

## Introduction to Population Genomics: The use of **Aristotelis HPC**

19-06-2025 Σαγώνας Κωνσταντίνος https://www.sagonaslab.com/



## **Population Genomics**

#### Main questions:

- Footprints of natural and artificial selection
- Local adaptation to changing environments
- Population structure, gene flow, (adaptive) introgression
- Evolution of genomic variation within and between species
- Deleterious mutations /Conservation Biology
- Speciation / Hybrid zones
- Methods in population genomics (demographic inferences, genome scans/GWAS, ...)
- Population Epigenetics
- Population Transcriptomics



### **Population Genomics**



SNP caller for low-coverage data, e.g. ANGSD

Traditional SNP caller, e.g. GATK, Stacks, ipyRAD, Vcftools, Bcftools Allele counts data, e.g. synchronized mpileup

One of the first things we are interested in is: How many sequences are present in my fastq file?

**Fastq.gz** files: zcat AE13\_R1.fastq.gz | wc –l #(divide by 4)

#### Trimming reads has the objectives:

- Removing adapters
- Removing low quality bases
- Excluding short reads after quality trimming

#### **Mapping Reads**

Read mapping corresponds to a balance between **speed and accuracy**:

- Faster algorithms find approximate positions (mapping)
- Slower, precise methods (alignment) match each base accurately

Software often allows fine-tuning of detection, e.g. --very-fast vs. --very-sensitive modes in Bowtie2

Different format used for the outputs: **SAM** (text-based) / **BAM** (binary) / **CRAM** (compressed)

Summarizing the results of the mapping (e.g., bamtools stats / samtools flagstat)

#### **Mapping Reads**

Mappingqualities(MAPQ;expressedin-10log10probabilitythat the mapping position is wrong)

Filtering low-confidence reads (e.g. MAPQ <5, <20) is generally performed since these mapping are more error-prone, excluding them therefore improve downstream analysis such as SNP calling



Mark (not remove) duplicates!

- **Picard** (Markduplicates)
  - Samtools (markdup)

#### Variant Calling

Multi-sample variant calling. By analyzing all individuals together, variant callers can:

- Increase sensitivity: Detect low-frequency variants that might be missed in single-sample analyses
- Improve accuracy: Use allele frequency data to distinguish true variants from errors
- Enable joint genotyping: Call genotypes across all samples consistently, which helps in downstream analyses like population genetics and GWAS

#### GATK, FreeBayes, Vcftools, Bcftools, etc



#### Variant Calling

#### Selecting most reliable variants

136 ##contig=<ID=chrUn g1000247,length=36422,assembly=hg19> 137 ##contig=<ID=chrUn g1000248,length=39786,assembly=hg19> 138 ##contig=<ID=chrUn g1000249,length=38502,assembly=hg19> 139 ##reference=file:///l3bioinfo/ucsc.hgl9.fasta 140 ##source=SelectVariants FILTER INFO FORMAT 15001711232757A 141 #CHROM POS ID REF ALT OUAL 142 chrl 14464 . A T 810.77 PASS AC=1;AF=0.500;AN=2;BaseQRankSum=0.152;ClippingRankSum=0.00;DP=56;Excess 143 chrl 15274 . A T 1084.77 PASS AC=2;AF=1.00;AN=2;DP=43;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC 144 chrl 28563 . A G 139.90 PASS AC=2;AF=1.00;AN=2;DP=5;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MO= 145 chrl 49298 . T C 515.77 PASS AC=2;AF=1.00;AN=2;DP=17;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC 146 chrl 52238 . T G 716.77 PASS AC=2;AF=1.00;AN=2;DP=22;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC 147 chrl 55926 . T C 120.90 PASS AC=2;AF=1.00;AN=2;DP=5;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MQ= 148 chrl 61442 . A G 314.77 PASS AC=2;AF=1.00;AN=2;DP=10;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC 149 chrl 61947 . C T 397.77 PASS AC=1;AF=0.500;AN=2;BaseQRankSum=3.01;ClippingRankSum=0.00;DP=33;ExcessF 150 chr1 61987 . A G 703.77 PASS AC=1;AF=0.500;AN=2;BaseQRankSum=0.426;ClippingRankSum=0.00;DP=42;Excess 151 chrl 61989 . G C 703.77 PASS AC=1;AF=0.500;AN=2;BaseQRankSum=0.125;ClippingRankSum=0.00;DP=41;Excess 152 chrl 69511 . A G 358.77 PASS AC=2;AF=1.00;AN=2;DP=13;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC PASS AC=2;AF=1.00;AN=2;DP=7;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MO= 153 chrl 83084 . T A 204.80 PASS AC=2;AF=1.00;AN=2;DP=19;ExcessHet=3.0103;FS=0.000;MLEAC=2;MLEAF=1.00;MC 154 chrl 89946 . A T 613.77

