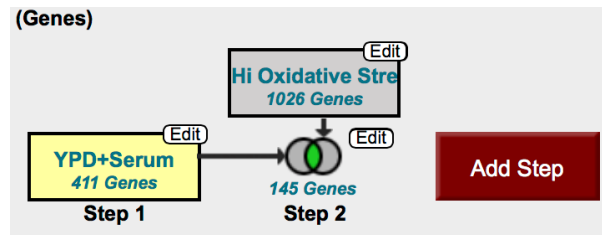
Add Step ←

Hint: To change the name of the first step to “YPD+Serum” access the “Edit” option in the search box and select “Rename”.

- Set up a search using “Hi Oxi stress” data set this time using the no oxidative stress control (“No Oxi”) as a reference sample.



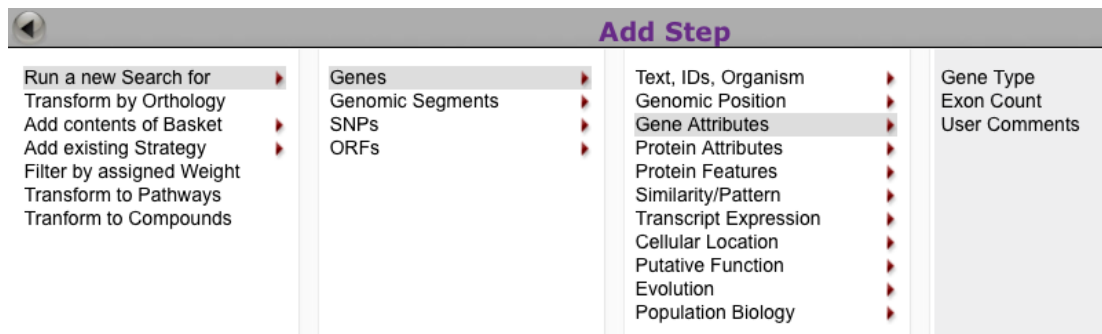
- Your search should return 145 genes:



c. **Identify genes with more than one exon.**

The Gene Attributes function offers an option to further customize your search to identify genes with more than one exon.

- Click “Add Step” and navigate to “Exon count” option:



- For this exercise leave the parameters of the search in *C. albicans* at default (between 2 and 20 exons):

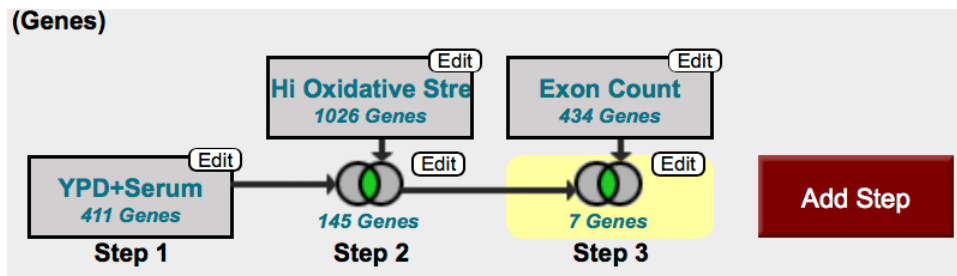
Organism [select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

- ☐ Oomycetes
- ☒ Fungi
 - ☐ Agaricomycetes
 - ☐ Blastocladiomycetes
 - ☐ Chytridiomycetes
 - ☐ Eurotiomycetes
 - ☐ Leotiomycetes
 - ☐ Pneumocystidomycetes
 - ☐ Pucciniomycetes
 - ☒ Saccharomycetes
 - ☒ Candida
 - ☒ Candida albicans
 - ☐ Candida glabrata
 - ☐ Saccharomyces
 - ☐ Yarrowia
 - ☐ Schizosaccharomycetes
 - ☐ Sordariomycetes
 - ☐ Tremellomycetes
 - ☐ Ustilaginomycetes
 - ☐ Zygomycetes

