

PO21: Precision Oncology Course

Variant Detection

Outline

- Definition and types of genomic variants
- Genomic variants in cancer
- Steps for variant detection after alignment
- Algorithms for variant calling
- Pipeline for targeted DNA sequencing: varca

Genomic Variants

Definition, relevance and types

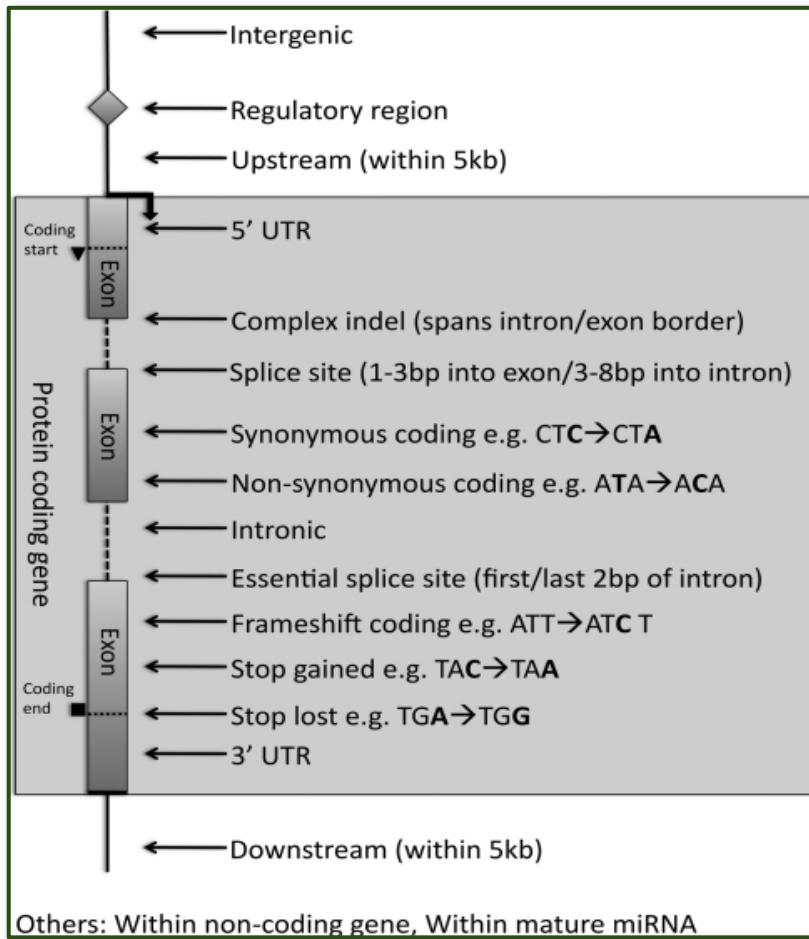
What are genomic variants?

- Genomic variants are **permanent** changes in the DNA sequence of an organism.
- They can **emerge by different mechanisms**:
 - Recombination during gametes formation.
 - Errors during the DNA replication.
 - External factors like radiation, viruses, transposons, tobacco, UV light.

Types of genomic variants

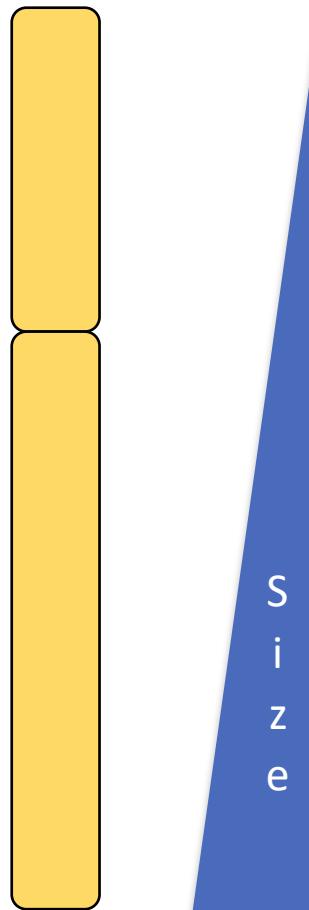
- Different types of variants according to different criteria:
 - Variant **size**
 - Variant **position** in the DNA sequence
 - **Consequence** of the variant in transcription and translation
 - **Clinical implication**

Classification of variants according to the sequence position



- Intergenic
- In regulatory regions
- Upstream
- Downstream
- In genes
 - Untranslated regions: 5'UTR y 3'UTR
 - Exons
 - Introns
 - Splicing regions

Classification of variants according to size



SNVs Single Nucleotide Variants (SNP)
Indels Small Insertion and Deletions

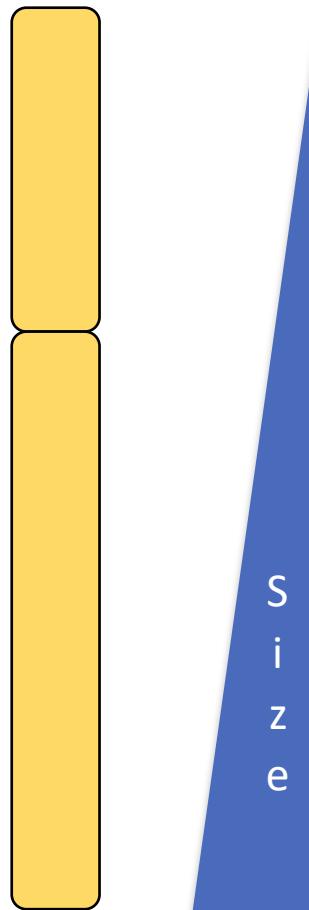
VNTRs (Micro, Minisatellites)
Variable Number of Tandem Repetitions

CNVs Copy Number Variations
Translocations, Inversions

Aneuploidies

Structural variants

Classification of variants according to size



SNVs Single Nucleotide Variants (SNP)
Indels Small Insertion and Deletions

Small scale variants

VNTRs (Micro, Minisatellites)
Variable Number of Tandem Repetitions

CNVs Copy Number Variations
Translocations, Inversions

Aneuploidies

Genomic Variants in Cancer

Intra-tumoral heterogeneity

Variant Allele Frequency

Proportion of DNA molecules in the sample carrying the variant

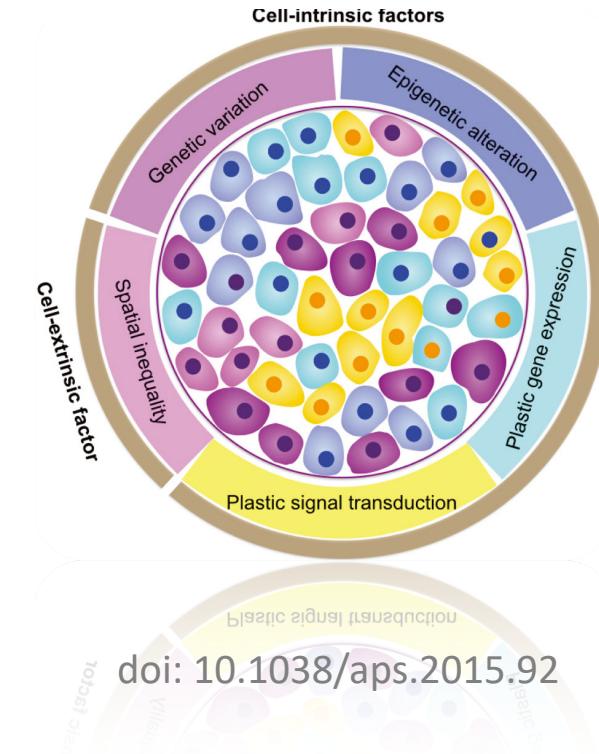
$$VAF = \frac{\text{sequence reads with a DNA variant}}{\text{overall coverage at that locus}}$$

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

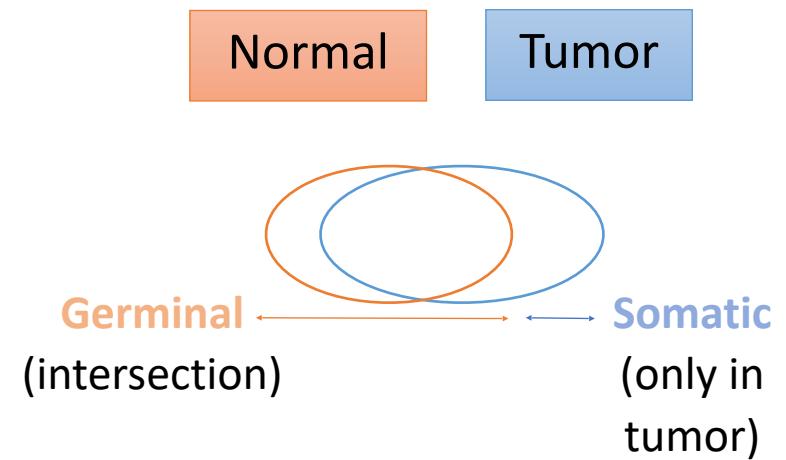


Somatic vs Germline variants

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

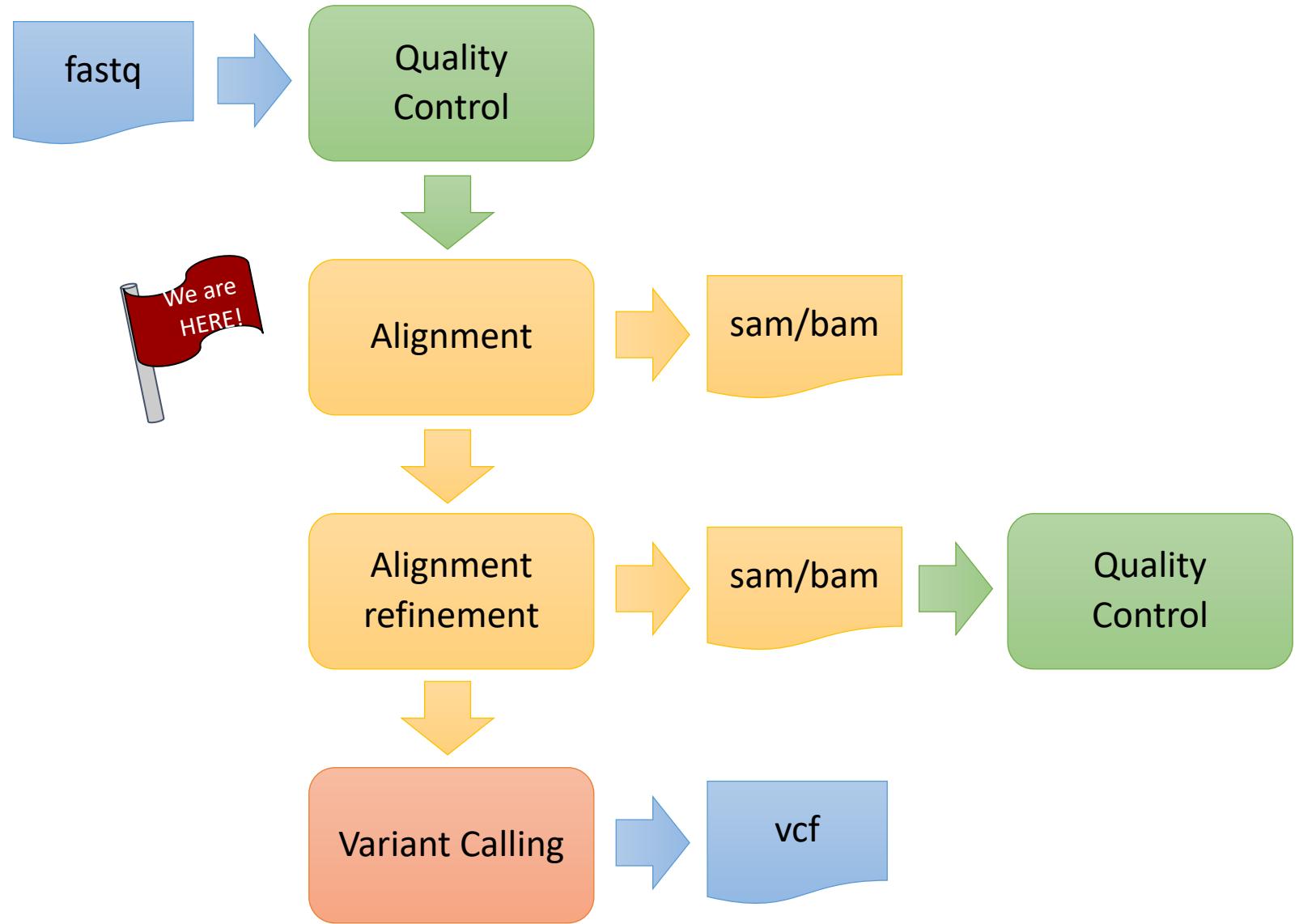
Identified by comparison:



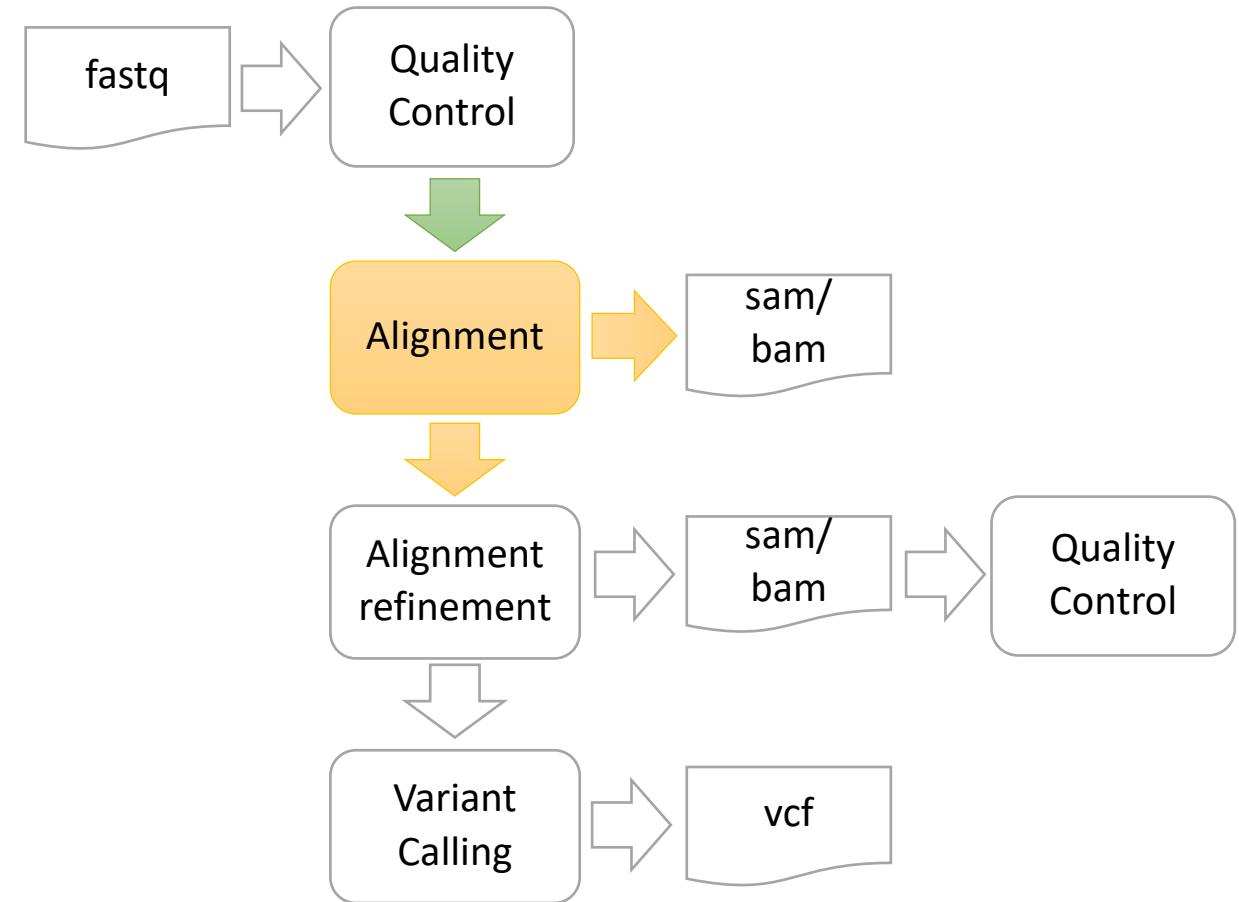
Variant detection

Sequencing and bioinformatics process

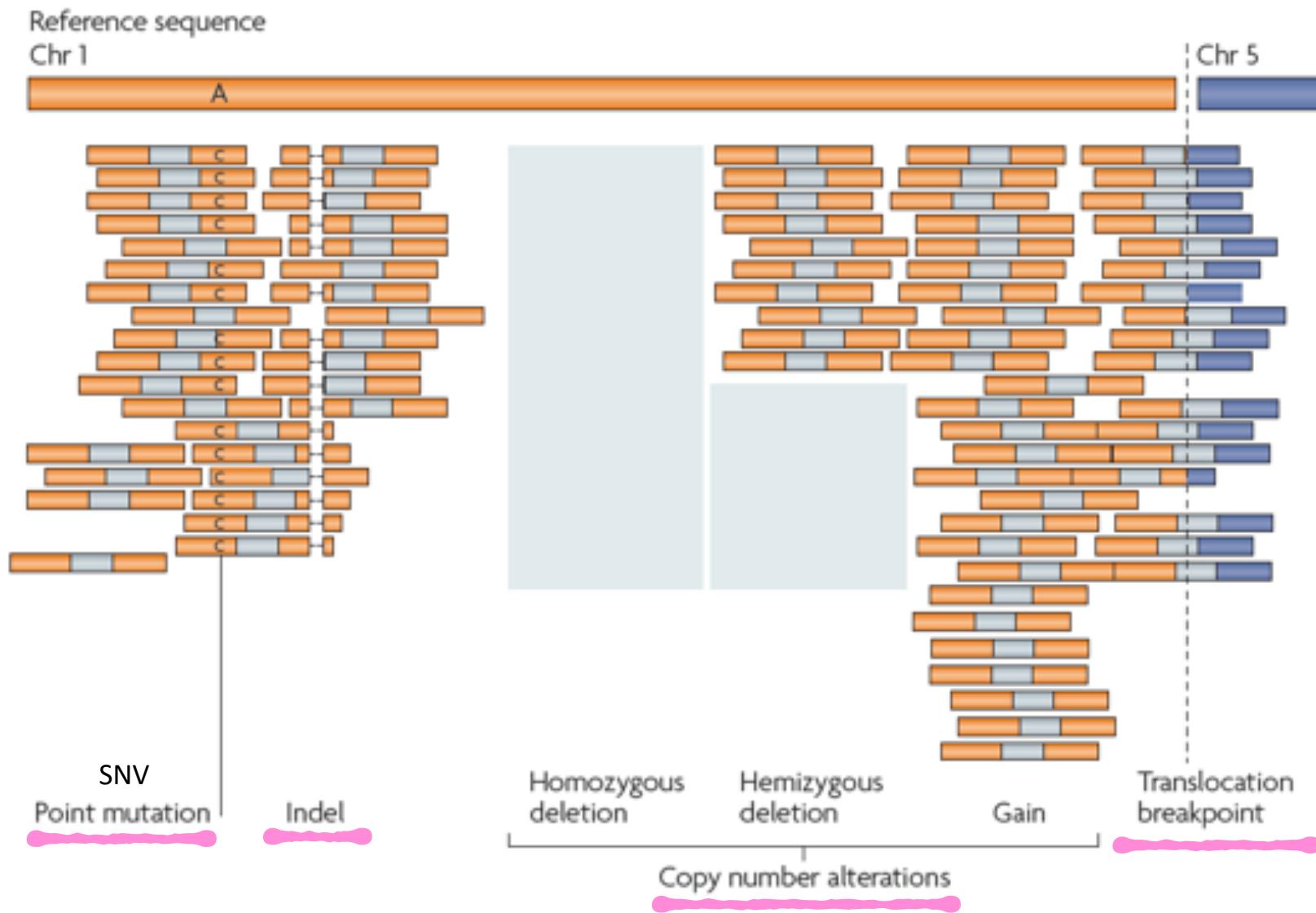
General steps



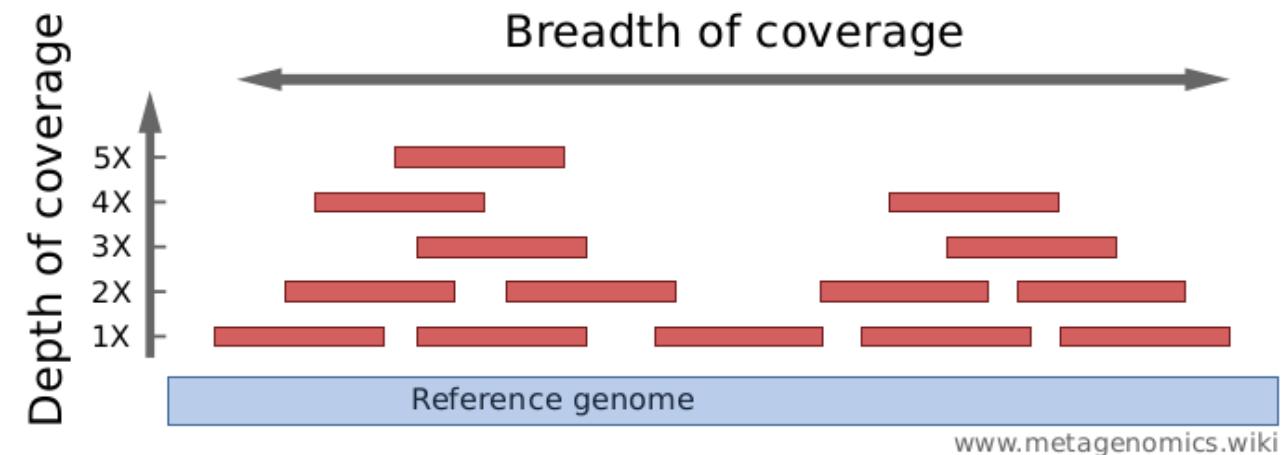
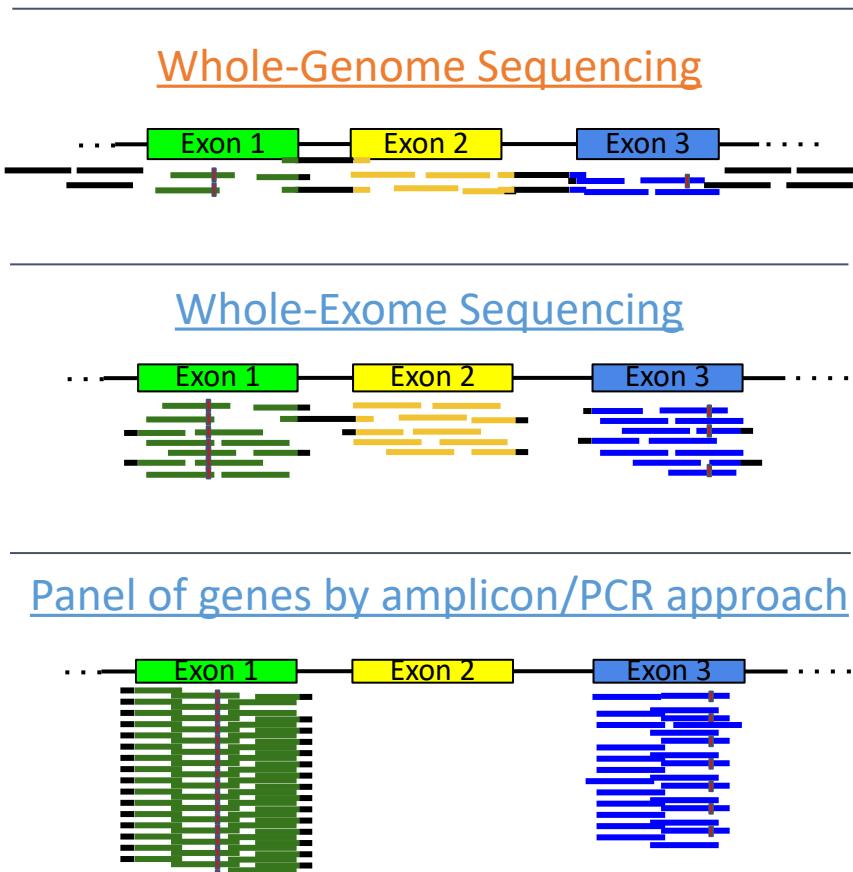
Alignment



Alignment of reads uncover potential variant sites



DNA-seq strategies – Sequencing coverage



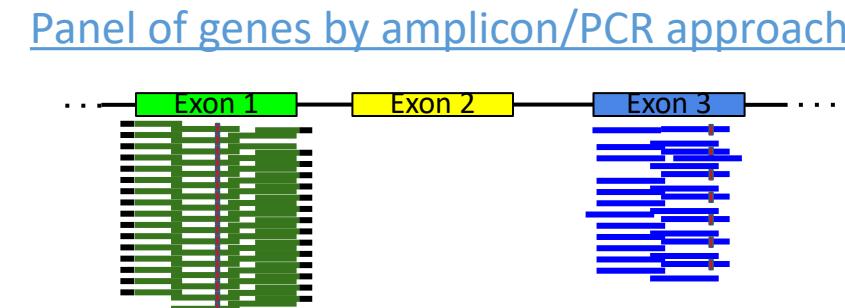
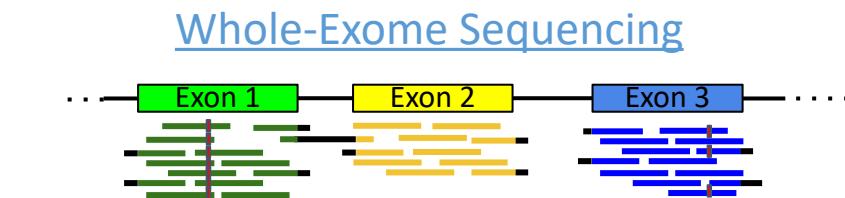
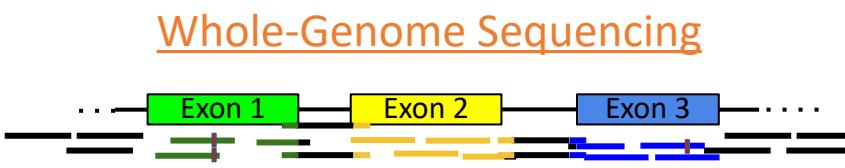
Depth: times that a base is sequenced

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)

