

Whole Exome Sequencing for SNP and Indel Discovery source code

Tools and Libraries Used

- **FastQC**

Tool used to perform quality control checks on raw sequence data.

```
sudo apt install fastqc
```

- **MultiQC**

Aggregates results from FastQC into a single HTML report.

```
pip install multiqc
multiqc .
```

- **GATK (Genome Analysis Toolkit)**

Toolkit for variant discovery and genotyping developed by the Broad Institute.

```
# Download GATK from https://github.com/broadinstitute/gatk/releases
# Make it executable:
chmod +x gatk
./gatk --help
```

- **Samtools**

Utilities for manipulating alignments in the SAM/BAM format.

```
sudo apt install samtools
```

- **BWA (Burrows-Wheeler Aligner)**

For aligning sequence reads to a large reference genome.

```
sudo apt install bwa
bwa index hg38.fa
bwa mem hg38.fa reads_1.fastq.gz reads_2.fastq.gz > aligned.sam
```

- **Picard**

Toolset for manipulating high-throughput sequencing data.

```
# Download from https://broadinstitute.github.io/picard/
java -jar picard.jar --help
```

- **Funcotator (Functional Annotation)**

GATK tool for annotating variants using known biological databases.

```
# Part of GATK bundle. Requires downloading data sources:
wget https://...funcotator_dataSources.v1.7.20200521g.tar.gz
tar -xvzf funcotator_dataSources.v1.7.20200521g.tar.gz
```

- **wget / curl**

Used to download required reference and known site files.

```
sudo apt install wget
```

1 Project Folder Structure

```
project_root/
  supporting_files/
    hg38/
      hg38.fa
      hg38.fai
      hg38.dict
      Homo_sapiens_assembly38.dbsnp138.vcf
      Homo_sapiens_assembly38.dbsnp138.vcf.idx

  reads/
    SRR062634_1.filt.fastq.gz
    SRR062634_2.filt.fastq.gz

  aligned_reads/
    SRR062634_sorted_dedup_bqsr_reads.bam

  data/
    recal_data.table

  results/
    raw_variants.vcf
    raw_snps.vcf
    raw_indels.vcf

    filtered_snps.vcf
    filtered_indels.vcf

    analysis-ready-snps.vcf
    analysis-ready-indels.vcf

    analysis-ready-snps-filteredGT.vcf
    analysis-ready-indels-filteredGT.vcf

    analysis-ready-snps-filteredGT-functotated.vcf
    analysis-ready-indels-filteredGT-functotated.vcf

  scripts/
    variant_calling.sh
    filter_and_annotation.sh
```

2 DataSets

```
1  #!/bin/bash
2  # Script: Variant_calling.sh
3  # Purpose: Download paired-end FASTQ files for sample SRR794247 (HG03012)
4  # Project: Whole Exome Sequencing (WES) Pipeline
5
6  # Define output directory
7  OUTPUT_DIR="/Users/aravindsudhakar/Desktop/BioInformatics/reads"
8
9  # Download paired-end FASTQ files
10 wget -P $OUTPUT_DIR
   ↪ ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR794/SRR794247/SRR794247_1.fastq.gz
11 wget -P $OUTPUT_DIR
   ↪ ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR794/SRR794247/SRR794247_2.fastq.gz
```

3 Reference Genome Preparation

```
1  #!/bin/bash
2  # -----
3  # Purpose: Prepare reference genome and supporting files for variant calling
4  # -----
5
6  # Define base directory
7  REF_DIR=~/Desktop/demo/supporting_files/hg38/
8
9  # Download reference genome (hg38) and unzip
10 wget -P $REF_DIR https://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/hg38.fa.gz
11 gunzip $REF_DIR/hg38.fa.gz
12
13 # Index the reference genome (.fai index) using samtools
14 samtools faidx $REF_DIR/hg38.fa
15
16 # Create sequence dictionary (.dict) using GATK
17 gatk CreateSequenceDictionary \
18   -R $REF_DIR/hg38.fa \
19   -O $REF_DIR/hg38.dict
20
21 # Download known sites VCF files for Base Quality Score Recalibration (BQSR)
22 wget -P $REF_DIR \
23   https://storage.googleapis.com/genomics-public-data/resources/broad/hg38/v0/\
24   Homo_sapiens_assembly38.dbsnp138.vcf
25
26 wget -P $REF_DIR \
27   https://storage.googleapis.com/genomics-public-data/resources/broad/hg38/v0/\
28   Homo_sapiens_assembly38.dbsnp138.vcf.idx
```

4 Variant Calling Setup: Directory Paths

```
1  #!/bin/bash
2  # -----
3  # Script: variant_calling_setup.sh
4  # Purpose: Define directory paths for variant calling pipeline (common paths)
5  # -----
6
7  # Path to reference genome (FASTA)
8  ref="./supporting_files/hg38/hg38.fa"
9
10 # Path to known variant sites for BQSR
11 known_sites="./supporting_files/hg38/Homo_sapiens_assembly38.dbsnp138.vcf"
12
13 # Directory containing aligned BAM files
14 aligned_reads="./aligned_reads"
15
16 # Directory containing input FASTQ read files
17 reads="./reads"
18
19 # Output directory for variant calling results
20 results="./results"
21
22 # Data directory for intermediate or other input files
23 data="./data"
```