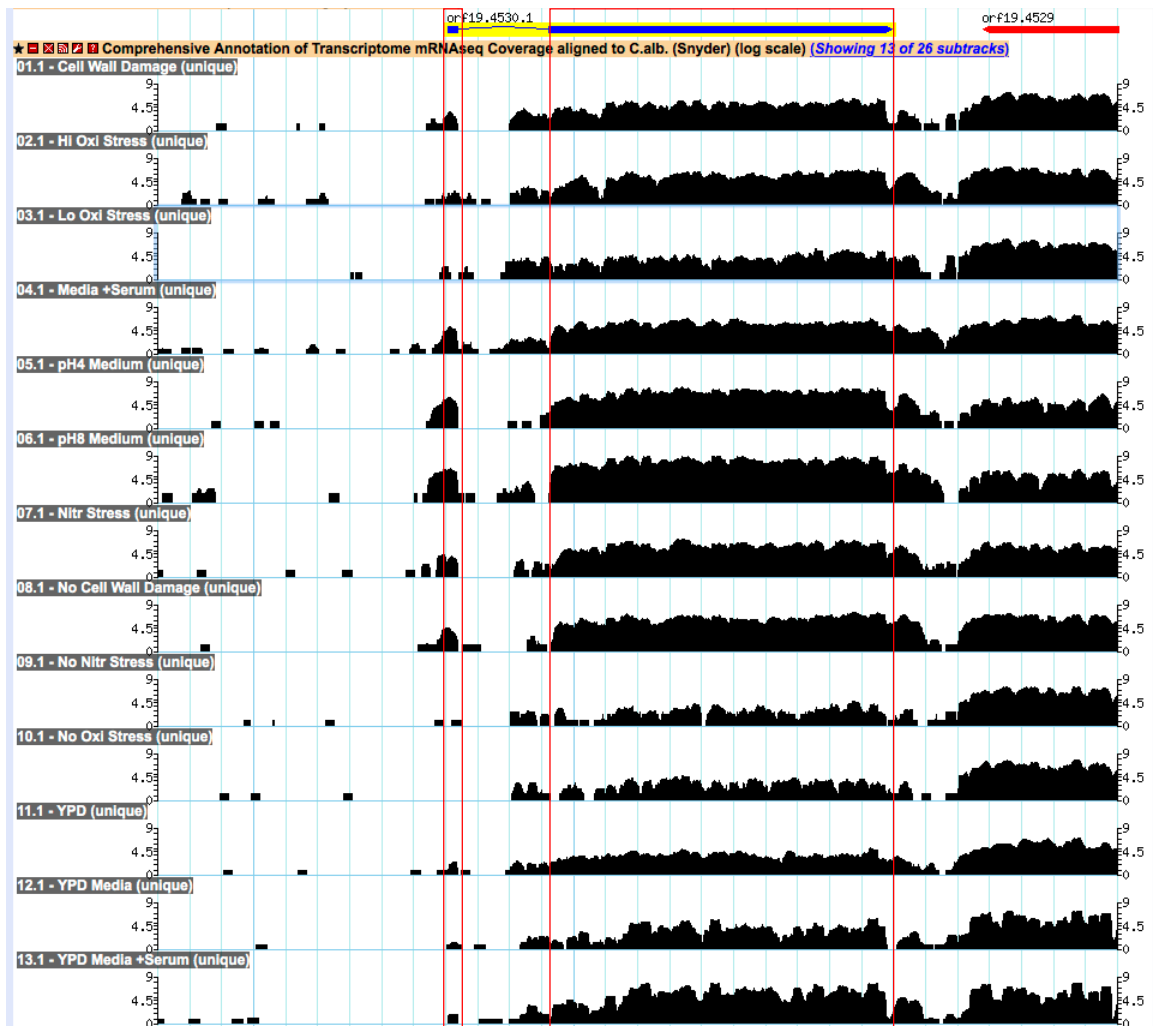
[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

Exon Count >=

Exon Count <=



- Explore the first gene in the list - **orf19.4530.1**.
- d. Explore RNA-seq tracks and genome annotations using Genome Browser.
- From the gene page navigate to the Genome Browser window.
 - Turn on the track called “Comprehensive Annotation of Transcriptome mRNAseq Coverage aligned to C.alb. (Snyder) (log scale)” by navigating to the Select tracks tab in the GBrowse window.
 - Examine how well RNA-seq tracks correspond to current gene model.



