

*This document mainly talks about Plexseq reads quality checking, processing VCF file from Mayo, and analysis of genetic data from both Mayo and Plexseq.*

*Alert: All the file path in the scripts (shell, python and others) are subject to my local path and should be adapted to the local path of whoever's using these scripts.*

## **Plexseq Results Quality Checking**

samtools-1.3 is used for all following process

Use *diff.py* to figure out sample ids in *plexseq\_diagnosis.xlsx* (information saved in *sample.txt*) that corresponding to diff-SNPs in *quality\_checking\_12\_13\_2017* (information saved in *difflist.txt*), saved in ***diffid.txt***

Use *filter.py* to filter out **102** samples out of 9024 samples we have that contain diff\_SNPs, saved in copy directory.

Use *bwa.sh* to align 102 samples to hg38.fa, save all the sam files in *sam* folder

Use *sam2bam.sh* to convert sam files to bam files

Use *sort.sh* to sort, use *index.sh* to index all bam files, save in *sorted* folder

Use *brc.sh* to count the reads at interests (*region.txt* includes the snps information obtained from ncbi snps database:

[https://www.ncbi.nlm.nih.gov/SNP/snp\\_ref.cgi?searchType=adhoc\\_search&type=rs&rs=rs4666451](https://www.ncbi.nlm.nih.gov/SNP/snp_ref.cgi?searchType=adhoc_search&type=rs&rs=rs4666451)) All the resulting txt files are saved in *sorted* folder

Use *interpret.py* to interpret the results, saved in ***213snps.xlsx***

## **Processing (Mayo)**

1. Use */Volumes/CORSAIR/Mayo/cleaning\_script/vcf2allele.py* to parse VCF into reads file
2. Use */Volumes/CORSAIR/Mayo/comparison/compare.py* to compare the reads of overlap SNPs from mayo and Plexseq
3. Use */Volumes/CORSAIR/Mayo/hwe/HWE\_recessive.py* or */Volumes/CORSAIR/Mayo/hwe/HWE.py* to calculate Hardy-Weinberg equilibrium (can choose to compute from dominant or recessive)

4. Use /Volumes/CORSAIR/Mayo/mayo\_data\_snp\_flip/vcf2allele\_fliped.py to flip SNPs that are on reverse strands from Plexseq, since Mayo always report SNPs in the forward strands and we want to be consistent about it.

**Analysis:**

5. /Volumes/CORSAIR/TFA (Complex folder):
  - a. **33 loci folder**: info from 33 loci paper
  - b. **GFL folder**: group fused lasso
  - c. **Old binarize method folder**: old binarization method input and results 0.5 for heterozygous
  - d. TWAS folder: TWAS paper information
  - e. **Status prediction folder**: inputs and results from snp\_algorithm.py and survey\_data\_algorithm.py
  - f. Binarize\_mayo.py including functions:
    - 1 zscore calculation
    - 2 binarize data
    - 3 make ped file (for plink)
    - 4 make map file(for plink)
    - 5 ped file transformation (get rid of "/")
  - g. Overlap.py including functions:
    - 1 mayo/plexseq overlap ID check
    - 2 overlap of our snps with 33 and cis-eQTL
    - 3 ...multiple overlap checkings upon requested
  - h. Vital\_match.py including functions:
    - 1 mayo-plexseq overlap checking
    - 2 vital status match
  - i. Plink\_mac folder: utilize Plink to calculate LD.