

Population data, SNPs and alleles

Exercise 4

4.1 Diversifying or purifying selection

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

- a. Find the *P. falciparum* genes that are the most diverse among sequenced strains and thus appear to be under diversifying selection.

Hint: for this exercise you need to use the “Genes by SNP Characteristics” search under “Identify Genes by:”.

- b. What strikes you about the known genes in this result set? What does the distribution of these genes on chromosomes look like? (hint: Click on the “Genome View” tab).

SNPs - Step 1 - 251 Genes Add 251 Genes to Basket | Download 251 Genes

Genes: Genome View (beta)

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

Gene ID	Product Description	Total SNPs	Non-synonymous SNPs	Synonymous SNPs	Nonsense SNPs	Non-coding SNPs	Non-syn/syn SNP ratio	SNPs per Kb (CDS)
PF3D7_1200600	erythrocyte membrane protein 1, PfEMP1 (VAR2CSA)	253	196	57	0	0	3.44	27.59
PF3D7_0402200	surface-associated	116	100	15	2	1	6.67	17.27

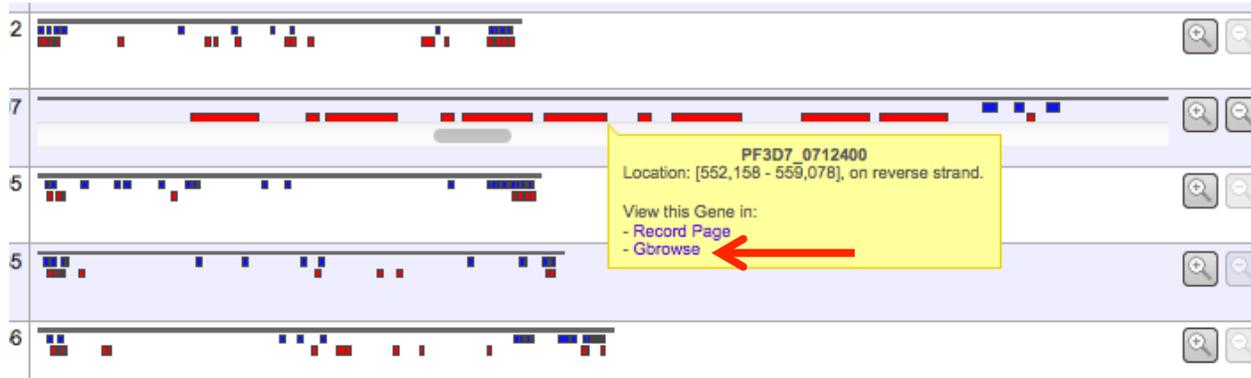
- c. Find all paralogs of these genes in *P. falciparum*. (hint: add an orthology transform step).
Now look at the distribution of the genes on chromosomes – do you notice a pattern?

Sequence	Organism	Chromosome	#Genes	Length	Gene Locations
PF3D7_03_v3	Plasmodium falciparum 3D7	3	30	1,067,971	[Gene locations on chromosome 3]
PF3D7_04_v3	Plasmodium falciparum 3D7	4	65	1,200,490	[Gene locations on chromosome 4]
PF3D7_05_v3	Plasmodium falciparum 3D7	5	30	1,343,557	[Gene locations on chromosome 5]
PF3D7_06_v3	Plasmodium falciparum 3D7	6	44	1,418,242	[Gene locations on chromosome 6]
PF3D7_07_v3	Plasmodium falciparum 3D7	7	62	1,445,207	[Gene locations on chromosome 7]
PF3D7_08_v3	Plasmodium falciparum 3D7	8	43	1,472,805	[Gene locations on chromosome 8]

- d. How many genes are in the middle cluster on chromosome 7? It might help to zoom in on this are (hint: use the magnifying glass icon on the left hand side).

PF3D7_06_v3	Plasmodium falciparum 3D7	6	44	1,418,242	[Gene locations on chromosome 6]
PF3D7_07_v3	Plasmodium falciparum 3D7	7	62	1,445,207	[Gene locations on chromosome 7]
PF3D7_08_v3	Plasmodium falciparum 3D7	8	43	1,472,805	[Gene locations on chromosome 8]
PF3D7_09_v3	Plasmodium	9	33	1,541,735	[Gene locations on chromosome 9]

