

Functional Genomics in EuPathDB Transcriptomics and Proteomics Exercise 3


- 1. Functional genomics data in EuPathDB databases includes transcription, protein and metabolic level data.**

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

- What kind of data types can be used to provide evidence of transcriptional activity? *Hint:* click on “Transcript Expression” to expand the list of possible searches.
- Explore organisms that have microarray data. What organisms have expressed sequence tag (EST), RNA sequence, ChIP-chip or SAGE tag data?
- What does RNA-seq data tell you that microarray data cannot? What does ChIP-chip data tell you about a gene?
- High throughput phenotyping is also a transcriptomic experiment. This data is located under putative function.
- What about protein expression data? What does quantitative proteomic data provide?

Identify Genes by:

Expand All | Collapse All

- ☒ Text, IDs, Organism
- ☒ Genomic Position
- ☒ Gene Attributes
- ☒ Protein Attributes
- ☒ Protein Features
- ☒ Similarity/Pattern
- ☒ Transcript Expression
 - EST Evidence
 - RT PCR Evidence
 - SAGE Tag Evidence
 - Microarray Evidence
 - RNA Seq Evidence
 - ChIP on Chip Evidence
 - TF Binding Site Evidence
- ☐ Protein Expression
 - Mass Spec. Evidence
 - Quantitative Mass Spec. Evidence
- ☒ Cellular Location
- ☐ Putative Function
 - GO Term
 - EC Number
 - Metabolic Pathway - MPMP
 - Y2H Protein Interaction
 - Predicted Functional Interaction
 - Phenotype
 -  High-Throughput Phenotyping
- ☒ Evolution
- ☒ Population Biology

- f. Go to the Data Summary Section, can you find the same information there? *Hint:* data summary table in on the left side of the home page.

EuPathDB Eukaryotic Pathogen Database Resources

Version 2.19 25 Sep 13

Gene ID: PF3D7_1133400 Gene Text Search: synth*

About EuPathDB | Help | EuPathDB Example's Profile | Logout | Contact Us

Home | New Search | My Strategies | My Basket (7) | Tools | Data Summary | Downloads | Community

Data Summary

EuPathDB Bioinformatics Resource Center for Biodefense and Emerging/Re-emerging Infectious Diseases is a portal for accessing genomic-scale datasets associated with the eukaryotic pathogens: (mouse over the logos: Acanthamoeba, Annacalia, Babesia, Cryptosporidium, Edhazardia, Eimeria, Encephalitozoon, Endotrypanum, Entamoeba, Enterocytozoon, Giardia, Gregarina, Hamiltosporidium, Leishmania, Nematocida, Neospora, Nosema, Plasmodium, Theileria, Toxoplasma, Trichomonas, Trypanosoma, Vivria, Vittiforma).

AmoeboDB CryptoDB GiardiaDB MicrosporidiaDB PiroplasmaDB PlasmoDB ToxoDB TrichDB TriTrypDB

Identify Genes by: Identify Other Data Types: Tools:

2. 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 upregulated based on RNA-seq data at late time points (30, 35 and 40-hours) compared to early time points in this experiment (1, 10, 15, 20, 25 hrs). *Hint:* for this exercise use a fold change search based on the "Transcriptome during intraerythrocytic development (Bartfai *et al.*)" experiment.

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 FCpV Fold Change... P Percentile

Organism	Data Set	Choose a search
P. falciparum 3D7	Transcriptome during intraerythrocytic development (Bartfai et al.)	FC FCpV P
P. falciparum 3D7	Transcriptomes of 7 sexual and asexual life stages (Lopez-Barragan et al.)	FC FCpV P
P. falciparum 3D7	Strand specific transcriptomes of 4 life cycle stages (Lopez-Barragan et al.)	FC P
P. falciparum 3D7	NSR-seq Transcript Profiling of malaria-infected pregnant women and children (Vignali et al.)	FC FCpV 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

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:

Fold Change

Fold Change with pValue

Percentile

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

Tutorial

For the Experiment

Post-Infection (RBC) RNA-Seq time Series

return

protein coding

Genes

that are

up-regulated

with a Fold change >=

12

between each gene's

average

expression value

in the following

Reference Samples

☒ Hour 5
☒ Hour 10
☒ Hour 15
☒ Hour 20
☒ Hour 25
☐ Hour 30
☐ Hour 35
☐ Hour 40

select all | clear all

and its

average

expression value

in the following

Comparison Samples

☐ Hour 15
☐ 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 >= 12. 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

Direction: the direction of change in expression. **Choose up-regulated.**

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

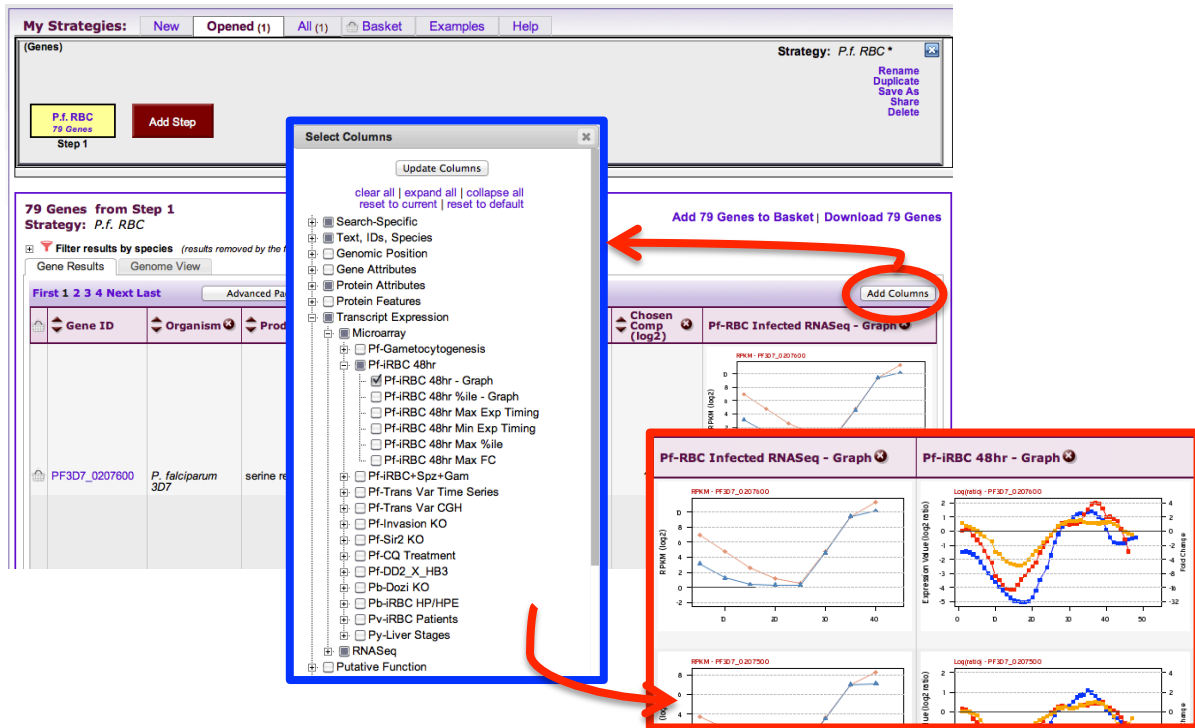
Operation Applied to Reference Samples: fold change is calculated as the ratio of two values (expression in reference)/(expression in comparison). When you choose multiple samples to serve as reference, we generate one number for the fold change calculation by using the minimum, maximum, or average. **Choose average**

