

## Affymetrix B Allele Frequency Calculation

**Author:** Greta Peterson, Golden Helix, Inc.

### Overview

Using Affymetrix CEL files as its source, this script combines quantile normalized SNP A and B probe intensities for each marker into a theta value, then calculates B-Allele Frequencies for each marker.

For information on accessing B Allele Frequencies for Illumina data, see [http://www.goldenhelix.com/Downloads/SVS/ExportingDSF-GenomeStudio\\_v4.doc](http://www.goldenhelix.com/Downloads/SVS/ExportingDSF-GenomeStudio_v4.doc)

### Recommended Directory Location

Save the script to the following directory:

\*..\Application Data\Golden Helix SVS\UserScripts\SVS\Tools\

**Note:** The **Application Data** folder is a hidden folder on Windows operating systems and its location varies between XP and Vista. The easiest way to locate this directory on your computer is to open SVS and select **Tools >Open Folder > User Scripts Folder**. If saved to the proper folder, this script will be accessible from the project navigator's Tools menu.

### Preparing to use the Script

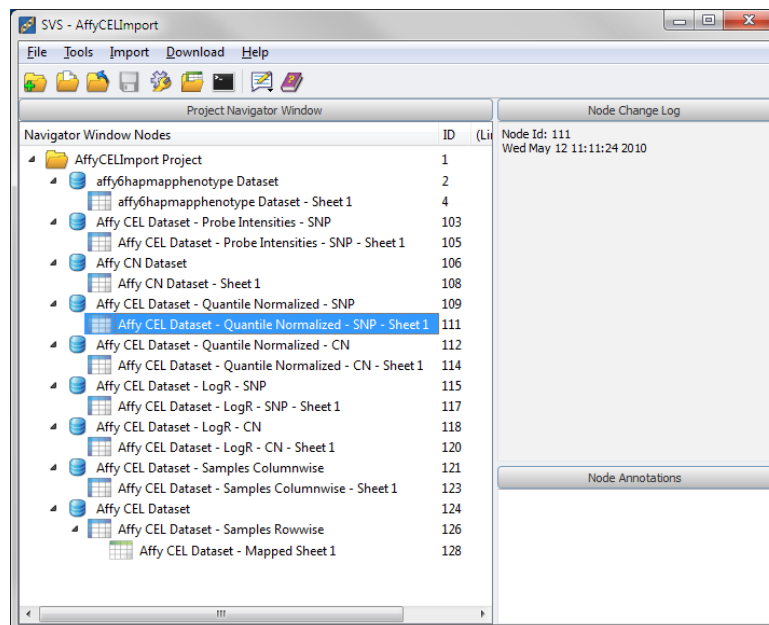
This script requires that CEL files be imported into SVS with intermediate spreadsheets. Follow the instructions below to obtain these intermediate files.

1. From an SVS project, a reference spreadsheet needs to be imported to use when importing CEL files unless all samples are used as reference, or one of the provided reference datasets are used.

**If a reference spreadsheet is used**, then the content of the reference spreadsheet depends on the Affymetrix Array type.

- a. For combining 250K NSP and 250K STY arrays, matching columns are needed. The row labels do not matter, but will be used in the spreadsheets on import of the data. The second column is the CEL file name of the NSP array for the sample. The third column is the CEL file name of the STY array for the sample. Subsequent columns can contain phenotype information.
  - b. For other Affymetrix arrays, the row labels need to match the CEL file names (with or without the .CEL extension), and one of the subsequent columns needs to be binary.
2. To import the CEL files, go to **Import > Affymetrix > CEL**. Select the CEL files. For combined Affymetrix 250K NSP and 250K STY arrays, select both sets of files. Click **Next**.

3. Choose the appropriate parameters in the Import CEL dialog, including selecting the **Quantile Normalized A/B intensities** under **Output Options**.
4. Click **Finish**. The CEL file import process might take a while, and when it is done the chosen spreadsheets are created. See **Figure 1** below. One spreadsheet should have “Quantile Normalized – SNP” in its name or for 500k data there should be two spreadsheets one with “Quantile Normalized – NSP” in its name, the other with “Quantile Normalized – STY” in its name. These are the spreadsheets to use for the B-Allele Frequency calculation.



**Figure 1: SVS Project Navigator after importing Affymetrix 6.0 CEL files with intermediate files**

## Using the Script

1. Go to **Tools > Affymetrix B Allele Frequency Calculation**.
2. You will be asked to select the spreadsheet containing quantile normalized A and B intensities and marker map. Select the appropriate spreadsheet and marker map file.