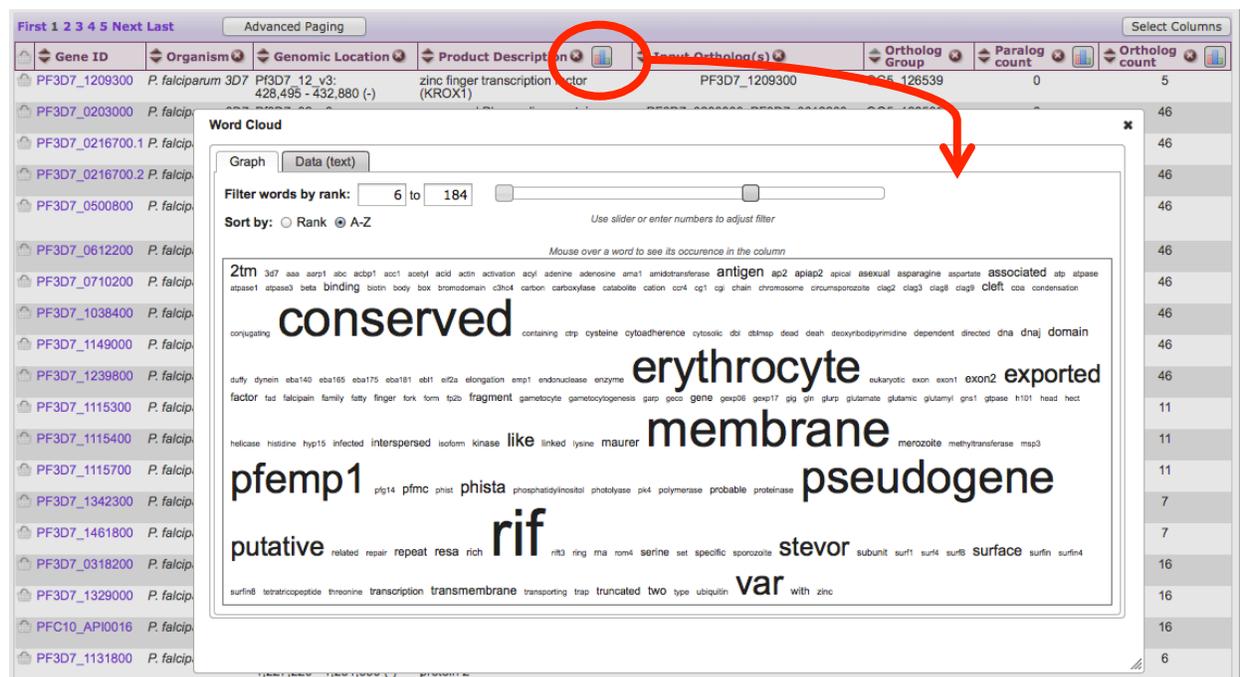
- e. What are these genes? Mouse over one of them then click on Gbrowse in the popup. When you get to the gbrowse view you can mouse over the genes in the region to see what they are:



f. Go back to your search strategy page (hint: you can easily get back to your strategies by clicking on the “My Strategies” link in the grey tool bar).



g. How can you quickly get an idea about which genes are represented in your results? (hint: click on the column analysis icon next to the product description). You can customize the filter and sorting of the resulting word cloud. You can also click on the “Data” tab in the popup to get a table of the results.



4.2 Isolate comparison

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

a. Go to the “Identify SNPs based on Isolate Comparison” search.

Hint: you can find this under “SNPs” in the “Identify other data types” section.

Identify Other Data Types:

- Expand All | Collapse All
- Isolates
- Genomic Sequences
- Genomic Segments (DNA Motif)
- SNPs
 - SNP ID(s)
 - Gene ID
 - Allele Frequency
 - Genomic Location
 - Isolate Comparison
 - Presence in isolate assay
- ESTs
- ORFs
- SAGE Tags

Identify SNPs based on Isolate Comparison

Set A isolate identifiers Enter list: CP3.478952 , CP3.478989, CP3.478990, CP3.478991 , CP3.478992, CP3.478993, CP3.478994 , CP3.478995

Copy Isolates from My Basket (0 Isolates)

Upload from a text file: Browse...

Minimum percentage of isolates in Set A with same allele >=

Set B isolate identifiers Enter list: CP3.478940, CP3.479001, CP3.479003, CP3.479004, CP3.479006, CP3.479007, CP3.479008, CP3.479009, CP3.479010, CP3.479011, CP3.479013, CP3.479014

Copy Isolates from My Basket (0 Isolates)

Upload from a text file: Browse...

Minimum percentage of isolates in Set B with same allele >=

Give this step a weight

b. What does this search do? What is in Set A and B. Run the query and look at your results. How many SNPs were identified between isolates from Brazil and Malawi? What could you use this information for?

c. Find SNPs that differentiate isolates from East Africa and those from West Africa.