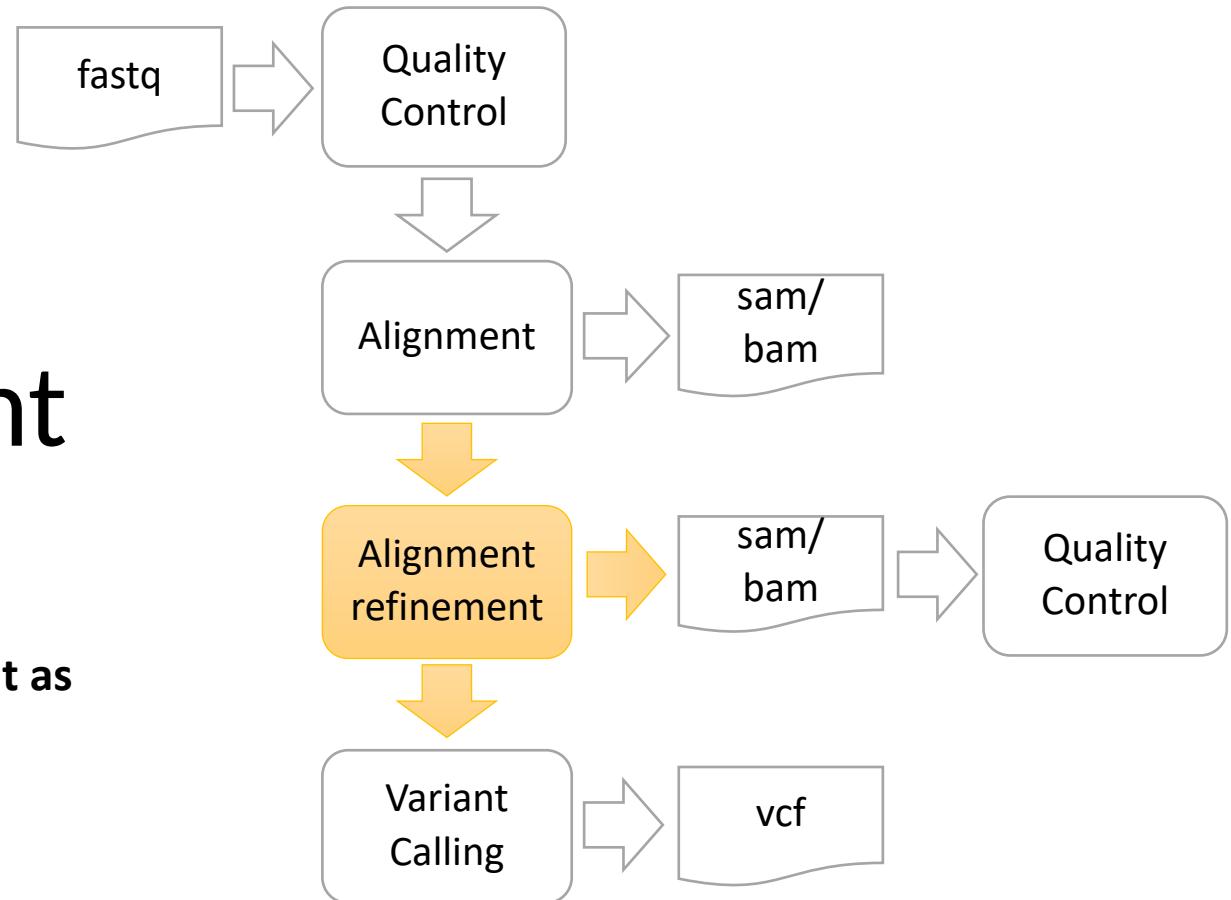
DNA-seq strategies – Type of variants



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.



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)



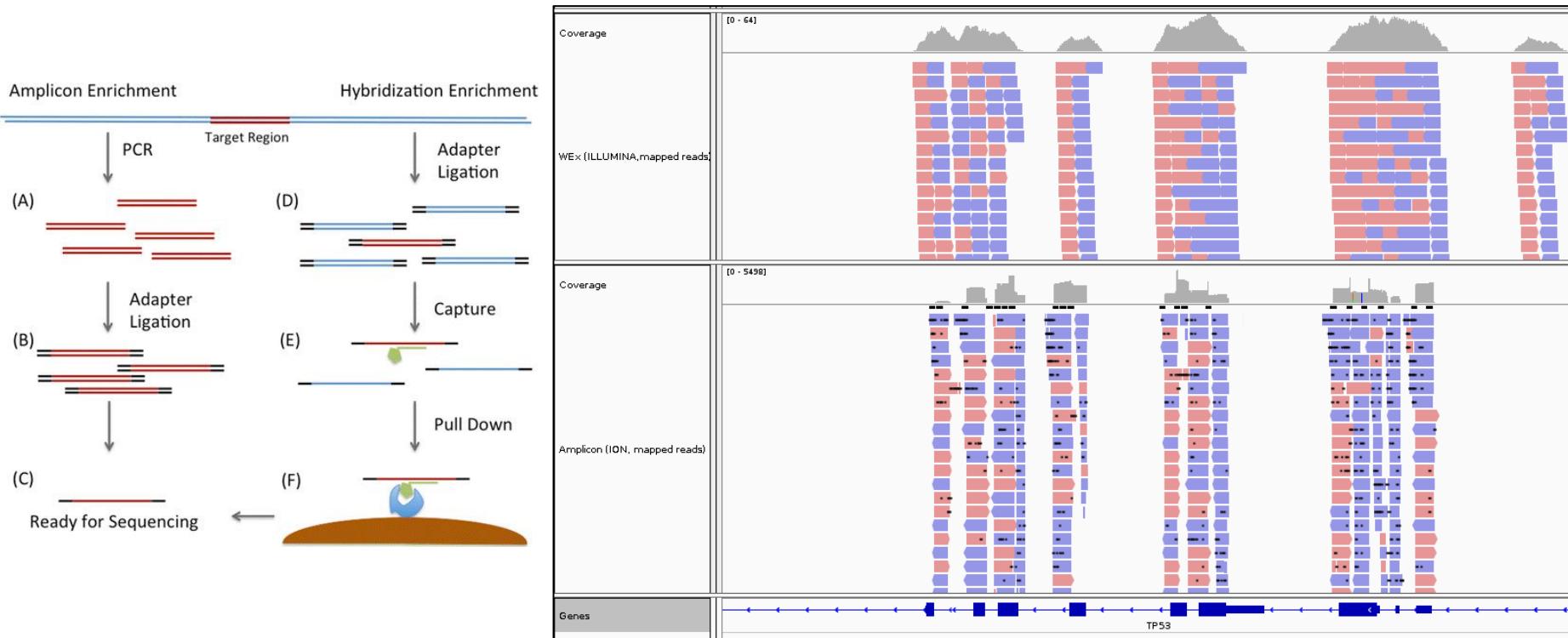
After marking/removing duplicates, the variant caller will only see :



... 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.

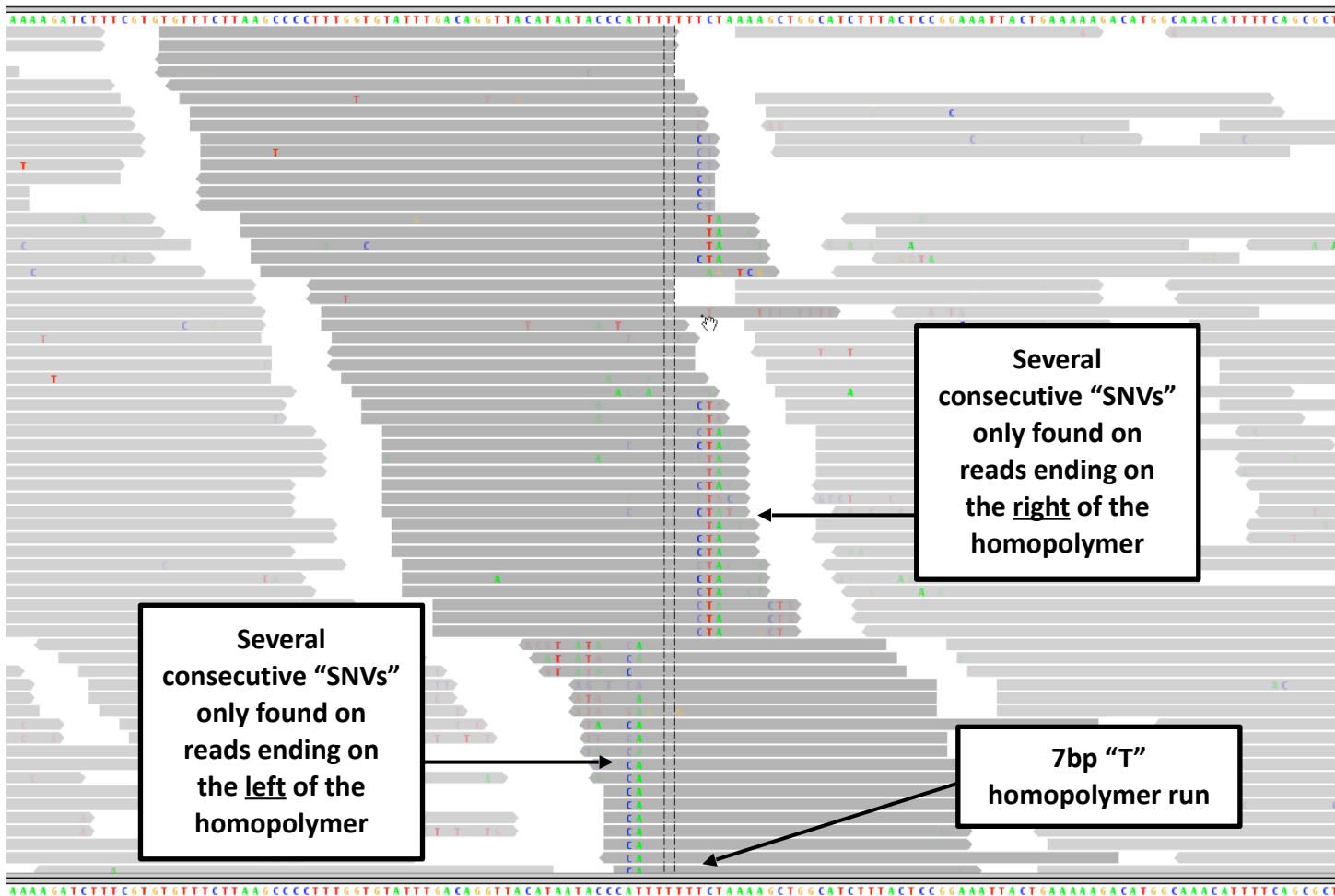
*During alignment, **penalties on mismatches are much cheaper than gaps (indels)**.*

- 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

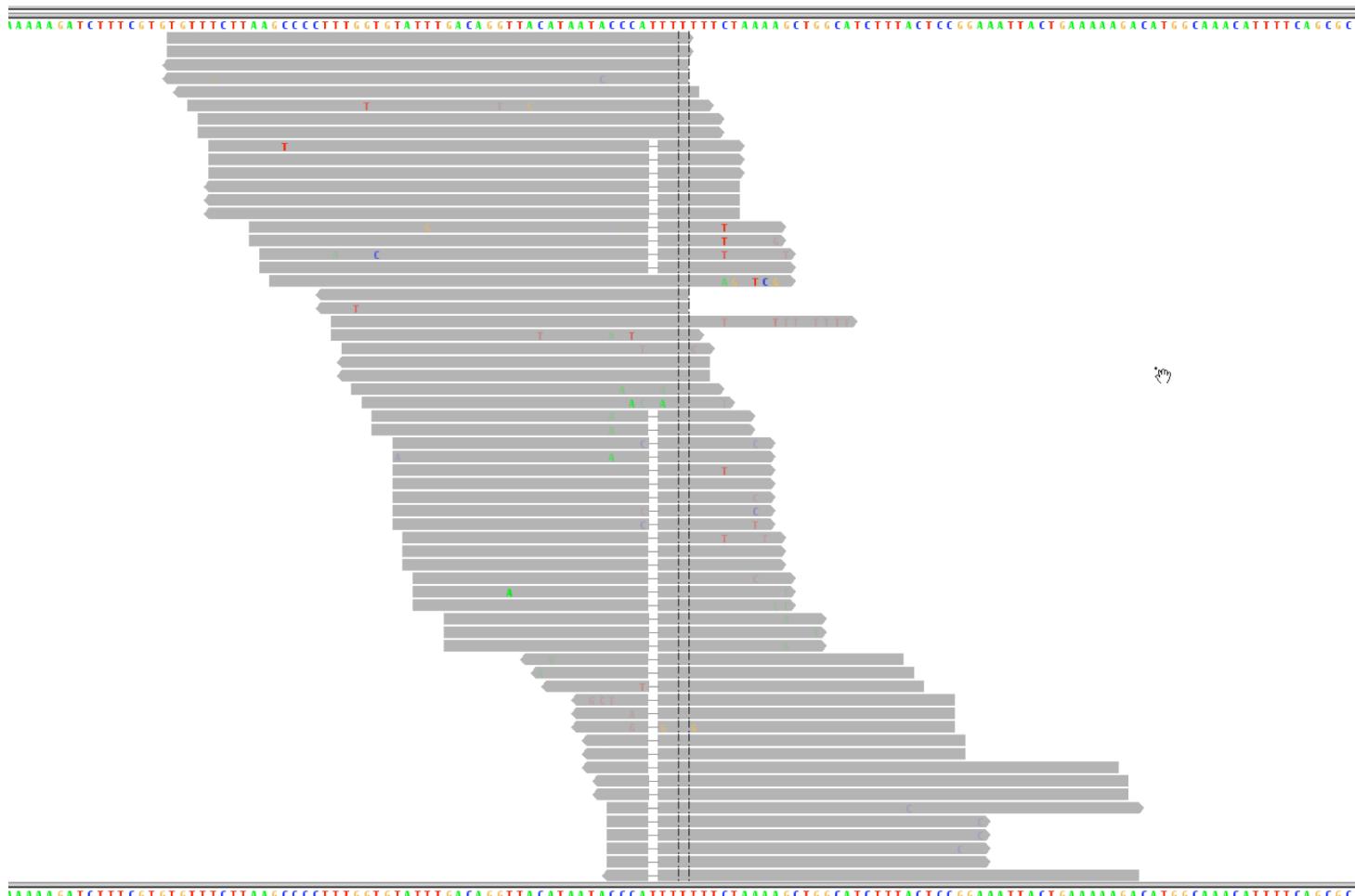
[https://github.com/broadinstitute/gatk-docs/blob/master/gatk3-tutorials/\(howto\)_Perform_local_realignment_around_indels.md](https://github.com/broadinstitute/gatk-docs/blob/master/gatk3-tutorials/(howto)_Perform_local_realignment_around_indels.md)

