# Exploring Transcriptomics Data
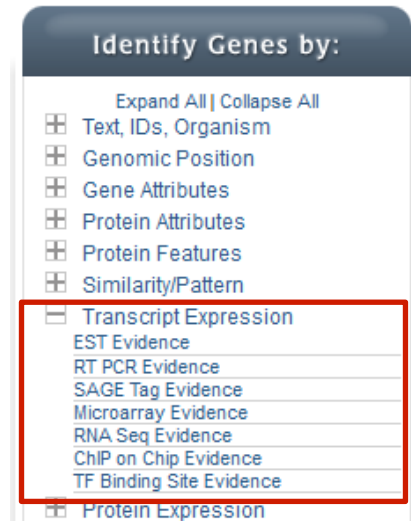## Exercise 13

## 13.1 Evidence of expression at the transcriptional level.
Note: For this exercise use http://www.eupathdb.org

a. 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.



b. Explore organisms that have microarray data. What organisms have expressed sequence tag (EST), RNA sequence, ChIP-chip or SAGE tag data?
c. What does RNA-seq data tell you that microarray data cannot? What does ChIP-chip data tell you about a gene?
d. 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.

## 13.2 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 "P.f. post infection (RBC) RNA-seq time series (fold change)".



> hint: there are several parameters to manipulate in this search:

**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.

**Global min / max in selected time points**: Choose whether the selected samples must contain the global minimum or maximum expression value of all samples. <mark>Choose maximum</mark> [Choosing minimum: if searching for up-regulated genes then the Reference values selected must be the global minimum. If searching for down-regulated genes, the comparison values selected must be the global minimum. Choosing maximum: if searching for up-regulated genes, then the Comparison values selected must be the global maximum. If searching for down-regulated genes, the reference values selected must be the global maximum.]



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 expr profile graph (GS array)" and compare the RNA-seq to the microarray graphs).

c. Which gene has 17 exons? (hint: use the columns)
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 "Pf RNA-Seq - intraerythrocytic cycle: 5-40hr post invasion (log2) [Stunnenberg lab]".  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 - Combined").
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).

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

1. Find all genes in *P. falciparum* that are upregulated at least 50-fold in ookinetes compared to other stages: P.f. seven stages - RNA Seq (fold change)



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 form 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.

## 13.3 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.

a.  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.



b.  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.)

## 13.4 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 (Boothroyd/Conrad)* microarray experiment.

- There are multiple parameters that need to be set.
- **Experiment**: choose <mark>Oocyst, Tachyzoite and Bradyzoite Development</mark>.
- **Direction**: choose <mark>down-regulated</mark> since we want to find things more highly expressed in oocysts than in other stages.
- Notice setting the Direction to down-regulated automatically changes the ***Operations Applied to Reference Samples*** from average to <mark>maximum</mark> and minimum for the comparator samples. This would enable you to find the genes with the maximum difference between these two sets of samples. Let's leave the reference set to maximum.
- **Reference Samples**: choose the <mark>three oocyst samples: (unsporulated, 4 days sporulated and 10 days sporulated</mark>.
- **Comparison Samples**: choose the <mark>4 non-oocyst samples: 2 days, 4 days, 8 days *in vitro*, and 21 days *in vivo*</mark>. (ie, tachyzoite and three bradyzoite samples)
- **Operation Applied to Comparison Samples:** choose <mark>maximum</mark> since the goal is to find genes with 10-fold higher expression in at least one of the oocyst samples compared to any of the non-oocyst samples.
- **Fold Change >= <mark>10</mark>**.
- **Global min/max in selected time points:** choose "<mark>don't care</mark>". Since we have selected all the samples between the reference and comparator time points, the global max and the global min will have to be within the selected time points. If we had not selected all the time points, then changing this parameter would make a difference as the global min or max could be in a time point that we didn't select.
- **Protein coding only** as <mark>yes</mark>. We want to only look at polyadenylated transcripts.

**Identify Genes based on T.g. Life Cycle Stages (fold change)**

Experiment ❓ ◉ Oocyst, Tachyzoite and Bradyzoite Development

Direction ❓ down-regulated ▾

Reference Samples ❓
- ☑ oocyst - d0 unsporulated
- ☑ oocyst - d4 sporulation
- ☑ oocyst - d10 sporulation
- ☐ tachyzoite - d2 in vitro
- ☐ bradyzoite - d4 in vitro
- ☐ bradyzoite - d8 in vitro
- ☐ bradyzoite - d21 IN VIVO
  select all | clear all

Operation Applied to Reference Samples ❓ maximum ▾

Comparison Samples ❓
- ☐ oocyst - d0 unsporulated
- ☐ oocyst - d4 sporulation
- ☐ oocyst - d10 sporulation
- ☑ tachyzoite - d2 in vitro
- ☑ bradyzoite - d4 in vitro
- ☑ bradyzoite - d8 in vitro
- ☑ bradyzoite - d21 IN VIVO
  select all | clear all

Operation Applied to Comparison Samples ❓ maximum ▾

Fold change >= ❓ 10

Global min / max in selected time points ❓ Don't care ▾

Protein Coding Only: ❓ yes ▾

⊞ Give this search a weight

⊞ Give this search a name

Get Answer

b. Add a step 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: use the ***Expression Profiling of T. gondii Oocyst/Tachyzoite/Bradyzoite stages (str M4) (Boothroyd/Conrad)*** -> **T.g. Life Cycle Stages (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***".
- Set ***Protein Coding Only*** to YES.
- Note: you can turn on the column for "M4 Life Cycle Stages – graph" to see the graphs in the final result table.