- For this exercise we are going to use the same SNPs by isolate comparison search as above. However, first we have to identify isolate IDs from West Africa and ones from East Africa. To do this use the Geographic location query under the isolates section (note that you will need to run this query twice, once for each set of countries):

Identify Other Data Types:

- Expand All | Collapse All
- Isolates
- Isolate ID(s)
- Taxon/Strain
- Host Name
- Isolation Source
- Isolate Sequence Name
- Geographic Location
- Text (search product name, notes, submitter etc.)
- Isolates Clustering
- Genomic Sequences
- Genomic Segments (DNA Motif)
- SNPs
- ESTs
- ORFs
- SAGE Tags

Identify Isolates based on Geographic Location

Geographic Locations

- Africa
 - Benin
 - Cameroon
 - Cape Verde
 - Central African Republic
 - Congo (Dem. Rep.)
 - Gabon
 - Gambia
 - Ghana
 - Guinea
 - Kenya
 - Liberia
 - Madagascar
 - Malawi
 - Mali
 - Mozambique
 - Niger
 - Nigeria
 - Republic Of Cote D Ivoire
 - Senegal
 - Sierra Leone
 - Sudan
 - Tanzania
 - Uganda
 - Zambia
- Asia
- Europe
- N. America
- Oceania/Australia
- S. America
- Unknown

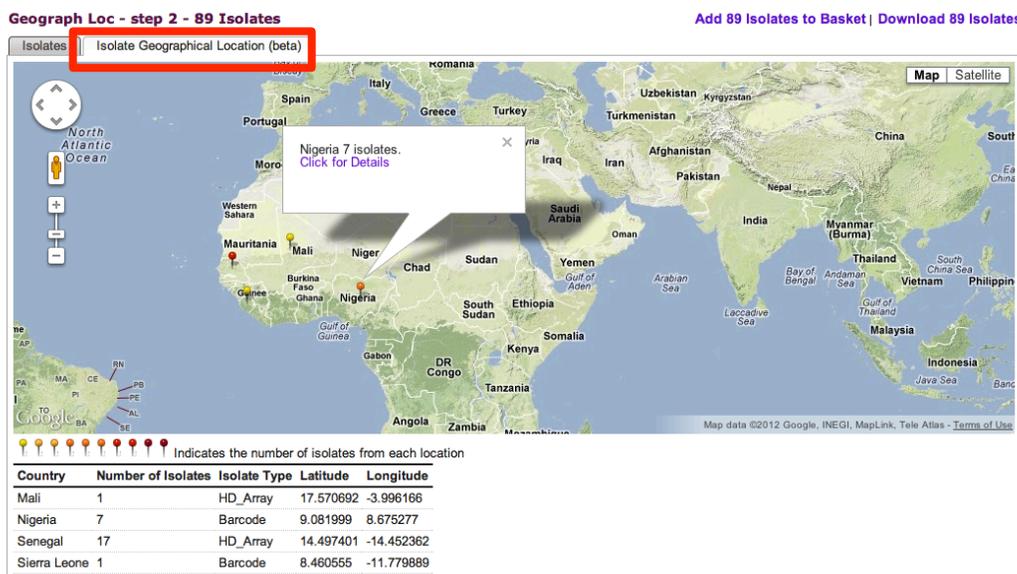
Isolate assay type

- 3kChIP
- Barcode
- HD_Array
- Sequencing Typed

Some East African countries: Kenya, Madagascar, Malawi, Mozambique, Tanzania, Sudan, Uganda, Zambia

Some West African Countries: Cameroon, Gabon, Liberia, Mali, Nigeria, Senegal, Sierra Leone

- For isolate assay type select HD_Array since this array has the most SNPs. You could also try the 3K_chip or even Barcode but shouldn't mix the assay types in one analysis.
- Confirm the distribution of the isolates you get by clicking on the "Geographic location" tab:



- Once you have isolates based on geographic location you will need to copy the IDs and paste them into the SNPs by isolate comparison query (make sure you put isolates from one set of countries into the input box for set A and the other set in the input box for Set B). You might find it useful to use the NotePad on your PC or open the query in another window or tab.
 - o To do this easily, click on "Download Results", select "Tab delimited (Excel):" then unselect all the columns and click on "Get Report". Now copy the list of IDs.
 - o If the above steps are taking too long, feel free to copy the IDs from the following link: <http://goo.gl/rhRdO>
- Once you have the isolate IDs pasted in the isolate comparison query, run it and examine your results. Did you get any results? Revise the query and change the minimum percentage parameters to 70 for both set A and B:

Revise Step

Revise Step 1 : Isolate Comparison

Set A isolate identifiers ?

BC.458086; BC.458090; BC.458091;
BC.458092; BC.458093; BC.458101;
BC.458105; BC.458120; BC.458124;
BC.458125; BC.458126; BC.458127;
BC.458128; BC.458129; BC.458130;

Enter list:
 Copy Isolates from My Basket (0 Isolates)

Minimum percentage of isolates in Set A with same allele >= ?

Set B isolate identifiers ?