Indel realignment



Taken from GATK team

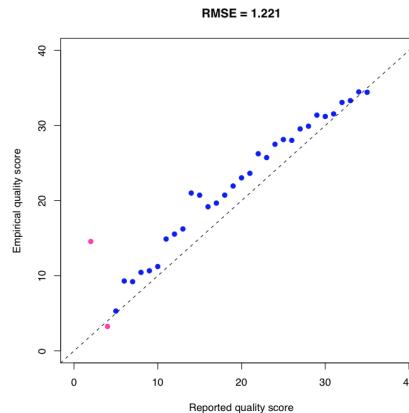
Indel realignment



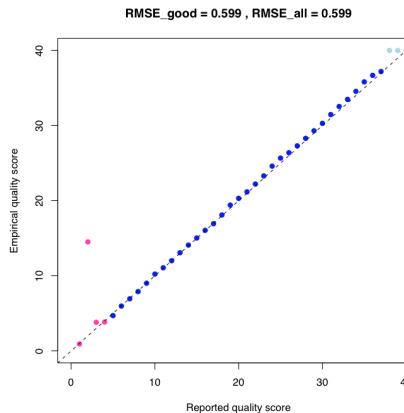
Taken from GATK team

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

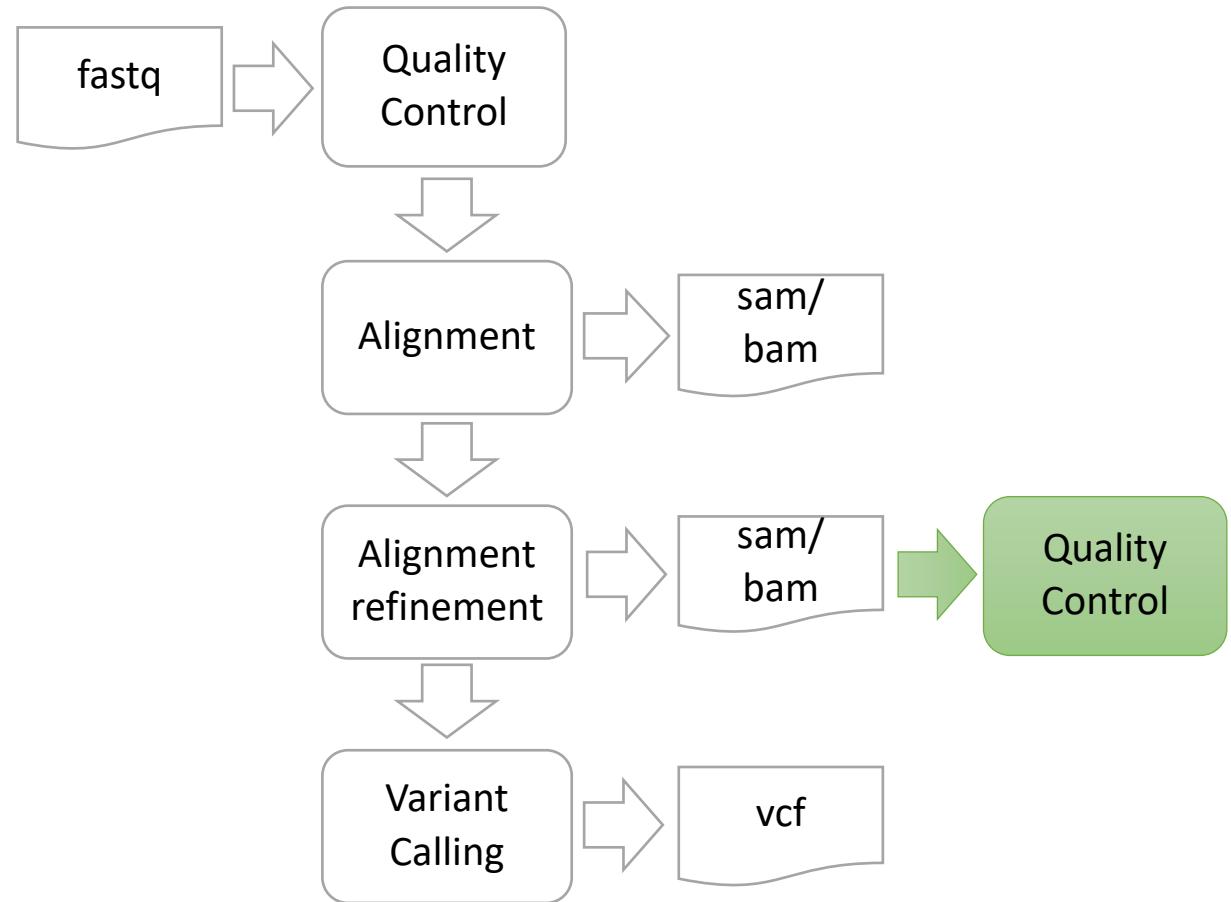


After BQSR recalibration

METHOD: by GATK

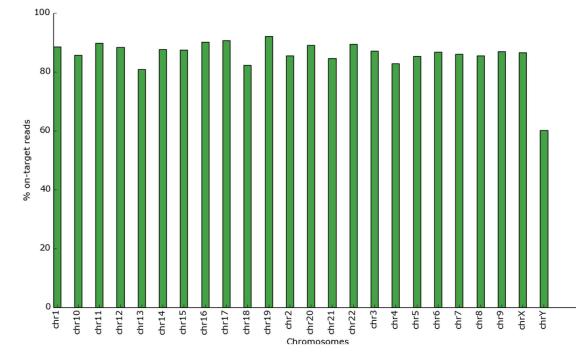
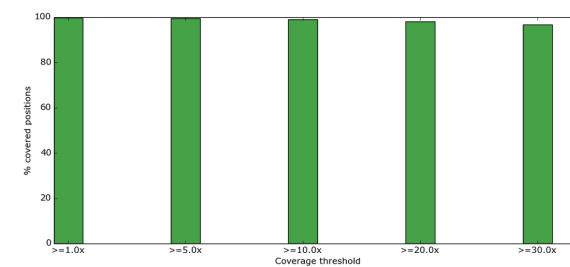
<https://gatk.broadinstitute.org/hc/en-us/articles/360035890531-Base-Quality-Score-Recalibration-BQSR->

Alignment Quality Control

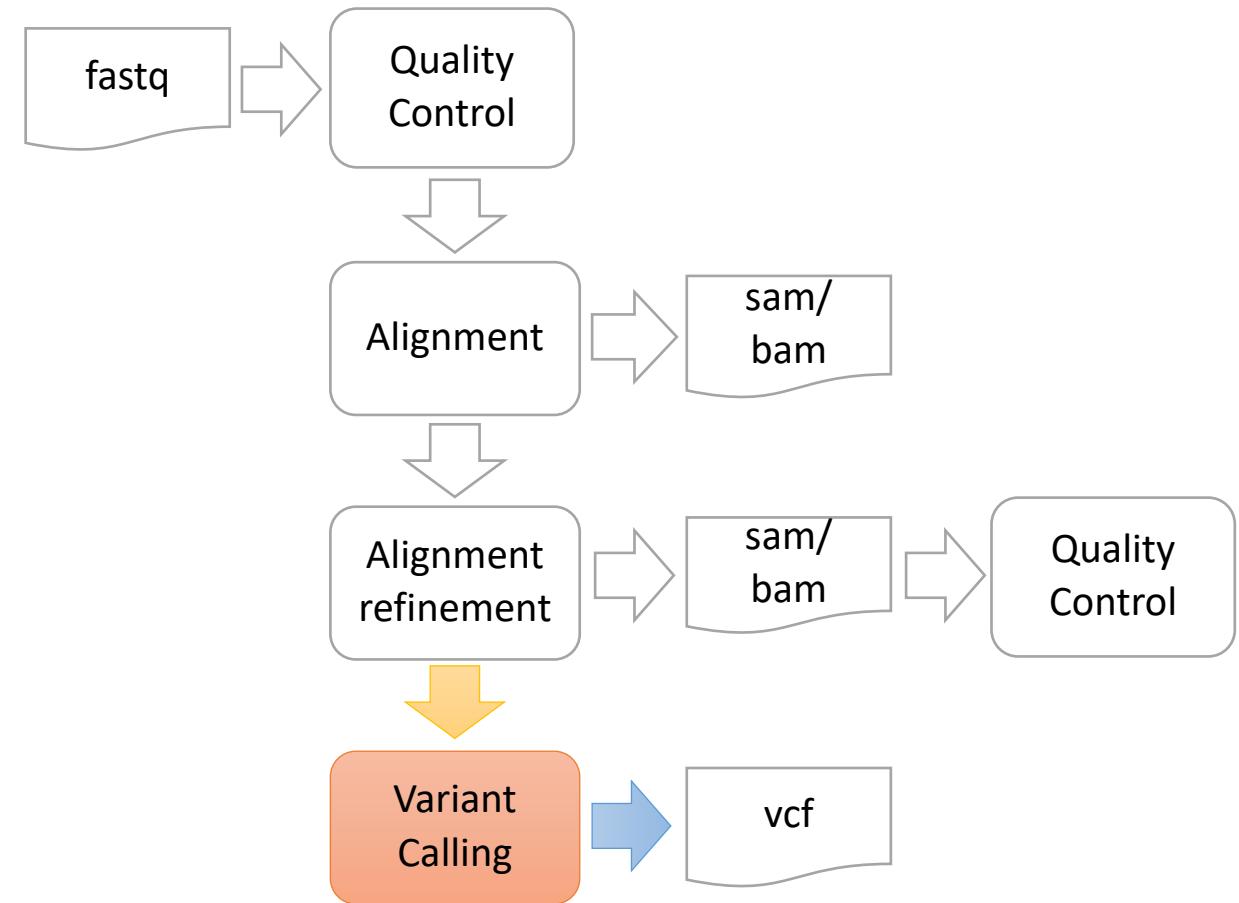


Alignment Quality Control

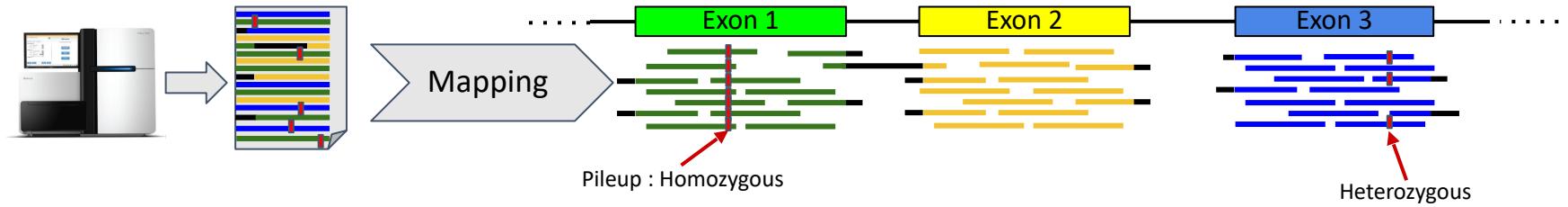
- Mean sequencing depth
 - Is there enough coverage in regions of interest?
 - Are the reads on-target?
-
- Software:
 - ngsCAT
 - QualiMap



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/>)

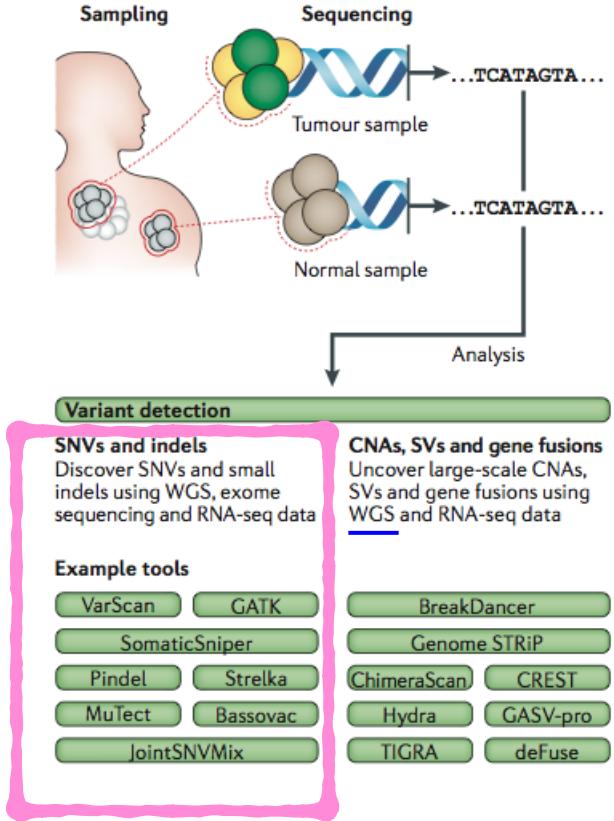
- Human assemblies (Versions):
 - + GRCh37/hg19 : former version. Released in 2012. It is still used for analysis.
 - + CRCh38/hg38 : current version (Sep. 2017). Released in 2014. More accurate, comprehensive (includes Haplotypes) and sophisticated.