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**

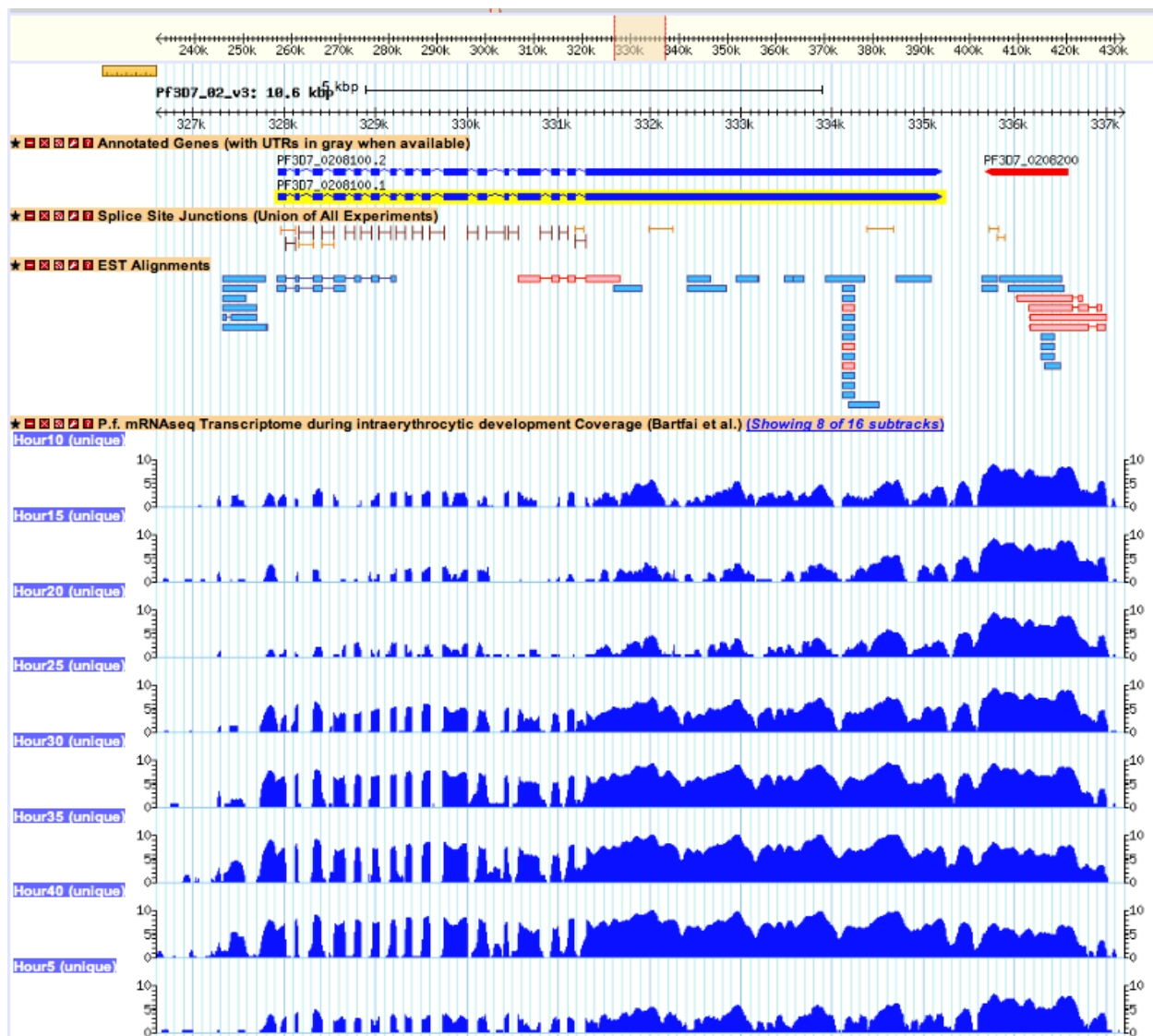
Operation Applied to Comparison Samples: see explanation above. **Choose average**

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.

- b. For the genes returned by the search, how does the RNA-sequence data compare to microarray data? (Hint: add the column called “Pf-iRBC 48hr - Graph” and compare the RNA-seq to the microarray graphs).

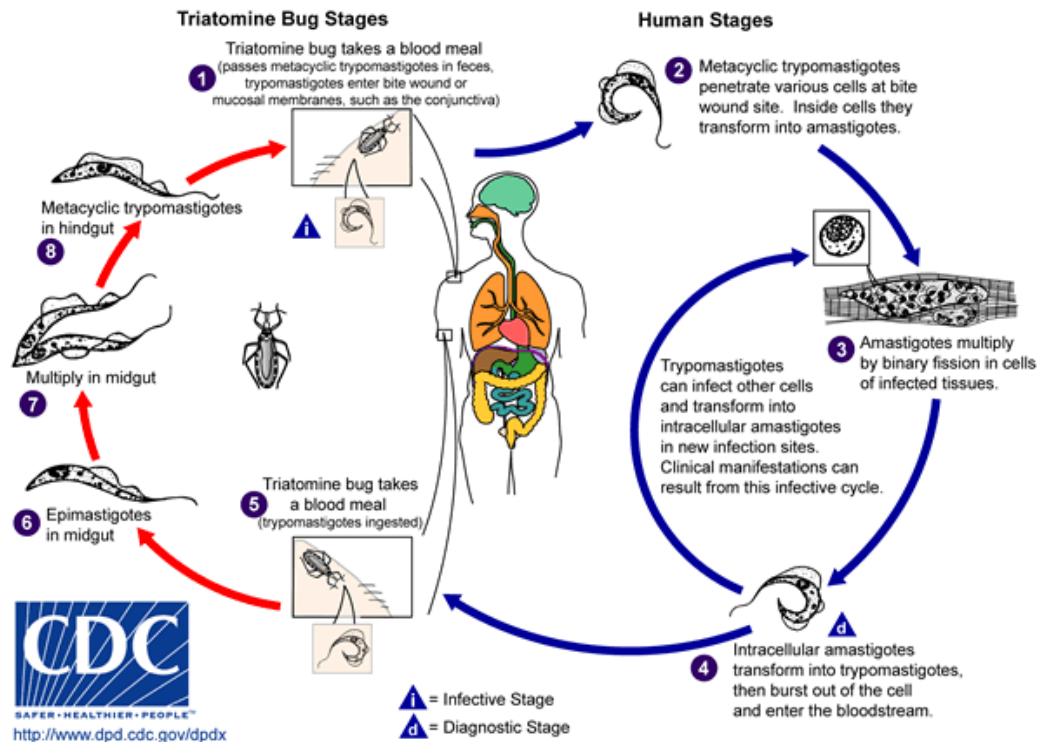


- c. Which gene has 16 exons? (Hint: add a column for number of exons)
- d. Is this gene alternatively spliced? Look at the gene page. Take note of the Gene ID.
- e. View this gene in the genome browser and load the RNA-seq tracks for this experiment “P.f. mRNAseq Transcriptome during intraerythrocytic development Coverage (Bartfai *et al.*)”. Do these tracks match the results you got above? (ie. is this gene differentially regulated between the early time points and the late ones?)
- f. 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)”).
- g. 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).



