

# PO: Precision Oncology Course

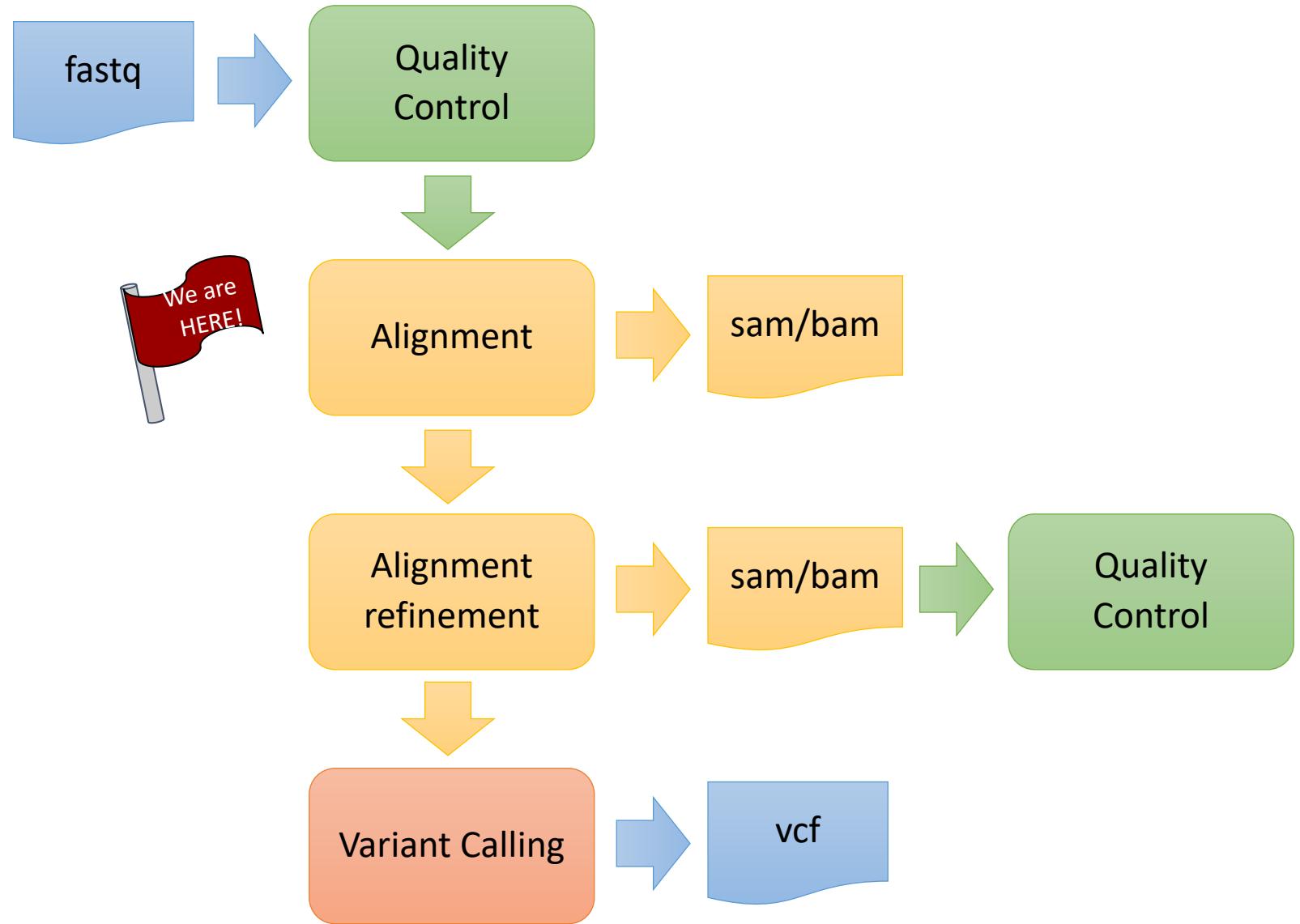
## Variant Detection

# Outline

- Steps for variant detection after alignment
- Algorithms for variant calling
- Pipeline for targeted DNA sequencing: varca

Steps for variant  
detection after alignment

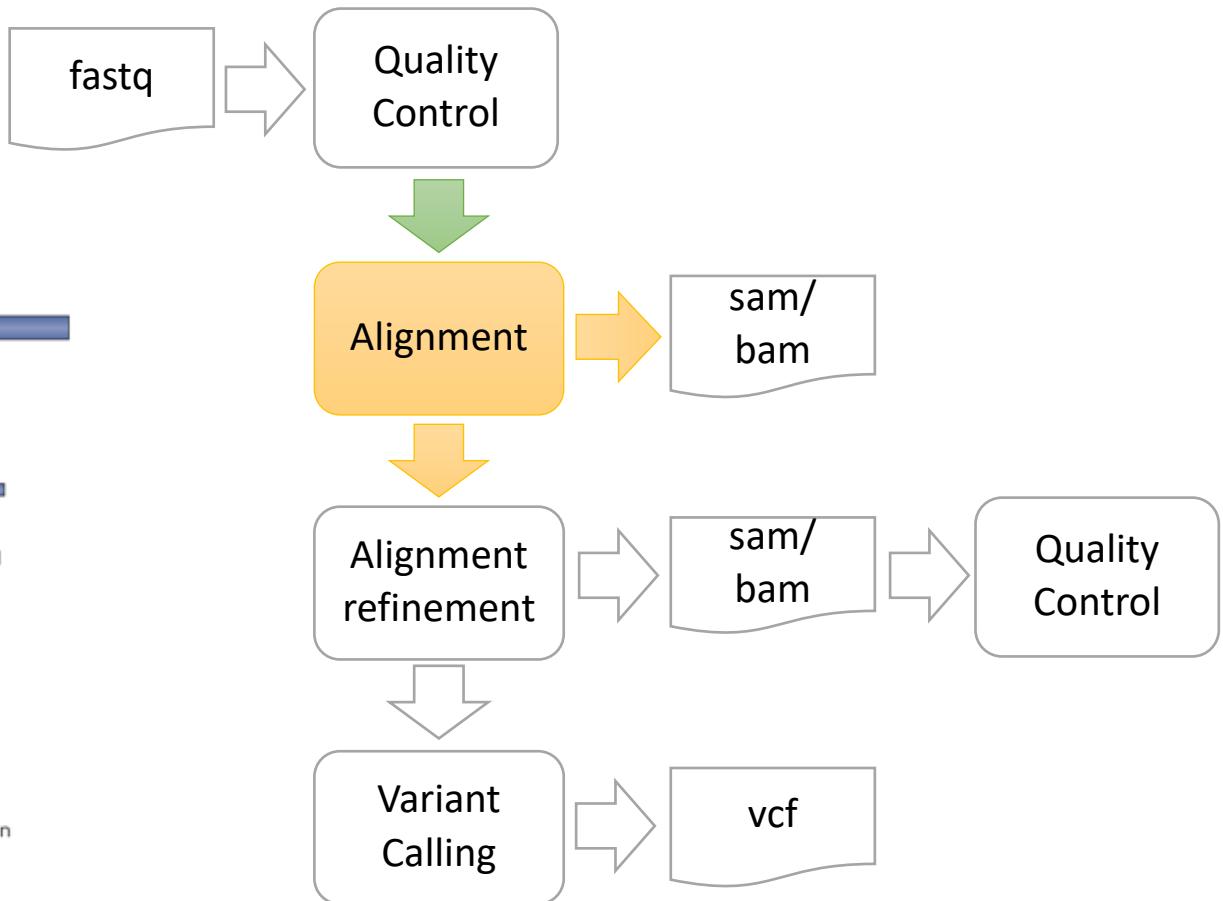
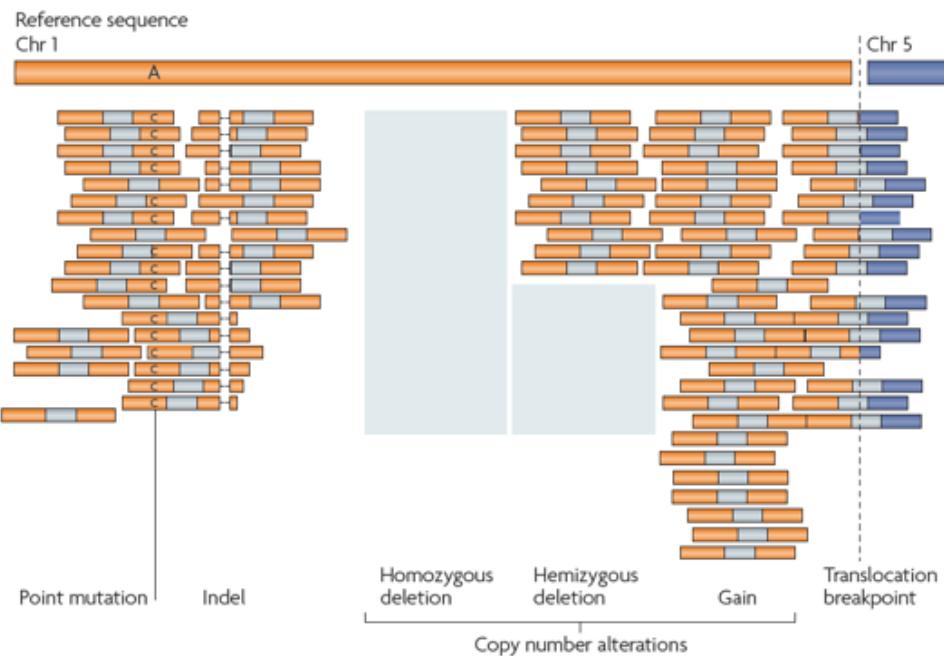
# General steps



