

# Variant calling

## Detecting variants in NGS data

Samtools and the Genome Analysis ToolKit (GATK)

University of Cambridge

Cambridge, UK

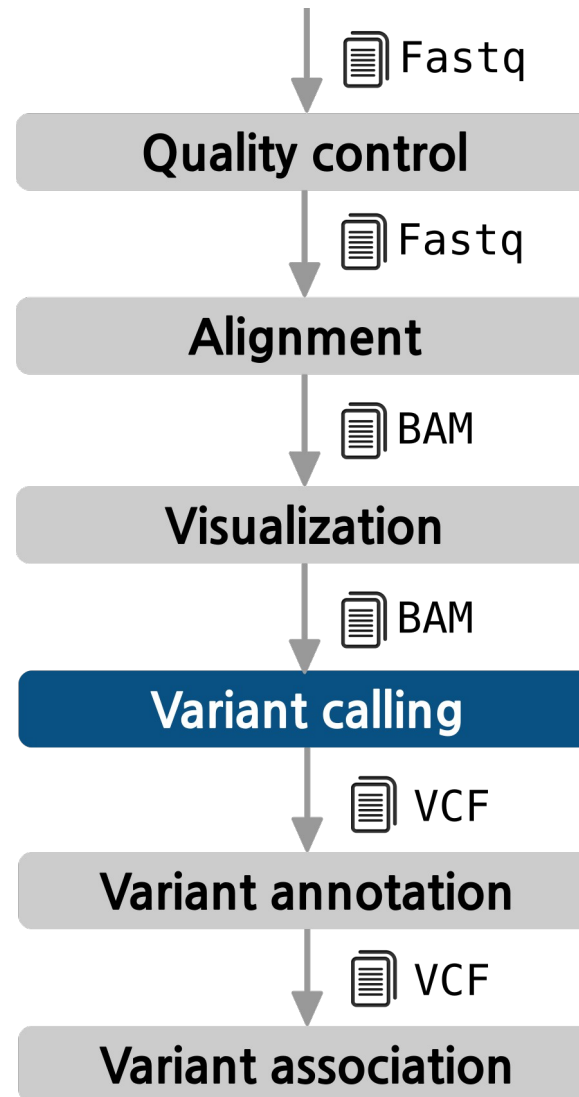
18<sup>th</sup> March 2016



Courtesy of Marta Bleda

# The pipeline

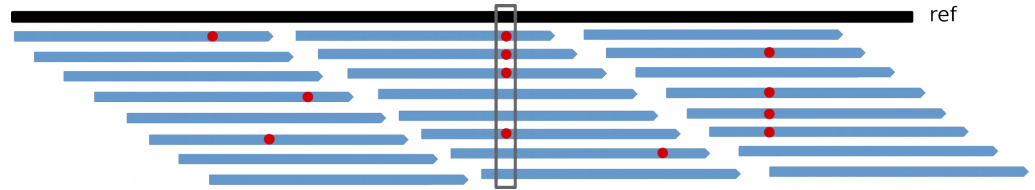
---



# Objective

---

Assign a genotype to each position



## Problems

Some variation observed in BAM files is caused by mapping and sequencing artifacts:

- **PCR artifacts:**
  - Mismatches due to errors in early PCR rounds
  - PCR duplicates
- **Sequencing errors:** erroneous call, either for physical reasons or to properties of the sequenced DNA
- **Mapping errors:** often happens around repeats or other low-complexity regions

Separate **true variation** from machine artifacts

# Variant calling process pipeline

---

## 1. Mark duplicates

Duplicates should not be counted as additional evidence

## 2. Local realignment around INDELS

Reads mapping on the edges of INDELS often get mapped with mismatching bases introducing false positives

## 3. Base quality score recalibration (BQSR)

Quality scores provided by sequencing machines are generally inaccurate and biased

## 4. Variant calling

Discover variants and their genotypes

# 1. Mark duplicates

---

- The same DNA molecule can be **sequenced several times during PCR**
- **Not informative**
- **Not** to be counted as **additional evidence** for or against a putative variant
- Can result in **false variant calls**

## Tools

- **Samtools:** `samtools rmdup`
- **Picard:** `MarkDuplicates`

# 1. Mark duplicates

---

- The same DNA molecule can be **sequenced several times during PCR**
- **Not informative**
- **Not** to be counted as **additional evidence** for or against a putative variant
- Can result in **false variant calls**

