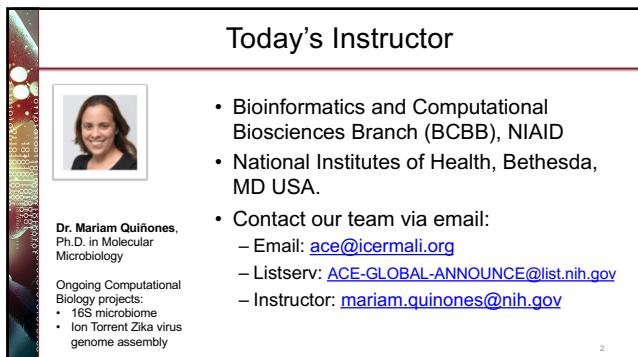


AFRICAN CENTERS OF EXCELLENCE
IN BIOINFORMATICS
KAMPALA, UGANDA

Introduction to DNA-Seq and Variant Calling
Mariam Quiñones, PhD

1



Today's Instructor

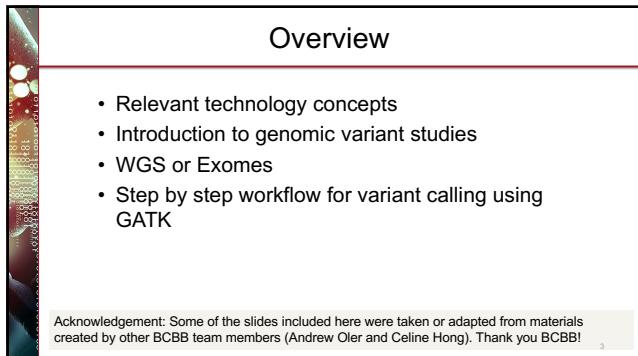


Dr. Mariam Quiñones,
Ph.D. in Molecular
Microbiology

Ongoing Computational
Biology projects:
• 16S microbiome
• Ion Torrent Zika virus
genome assembly

- Bioinformatics and Computational Biosciences Branch (BCBB), NIAID
- National Institutes of Health, Bethesda, MD USA.
- Contact our team via email:
 - Email: ace@icermali.org
 - Listserv: ACE-GLOBAL-ANNOUNCE@list.nih.gov
 - Instructor: mariam.quinones@nih.gov

2



Overview

- Relevant technology concepts
- Introduction to genomic variant studies
- WGS or Exomes
- Step by step workflow for variant calling using GATK

Acknowledgement: Some of the slides included here were taken or adapted from materials created by other BCBB team members (Andrew Oler and Celine Hong). Thank you BCBB!

3

**Important developments in studying genomic variations:
Timeline to Human Genome Project**

https://www.mun.ca/biology/scarr/Human_Genome_Project_timeline.html

Genomic variations have been studied in large populations studies

International HapMap Project

- Project that started on 2002 with the goal of describing patterns of human genetic variation and create a haplotype map using SNPs present in at least 1% of the population, which were deposited in dbSNPs.

1000 Genomes
A Deep Catalog of Human Genetic Variation

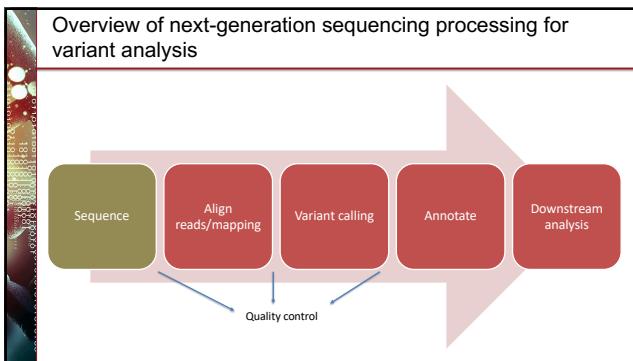
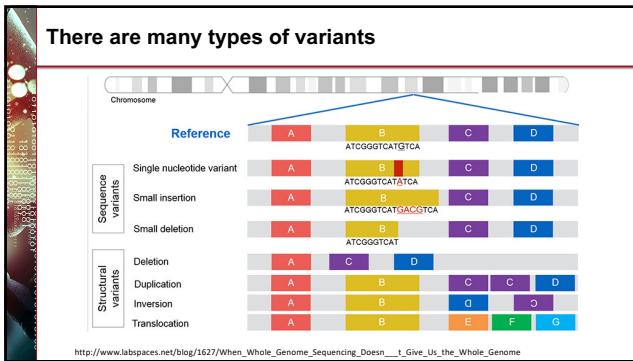
- Started in 2008 with a goal of using at least 1000 individuals (about 2,500 samples at 4X coverage), interrogate 1000 gene regions in 900 samples (exome analysis), find most genetic variants with allele frequencies above 1% and to a 0.1% if in coding regions as well as indels and structural variants
- Make data available to the public <ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/> or via Amazon Cloud <http://s3.amazonaws.com/1000genomes>