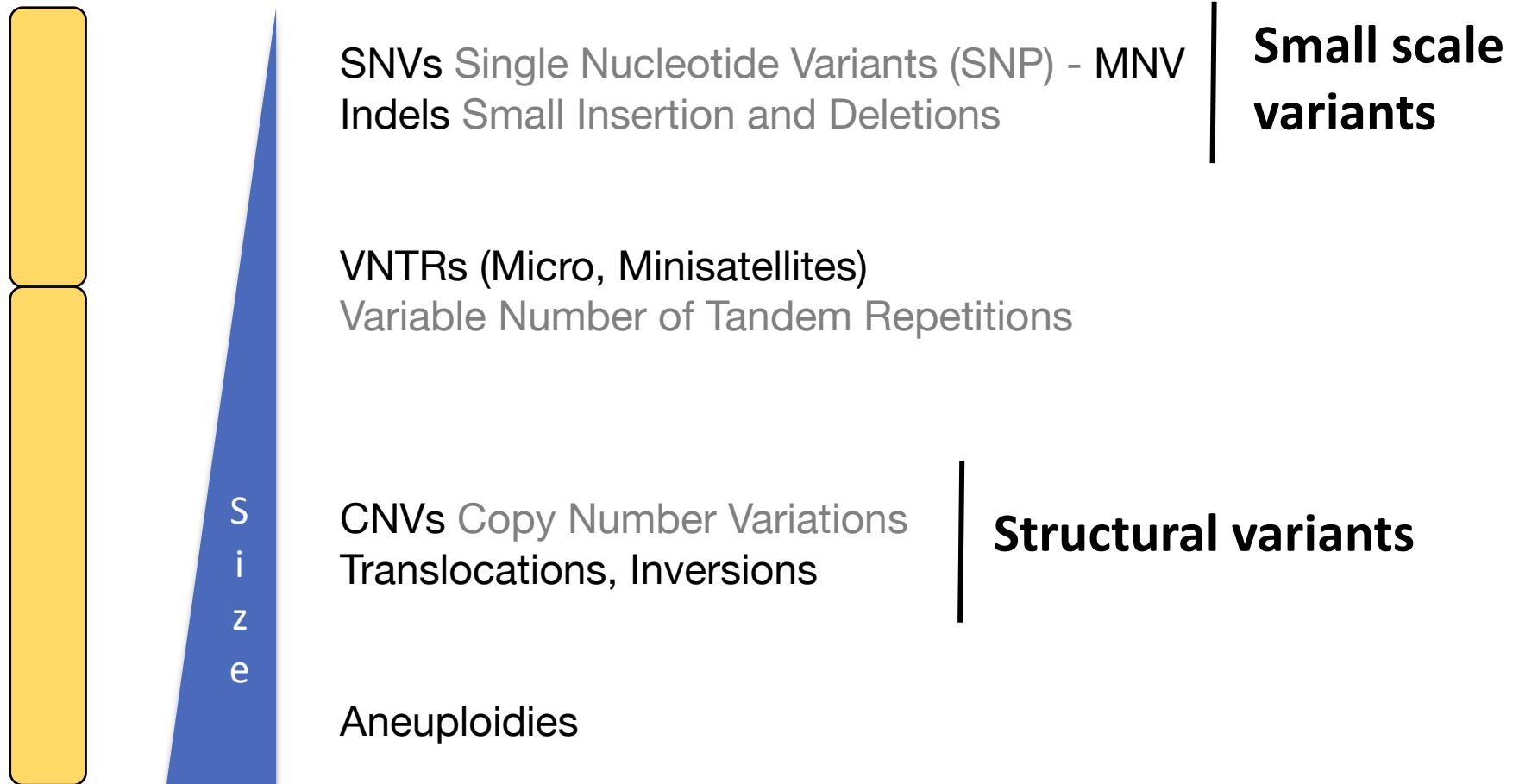
# Alignment

Reference

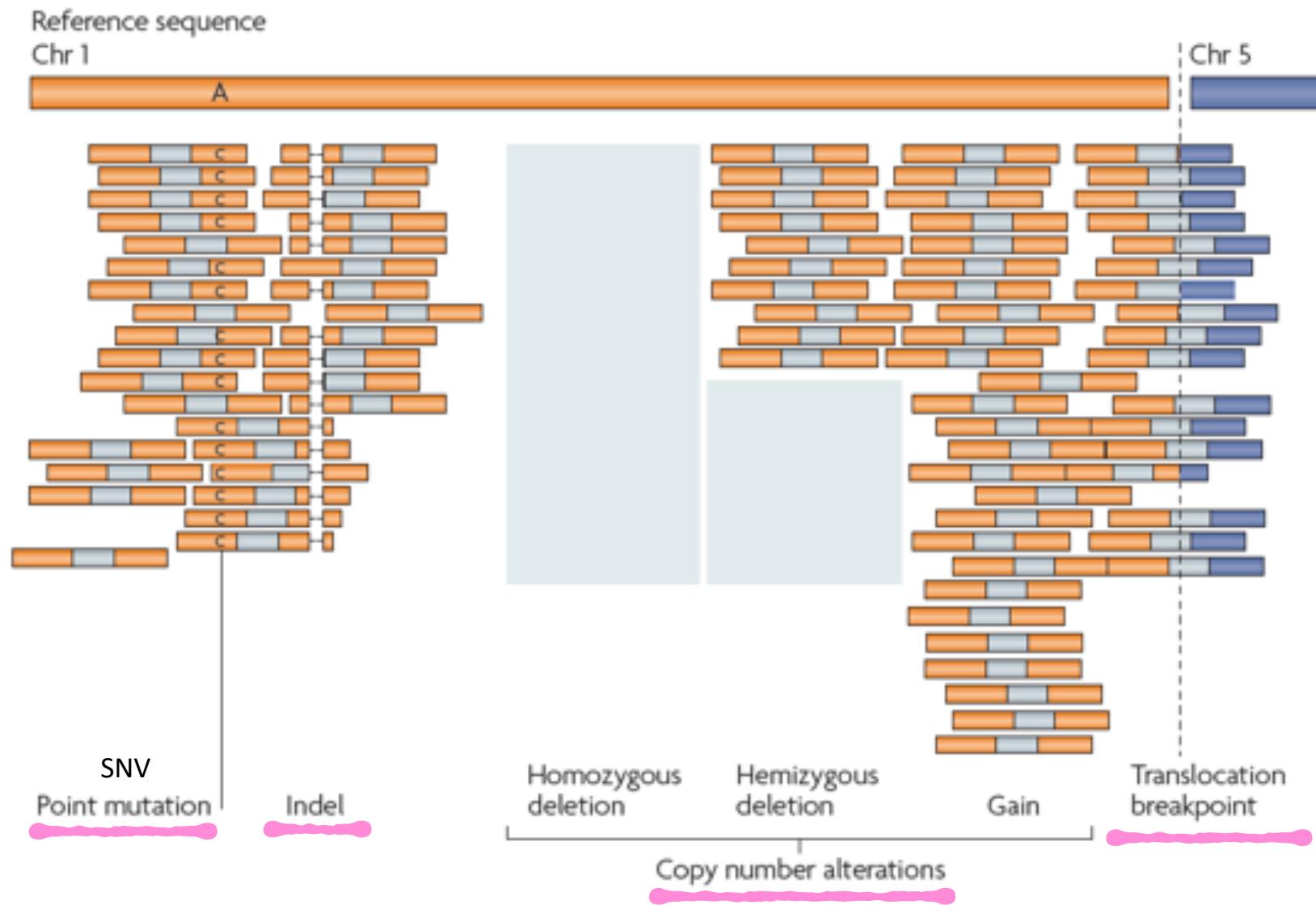
Sample



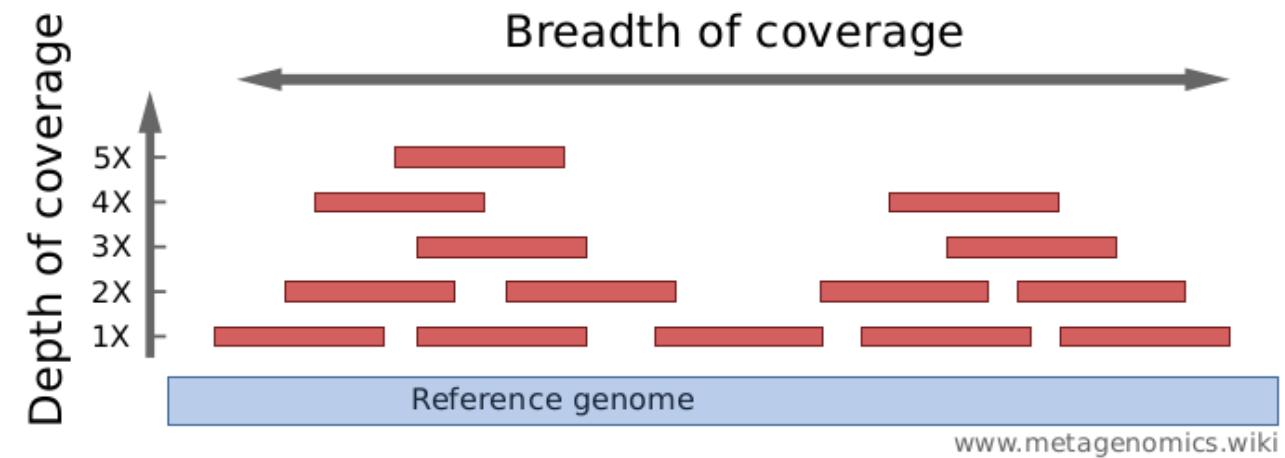
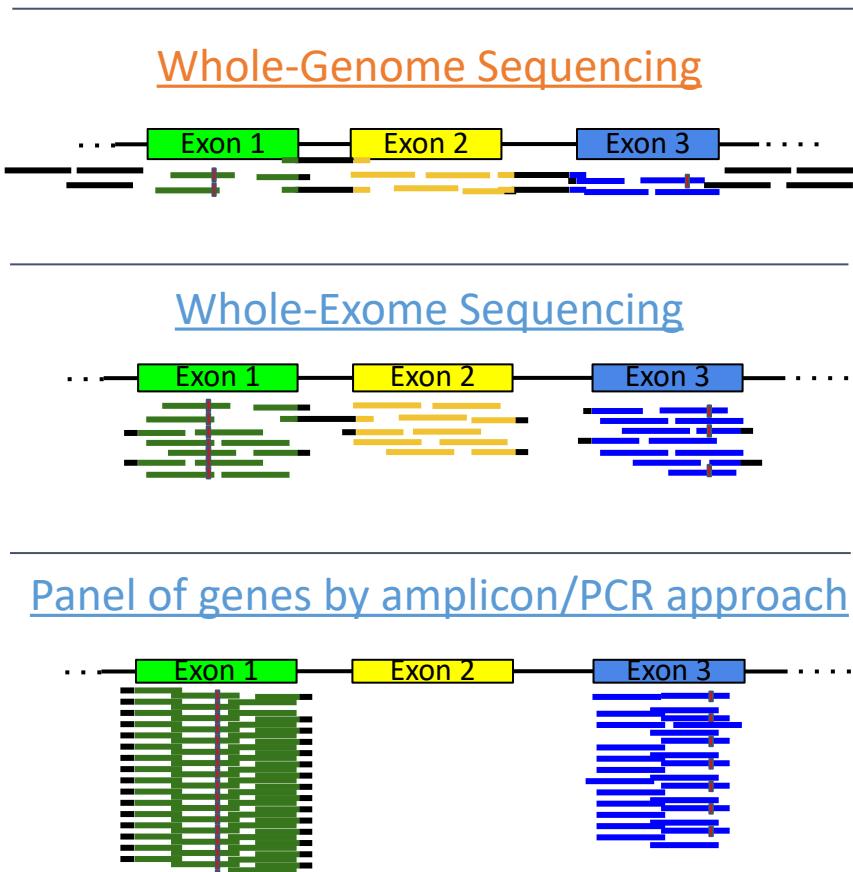
# Classification of variants according to size



# Alignment of reads uncovers potential variant sites



# DNA-seq strategies – Sequencing coverage



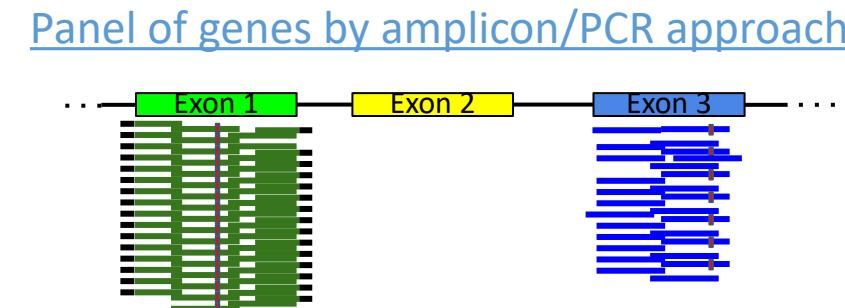
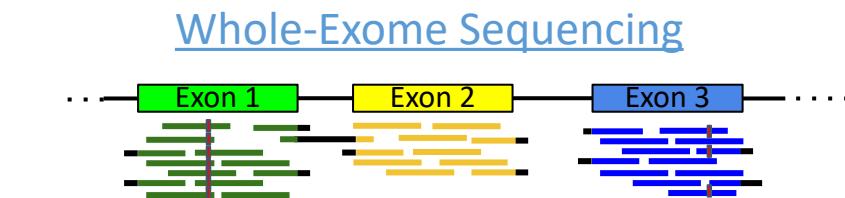
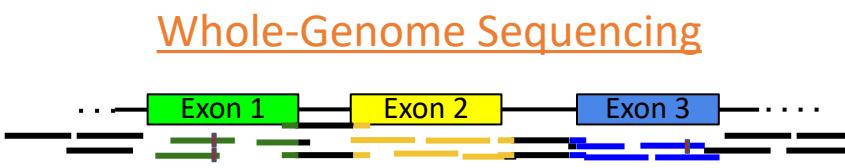
**Depth:** times that a base is sequenced

On average:

$$\frac{\sum \text{Number of times a sequenced base is covered by reads}}{\text{Length of the sequenced genome}}$$