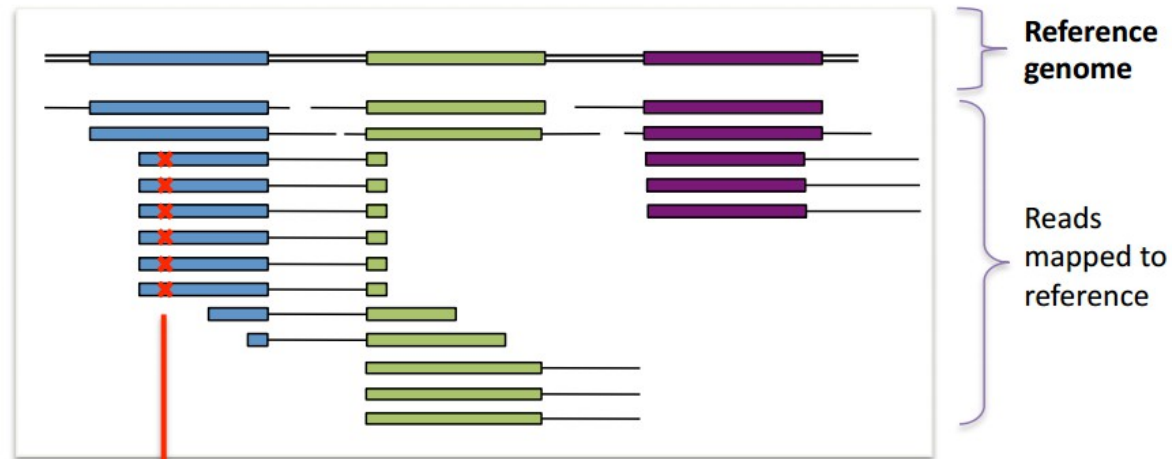
## Tools

- **Samtools:** `samtools rmdup`
- **Picard:** `MarkDuplicates`

# 1. Mark duplicates

The reason why duplicates are bad

✘ = sequencing error propagated in duplicates



After marking duplicates, the GATK will only see :