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



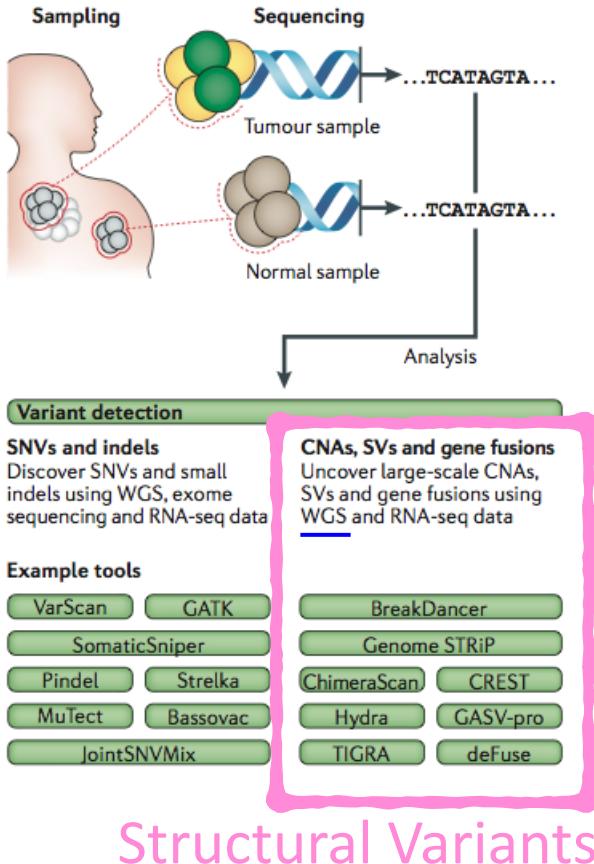
Several Methods have been published.

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



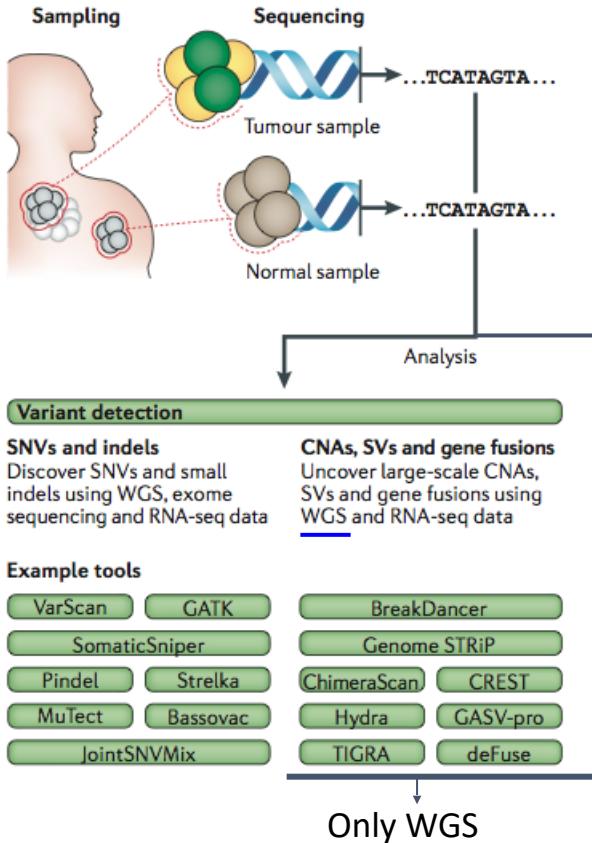
Several Methods have been published.

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

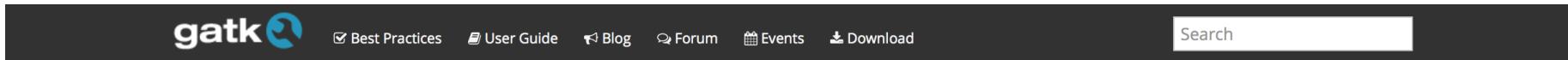
Appropriate methods for Whole-Exome seq

Further reading:

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

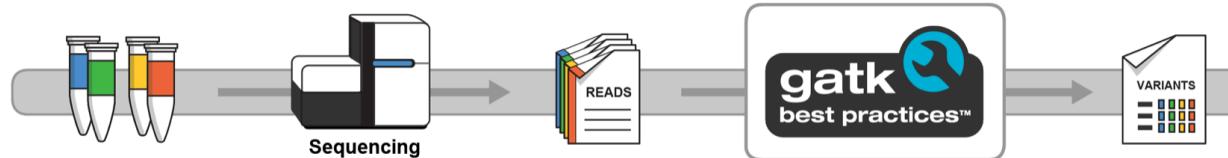
Nam J.N. et al. Evaluation of somatic copy number estimation tools for whole-exome sequencing data. Brief. Bioinformatics (2015)

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

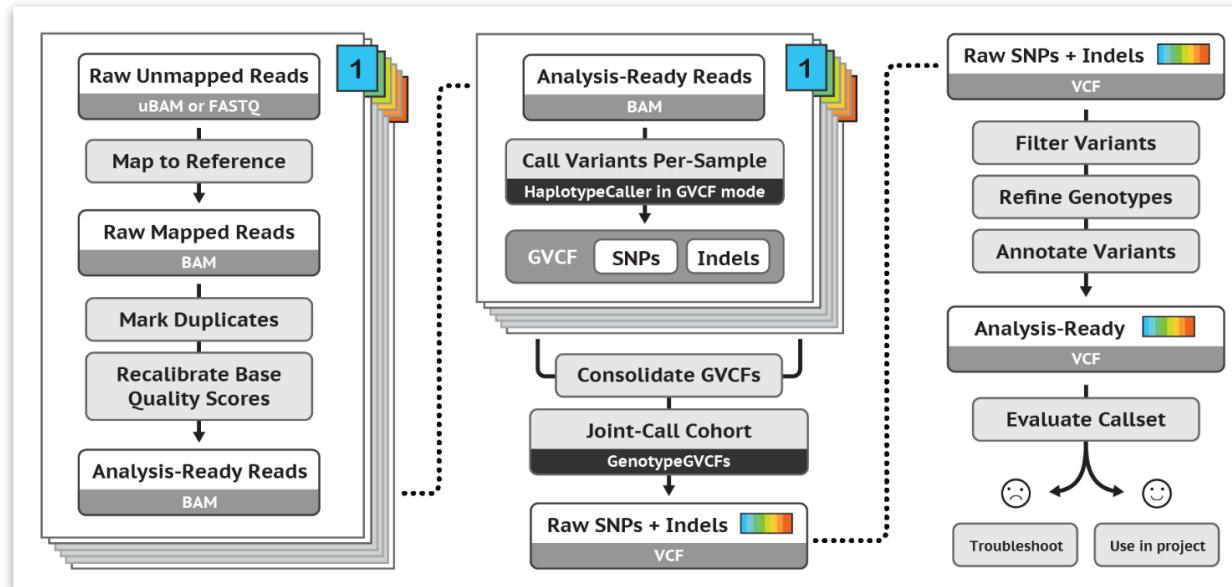


Forum

Ask our team for help and report issues

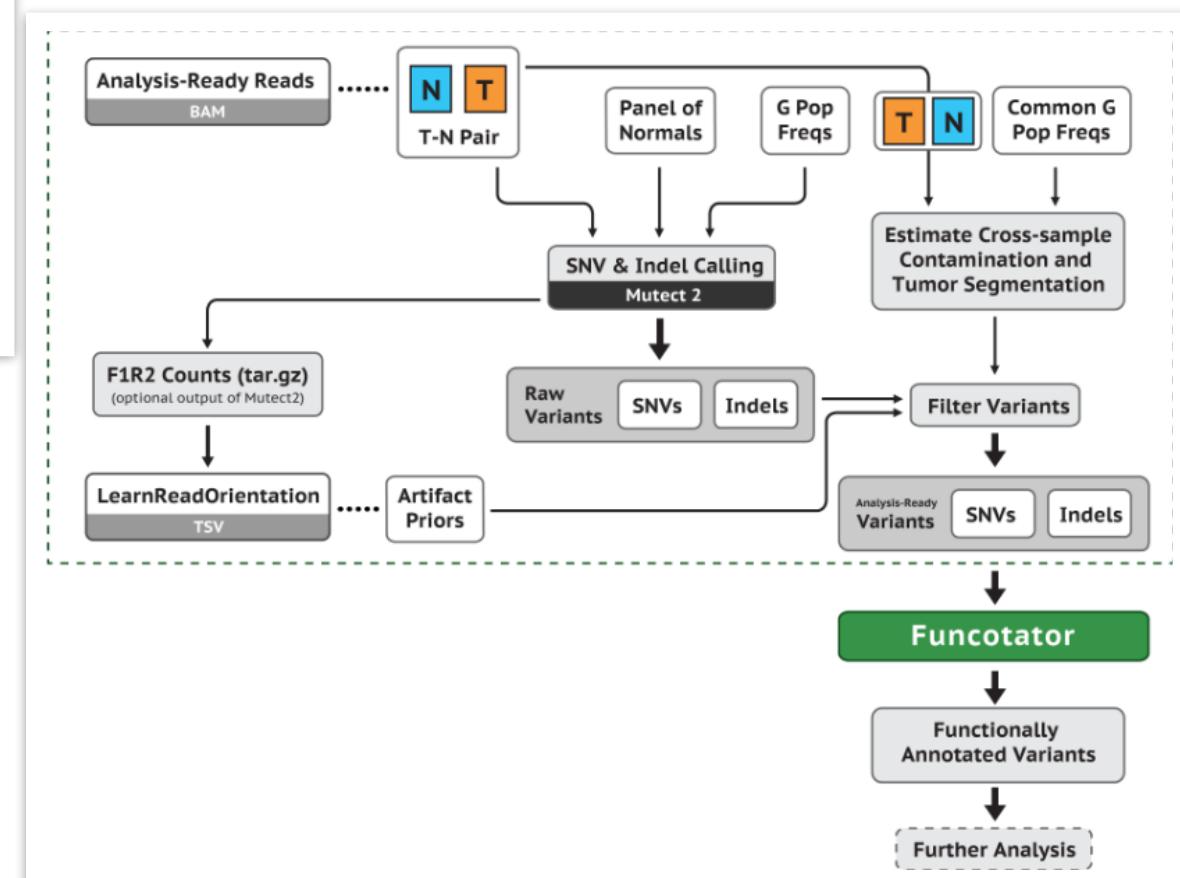


GATK Best Practices



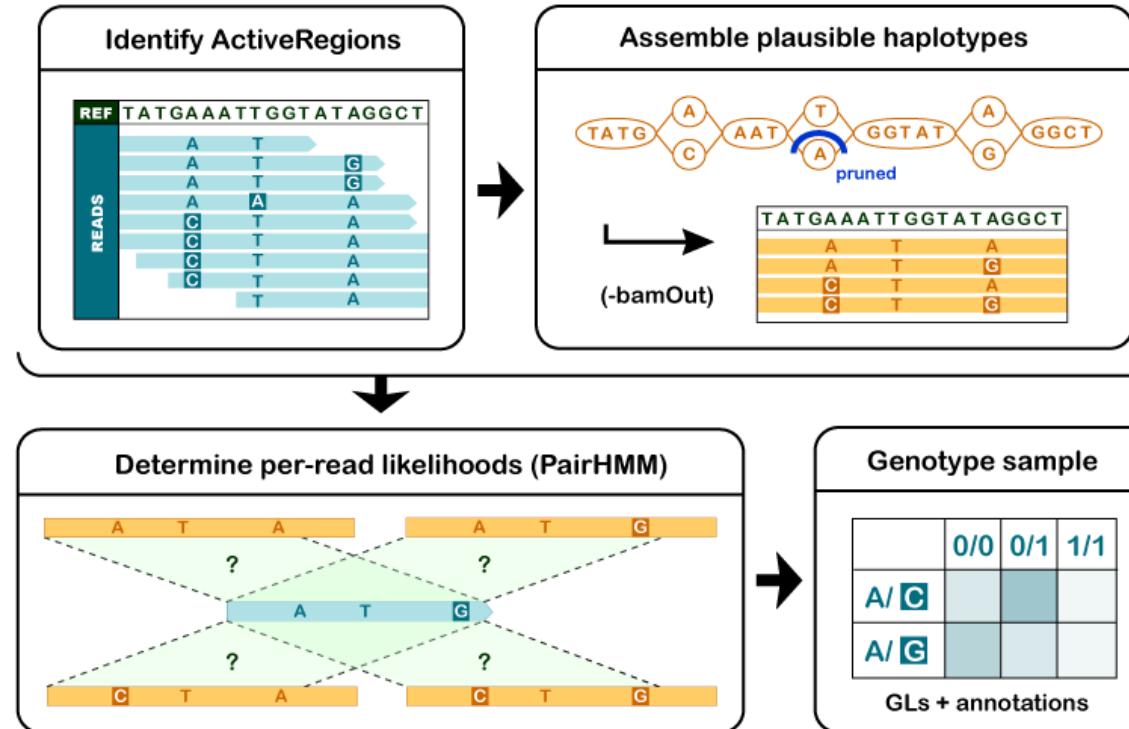
Somatic small-scale variants

Germline small-scale variants



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://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>

