


**FACULTAD DE MEDICINA**  
UNIVERSIDAD DE CHILE


**GENOMED-Lab**  
<http://genomed.med.uchile.cl>



**ICBM**  
INSTITUTO DE CIENCIAS BIOMÉDICAS

## Análisis de datos NGS: Llamado de variantes

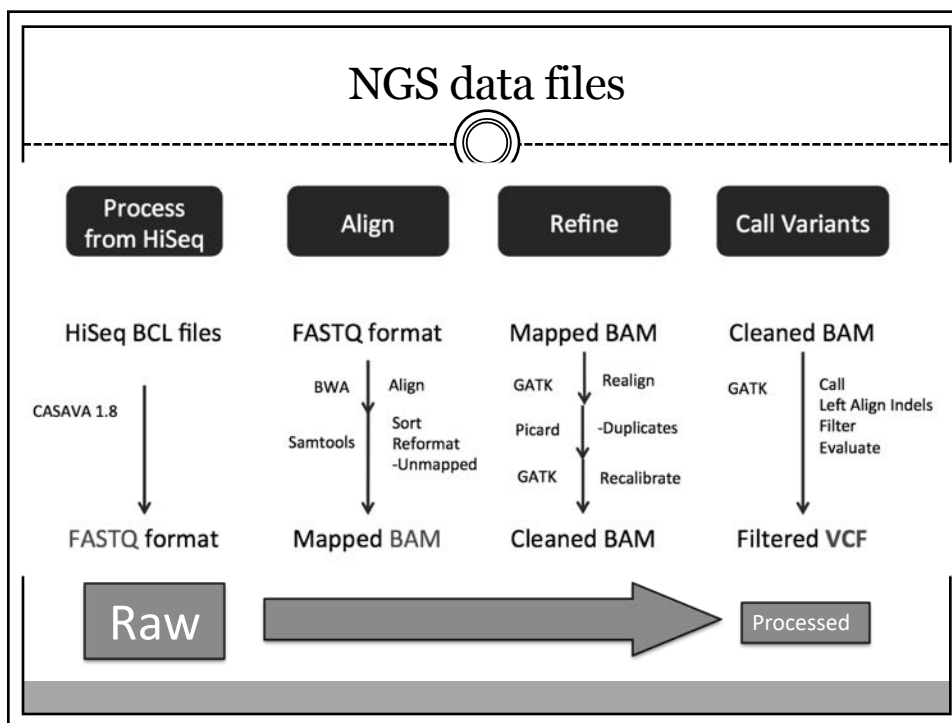
---



RICARDO A. VERDUGO, Ph.D.

Programa de Genética Humana, ICBM  
Facultad de Medicina, U. de Chile

Abril 2019



--

\_\_\_\_\_

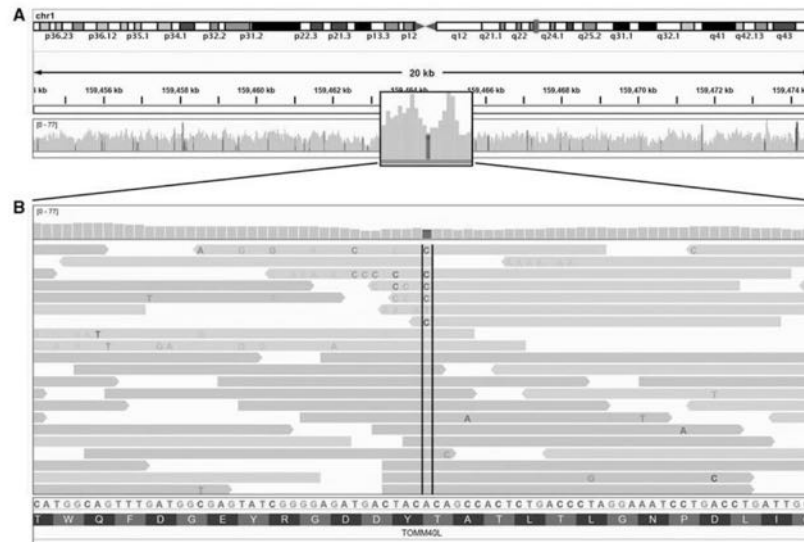
## Archivo SAM: registros



**HEADER** containing metadata (sequence dictionary, read group definitions etc)  
**RECORDS** containing structured read information (1 line per read record)