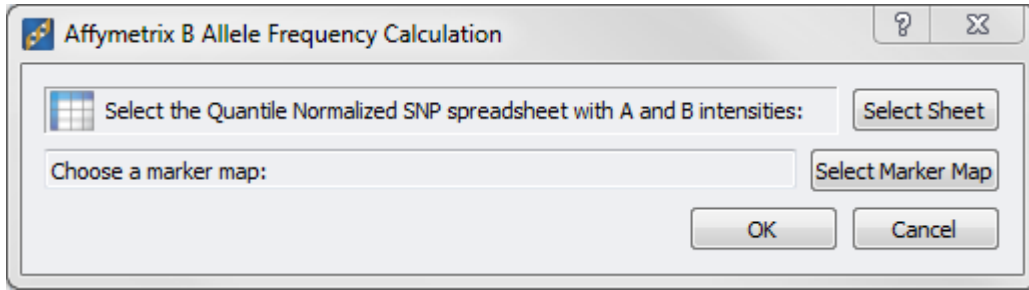


Figure 2: Prompt to select the spreadsheet with A & B intensities and Marker Map

3. Selecting a marker map is optional. If you do not select a marker map, then the final spreadsheets will not be marker mapped. Otherwise, select a marker map and press **OK**.

### How the Script Works

The script will first use the A (X coordinate) and B (Y coordinate) intensities to calculate the theta coordinate in a polar coordinates transformation. The R coordinate is not needed for the B Allele Frequency computation. The theta coordinates are output in a spreadsheet for later reference. Plotting histograms of theta coordinates per SNP can yield important clustering information.

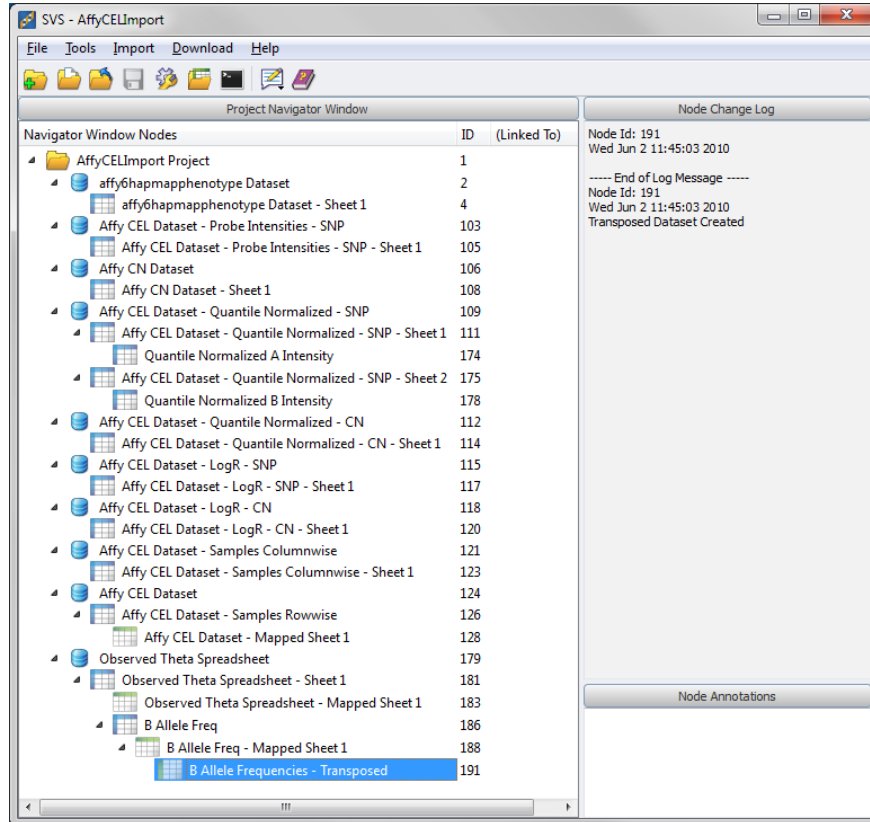
The formula (Wang, 2008) for the theta coordinate of the  $i^{\text{th}}$  sample and  $j^{\text{th}}$  marker is defined as:

$$\theta_{ij} = \frac{\arctan\left(\frac{Y_{ij}}{X_{ij}}\right)}{\pi/2}$$

Next, the script will use the theta values and the assumption that ideally the AA cluster will be near a theta value of 0.1, the AB cluster will be near a theta value of 0.5 and the BB cluster will be near a theta value of 0.9 to calculate approximate cluster means for each marker ( $\theta_{AA}$ ,  $\theta_{AB}$ ,  $\theta_{BB}$ ). These cluster means will be used to calculate the B-Allele Frequency with the following formula (Wang, 2008):

$$B_{ij} = \begin{cases} 0 & \text{if } \theta_{ij} < \theta_{AA} \\ \frac{0.5(\theta_{ij} - \theta_{AA})}{\theta_{AB} - \theta_{AA}} & \text{if } \theta_{AA} \leq \theta_{ij} < \theta_{AB} \\ 0.5 + \frac{0.5(\theta_{ij} - \theta_{AB})}{\theta_{BB} - \theta_{AB}} & \text{if } \theta_{AB} \leq \theta_{ij} < \theta_{BB} \\ 1 & \text{if } \theta_{BB} \leq \theta_{ij} \end{cases}$$