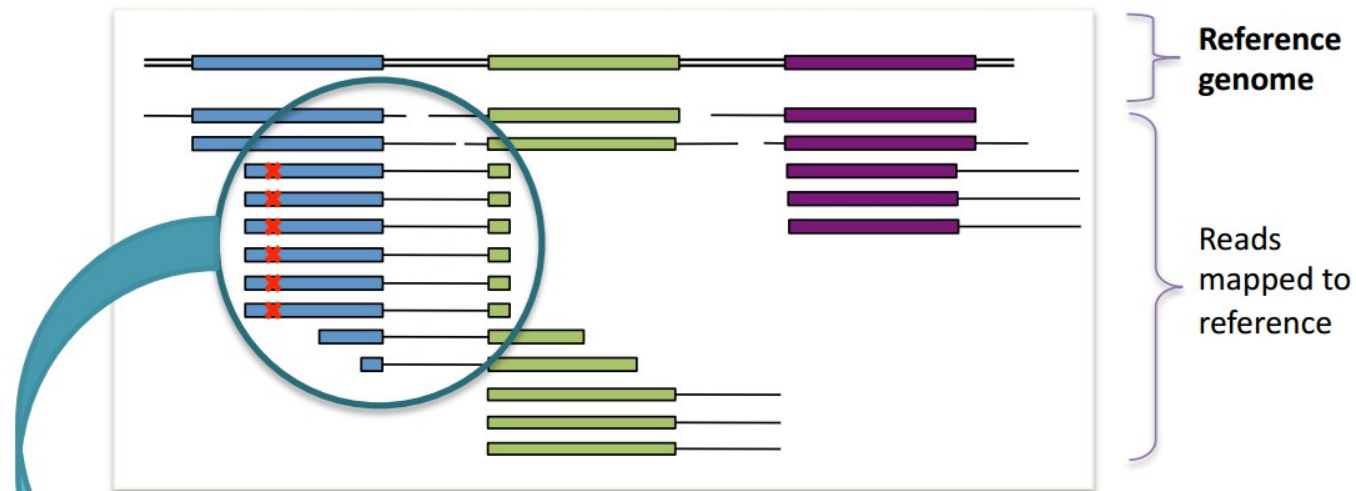


... and thus be more likely to make the right call

# 1. Mark duplicates

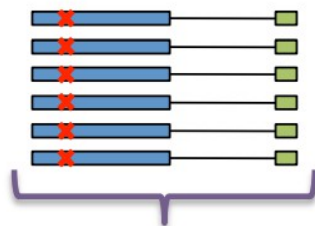
## Duplicate identification

Duplicates have the **same starting position** and the **same CIGAR string**



Reference genome

Reads mapped to reference



Easy to bag & tag

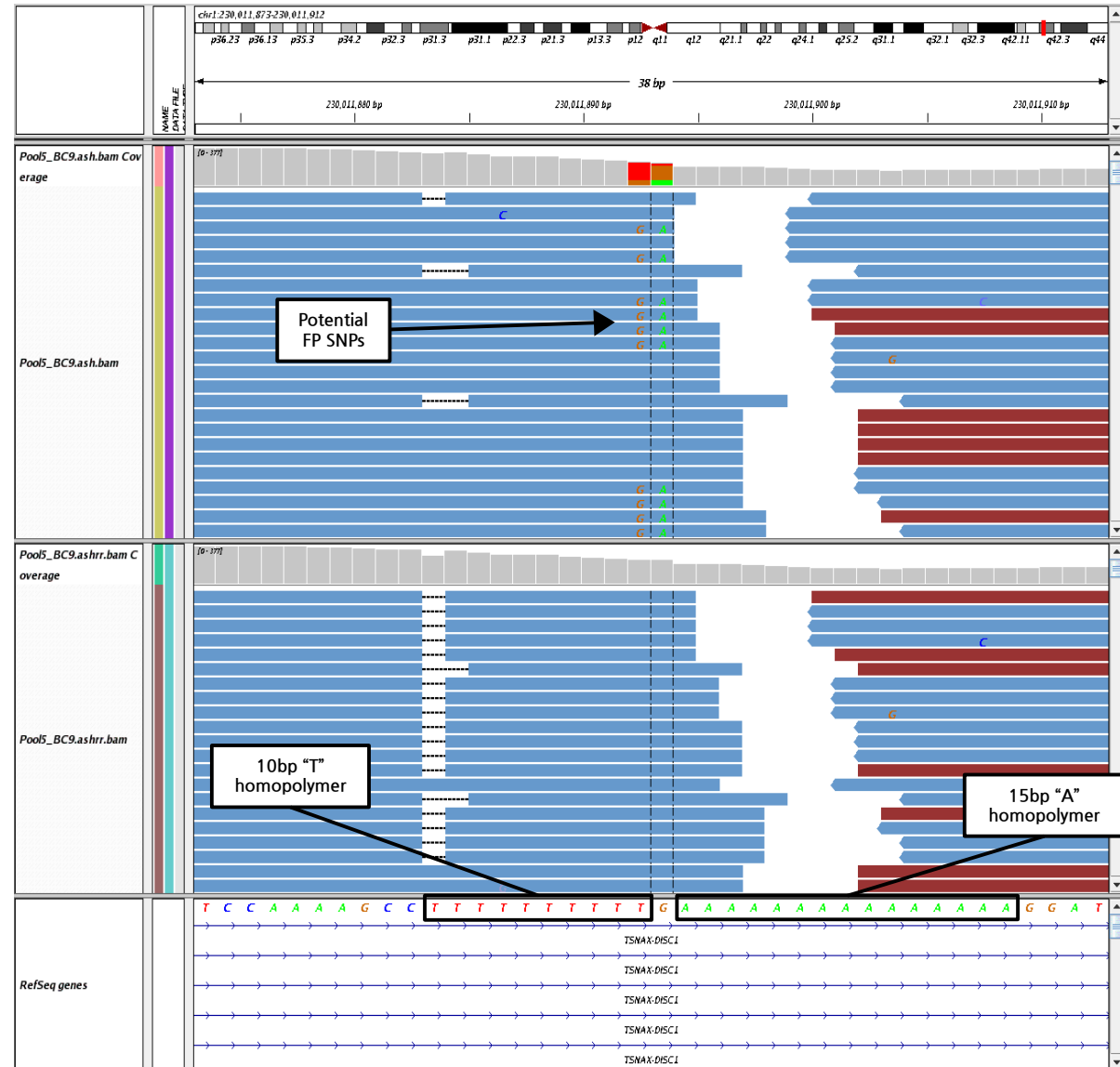
```
POS: 340  
CIGAR: 42M1D38M3I18M
```

Hey, Picard has an app for that!