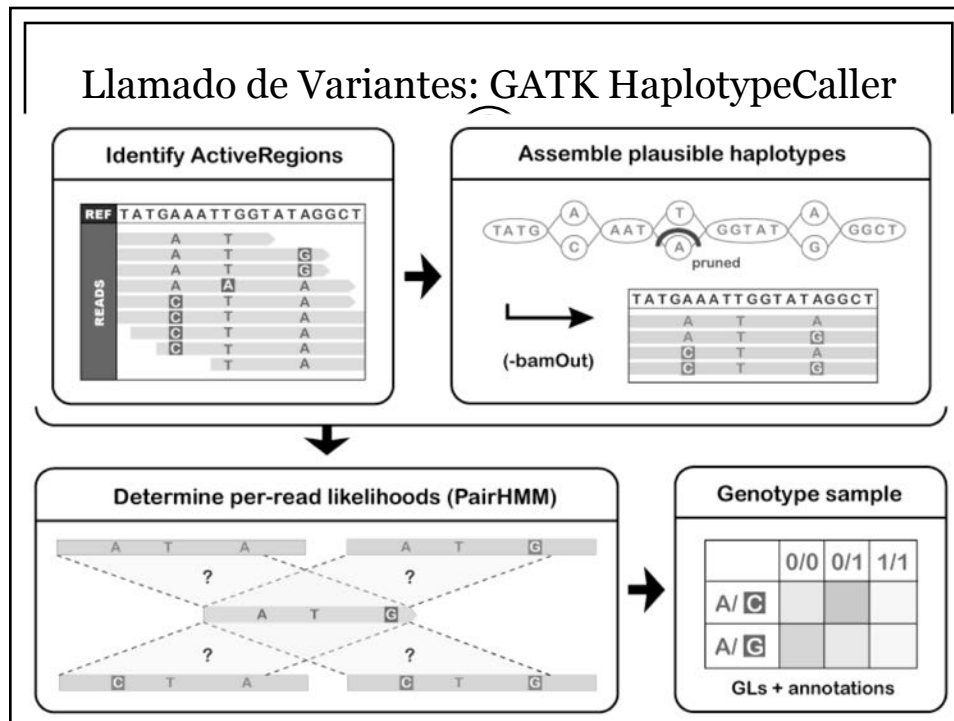
read name	position	CIGAR	read sequence	metadata
SLX1:1:127:63:4	99	1	10052169	60
		23M6N10M	=	14
			10	GAAGATACTGGT
			768832	'48:::
				SM:Z:JPTGBMN01 ...
	flags	MAPQ	mate information	quality scores