- Do these results make sense?
- Do you think there is evidence for alternative splicing or intron read through due to various growth conditions?
- Can you explain transcript mapping outside the current annotation coordinates?
Hint: You can choose which tracks to view (turn off / turn on tracks) by clicking on the blue link at the top "Showing 13 of 26 sub-tracks" and further adjusting the view options.

3. Find sporozoite-specific *Cryptosporidium* genes that are expressed at the protein level and are likely secreted.

For this exercise use <http://cryptodb.org>

- 3 a. Find *Cryptosporidium* genes that are expressed at the protein level with evidence from any of the sporozoite proteomics experiments available in CryptoDB. Explore the available proteomics data and select samples that make sense. You may need to click on the '+' sign to expand experiments to see the underlying samples.

Identify Genes based on Mass Spec. Evidence

Experiment/Samples [?](#) [select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

- ☒ **Cryptosporidium**
 - ☒ ***Cryptosporidium parvum***
 - ☐ Enriched cytoskeletal and membrane fractions (Madrid-Aliste et al.)
 - ☐ Mitochondrial Fraction Proteomics (Putignani)
 - ☒ **Oocyst Wall Proteome (lowall) (Ferrari)**
 - ☐ Intact Oocysts
 - ☐ Oocyst walls
 - ☒ Sporozoites
 - ☒ **Proteome during Sporozoite Excystation (ISSC162) (Snelling et al.)**
 - ☒ Insoluble Excysted Fraction LC-MS/MS
 - ☐ Insoluble Non-excysted Fraction LC-MS/MS
 - ☐ Soluble Excysted and Non-excysted Fraction LC-MS/MS
 - ☒ **Sporozoite Proteome (lowall) (Sanderson et al.)**
 - ☒ 1D Gel LC-MS/MS
 - ☒ 2D Gel LC-MS/MS
 - ☒ MudPit Insoluble fractions
 - ☒ MudPit Soluble fractions

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

Minimum Number of Unique Peptide Sequences [?](#)

Minimum Number of Spectra [?](#)

[Advanced Parameters](#)

[Get Answer](#)

- 3 b. Remove any gene with peptide evidence from non-sporozoite samples
Hint: add a step for mass spec data and think about how you will combine your results.

Add Step 2 : Mass Spec. Evidence

Experiment/Samples [select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

- ☒ **Cryptosporidium**
 - ☒ **Cryptosporidium parvum**
 - ☐ Enriched cytoskeletal and membrane fractions (Madrid-Aliste et al.)
 - ☐ Mitochondrial Fraction Proteomics (Putnam)
 - ☒ Oocyst Wall Proteome (Iowall) (Ferrari)
 - ☒ Intact Oocysts
 - ☒ Oocyst walls
 - ☐ Sporozoites
 - ☐ Proteome during Sporozoite Excystation (ISSC162) (Snelling et al.)
 - ☐ Insoluble Excysted Fraction LC-MS/MS
 - ☒ Insoluble Non-excysted Fraction LC-MS/MS
 - ☐ Soluble Excysted and Non-excysted Fraction LC-MS/MS
 - ☐ Sporozoite Proteome (Iowall) (Sanderson et al.)
 - ☐ 1D Gel LC-MS/MS
 - ☐ 2D Gel LC-MS/MS
 - ☐ MudPit Insoluble fractions
 - ☐ MudPit Soluble fractions

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

Minimum Number of Unique Peptide Sequences

Minimum Number of Spectra

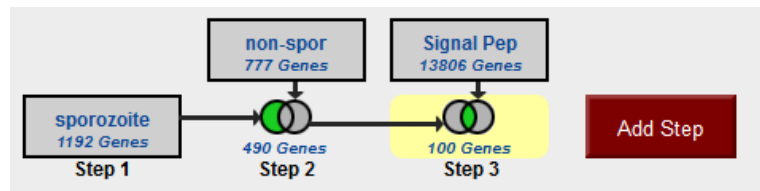
Advanced Parameters

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](#)

3 c. How many of these genes are also predicted to be secreted?



3 d. So far you have been searching for *C. parvum* genes because we only have proteomics data from this species. However, what if you are studying *C. muris*? How can you garner information about the protein expression of *C. muris* genes from your *C. parvum* results? (Hint: add a step then select the “Transform by Orthology” option).