## 2. Local realignment around INDELS

- Reads **near INDELS** are mapped with mismatches
- **Realignment** can identify the most consistent placement for these reads
  1. **Identify** problematic regions
  2. **Determine the optimal** consensus sequence
- **Minimizes mismatches** with the reference sequence
- **Refines** location of **INDELS**



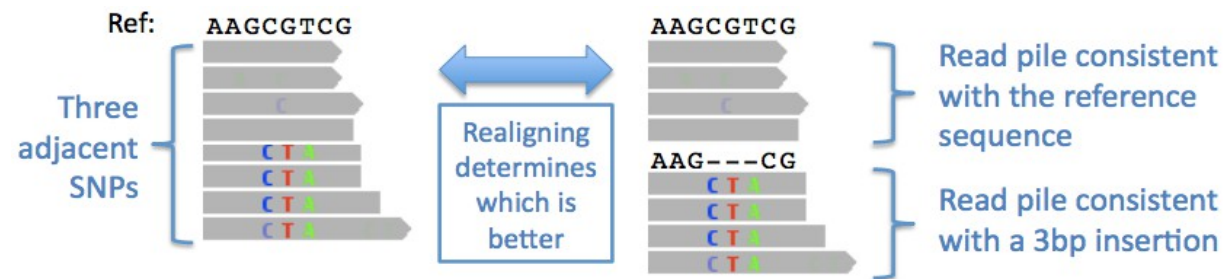
## 2. Local realignment around INDELS

- Reads **near INDELS** are mapped with mismatches
- **Realignment** can identify the most consistent placement for these reads

1. **Identify** problematic regions

2. **Determine the optimal** consensus sequence

- **Minimizes mismatches** with the reference sequence
- **Refines** location of **INDELS**



# 3. Base quality score recalibration

---

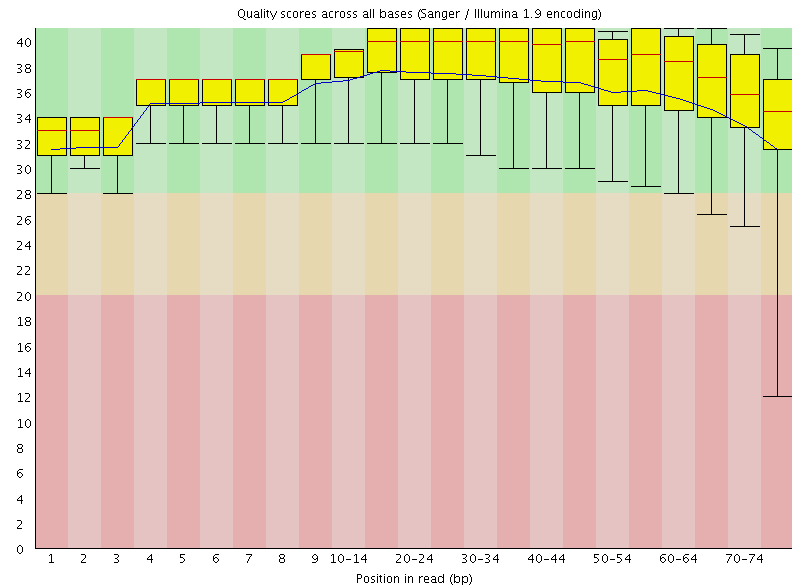
- **Calling algorithms rely** heavily on the **quality scores** assigned to the individual base calls in each sequence read
- Unfortunately, the scores produced by the machines are subject to various sources of **systematic error**, leading to over- or under-estimated base quality scores in the data

## How?

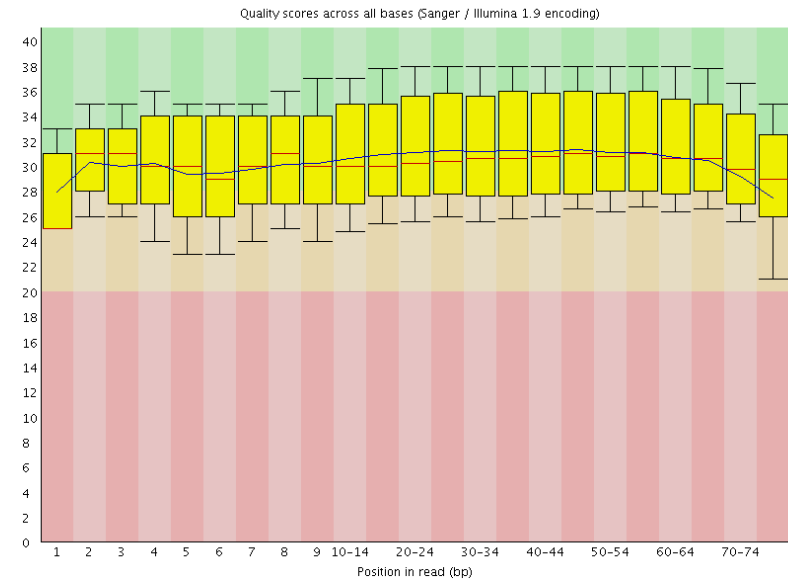
1. **Analyze covariation** among several features of a base:
  - Reported quality score
  - Position within the read
  - Preceding and current nucleotide
2. Use a set of **known variants** (i.e.: dbSNP) to model error properties of real polymorphism and determine the **probability that novel sites are real**
3. **Adjust** the quality scores of all reads in a BAM file