## Integrative Genomic Viewer

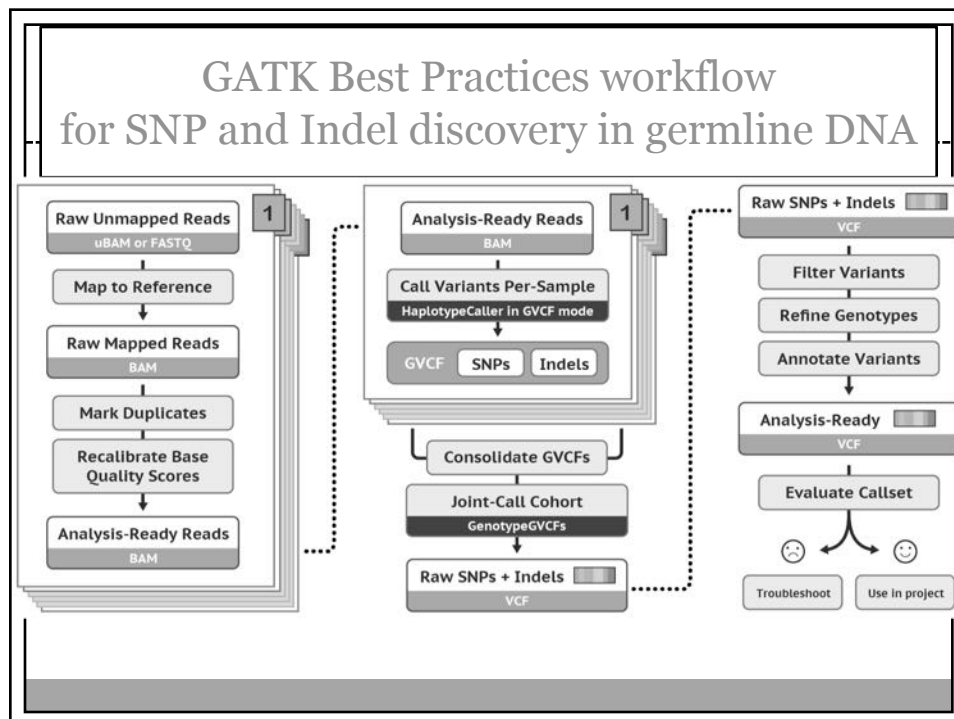


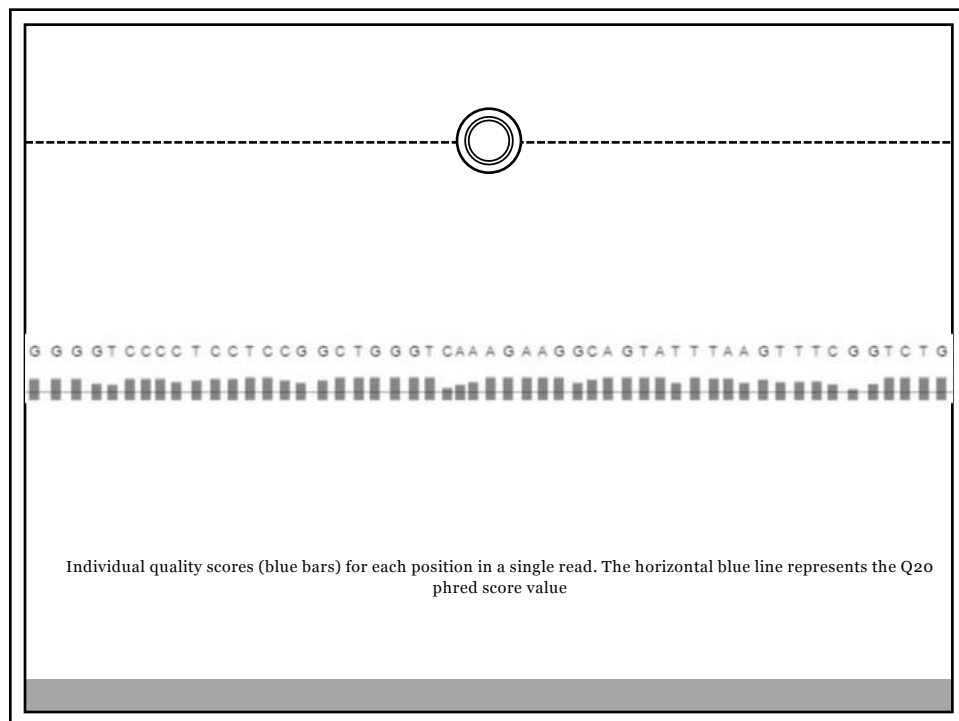
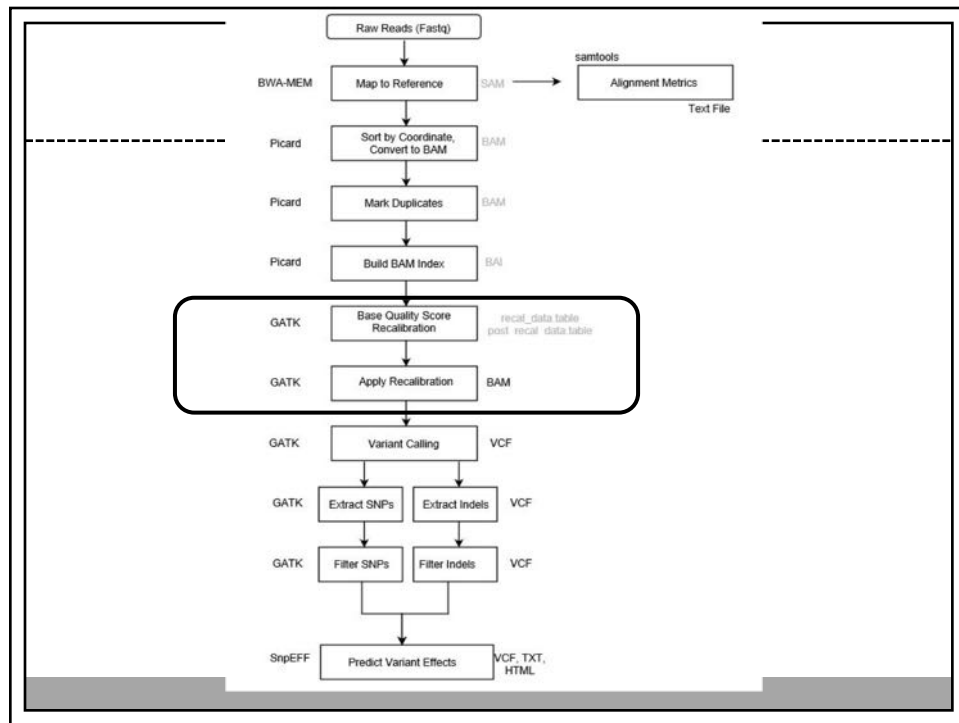
<http://software.broadinstitute.org/software/igv/>

## Llamado de Variantes: GATK HaplotypeCaller



## GATK Best Practices workflow for SNP and Indel discovery in germline DNA





## Why do we care about quality scores so much?



- Variant calling algorithms rely on the quality score assigned to the individual base calls
- Tells us how much we can trust that particular observation to inform us about the biological truth of the site
- If we have a basecall that has a low quality score, that means we're not sure we actually read that A correctly, and it could actually be something else
- So we won't trust it as much as other base calls that have higher qualities
- We use that score to weigh the evidence that we have for or against a variant existing at a particular site

## Why Recalibrate?



- Scores produced by the machines are subject to various sources of systematic technical error
- Leads to over- or under-estimated base quality scores in the data.
- Errors can arise due to the physics or the chemistry of how the sequencing reaction works, possibly manufacturing flaws in the equipment.

## Why Recalibrate?

Base quality score recalibration (BQSR) is a process in which we apply machine learning to model these errors empirically and adjust the quality scores accordingly.

Raw, high-sensitivity callsets contain many false positives

- Mutation calling algorithms are very permissive by design
  - How to filter?
    - Hand-tuned hard-filtering requires time and expertise
    - Better to learn what the filters should be from the data itself
  - Must enable analysts to trade off sensitivity and specificity depending on project goals
- ☒ **Building a model of what true genetic variation looks like will allow us to rank-order variants based on their likelihood of being real**

## How does BQSR work?



1. You provide GATK Base Recalibrator with a set of known variants.
2. GATK Base Recalibrator analyzes all reads looking for mismatches between the read and reference, skipping those positions which are included in the set of known variants (from step 1).
3. GATK Base Recalibrator computes statistics on the mismatches (identified in step 2) based on the reported quality score, the position in the read, the sequencing context (ex: preceding and current nucleotide).
4. Based on the statistics computed in step 3, an empirical quality score is assigned to each mismatch, overwriting the original reported quality score.

