

## Treat Workflow scripts and Files Description

Treat workflow is executed using one standalone script which in turn calls different scripts to get to the final results. To execute the TREAT it requires three reference files. Example reference file can be found in TREAT scripts folder named examples.

- Configuration file (path to various tools and reference datasets)

```
## path to all the scripts and reference genomes and tools

script_path=path_to_treat_workflow_files
bwa_ref=path_to_bwa_ref
bowtie_ref=path_to_bowtie_ref
human_ref=path_to_ref
sift_ref=path_to_sift_ref
dbSNP_ref=path_to_dbSNP_ref
bwa_path=path_to_bwa
bowtie_path=path_to_bowtie
java_path=path_to_java
samtools_path=path_to_samtools
gatk_path=path_to_GATK
sift_path=path_to_SIFT
OnTarget=OnTarget_file
bed_path=path_to_bedtools
dbSNP_rsids=dbSNP_rsids_file
snvmix_path=path_to_SNVmix_tool
kgenome=path_to_kgenome_data
picard_path=path_to_picard
GeneIdMap=geneid_to_gene_name_mapping_file
hapmap=path_to_hapmap_data
ref_flat=ref_flat_file
codon_ref=codon_usage_file
UCSC_tracks=ucsc_tracks_references
```

- Run information file (meta data information about the project, how many samples need to be processed, parameters user want to use in the workflow)

```
## information about the run and parameter user want to use to run TREAT

tool=type_of_analysis
Date=mm/dd/yyyy
Aligner=BWA/BOWTIE
SNV_caller=SNVMix/GATK
paired=1/0
ReadLength=XXX
Disease=XXX
queue=queue_for_the_cluster
variant_type=BOTH/SNV/INDEL
PI=lastname_firstname_lanID
output_dir=Output_directory_location
input_dir=path_to_input_directory
email=lastname.firstname@mayo.edu
sampleNames=sampleA:sampleB
config=path_to_configuration_file
sample_info=path_to_sample_info_file
analysis=mayo/all/variant/annotation
output_folder=name_for_the_output_folder
center=MAYO/TGEN/BMC/XXX
platform=illumina
GenomeBuild=hg18/hg19
SampleInformation=Free Text for HTML report (information about the samples)
```

- Sample information file (names of the files user want to use in the workflow, as this tool is independent of the naming convention of the input files)

There are different formats for the sample info files, depends on which module user wants to use, refer to the user manual for all the sample info formats.

```
$script_path/treat.sh
Usage:
```

1. <run\_info file>

```
treat.sh (input run_info file)

{
    i. alignment.sh
    ii. variant.sh
    iii. annotation.sh
    iv. numbers.sh
    v. generate.html.sh
}
```

All the steps are interdependent, and each is executed after the completion of the last one executed.

The Workflow is divided into four (4) modules:

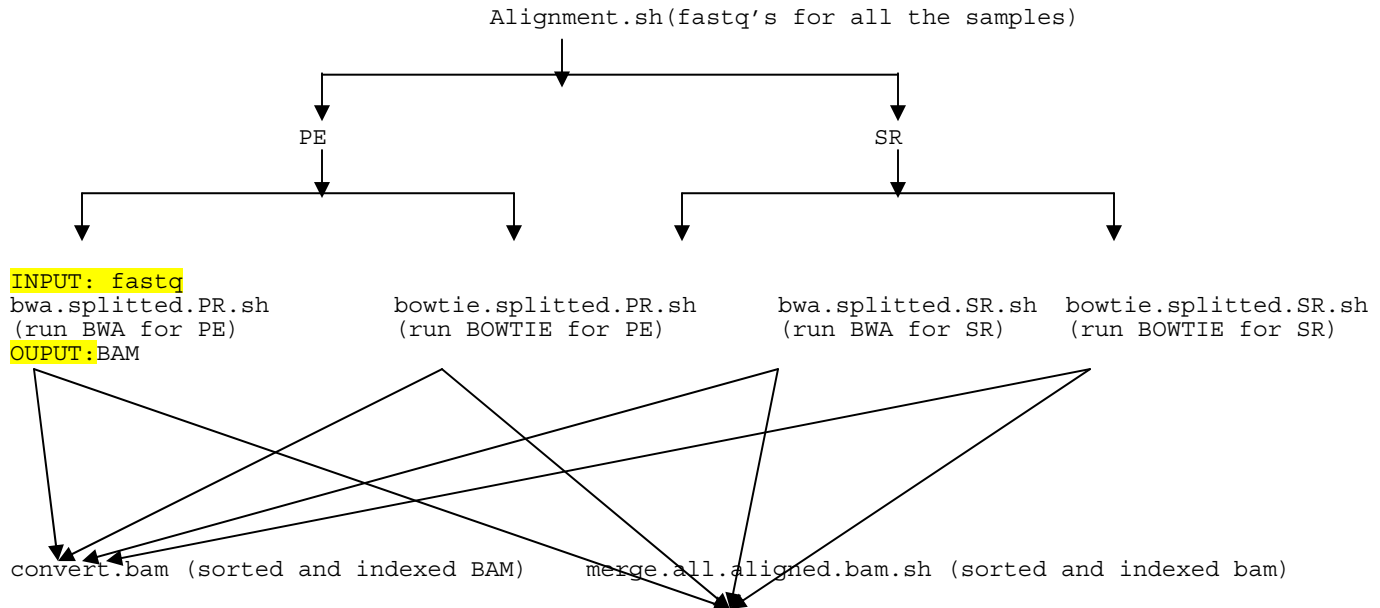
1. Alignment
2. Variant
3. Annotation
4. Numbers