# 3. Base quality score recalibration

## Before



## After



Phred Quality score:

$$Q = -10 \log_{10} P(\text{error})$$

| Phred Quality Score | Probability of Incorrect Base Call | Base Call Accuracy |
|---------------------|------------------------------------|--------------------|
| 10                  | 1 in 10                            | 90%                |
| 20                  | 1 in 100                           | 99%                |
| 30                  | 1 in 1,000                         | 99.9%              |
| 40                  | 1 in 10,000                        | 99.99%             |
| 50                  | 1 in 100,000                       | 99.999%            |

# 4. Variant calling

## Variant discovery process

---

### Steps

1. **Variant calling:** Identify the positions that differ from the reference
2. **Genotype calling:** calculate the genotypes for each sample at these sites

### Initial approach

**Independent** base assumption

Counting the number of times each allele is observed

### Evolved approach

**Bayesian inference** → Compute genotype likelihood

Advantages:

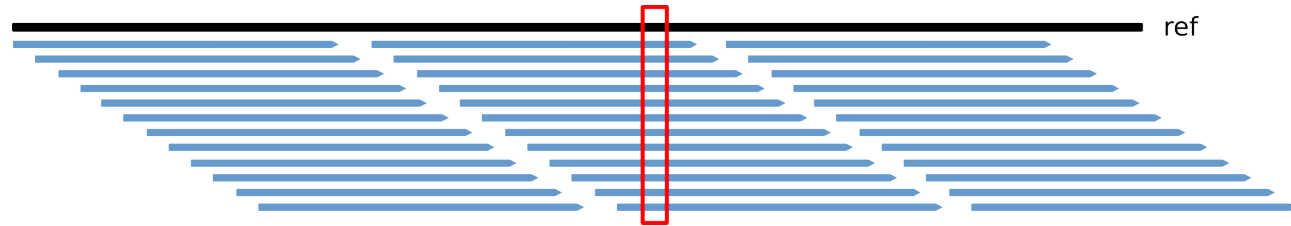
Provide statistical measure of **uncertainty**

Lead to **higher accuracy** of genotype calling

# 4. Variant calling

## Variant discovery process

---



Reference = A

# 4. Variant calling

## Variant discovery process

---



Reference =  $A$

AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0$

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV

# 4. Variant calling

## Variant discovery process



Reference = A

AAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0$

GGGGGGGGGGGGGGGGGGGGGGGGGG

$N=30, X=30$

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV



# 4. Variant calling

## Variant discovery process



Reference = **A**

AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0$

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

$N=30, X=30$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGG

$N=30, X=15$

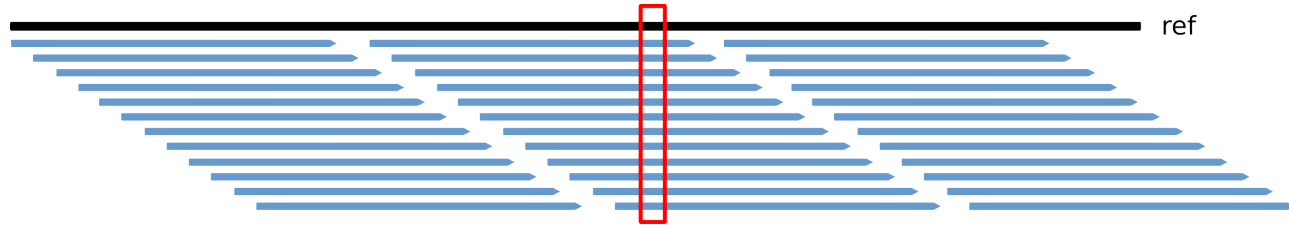
$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV



