

# NGS I : VARIANT DETECTION

---

Javier Perales-Patón  
[jperales@cnio.es](mailto:jperales@cnio.es)



Translational Bioinformatics Unit  
CNIO. Madrid, Spain.

Fátima Al-Shahrour  
[falshahrour@cnio.es](mailto:falshahrour@cnio.es)

Elena Piñeiro-Yáñez  
[epineiro@cnio.es](mailto:epineiro@cnio.es)

Pedro Fernandes  
[pfern@igc.gulbenkian.pt](mailto:pfern@igc.gulbenkian.pt)

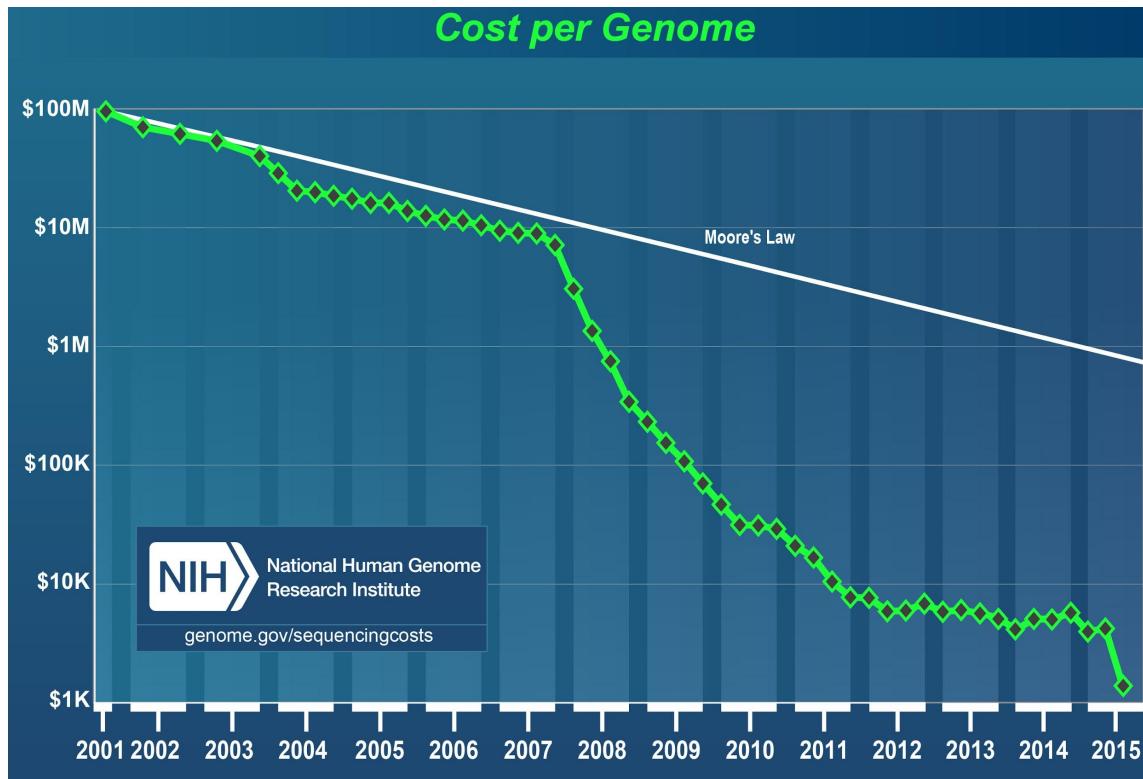


INSTITUTO  
GULBENKIAN  
DE CIÊNCIA

# Today

09:30 - 10:00	Introduction to the course and self presentation of the participants. Personalized medicine.
<b>11:30 - 12:30</b>	<b>NGS I : Variant detection.</b>
14:00 - 16:00	Playing with the data and the methods.
16:30 - 18:00	Practical : Running the pipeline.

# Sequencing cost has been coming down



# Sequencing cost has been coming down

**Cost per Genome**

\$10 Mardis *Genome Medicine* 2010, 2:84  
<http://genomemedicine.com/content/2/11/84>

 Genome **Medicine**

**MUSINGS**

## The \$1,000 genome, the \$100,000 analysis?

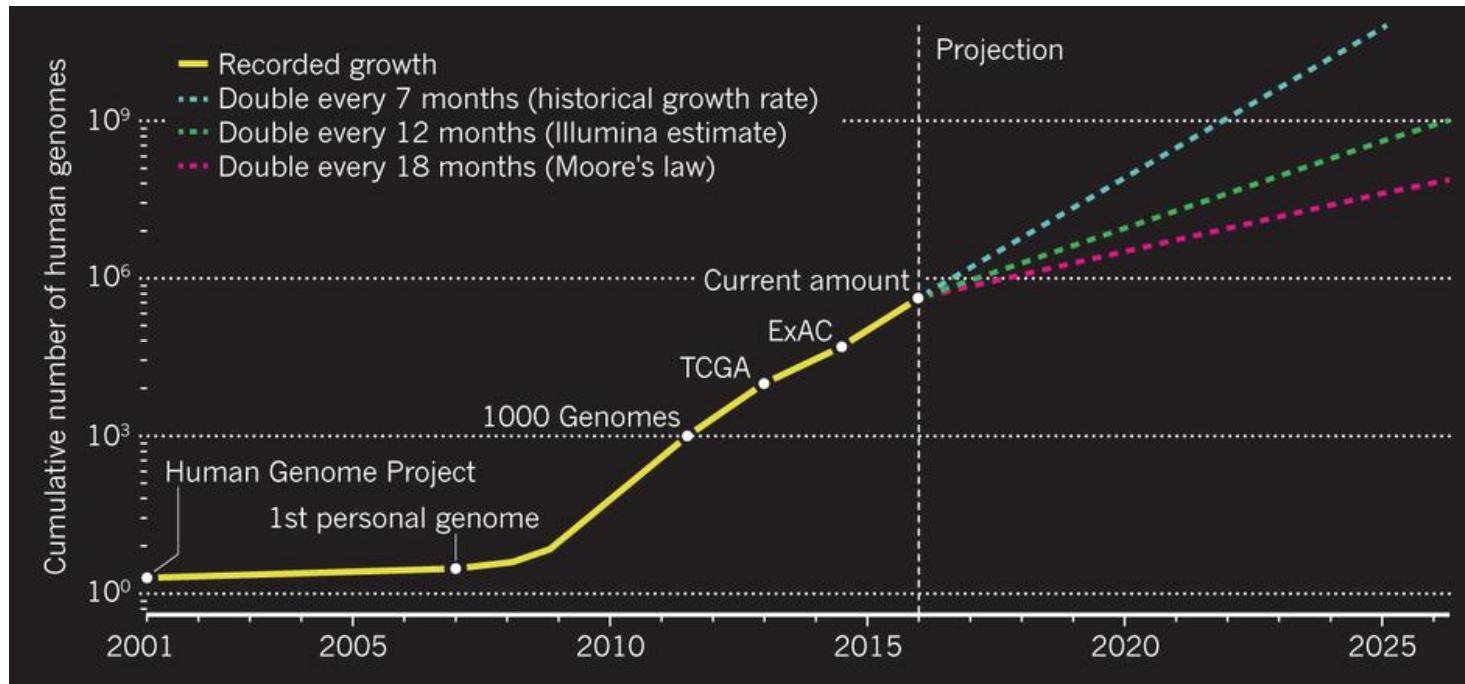
Elaine R Mardis\*

Although each presenter emphasized the rapidity with which these data can now be generated using next-generation sequencing instruments, they also listed the large number of people involved in the analysis of these datasets.

[...]

The large number of specialists was critical for the completion of the data analysis, the annotation of variants, the interpretive 'filtering' necessary to deduce the causative or 'actionable' variants, the clinical verification of these variants, and the communication of results and their ramifications to the treating physician, and ultimately to the patient. At the end of the day, although the idea of clinical whole-genome sequencing for diagnosis is exciting and potentially life-changing for these patients, one does wonder how, in the clinical translation required for this practice to become common-place, such a 'dream team' of specialists would be assembled for each case.

# DNA sequencing soars



- + 1000 Genomes Project : hundreds of genomes.
- + TCGA : thousands (genome & exomes).
- + ExAC : > 60,000 exomes.

Stephens ZD et al. **Big Data: Astronomical or Genomical?**. PLoS Biol. 2015 Jul 7;13(7)

Eisenstein M. **Big data: The power of petabytes**. Nature. 2015 Nov 5;527(7576):S2-4.

# 1000 Genomes Project

## ARTICLE

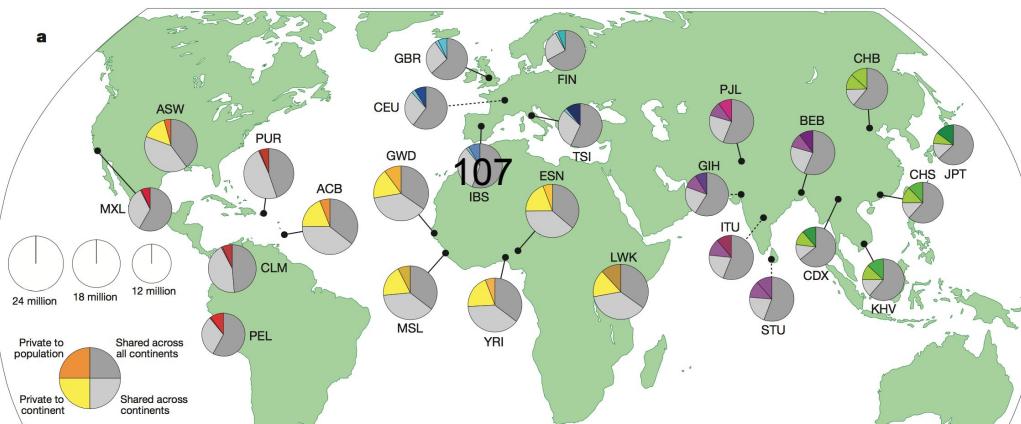
OPEN

doi:10.1038/nature15393

### A global reference for human genetic variation

The 1000 Genomes Project Consortium\*

The 1000 Genomes Project set out to provide a comprehensive description of common human genetic variation by applying whole-genome sequencing to a diverse set of individuals from multiple populations. Here we report completion of the project, having reconstructed the genomes of 2,504 individuals from 26 populations using a combination of low-coverage whole-genome sequencing, deep exome sequencing, and dense microarray genotyping. We characterized a broad spectrum of genetic variation, in total over 88 million variants (84.7 million single nucleotide polymorphisms (SNPs), 3.6 million short insertions/deletions (indels), and 60,000 structural variants), all phased onto high-quality haplotypes. This resource includes >99% of SNP variants with a frequency of >1% for a variety of ancestries. We describe the distribution of genetic variation across the global sample, and discuss the implications for common disease studies.