- Did the number of *C. parvum* genes increase or decrease? Why?

Add Step 4 : Transform by Orthology

Organism [select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

- ☒ **Apicomplexa**
 - ☒ **Cryptosporidium**
 - ☒ Cryptosporidium hominis
 - ☒ Cryptosporidium hominis TU502
 - ☒ Cryptosporidium muris
 - ☒ Cryptosporidium muris RN66
 - ☒ Cryptosporidium parvum
 - ☒ Cryptosporidium parvum Iowa II
 - ☒ Gregarina
 - ☒ Gregarina niphandrodes Unknown strain
 - ☒ Chromera
 - ☒ Chromera vella CCMP2878
 - ☒ Vitrella
 - ☒ Vitrella brassicaformis CCMP3155

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

Syntenic Orthologs Only? ☐ no ☒ yes

Advanced Parameters

[Run Step](#)

My Strategies: [New](#) | [Opened \(1\)](#) | [All \(27\)](#) | [Basket](#) | [Public Strategies](#)

(Genes)

399 Genes from Step 4
Strategy: sporozoite

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

All Results	Ortholog Groups	Apicomplexa			Chromerida		
		Cryptosporidium	Gregarina	Chromera	Vitrella		
		C.hominis	C.muris	C.parvum	G.niphandrodes	C.vella	V.brassicaformis
		TU502	RN66	Iowa II	Unknown strain	CCMP2878	CCMP3155
399	100	101	75	109	26	44	44

4. Comparing RNA abundance and Protein abundance data.

Note: for this exercise use <http://TriTrypDB.org>.

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

- 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.

Identify Genes by:

- Expand All | Collapse All
- Text, IDs, Organism
- Genomic Position
- Gene Attributes
- Protein Attributes
- Protein Features
- Similarity/Pattern
- Transcript Expression
- EST Evidence
- Microarray Evidence**
- Protein Expression
- Cellular Location
- Putative Function
- Evolution
- Population Biology

Identify Genes based on Microarray Evidence

Filter Data Sets: Type keyword(s) to filter

Legend: DC Direct Comparison FC Fold Change P Percentile

Organism	Data Set	Choose a search
<i>L. infantum</i> JPCM5	Expression profiling of the promastigote time-course (L.d. Samples) (Peter Myler)	FC P
<i>L. infantum</i> JPCM5	axenic and intracellular amastigote profiles (Barbara Papadopolou)	FC P
<i>L. major</i> strain Friedlin	Three Developmental Stages (Stephen M. Beverley)	DC FC P
<i>T. brucei</i> TREU927	Dynamic mRNA Expression analysis of cells undergoing synchronous life-cycle differentiation (Keith R. Matthews)	FC P
<i>T. brucei</i> TREU927	Expression profiling of five life cycle stages (Marilyn Parsons)	FC P
<i>T. brucei</i> TREU927	Procyclic TbDRBD3 Depletion (Antonio Estevez)	DC FC P
<i>T. brucei</i> TREU927	Expression profiling of in vitro differentiation time series (Christine Clayton)	FC P

Identify Genes based on T.b. Expression profiling of five life cycle stages Microarray (fold change)

For the Experiment (Expression profiling of five life cycle stages: 3) return genes that are (Down-regulated) with a Fold change >= 2 between each gene's (average) expression value in the following (Reference Samples):

- ☒ Blood Form
- ☒ Stumpy
- ☐ PCF Log
- ☐ PCF Stat

and its (average) expression value in the following (Comparison Samples):

- ☐ Blood Form
- ☐ Stumpy
- ☒ PCF Log
- ☒ PCF Stat

Protein Coding Only: ☒ protein coding

Example showing one gene that would meet search criteria

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

Down-regulated

Expression

Average Reference

Average Comparison

2 fold

Reference Comparison Samples Samples

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

For each gene, the search calculates:

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

and returns genes when fold change >= 2. To narrow the window, use the minimum reference value, or maximum comparison value. To broaden the window, use the maximum reference value, or minimum comparison value.

See the detailed help for this search.

Advanced Parameters

Get Answer

- 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.