c. Revise the first step of this strategy to find genes where all oocyst stages (d0, 4, 10) are 10 fold higher than any of the non-oocyst stages.
   - Hint, change the "***Operation applied to reference samples***" to <mark>minimum</mark>.
   - Does this result in cleaner, more convincing looking graphs?  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 this afternoon**.



d. Revise the first step of this strategy to find genes that are 3 fold higher in d4 oocysts than any other life cycle stage in this experiment.
   - Do all these genes have d4 oocysts as the global maximum time point?
   - Note that we still have the step to limit the percentile of non-oocyst samples to <= 50[th] percentile.  What happens if you revise this step to also include the d0 and d10 oocyst samples in this percentile range?  Do you get more of fewer results back?  Why?

### 13.5  Exploring EST evidence in *Entamoeba.*
   **Note: For this exercise use** http://amoebadb.org

a. Find all *Entamoeba* genes that have EST evidence.
b. Which gene has the highest number of ESTs?
c. Go to the following *E. dispar* gene in the genome browser: **EDI_145400**
   Hint:  from the home page click on "Genome Browser" under the tools section, enter the ID in the "Landmark or Region" box.

d.  Look at the EST evidence for this gene.  Does it support the gene model?  Does the EST data tell you something else about the gene?
Hint: it would be easier to view this if you only have the gene model and EST tracks on, also you may have to zoom in or out to get a good view of the gene.



e.  Now, go to the following *E. histolytica* gene in the genome browser: **EHI_163570** (just like you did in step 'c' above).  What does the EST evidence look like?  Are there any ESTs that support an alternative gene model?  Look to the left of this gene. What do these ESTs mean?



### 13.6   Finding all ESTs that do not coincide with a gene model in *Entamoeba.*
   **Note: For this exercise use http://www.amoebadb.org**

a.  Find all ESTs that do not overlap with genes.  Hint:  Use the "Extent of Gene Overlap" search under the heading "Identify Other Data Types".
b.  How many ESTs did you get?  Hint: make sure you changed the **base overlap to 'O'** and selected "does not overlap with a gene".
c.  Visit one of the EST pages and explore it.  Can you get to the genome browser from here to see this EST? Hint: Look at the "Alignments to genomic sequence" section, click on "view".

### 13.7   Exploring genes expressed during encystation in *Giardia* based on SAGE tag evidence.
   **Note: For this exercise use http://www.giardiadb.org**

a. Find all genes on chromosome " GLCHR05". Hint: search for genes based on genomic location, under the menu "Genomic Position" in the "Identify Genes by" section.
b. How many of those genes have SAGE tag evidence during encystation? Hint: add a SAGE tag evidence (under Identify Genes By Transcript Expression) step using the following parameters: allow the tag to align 20 bp from either end and align to only one place in the genome. Find only genes with a tag count >= 5.
c. Do any of these genes have nucleic acid binding activity? Can you tell this from the product names? What other information can you add that would help? Hint: add a "Predicted GO Function" column to the list of your results.
d. How would you identify other genes in the *Giardia* that have a nucleic acid binding function and are expressed during encystation?



## 13.8  Exploring ChIP-chip data in *Toxoplasma*.
### Note:  For this exercise use http://www.toxodb.org

a. Use ChIP-chip searches to identify all genes that are differentially expressed between Type I and Type II strains of *T. gondii*.
How many genes are likely transcriptionally active in Type I but not Type II strains?
b. Go to the gbrowse view of one of those genes. Hint: use the column to add the gbrowse link to your list of results.
c. Turn on the tracks for ChIP-chip data.  Does the data agree with your search results?
d. Zoom out and explore neighboring genes.
e. Can you find centromeres?  Turn off all tracks except the ChIP-chip Centromeres track.  Hint: zoom out so that you have a view of the entire chromosome.
f. Once you locate the centromeric peaks, zoom in to explore what they look like.

g. Turn on one of the ChIP-chip graphs.  What do you notice?

**13.9   Exploring microarray data in TriTyrpDB.**
         **Note:  For this exercise use http://www.tritrypdb.org**

a. Find all genes in *T. brucei* that are up-regulated 48 hours (as compared to the 0 time point) post induction of differentiation (look for at least 4-fold induction).
   Hint: notice that there are two experiments in the "**T.b. differentiation time series (fold change)**" search.  You need to combine results from both experiments.  Also, in one of the experiments the "blood form – high density" is equal to the 0 time point. How did you combine the above two experiments?  Union or intersect?
b. How do these results compare to the RNA-seq data?  Can you view RNA-seq data for all of these results? Hint:  explore the columns in your result list.  Add the column that corresponds to the RNA-seq data.
c. Start a new search and look for genes that are greater than 4-fold down regulated in amasitgotes compared to metacyclics.  How many genes did you get?  Which gene is the most repressed gene? Hint: sort the "fold change avg." column in the list of results.
d. Go to the gene page of the top gene from step 'c'.  This gene is annotated as hypothetical – but can you get some hints to its possible function?  Hint: look at the protein feature section.

**13.10:   Comparing RNA abundance and Protein abundance data.   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.

   a. Go to the genes by microarray expression -> T.b. Expression profiling of five life cycle stages (fold change).  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.
   b. Add a step to compare with protein expression.  Genes by protein expression -> 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 (*it can be quite confusing* ☺).
   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.*