Variants are not only found in coding sequences

ENCODE

Main Goal:

- Find all functional elements in the genome



Sequencers

Machine	Cost	methods	Throughput per run	Read length	Error rate
Illumina MiSeq	\$128K	Small genomes, targeted gene	1.5-2Gb	2X300	0.8%
Ion Torrent	\$80K	Small genomes, targeted gene	1Gb	400	1.71%
Illumina NextSeq	\$250K	Exomes/transcriptome	120Gb	2x150	0.8%
Illumina HiSeq	\$654K	Genomes/exomes/transcriptomics	600Gb	2X150	0.76%
Illumina X Ten	\$10 Mil.	Genomes	1.6Tb	2x150	0.5%
PacBio	\$695K	Genomes	100Mb	15K	12.86%



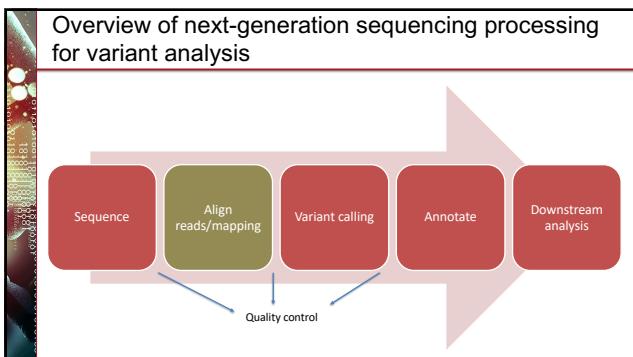
Next-generation sequencing overview		
	Exomes (ES/WES)	Genomes (GS/WGS)
Cost	~\$1000	~\$2000 (~\$1000 with HiSeq X Ten)
Size of bam files	~10 Gb	~200 Gb
DNA	Targeted and captured	Sheared DNA
What you can get	Most coding regions (+UTR)	coding and non-coding
Variants that can be examined	SNVs, indels (CNVs)	SNVs, indels, CNVs, structural variations