# 4. Variant calling

## Variant discovery process



Reference = **A**

AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0$

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

$N=30, X=30$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGGG

$N=30, X=15$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGCT

$N=30, X=12$

AAAGGGCCTT

$N=10, X=3$

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV

# 4. Variant calling

## Variant discovery process



Reference = **A**

AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0$

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

$N=30, X=30$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGGGGGGG

$N=30, X=15$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGGGCT

$N=30, X=12$

AAAGGGCCCTT

$N=10, X=3$

Cutoff for  $X \rightarrow$  value or proportion

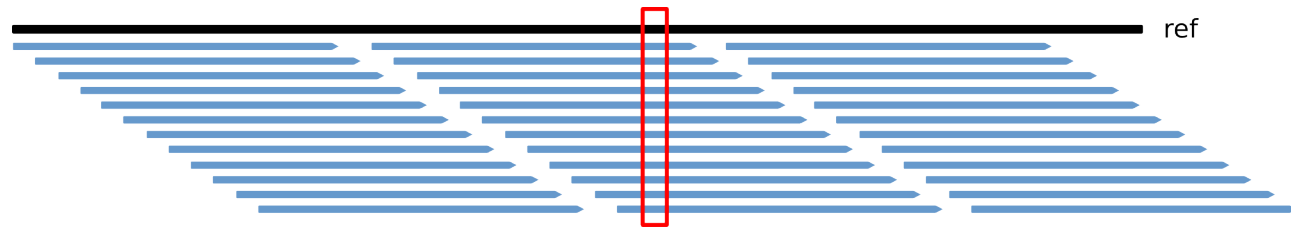
- $c_1 = 10\%, c_2 = 30\%$ 
  - $X \leq c_1 \rightarrow \mathbf{RR}$
  - $c_1 < X < c_2 \rightarrow \mathbf{RV}$
  - $X \geq c_2 \rightarrow \mathbf{VV}$

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
 $RR \quad RV \quad VV$

# 4. Variant calling

## Variant discovery process



Reference = **A**

AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

$N=30, X=0 \rightarrow \mathbf{RR}$

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

$N=30, X=30 \rightarrow \mathbf{VV}$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGGGGG

$N=30, X=15 \rightarrow \mathbf{RV}$

AAAAAAAAAAAAAAAAAGGGGGGGGGGGGGGGCT

$N=30, X=12 \rightarrow \mathbf{RV}$

AAAGGGCCTT

$N=10, X=3 \rightarrow \mathbf{RV?}$

Cutoff for  $X \rightarrow$  value or proportion

- $c_1 = 10\%, c_2 = 30\%$ 
  - $X \leq c_1 \rightarrow \mathbf{RR}$
  - $c_1 < X < c_2 \rightarrow \mathbf{RV}$
  - $X \geq c_2 \rightarrow \mathbf{VV}$

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV

# 4. Variant calling

## Variant discovery process



### Bayesian approximation

$\alpha$  = nucleotide-base error rate

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:

RR RV VV

$P(G=RR, X|N, \alpha)$  = P of all R calls being correct and all V calls being wrong

$P(G=VV, X|N, \alpha)$  = P of all V calls being correct and all R calls being wrong

$P(G=RV, X|N, \alpha)$  = P of all R and V calls being correct

# 4. Variant calling

## Variant discovery process



### Bayesian approximation

$\alpha$  = nucleotide-base error rate

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:

RR RV VV

$$P(G=RR, X|N, \alpha) = \binom{N}{X} \alpha^X (1-\alpha)^{N-X}$$

$$P(G=VV, X|N, \alpha) = \binom{N}{X} (1-\alpha)^X \alpha^{N-X}$$

$$P(G=RV, X|N, \alpha) = \binom{N}{X} \left(\frac{1}{2}\right)^N$$

# 4. Variant calling

## Variant discovery process



### Bayesian approximation

$\alpha$  = nucleotide-base error rate

$p_{VV}$  }  
 $p_{VR}$  } Prior probabilities

$N$  = nucleotides  
 $G$  = true genotype  
 $R$  = reference base  
 $V$  = variant base  
 $X$  = variant nucleotides