5 Variant Calling Setup

```
1  #!/bin/bash
2  # -----
3  # Purpose: Define directory paths and input files for variant calling steps
4  # -----
5
6  # Define base project directory (modify as needed)
7  BASE_DIR="./project_root"
8
9  # Reference genome FASTA file
10 ref="${BASE_DIR}/supporting_files/hg38/hg38.fa"
11
12 # Known variant sites for Base Quality Score Recalibration (BQSR)
13 known_sites="${BASE_DIR}/supporting_files/hg38/Homo_sapiens_assembly38.dbsnp138.vcf"
14
15 # Directory containing aligned BAM files
16 aligned_reads="${BASE_DIR}/VC/aligned_reads"
17
18 # Directory containing raw paired-end FASTQ files
19 reads="${BASE_DIR}/VC/reads"
20
21 # Directory for storing result files
22 results="${BASE_DIR}/VC/results"
23
24 # Directory for storing intermediate data (e.g., metrics, temp files)
25 data="${BASE_DIR}/VC/data"
26
```

Step 1: Quality Control (FastQC)

```
1  # -----
2  # STEP 1: Quality Control using FastQC
3  # -----
4
5  echo "STEP 1: Quality Control - Running FastQC on raw reads"
6
7  # Run FastQC on forward and reverse reads
8  fastqc ${reads}/SRR794247_1.filt.fastq.gz -o ${reads}/
9  fastqc ${reads}/SRR794247_2.filt.fastq.gz -o ${reads}/
```

Read Quality and Trimming

Initial quality assessment with FastQC indicated that the sequencing reads were of high quality, with no significant adapter contamination or low-quality bases. Therefore, no trimming was performed prior to alignment.

Note: If trimming were necessary, the following command using Trimmomatic could be applied to remove adapters and low-quality bases:

```
1  # trimming command using Trimmomatic for paired-end reads
2  trimmomatic PE -phred33 \
3    ${reads}/SRR794247_1.fastq.gz ${reads}/SRR794247_2.fastq.gz \
4    ${reads}/SRR794247_1.trimmed.fastq.gz ${reads}/SRR794247_1.unpaired.fastq.gz \
5    ${reads}/SRR794247_2.trimmed.fastq.gz ${reads}/SRR794247_2.unpaired.fastq.gz \
6    ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 \
7    LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

Step 2: Mapping Reads to Reference Genome using BWA-MEM

```
1 echo "STEP 2: Map to reference using BWA-MEM"
2
3 # BWA index reference
4 bwa index ${ref}
5
6 # BWA alignment
7 bwa mem -t 4 -R "@RG\tID:SRR794247\tPL:ILLUMINA\tSM:SRR794247" \
8   ${ref} ${reads}/SRR794247_1.filt.fastq.gz ${reads}/SRR794247_2.filt.fastq.gz \
9   > ${aligned_reads}/SRR794247.paired.sam
```

Step 3: Mark Duplicates and Sort BAM File using GATK

```
1 echo "STEP 3: Mark Duplicates and Sort - GATK"
2
3 gatk MarkDuplicatesSpark -I ${aligned_reads}/SRR794247.paired.sam \
4   -O ${aligned_reads}/SRR794247_sorted_dedup_reads.bam
```

Step 4: Base Quality Recalibration

```
1 # -----
2 # STEP 4: Base quality recalibration
3 # -----
4
5 echo "STEP 4: Base quality recalibration"
6
7 # 1. Build the recalibration model
8 gatk BaseRecalibrator -I ${aligned_reads}/SRR794247_sorted_dedup_reads.bam \
9   -R ${ref} \
10  --known-sites ${known_sites} \
11  -O ${data}/recal_data.table
12
13 # 2. Apply the recalibration to adjust base quality scores
14 gatk ApplyBQSR -I ${aligned_reads}/SRR794247_sorted_dedup_reads.bam \
15   -R ${ref} \
16   --bqsr-recal-file ${data}/recal_data.table \
17   -O ${aligned_reads}/SRR794247_sorted_dedup_bqsr_reads.bam
```

Step 5: Collect Alignment & Insert Size Metrics

```
1 # -----
2 # STEP 5: Collect Alignment & Insert Size Metrics
3 # -----
4
5 echo "STEP 5: Collect Alignment & Insert Size Metrics"
6
7 gatk CollectAlignmentSummaryMetrics \
8   R=${ref} \
9   I=${aligned_reads}/SRR794247_sorted_dedup_bqsr_reads.bam \
10  O=${aligned_reads}/alignment_metrics.txt
11
12 gatk CollectInsertSizeMetrics \
13   INPUT=${aligned_reads}/SRR794247_sorted_dedup_bqsr_reads.bam \
14   OUTPUT=${aligned_reads}/insert_size_metrics.txt \
15   HISTOGRAM_FILE=${aligned_reads}/insert_size_histogram.pdf
```

Step 6: Call Variants - GATK HaplotypeCaller