**Breadth:** percentage of the sequenced genome covered by the reads (at a certain depth)

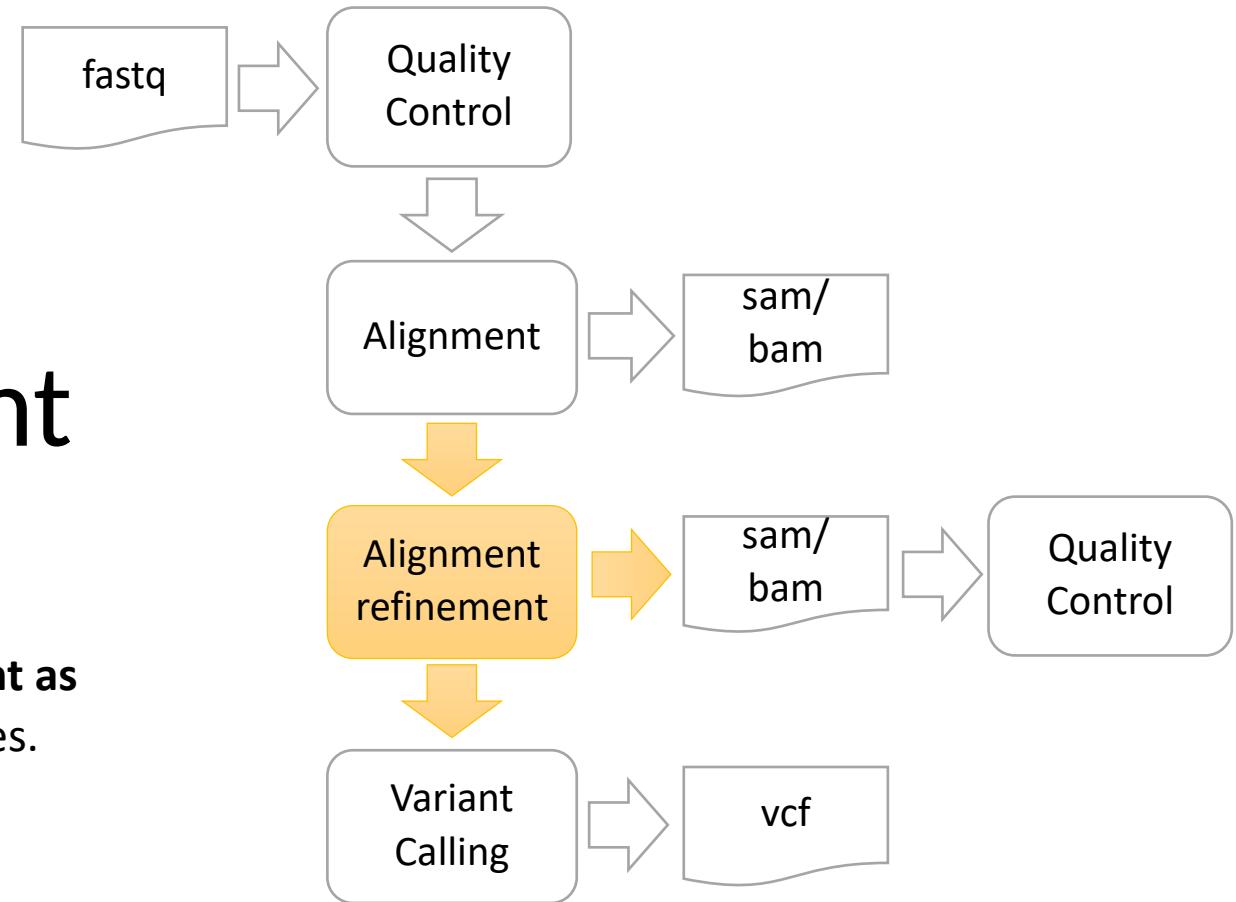
# DNA-seq strategies – Type of variants



Target	Type of variants discovered	Avg depth per pos	Cost
Entire genome	Coding variants, intronic and regulatory sites. Structural variants #Variants= 3M - 4M.	> 30x Most uniform	High
2% of the genome	Coding variants Some intronic and regulatory sites. Issues in the detection of structural variants #Variants= 20k - 60k.	> 50x - 100x	Low
Variable	Depends on the design (customizable) Challenging detection of structural variants # variants = ND	> 500x	Lower

# Alignment refinement

**Variant calling requires the most perfect alignment as possible** to avoid False Positives and False Negatives.



# Mark/remove duplicates

- Duplicates derive from **PCR amplification** (library preparation): one fragment is sequenced multiple times.
- An error at the beginning of the PCR (first steps) is propagated.
- Duplicates in hybrid-seq are **worthless** for the subsequent analysis:  
*Duplicates are source of False Positives calls while only provide redundancy.*

**Solution: retrieve the best one, discard the duplicates:**

Duplicates share the  
same alignment  
properties : sequence,  
start and end positions



\* = sequencing error propagated in duplicates

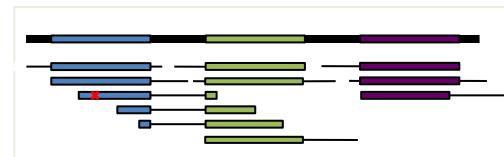
**METHOD:** by Picard-tools

[http://broadinstitute.github.io/  
picard/](http://broadinstitute.github.io/picard/)

(alternatives : samtools)



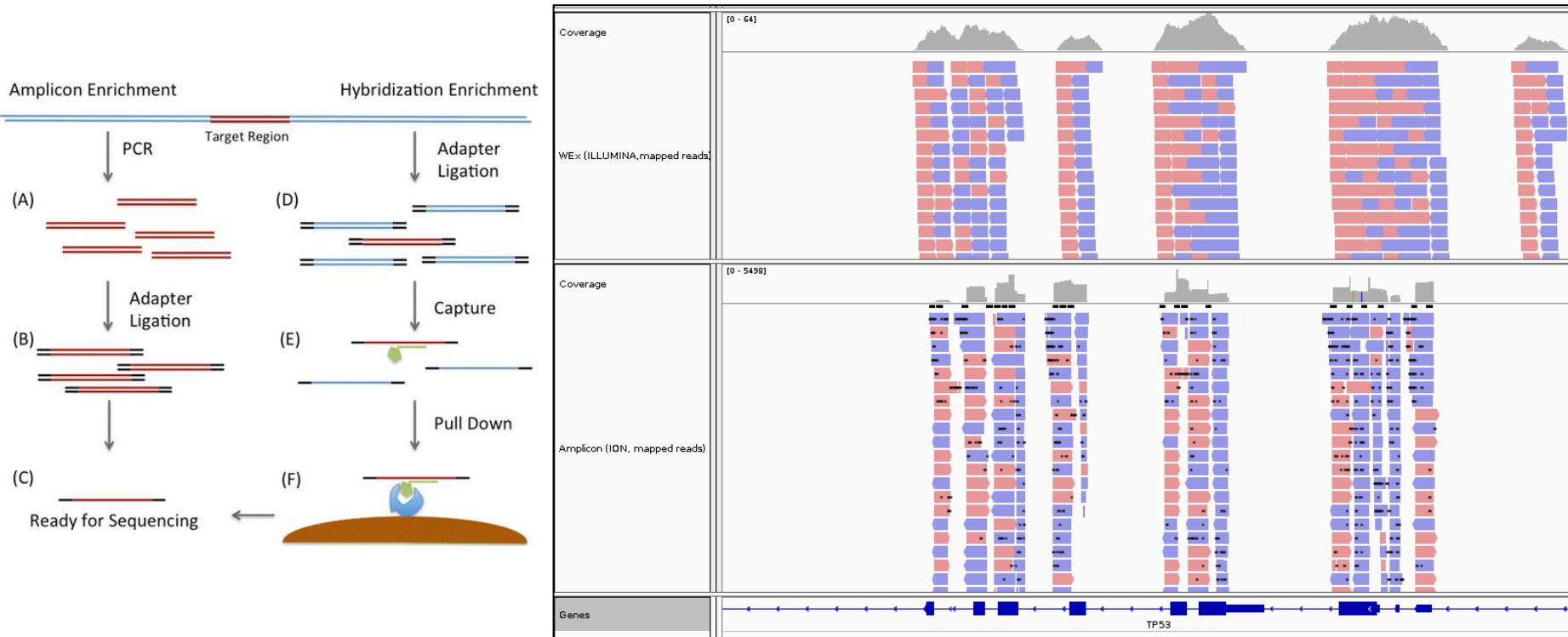
After marking/removing duplicates, the variant caller will only see :



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

Adapted from GATK

# Mark/remove duplicates: Amplicon seq



**WARNING: Do NOT remove duplicates in data derived from amplicon techniques (**Ion Torrent**).**

More info.: [https://github.com/broadgsa/gatk/blob/master/doc\\_archive/tutorials/\(How\\_to\)\\_Mark\\_duplicates\\_with\\_MarkDuplicates\\_or\\_MarkDuplicatesWithMateCigar.md](https://github.com/broadgsa/gatk/blob/master/doc_archive/tutorials/(How_to)_Mark_duplicates_with_MarkDuplicates_or_MarkDuplicatesWithMateCigar.md)

# Indel realignment

- Algorithms align reads very fast with high accuracy, but not perfectly.

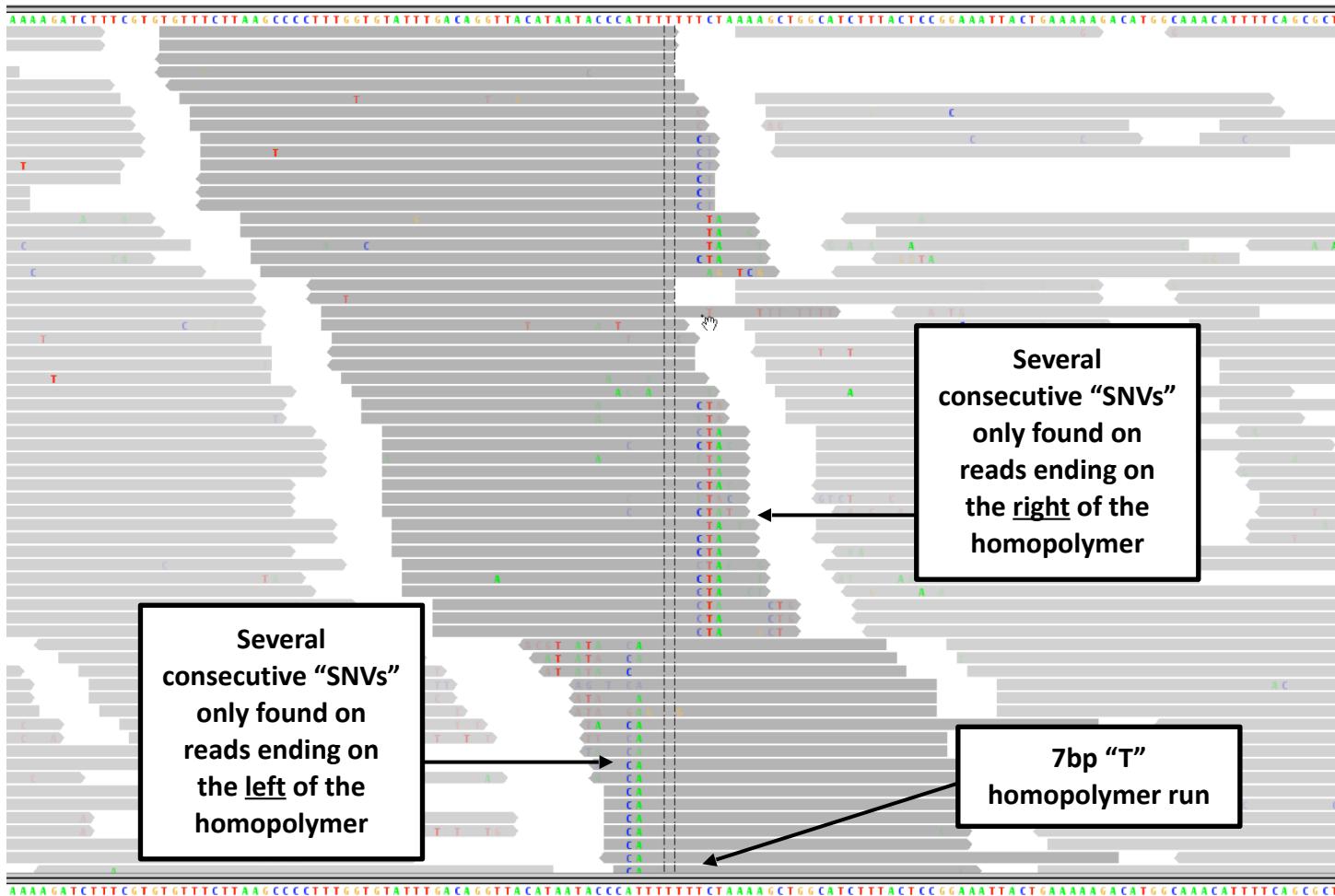
*During alignment, **penalties on mismatches are much cheaper than gaps (indels)**.*

- Also, there are sometimes multiple solutions (alignments) for a given read. **Aligners can choose one randomly.**
- Reads are aligned separately (one by one).
- **Indels can be no properly identified** in the alignment of the read.