Outcomes:  
RR RV VV

$$P(G = RR, X|N, \alpha) = \binom{N}{X} \alpha^X (1 - \alpha)^{N - X} (1 - p_{VV} - p_{RV})$$

$$P(G = VV, X|N, \alpha) = \binom{N}{X} (1 - \alpha)^X \alpha^{N - X} p_{VV}$$

$$P(G = RV, X|N, \alpha) = \binom{N}{X} \left(\frac{1}{2}\right)^N p_{RV}$$



# VCF file format

- Specification defined by the 1000 genomes (current version **4.2**):  
<http://www.1000genomes.org/wiki/Analysis/Variant%20Call%20Format/vcf-variant-call-format-version-41>
- Commonly **compressed and indexed** with bgzip/tabix
- Single-sample or multi-sample VCF

```
##fileformat=VCFv4.1
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTC G,GTCT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3
```

# VCF file format

```
#CHROM POS ID REF ALT QUAL FILTER INFO
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2

FORMAT NA00001 NA00002 NA00003
GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
```

genotype genotype quality read depth haplotype qualities

- **CHROM**: chromosome
- **POS**: position
- **ID**: identifier
- **REF**: reference base(s)
- **ALT**: non-reference allele(s)
- **QUAL**: quality score of the calls (phed scale)
- **FILTER**: "PASS" or a filtering tag
- **INFO**: additional information
- **FORMAT**: describes the information given by sample

# Software

| Software | Available from  | Calling method  | Prerequisites   | Comments   | Refs  |
|----------|---|-----------------|---|--|-------|
| SOAP2    | <a href="http://soap.genomics.org.cn/index.html">http://soap.genomics.org.cn/index.html</a>   | Single-sample   | High-quality variant database (for example, dbSNP)                  | Package for NGS data analysis, which includes a single individual genotype caller (SOAPSnp)  | 15    |
| realSFS  | <a href="http://128.32.118.212/thorfinn/realSFS/">http://128.32.118.212/thorfinn/realSFS/</a>   | Single-sample   | Aligned reads   | Software for SNP and genotype calling using single individuals and allele frequencies. Site frequency spectrum (SFS) estimation  | -     |
| Samtools | <a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>   | Multi-sample    | Aligned reads   | Package for manipulation of NGS alignments, which includes a computation of genotype likelihoods (samtools) and SNP and genotype calling (bcftools)  | 53    |
| GATK     | <a href="http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit">http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit</a> | Multi-sample    | Aligned reads   | Package for aligned NGS data analysis, which includes a SNP and genotype caller (Unified Genotyper), SNP filtering (Variant Filtration) and SNP quality recalibration (Variant Recalibrator)   | 32,33 |
| Beagle   | <a href="http://faculty.washington.edu/browning/beagle/beagle.html">http://faculty.washington.edu/browning/beagle/beagle.html</a>                                       | Multi-sample LD | Candidate SNPs, genotype likelihoods                                | Software for imputation, phasing and association that includes a mode for genotype calling   | 42    |
| IMPUTE2  | <a href="http://mathgen.stats.ox.ac.uk/impute/impute_v2.html">http://mathgen.stats.ox.ac.uk/impute/impute_v2.html</a>   | Multi-sample LD | Candidate SNPs, genotype likelihoods                                | Software for imputation and phasing, including a mode for genotype calling. Requires fine-scale linkage map  | 44    |
| QCall    | <a href="ftp://ftp.sanger.ac.uk/pub/rd/QCALL">ftp://ftp.sanger.ac.uk/pub/rd/QCALL</a>   | Multi-sample LD | 'Feasible' genealogies at a dense set of loci, genotype likelihoods | Software for SNP and genotype calling, including a method for generating candidate SNPs without LD information (NLDA) and a method for incorporating LD information (LDA). The 'feasible' genealogies can be generated using Margarita ( <a href="http://www.sanger.ac.uk/resources/software/margarita">http://www.sanger.ac.uk/resources/software/margarita</a> ) | 54    |
| MaCH     | <a href="http://genome.sph.umich.edu/wiki/Thunder">http://genome.sph.umich.edu/wiki/Thunder</a>   | Multi-sample LD | Genotype likelihoods  | Software for SNP and genotype calling, including a method (GPT_Freq) for generating candidate SNPs without LD information and a method (thunder_glf_freq) for incorporating LD information   | -     |

A more complete list is available from <http://seqanswers.com/wiki/Software/list>. LD, linkage disequilibrium; NGS, next-generation sequencing.