#### Variant Calling

#### Selecting most reliable variants

Quality by Depth (QD): ensures that variants have sufficient supporting read depth Mapping Quality (MQ): filtering out variants associated with mapping errors Base Quality (BQ): average base quality of the reads sporting the variant, filtering out variants with low base quality

Minimum Allele Frequency (AF): filtering out variants with low allele frequency Depth of Coverage (DP): ensures that the variant is supported by a large number of reads Variant type (SNPs vs. indels): calling INDELs is more challenging than SNPs Missing rate: filtering out variants with excessive missing data is a standard to reduce false positives (true also at the sample level)

### **Population Genomics**

The study of changes in allele frequency over time and space

With new technologies we can examine millions of "independent" markers along the genome. We can further detect which polymorphisms segregate in a way that is not consistent with others (e.g. selection)



### **Population Genomics**

How to deal with thousands or millions of markers?

- Measures of pairwise relatedness, pedigrees, genetic distances, summary statistics (Fstatistics, ABBA-BABA...)
- Dimensionality-reduction (e.g., DAPC, PCA, NMDS)
- Clustering of individuals (can be modelbased; e.g., ADMIXTURE)
- Tree based methods (phylogenetics)



# Describing population structure using allele frequencies

**F-statistics** measure the reduction in heterozygosity relative to HW expectations in specific populations (subpopulation).

$$F = 1 - (H_{\odot} / H_{E})$$
 or  $F = (H_{E} - H_{\odot}) / H_{E}$ 



Fis: probability that two alleles in an individual are identical by descent ( $\approx F$  averaged across all individuals) - intra-population

**Fsr:** fixation index - probability that two alleles from two populations are identical by descent

Fir: general genetic structure

In real life, we use different estimators (Weir an Cockerham, Bhatia, GST...)

## Integrate other data on population genomics



# GATK

### **Genome Analysis Toolkit**

#### Variant Discovery in High-Throughput Sequencing Data







## Step 1: HaplotypeCaller

- Command: gatk HaplotypeCaller
- Generates a GVCF per-sample for joint genotyping.
- Important options:
- -R: Reference genome (FASTA)
- -I: Input BAM file
- -O: Output GVCF file
- -ERC GVCF: Emit reference confidence mode



## Step 2: GenomicsDBImport

- Command: gatk GenomicsDBImport
- Converts per-sample GVCFs into a database format for efficient joint calling.
- Important options:
- --genomicsdb-workspace-path: Output database path
- --batch-size: Number of samples per batch
- --intervals: List of genomic intervals (chromosomes)
- -V: Input GVCFs



## Step 3: GenotypeGVCFs

- Command: gatk GenotypeGVCFs
- Performs joint genotyping on the database produced in previous step.
- Important options:
- -R: Reference genome
- -V: gendb:// path to GenomicsDB
- -O: Output VCF file



## Step 4: Variant Filtration

- Two main options: Hard Filtering or Variant Quality Score Recalibration (VQSR)
- Command (Hard Filtering): gatk VariantFiltration
- Important expressions:
- - QD < 2.0
- - FS > 60.0
- - MQ < 40.0
- - SOR > 3.0
- Use SelectVariants before this to filter SNPs or INDELs specifically



## Hard Filtering vs. VQSR

- Hard Filtering: Uses fixed thresholds for filtering variants.
- Pros: Simple, no external data required.
- Cons: Less flexible, risk of over/under-filtering.
- VQSR: Builds statistical models using known, high-quality variant datasets.
- Pros: More accurate, adaptive.
- Cons: Requires large datasets and known resources (truth sets).



## Post-Variant Calling Analyses in R

#### Population Genetic Indices

- adegenet diversity, clustering, PCA
- hierfstat F-statistics
- pegas haplotype and allele frequency analysis
- poppr clone correction, AMOVA

#### Landscape Genomics

- LEA genotype-environment associations (LFMM)
- ResistanceGA isolation-by-resistance modeling
- gdm generalized dissimilarity modeling
- raster / terra spatial data handling

#### **11** Demographic Inference

- fastsimcoalR simulate demographic scenarios
- abc approximate Bayesian computation
- msprime, SLiM demographic simulations

#### 🛷 Structure & Relatedness

- SNPRelate PCA, IBS, LD analysis
- adegenet DAPC for population structure
- pcadapt outlier detection for structure
- starmie STRUCTURE plot visualization

#### 🐗 Visualization

- ggplot2 / ggpubr elegant data plotting
- ggmap / sf spatial visualizations
- plotly interactive PCA, STRUCTURE plots



Registration Open for GBCC2025:         Joint Galaxy/Bioconductor Conference!						
Bioconductor Open source software for bioinformatics	About	Learn	Packages	Developers	Q Search	Get Started >

## Open source software for Bioinformatics

The Bioconductor project aims to develop and share open source software for precise and repeatable analysis of biological data. We foster an inclusive and collaborative community of developers and data scientists.

if (!require("BiocManager", quietly = TRUE))
install.packages("BiocManager")
BiocManager::install(version = "3.21")

BiocManager::install("packagename")