**METHOD:** by GATK

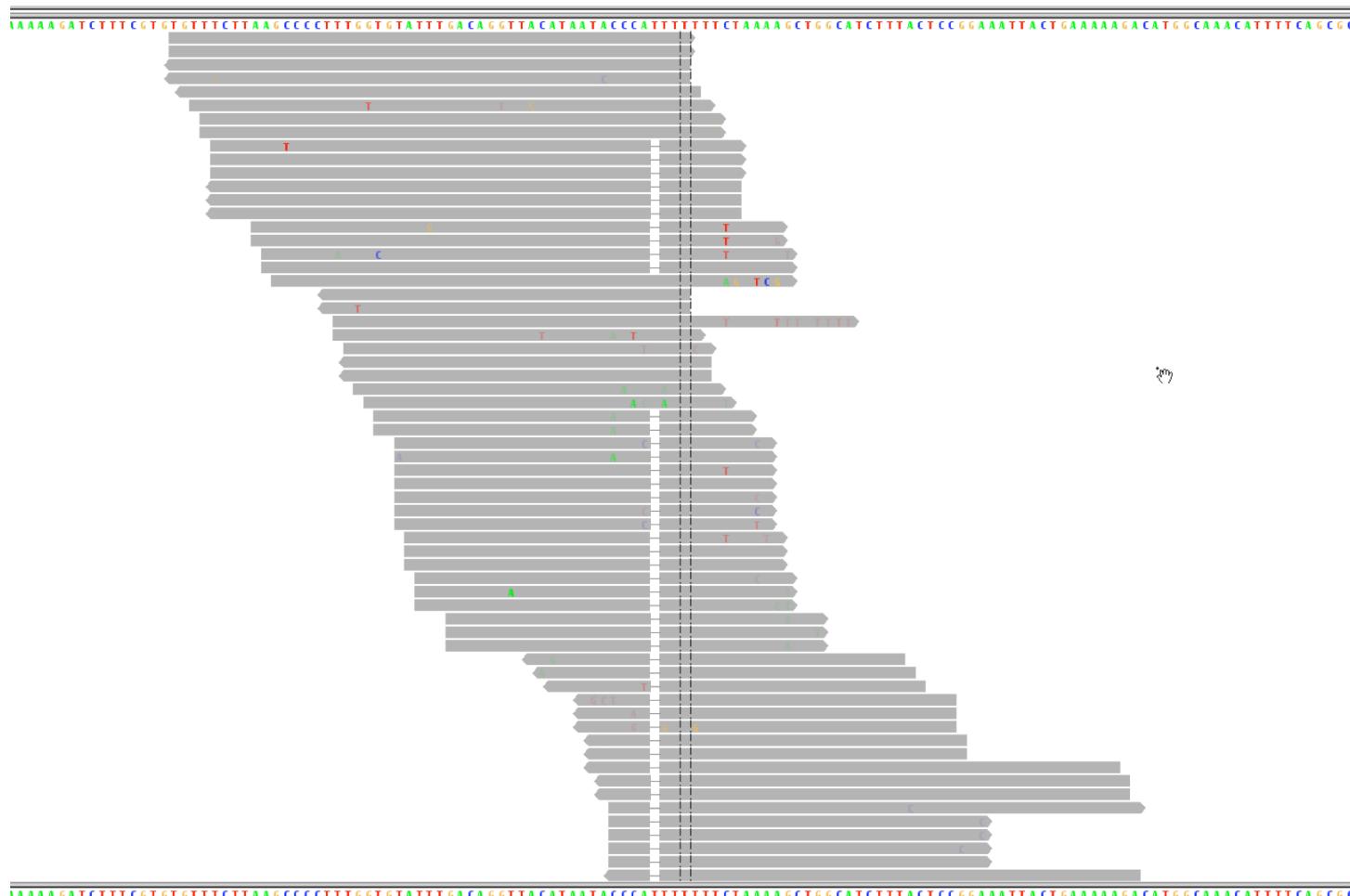
[https://github.com/broadinstitute/gatk-docs/blob/master/gatk3-tutorials/\(howto\)\\_Perform\\_local\\_realignment\\_around\\_indels.md](https://github.com/broadinstitute/gatk-docs/blob/master/gatk3-tutorials/(howto)_Perform_local_realignment_around_indels.md)

# Before indel realignment



Taken from GATK team

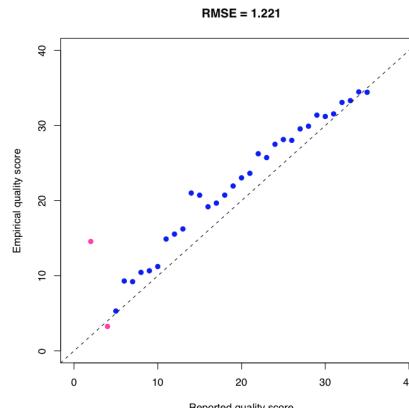
# After indel realignment



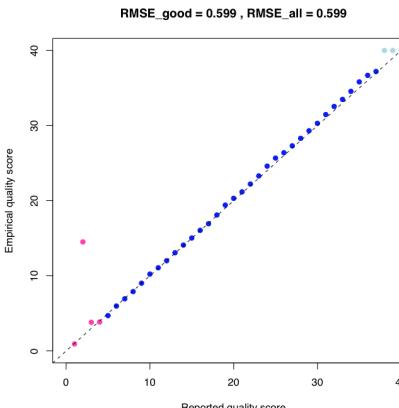
Taken from GATK team

# Base Quality Score Recalibration

- **Phred Quality score:** each position of the sequence has its particular **base quality score**.
- The individual quality measures are crucial during **Variant calling**.
- Different NGS technologies have their particular **bias in Quality Score** depending on the context. Recalibration **correct empirically** these biases.



Original

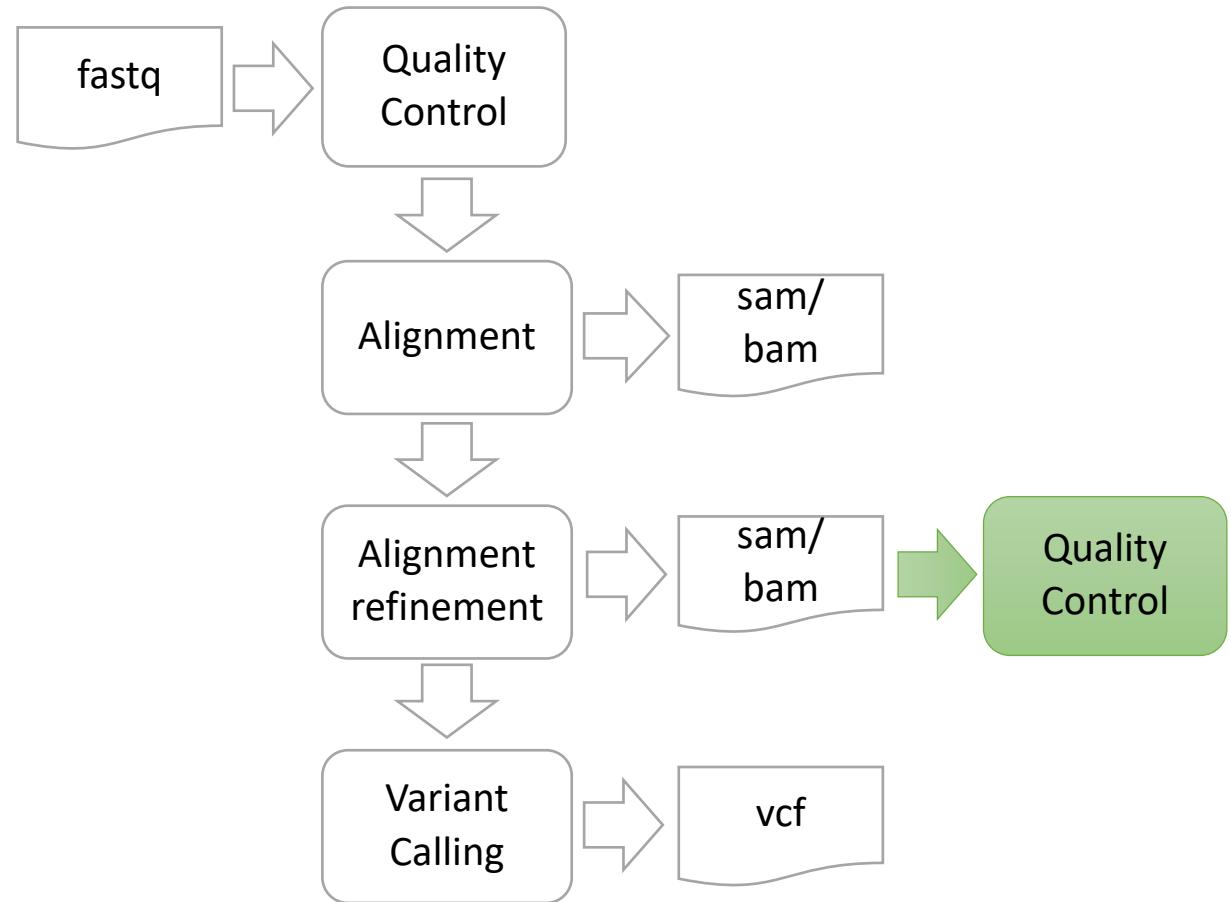


After BQSR recalibration

**METHOD:** by GATK

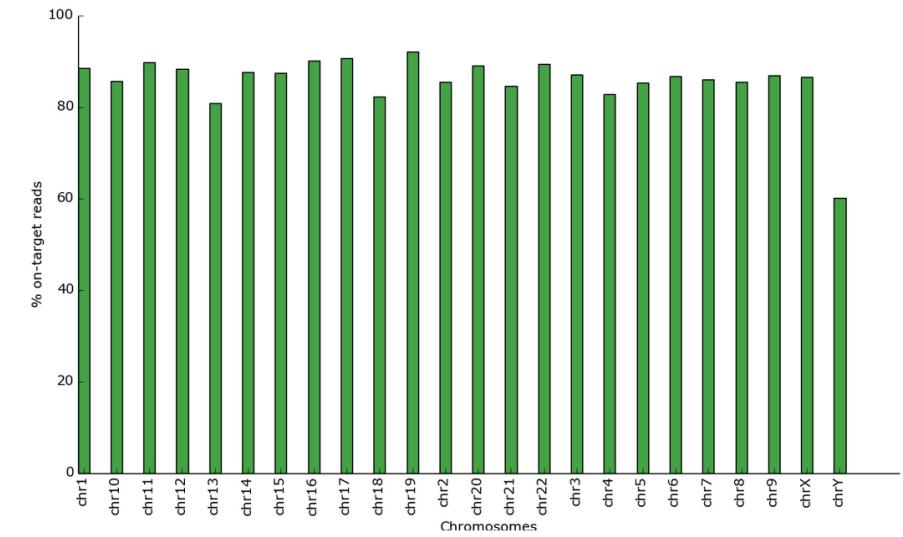
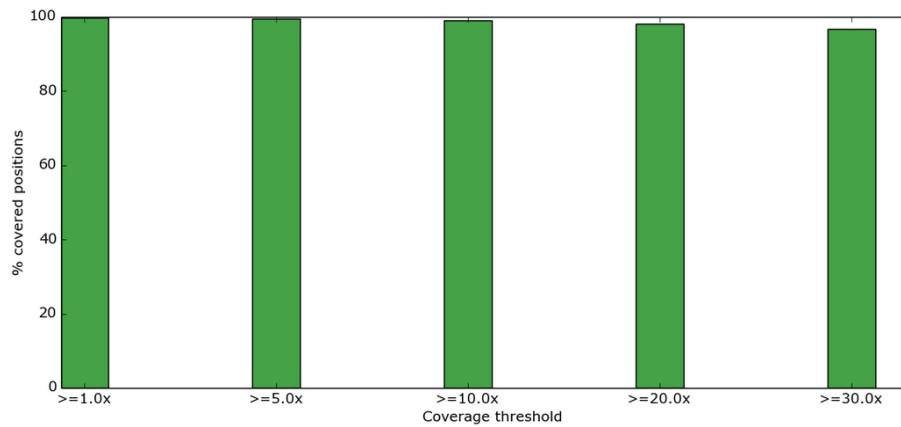
<https://gatk.broadinstitute.org/hc/en-us/articles/360035890531-Base-Quality-Score-Recalibration-BQSR->

# Alignment Quality Control



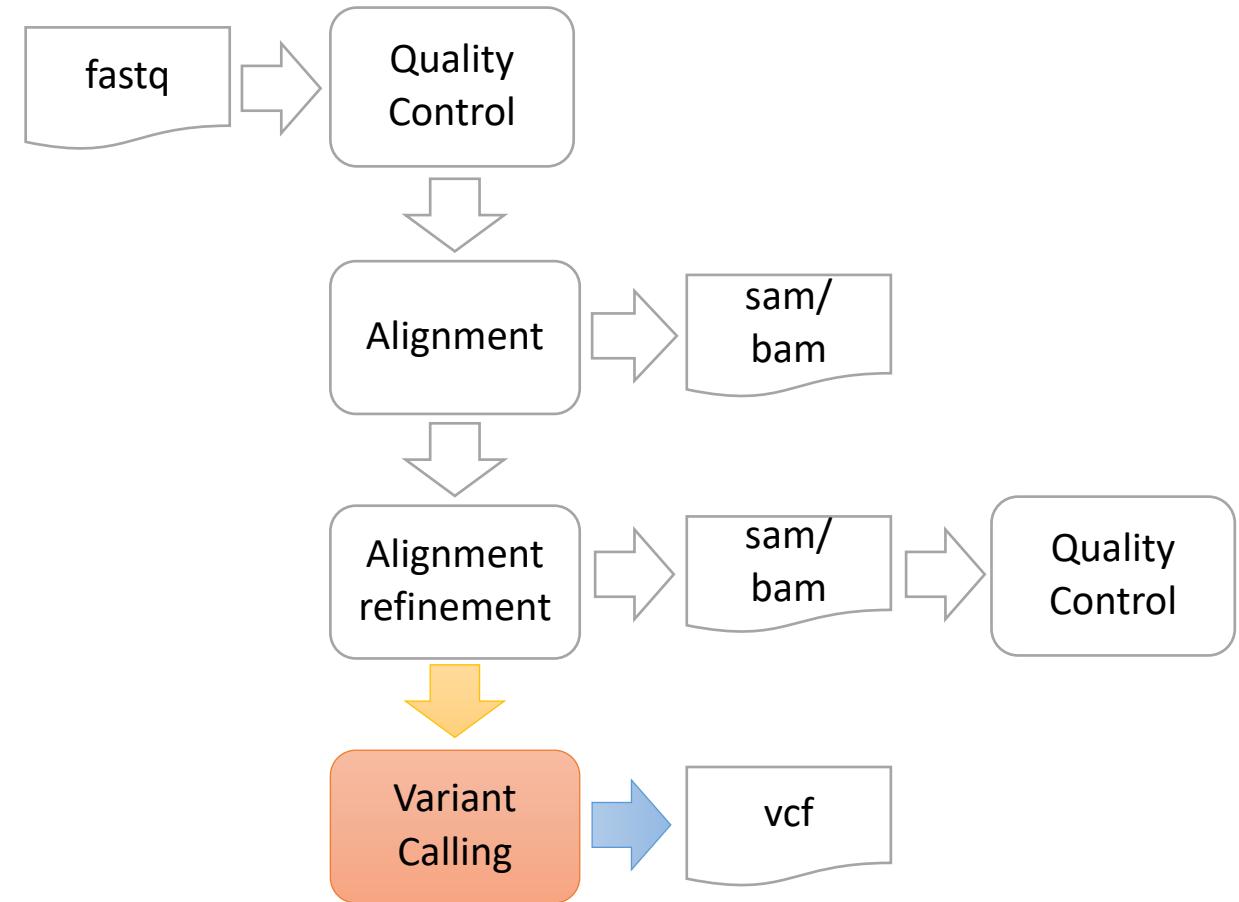
# Alignment Quality Control

- Mean sequencing depth
- Is there enough coverage in regions of interest?
- Are the reads on-target?