# Software

| Software | Available from  | Calling method  | Prerequisites                                      | Comments   | Refs  |
|----------|---|-----------------|--|--|-------|
| SOAP2    | <a href="http://soap.genomics.org.cn/index.html">http://soap.genomics.org.cn/index.html</a>   | Single-sample   | High-quality variant database (for example, dbSNP) | Package for NGS data analysis, which includes a single individual genotype caller (SOAPSnp)  | 15    |
| realSFS  | <a href="http://128.32.118.212/thorfinn/realSFS/">http://128.32.118.212/thorfinn/realSFS/</a>   | Single-sample   | Aligned reads                                      | Software for SNP and genotype calling using single individuals and allele frequencies. Site frequency spectrum (SFS) estimation  | -     |
| Samtools | <a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>   | Multi-sample    | Aligned reads                                      | Package for manipulation of NGS alignments, which includes a computation of genotype likelihoods (samtools) and SNP and genotype calling (bcftools)  | 53    |
| GATK     | <a href="http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit">http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit</a> | Multi-sample    | Aligned reads                                      | Package for aligned NGS data analysis, which includes a SNP and genotype caller (Unified Genotyper), SNP filtering (Variant Filtration) and SNP quality recalibration (Variant Recalibrator)   | 32,33 |
| Beagle   | <a href="http://faculty.washington.edu/browning/beagle/beagle.html">http://faculty.washington.edu/browning/beagle/beagle.html</a>                                       | Multi-sample LD | Candidate SNPs, genotype likelihoods               | Software for imputation, phasing and association that includes a mode for genotype calling   | 42    |
| IMPUTE2  | <a href="http://mathgen.stats.ox.ac.uk/impute/impute_v2.html">http://mathgen.stats.ox.ac.uk/impute/impute_v2.html</a>   | Multi-sample LD | Candidate SNPs, genotype likelihoods               | Software for imputation and phasing, including a mode for genotype calling. Requires fine-scale linkage map  | 44    |
| QCall    | <a href="ftp://ftp.sanger.ac.uk/pub/">ftp://ftp.sanger.ac.uk/pub/</a>   | Multi-sample LD | 'Feasible'   | Software for SNP and genotype calling, including a method for generating candidate SNPs without LD information (NLDA) and a method for incorporating LD information (LDA). The 'feasible' genealogies can be generated using Margarita ( <a href="http://www.sanger">http://www.sanger</a> | 54    |

## Platypus:

Andy Rimmer, Hang Phan, Iain Mathieson, Zamin Iqbal, Stephen R. F. Twigg, WGS500 Consortium, Andrew O. M. Wilkie, Gil McVean, Gerton Lunter. **Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications.** *Nature Genetics* (2014) doi:10.1038/ng.3036

## FreeBayes:

Garrison, E. & Marth, G. **Haplotype-based variant detection from short-read sequencing.** arXiv <http://arxiv.org/abs/1207.3907> (2012).

# Prerequisites: JAVA and Picard tools

---

- **Requires Java** (<http://www.oracle.com/technetwork/java/javase/downloads/index.html>)
  - Check your java version

```
java -version
```

GATK  $\geq$  2.6 → Requires Java version 1.7

- **Picard (current version 1.130)**

- Website: <http://broadinstitute.github.io/picard/>
- For a compiled version click on “Latest Release” and download [picard-tools-1.130.zip](#)
- Testing:



```
java -jar picard.jar -h
```

- Usage




```
java -jar picard.jar <ToolName> [options]
```

# Samtools installation

---

- Samtools 1.2 download

<http://www.htslib.org/download/>

- Download  [samtools-1.2](#)  [bcftools-1.2](#)  [htslib-1.2.1](#)
- Uncompress each of the files and inside the uncompressed folder execute:

```
make  
make install
```

- Check if Samtools is working

```
samtools
```

- Usage

```
samtools <command> [options]
```

# GATK installation

- GATK 3.3-0 download

<http://www.broadinstitute.org/gatk/>

- We need to register before download
- Go to Downloads and click
- Accept the license agreement
- Extract the file in the applications folder



You must be logged into the forums to proceed

You do not seem to be logged into the forums

Register

Login Here »

- Check if GATK is working

```
java -jar GenomeAnalysisTK.jar -h
```

Show GATK help

- Usage

```
java -jar GenomeAnalysisTK.jar -T <ToolName> [arguments]
```

# Filtering recommendations

---

## Filtering recommendations for SNPs:

- $QD < 2.0$
- $MQ < 40.0$
- $FS > 60.0$
- $HaplotypeScore > 13.0$
- $MQRankSum < -12.5$
- $ReadPosRankSum < -8.0$

## Filtering recommendations for indels:

- $QD < 2.0$
- $ReadPosRankSum < -20.0$
- $InbreedingCoeff < -0.8$
- $FS > 200.0$