Phase 3	WGS	WExS
Raw bases	89 Tb	18 Tb
Samples	2,504	2,504
Region	Genome	Exome
Mean Depth	8.45x	75x
SNPs	85M	1.5M
Indels	3.6M	22K
Structural Variants	60K	6.5K
Het. Concordance (SNPs)	99.4%	99.8%

<http://www.1000genomes.org/about#ProjectSamples> ; Phase 1 n=1092 → Phase 3 n=2504

**ARTICLE**

**OPEN**  
doi:10.1038/nature19057

## Analysis of protein-coding genetic variation in 60,706 humans

Large-scale reference data sets of human genetic variation are critical for the medical and functional interpretation of DNA sequence changes. Here we describe the aggregation and analysis of high-quality exome (protein-coding region) DNA sequence data for 60,706 individuals of diverse ancestries generated as part of the Exome Aggregation Consortium (ExAC). This catalogue of human genetic diversity contains an average of one variant every eight bases of the exome, and provides direct evidence for the presence of widespread mutational recurrence. We have used this catalogue to calculate objective metrics of pathogenicity for sequence variants, and to identify genes subject to strong selection against various classes of mutation; identifying 3,230 genes with near-complete depletion of predicted protein-truncating variants, with 72% of these genes having no currently established human disease phenotype. Finally, we demonstrate that these data can be used for the efficient filtering of candidate disease-causing variants, and for the discovery of human 'knockout' variants in protein-coding genes.



An average of 85 heterozygous and 35 homozygous Protein-truncating variants per individual. The majority found in any one person are common, and only 2 are singleton.

**Blog Post:** <http://blogs.nature.com/freeassociation/2016/08/joint-calling-of-the-exac-publications.html>

**Paper :** [Lek M et al. Analysis of protein-coding genetic variation in 60,706 humans. \(2016\) Nature.](#)

# The objective of the Variant Detection:

---

Identify the most likely **genotype** for each genomic position from the individual.

- - -

In Cancer genomics, if there is a **matched-normal sample** to be compared against the tumour sample:

- + Identify **somatic variants** (i.e. only in tumour sample).
- + Identify **copy-number alterations** (large genomic aberrations).

# Some concepts

- What is a genetic **variant** ?

Genetic differences in individuals as compared to a reference genome (built from a population).

**Nomenclature:**

- **First level:** Genomic position and nucleotide change.

**Chromosome Name:** Genomic position (coordinates): Reference allele > Alternative allele  
Chr12:25398284-25398284:G>A

- Classes of variants :

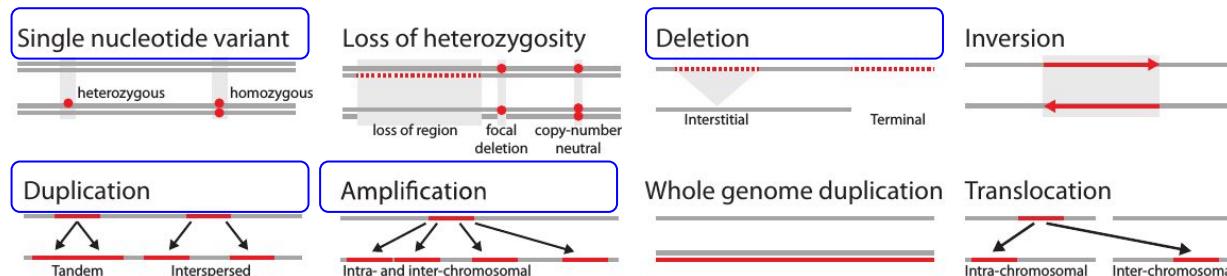
- **Germline** : inherited. E.g. a SNP, or a SNV related to a rare disease.
- **Somatic** : acquired within a cell lineage. E.g. Cancer mutations.

- Polymorphism

common variant in a given population (SNP, Single Nucleotide Polymorphism).

Present in at least 1% in a population.

- Types of genomic variants:



# Variants need a context

Example:

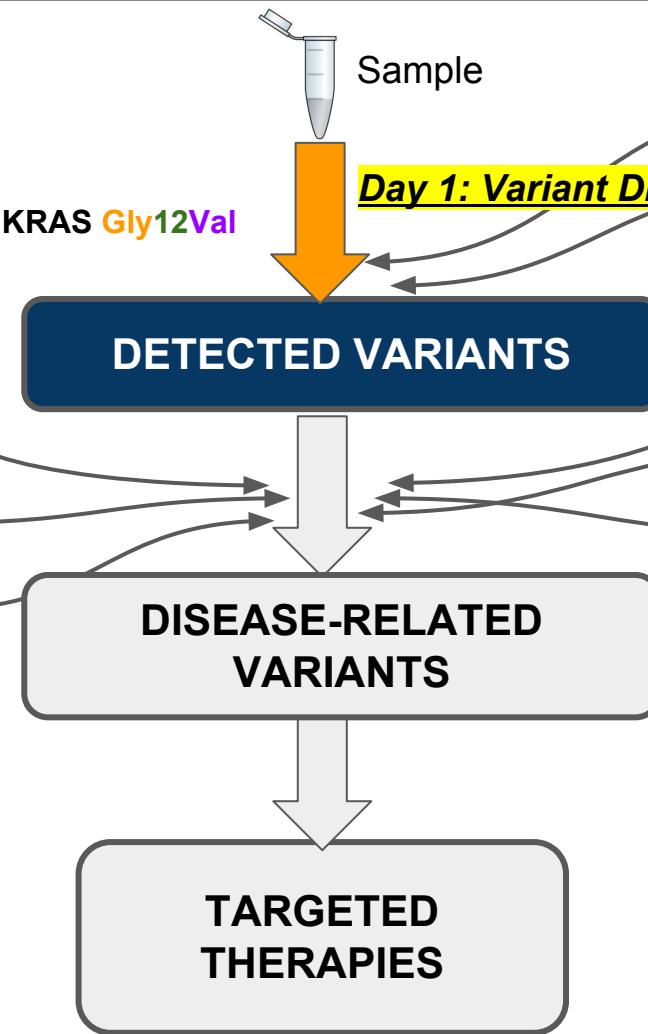
chr12:25398284-25398284:G>A → KRAS Gly12Val

White list



ClinVar  
CTGAGGAGAACT  
TACAAGACAGGT

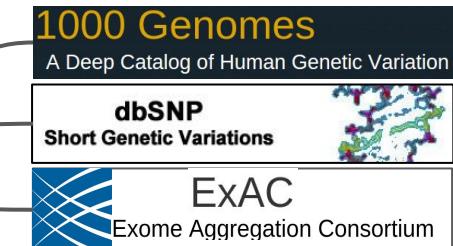
OMIM  
Online Mendelian Inheritance in Man



Prior-knowledge for modeling

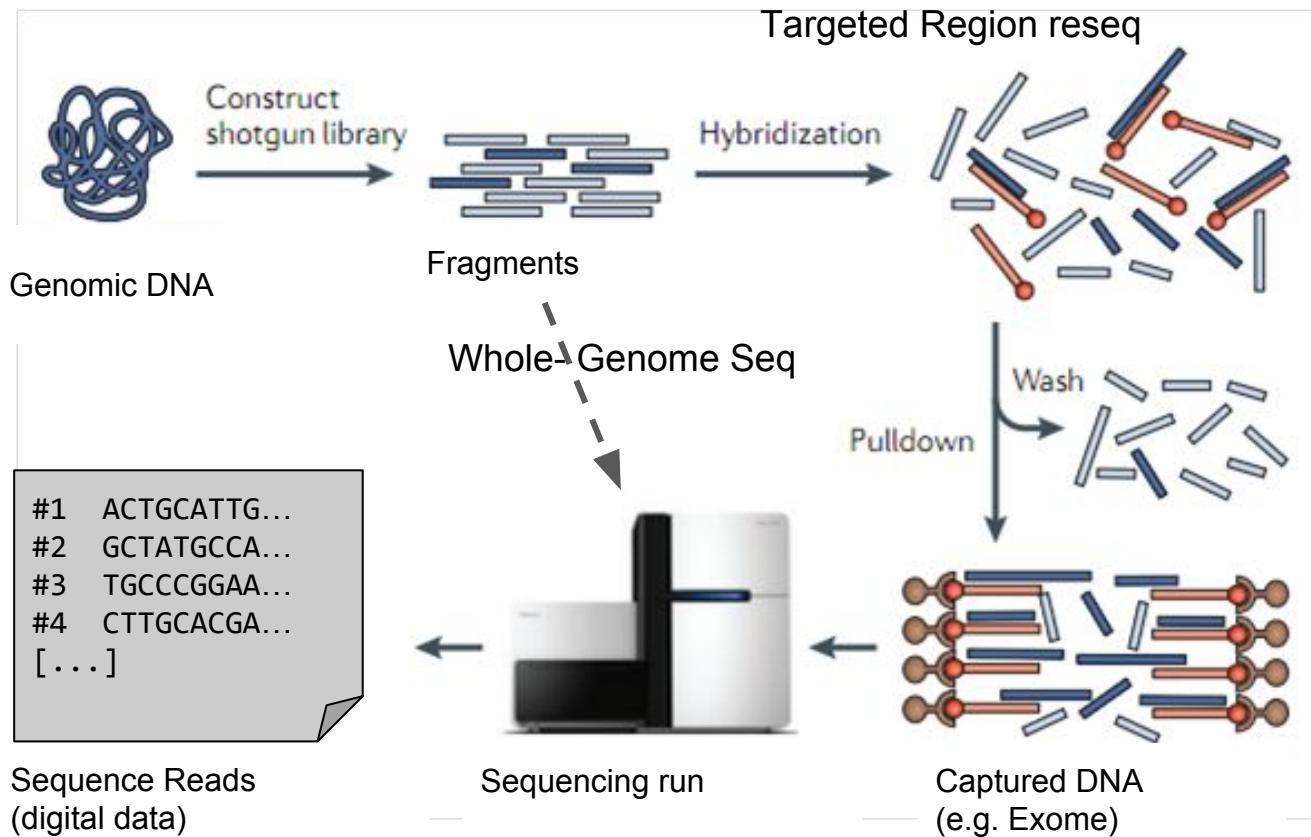


Black list



Day 4: real cases studies!

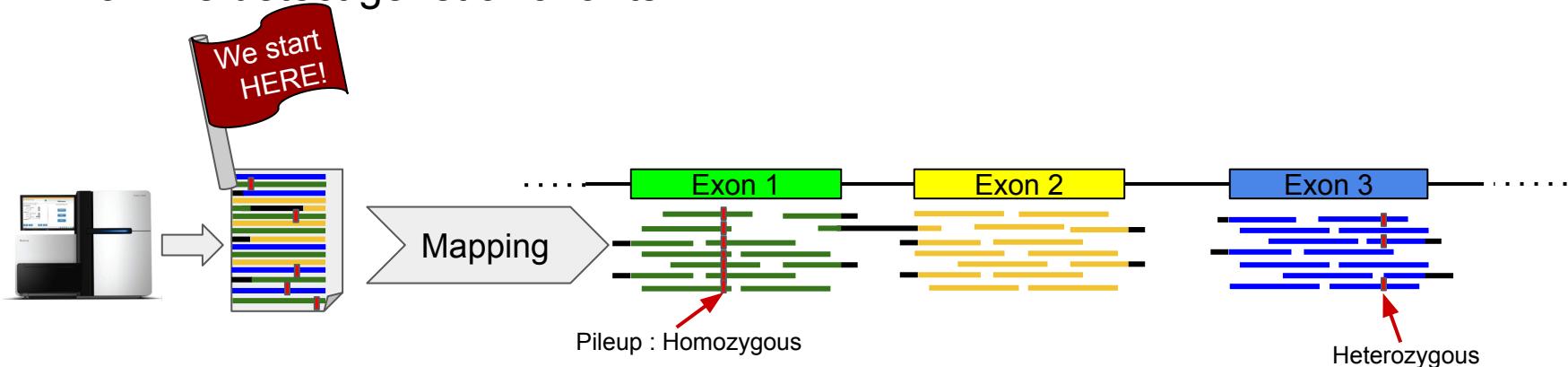
# DNA Sequencing data generation



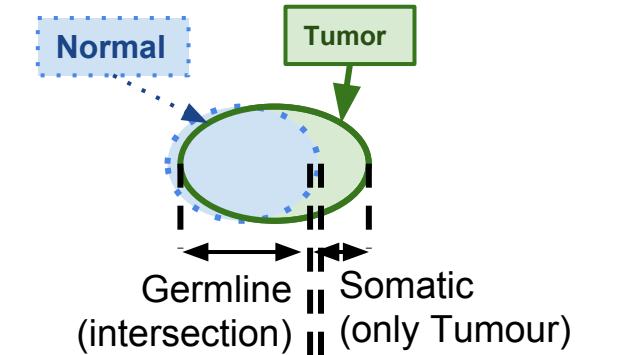
The sequence reads belong from the ends of the original fragment.