3. Exploring microarray data in TriTrypDB.

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



Find *T. cruzi* genes that are upregulated in amastigotes compared to trypomastigotes. Go to the transcript expression section then select microarray.

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
- ☐ Protein Expression
- ☐ Cellular Location
- ☐ Putative Function
- ☐ Evolution
- ☐ Population Biology

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)	P
<i>L. major</i> strain Friedlin	Three Developmental Stages (Stephen M. Beverley)	DC 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
<i>T. brucei</i> TREU927	Expression profiling of in vitro differentiation time series (Christine Clayton)	FC
<i>T. brucei</i> TREU927	induced DHH1 in wild type and DEAD:DQAD mutant (Mark Carrington)	P
<i>T. brucei</i> TREU927	Procyclic trypanosomes treated with heat shock (Mark Carrington)	DC P
<i>T. cruzi</i> CL Brener Esmeraldo-like	Life-Cycle Stages (Rick Tarleton)	FC P

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

Identify Genes based on T.c. Life-Cycle Stages Microarray (fold change) [Tutorial](#) [YouTube](#)

For the Experiment **Life-Cycle Stages tcruCLBrennerEsmeraldo-like**

return genes that are **up-regulated**

with a Fold change \geq **2**

between each gene's expression value

in the following **Reference Samples**

☐ amastigotes
☒ trypomastigotes
☐ epimastigotes
☐ metacyclics

[select all](#) | [clear all](#)

and its expression value

in the following **Comparison Samples**

☒ amastigotes
☐ trypomastigotes
☐ epimastigotes
☐ metacyclics

[select all](#) | [clear all](#)

Protein Coding Only: ☐ protein coding

[Advanced Parameters](#)

[Get Answer](#)

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

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

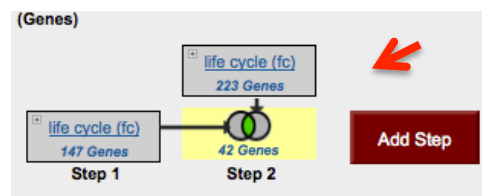
For each gene, the search calculates:

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

and returns genes when fold change \geq 2.

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

- How many genes did you find? Do the results seem plausible?
- Are any of these genes also upregulated in the replicative insect stage (epimastigotes)? 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 run an ortholog transform on your results).
- How many orthologs exist in *L. braziliensis*? (*Hint*: look at the filter table right above your results. Click on the number in of gene to view results from a specific species).



My Strategies: New Opened (1) All (3) Basket Examples Help

(Genes) Strategy: life cycle (fc) *

life cycle (fc) 223 Genes

life cycle (fc) 147 Genes Step 1

42 Genes Step 2

Orthologs 55 Genes Step 3

Add Step

55 Genes from Step 3 Strategy: life cycle (fc)

Add 55 Genes to Basket | Download 55 Genes

Filter results by species (results removed by the filter will not be combined into the next step.)

All Results	Ortholog Groups	Leishmania							Trypanosoma brucei				Trypanosoma cruzi						Trypanosoma vivax
		braziliensis	donovani	infantum	major	mexicana	tarentolae	Distinct genes	TREU927	strain 427	gambiense	Trypanosoma congolense	Distinct genes	esmeraldo	non-esmeraldo	unassigned	marinkellei	Sylvio	
1523	37	55	52	57	59	56	59	38	38	36	36	27	1017	330	316	194	94	83	31

Gene Results Genome View

First 1 2 3 Next Last Advanced Paging Add Columns

Gene ID	Organism	Genomic Location	Product Description	Input Ortholog(s)	Ortholog Group	Paralog count	Ortholog count
LbrM.02.0350	L. braziliensis MHOM/BR/75/M2904	LbrM.02: 147,781 - 154,645 (-)	ABC1 transporter, putative	TcCLB.510149.80	OG5_126568	8	76
LbrM.11.0960	L. braziliensis MHOM/BR/75/M2904	LbrM.11: 439,107 - 444,425 (+)	ABC transporter, putative	TcCLB.510149.80	OG5_126568	8	76
LbrM.11.1000	L. braziliensis MHOM/BR/75/M2904	LbrM.11: 458,406 - 464,144 (+)	ABC1 transporter, putative	TcCLB.510149.80	OG5_126568	8	76
LbrM.11.1010	L. braziliensis MHOM/BR/75/M2904	LbrM.11: 470,736 - 476,186 (+)	ABC1 transporter, putative	TcCLB.510149.80	OG5_126568	8	76

- Explore your results. Did you find anything interesting?

4. Finding genes based on RNAseq evidence and inferring function of hypothetical genes.

Note: Use <http://plasmodb.org> for this exercise.

- Find all genes in *P. falciparum* that are upregulated at least 50-fold in ookinetes compared to other stages: “Transcriptomes of 7 sexual and asexual life stages (Lopez-Barragan et al.)”

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

Organism **Data Set** **Choose a search**

P. falciparum 3D7	Transcriptome during intraerythrocytic development (Bartfai et al.)	FC FcPv P
P. falciparum 3D7	Transcriptomes of 7 sexual and asexual life stages (Lopez-Barragan et al.)	FC FcPv P
P. falciparum 3D7	Strand specific transcriptomes of 4 life cycle stages (Lopez-Barragan et al.)	FC FcPv P
P. falciparum 3D7	NSR-seq Transcript Profiling of malaria-infected pregnant women and children (Vignali et al.)	FC FcPv P
P. falciparum 3D7	Blood sta	FC FcPv P
P. yoelii yoelii 17XNL	Salivary	FC FcPv P

Revise Step 2 : P.f. seven stages - RNA Seq (percentile)

Experiment: percentile - P. falciparum Su Seven Stages RNA Seq data

Samples:

- ☒ Ring
- ☒ Early Trophozoite
- ☒ Late Trophozoite
- ☒ Schizont
- ☒ Gametocyte II
- ☒ Gametocyte V
- ☐ Ookinete

select all | clear all

Minimum expression percentile: 20

Maximum expression percentile: 50

Matches Any or All Selected Samples?: any

Protein Coding Only?: yes

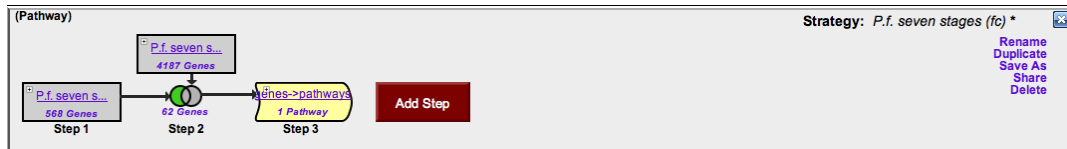
Advanced Parameters

Combine Genes in Step 1 with Genes in Step 2:

- ☐ 1 Intersect 2
- ☐ 1 Union 2
- ☐ 1 Relative to 2, using genomic colocation
- ☐ 1 Minus 2
- ☐ 2 Minus 1

b.