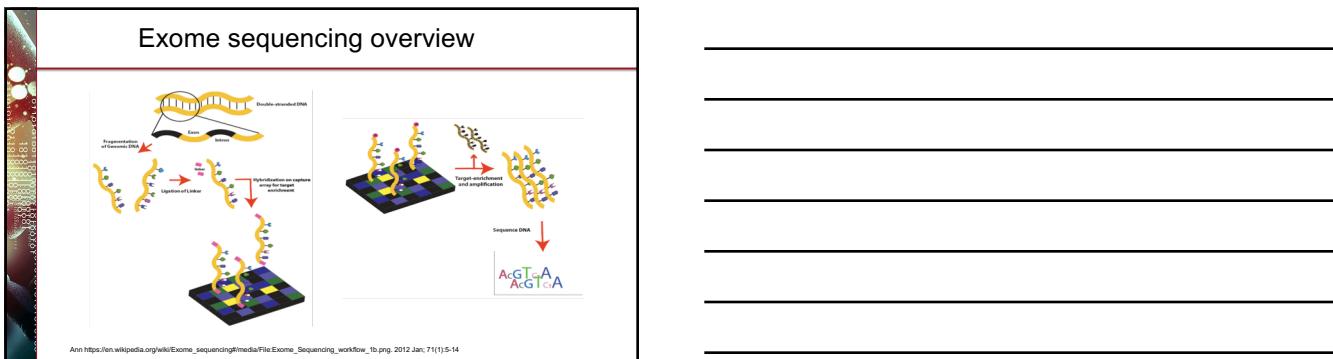
Processing reads for Variant Analysis

- Trimming:
 - Unless quality of read data is poor, low quality end trimming of short reads (e.g. Illumina) prior to mapping is usually not required.
 - If adapters are present in a high number of reads, these could be trimmed to improve mapping.
 - Recommended tools: *Prinseq*, *Btrim64*, *FastQC*
- Marking duplicates:
 - It is not required to remove duplicate reads prior to mapping but instead it is recommended to mark duplicates after the alignment.
 - Recommended tool: *Picard*.
 - A library that is composed mainly of PCR duplicates could produce inaccurate variant calling.

<http://www.biomedcentral.com/1471-2164/13/S8/S8>



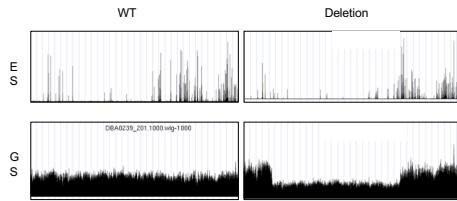
	<p>Variant analysis can be done using the entire genome (WGS) or only coding sequences (exome)</p>
<p>Whole Genome</p> <ul style="list-style-type: none"> • Useful for identifying SNPs but also structural variants in any genomic region • It is routinely done when studying genomes from bacteria or viruses 	<p>Exome-Seq</p> <ul style="list-style-type: none"> • Interrogates variants in specific genomic regions, usually coding sequences. • Cheaper and provides high coverage



	<p>Partial list of capture enrichment kits</p> <table border="1"> <thead> <tr> <th>Manufacturer</th><th>Kits</th><th>Regions targeted</th><th>Bases covered</th></tr> </thead> <tbody> <tr> <td>Illumina</td><td>Nextera Rapid capture</td><td>Exons + UTRs+ miRNA</td><td>62 Mb</td></tr> <tr> <td>Nimblegen</td><td>SeqCap EZ Exome</td><td>Exons + UTR</td><td>96 Mb</td></tr> <tr> <td></td><td>SeqCap EZMedExome</td><td>Disease-associated regions</td><td>47 Mb</td></tr> <tr> <td rowspan="2">Agilent</td><td>SureSelect Human All Exon V6</td><td>Exons+UTRs</td><td>60Mb</td></tr> <tr> <td>Clinical Research Exome</td><td>Disease-relevant targets</td><td>51Mb</td></tr> </tbody> </table>	Manufacturer	Kits	Regions targeted	Bases covered	Illumina	Nextera Rapid capture	Exons + UTRs+ miRNA	62 Mb	Nimblegen	SeqCap EZ Exome	Exons + UTR	96 Mb		SeqCap EZMedExome	Disease-associated regions	47 Mb	Agilent	SureSelect Human All Exon V6	Exons+UTRs	60Mb	Clinical Research Exome	Disease-relevant targets	51Mb
Manufacturer	Kits	Regions targeted	Bases covered																					
Illumina	Nextera Rapid capture	Exons + UTRs+ miRNA	62 Mb																					
Nimblegen	SeqCap EZ Exome	Exons + UTR	96 Mb																					
	SeqCap EZMedExome	Disease-associated regions	47 Mb																					
Agilent	SureSelect Human All Exon V6	Exons+UTRs	60Mb																					
	Clinical Research Exome	Disease-relevant targets	51Mb																					

How exomes and genomes look

- Genome has even coverage
 - Even a deletion is observed by eye



All about exomes

- What is exome?
 - Sequencing targeted exonic regions
 - ~2% of genome
 - Important to know
 - You will NOT get a whole exome!
 - Not all exons in all genes are captured!
 - Important to know the negative results vs no data
 - Coverage will vary in targets

What to know about your data

- What is the sequence depth?
 - The depth your sequence
 - 10x, 30x, 50x, 100x
 - Read length
 - How long is your read length?
 - What software was used to align?
 - Which reference was used?
 - Which capture kit was used? What is covered?
 - What variant calling was used to call variants?