## ALIGNMENT

### Input: FASTQ

#### Example for fastq:

```
@R0211812_0096:1:1:1223:2139#0/1
NGTGAGCCACAGCACCAGCTGCTTCATCTGTAAATAAGGTTTACAGCATCTGCAACCCCTTAACCTCTGGTTTATGGTTATTCCATAACCCAAACAA
+R0211812_0096:1:1:1223:2139#0/1
BSUUJPTXNTPVWPW\ [U] \c_cccccPcc [ ] [ [GIKGGVWXX [SVXVY_X___BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
```



### Output: Sorted BAM & mapping numbers for the BAM

#### Note:

- BAM includes all the reads (mapped or unmapped), which we achieve for longer term. BAM is binary and compressed file sorted in coordinate order. There is Samtools command to extract header of Samtools which contains information about the BAM.

```
/projects/bsi/bictools/apps/alignment/samtools/samtools-0.1.12a/samtools view -H $bam
@SQ      SN:chrM  LN:16571
@SQ      SN:chr1  LN:247249719
@SQ      SN:chrX  LN:154913754
@SQ      SN:chrY  LN:57772954
@PG      ID:bwa   PN:bwa   VN:0.5.9-r16
@PG      ID:BWA
@RG      ID:6-HBPN-Land-JL      SM:6-HBPN-Land-JL
LB:/data1/bsi/refdata/bictools/sequence/human/ncbi/36.49/indexed/BWA_Indexed_reference_including_ChRM/allChr.fa
PL:illumina      CN:MAYO
```

(GenomeBuild, aligner used, sample name, platform information etc.)

- Fastq's can be extracted from BAM using Picard tool

```
$ /data2/bsi/RandD/exome/bam2fastq/bam2fastq.sh
Usage: for PE: bam2fastq.sh <input bam> (provide full path) <name for read1> <name for read2>
        for SR: bam2fastq.sh <input bam> (provide full path) <name for read1>
```

#### Folder to look at:

Output\_folder/alignement/sample/

Each sample folder should includes important file

- Sorted bam (\*.bam) - aligned BAM file (no filtering)

- Indexed bam (\*.bai) - indexed BAM

Important script we use to add platform, read group, sample. GenomeBuild information to the sam files

```
Perl $script_path/add.read_group.platform.pl -i input_sam -o output_sam -r $sample -s $sample -p $platform -a $Aligner -c $center -l $GenomeBuild
```