# Fundamentals of variant detection

- How we detect genetic variants ?

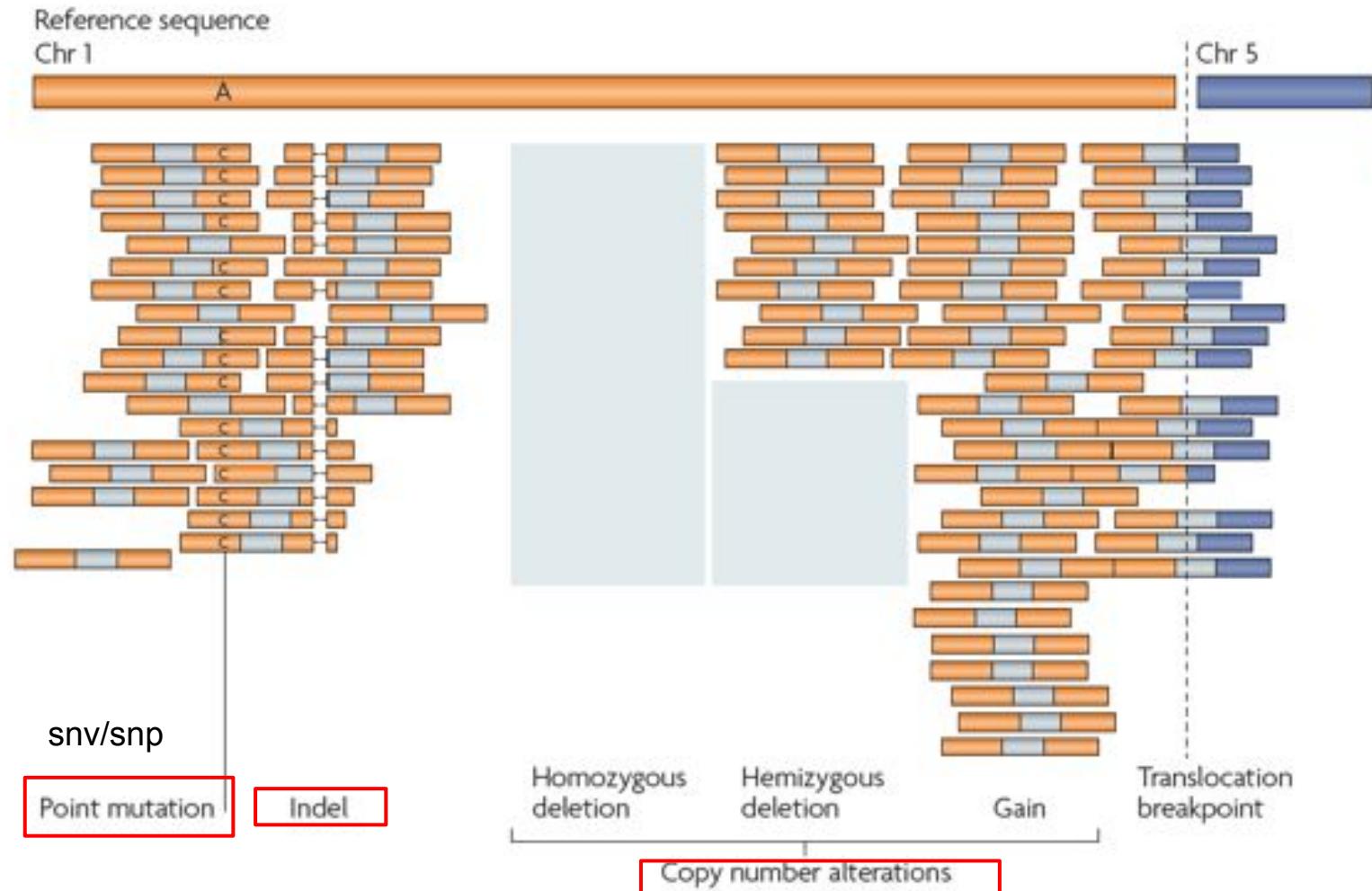


Discern somatic mutations by comparison:

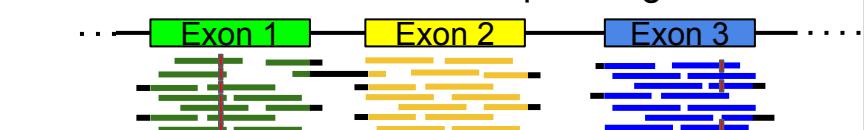
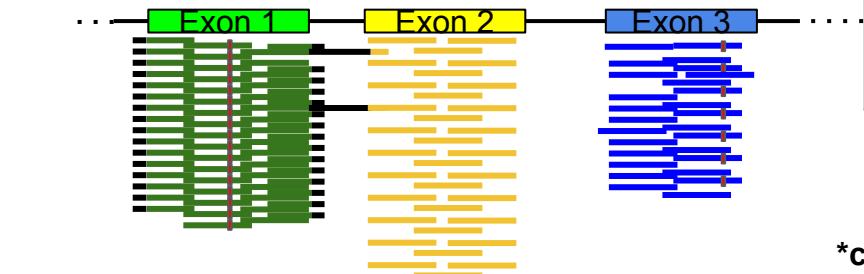


**Are these differences true calls?**  
Statistical method to estimate the most likely genotype.

# Different types of variants detected by mapping reads

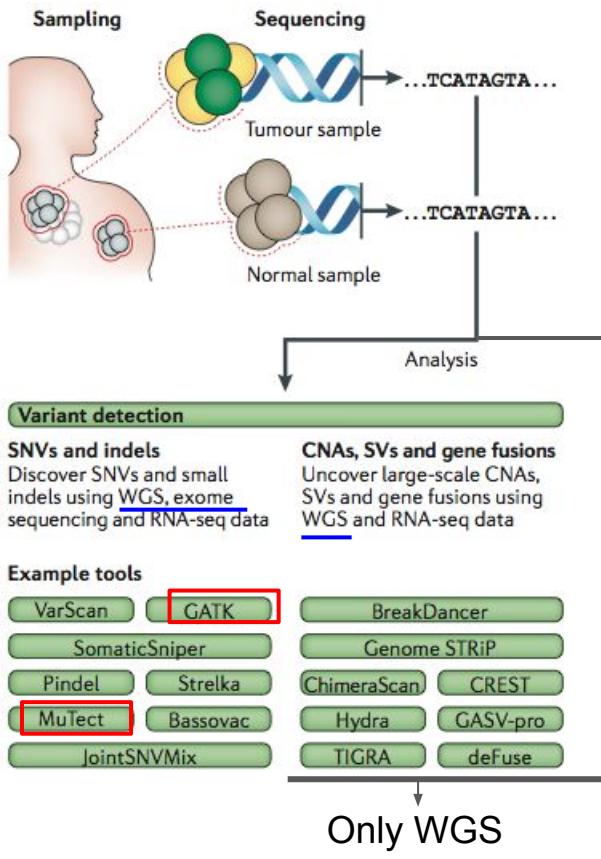


# Whole-Genome, and targeted resequencing

	# bp pos seq	Type of variants discovered	Avg Coverage per pos	Cost
<u>Whole-Genome Sequencing</u> 	~100 Gb	- <b>coding variants*</b> , intronic and regulatory sites. - <b>Structural variants</b> - <b>CNA</b> #Variants= 3M - 4M.	30x	High
<u>Whole-Exome Sequencing</u> 	~32Mb 50Mb	- <b>coding variants*</b> . - Some intronic and regulatory sites. - <b>CNA</b> (challenging). #Variants= 20k - 60k.	20x - 80x	Low
<u>Panel of genes by amplicon/PCR approach</u> 	ND	Depends on the design - Particular <b>coding variants*</b> - <b>CNA</b> (challenging) # variants = ND	1000x - 5000x	Low

\***coding variants**: missense, stop gained, stop lost, frameshift, splice region...