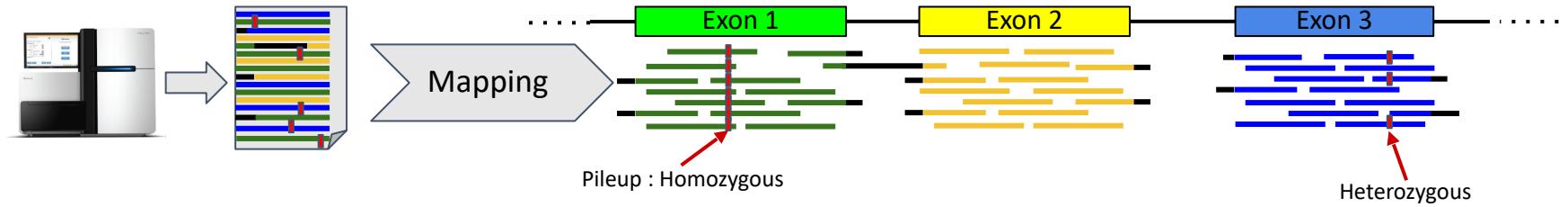


- Examples of **software**: ngsCAT, QualiMap

# Variant calling



# Fundamentals of Variant Calling



1

Identify the most likely genotype for each genomic position using statistical methods.

2

Identify the differences by comparing with the reference genome.

# What is Crucial in Variant calling

- For clinical practices, the use of **gold standard methods and reproducible analysis** are mandatory.
- The analysis is based on the comparison against the reference genome:

*A single consensus sequence for the whole genome. It was built up from a high quality set of representative samples of the species (from different populations).*

*It is the first-line comparison during analysis.*

By Genome Reference Consortium (GRC) (<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>)

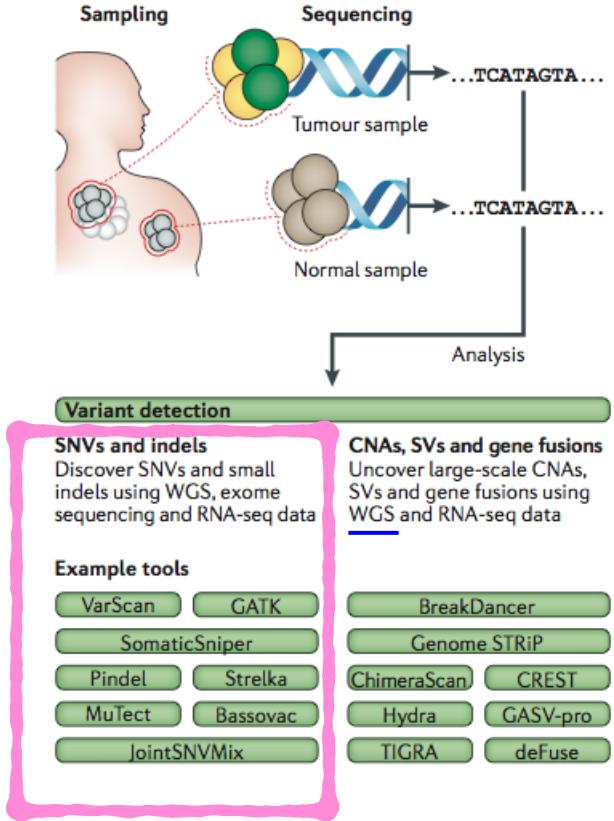
- Human assemblies (Versions):
  - + GRCh37/hg19 : former version. Released in 2012. It is still used for analysis.
  - + CRCh38/hg38 : current version (Sep. 2017). Released in 2014. More accurate, comprehensive (includes Haplotypes) and sophisticated.

We must **keep consistency in the Genome Reference Version** through the variant analysis.

- We must know what **regions along the genome were sequenced** in the experiment, that is, the sequencing library.

# Algorithms for Variant Calling

SNVs and  
Indels



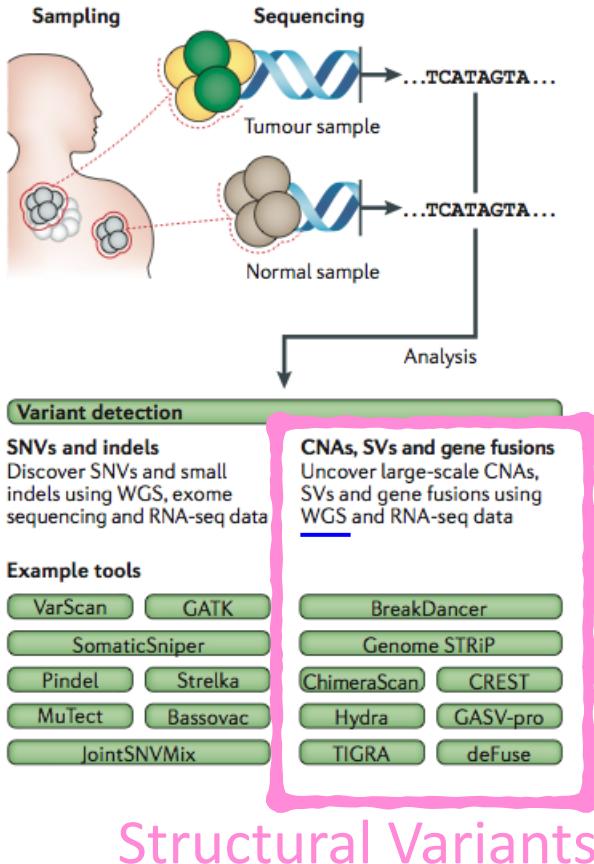
Several Methods have been published.

## Further reading:

Ding, L., Wendl, M. C., McMichael, J. F. & Raphael, B. J. Nat Rev Genet – (2014). doi:10.1038/nrg3767

Nam J.N. et al. Evaluation of somatic copy number estimation tools for whole-exome sequencing data. Brief. Bioinformatics (2015)

# Algorithms for Variant Calling



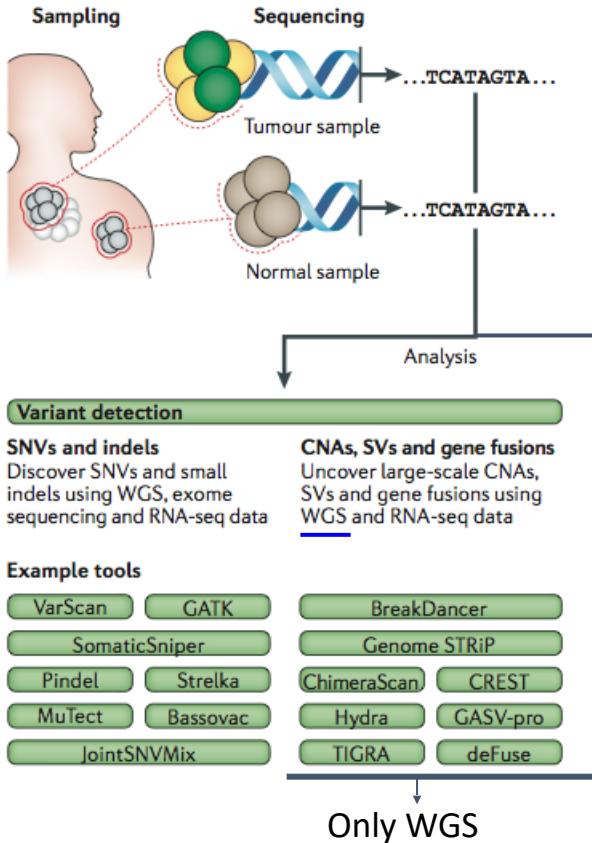
Several Methods have been published.

## Further reading:

Ding, L., Wendl, M. C., McMichael, J. F. & Raphael, B. J. Nat Rev Genet – (2014). doi:10.1038/nrg3767

Nam J.N. et al. Evaluation of somatic copy number estimation tools for whole-exome sequencing data. Brief. Bioinformatics (2015)

# Algorithms for Variant Calling



Several Methods have been published.

Tool	Year	Language	Paired or pooled data	Segmentation	Feature
ADTEX	2014	Python, R	Both	HMM	Noise reduction Ploidy estimation
CONTRA	2012	Python, R	Both	CBS	GC correction
Control-FREEC	2011	C++, R	Paired	LASSO	GC correction, mappability
EXCAVATOR	2013	Perl, R	Both	HSLM	GC correction, mappability, exon-size correction
ExomeCNV	2011	R	Paired	CBS	GC correction, mappability
Varscan2	2012	Java, Perl, R	Paired	CBS	GC correction

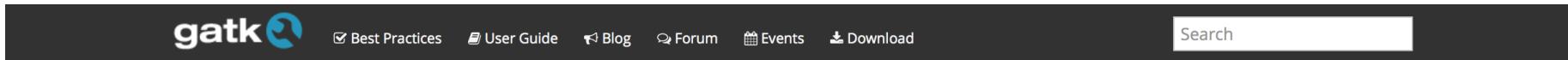
Appropriate methods for Whole-Exome seq

## Further reading:

Ding, L., Wendl, M. C., McMichael, J. F. & Raphael, B. J. Nat Rev Genet – (2014). doi:10.1038/nrg3767

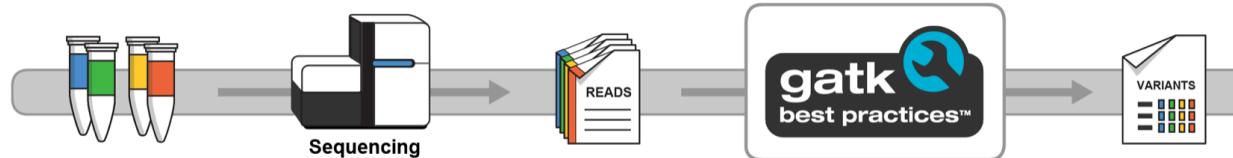
Nam J.N. et al. Evaluation of somatic copy number estimation tools for whole-exome sequencing data. Brief. Bioinformatics (2015)

# GATK for variant calling analysis



## Genome Analysis Toolkit

Variant Discovery in High-Throughput Sequencing Data



Developed in the Data Sciences Platform at the [Broad Institute](#), the toolkit offers a wide variety of tools with a primary focus on variant discovery and genotyping. Its powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

[Learn More](#)