## From annotations to mixture models

- Each variant has a diverse set of statistics associated with it.
- These annotations tend to form Gaussian clusters
- We can fit a “Gaussian mixture model” to the annotations known variants in our dataset.
- Any new variant can be scored by evaluating the associated annotations in this model.



Variant annotations are the “features” of the model

### VCF record for an A/G SNP at 22:49582364

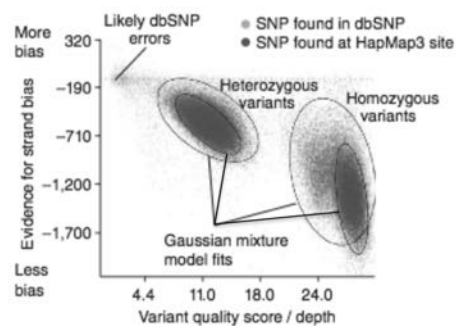
22	49582364	.	A	G	198.96	.
AC=3; AF=0.50; AN=6; DP=87; MLEAC=3; MLEAF=0.50; MQ=71.31; MQ0=22; QD=2.29; SB=-31.76 GT:DP:GQ						
INFO field		AC	No. chromosomes carrying alt allele		MLEAF	Max likelihood AF
		AN	Total no. of chromosomes		MQ	RMS MAPQ of all reads
		AF	Allele frequency		MQ0	No. of MAPQ 0 reads at locus
		DP	Depth of coverage		QD	QUAL score over depth
		MLEAC	Max likelihood AC			
		0/1:12:99	0/1:11:89	0/1:28:37		

Note that VQSR will only look at INFO annotations;

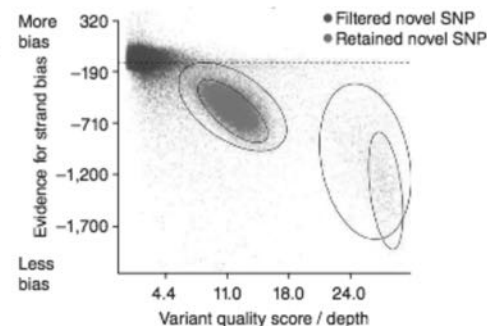
Two steps: (1) train a model then (2) apply to callset

**Basic idea: training on high-confidence known sites to determine the probability that other sites are true**

(1) Train model using HapMap

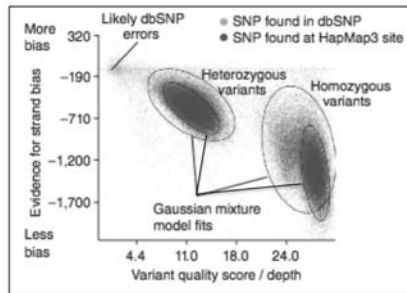


(2) Apply model to callset



## (1) Training the model

### (1) Train model using e.g. HapMap



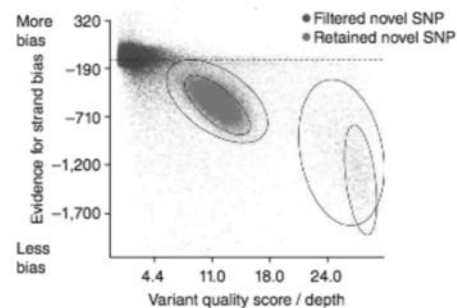
- We choose a training set
- Variants that are both in the training set and in our callset are selected.
- We train the model using the annotations of the selected variants

- This tells us **what good variants look like**
- A similar model for the variants in our callset that least look like good variants is also created (bad model, no biscuit!)
- All variants can now be ranked based on the ratio between their scores in the good model and the bad model (= VQSLOD)

## (2) Applying the model to our callset

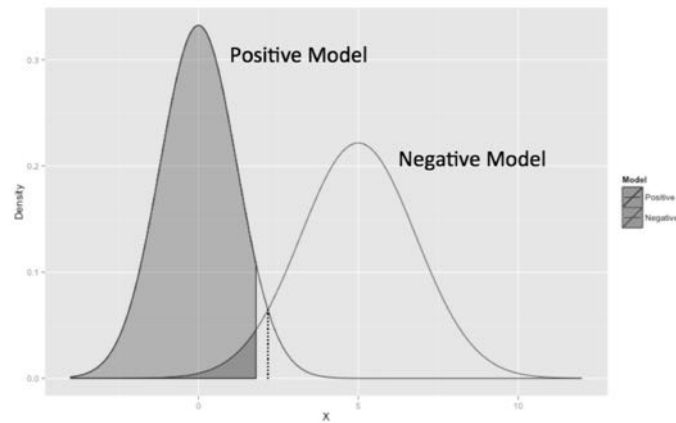
- Using the ranking produced by the model, filtering variants is as easy as setting a single threshold value
- Any variants whose score falls below the threshold is filtered out