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 ~ 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 ¹	Sequenom	97	92	5	94.85%
Ovarian ²	Sequenom/PCR/454	1655	1483	172	89.61%
Ovarian ²	Capture/Illumina	6497	6232	265	95.92%
Head and Neck ³	Sequenom	321	288	33	89.72%
Breast ⁴	Sequenom/PCR/454	455	428	0	94.07%

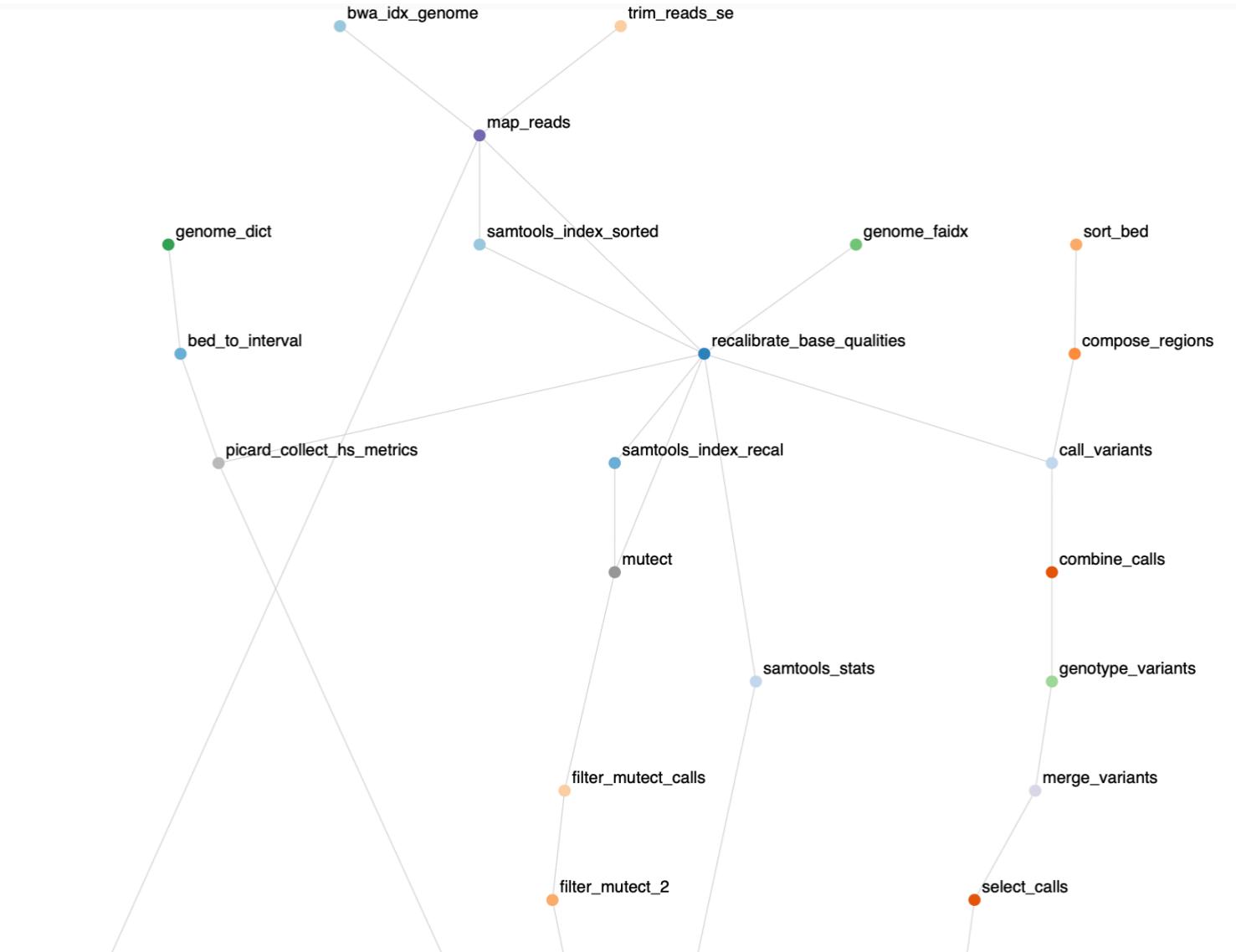
<https://gatk.broadinstitute.org/hc/en-us/articles/360037593851-Mutect2>