- c. The above search will give you all genes that are upregulated by 50 fold in ookinetes compared to the other stages. However, this does not mean that these genes are not expressed in the other stages. How can you remove genes from the list that are likely not expressed in the other stages? (*hint: run a search for genes based on RNAseq evidence from the same experiment, but this time select the percentile search*): **P.f. seven stages - RNA Seq (percentile)**
- d. Which metabolic pathways are represented in this gene list? (*hint: transform results to metabolic pathways*).



4. Finding genes that are essential in procyclics but not in blood form *T. brucei*.
Note: for this exercise use <http://TriTrypDB.org>.

- Find the query for high throughput Phenotyping.

Identify Genes by:

Expand All | Collapse All

- ☐ Text, IDs, Organism
- ☐ Genomic Position
- ☐ Gene Attributes
- ☐ Protein Attributes
- ☐ Protein Features
- ☐ Similarity/Pattern
- ☐ Transcript Expression
- ☐ Protein Expression
- ☐ Cellular Location
- ☐ Putative Function
- ☐ GO Term
- ☐ EC Number
- ☐ Phenotype
- ☐ High-Throughput Phenotyping
- ☐ Evolution
- ☐ Population Biology

Identify Genes based on High-Throughput Phenotyping

Experiment ☒ Quantitated from the CDS Sequence
☐ Quantitated from gene model (5 prime UTR + CDS)

Direction

Reference Sample(s) ☒ Uninduced sample

Comparison Sample(s) ☐ Induced bloodstream form (day 3)
☐ Induced bloodstream form (day 6)
☐ Induced procyclics
☐ DIF (induced throughout growth) form[†]
[select all](#) | [clear all](#)

fold difference

P value less than or equal to

Apply to Any or All Selected Samples?

Protein Coding Only:

- Think about how to set up this query. (*hint: you will have to setup a two step strategy*).
- Remember you can play around with the parameters but there is no one correct way of setting them up – try the default parameters first and select the “induced procyclics” as the comparison sample.

Identify Genes based on High-Throughput Phenotyping

Experiment ☒ Quantitated from the CDS Sequence
☐ Quantitated from gene model (5 prime UTR + CDS)

Direction Decrease in coverage

Reference Sample(s) ☒ Uninduced sample

Comparison Sample(s) ☐ Induced bloodstream form (day 3)
☐ Induced bloodstream form (day 6)
☒ Induced procyclics
☐ DIF (induced throughout growth) form'
[select all](#) | [clear all](#)

fold difference 1.5

P value less than or equal to 1E-6

Apply to Any or All Selected Samples? any

Protein Coding Only: yes

My Strategies: New Opened (1) All (1) Basket Examples Help

Add Step

Strategy: *T.b. RNAi fc* *
[Rename](#)
[Duplicate](#)
[Save As](#)
[Share](#)
[Delete](#)

- 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.

Add Step 2 : High-Throughput Phenotyping

Experiment ☒ Quantitated from the CDS Sequence
☐ Quantitated from gene model (5 prime UTR + CDS)

Direction Decrease in coverage

Reference Sample(s) ☒ Uninduced sample

Comparison Sample(s) ☒ Induced bloodstream form (day 3)
☒ Induced bloodstream form (day 6)
☐ Induced procyclics
☐ DIF (induced throughout growth) form'
[select all](#) | [clear all](#)

fold difference 1.5

P value less than or equal to 1E-6

Apply to Any or All Selected Samples? any

Protein Coding Only: yes

Advanced Parameters

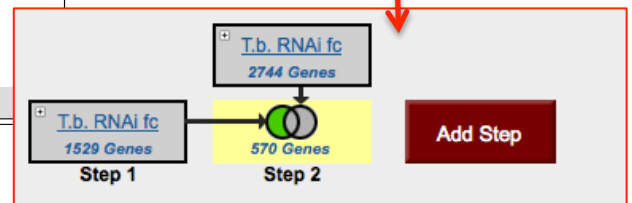
Combine Genes in Step 1 with Genes in Step 2:

☐ 1 Intersect 2
 ☒ 1 Minus 2

☐ 1 Union 2
 ☐ 2 Minus 1

☐ 1 Relative to 2, using genomic colocation

Run Step




5. Exploring Expression Quantitative Trait Locus (eQTL) data in PlasmoDB.

Genetic crosses were instrumental in implicating the PfCRT gene in chloroquine resistance. PlasmoDB contains expression quantitative trait locus data from Gonzales *et. al.* PLoS Biol 6(9): e238. The trait that was examined in this study was gene expression using microarray experiments.

- Go to the gene page for the gene with the ID PF3D7_0630200. Can you identify the genomic region (haplotype block) that is “most” associated with this gene, ie. has the highest LOD score? (Hint: examine the table called “Regions/Spans associated by eQTL experiment on HB3 x DD2 progeny” on the gene page.

Regions/Spans associated by eQTL experiment on HB3 x DD2 progeny (LOD cut off = 1.5) [Hide](#)