### (2) Apply model to callset



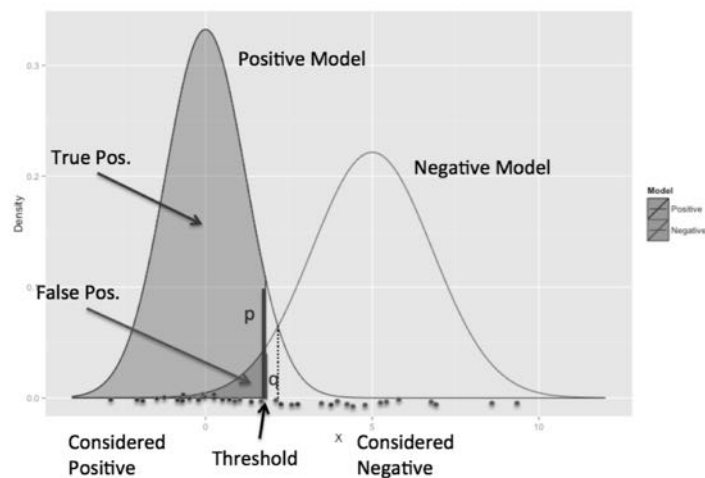
**But how do we set that threshold?**

There are in fact two components to the model



- A **negative model** is also built during training
- It represents the probability of variants to be **false positives**

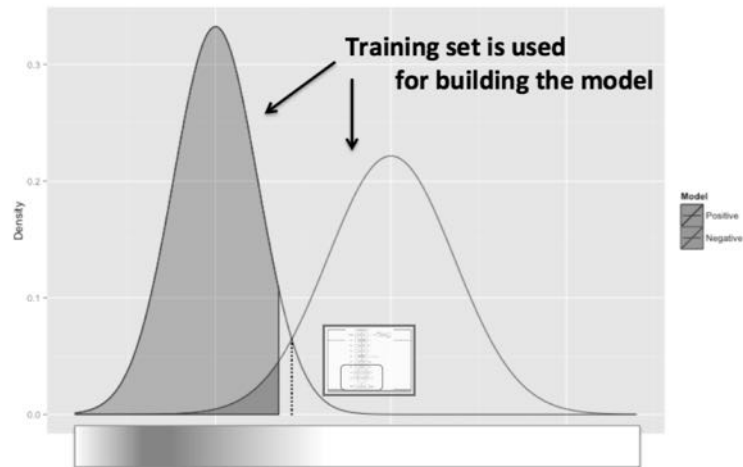
The VQSLOD threshold is a tradeoff between TP and FP



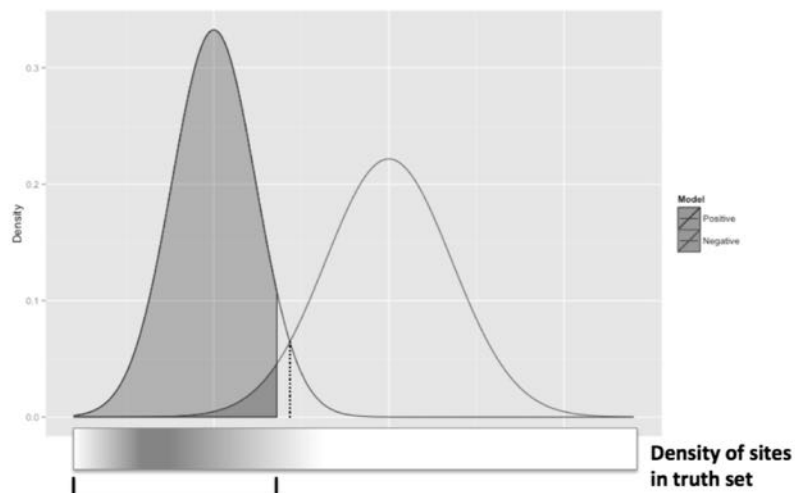
$$\text{VQSLOD}(x) = \text{Log}(p(x)/q(x))$$

(VQSLOD is distinct from QUAL!)

## Role of training and truth resources



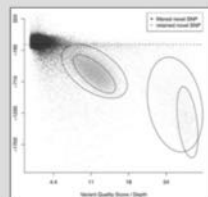
## We set the threshold based on sensitivity to truth data



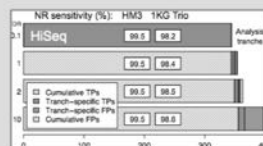
What threshold do we need to set to capture X % of the sites in the truth set?

## Variant Recalibration steps & tools

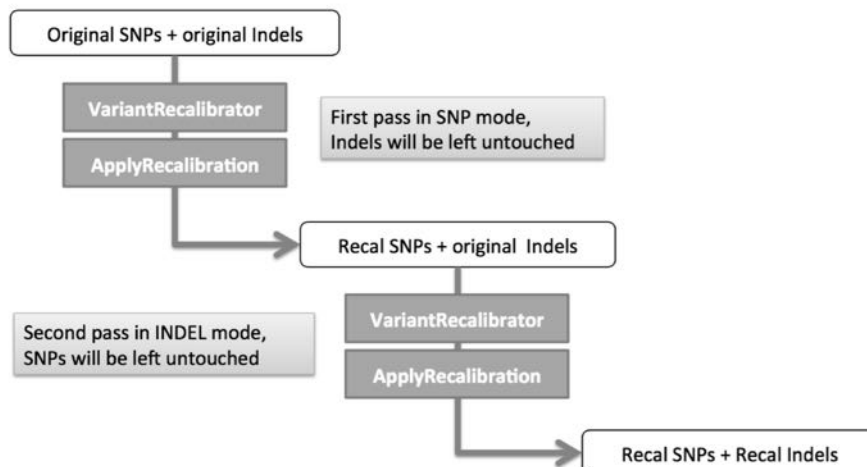
- Build and Apply the models (from resources and callset)  
→ **VariantRecalibrator**