The B-Allele Frequencies are output in a spreadsheet with markers in columns and samples in rows (see **Figure 3**). Then this spreadsheet is transposed to have samples in columns and markers in rows (see **Figure 4**) for plotting of the B-Allele Frequencies for particular samples (see **Figure 5**).



**Figure 3: Project navigator after running script**

Map	Columns	R 1	R 2	R 3	R 4	R 5	R 6
1	SNP A-8575125	0.00584935	0	0	0	0.0108765	0
2	SNP A-8575115	1	1	1	1	0.993548	1
3	SNP A-8575371	0.751008	0.977827	0.860242	1	1	0.999194
4	SNP A-8709646	0.253214	0.0732964	0.163106	0	0	0.0110274
5	SNP A-8497791	0.556183	0	0.983464	0.0694269	0	0.41357
6	SNP A-1909444	0.503106	0.464248	0.444887	0	0.493967	0
7	SNP A-8358063	0.170745	0.137749	0	0	0	0
8	SNP A-8329892	0.0352828	0	0.391298	0	0	0
9	SNP A-8408912	0.509808	0.440435	0.486324	0.0115997	0.502856	0
10	SNP A-8294056	0.394919	0.626142	0.366547	0.596869	0.662869	0.521342
11	SNP A-1886933	0.265613	0.302708	0	0	0	0
12	SNP A-2236359	0.547096	0.584492	0.608295	1	1	1
13	SNP A-8515688	0.549349	1	0.00810824	0	0	0
14	SNP A-2205441	0	0	0.00750051	0	0.00499782	0.00305665
15	SNP A-8524447	0.628757	0.374134	0.378954	0.434353	0	0.392218
16	SNP A-8573955	0.0584712	0.322691	0.273105	0.0100979	0	0
17	SNP A-8530278	0.94923	1	0.968932	0.365925	1	0
18	SNP A-8573668	0.101941	0	0.490766	0.430632	1	0
19	SNP A-8573414	0.338282	0.459666	0.0751788	0.493358	0.38542	1
20	SNP A-8531044	1	1	1	0.478399	1	0.476151
21	SNP A-8530320	0.680246	0.837322	0.667639	0.359571	0.959486	0.399612
22	SNP A-8572481	1	1	1	0.99852	0.986505	0.987481
23	SNP A-2116190	0.958477	0.527655	0.989442	0.471812	0.550019	0.426353
24	SNP A-8325638	0	0.907663	0.484027	0.910687	1	1

**Figure 4: B Allele Frequency - Transposed spreadsheet, ready for plotting**

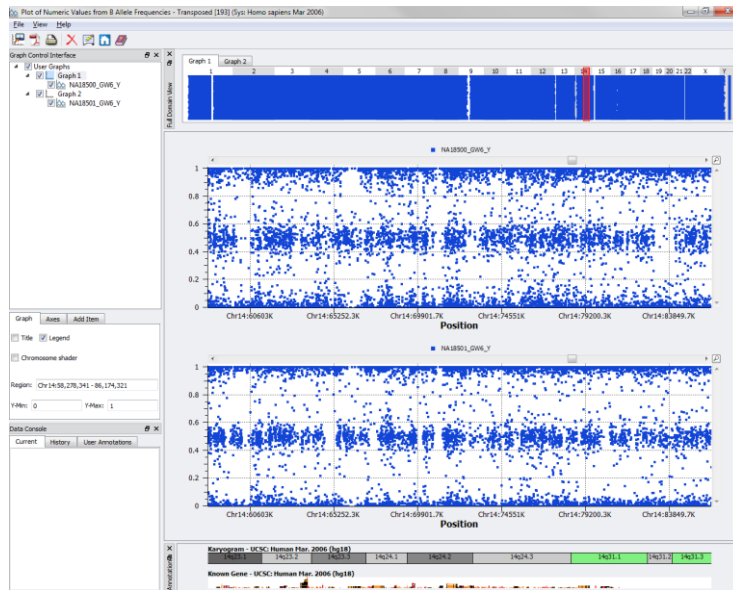


Figure 5: Two B-Allele Frequency plots

## Works Cited

Wang, K. a. (2008). Copy Number Variation Detection via High-Density SNP Genotyping. *Cold Spring Harb. Protoc.*; doi:10.1101/pdb.top46.