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

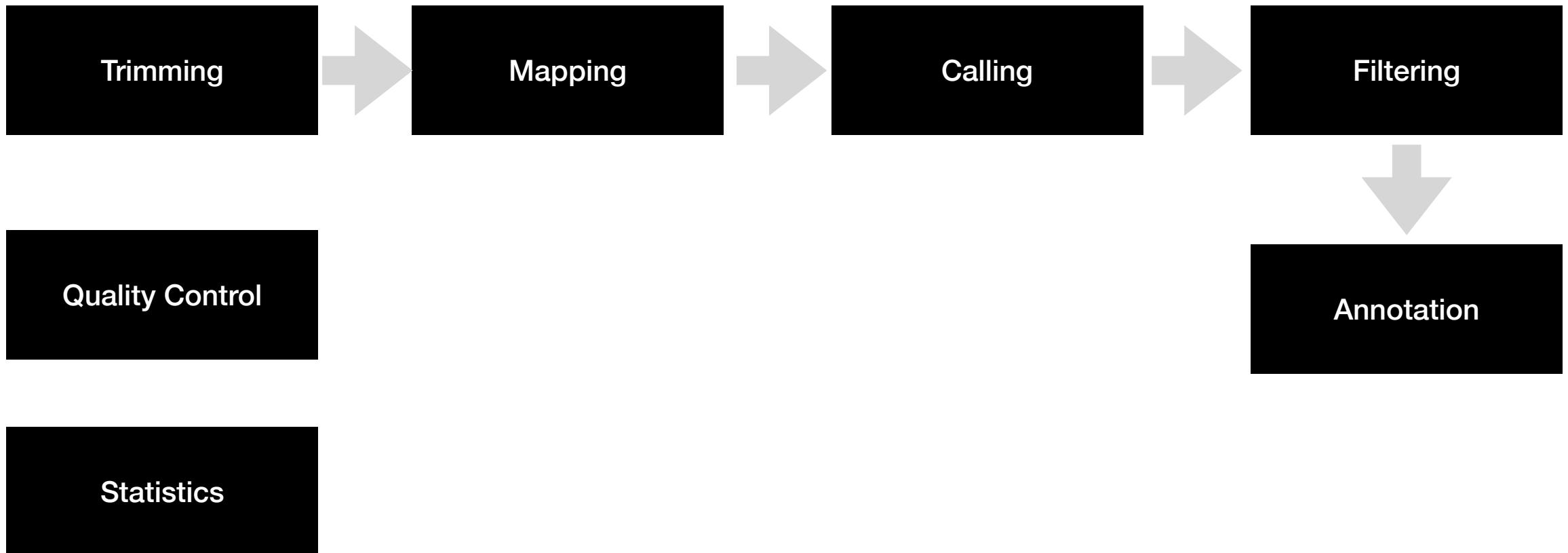
varca



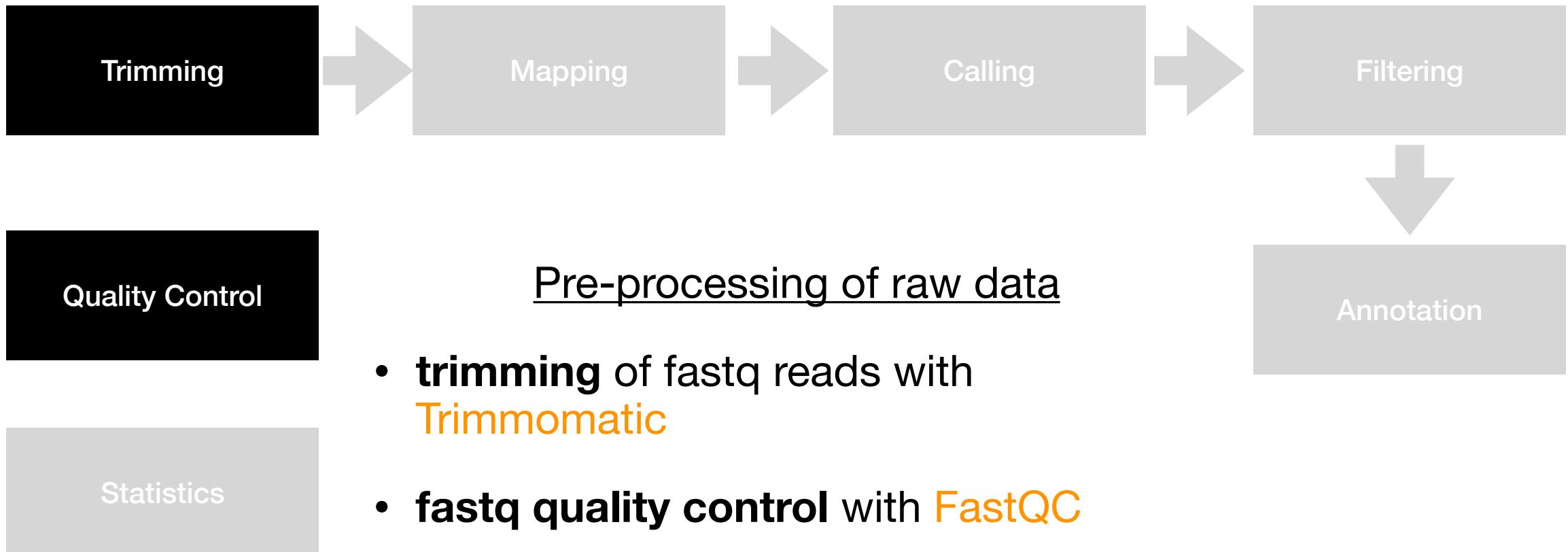
https://gitlab.com/bu_cnio/varca



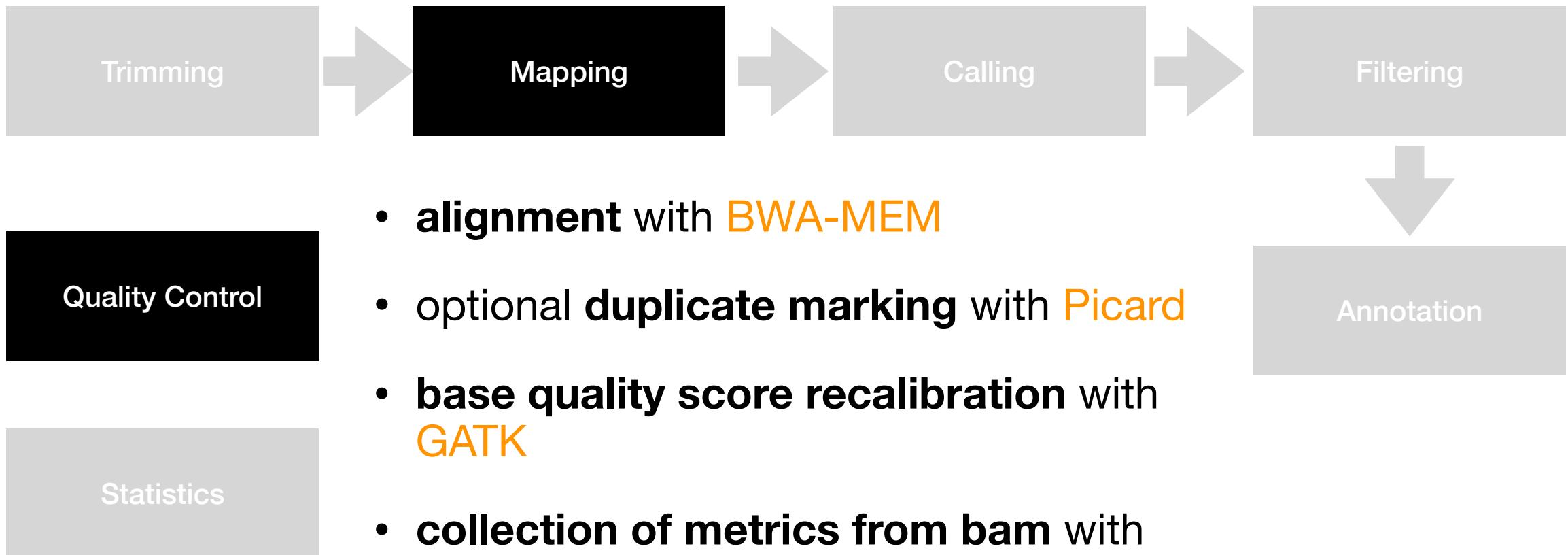
Steps in the pipeline



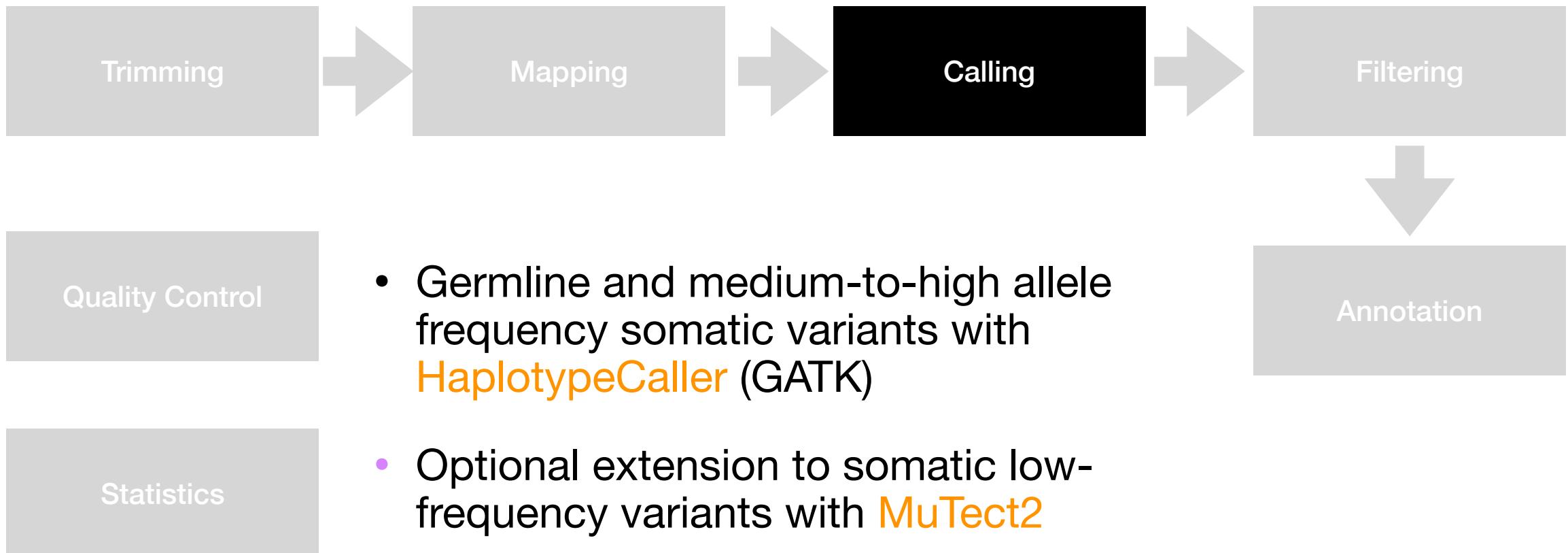
Steps in the pipeline



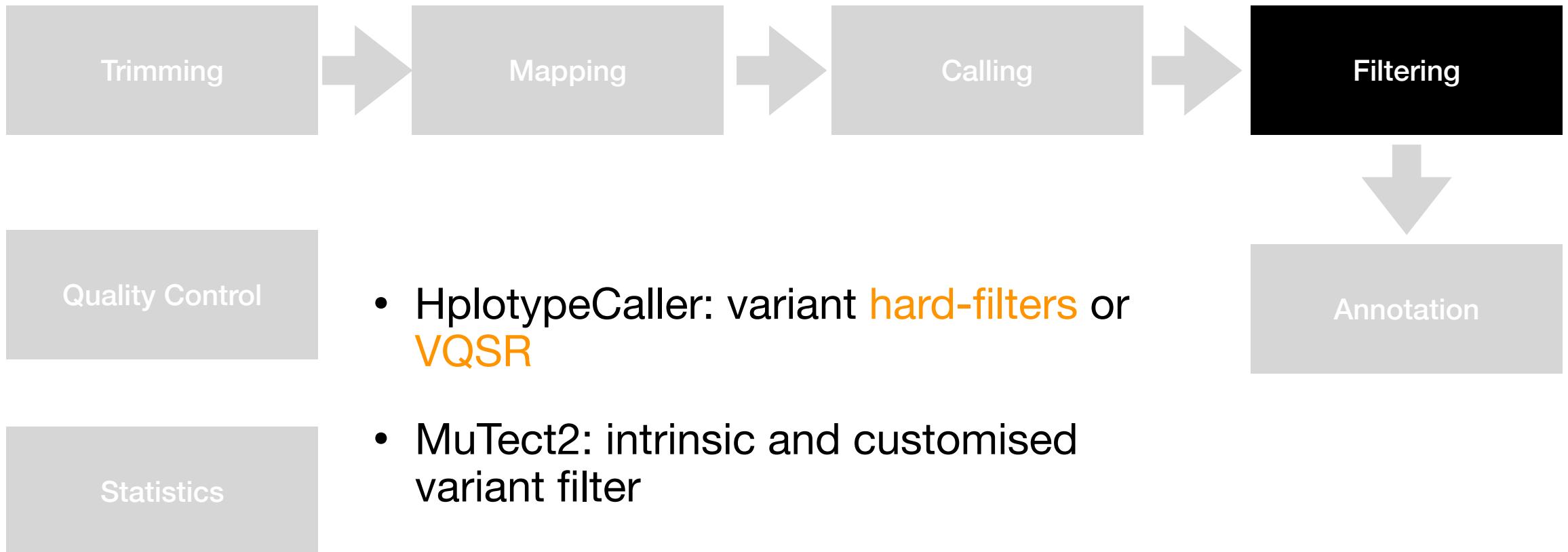
Steps in the pipeline



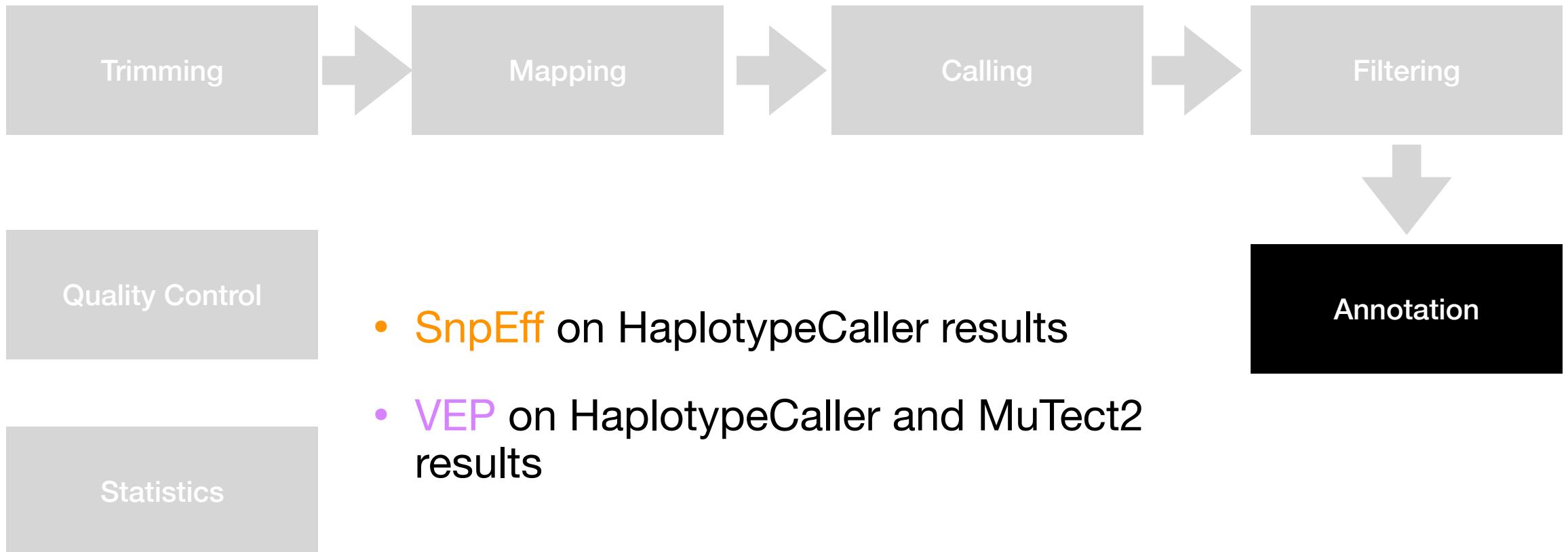
Steps in the pipeline



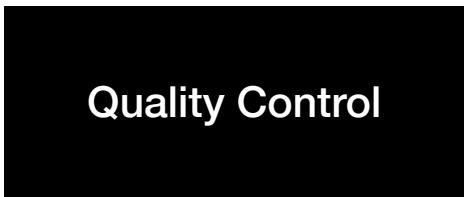
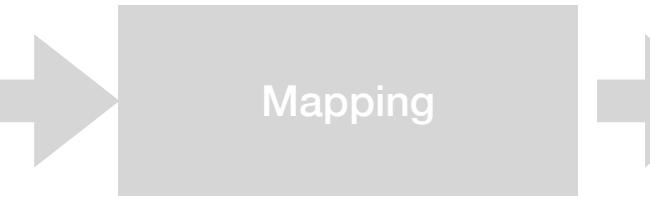
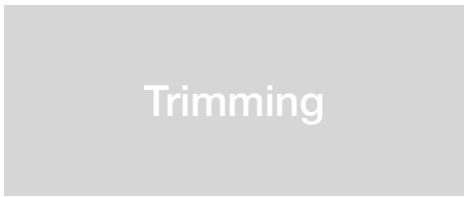
Steps in the pipeline



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 various analysis modules: General Stats, SnpEff, Picard, HSMetrics, Target Region Coverage, Mark Duplicates, Samtools, Percent Mapped, and Alignment metrics. The main content area has tabs for "General Statistics" and "SnpEff". The "General Statistics" tab displays a table of general metrics for samples "all", "normal-1", and "tumor-1". The "SnpEff" tab provides a brief description of the tool and links to its documentation.

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

Configuration files

2 configuration files with the **samples** definition

1 configuration file with the **contigs** to analyse

1 configuration file with **paths** to additional files used in the process **and parameters**

Configuration files

samples.tsv and **units.tsv**

```
1 sample control  
2 A B  
3 B -
```

unable MuTect2 calling on the sample

samples

list of samples (sample column)

establishes the MuTect2 calling mode (control column)

- No MuTect2 execution:

sample	control
A	-

sample	control
A	A

- MuTect2 execution in tumor-only mode:

sample	control
A	B
B	-

being A:tumor, B:normal

Configuration files

samples.tsv and **units.tsv**

1	sample	control
2	A	B
3	B	-

unable MuTect2 calling on the sample

samples

list of samples (sample column)
establishes the MuTect2 calling mode (control column)

paired-end and single-end definition

1	sample	unit
2	A	1
3	B	1
4	B	2

platform
ILLUMINA

fq1	fq2
data/reads/a.chr21.1.fq	data/reads/a.chr21.2.fq
data/reads/b.chr21.1.fq	data/reads/b.chr21.2.fq
data/reads/b.chr21.1.fq	

units

sequential to identify multiple sequences for the same sample

Configuration files

contigs.tsv

1	chr1			
2	chr2			
3	chr3			
4	chr4			
5	chr5			
6	chr6			
7	chr7	@SQ	SN:chr1_KI270708v1_random	LN:127682 M5:1e95e047b98ed92148dd84d6c037158c
8	chr8	@SQ	SN:chr1_KI270709v1_random	LN:66860 M5:4e2db2933ea96aee8dab54af60ecb37d
9	chr9	@SQ	SN:chr1_KI270710v1_random	LN:40176 M5:9949f776680c6214512ee738ac5da289
10	chr10	@SQ	SN:chr1_KI270711v1_random	LN:42210 M5:af383f98cf4492c1f1c4e750c26ccb40
11	chr11	@SQ	SN:chr1_KI270712v1_random	LN:176043 M5:c38a0fecae6a1838a405406f724d6838
12	chr12	@SQ	SN:chr1_KI270713v1_random	LN:40745 M5:cb78d48cc0adbc58822a1c6fe89e3569
13	chr13	@SQ	SN:chr1_KI270714v1_random	LN:41717 M5:42f7a452b8b769d051ad738ee9f00631
14	chr14	@SQ	SN:chr2_KI270715v1_random	LN:161471 M5:b65a8af1d7bbb7f3c77eea85423452bb
15	chr15	@SQ	SN:chr2_KI270716v1_random	LN:153799 M5:2828e63b8edc5e845bf48e75fbad2926
16	chr16	@SQ	SN:chr3_GL000221v1_random	LN:155397 M5:3238fb74ea87ae857f9c7508d315babb
17	chr17	@SQ	SN:chr4_GL000008v2_random	LN:209709 M5:a999388c587908f80406444cebe80ba3
18	chr18	@SQ	SN:chr5_GL000208v1_random	LN:92689 M5:aa81be49bf3fe63a79bdc6a6f279abf6
19	chr19	@SQ	SN:chr9_KI270717v1_random	LN:40062 M5:796773a1ee67c988b4de887addbed9e7
20	chr20	@SQ	SN:chr9_KI270718v1_random	LN:38054 M5:b0c463c8efa8d64442b48e936368dad5
21	chr21	@SQ	SN:chr9_KI270719v1_random	LN:176845 M5:cd5e932cfc4c74d05bb64e2126873a3a
22	chr22	@SQ	SN:chr9_KI270720v1_random	LN:39050 M5:8c2683400a4aeeb40abff96652b9b127
23	chrX	@SQ	SN:chr11_KI270721v1_random	LN:100316 M5:9654b5d3f36845bb9d19a6dbd15d2f22
24	chrY	@SQ	SN:chr14_GL000009v2_random	LN:201709 M5:862f555045546733591ff7ab15bcecbe
25	chrM	@SQ	SN:chr14_GL000225v1_random	LN:211173 M5:63945c3e6962f28ffd469719a747e73c
		@SQ	SN:chr14_KI270722v1_random	LN:194050 M5:51f46c9093929e6edc3b4dfd50d803fc
		@SQ	SN:chr14_GL000194v1_random	LN:191469 M5:6ac8f815bf8e845bb3031b73f812c012
		@SQ	SN:chr14_KI270723v1_random	LN:38115 M5:74a4b480675592095fb0c577c515b5df