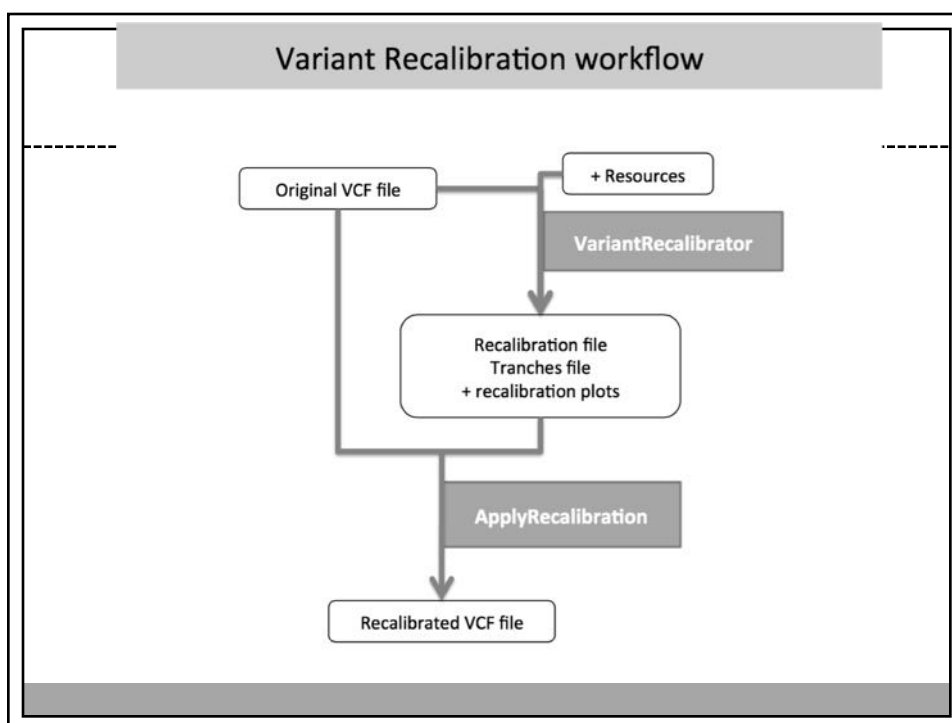
- Use VQSLOD to filter variants and write a new annotated VCF  
→ **ApplyRecalibration**



**NOTE: SNPs and Indels must be recalibrated separately!**



Pro-tip: Run VQSR twice in succession according to this workflow. That way you avoid having to split them, recalibrate and combine them again.



TOOL TIPS

## VariantRecalibrator

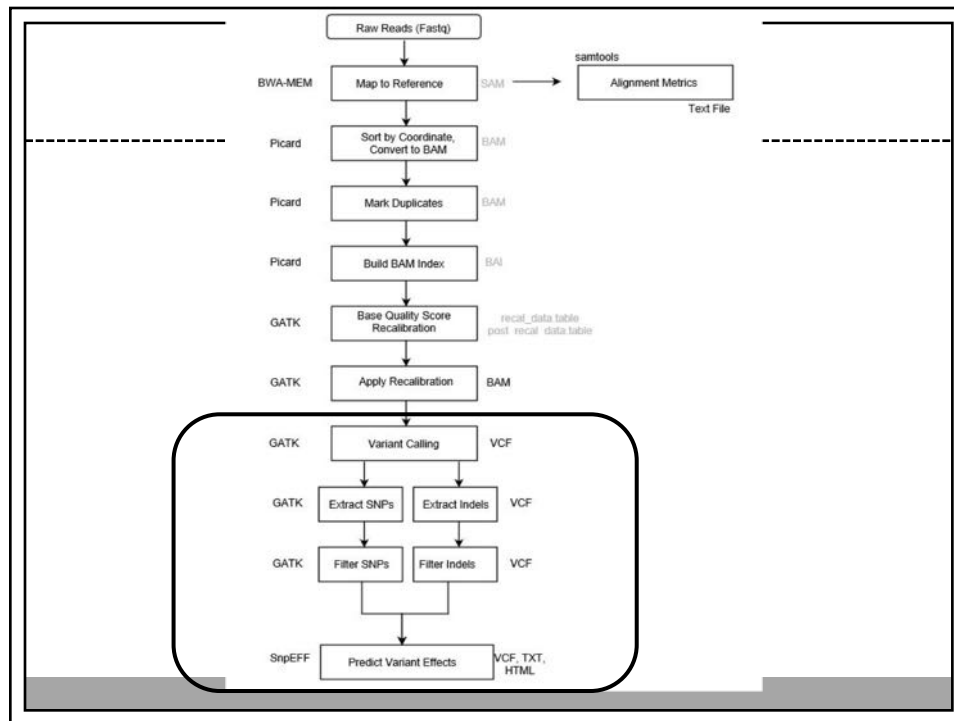
- Build the Gaussian mixture model using the variants in the input callset which overlap the training data

```

java -jar GenomeAnalysisTK.jar -T VariantRecalibrator \
  -R human.fasta \
  -input raw.SNPs.vcf \
  -resource: {see next slide} \
  -an DP -an QD -an FS -an MQRankSum {...} \
  -mode SNP \
  -recalFile raw.SNPs.recal \
  -tranchesFile raw.SNPs.tranches \
  -rscriptFile recal.plots.R