### Best Practices

Pipelines optimized for accuracy and performance



### Blog

Announcements and progress updates



### User Guide

Detailed documentation, tutorials and resources

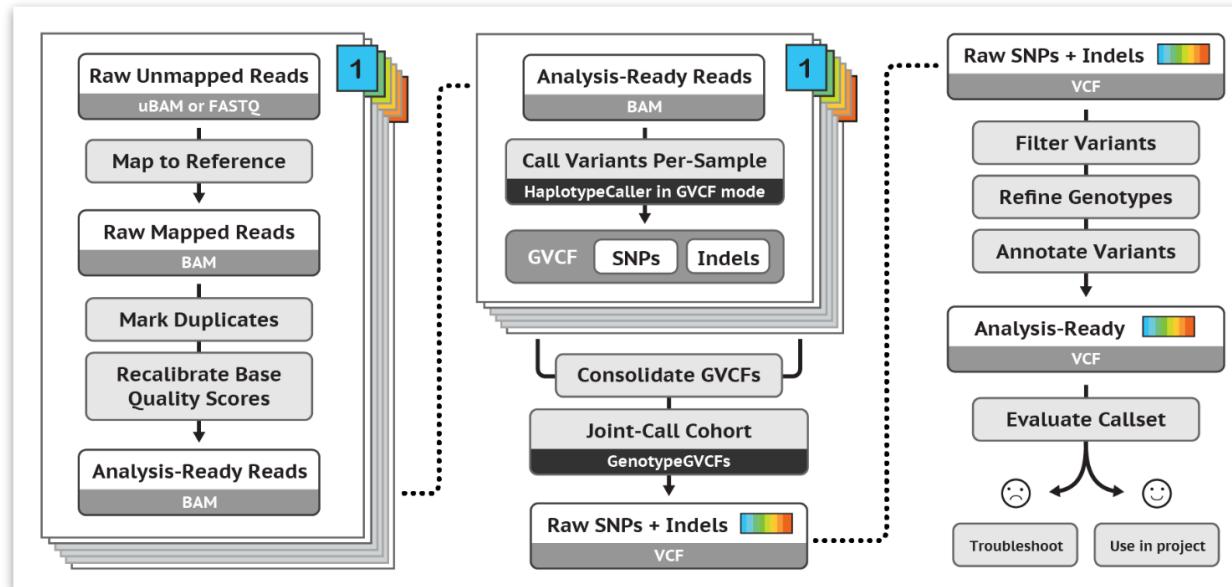


### Forum

Ask our team for help and report issues

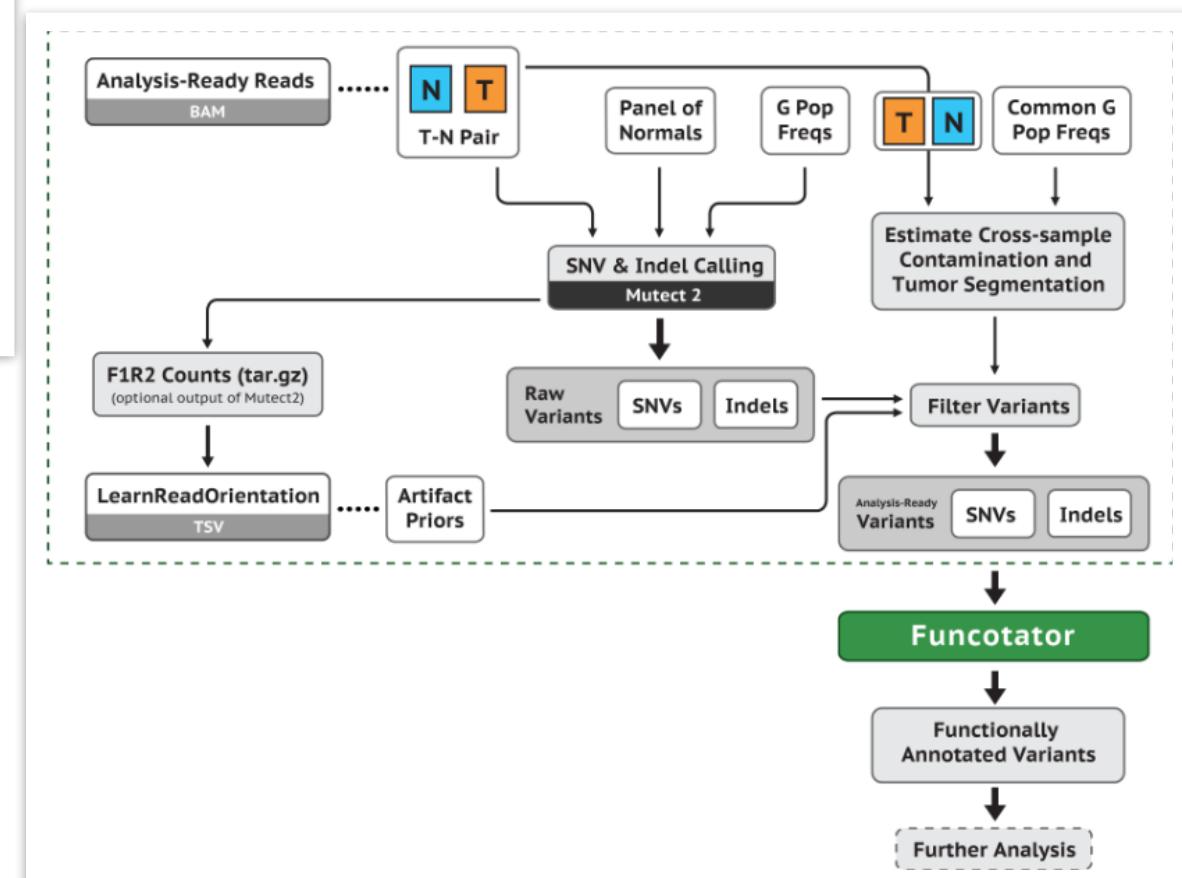


# GATK Best Practices



Somatic small-scale variants

Germline small-scale variants

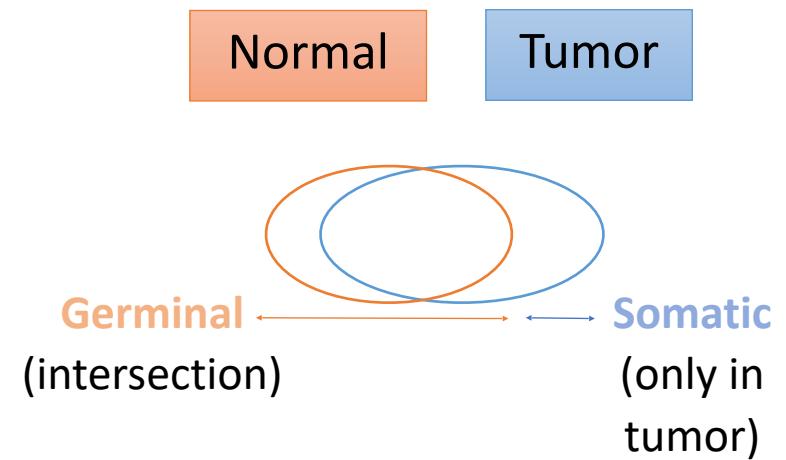


# Somatic vs Germline variants

**Germline:** appear in gametes  
Inheritable  
Affect to future generations  
e.g.: variants involved in rare diseases

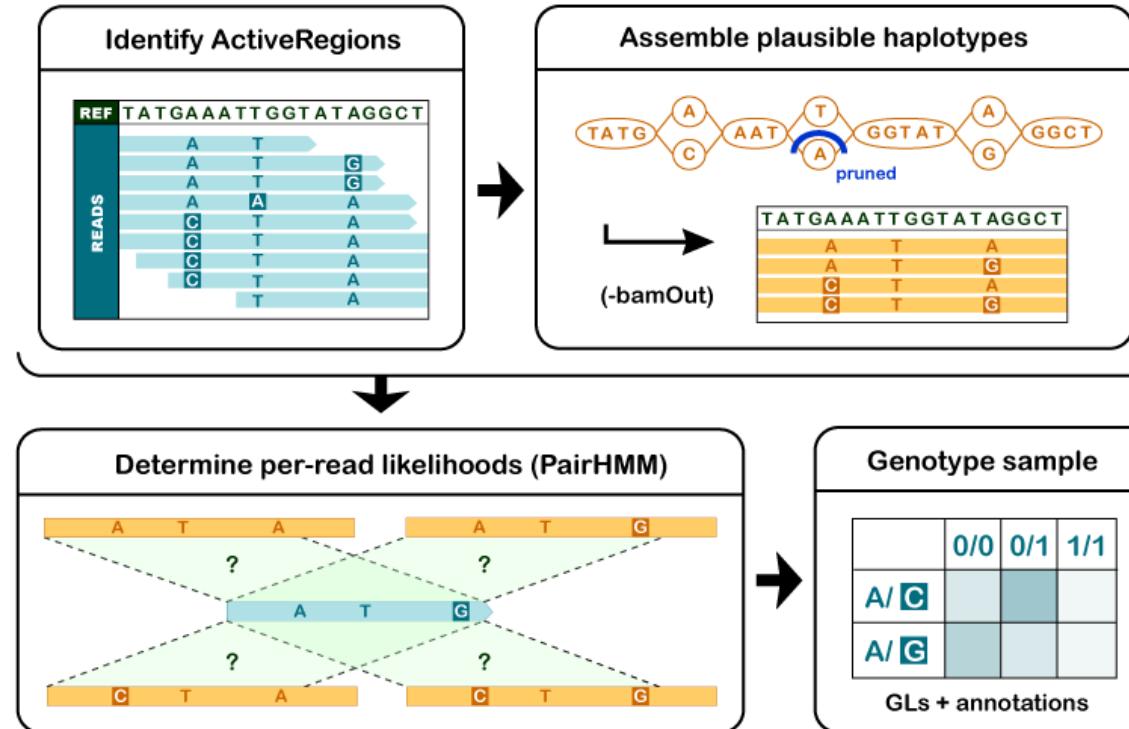
**Somatic:** appear in different from germline cells  
Acquired  
Only affect to the lineage of the affected cell  
e.g.: variants causing cancer

**Identified by comparison:**



# Variant Calling for SNVs and Indels

**Haplotype Caller :** Variant calling based on the calculation of genotype likelihoods:



**Assumptions:** It bases the calling in the indicated ploidy (e.g. 2n)  
**Limited detection of low allele frequencies.**

Further reading:

[https://github.com/broadgsa/gatk/blob/master/doc\\_archive/methods/HC\\_overview:\\_How\\_the\\_HaplotypeCaller\\_works.md](https://github.com/broadgsa/gatk/blob/master/doc_archive/methods/HC_overview:_How_the_HaplotypeCaller_works.md)  
<https://gatk.broadinstitute.org/hc/en-us/sections/360007226771?name=methods>

# VCF file

## HEADER										
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SampleName	
chr17	87234	.	G	A	2000	PASS	DP=80	GT:PL	1/1:3000,220,0	
chr17	98764	.	T	C	340	PASS	DP=30	GT:PL	0/1:1200,0,200	
chr17	108764	.	G	C	10	FILTERED	DP=7	GT:PL	0/1:37,0,200	

Genomic coordinates      Nucleotide change      score  
(higher → better)      filtered?

Allele1 / Allele2 (diploid)  
1/1 → homozygous mutant  
0/1 → heterozygous mutant  
0/0 → homozygous reference

Likelihood for each GT:  
0/0, 0/1, 1/1.  
(lower → better)  
0 is the best score.

More info.:

<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

<https://gatk.broadinstitute.org/hc/en-us/articles/360035531692-VCF-Variant-Call-Format>

# Intra-tumoral heterogeneity

## Variant Allele Frequency

Proportion of DNA molecules in the sample carrying the variant

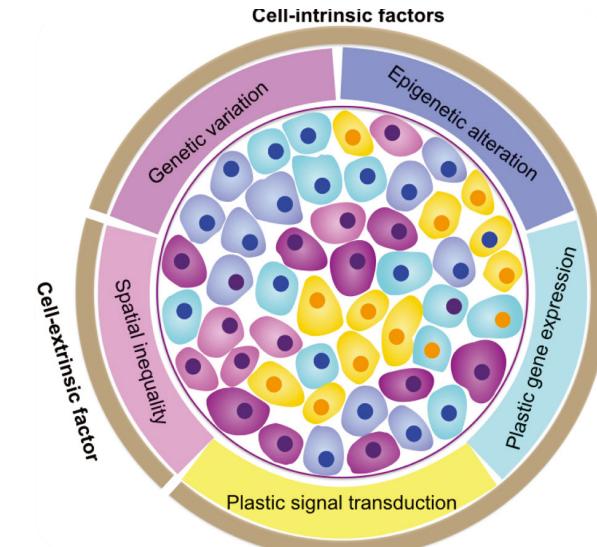
$$VAF = \frac{\text{sequence reads with a DNA variant}}{\text{overall coverage at that locus}}$$

For a diploid organism:

- **heterozygous loci** should be near 0.5 VAF
- **homozygous loci** should be near 1 VAF
- **reference loci** should be near 0 VAF

doi: [10.28092/j.issn.2095-3941.2016.0004](https://doi.org/10.28092/j.issn.2095-3941.2016.0004)

**Clonal composition in cancer**  
changes 0.5/1.0 diploid  
variant allele frequencies



doi: [10.1038/aps.2015.92](https://doi.org/10.1038/aps.2015.92)

# Variant Calling for somatic variants: MuTect2

**SNV and Indel caller.**

Similar logic to Haplotype Caller but:

- It allows variable allele frequencies.
- It includes logic to avoid germline variants.

The screenshot shows two adjacent web pages. The left page is the CGA homepage, featuring a search bar, a login button, and a sidebar with links to various genomic analysis tools like ABSOLUTE, BreakPointer, and MuTect. The right page is the MuTect page, which includes a brief introduction, a section on how it works, and a table of validation rates from cancer studies. The MuTect page also contains mathematical formulas for LOD scores.

**What does MuTect do?**

MuTect is a method developed at the Broad Institute for the reliable and accurate identification of somatic point mutations in next generation sequencing data of cancer genomes.

For complete details, please see our publication in *Nature Biotechnology*:

Cibulskis, K. et al. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat Biotechnology* (2013).doi:10.1038/nbt.2514

**How does it work?**

In brief, muTect consists of three steps.