Files and tools used for SNPs and Small INDELS

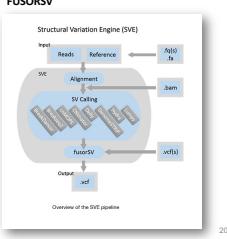
File type	Origin
FASTQ	Raw reads from sequencer
SAM	Sequence Alignment/Map
BAM	Binary version of SAM
gVCF/VCF	Variant call format

Tool	Purpose
BWA mem	Read alignment to reference
Picard	Mark duplicates
GATK (Haplotype caller)	QC/variant calling
Samtools	Sort sam/bam, convert sam<->bam

How to identify structural variants?

Try using tools which leverage information on read pairs or depth.

- 1- Read Pair – use pair ends.
GASVPro ([Sindi 2012](#)), BreakDancer (Chen 2009)
- 2- Read depth – CNVnator (Abyzov, 2010)
- 3- Split read – PinDEL (Ye 2009)
- 4- Denovo assembly



VCF format (version 4.3)

- Format used to report information about a position in the genome
- Used by 1000 genomes to report all variants

VCF FORMAT

The Variant Call Format (VCF) is a TAB-delimited format with each data line consisting of the following fields:

Col	Field	Description
1	CHROM	Chromosome name
2	POS	the left-most POSITION of the variant
3	ID	unique variant identifier
4	REF	the REFERENCE allele
5	ALT	the ALTERNATE allele(s); separated by comma
6	QUAL	variant/reference QUALITY
7	FILTER	FILTERs applied
8	INFO	INFORMATION related to the variant; separated by semi-colon
9	FORMAT	FORMAT of the genotype fields; separated by colon (optional)
10+	SAMPLE	SAMPLE genotypes and per-sample information (optional)

VCF format

See <http://samtools.github.io/hts-specs/VCFv4.3.pdf>

```

@HEADER LN:NA12878
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT
chr1 873762 rs3828047 T A 5231.78 PASS [AMPLIFICATION] 0F AD 87 90 PL 1/1:0.141:289:99:255:0,255
chr1 877644 rs2858421 A T 3931.66 PASS [AMPLIFICATION] 0F AD 38 90 PL 1/1:0.105:94:99:255:255,0
chr1 899282 rs28548431 C T 71.77 LowQual [AMPLIFICATION] 0F AD 39 90 PL 0/1:14.4:14:60:91:61,0,255
chr1 974165 rs9442791 T C 29.84 LowQual [AMPLIFICATION] 0F AD 39 90 PL 0/1:14.4:14:60:91:61,0,255

```

How variation is represented in a VCF

Each line represents one variant (here everything is a SNP, but some could be indels or CNVs) as well as the genotype of our sample, NA12878, at that variant. I've chosen these four variants because they each represent an important aspect in interpreting a VCF file:

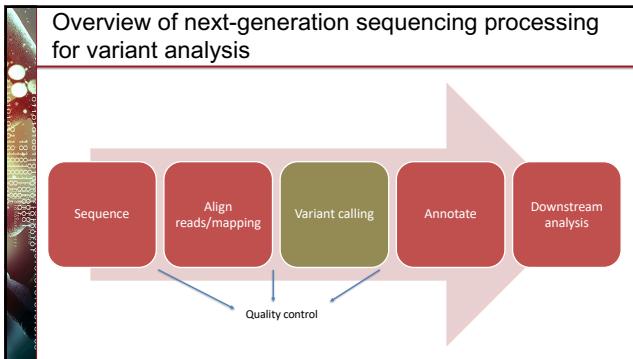
- chr1:873762 is a novel T/G polymorphism, found with very high confidence (QUAL = 5231.78).
- chr1:877644 is a known A/G SNP (rs3828047), found with very high confidence (QUAL = 3931.66)
- chr1:899282 is a known C/T SNP (rs28548431), but has a relative low confidence (QUAL = 71.77)
- chr1:974165 is a known T/C SNP but we have so little evidence for this variant in our data that although we write out a record for it (book keeping, really) our statistical evidence is so low that we filter the record out as a bad site "LowQual".