```

SNP example – see documentation for indel recommendations



## 1) Call Variants

- We use the GATK HaplotypeCaller tool
- This step is designed to maximize sensitivity in order to minimize false negatives, i.e. failing to identify real variants
- Creates a single file with both SNPs and indels
- We extract each type of variant into it's own file so we can process them individually



## 2) Filter Variants

- The first step is designed to maximize sensitivity and is thus very lenient in calling variants
- Good because it minimizes the chance of missing real variants
- But means that we need to filter the raw call set in order to reduce the amount of false positives
- Important in order to obtain the the highest-quality call set possible



## 3) Annotation

- We use SnpEff
- Annotates and predicts the effects of variants on genes
  - Codon changes
  - Amino acid changes
  - Genomic region
  - Functional effect (silent, missense)
- SnpEff has pre-built databases for thousands of genomes



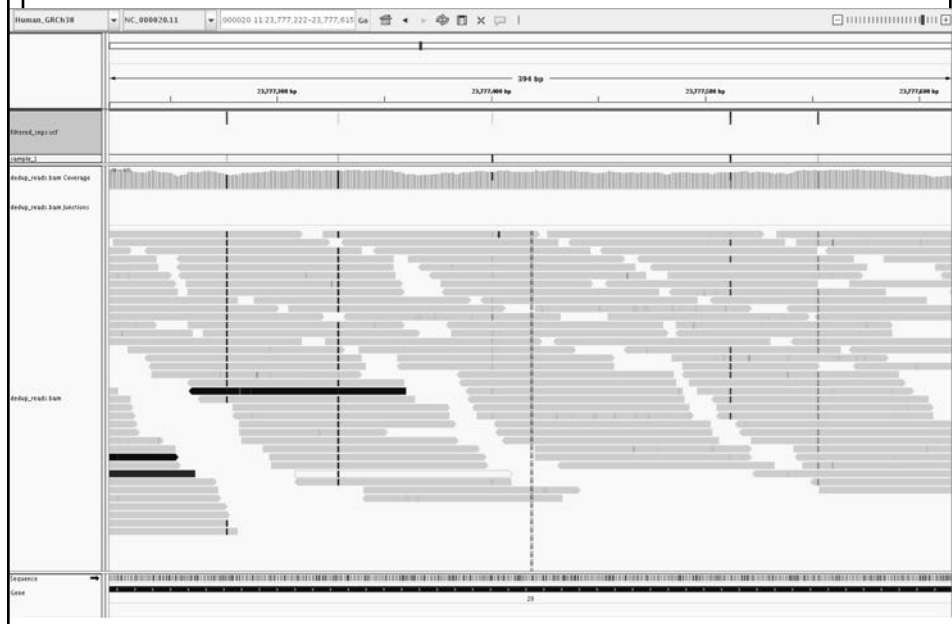
## Archivo VCF



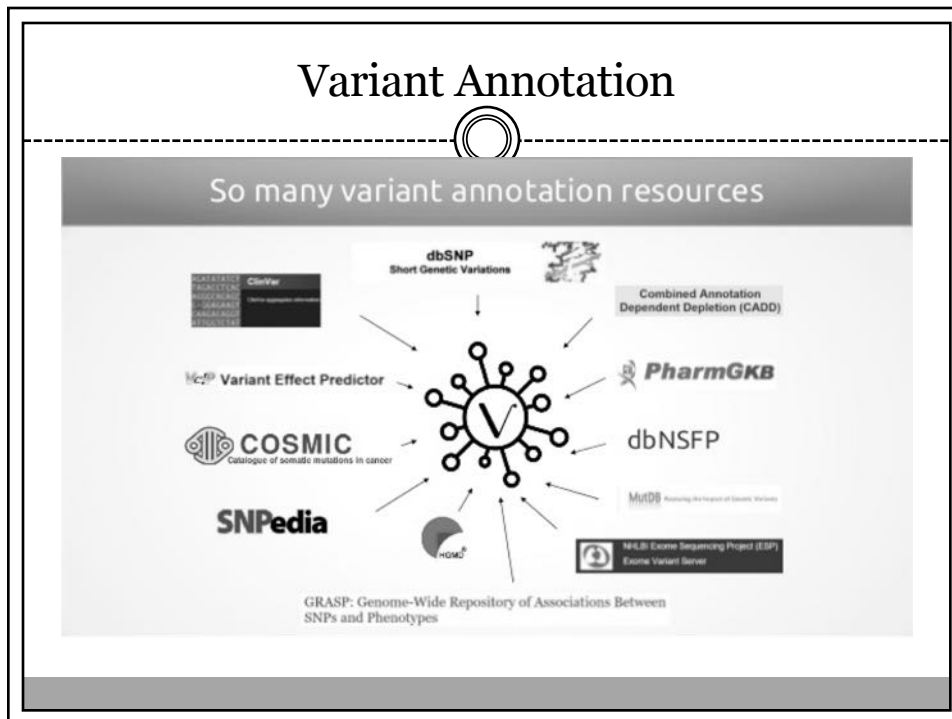
```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=1,Type=Float,Description="Allele 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 microsat GTCT G,GTACT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0|1:35:4 0|2:17:2 1|1:40:3
```