Haplotype Block	Genomic Segment (Liberal)	Genomic Segment (Conservative)	LOD Score (opens a haplotype plot)	Search for Genes (Liberal by Default)	Search for Genes (Liberal by Default)
PF3D7_05_v3_68.8	PF3D7_05_v3:1010972-1040241	PF3D7_05_v3:1018620-1018825	4.94	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_68.8	PF3D7_05_v3:959929-1010786	PF3D7_05_v3:1007897-1008018	4.94	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_65.9	PF3D7_05_v3:870388-1007896	PF3D7_05_v3:918503-959928	4.9	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_25.8	PF3D7_05_v3:389050-493947	PF3D7_05_v3:398963-405946	3.29	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_48.7	PF3D7_05_v3:683733-732922	PF3D7_05_v3:686437-693079	3.2	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_45.8	PF3D7_05_v3:628981-686436	PF3D7_05_v3:683548-683732	3.2	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_42.9	PF3D7_05_v3:555274-683547	PF3D7_05_v3:628753-628980	3.2	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_31.5	PF3D7_05_v3:405947-628752	PF3D7_05_v3:493948-55273	2.99	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_20	PF3D7_05_v3:260855-355367	PF3D7_05_v3:304284-325885	2.87	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_22.9	PF3D7_05_v3:325886-398962	PF3D7_05_v3:355368-389049	2.81	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_60.2	PF3D7_05_v3:770125-918502	PF3D7_05_v3:814427-870387	2.18	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_54.4	PF3D7_05_v3:693080-769886	PF3D7_05_v3:732923-733046	2.15	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_11.4	PF3D7_05_v3:252443-304283	PF3D7_05_v3:260710-260854	2.14	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_5.7	PF3D7_05_v3:166792-260709	PF3D7_05_v3:225881-252442	2.13	Genes Contained in this Region	Genes Associated to this Region
PF3D7_08_v3_57.5	PF3D7_08_v3:408724-684033	PF3D7_08_v3:570281-647334	2.11	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_28.9	PF3D7_07_v3:496401-694858	PF3D7_07_v3:611138-611341	1.98	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_57.3	PF3D7_05_v3:733047-814426	PF3D7_05_v3:769887-770124	1.98	Genes Contained in this Region	Genes Associated to this Region
PF3D7_08_v3_40.3	PF3D7_08_v3:768381-783997	PF3D7_08_v3:768494-768653	1.97	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_20.2	PF3D7_07_v3:391071-427528	PF3D7_07_v3:392209-425264	1.79	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_17.3	PF3D7_07_v3:371129-392208	PF3D7_07_v3:377646-391070	1.69	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_0	PF3D7_05_v3:86612-225880	PF3D7_05_v3:140933-166791	1.67	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_26	PF3D7_07_v3:451719-611137	PF3D7_07_v3:463358-496400	1.65	Genes Contained in this Region	Genes Associated to this Region
PF3D7_08_v3_91.8	PF3D7_08_v3:1-230964	PF3D7_08_v3:122068-122241	1.64	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_23.1	PF3D7_07_v3:425265-463357	PF3D7_07_v3:427529-451718	1.64	Genes Contained in this Region	Genes Associated to this Region
PF3D7_08_v3_48.9	PF3D7_08_v3:647335-751204	PF3D7_08_v3:684034-725296	1.6	Genes Contained in this Region	Genes Associated to this Region
PF3D7_07_v3_14.4	PF3D7_07_v3:358161-377645	PF3D7_07_v3:370990-371128	1.57	Genes Contained in this Region	Genes Associated to this Region
PF3D7_05_v3_83.1	PF3D7_05_v3:1018826-1095899	PF3D7_05_v3:1040242-1045759	1.53	Genes Contained in this Region	Genes Associated to this Region

Other genes that have similar associations based on eQTL experiments

- What kinds of genes do you find in this region? Click on the first link in the column “Genomic segment (liberal)”. Now examine the gene table on the genomic segment page.

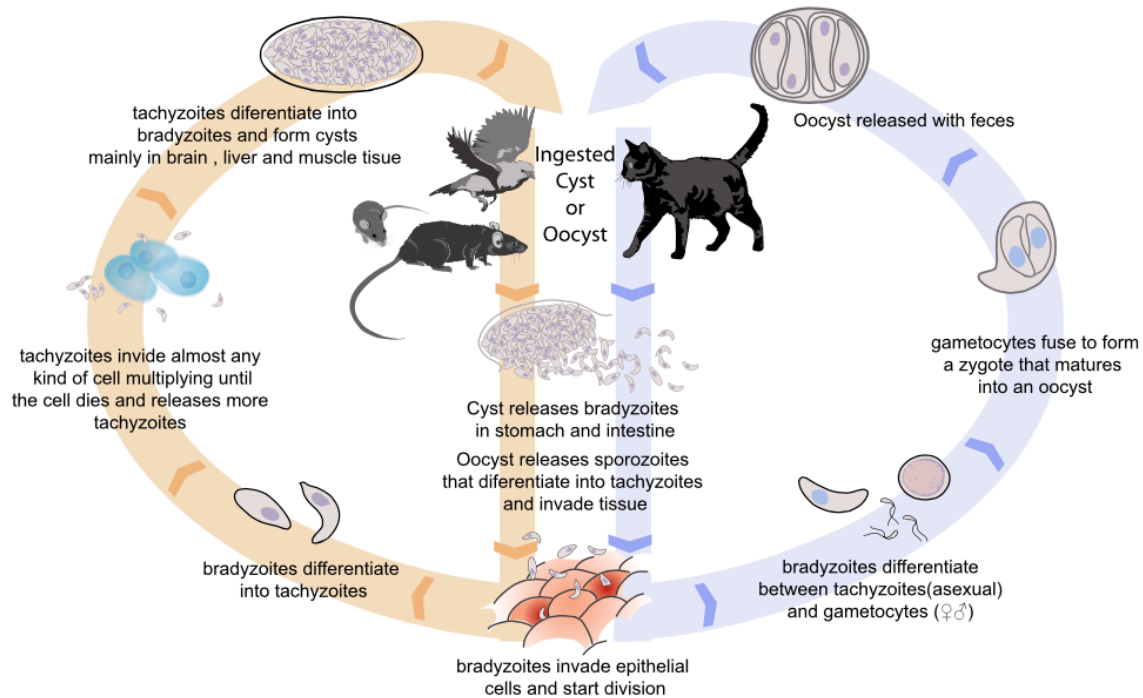
Genes [Hide](#)

Gene ID	Start	End	Strand	Product Description
PF3D7_0523000	957890	962149	forward	multidrug resistance protein (MDR1)
PF3D7_0523100	963227	965044	reverse	mitochondrial processing peptidase alpha subunit, putative
PF3D7_0523200	966123	969737	forward	conserved Plasmodium protein, unknown function
PF3D7_0523300	970266	970962	reverse	conserved Plasmodium protein, unknown function
PF3D7_0523400	973518	975876	forward	DnaJ protein, putative
PF3D7_0523500	976690	977815	reverse	outer arm dynein lc3, putative
PF3D7_0523600	978665	979870	forward	conserved Plasmodium protein, unknown function
PF3D7_0523700	980754	985354	reverse	conserved Plasmodium membrane protein, unknown function
PF3D7_0523800	990005	992059	forward	transporter, putative
PF3D7_0523900	993433	994607	reverse	conserved Plasmodium membrane protein, unknown function
PF3D7_0524000	998753	1002124	forward	karyopherin beta (KASbeta)
PF3D7_0524100	1004237	1008108	forward	conserved Plasmodium protein, unknown function
PF3D7_0524200	1008636	1009404	reverse	conserved Plasmodium membrane protein, unknown function

- c. What other genes are associated with this block?
(Hint: go back to the gene page eQTL table, and click the “genes associated with this region” link. Run the search on the next page and examine the list of genes. It might be useful to sort this list based on the LOD scores.)

6. 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 “Expression Profiling of *T. gondii* Oocyst, Tachyzoite, Bradyzoite stages (John Boothroyd)” microarray experiment.