# Methods for Variant Detection



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

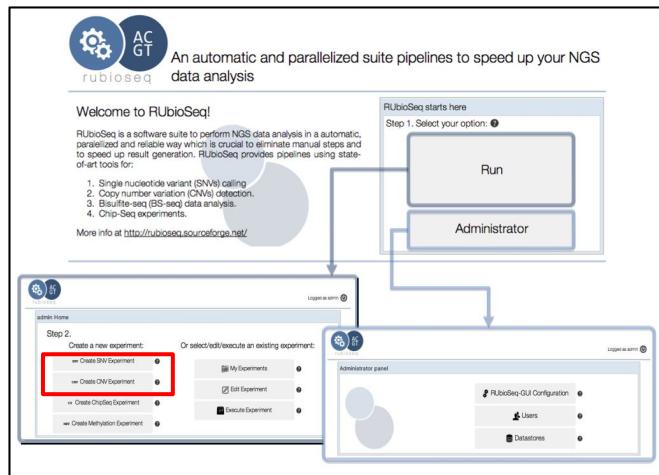
Appropriate methods for Whole-Exome seq

## 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)

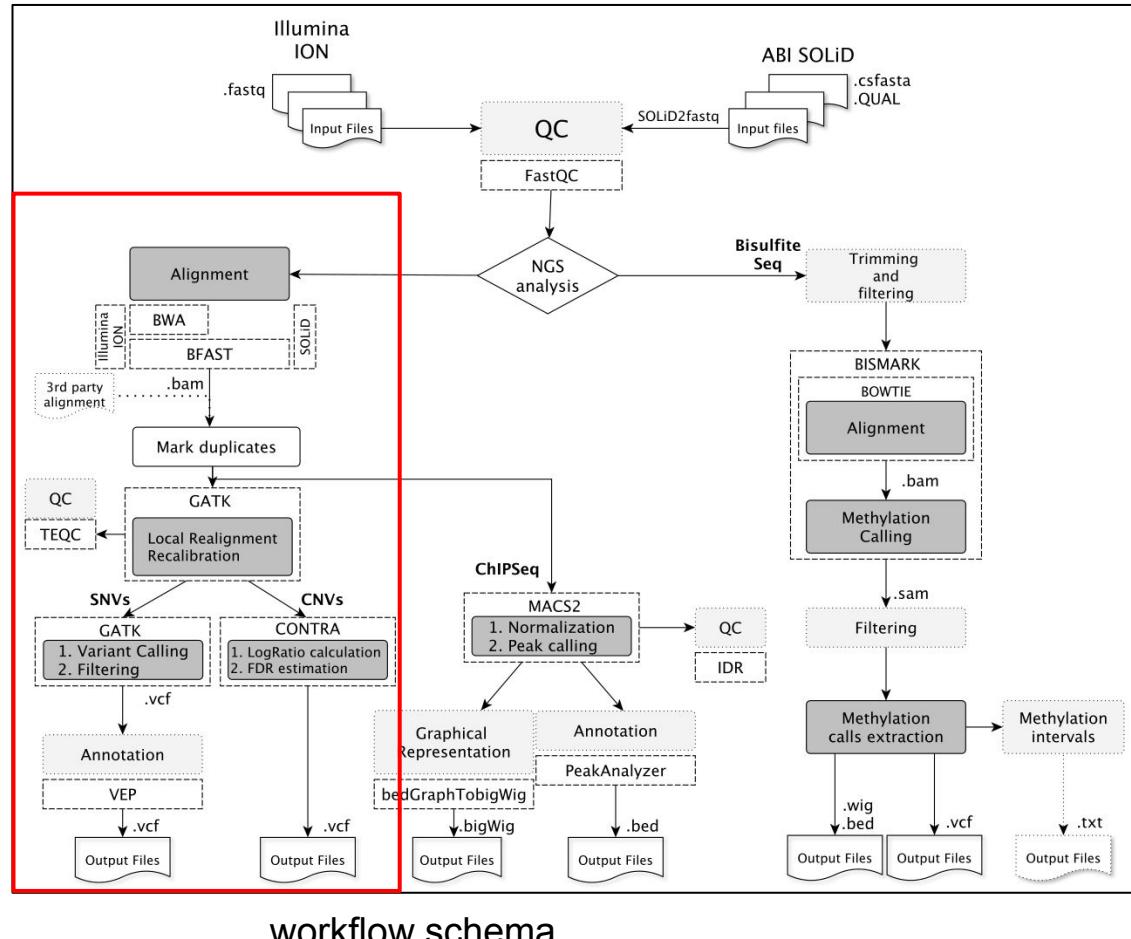
# Our proposed Variant Calling Pipeline



Graphic User Interface



<http://rubioseq.bioinfo.cnio.es/>



workflow schema

Developed by the **Bioinformatics Unit** at the **Spanish National Cancer Research Centre** (Madrid, Spain).

Rubio-Camarillo *et al.* Comput Methods Programs Biomed. 2017 Jan;138:73-81

Rubio-Camarillo *et al.* Bioinformatics (2013) 29 (13), 1687-1689

# 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 specie (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 the preference for analysis.
  - + **CRCh38/hg38** : current version (Sep. 2017). Released in 2014. More accurate, comprehensive (includes Haplotypes) and sophisticated.

*“CRCh38 is here now, but still waiting.”*

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

# Bundle of files for Variant Detection

---

1. **Genome Reference** (standard 1000 Genomes, fasta).
2. List of **Target beats or intervals** of genomic regions sequenced by the Library protocol.
3. **dbSNP** (VCF file) for a recent dbSNP release (build 138, it includes the 1000 Genomes).
4. HapMap genotypes and sites VCFs
5. **OMNI 2.5 genotypes for 1000 Genomes samples** (VCF).
6. The current best set of **known indels** to be used for local realignment); use both files:
  - 1000G\_phase1.indels.b37.vcf (currently from the 1000 Genomes Phase I indel calls)
  - Mills\_and\_1000G\_gold\_standard.indels.b37.sites.vcf

Q: How you can get this bundle of files?

A: you could get them from the **Broad Institute's FTP**

<ftp://gsapubftp-anonymous@ftp.broadinstitute.org/bundle/hg19/>

You also need the Intervals file from your NGS provider (Illumina, Ion Torrent,...)

For this workshop, we got these files for you.

More info.:

<https://software.broadinstitute.org/gatk/download/bundle>

# Bundle from the GATK's repository



ftp://ftp.broadinstitute.org/bundle/2.8/hg19/

## Index of /bundle/2.8/hg19/

Name	Size	Date Modified
[parent directory]		
1000G_omni2.5.hg19.sites.vcf.gz	49.4 MB	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.gz.md5	97 B	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.idx.gz	464 kB	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.idx.gz.md5	101 B	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.gz	42.9 MB	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.gz.md5	103 B	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.idx.gz	326 kB	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.idx.gz.md5	107 B	12/8/13, 1:00:00 AM
1000G_phase1.snp5.high_confidence.hg19.sites.vcf.gz	1.7 GB	12/8/13, 1:00:00 AM
1000G_phase1.snp5.high_confidence.hg19.sites.vcf.gz.md5	117 B	12/8/13, 1:00:00 AM
1000G_phase1.snp5.high_confidence.hg19.sites.vcf.idx.gz	3.4 MB	12/8/13, 1:00:00 AM
1000G_phase1.snp5.high_confidence.hg19.sites.vcf.idx.gz.md5	121 B	12/8/13, 1:00:00 AM
CEUTrio.HISeq.WGS.b37.bestPractices.hg19.vcf.gz	407 MB	12/8/13, 1:00:00 AM
CEUTrio.HISeq.WGS.b37.bestPractices.hg19.vcf.gz.md5	119 B	12/8/13, 1:00:00 AM
CEUTrio.HISeq.WGS.b37.bestPractices.hg19.vcf.idx.gz	3.2 MB	12/8/13, 1:00:00 AM
CEUTrio.HISeq.WGS.b37.bestPractices.hg19.vcf.idx.gz.md5	123 B	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.gz	19.1 MB	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.gz.md5	120 B	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.idx.gz	426 kB	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.idx.gz.md5	124 B	12/8/13, 1:00:00 AM
dbSNP_138.hg19_excluding_sites_after_129.vcf.gz	334 MB	12/8/13, 1:00:00 AM
dbSNP_138.hg19_excluding_sites_after_129.vcf.gz.md5	119 B	12/8/13, 1:00:00 AM
dbSNP_138.hg19_excluding_sites_after_129.vcf.idx.gz	3.6 MB	12/8/13, 1:00:00 AM
dbSNP_138.hg19_excluding_sites_after_129.vcf.idx.gz.md5	123 B	12/8/13, 1:00:00 AM
dbSNP_138.hg19.vcf.gz	1.4 GB	12/8/13, 1:00:00 AM
dbSNP_138.hg19.vcf.gz.md5	93 B	12/8/13, 1:00:00 AM
dbSNP_138.hg19.vcf.idx.gz	3.8 MB	12/8/13, 1:00:00 AM
dbSNP_138.hg19.vcf.idx.gz.md5	97 B	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.gz	58.0 MB	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.gz.md5	94 B	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.idx.gz	807 kB	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.idx.gz.md5	98 B	12/8/13, 1:00:00 AM

**Everything** you need for the variant calling is stored in the FTP of the GATK team.

It includes the **human genome reference** too!

There are independent bundles for each assembly of the human genome (hg19, hg38,etc). The difference across them is the coordinate system for the annotations, but it is the same data.

You can download the different bundles from GATK's FTP (Broad Institute) visiting this URL with your Internet Browser:

<ftp://gsapubftp-anonymous@ftp.broadinstitute.org/bundle/2.8/>



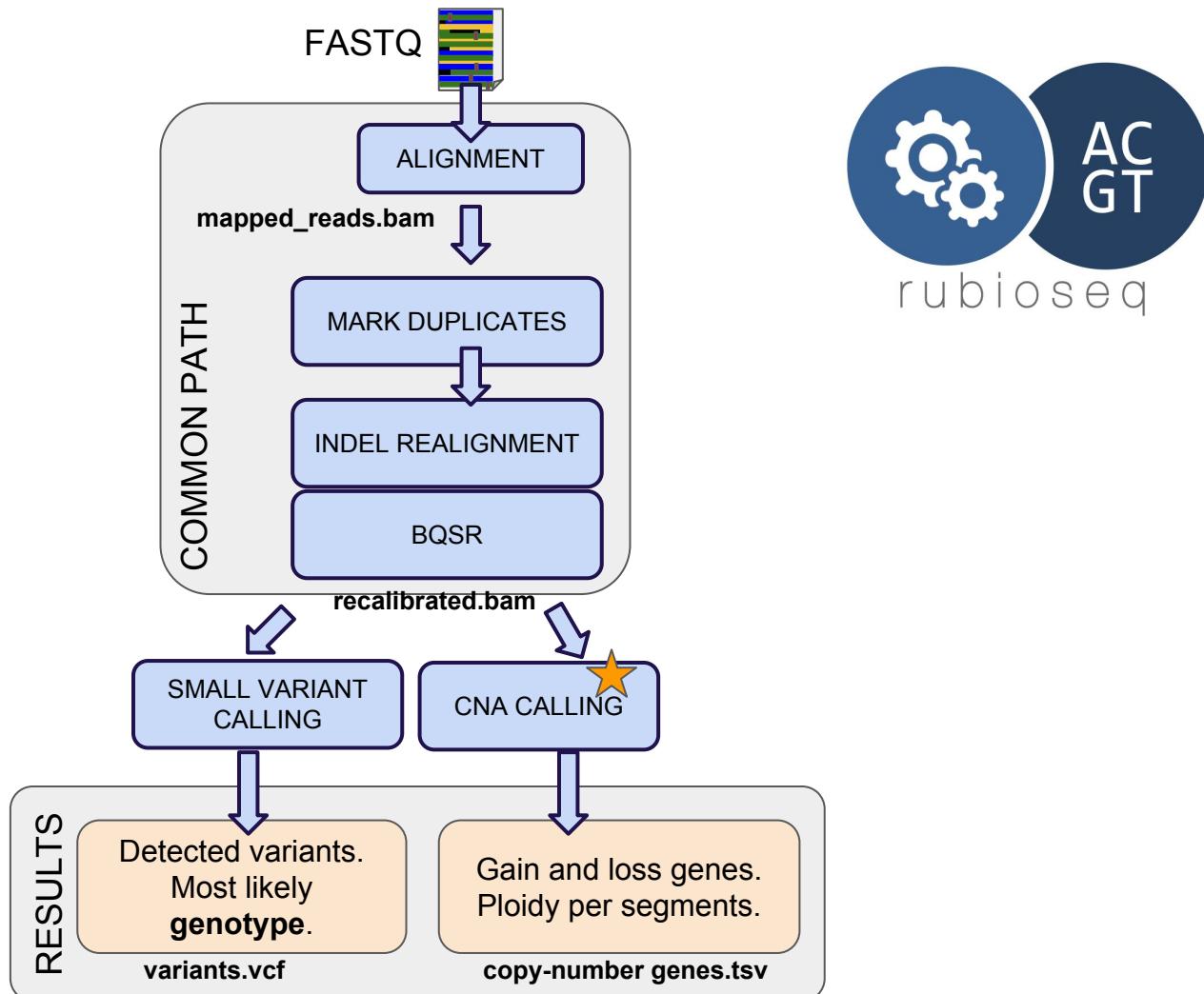
Tip for home: The following UNIX command downloads the whole bundle for **hg19** in one step (~hrs) :

```
$ wget -r -nH --cut-dirs=2 --reject-regex "NA12878|CEUTrio" \
-P /path/to/your_directory/ ftp://gsapubftp-anonymous@ftp.broadinstitute.org/bundle/hg19/*
```

More info.:

<http://gatkforums.broadinstitute.org/discussion/1213/whats-in-the-resource-bundle-and-how-can-i-get-it>