**Output:** will add the above said information to the sam and which can be use further

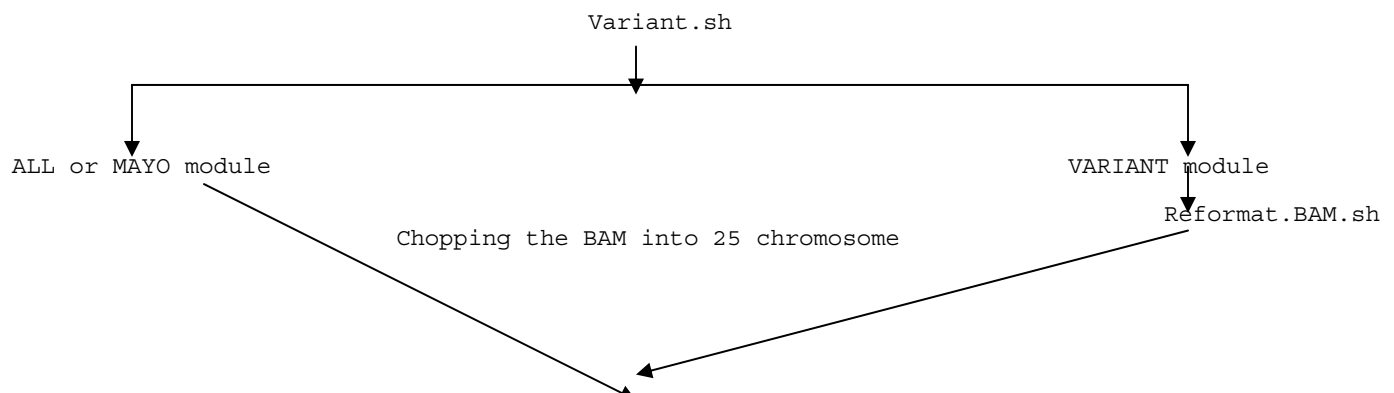
Tools used for this module:

- SAMTOOLS
- BWA
- BOWTIE
- PERL SCRIPTING
- SHELL SCRIPTING

## VARIANT

**Input:** BAM

This module includes Realignment, Variant calling, On Target calling. This is the slowest step in the workflow, so we chop the input BAM into 25 chromosomes (Homosapiens) and run them in parallel to get the results faster.



```

variant_per_chr.sh (parallel processing)
• Realign.sh (realignment using GATK)
Input: per Chr BAM
Output: realigned and recalibrated BAM (Q>= 20 as a filter to BAM's to call variants)
• Call.indels.sh (calling Indels using GATK)
Input: per Chr realigned and recalibrated BAM
Output: VCF output INDEL calls (per Chr)
• Call.snvs.sh (calling SNVs using SNVMix/GATK)
Input: per Chr realigned and recalibrated BAM
Output: outputfile (per Chr)
• Format.variants.sh (format the variants in tab delimited format)
Input: output files from call.snvs and call.indels to format in tab delimited format.
Output: Tab delimited format with required columns
• OnTarget.SNV.sh (intersect the SNV with target kit provided)
Input: SNV formatted file
Output: OnTarget SNVs
• OnTarget.INDEL.sh (intersect the INDEL with target kit provided)
Input: INDEL formatted file
Output: OnTarget INDELs
• OnTarget.BAM.sh (intersect the BAM with target kit provided)
Input: realigned and recalibrated BAM
Output: OnTarget BAM
• OnTarget.PILEUP.sh (intersect the PILEUP with target kit provided)
Input: pileup from realigned and recalibrated BAM
Output: On Target pileup to get coverage numbers

```

## **OUTPUT: variant calls (On Target SNVs and INDELs)**

### **Folder to look for realigned and recalibrated BAM**

Output\_folder/realigned\_data/sample  
 Realigned BAM for IGV visualization

### **FOR VARIANT CALLS**

Output\_folder/variants/SNV  
 Output\_folder/variants/INDEL  
 (all the variants calls are per chromosome)

### **RAW variant calls:**

- SNV \$sample.chr\*.snps.raw.gatk.vcf / \$sample.chr\*.snps.raw.snvmix
- INDEL \$sample.chr\*.indel.gatk.vcf

### **FILTERED variant calls:**

- SNV \$sample.chr\*.raw.snvs
  - SNVMix
    - Probability greater or equal to 0.8 ( $\geq 0.8$ )
    - Discard reference homozygous calls
  - GATK
    - Discard reference homozygous calls
    - Quality filters
- INDEL \$sample.chr\*.raw.indels

### **OnTarget Variant calls**

Output\_folder/OnTarget/

- SNV `$sample.chr*.raw.snvs.bed.i.ToMerge`
- INDEL `$sample.chr*.raw.indels.bed.i.ToMerge`