Identify Genes by:

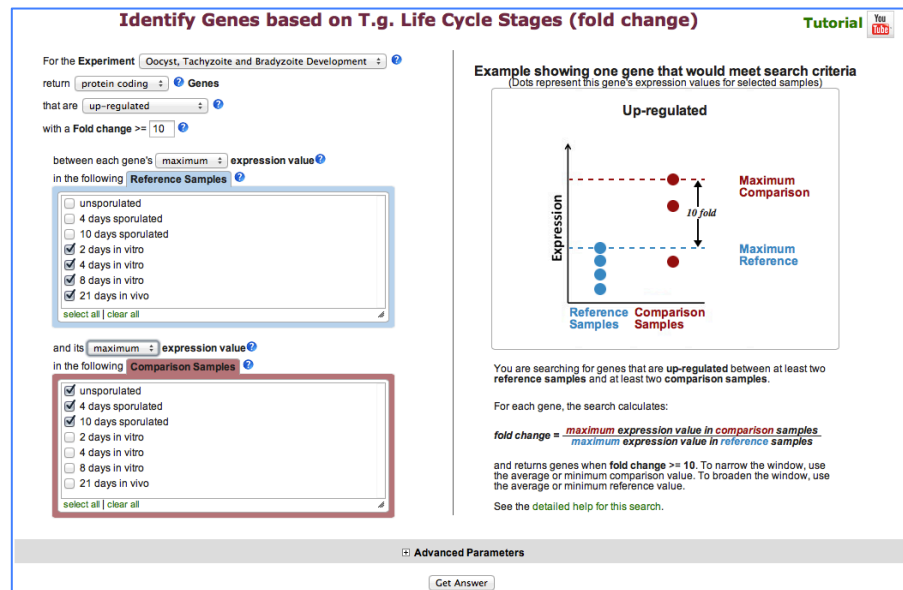
- Expand All | Collapse All
- Text, IDs, Organism
- Genomic Position
- Gene Attributes
- Protein Attributes
- Protein Features
- Similarity/Pattern
- Transcript Expression
- EST Evidence
- Microarray Evidence**
- ChIP on Chip Evidence
- Protein Expression
- Cellular Location
- Putative Function
- Evolution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

Identify Genes based on Microarray Evidence

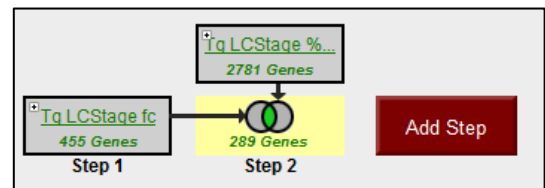
Filter Data Sets: Legend: FC Fold Change FCC Fold Change wit... P Percentile S Similarity

Organism	Data Set	Choose a search
<i>T. gondii</i> ME49	Differential Expression Profiling GCN5-A mutant (William Sullivan)	FC FCC P
<i>T. gondii</i> ME49	Bradyzoite Differentiation (Multiple 6-hr time points and Extended time series) (Paul H. Davis)	FC P
<i>T. gondii</i> ME49	Expression profiling of the 3 archetypal <i>T. gondii</i> lineages (David S. Roos)	FCC P
<i>T. gondii</i> ME49	Transcript Profiling Infection (Vern B. Carruthers)	FC FCC P
<i>T. gondii</i> ME49	Mutants and wild-type during bradyzoite differentiation in vitro (Mariana Matrajt)	FC FCC P
<i>T. gondii</i> ME49	Bradyzoite Differentiation (Single Time-Point) (Michael W White)	P
<i>T. gondii</i> ME49	Cell Cycle Expression Profiles (Michael W White)	FC P S
<i>T. gondii</i> ME49	Expression Profiling of oocyst, tachyzoite, and bradyzoite development in strain M4 (John Boothroyd)	FC P

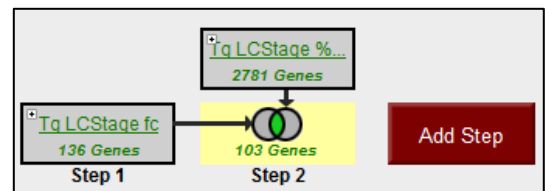
In this example the maximum expression value between genes in the reference and comparison groups was used to determine the fold difference.



- b. Add a step to limit this set of genes to only those for which all the non-oocyst stages are expressed below 50th 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 50th 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**”.
 - Note: you can turn on the columns called “Tg-M4 Life Cycle Stages – graph” and “Tg-M4 Life Cycle Stage %ile- graph” to view the graphs in the final result table.



- c. Revise the first step of this strategy and compare the maximum expression of the reference samples to the minimum of the comparison samples.
- Does this result look cleaner/more convincing? Why?
 - Would you consider these genes to be oocyst specific?
 - **Save this strategy as we'll use this strategy for an exercise we are doing later during the course.**



- d. Revise the first step of this strategy to find genes that are 3 fold higher in day 4 oocysts than any other life cycle stage in this experiment.
- Do all these genes have day 4 oocysts as the global maximum time point?
 - Note that we still have the step to limit the percentile of non-oocyst samples to $\leq 50^{\text{th}}$ percentile. What happens if you revise this step to also include the unsporulated and day 10 oocyst samples in this percentile range? Do you get more or fewer results back? Why?

My Strategies: [New](#) [Opened \(1\)](#) [All \(1\)](#) [Basket](#) [Examples](#) [Help](#)

Strategy: **Tg LCStage fc***

(Genes)

Tg LCStage %tilt
1920 Genes

Tg LCStage fc
67 Genes
Step 1

4 Genes
Step 2

[Add Step](#)

[Rename](#)
[Duplicate](#)
[Save As](#)
[Share](#)
[Delete](#)

4 Genes from Step 2
Strategy: **Tg LCStage fc**

[Add 4 Genes to Basket](#) | [Download 4 Genes](#)

[Filter by organism or strain](#) (results removed by the filter will not be combined into the next step.)
[Filter by strains \(advanced\)](#) (results removed by the filter will not be combined into the next step.)

[Gene Results](#) [Genome View](#)

Advanced Paging [Add Columns](#)

Gene ID	Gene Group (representative gene)	Genomic Location	Product Description	Tg-M4 Life Cycle Stages - graph	Tg-M4 Life Cycle Stage %ile - graph
TGME49_258800	TGGT1_258800	TGME49_chrVIIb: 3,177,135 - 3,178,728 (+)	roptry kinase family protein ROP31 (ROP31)		
TGME49_233300	TGGT1_233300	TGME49_chrVIII: 2,569,523 - 2,577,098 (-)	RhoGAP domain-containing protein		

7. Comparing RNA abundance and Protein abundance data.

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

In this exercise we want to 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.

- Go to the genes by microarray expression and select the fold change search for the “Expression profiling of five life cycle stages (Marilyn Parsons)” experiment.

The image shows the TriTrypDB interface. On the left, the 'Identify Genes by:' sidebar has 'Microarray Evidence' highlighted with a red box and a red arrow pointing to the main panel. The main panel, titled 'Identify Genes based on Microarray Evidence', has a 'Filter Data Sets' field and a 'Legend' with 'DC' (Direct Comparison), 'FC' (Fold Change), and 'P' (Percentile). Below is a table of data sets:

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 Papadopoulou)	P
<i>L. major</i> strain Friedlin	Three Developmental Stages (Stephen M. Beverley)	DC 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
<i>T. brucei</i> TREU927	Expression profiling of in vitro differentiation time series (Christine Clayton)	FC
<i>T. brucei</i> TREU927	induced DHH1 in wild type and DEAD:DQAD mutant (Mark Carrington)	P
<i>T. brucei</i> TREU927	Procyclic trypanosomes treated with heat shock (Mark Carrington)	DC P
<i>T. cruzi</i> CL Brener Esmeraldo-like	Life-Cycle Stages (Rick Tarleton)	FC P

Configure the search to return protein-coding genes that are down-regulated 2 fold in procyclic form (PCF) (I chose both log and Stat and averaged them) relative to the Blood Form reference sample.