# Point mutations and CNV Calling Process

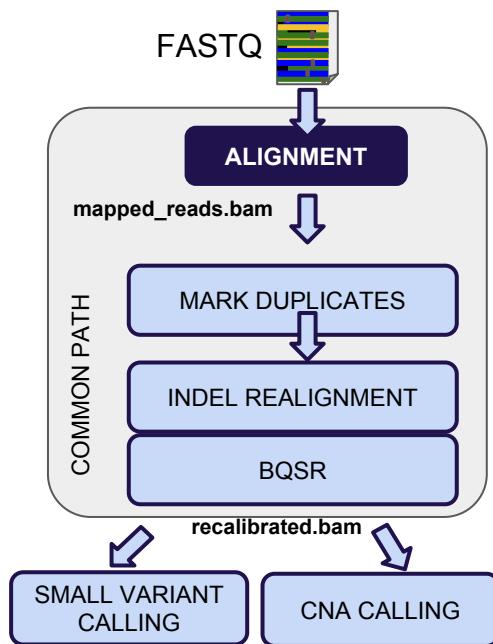


★ It requires **tumour and matched normal sample (or a panel of normals)**

# Variant Calling Pipeline

## 1. Alignment

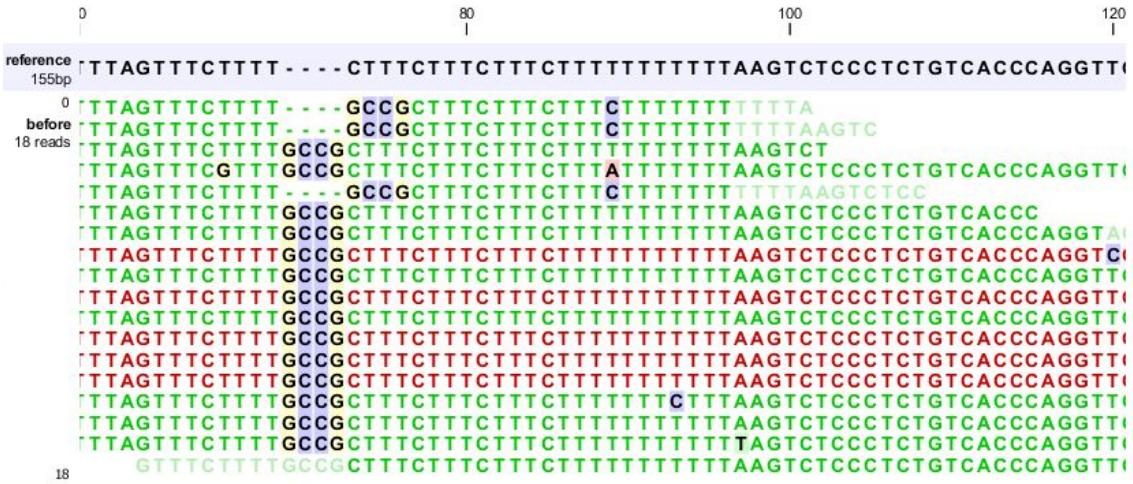
### WORKFLOW:



### METHOD: by BWA & Bfast+BWA

<https://github.com/lh3/bwa#citing-bwa>

<http://sourceforge.net/projects/bfast/files/bfast%2Bbwa/>

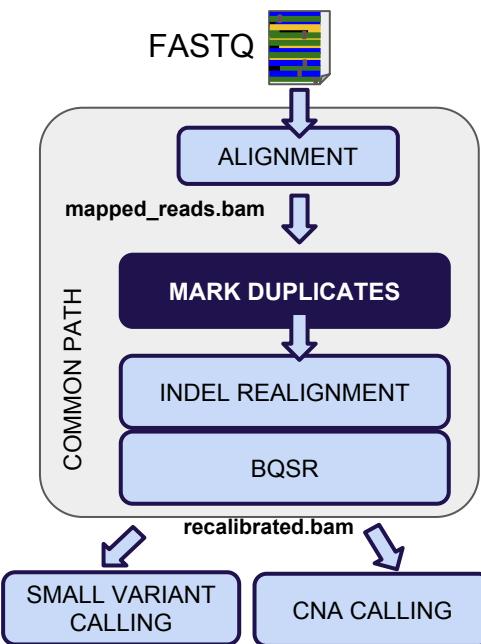


- Fast mapping on the reference genome by creating indexes. It is computationally intensive, but it is done only once.
- Search for candidate sites to align a given read by using seeds (fragments of a read).

# Variant Calling Pipeline

## 2. Mark duplicates

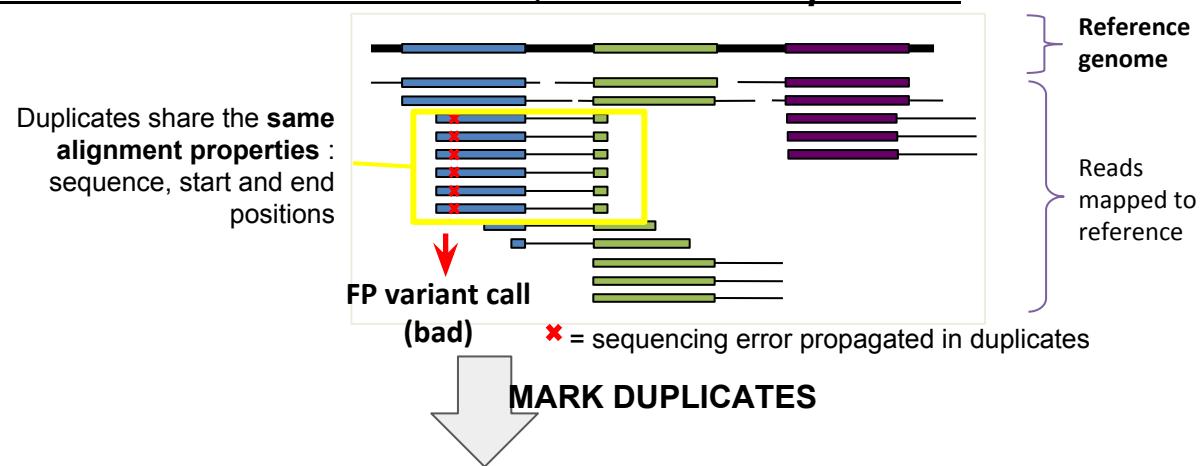
### WORKFLOW:



### Under the hood:

- 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.
- Therefore, duplicates are **worthless** for the analysis:  
*Duplicates are source of False Positives calls while only provide redundancy.*

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



**METHOD:** by Picard-tools

<http://broadinstitute.github.io/picard/>

(alternatives : samtools)

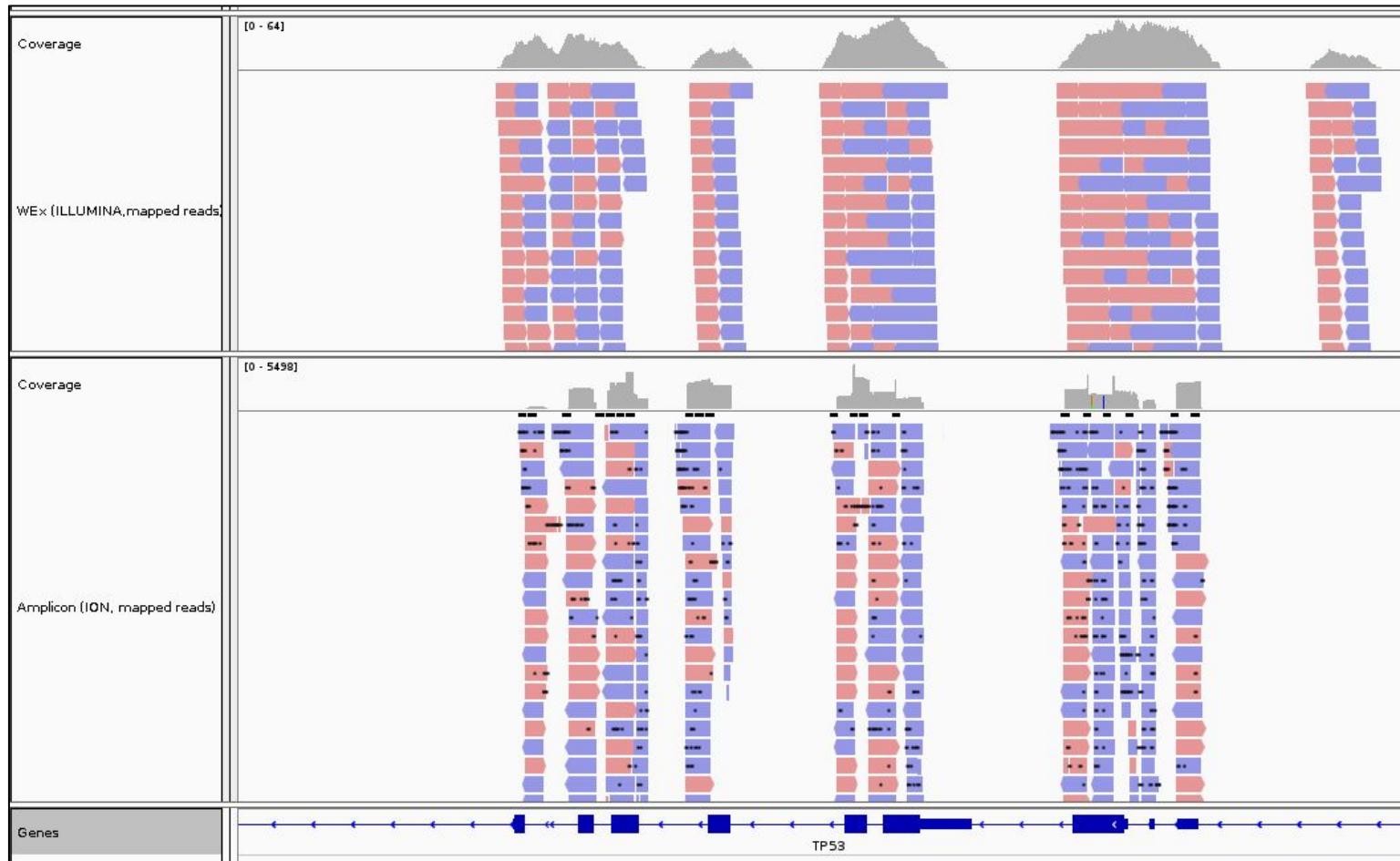
Adapted from GATK

After marking duplicates, the variant caller will only see :



## Variant Calling Pipeline

# 2. Mark duplicates: WEx Vs. Amplicon



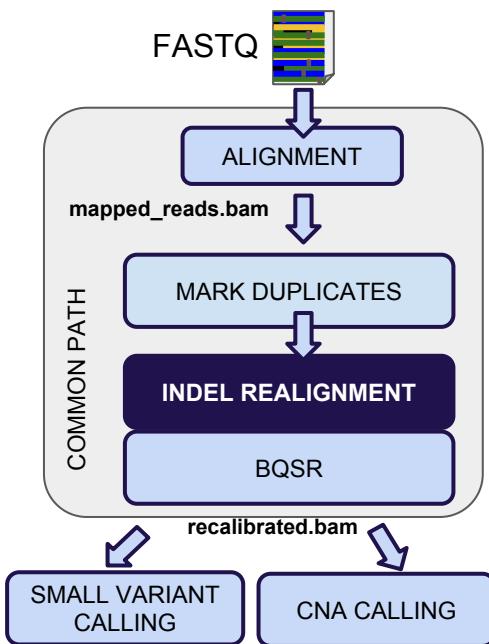
**WARNING:** Do NOT mark duplicates in data derived from amplicon techniques (**Ion Torrent**).