Popular tools for mapping, variant calling and annotation

Task	Popular tools
Align reads/mapping	BWA-mem Novoalign Isaac
Variant calling	GATK (Broad) Platypus (Wellcome Trust) Starling (Illumina)
Annotating variants	AnnoVar VEP SnpEff Gemini

Commonly used sources for reference genomes and annotations

Task	Popular source
Control population frequency	ExAC 1000GP ESP
Annotation	RefSeq Ensemble UCSC genes GENCODE
Visualization	UCSC genome browser IGV
Clinical relevance	HGMD OMIM CGD ClinVar
Variants and effect on gene expression	GTEX



Example of variant calling with GATK

Variation in coagulation and fibrinolysis genes evaluated for their contribution to cerebrovascular complications in adults with bacterial meningitis in the Netherlands

A.T. Kloek¹, H.N. Khan^{1,2}, M. Valls Serra³, A. Jongman¹, A.H. Zwinkelman³, F. Baas⁴, A. van der Ende², D. van de Beek¹, X. Bi², B. Ferwerda^{4,5}, M.C. Brouwer^{1,6}

- Sequencing: Solid 5500xl sequencer (Life technologies)
- Mapping: Alignment to GRCh37/HG19 reference genome using the Lifescope aligner (version 2.5.1, Applied Biosystems)
- Preprocessing and QC: Reads were realigned, recalibrated and called using GATK HaplotypeCaller (version 3.3-0). Variants were filtered on minimal mapping quality score ≥30, minimal read depth of 20, allele balance between 0.2 and 0.8 for heterozygotes, >90 for homozygotes and genotype quality score of 99.
- Additional QC and Filtering: All samples were combined and converted to PLINK format.¹⁹ Final QC was performed to exclude related individuals and samples with non-European ancestry based on a previous multidimensional scaling test.²⁰ Markers with a missingness of 5% and Hardy-Weinberg equilibrium (HWE) p-value of <1.0⁻⁵ were excluded from the analysis. For the regression analysis only markers with minor allele frequency (MAF) ≥ 0.01 were included.

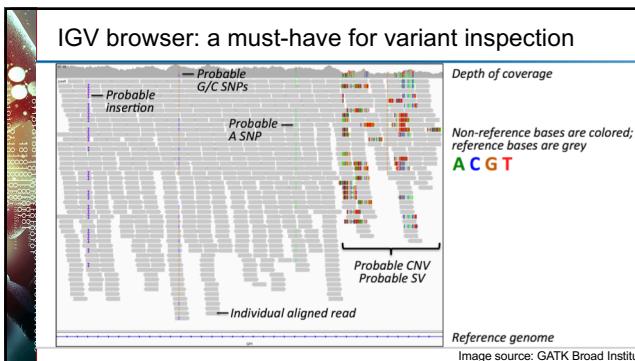
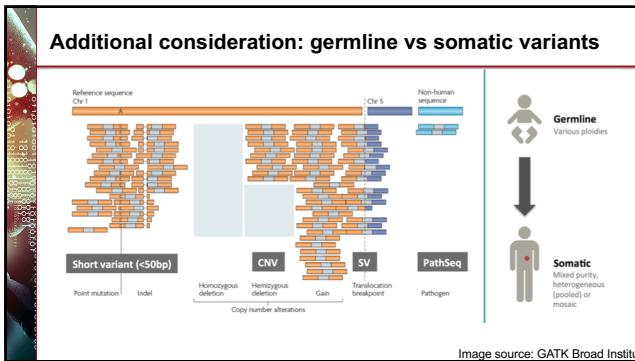
Example of variant calling and filtering with FreeBayes / samtools

Front. Genet. | www.frontiersin.org | Volume 9 | Article 147 | April 24, 2018 | DOI: 10.3389/fgene.2018.00147 | PMID: 29500203 | PMC ID: PMC5898203

A New Single Nucleotide Polymorphism Database for Rainbow Trout Generated Through Whole Genome Resequencing

Guanqiu Ge,¹ Torben Nebe,² Devon E. Pearce,³ Thomas Moen,⁴ Kerry A. Neale,⁵ Gary H. Thorpe,⁶ System Lines,⁷ and Yizhu Pei^{1*}

- SNP quality filter (QUAL):** SNPs not located within 4 bases distance to an indel, and have the phred-scaled variant quality score, QUAL, larger than 30.
- Low-complexity filter (LC):** We removed the SNPs that are in the low-complexity regions as using mdust (<https://github.com/lh3/mdust>)
- Maximum depth filter (DP):** maximum depth value of 1,500, corresponding to a coverage of 24.6 reads per sample
- Minimum sample filter (NS):** Only SNP sites with sequence reads coverage from at least 58 of the 61 samples passed this filter,
- Double haploid filter (DH):** Briefly, we filtered out SNP sites with heterozygous genotypes in at least two of the 11 DH lines that we re-sequenced in this study.