```
1 # -----
2 # STEP 6: Call Variants - gatk haplotype caller
3 # -----
4
5 echo "STEP 6: Call Variants - gatk haplotype caller"
6
7 gatk HaplotypeCaller \
8     -R ${ref} \
9     -I ${aligned_reads}/SRR794247_sorted_dedup_bqsr_reads.bam \
10    -O ${results}/raw_variants.vcf
```

Step 7: Extract SNPs and Indels

```
1 # Extract SNPs
2 gatk SelectVariants \
3     -R ${ref} \
4     -V ${results}/raw_variants.vcf \
5     --select-type SNP \
6     -O ${results}/raw_snps.vcf
7
8 # Extract Indels
9 gatk SelectVariants \
10    -R ${ref} \
11    -V ${results}/raw_variants.vcf \
12    --select-type INDEL \
13    -O ${results}/raw_indels.vcf
```

Annotation and Filtering

```
# Set base working directory
BASE_DIR="/workspace/project"

# Define reference genome and known sites
ref="${BASE_DIR}/supporting_files/hg38/hg38.fa"

# Input and output directories
results="${BASE_DIR}/results"
data="${BASE_DIR}/data"
```

Variant Filtering

```
# -----
# Filter Variants - GATK
# -----

# Filter SNPs
gatk VariantFiltration \
    -R ${ref} \
    -V ${results}/raw_snps.vcf \
    -O ${results}/filtered_snps.vcf \
    -filter-name "QD_filter" -filter "QD < 2.0" \
    -filter-name "FS_filter" -filter "FS > 60.0" \
    -filter-name "MQ_filter" -filter "MQ < 40.0" \
    -filter-name "SOR_filter" -filter "SOR > 4.0" \
    -filter-name "MQRankSum_filter" -filter "MQRankSum < -12.5" \
    -filter-name "ReadPosRankSum_filter" -filter "ReadPosRankSum < -8.0" \
```

```

-genotype-filter-expression "DP < 10" \
-genotype-filter-name "DP_filter" \
-genotype-filter-expression "GQ < 10" \
-genotype-filter-name "GQ_filter"

# Filter INDELS
gatk VariantFiltration \
  -R ${ref} \
  -V ${results}/raw_indels.vcf \
  -O ${results}/filtered_indels.vcf \
  -filter-name "QD_filter" -filter "QD < 2.0" \
  -filter-name "FS_filter" -filter "FS > 200.0" \
  -filter-name "SOR_filter" -filter "SOR > 10.0" \
  -genotype-filter-expression "DP < 10" \
  -genotype-filter-name "DP_filter" \
  -genotype-filter-expression "GQ < 10" \
  -genotype-filter-name "GQ_filter"

```

Selecting Variants That Pass Filters

Select Variants that PASS filters

```

gatk SelectVariants \
  --exclude-filtered \
  -V ${results}/filtered_snps.vcf \
  -O ${results}/analysis-ready-snps.vcf

```

```

gatk SelectVariants \
  --exclude-filtered \
  -V ${results}/filtered_indels.vcf \
  -O ${results}/analysis-ready-indels.vcf

```

Remove variants that failed genotype filters

```

cat analysis-ready-snps.vcf | grep -v -E "DP_filter|GQ_filter" >
↪ analysis-ready-snps-filteredGT.vcf
cat analysis-ready-indels.vcf | grep -v -E "DP_filter|GQ_filter" >
↪ analysis-ready-indels-filteredGT.vcf

```

Annotating Variants - GATK Funcotator

Annotate SNPs using Funcotator

```

gatk Funcotator \
  --variant ${results}/analysis-ready-snps-filteredGT.vcf \
  --reference ${ref} \
  --ref-version hg38 \
  --data-sources-path /path/to/funcotator_dataSources.v1.7.20200521g \
  --output ${results}/analysis-ready-snps-filteredGT-funcotated.vcf \
  --output-file-format VCF

```

Annotate INDELS using Funcotator

```

gatk Funcotator \
  --variant ${results}/analysis-ready-indels-filteredGT.vcf \
  --reference ${ref} \
  --ref-version hg38 \
  --data-sources-path /path/to/funcotator_dataSources.v1.7.20200521g \
  --output ${results}/analysis-ready-indels-filteredGT-funcotated.vcf \
  --output-file-format VCF

```

```
echo "Annotation completed successfully."

# Optional: summarize the number of annotated variants
echo "Summary of annotated SNP variants:"
grep -v "^#" ${results}/analysis-ready-snps-filteredGT-functotated.vcf | wc -l

echo "Summary of annotated INDEL variants:"
grep -v "^#" ${results}/analysis-ready-indels-filteredGT-functotated.vcf | wc -l
```