~3300 contigs in hg38

Configuration files

config.yaml

```
1 samples: samples.tsv
2 units: units.tsv
3 contigs: contigs.tsv
4
5 outdir: "out"
6 logdir: "log"
7
8 ref:
9 # Genome database of snpeff to be used in the annotation with this resource. Available databases can be checked with java -jar
10 name: GRCh38.86
11 # Path to the reference genome, ideally as it is provided by the GATK bundle.
12 genome: /home/epineiro/REFERENCES/Homo_sapiens/HG38/hg38.fa
13 # Path to any database of known variants, ideally as it is provided by the GATK bundle.
14 known-variants: /home/epineiro/REFERENCES/Homo_sapiens/HG38/dbsnp_151.hg38.vcf
15
16 filtering:
17 # Set to true in order to apply machine learning based recalibration of
18 # quality scores instead of hard filtering.
19 vqsr: false
20 hard:
21 # hard filtering as outlined in GATK docs
22 # (https://gatkforums.broadinstitute.org/gatk/discussion/2806/howto-apply-hard-filters-to-a-call-set)
23 snvs:
24 "QD < 2.0 || QUAL < 100.0 || DP < 50.0 || SOR > 3.0 || FS > 60.0 || MQ < 40.0 || MQRankSum < -12.5 || ReadPosRankSum < -8.0"
25 indels:
26 "QD < 2.0 || QUAL < 100 || DP < 10.0 || FS > 200.0 || ReadPosRankSum < -20.0"
27 #depth of coverage threshold to apply to variants identified with MuTect2
28 depth: "DP < 30"
```

Reference Genome and known variants file

Variant filtering system and parameters for hard filtering

Configuration files

config.yaml

```
30 processing:  
31   remove-duplicates: true  
32  
33   # see https://gatkforums.broadinstitute.org/gatk/discussion/4133/when-should-i-use-l-to-pass-in-a-list-of-intervals.  
34   restrict-regions: /home/epineiro/REFERENCES/Homo_sapiens/HG38/S04380219_Padded_v5_hg38_varca.bed  
35   # If regions are restricted, uncomment this to enlarge them by the given value in order to include  
36   # flanking areas.  
37   # region-padding: 100  
38  
39 params:  
40   gatk:  
41     HaplotypeCaller: ""  
42     BaseRecalibrator: ""  
43     GenotypeGVCFs: ""  
44     #If vqsr set to true, fill the following section with the required information (see VariantRecalibrator documentation in GATK)  
45     VariantRecalibrator:  
46       #paths to the necessary resources  
47       hapmap: "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_hapmap_3.3.hg38.vcf.gz"  
48       omni: "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_1000G_omni2.5.hg38.vcf.gz"  
49       g1k: "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_1000G_phase1.snps.high_confidence.hg38.vcf.gz"  
50       dbsnp: "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_Homo_sapiens_assembly38.dbsnp138.vcf"  
51       #paths to the resource indexes  
52       aux: ["/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_hapmap_3.3.hg38.vcf.gz.tbi",  
53           "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_1000G_omni2.5.hg38.vcf.gz.tbi",  
54           "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_1000G_phase1.snps.high_confidence.hg38.vcf.gz.tbi",  
55           "/home/epineiro/REFERENCES/Homo_sapiens/HG38/resources_broad_hg38_v0_Homo_sapiens_assembly38.dbsnp138.vcf.tbi"]  
56       #set definition for each resource  
57       parameters: {"hapmap": {"known": False, "training": True, "truth": True, "prior": 15.0},  
58                     "omni": {"known": False, "training": True, "truth": False, "prior": 12.0},  
59                     "g1k": {"known": False, "training": True, "truth": False, "prior": 10.0},  
60                     "dbsnp": {"known": True, "training": False, "truth": False, "prior": 2.0}}  
61  
#Names of the annotations used for calculations
```

Remove duplicates in the alignment?

Restrict variants to specific regions?

Configuration parameters for trimming, deduplication and GATK variant calling, vqsr

Configuration files

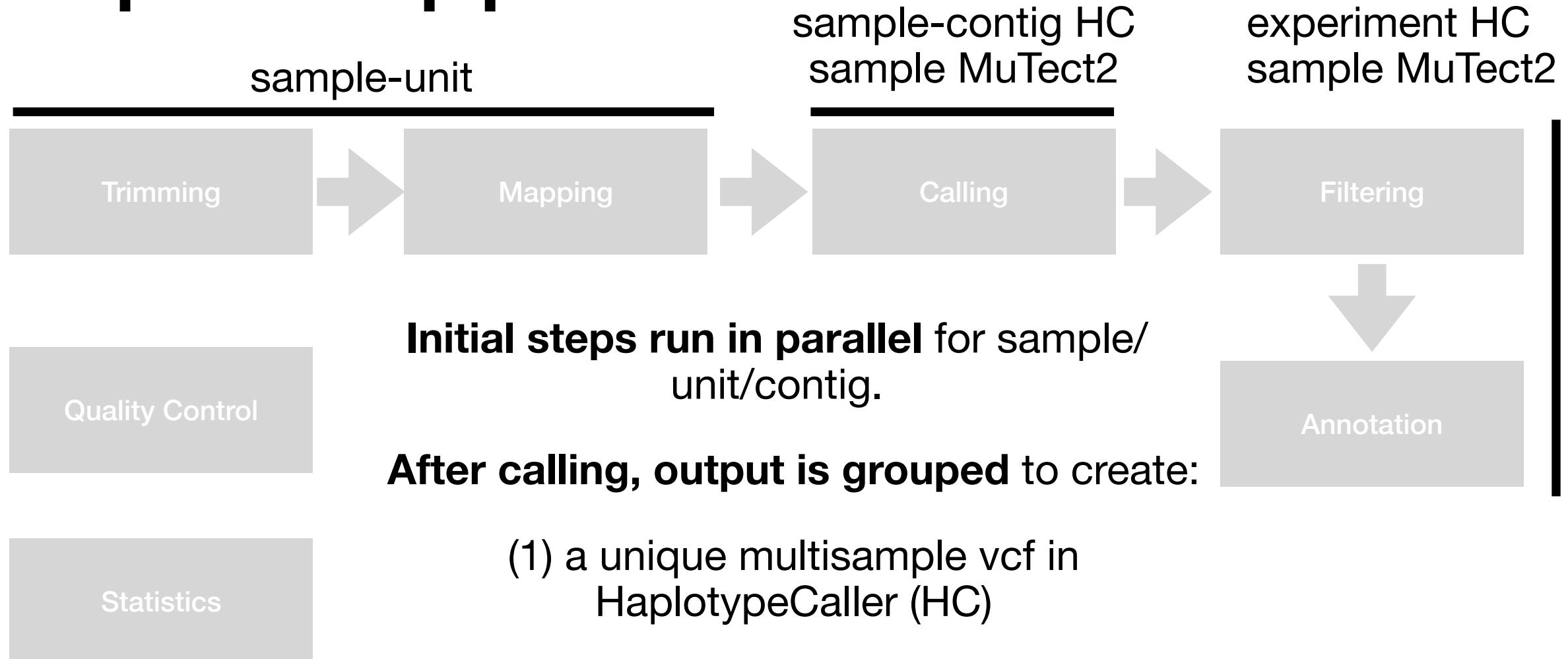
config.yaml

```
82 annotation:  
83   vep:  
84     # If set to true, a cache directory must be created using the instructions of VEP. Using a cache allows for a fastest and  
85     cache: true  
86     #Indicate the version of the cache if cache is set to true. e.g. 100. This version must match the one indicated in the vep  
87     cache_version: 100  
88     #Indicate the path to the cache files if cache is set to true  
89     cache_directory: /home/epineiro/Programs/VEP/VEP100/.vep/  
90     #Indicate the assembly version of the reference genome  
91     assembly: GRCh38  
92     #Indicate the annotations to include with the VEP execution and additional parameters. See the documentation (https://www.ensembl.org/papers/vep/)  
93     annotations: "--sift b --polyphen b --ccds --uniprot --hgvs --symbol --numbers --domains --regulatory --canonical --protein  
94  
95 resources:  
96   default:  
97     threads: 1  
98     mem: 32000  
99     walltime: 720  
100  snpeff:  
101    mem: 8000  
102    walltime: 480  
103  call_variants:  
104    mem: 10000  
105    threads: 4  
106    walltime: 480  
107  combine_calls:  
108    mem: 8000  
109    threads: 1  
110    walltime: 480  
111  recalibrate_calls:  
112    mem: 8000  
113    walltime: 480
```

Parameters for the annotation with VEP

Machine requirements

Steps in the pipeline



Result's files

VCF multisample

Output directory

trimmed
mapped
dedup
recal -> bam files

called
genotyped
filtered -> vcf file with all HC variants and PASS label
snpeff
annotated -> annotated vcf

qc
mutect
mutect_filter-> vcf files with MuTect2 variants and PASS label

```
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT normal tumor
chr1 13656 . CAG C 573.29 PASS AC=4;AF=1.00;AN=4;DP=14;ExcessHet=3.0103;FS=0.000;MLEAC=4;MLEAF=1.00;MQ=22.61;QD=25.36;SOR=5.421
GT:AD:DP:GQ:PL 1/1:0,10:10:30:450,30,0 1/1:0,4:4:12:159,12,0
chr1 14653 . C T 75.33 snv-hard-filter
AC=1;AF=0.500;AN=2;BaseQRankSum=-2.530e-01;ClippingRankSum=0.00;DP=5;ExcessHet=3.0103;FS=0.000;MLEAC=1;MLEAF=0.500;MQ=24.61;MQRankSum=-2.530e-01;QD=15.07;ReadPosRankSum=0.524;SOR=0.693 GT:AD:DP:GQ:PL 0/1:1,4:5:15:102,0,15 ./:0,0:0::0,0,0
chr1 14907 . A G 2162.44 PASS
AC=2;AF=0.500;AN=4;BaseQRankSum=1.26;ClippingRankSum=0.00;DP=117;ExcessHet=4.7712;FS=4.047;MLEAC=2;MLEAF=0.500;MQ=44.68;MQRankSum=-2.370e-01;QD=18.64;ReadPosRankSum=0.183;SOR=1.129 GT:AD:DP:GQ:PL 0/1:17,47:64:99:1169,0,391 0/1:12,40:52:99:1022,0,201
chr1 14930 . A G 3506.44 PASS
AC=2;AF=0.500;AN=4;BaseQRankSum=0.814;ClippingRankSum=0.00;DP=160;ExcessHet=4.7712;FS=6.365;MLEAC=2;MLEAF=0.500;MQ=44.54;MQRankSum=-2.025e+00;QD=21.92;ReadPosRankSum=0.461;SOR=0.709 GT:AD:DP:GQ:PL 0/1:19,64:83:99:1814,0,306 0/1:12,65:77:99:1721,0,145
chr1 14933 . G A 10.92 snv-hard-filter
AC=2;AF=0.500;AN=4;BaseQRankSum=-5.500e-01;ClippingRankSum=0.00;DP=158;ExcessHet=4.7712;FS=4.667;MLEAC=2;MLEAF=0.500;MQ=44.88;MQRankSum=-2.764e+00;QD=0.07;ReadPosRankSum=0.629;SOR=1.231 GT:AD:DP:GQ:PL 0/1:70,9:79:15,0,2269 0/1:69,10:79:24:24,0,2001
chr1 14976 . G A 16.95 snv-hard-filter
AC=1;AF=0.250;AN=4;BaseQRankSum=0.210;ClippingRankSum=0.00;DP=123;ExcessHet=3.0103;FS=3.614;MLEAC=1;MLEAF=0.250;MQ=41.60;MQRankSum=-1.883e+00;QD=0.31;ReadPosRankSum=-1.300e-02;SOR=0.697 GT:AD:DP:GQ:PL 0/1:48,7:55:46:46,0,1362 0/0:68,0:68:9:0,9,1771
chr1 16949 . A C 59.44 snv-hard-filter
AC=2;AF=0.500;AN=4;BaseQRankSum=2.74;ClippingRankSum=0.00;DP=18;ExcessHet=4.7712;FS=0.000;MLEAC=2;MLEAF=0.500;MQ=23.95;MQRankSum=0.728;QD=3.30;ReadPosRankSum=-7.890e-01;SOR=0.169 GT:AD:DP:GQ:PL 0/1:5,2:7:35:35,0,125 0/1:8,3:11:53:53,0,195
chr1 17020 . G A 77.44 snv-hard-filter
AC=2;AF=0.500;AN=4;BaseQRankSum=-1.440e+00;ClippingRankSum=0.00;DP=12;ExcessHet=4.7712;FS=3.282;MLEAC=2;MLEAF=0.500;MQ=23.53;MQRankSum=0.253;QD=6.45;ReadPosRankSum=-4.310e-01;SOR=2.303 GT:AD:DP:GQ:PL 0/1:4,2:6:38:38,0,92 0/1:3,3:6:68:68,0,68
chr1 17385 . G A 91.44 snv-hard-filter
AC=2;AF=0.500;AN=4;BaseQRankSum=-8.760e-01;ClippingRankSum=0.00;DP=18;ExcessHet=4.7712;FS=0.000;MLEAC=2;MLEAF=0.500;MQ=35.20;MQRankSum=0.660;QD=5.08;ReadPosRankSum=1.07;SOR=0.602 GT:AD:DP:GQ:PL 0/1:5,2:7:37:37,0,161 0/1:8,3:11:83:83,0,258
```

Elena Piñeiro-Yáñez

epineiro@cnio.es

Coral Fustero-Torre

cfustero@cnio.es

María José Jiménez-Santos

mjjimenez@cnio.es

Fátima Al-Shahrour

falshahrour@cnio.es