More info.: <http://gatkforums.broadinstitute.org/discussion/5847/remove-duplicates-from-targetted-sequencing-using-amplicon-approach>

## Variant Calling Pipeline

### 3. Indel realignment

#### WORKFLOW:



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

*During alignment, penalties on mismatches are much cheaper than gaps (indels). Aligners will tend to choose Mismatches at the beginning, and locate indels in the rest.*

Also, there are sometimes multiple solution (alignments) for a given read. Aligners choose one randomly.

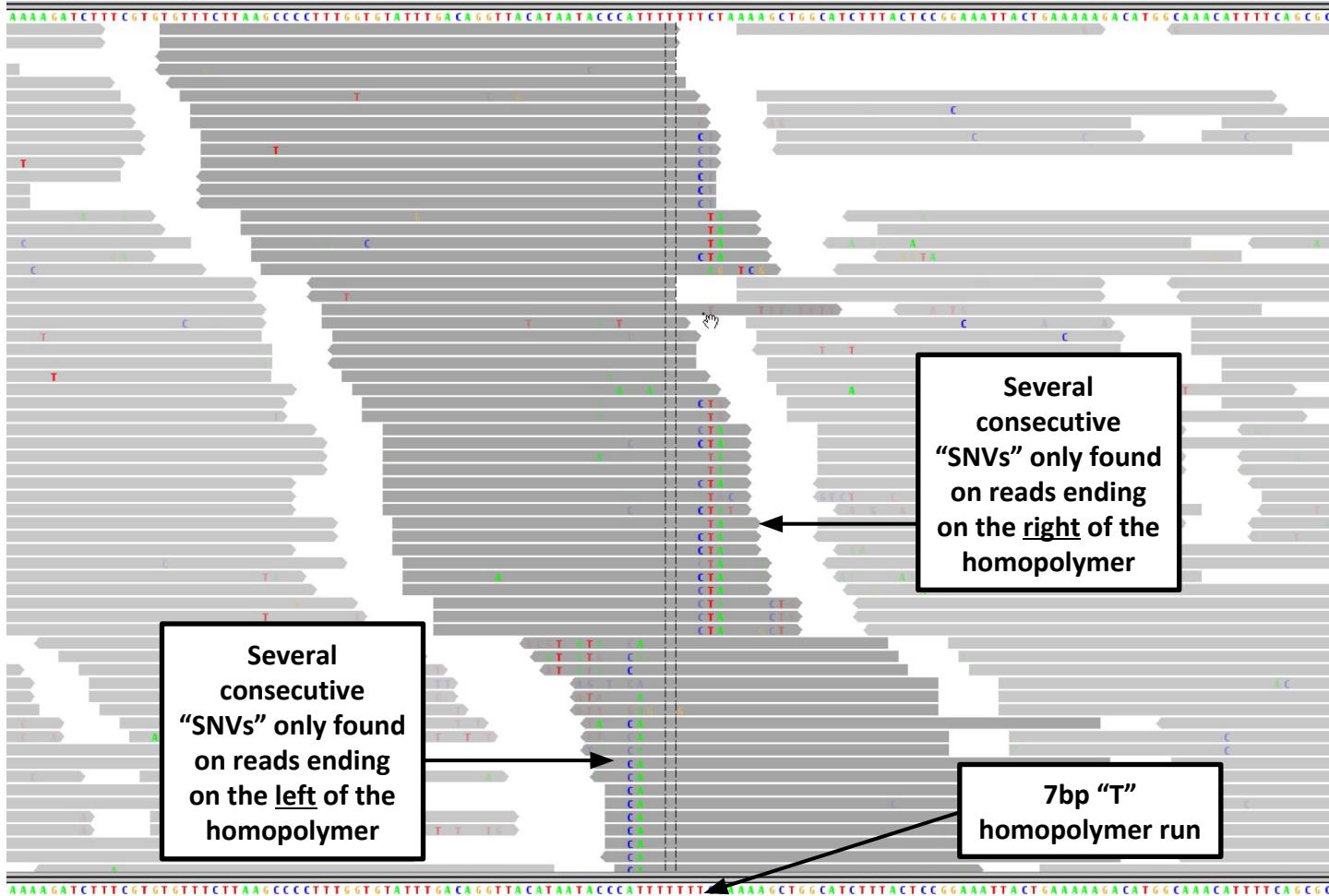
Variant calling requires the most perfect alignment as possible to avoid False Positives.

#### METHOD: by GATK

[https://www.broadinstitute.org/gatk/gatkdocs/org\\_broadinstitute\\_gatk\\_tools\\_walkers\\_indels\\_IndelRealigner.php](https://www.broadinstitute.org/gatk/gatkdocs/org_broadinstitute_gatk_tools_walkers_indels_IndelRealigner.php)

## Variant Calling Pipeline

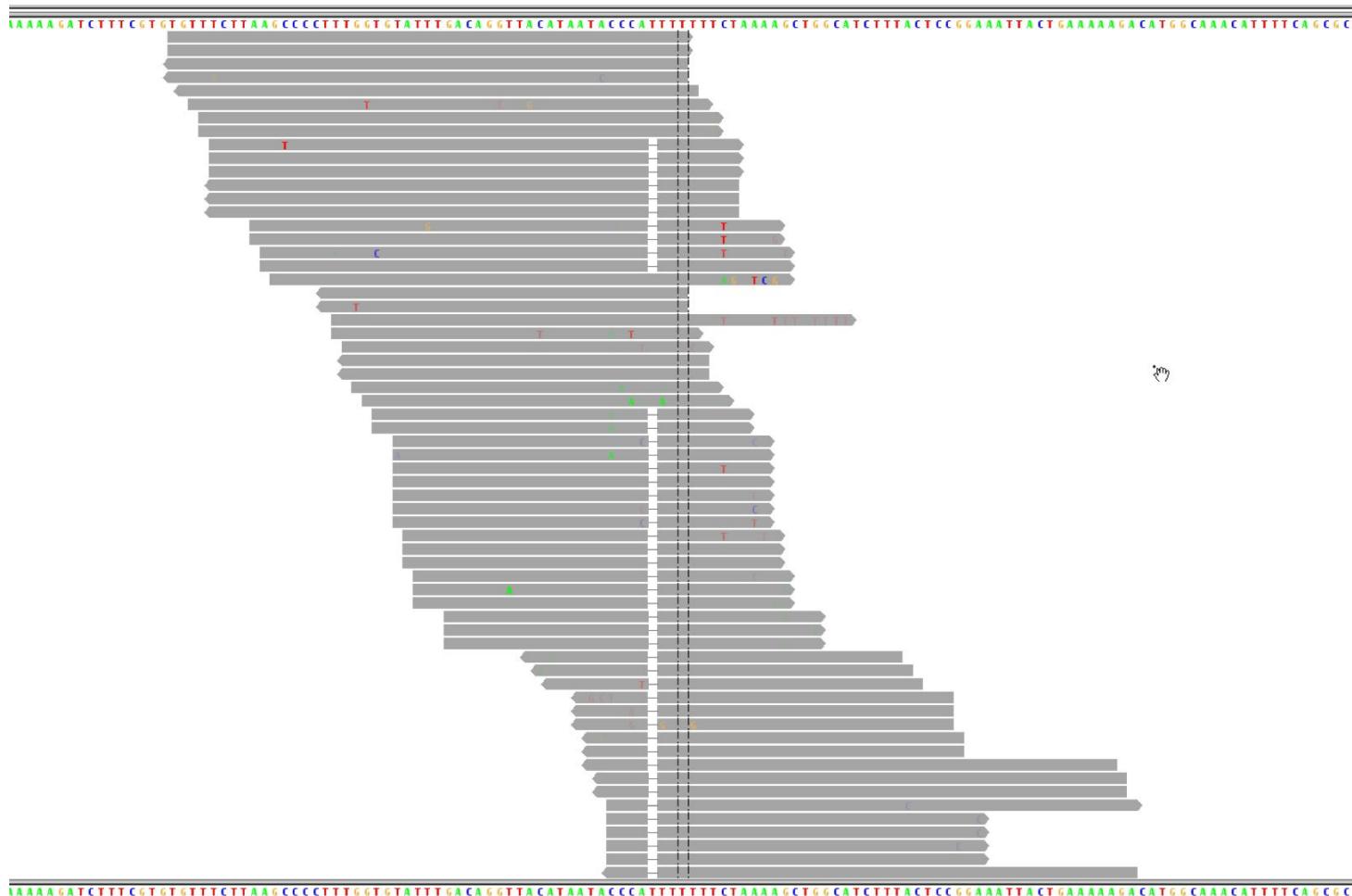
### 3. Indel realignment



Taken from GATK team

## Variant Calling Pipeline

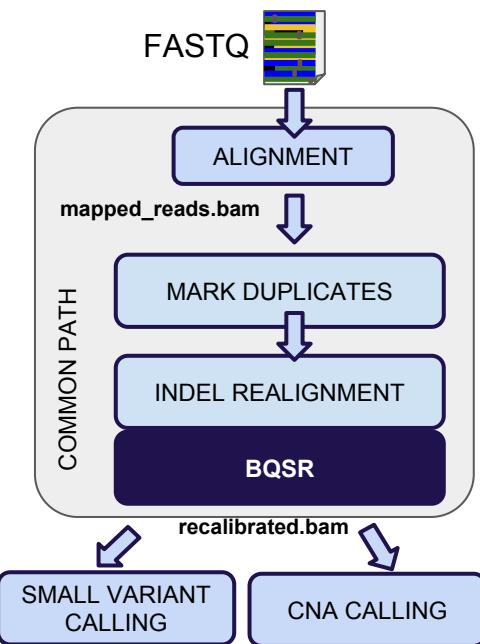
### 3. Indel realignment



Taken from GATK team

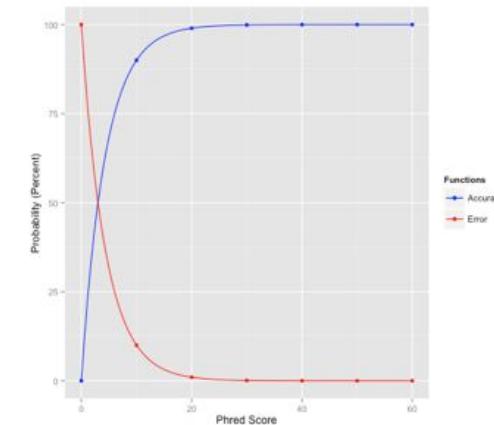
## 4. Base Quality Score Recalibration

### WORKFLOW:

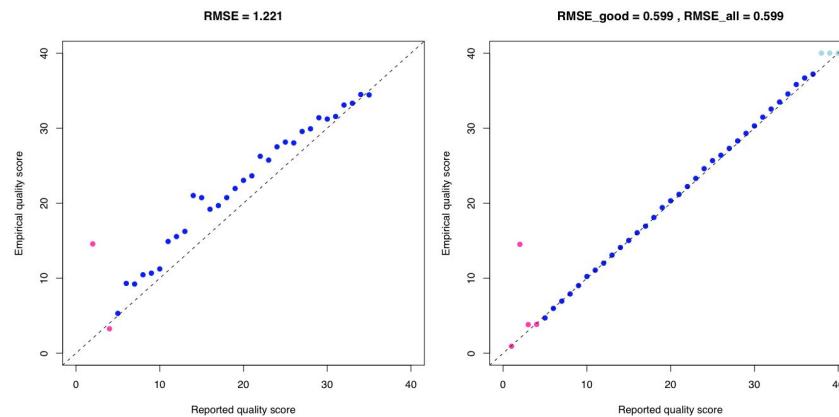


Phred Quality : each position of the sequence has its particular **base Quality score**.