BC.458098; BC.458110; BC.458111;
BC.458112; BC.458113; BC.458114;
BC.458115; BC.458116; BC.458117;
BC.458118; BC.458119; BC.458150;
BC.458168; BC.458169; CP3.273609;

Enter list:
 Copy Isolates from My Basket (0 Isolates)

Minimum percentage of isolates in Set B with same allele >= ?

- What do your results look like now?
 - o Which SNP differentiates more isolates (hint: look at the numbers in the columns for SetA and SetB)?
 - o Do you think these SNPs are synonymous or non-synonymous? (hint: click on “select columns” and add the column called “non-synonymous”).
 - o What are the genes that include these SNPs? (hint: click on the gene IDs in the “Gene ID” column).

4.3 Analyzing SNPs on a defined list of genes.

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

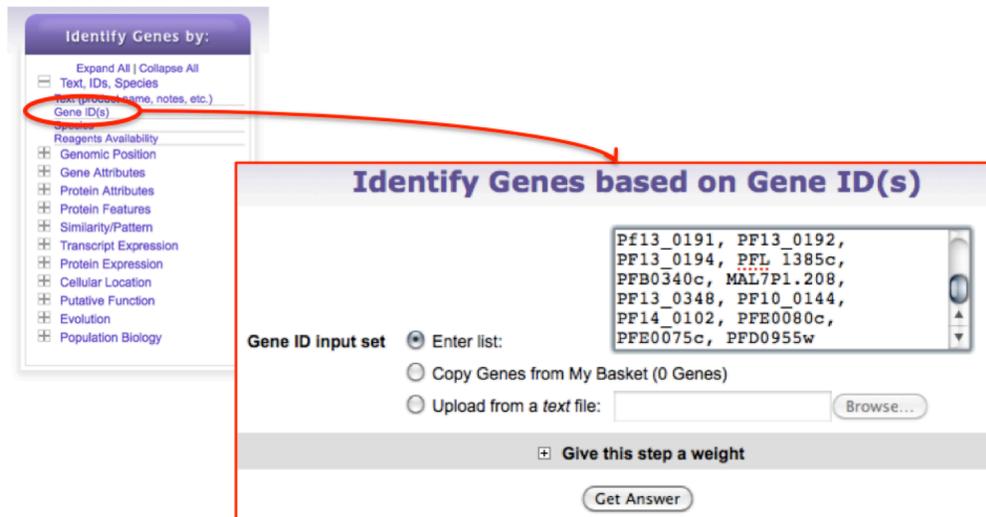
You just read the recent paper by Tetteh *et.al.* (<http://www.ncbi.nlm.nih.gov/pubmed/19440377>) where they perform an analysis of SNPs on a set of *P. falciparum* genes. Their conclusion is that these genes are under “balancing” selection – under diversifying selection due to their exposure to the host’s immune pressure. You decide you would like to analyze their list of genes in PlasmoDB.

Here is the list of gene IDs from their paper:

PFF0615c, Pf13_0338, PFE0395c, PF14_0201, PFF0995c, PF10_0346, PF10_0347, PF10_0348, PF10_0352, PF13_0197, PF13_0196, MAL13P1.174, PF13_0193, MAL13P1.173, Pf13_0191, PF13_0192, PF13_0194, PFL1385c, PFB0340c, MAL7P1.208, PF13_0348, PF10_0144, PF14_0102, PFE0080c, PFE0075c, PFD0955w

- **Put the above list into a step.**

Hint: use the Identify Genes based on Gene ID(s) search option.



Add a step to your strategy to identify how many of these genes are under diversifying selection. Hint: the “Identify Genes based on SNP Characteristics” is found under the population biology menu (see figure on the next page).

- What parameters would you chose?
- Would you expect genes under balancing selection to have a high or low non-synonymous/synonymous SNP ration?
- How many genes were returned by your search? Of these, now many intersect with the set of genes from the paper?

Click on the result for your ID search in the first step (25 genes) and add columns for SNP characteristics (under population biology). Do all these genes appear to be under balancing selection? Is this consistent with the results of your strategy?

(Genes)

Gene ID(s)
26 Genes

Step 1

Add Step

Add Step

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

Genes
Genomic Segments (DNA Motif)
SNPs
ORFs
SAGE Tags

Text, IDs, Species
Genomic Position
Gene Attributes
Protein Attributes
Protein Features
Similarity/Pattern
Transcript Expression
Protein Expression
Cellular Location
Putative Function
Evolution
Population Biology

SNP Characteristics

Close

Add Step

Add Step 2 : SNP Characteristics

Organism

Reference

Comparator

SNP Class

Number of SNPs of above class >=

Number of SNPs of above class <=

Non-synonymous / synonymous SNP ratio >=

Non-synonymous / synonymous SNP ratio <=

SNPs per KB (CDS) >=

SNPs per KB (CDS) <=

Give this search a weight

Give this search a name

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