1. Preprocessing the aligned reads in the tumor and normal sequencing data. In this step we ignore reads with too many mismatches or very low quality scores since these represent noisy reads that introduce more noise than signal.
2. A statistical analysis that identifies sites that are likely to carry somatic mutations with high confidence. The statistical analysis predicts a somatic mutation by using two Bayesian classifiers – the first aims to detect whether the tumor is non-reference at a given site and, for those sites that are found as non-reference, the second classifier makes sure the normal does not carry the variant allele. In practice the classification is performed by calculating a LOD score (log odds) and comparing it to a cutoff determined by the log ratio of prior probabilities of the considered events. For the tumors we calculate

$$LOD_t = \log_{10} \left( \frac{P(\text{observed data in tumor|site is mutated})}{P(\text{observed data in tumor|site is reference})} \right)$$

$LOD_n = \log_{10} \left( \frac{P(\text{observed data in normal|site is reference})}{P(\text{observed data in normal|site is mutated})} \right)$

Since we expect somatic mutations to occur at a rate of  $\sim 1$  in a Mb, we require  $LOD_t > \log_{10}(0.5 \times 10^{-1}) = -6.3$  which guarantees that our false positive rate, due to noise in the tumor, is less than half of the somatic mutation rate. In the normal, not in dbSNP sites, we require  $LOD_n > \log_{10}(0.5 \times 10^{-1}) = -2.3$  since non-dbSNP germline variants occur roughly at a rate of 100 in a Mb. This cutoff guarantees that the false positive somatic call rate, due to missing the variant in the normal, is also less than half the somatic mutation rate.

3. Post-processing of candidate somatic mutations to eliminate artifacts of next-generation sequencing, short read alignment and hybrid capture. For example, sequence context can cause hallucinated alternate alleles but often only in a single direction. Therefore, we test that the alternate alleles supporting the mutations are observed in both directions.

As muTect attempts to call mutations it also generates a coverage file (in a wiggle file format, which indicates for every base whether it is sufficiently covered in the tumor and normal to be sensitive enough to call mutations). We currently use cutoffs of at least 14 reads in the tumor and at least 8 in the normal (these cutoffs are applied after removing noisy reads in the preprocessing step). In addition, wiggle files can also be generated of the observed depth in the tumor and in the normal.

Most cancer genome studies at the Broad Institute have made use of MuTect and have validated the mutation calls as a part of their cancer biology papers, showing that MuTect has a very low false positive rate. A summary of validation rates from these papers are show below:

publication	technology	candidates	validated	no result	validation rate
Multiple Myeloma <sup>1</sup>	Sequenom	97	92	5	94.85%
Ovarian <sup>2</sup>	Sequenom/PCR/454	1655	1483	172	89.61%
Ovarian <sup>2</sup>	Capture/Illumina	6497	6232	265	95.92%
Head and Neck <sup>3</sup>	Sequenom	321	288	33	89.72%
Breast <sup>4</sup>	Sequenom/PCR/454	455	428	0	94.07%

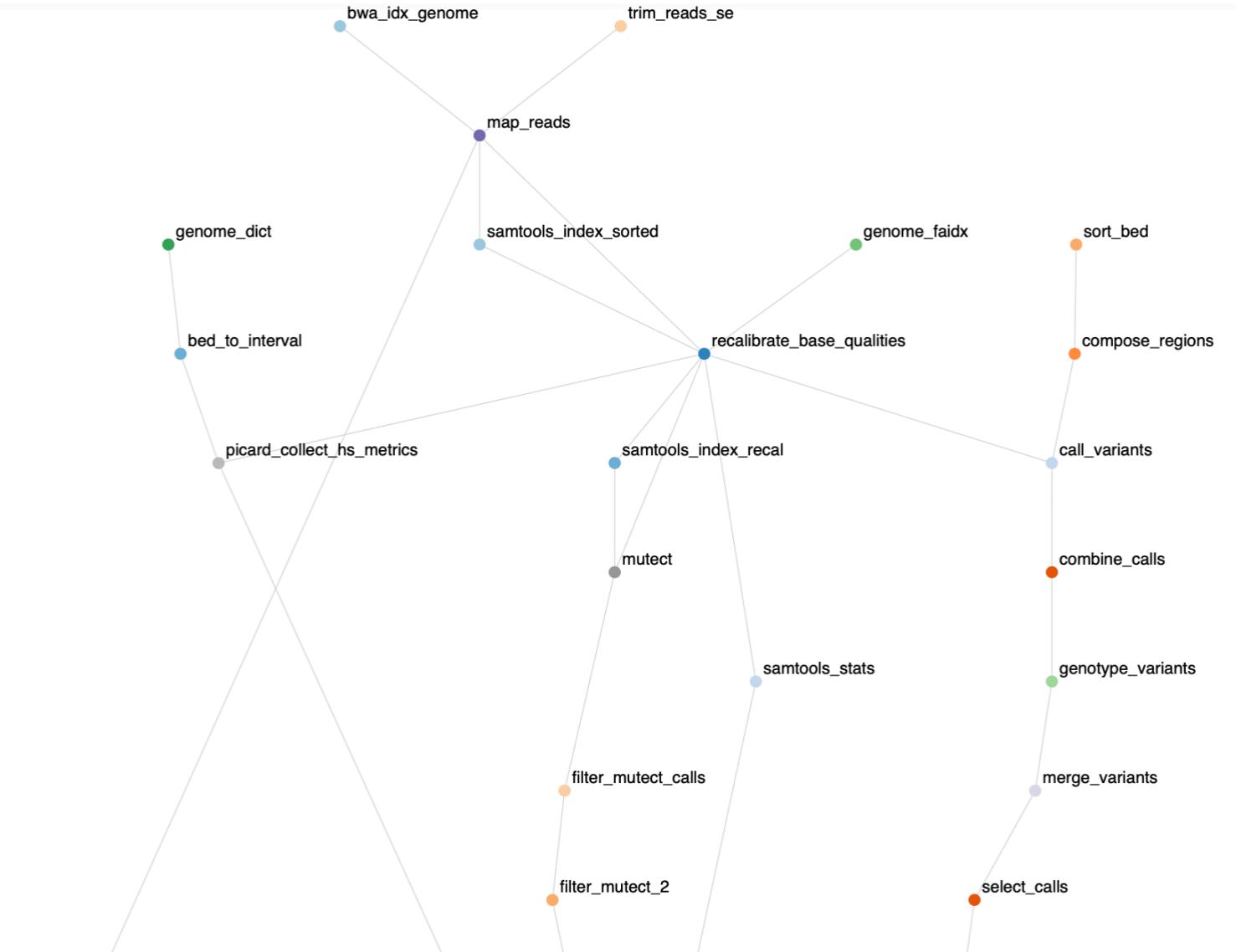
<https://gatk.broadinstitute.org/hc/en-us/articles/360037593851-Mutect2>

Cibulskis, K. et al.  
Nat Biotechnology (2013).doi:10.1038/nbt.2514

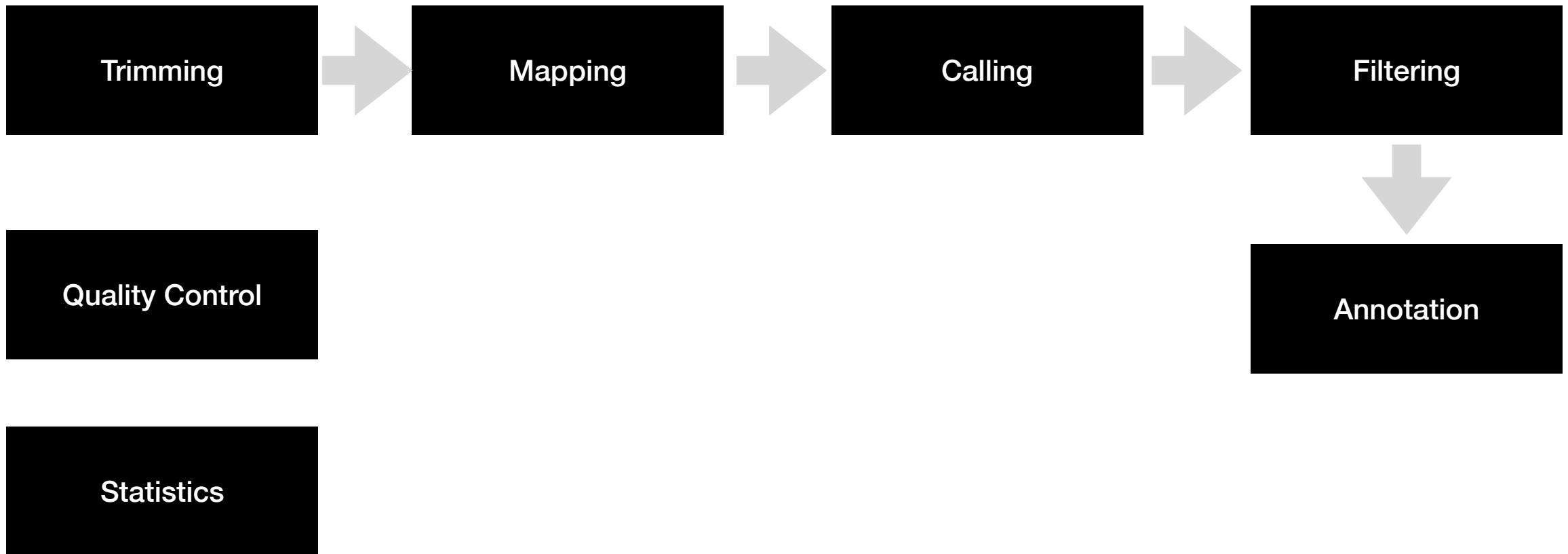
# varca



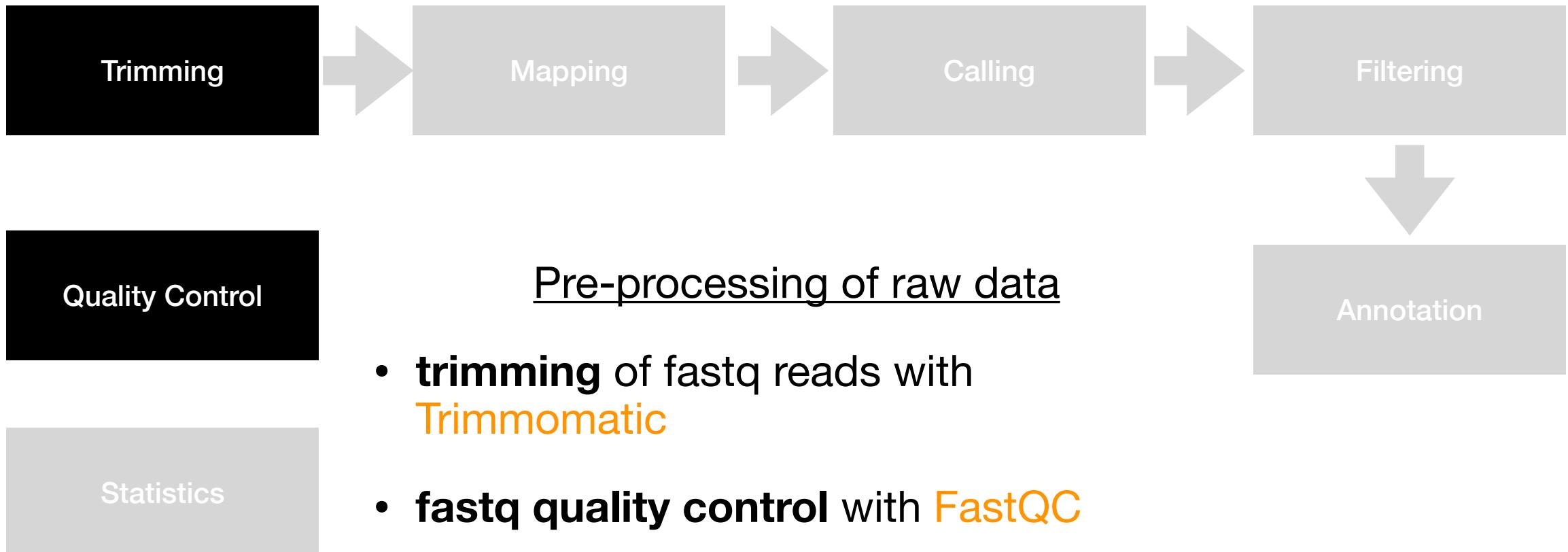
[https://gitlab.com/bu\\_cnio/varca](https://gitlab.com/bu_cnio/varca)