The image shows the search configuration page for 'Identify Genes based on T.b. Expression profiling of five life cycle stages Microarray (fold change)'. The 'For the Experiment' dropdown is set to 'Expression profiling of five life cycle stages'. The 'return genes that are' dropdown is set to 'down-regulated'. The 'with a Fold change >= 2' dropdown is set to '2'. The 'between each gene's average expression value' dropdown is set to 'average'. The 'in the following Reference Samples' dropdown is set to 'Blood Form'. The 'and to average expression value' dropdown is set to 'average'. The 'in the following Comparison Samples' dropdown is set to 'PCF Log'. The 'Protein Coding Only' dropdown is set to 'protein coding'. The 'Advanced Parameters' section is expanded, showing the 'Get Answer' button.

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

Down-regulated

Expression

Average Reference

Average Comparison

Reference Comparison 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:

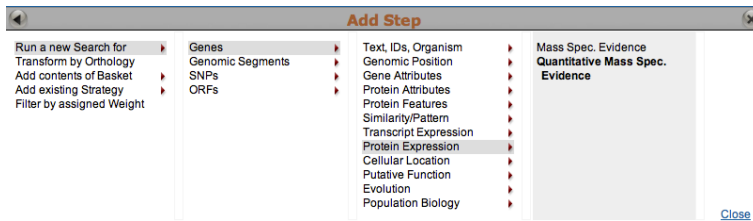
$$\text{fold change} = \frac{\text{average expression value in reference samples}}{\text{average expression value in comparison samples}}$$

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.

- Add a step to compare with quantitative protein expression. Select protein expression then “Quantitative Mass Spec Evidence”. Configure this search to return genes that are downregulated in procyclic form relative to Blood form.

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

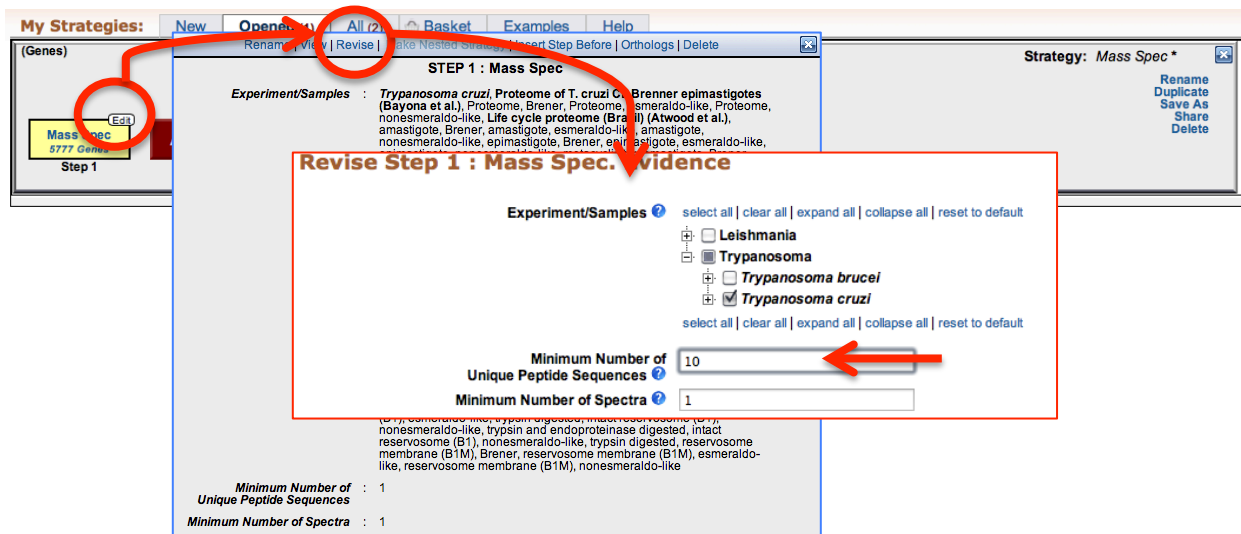


- d. Try changing directions and compare up-regulated genes/proteins. (*hint, revise the existing strategy ... you might want to duplicate it so you can keep both*). When you change one of the steps but not the other do you have any genes in the intersection? Why might this be??
- e. Can you think of 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 RNASequencing step that has the same samples*).

8. Finding all genes with mass spec evidence in *T. cruzi*.

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

- a. How many genes in *T. cruzi* have mass spec evidence?
- b. How many genes from the results in a. have at least 10 unique peptide hits? (*hint: try revising the step in 'a' and change the "minimum number of unique peptide sequences" option to 10.*



- c. Can you expand the list of results in 'b' to include possible orthologs/paralogs in *T. cruzi*?

Hint: you will have to use the ortholog transform option when adding a step and select only *T. cruzi*. Explore the columns in your result set.

My Strategies: New Opened (1) All (2) Basket Examples Help

(Genes) Strategy: Mass Spec *
Rename Duplicate Save As Share Delete

Mass Spec 634 Genes Step 1 → Orthologs 3124 Genes Step 2 Add Step

3124 Genes from Step 2 Strategy: Mass Spec Add 3124 Genes to Basket | Download 3124 Genes

Filter results by species (results removed by the filter will not be combined into the next step.)

All Results	Ortholog Groups	Leishmania						Trypanosoma brucei			Trypanosoma cruzi				Trypanosoma					
		braziliensis	donovani	infantum	major	mexicana	tarentolae	Distinct genes	TREU927	strain 427	gambiense	Distinct genes	esmeraldo	non-esmeraldo	unassigned	marinkellei	Sylvio	cogolense	evansi	vivax
3124	376	0	0	0	0	0	0	0	0	0	0	3113	637	690	207	613	977	0	0	0

Gene Results Genome View

First 1 2 3 4 5 Next Last Advanced Paging Add Columns

Gene ID	Organism	Genomic Location	Product Description	Input Ortholog(s)	Ortholog Group	Paralog count	Ortholog count
TCSYLVIQ_000024	<i>T. cruzi</i> Sylvio X10/1	ADWP02000075: 3 - 665 (-)	retrotransposon hot spot (RHS) protein, putative	TcCLB.410923.20, TcCLB.459199.10, TcCLB.463155.20, TcCLB.503483.9, TcCLB.503607.4, TcCLB.503861.10, ...	OG5_126555	443	772
TCSYLVIQ_000061	<i>T. cruzi</i> Sylvio X10/1	ADWP02000144: 81 - 671 (+)	retrotransposon hot spot (RHS) protein, putative	TcCLB.410923.20, TcCLB.459199.10, TcCLB.463155.20, TcCLB.503483.9, TcCLB.503607.4, TcCLB.503861.10, ...	OG5_126555	443	772
TCSYLVIQ_000111	<i>T. cruzi</i> Sylvio X10/1	ADWP02000314: 1 - 663 (+)	retrotransposon hot spot (RHS) protein, putative	TcCLB.410923.20, TcCLB.459199.10, TcCLB.463155.20, TcCLB.503483.9, TcCLB.503607.4, TcCLB.503861.10, ...	OG5_126555	443	772
TCSYLVIQ_000114	<i>T. cruzi</i> Sylvio X10/1	ADWP02000331: 2 - 661 (+)	retrotransposon hot spot (RHS) protein, putative	TcCLB.410923.20, TcCLB.459199.10, TcCLB.463155.20, TcCLB.503483.9, TcCLB.503607.4, TcCLB.503861.10, ...	OG5_126555	443	772
TCSYLVIQ_000134	<i>T. cruzi</i> Sylvio X10/1	ADWP02000370: 140 - 1,198 (+)	retrotransposon hot spot (RHS) protein, putative	TcCLB.410923.20, TcCLB.459199.10, TcCLB.463155.20, TcCLB.503483.9, ...	OG5_126555	443	772

9. Finding genes with mass spec evidence in *P. berghei* gametocytes.

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

- a. Find all *P. berghei* genes that have mass spec evidence in either or both male and female gametocytes.