The screenshot shows the 'My Strategies' interface. On the left, 'Step 1' is titled 'Tb LifeCyc Marra' and contains '553 Genes'. A red arrow points from the 'Add Step' button in Step 1 to the 'Add Step' button at the top of the interface. Another red arrow points from the 'Add Step' button at the top to the 'Add Step 2 : Quantitative Mass Spec. Evidence' window. This window shows a table of data sets for *T. brucei* TREU927. A red circle highlights the 'DC' button for the second data set, 'Quantitative phosphoproteomes of bloodstream and procyclic forms (Tb427) (Urbanik et al.)'. A third red arrow points from this circle to the 'Add Step 2 : T.bru. Quantitative phosphoproteomes of bloodstream and procyclic forms (Tb427) Proteomics (direct comparison)' window. This window shows 'Direction' set to 'down-regulated', 'Samples' set to 'Pcf-Bsf ratio', and 'Fold difference >= 2'. Below this, there are options for combining genes from Step 1 and Step 2: 'Intersect 2', 'Union 2', 'Relative to 2, using genomic colocation', '1 Minus 2', and '2 Minus 1'. A 'Run Step' button is at the bottom.

- 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).

5. OPTIONAL: Find genes with evidence of 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.

5a. Find all genes with evidence of phosphorylation in intracellular tachyzoites. 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.)”

Identify Genes based on Mass Spec. Evidence

Experiment/Samples [select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

- ☐ Eimeria
- ☒ Toxoplasma
 - ☒ *Toxoplasma gondii*
 - ☐ Oocyst Partially Sporulated Proteome (VEG) (Possenti, et al.)
 - ☐ Oocyst proteome (M4 Typell) (Wastling)
 - ☐ Oocyst proteome - Fractionated (M4 type II) (Fritz et al.)
 - ☐ Proteome During Infection in H. sapiens (Wastling)
 - ☐ Tachyzoite Intra- and Extracellular Lysine-Acetylomes (RH) (Jeffers and Xue)
 - ☐ Tachyzoite Rhoptry proteome (RH) (Bradley et al.)
 - ☐ Tachyzoite conoid proteome (RH) (Hu et al.)
 - ☐ Tachyzoite membrane and cytosolic fractions (RH) (Dybas et al.)
 - ☐ Tachyzoite phosphoproteome - Calcium dependent (RH) (Nebl et al.)
 - ☒ Tachyzoite phosphoproteome from purified parasite or infected host cell (RH) (Treeck et al.)
 - ☐ Infected host cell, phosphopeptide-depleted (peptide discovery against TgME49)
 - ☐ Infected host cell, phosphopeptide-depleted (peptide discovery against TgGT1)
 - ☒ Infected host cell, phosphopeptide-enriched (peptide discovery against TgME49)
 - ☐ Infected host cell, phosphopeptide-enriched (peptide discovery against TgGT1)
 - ☐ Purified tachyzoites phosphopeptide-depleted (peptide discovery against TgGT1)
 - ☐ Purified tachyzoites phosphopeptide-depleted (peptide discovery against TgME49)
 - ☐ Purified tachyzoites phosphopeptide-enriched (peptide discovery against TgGT1)
 - ☐ Purified tachyzoites phosphopeptide-enriched (peptide discovery against TgME49)
 - ☐ Tachyzoite secretome (RH) (Zhou et al.)
 - ☐ Tachyzoite subcellular fractions (Moreno)
 - ☐ Tachyzoite total proteome (RH) (Wastling)

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#) | [reset to default](#)

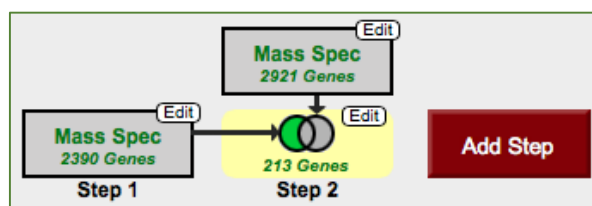
Minimum Number of Unique Peptide Sequences [?](#)

Minimum Number of Spectra [?](#)

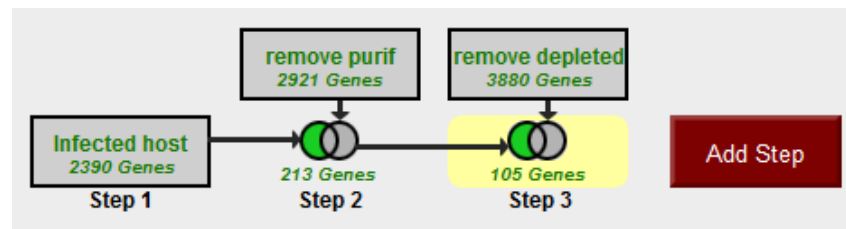
[Advanced Parameters](#)

[Get Answer](#)

5b. Remove all genes with phosphorylation evidence from purified tachyzoites.