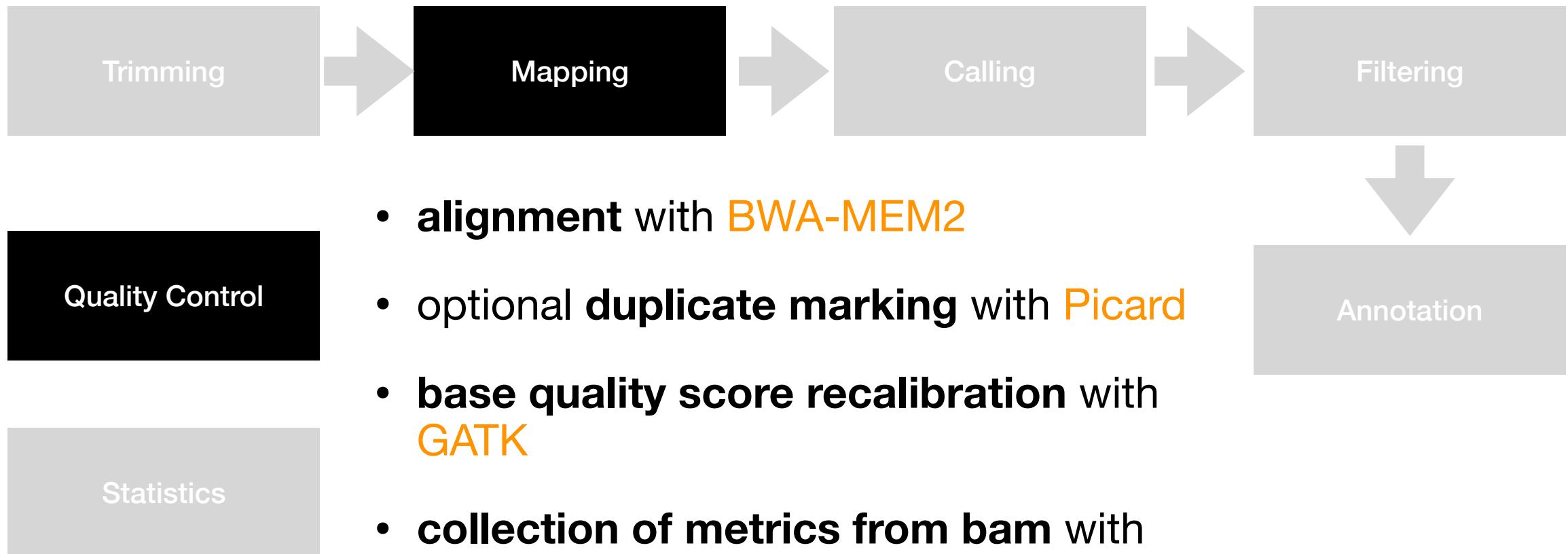
# Steps in the pipeline



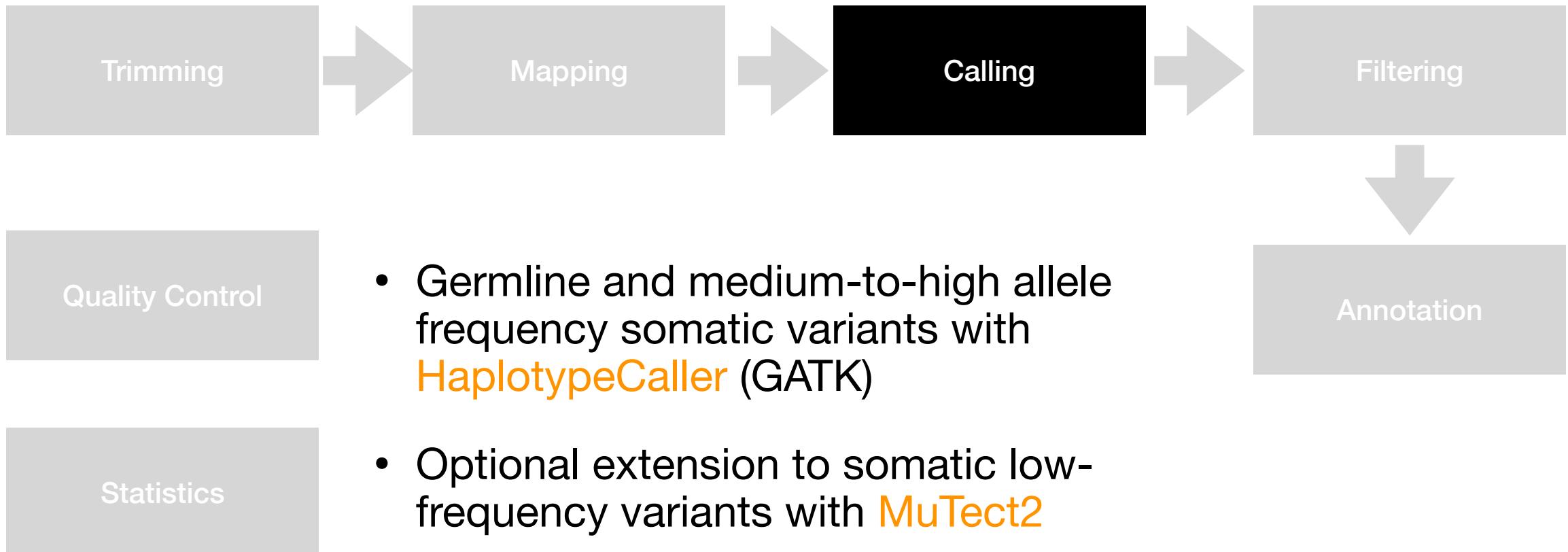
# Steps in the pipeline



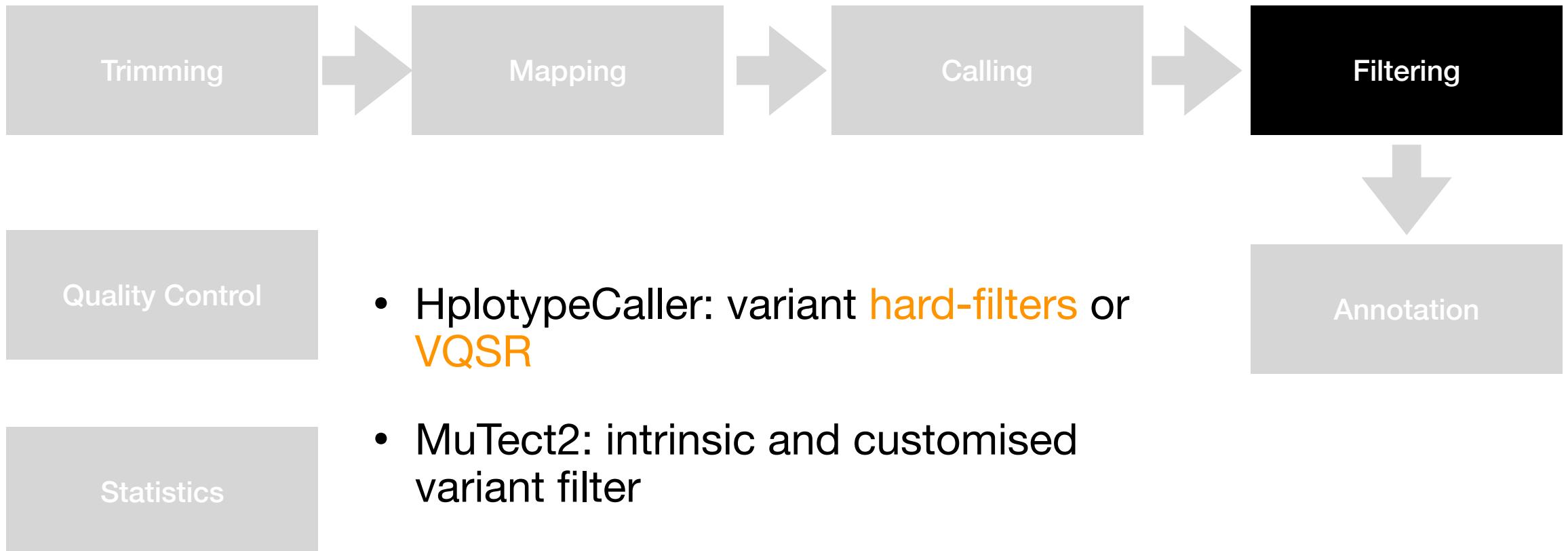
# Steps in the pipeline



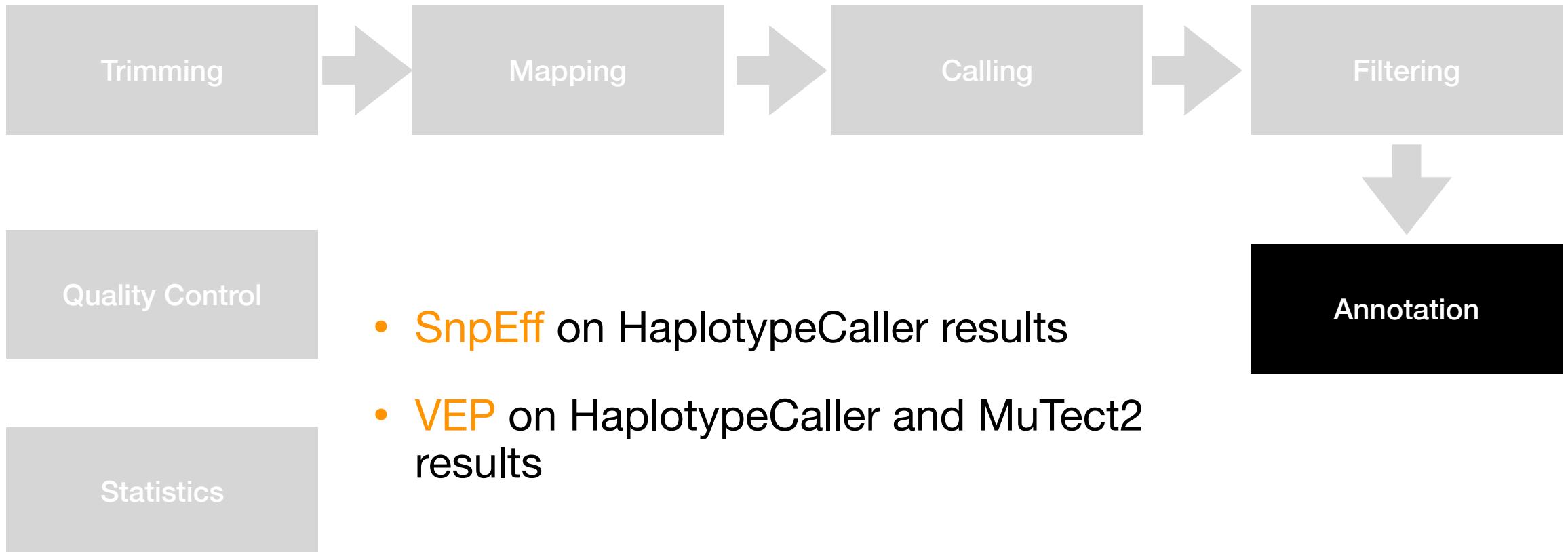
# Steps in the pipeline



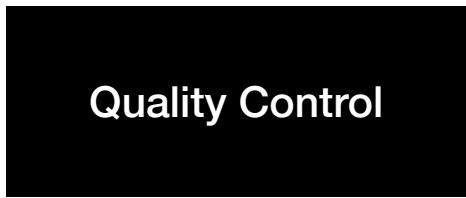
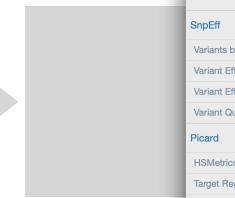
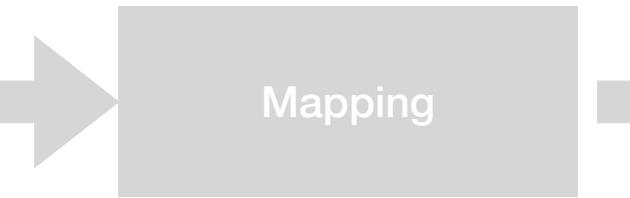
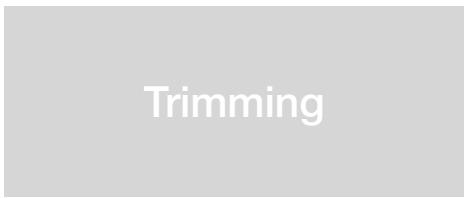
# Steps in the pipeline



# Steps in the pipeline



# Steps in the pipeline



- Quality report with **MultiQC**
- Statistic Plots: allele frequencies and depth of sequencing

The screenshot shows the MultiQC interface version 1.7. The left sidebar lists modules: General Stats, SnpEff, Picard, HSMetrics, Target Region Coverage, Mark Duplicates, Samtools, Percent Mapped, and Alignment metrics. The main area displays 'General Statistics' for samples 'all', 'normal-1', and 'tumor-1'. The 'SnpEff' section is also visible.

Sample Name	Change rate	Ts/Tv	M Variants	Fold Enrichment	Target Bases 30X	% Dups	Error rate	M Non-Primary
all	21 724	2.245	0.14			3.5%	0.37%	0.0
normal-1			7	24%				
tumor-1			5	6%		4.3%	0.73%	0.0

# Steps in the pipeline: parallelization