The individual quality measures are NOT very important during the alignment step (mapping), but crucial during Variant calling.



Different NGS technologies have their particular bias in QS depending on the context. **They could correct empirically** these biases.



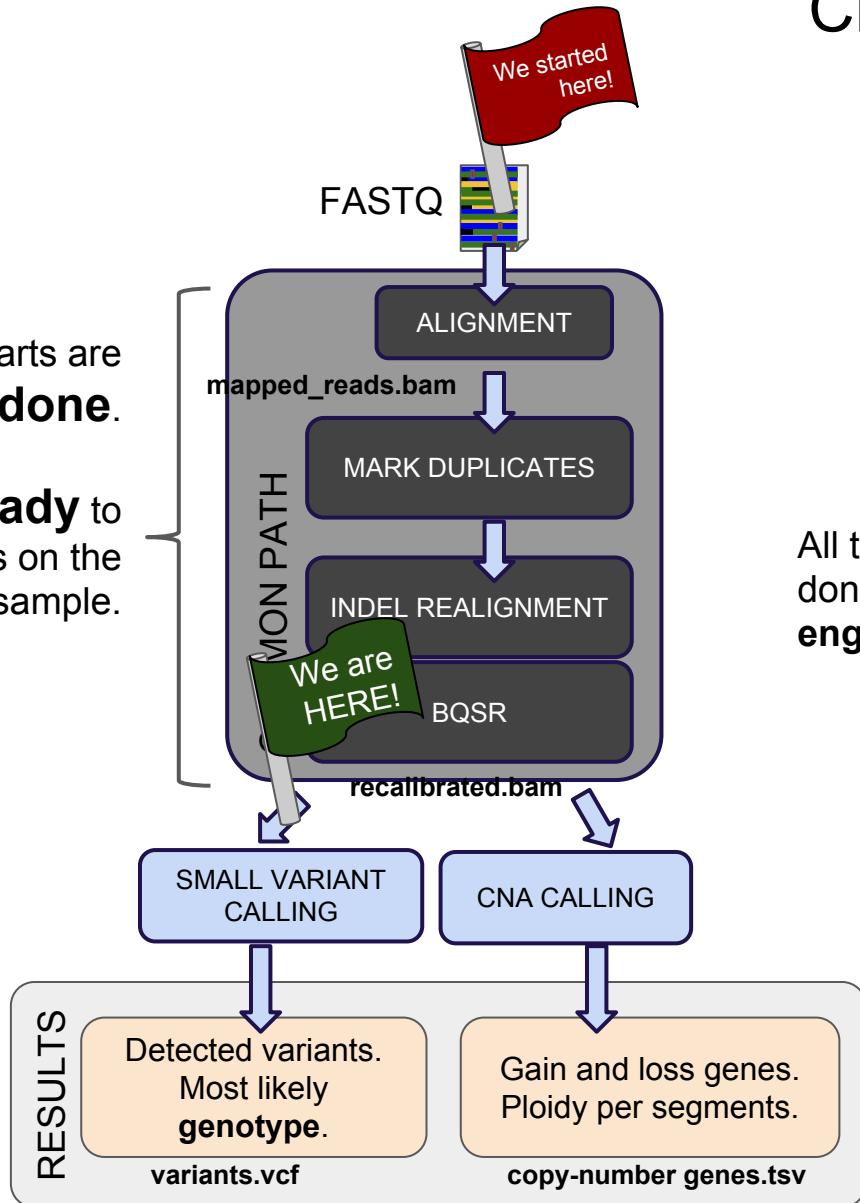
### METHOD: by GATK

<http://gatkforums.broadinstitute.org/discussion/44/base-quality-score-recalibration-bqsr>

# Point mutations and CNV Calling Process



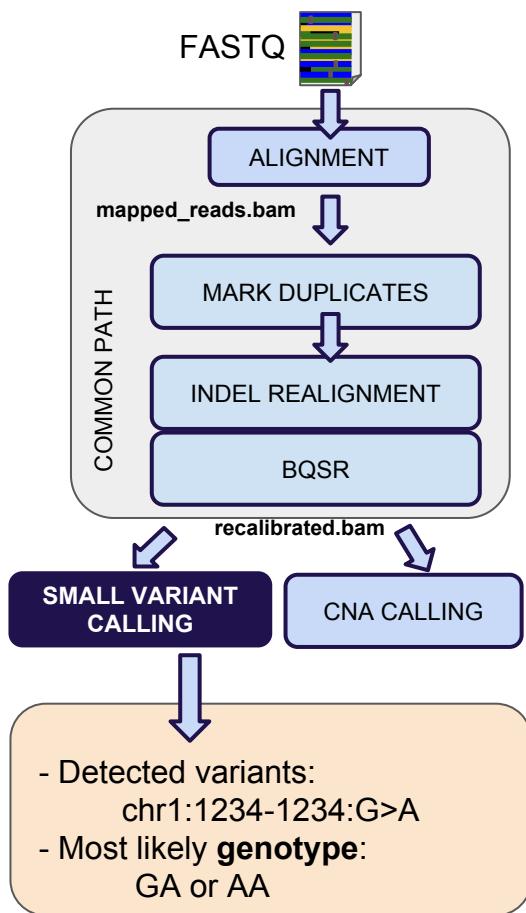
The first two parts are **done**.  
We are **ready** to discover variants on the sample.



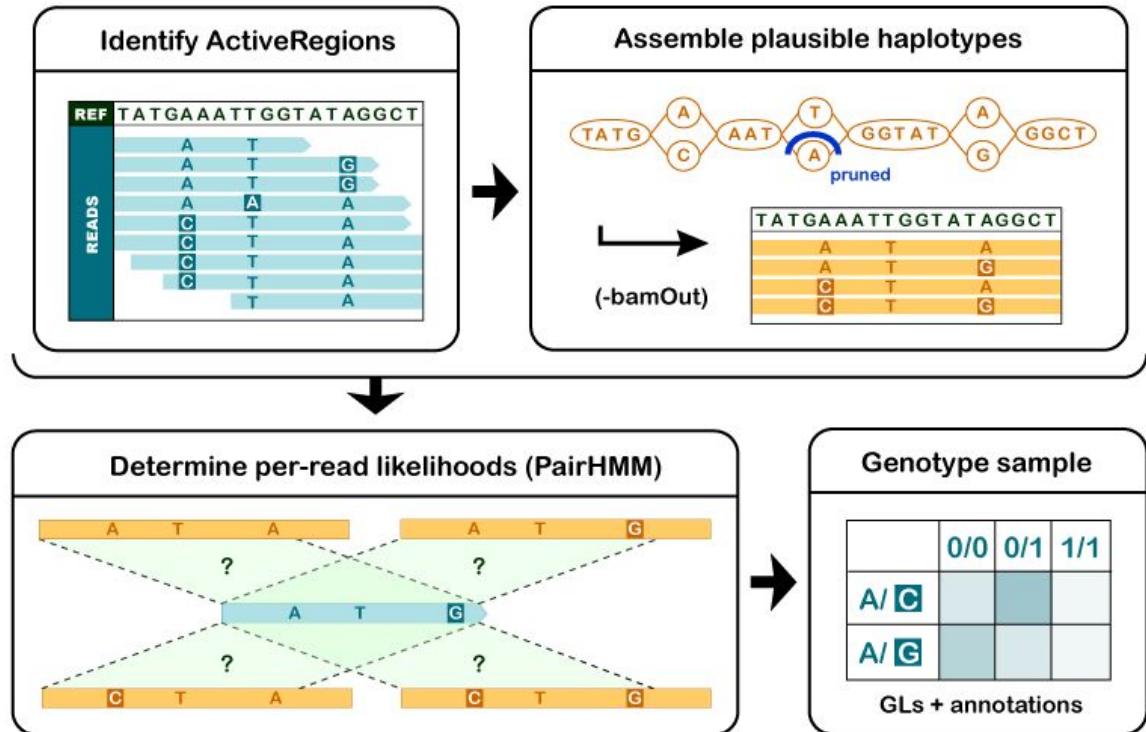
All these steps are **automatically** done by the pipeline (**RUBioSeq's engine**).

# 5. GATK Variant Calling Process : SNV & Indels

## WORKFLOW:



**Haplotype Caller** (new, included in **RUBioSeq v3.8.1**) : Variant calling based on the calculation of genotype likelihoods:



**Assumptions:** Diploid genome (2n).  
**Limitation:** Allele freq > 20%.

Further reading:

<http://gatkforums.broadinstitute.org/discussion/4148/hc-overview-how-the-haplotypecaller-works>

HC steps 1-4: <https://software.broadinstitute.org/gatk/documentation/topic?name=methods>

# GATK is in active development

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

gatk Home About Guide Blog Forum Download Search

The Genome Analysis Toolkit or GATK is a software package for analysis of high-throughput sequencing data, developed by the [Data Science and Data Engineering](#) group at the [Broad Institute](#). The toolkit offers a wide variety of tools, with a primary focus on variant discovery and genotyping as well as strong emphasis on data quality assurance. Its robust architecture, powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

[Learn more »](#)



## About

Overview of the tools and development team



## Guide

Detailed documentation, guidelines and tutorials



## Forum

Ask here for help with questions and bug reports



## Blog

Announcements and progress updates

**Latest version: 3.4-46**

[Release Notes](#)

[Download now ↗](#)

[For-profit users: click here](#)

# GATK is in active development

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

The screenshot shows the GATK website with several highlighted sections:

- GATK Best Practices**: A detailed page about workflow recommendations for variant discovery analysis.
- Support Forum**: A red-bordered section containing the "Ask the GATK team" and "Gziped gVCF files" threads.
- Forum**: A red-bordered section containing the "About variant calling of single sample .." thread.
- Blog**: A red-bordered section containing the "How can I use command to genotype refinement of population priors(no family groups)" thread.

**Best Practices Workflow Diagram:**

```
graph LR; A[Sequencing] --> B[FASTQ]; B --> C[GATK Best Practices]; C --> D[VCF]
```

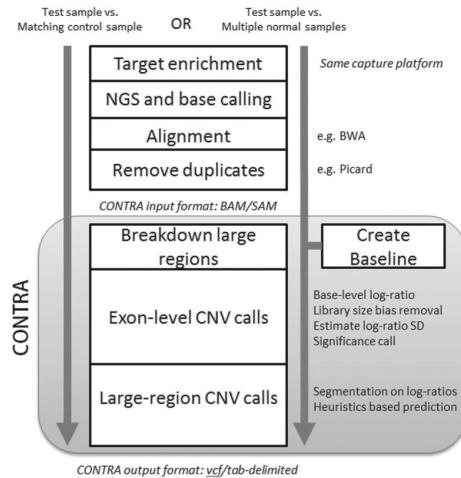
**Support Forum Threads:**

- Ask the GATK team**: Errors, bugs, problems and usage questions for the developers of the GATK or the community at large.
- Gziped gVCF files**: Tools built on top of GATK.
- About variant calling of single sample ..**: Questions about variant calling.
- GATK HaplotypeCaller run**: Questions about the HaplotypeCaller tool.
- How can I use command to genotype refinement of population priors(no family groups)**: Questions about command-line tools for genotype refinement.