## FORMAT for the SNV files

### GATK

Chr	Position	Ref	Alt	GenotypeClass	Alt-SupportedReads	Ref-SupportedReads	ReadDepth	Quality
Example:								
chr1	4770	A	G	GG	4	0	4	345.63
chr1	4793	A	G	GG	5	0	5	34.67
chr1	42101	T	G	GG	3	0	3	880.56
chr1	53560	T	C	TC	18	35	53	1003.32

### SNVmix

Chr	Position	Ref	Alt	GenotypeClass	Alt-SupportedReads	Ref-SupportedReads	ReadDepth	Probability
Example:								
chr1	4770	A	G	GG	4	0	4	0.93
chr1	4793	A	G	GG	5	0	5	0.96
chr1	42101	T	G	GG	3	0	3	0.88
chr1	53560	T	C	TC	18	35	53	1.00

## FORMAT for the INDEL files

### GATK

Chr	Start	Stop	INDEL information
Example:			
chr1	53598	53601	-CTA:12/36
chr1	1578607	1578610	-GCG:90/122
chr1	1637509	1637509	+T:18/18
chr1	1637753	1637753	+TTTCTT:70/70

Tools used for this module:

- GATK
- SNVmix
- SAMTOOLS
- PERL scripting
- SHELL scripting
- JAVA 64BIT
- BED TOOLS
- PICARD

## NOTE:

- **SAMTOOL has a bug that even we sort the file as coordinate sort it won't change the sort tag to coordinate. So we use PICARD to fix that header tag.**
- **SNVmix accepts the pileup which can be generated using few parameters**  
`$samtools/samtools pileup -s -f $ref bam > pileup`
- **SNVmix applies quality filters**
  - Min mapping quality 20
  - Minimum base quality 20
- **GATK quality filters (BOTH SNVs AND INDELs)**
  - Min mapping quality 20
  - Min base quality 20
  - Number of threads 4

## IMPORTANT SCRIPTS

To reformat the external Bam

```
perl $script_path/add.read_group.platform.pl -i input_sam -o output_sam -r $sample -s $sample -p $platform -a $Aligner -c $center -l $GenomeBuild
```