5c. Remove all genes that are also present in the phosphopeptide-depleted fractions (select both intracellular and extracellular).



5d. 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.* Could you achieve this same 105 genes with a two-step strategy? *Hint: remove depleted and tachyzoite proteins in one step rather than two.*

5e. Are any of these genes likely to be secreted? *Hint: add a step searching for genes with secretory signal peptides.*

My Strategies: New Opened (3) All (3) Basket Public Strategies (7) Help

(Genes)

22 Genes from Step 4
Strategy: Infected host

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

All Results	Ortholog Groups	<i>E.acervulina</i>	<i>E.brunetti</i>	<i>E.falciformis</i>	<i>E.maxima</i>	<i>E.mitis</i>	<i>E.necatrix</i>	<i>E.praecox</i>	<i>E.tenella</i>	<i>H.hammondi</i>	<i>N.</i>
		Houghton	Houghton	Bayer Haberkorn 1970	Weybridge	Houghton	Houghton	Houghton	strain Houghton	strain H.H.34	L
22	22	0	0	0	0	0	0	0	0	0	

Filter by strains (advanced)

Gene Results Genome View Analyze Results **BETA**

First 1 2 Next Last Advanced Paging

Gene ID	Gene Group (representative gene)	Genomic Location	Product Description
TGME49_294940	TGGT1_294940	TGME49_chrla: 1,282,608 - 1,287,925 (-)	hypothetical protein
TGME49_222870	TGGT1_222870	TGME49_chrlt: 1,271,864 - 1,275,140 (+)	hypothetical protein
TGME49_320150	TGGT1_320150	TGME49_chrlv: 464,394 - 473,129 (-)	elongation factor Tu GTP binding domain-containing protein

5f. 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.*

5g. 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 most number of non-synonymous SNPs. Hint: you can sort the columns by clicking on the up/down arrows next to the column names.

Gene Results

Genome View

Analyze Results

BETA

First 1 2 Next Last

Advanced Paging

Add Columns

Gene ID	Product Description	Total SNPs All Strains	NonSynonymous SNPs All Strains	Synonymous SNPs All Strains	Non-Coding SNPs All Strains	SNPs with Stop Codons All Strains	NonSyn/Syn SNP Ratio All Strains
TGME49_271110	hypothetical protein	890	157	44	679	10	3.57
TGME49_257595	hypothetical protein	317	123	51	131	12	2.41
TGME49_219640	hypothetical protein	382	85	34	263	0	2.5
TGME49_288370	hypothetical protein	224	82	35	105	2	2.34
TGME49_216840	hypothetical protein	189	75	23	89	2	3.26
TGME49_257640	hypothetical protein	110	66	12	31	1	5.5
TGME49_320150	elongation factor Tu GTP binding domain-containing protein	378	65	22	286	5	2.95
TGME49_235960	hypothetical protein	155	58	14	77	6	4.14
TGME49_288880	hypothetical protein	220	56	17	147	0	3.29
TGME49_269750	CrcB family protein	95	54	20	18	3	2.7
TGME49_315700	hypothetical protein	338	54	14	265	5	3.86
TGME49_308070	hypothetical protein	188	43	22	123	0	1.95
TGME49_269420	hypothetical protein	45	37	8	0	0	4.63
TGME49_200440	hypothetical protein	72	35	11	24	2	3.18
TGME49_259830	diacylglycerol kinase catalytic domain-containing protein	176	32	3	139	2	10.67
TGME49_236220	PCI domain-containing protein	383	28	18	332	5	1.56
TGME49_231180	hypothetical protein	54	25	9	18	2	2.78
TGME49_294940	hypothetical protein	137	16	7	111	3	2.29