**Bottom Navigation:**

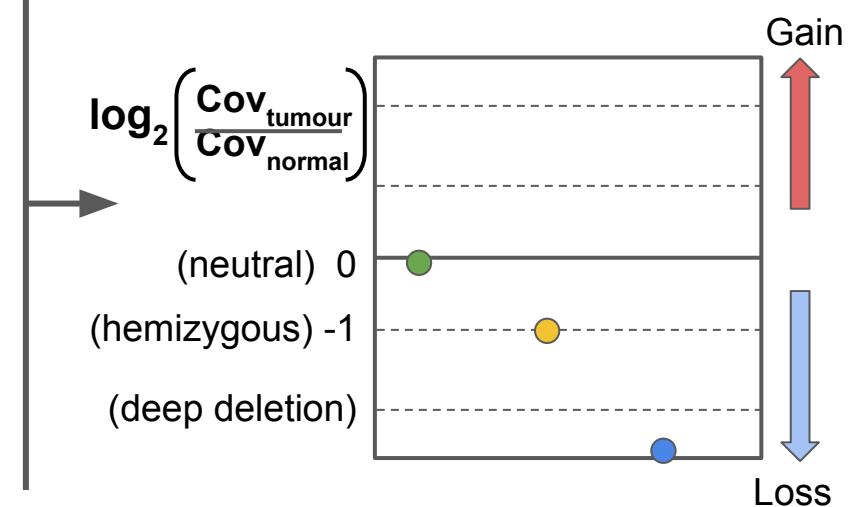
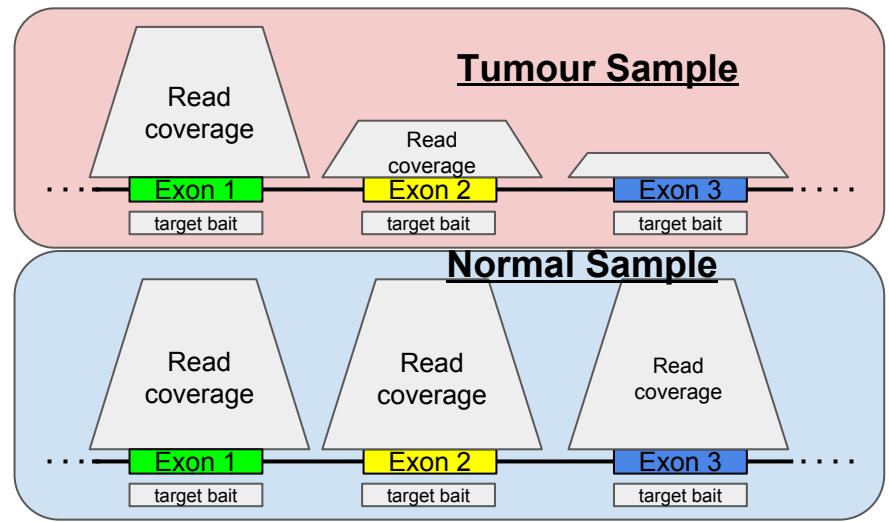
PM17 @GTPB Day #1 - NGS I : Variant Detection

# 6. CNA Variant Calling



- **Normalization:** Split large regions. GC-content bias, unbalanced library size effect on log-ratios.
- Read-depth coverage & log<sub>2</sub> CN ratio are corrected.
- Significance:  
**Assumption:** log<sub>2</sub>-transformed coverage fits a normal distribution:

$RLR \sim N(\mu_d, \sigma_d)$  ; Two-tailed P-value.  
multiple testing correction (FDR).



Li J et al. CONTRA: copy number analysis for targeted resequencing. (2012) Bioinformatics

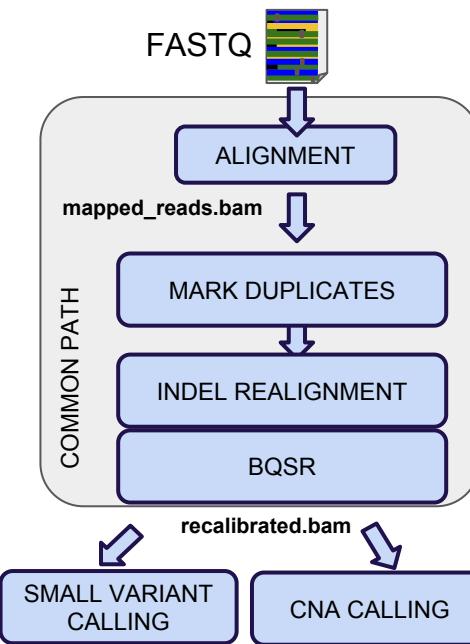
# ● What have we learnt?

## Main concepts:

- Variants.
- How to detect them.
- Differences between platforms.

## Requirements:

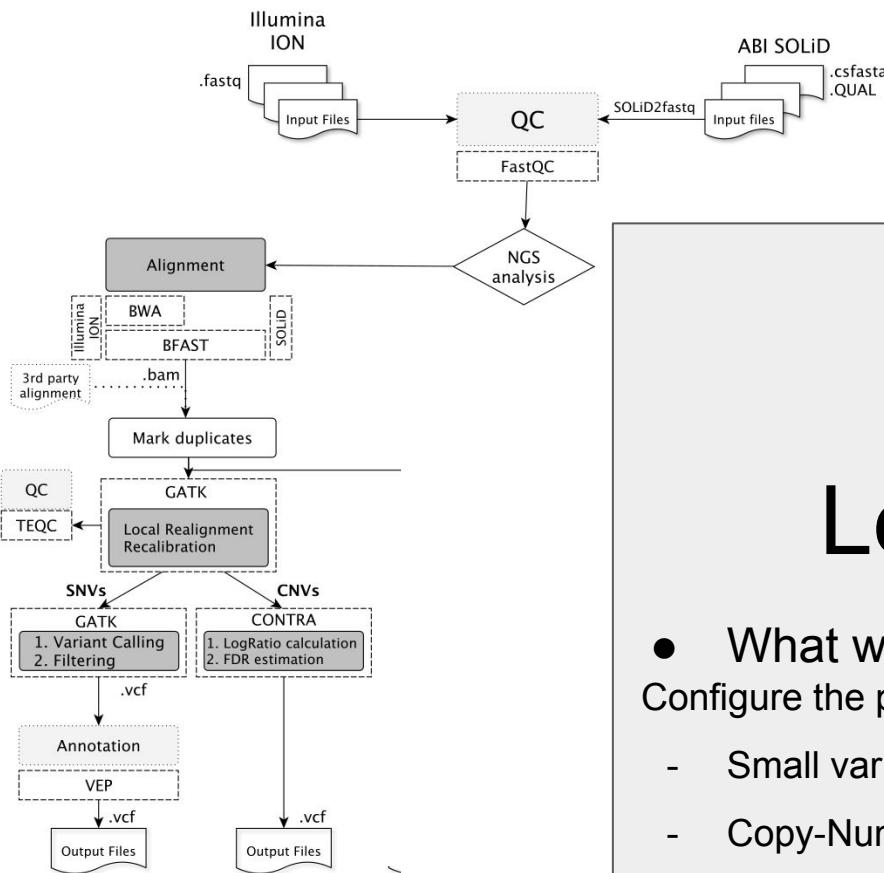
- Files.
- Methods



# ● Questions?

Thanks to Miriam Rubio-Camarillo & Gonzalo Gómez for the support and development of RUBioSeq.

# First hands-on :: Configure RUbioSeq+ Configuration files



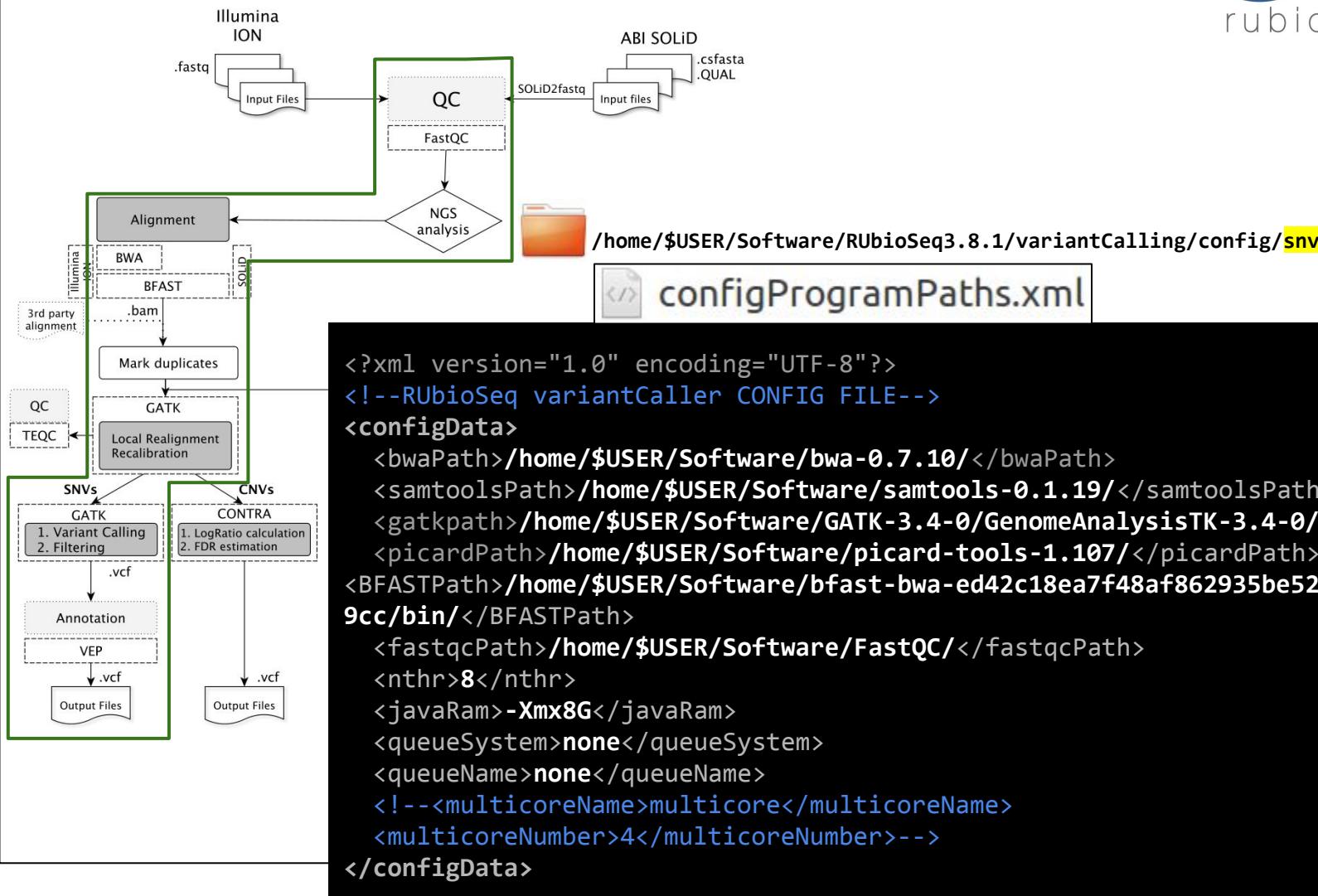
It is already installed.  
Let's configure it!

- What will we do?

Configure the pipeline for the detection of:

- Small variants (SNV + indels)
- Copy-Number Variants

# First hands-on :: Configure RUbioSeq+ Configuration for Small variant analysis



# First hands-on :: Configure RUBioSeq+ Configuration for CNV analysis