**To parse the VCF SNV from GATK**  
perl \$script\_path/parse.vcf.SNV.pl input\_snv\_vcf > output\_file

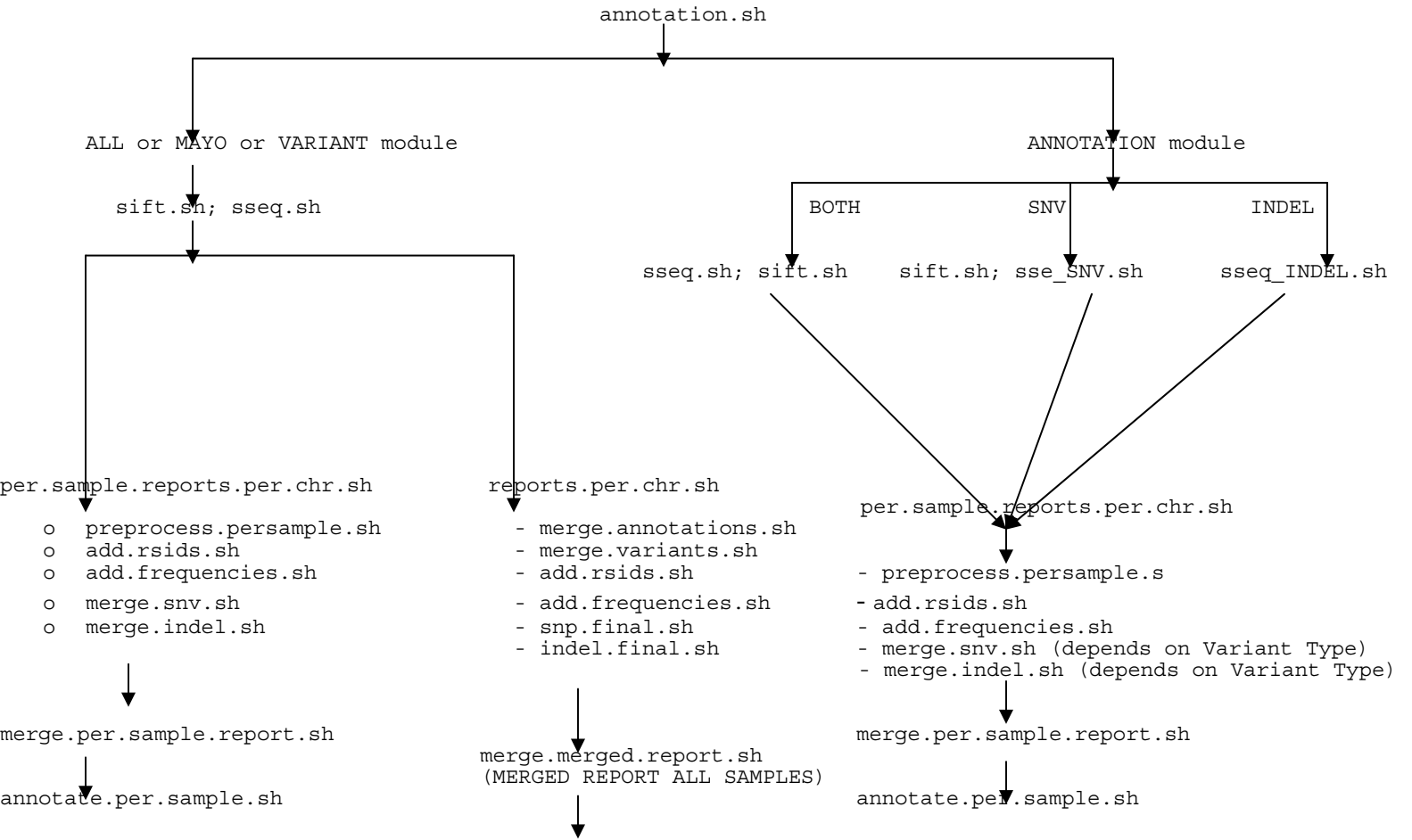
**To parse the VCF INDEL from GATK**  
perl \$script\_path/parse.vcf.INDEL.pl input\_indel\_vcf> output\_file

**To parse the SNVMix SNV output**  
perl \$script\_path/parse.snvmix.to.snvs.pl input\_snvmix\_file > output\_file

**To add rsIDs to the SNV on file**  
perl \$script\_path/add\_dbsnp\_ucsc.pl  
Raw parameters:  
Usage: /projects/bsi/bictools/scripts/dev/exome/TREAT1.0\_hold/add\_dbsnp\_ucsc.pl [-i input file containing chr and pos in tab delimited format -full path] [-b pos base in input file, 1 for 1-based, 0 for 0-based] [-s dbSNP UCSCChg##dbSNP### file use chromStart 0-based, current rsids only] [-c column# of chr in the input file] [-p column# of pos ] [-o output file path] [-r chromosome ]

ANNOTATION

Input: OnTarget Variant calls (per Chr)



(PER SAMPLE REPORTS)

variant.distance.sh  
(Splice variants +/- from Exons)

(PER SAMPLE REPORTS)

Output: Per sample and Merged Variant reports

### Folders to look for:

Output folder/Reports per Sample

(Per sample reports are available)

- SNV
  - \$sample.SNV.cleaned\_annot.xls ( All the OnTarget variant calls)
  - \$sample.SNV.cleaned\_annot\_filtered.xls (default filtered variant calls)
- INDEL
  - \$sample.INDEL.cleaned\_annot.xls
  - \$sample.INDEL.cleaned\_annot\_filtered.xls

Output folder/Reports

(Merged reports)

- SNV
  - SNV.cleaned\_annot.xls
  - SNV.cleaned\_annot\_filtered.xls
- INDEL
  - INDEL.cleaned\_annot.xls
  - INDEL.cleaned\_annot\_filtered.xls
- Variant Distance
  - variantLocation\_SNVs
  - variantLocation\_INDELs

Tools used for this module:

1. JAVA scripts
2. perl and shell scripting
3. SIFT
4. SSEQ web submission script
5. JAVA 64 Bit

### Important scripts used:

- Error with SIFT that it flips the input allele to something , so this script help us to discard those annotation  
perl \$script\_path/sift.inconsistent.pl \$id \$sift/\$snv\_file.sift \$sift
- To add MAYO annotation to the reports  
java -jar /projects/bsi/bictools/scripts/dev/exome/TREAT2.0/exome\_annot.jar  
Usage: java -jar exome\_annot.jar annotate <Folder with SNV and INDEL File (must end with INDEL.report and SNV.report)>java -jar exome\_annot.jar tissue\_files
- To get the variant distance  
java -jar /projects/bsi/bictools/scripts/dev/exome/TREAT2.0/exonvariantlocation.jar  
USAGE: java -jar exonvariantlocation.jar <ref\_seq\_file> <variant\_file> <variant\_type(snp or indel)>
- To add dbsnp id to the list of variants (using reference dbSNP file)  
\$script\_path/add.rsids.sh  
Usage: <TempReportDir> <configuration file> <variant file><chromosome>



- To add allele frequencies from different populations from Hapmap and 1kgenome  
\$script\_path/add.frequencies.sh  
Usage: <TempReportDir> <configuration file> <variant file with rsids ><chromosome>
- To preprocess the input file so as to merge the input variants and annoattaions  
\$script\_path/preprocess.sh  
Usage:<sample> <sift dir> <sseq dir> <configuration file> <tempReport dir> <input variant folder><email><analysis type><run info><reports per sample> <input variant folder><chromosome>
- To merge per sample SNV and INDEL files with all the annotation  
\$script\_path/merge.snv  
Usage:<Tempreports><configuration><sample name><chromosome><sseq dir><sift dir><variant file>  
\$script\_path/merge.indel.sh  
Usage <TempReports> <config> <sample> <which\_chr> <sseq> <indel\_file><input\_dir>
- All the processing is done on per chr variant per sample so to merge all the chr in one file  
\$script\_path/merge.per.sample.report.sh  
Usage: <config> <output\_dir> <TempReports> <sample> <run\_info>
- To annotate using MAYO annoattaions for per sample SNV and INDEL reports  
\$script\_path/annotate.per.sample.sh  
Usage : <config> < output\_dir>
- To get variant splice distance  
\$script\_path/variant.distance.sh  
usage:<config><TempReports><output\_dir>
- To merge all the annotations from different samples so that we can use same annotations with the merged report  
\$script\_path/merge.annotations.sh  
Usage: <sift dir><sseq dir><config><chromosome>
- To merge all the variant calls for all the samples  
\$script\_path/ merge.variants.sh  
Usage: <output dir><temp reports><config><chromossome><run\_info>
- To merge all the annotation with the variant call (merged report)  
\$script\_path/snp.final.sh  
Usage: <Tempreports dir><sift dir><sseq dir><config><chromosome><variant file>  
\$script\_path/indel.final.sh  
Usage:<Tempreports><sseq dir><configuration file><chromosome><variant file>
- To add MAYO annotation to the merged report  
\$script\_path/merge.merged.report.sh  
Usage: <output\_dir><TempReports><config>

## NUMBERS

numbers.sh  
Run in parallel for x samples  
↓  
sample.numbers.sh

### Folder to look for:

Output\_folder/numbers  
There are two files per sample  
1. \$sample.out  
2. \$sample.coverage.out

Example:

Total Reads  
2915968  
Mapped Reads  
2838835  
.....

Each number is preceded by a description and files are used to create HTML report

### Tools used for this module:

1. shell and perl scripting

Important Scripts:

To get the value for transition to transversion ratio

```
perl $script_path/transition.transversion.pl
```

To get one annotation call per variant from Seattle seq (INDELS) with default priority list provided

```
perl $script_path/to.parse.sseq.result.indel.per.sample.pl
```

```
frameshift > coding > splice-5=splice3 > others
```

To get one annotation call per variant from Seattle seq (SNVs) with default priority list provided

```
perl $script_path/to.parse.sseq.result.per.sample.pl
```

```
nonsense > missense > coding-synonymous > coding-notMod3 > others
```

To generate HTML report

```
$script_path/generate.html.sh $output_dir $run_info
```

This script generates HTML report and send out an email to the person running the TREAT stating that the workflow is complete.

To clean the space on secondary space

```
$script_path/cleanspace.sh $output_dir
```