(hint: mass spec searches are in the "protein expression" expression section. Either or both is the Union of both results, not the intersection).

Identify Genes by:

Expand All | Collapse All

- ☐ Text, IDs, Organism
- ☐ Genomic Position
- ☐ Gene Attributes
- ☐ Protein Attributes
- ☐ Protein Features
- ☐ Similarity/Pattern
- ☐ Transcript Expression
- ☐ Protein Expression
- ☒ Mass Spec. Evidence
- ☐ Cellular Location

- How many genes did you get? How did you get to this number?
- Try running this search in two different ways:

- i. Select both male and female gametocyte options and run the search.
 - ii. Select one of them first, run the search then add the other one using the add step button. How did you combine the two steps? Do you get the same results as in (i)?
- b. **Find all genes that have mass spec evidence in both male and female gametocytes.**
(*hint*: use the strategy you developed in (ii) to get this answer, but change the union into an intersection).
- c. **Find genes that have mass spec evidence only in male gametocytes and not in female ones.**
(*hint*: modify the set operation in b).
- d. **Find genes that have mass spec evidence only in female gametocytes and not in male ones.**
(*hint*: modify the set operation in b).
- e. **Which female gametocyte gene has the highest number of peptide sequences?**
(*hint*: look at the “number of peptide sequences” column in the list of results).
- f. **What does the distribution of peptides in the gene from ‘e’ look like?**
(*hint*: go to the gene page and look at the “Protein features” section, or go to the genome browser from the gene page and turn on the right tracks).

10. Finding genes with evidence of phosphorylation in intracellular *Toxoplasma* tachyzoites.

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

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

- a. 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](#)

☒ **Toxoplasma**

☒ ***Toxoplasma gondii***

☒ **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)

☐ Oocyst proteome (M4 type II) (Fritz et al.)

☐ Tachyzoite secretome (RH) (Zhou et al.)

☐ Tachyzoite membrane and cytosolic fractions (RH) (Dybas et al.)

☐ Tachyzoite subcellular fractions (Moreno)

☐ Tachyzoite conoid proteome (RH) (Hu et al.)

☐ Intra- and Extracellular Tachyzoite Lysine-Acetylomes (RH) (Jeffers and Xue)

☐ Calcium dependent tachyzoite phosphoproteome (RH) (Nebl et al.)

☐ Tachyzoite Rhoptyr proteome (RH) (Bradley et al.)

☐ Tachyzoite total proteome (RH) (Wastling)

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

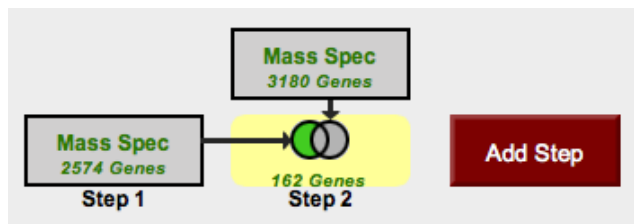
Minimum Number of Unique Peptide Sequences [?](#)

Minimum Number of Spectra [?](#)

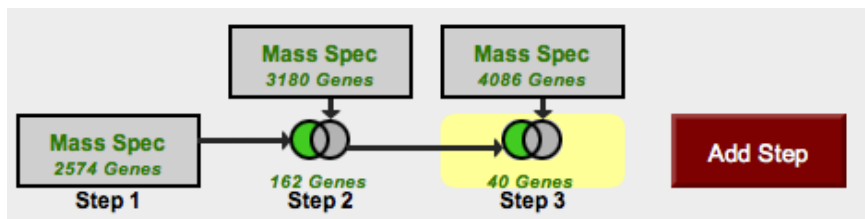
[Advanced Parameters](#)

[Get Answer](#)

- b. Remove all genes with phosphorylation evidence from purified extracellular tachyzoites.



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



d. Explore your results. What kinds of genes did you find? Are any of these results to be secreted? (Hint: add a step searching for genes with secretory signal peptides).

My Strategies: New Opened (1) All (2) Basket Examples Help

(Genes) Strategy: Mass Spec *

Mass Spec 2574 Genes Step 1 → Mass Spec 3180 Genes Step 2 → Mass Spec 4086 Genes Step 3 → Signal Pep 9872 Genes Step 4

9 Genes from Step 4 Strategy: Mass Spec Add 9 Genes to Basket | Download 9 Genes

Filter by organism or strain (results removed by the filter will not be combined into the next step.)

All Results	Ortholog Groups	Toxoplasma gondii						Neospora caninum	Eimeria tenella	
		All	Non-redundant	GT1	ME49	VEG	RH			
9	9		9	0	0	9	0	0	0	0

Filter by strains (advanced) (results removed by the filter will not be combined into the next step.)

Gene Results Genome View

Advanced Paging Add Columns

Gene ID	Product Description
TGME49_28880	hypothetical protein
TGME49_257568	hypothetical protein
TGME49_229680	hypothetical protein
TGME49_231180	hypothetical protein
TGME49_269420	hypothetical protein
TGME49_200440	hypothetical protein
TGME49_216840	hypothetical protein
TGME49_308070	hypothetical protein
TGME49_219640	hypothetical protein

Advanced Paging

e. 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).

f. What about polymorphism data? Go back to your strategy and add columns for SNP data found under the population section. Explore the gene page for the gene that has the most number of nonsynonymous SNPs.

Gene Results Genome View

Advanced Paging Add Columns

Gene ID	Product Description	Nonsynonymous SNPs All Strains	Synonymous SNPs All Strains	Total HTS SNPs All Strains	Total HTS Non-Synonymous SNPs	Total HTS Synonymous SNPs
TGME49_219640	hypothetical protein	33	60	383	81	302
TGME49_308070	hypothetical protein	20	34	188	42	146
TGME49_28880	hypothetical protein	17	27	221	52	169
TGME49_200440	hypothetical protein	14	13	72	36	36
TGME49_216840	hypothetical protein	13	42	189	71	118
TGME49_269420	hypothetical protein	9	27	45	31	14
TGME49_257568	hypothetical protein	5	16	30	20	10
TGME49_231180	hypothetical protein	3	15	54	24	30
TGME49_229680	hypothetical protein	0	6	33	2	31

Advanced Paging