<http://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

## 4) Visualization - IGV



## Variant Annotation



## Variant Annotation: SnpEff

- Variant annotation and effect prediction tool. It annotates and predicts the effects of genetic variants (such as amino acid changes).
- Many effects are calculated: such as SYNONYMOUS\_CODING, NON\_SYNONYMOUS\_CODING, FRAME\_SHIFT, STOP\_GAINED just to name a few.

## SnpEff: Public databases



- **ENCODE** datasets are supported by SnpEff (by means of BigWig files provided by ENCODE project).
- **Epigenome Roadmap** provides data-sets that can be used with SnpEff.
- **TFBS** Transcription factor binding site predictions can be annotated. Motif data used in this annotations is generated by Jaspar and ENSEMBL projects
- **NextProt** database can be used to annotate protein domains as well as important functional sites in a protein (e.g. phosphorylation site)

## CADD - Combined Annotation Dependent Depletion



- Framework that integrates multiple annotations into one metric by contrasting variants that survived natural selection with simulated mutations
- C-scores strongly correlate with allelic diversity, pathogenicity of both coding and non-coding variants, and experimentally measured regulatory effects, and also highly rank causal variants within individual genome sequences.
- C-scores of complex trait-associated variants from genome-wide association studies (GWAS) are significantly higher than matched controls and correlate with study sample size, likely reflecting the increased accuracy of larger GWAS.

<https://cadd.gs.washington.edu/>

## Puntajes CADD

Filter

SIFT: All

PolyPhen: All

Consequences: missense variant

Source: dbSNP

CADD: All

Filter Other Columns

Show/hide columns

CADD

0 - 100

0

1

include blank ☐

Apply

Cancel

Mutation Assessor

0.812

0.968

Variant ID	Chr: bp	Alleles	Evidence	AA	AA coord	SIFT							
rs1437042753	3:25609282	G/A	AD	R/W	1327	0							
rs1246819283	3:25624786	G/A	AD	R/W	743	0							
rs1173166728	3:25627237	G/A	AD	R/C	651	0	0.988	34	0.52	0.415	0.936		
rs1357930313	3:25632778	T/C	AD	D/G	343	0	0.949	34	0.789	0.39	0.837		
rs776876327	3:25634002	G/A/T	AD	R/C	284	0.01	0.875	34	0.434	0.259	0.673		
rs376036396	3:25615236	T/A	AD	E/V	1182	0.01	0.786	33	0.362	0.295	0.709		
rs1349980311	3:25615269	A/G	AD	L/P	1171	0	0.998	33	0.836	0.801	0.967		

## Software Libre

**Galaxy**  
COMMUNITY HUB

**Galaxy**  
Analyze Data Workflow Shared Data Visualizations Help Login or Register

Tools

search tools

Get Data  
Lift-Over  
Collection Operations  
Text Manipulation  
Datamash  
Convert Formats  
Filter and Sort  
Join, Subtract and Group  
Fetch Alignments/Sequences  
NGS: QC and manipulation  
NGS: DeepTools  
NGS: Mapping  
NGS: RNA Analysis  
NGS: SAMtools  
NGS: BamTools  
NGS: Picard  
NGS: VCF Manipulation  
NGS: Peak Calling  
NGS: Variant Analysis  
NGS: RNA Structure  
NGS: Du Novus  
NGS: Gemini  
NGS: Assembly  
NGS: Chromosome Conformation  
NGS: Motif  
Operate on Genomic Intervals  
Statistics  
Graph/Display Data

**Data intensive for everyone**

Galaxy is an open, web-based platform for data intensive biomedical research. If you are new to Galaxy start here or consult our help resources. You can install your own Galaxy by following the tutorial and choose from thousands of tools from the Tool Shed.

Galaxy is an open, web-based platform for data intensive biomedical research. If you are new to Galaxy start here or consult our help resources. You can install your own Galaxy by following the tutorial and choose from thousands of tools from the Tool Shed.

**Try Galaxy on the Cloud**

Now you can have a personal Galaxy within the infinite Universe

Tweets by @galaxyproject

**Galaxy Project** @galaxyproject  
One of 16 highlighted pubs in this month's Galaxy News galaxyproject.org/galaxy-updates... #usegalaxy

History

search datasets

Unnamed history  
6 shared  
53.7 MB

6: Base Coverage on data 5  
5: UCSC Main on Human: ncbiRefSeq (genome)  
4: Base Coverage on data 3  
3: Merge on data 1  
2: Base Coverage on data 1  
1: UCSC Main on Human: ncbiRefSeq (genome)

**Galaxy**  
Aplicación web gratuita para análisis de datos NGS  
<https://usegalaxy.org/>

## Resources in NGS data analysis

### Public forums:



BioBits: A Blog about bioinformatics