Reference genome

- There are different versions of human reference genome

Reference Name	Chr notation	Mitochondrial sequence	Additional sequences included
GRCh38 (Genome Reference Consortium)	1, 2... XY, MT	Yes	<ul style="list-style-type: none"> Unlocalized Unplaced Alternate loci Centromere sequences
hg19 (UCSC genome browser)	Chr1, Chr2... ChrX, ChrY, ChrM	Copied from previous release	<ul style="list-style-type: none"> Unlocalized Unplaced Alternate loci
b37/b37+decoy/hs37d5 (1000GP)	1, 2... XY, MT	Yes	<ul style="list-style-type: none"> Unlocalized: chromosome known, exact location unknown Unplaced: known to originate from human genome, chromosome unknown Alternate loci: alternate representation of specific human regions

Index files

- Index files are needed for files in next-generation analysis, as file sizes are big!
- Enables program to efficiently access the data, rather than having to read the whole file

File type	Index
FASTA	*fai
BAM	*bai
VCF	*vcf.idx

Reference FASTA index file

1	242529621	52	60	61
2	242529621	253464986	60	61
3	19882438	58057651	60	61
4	191154279	78198587	60	61
5	18982438	88057651	60	61
6	171115867	1888251307	60	61
7	159138663	1254218344	60	61
8	14892438	14892438	60	61
9	141213433	156482846	60	61
10	139534747	1788379889	60	61
11	138518989	18882438	60	61
12	138518989	1983430282	60	61
13	115169870	2119513996	60	61
14	18982438	22882438	60	61
15	182531392	2345741279	60	61
16	98354753	2449981581	60	61
17	83120528	2649981581	60	61
18	78877248	2624399817	60	61
19	59128983	278379886	60	61
20	63828983	278379886	60	61
21	48129895	282759925	60	61
22	51289866	282759925	60	61
X	1512728560	2929851273	60	61
Y	59373566	3886910193	60	61
MT	1493569	2147220000	60	61
GL008297..1	4262	3147239025	60	61
GL008226..1	15896	3147294651	60	61
GL008226..2	15896	3147294651	60	61
GL008231..1	27386	3147330289	60	61
GL008218..1	27682	3147358194	60	61
GL008231..2	27682	3147358194	60	61
GL008235..1	34474	314720850	60	61

GATK Best Practices Website

The GATK Best Practices website provides detailed instructions for variant discovery analysis. It includes recommended workflows for variant discovery analysis.

What is in the best practices?

The GATK Best Practices workflow provides the most up-to-date recommendations for variant discovery analysis. It includes a detailed description of each step in the process, including pre-processing, alignment, and variant calling. It also includes a section on quality control and reporting.

Current version is 4.2

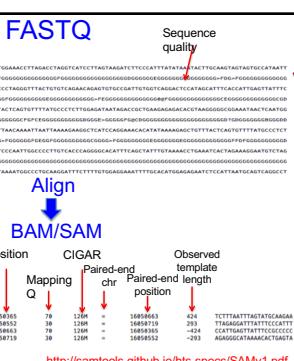
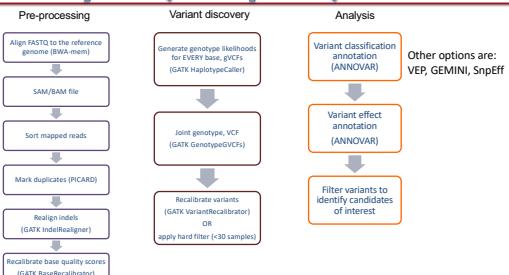
The Best Practices have been updated to reflect the latest developments in the field. This includes new recommendations for pre-processing, alignment, and variant calling. It also includes a section on quality control and reporting.

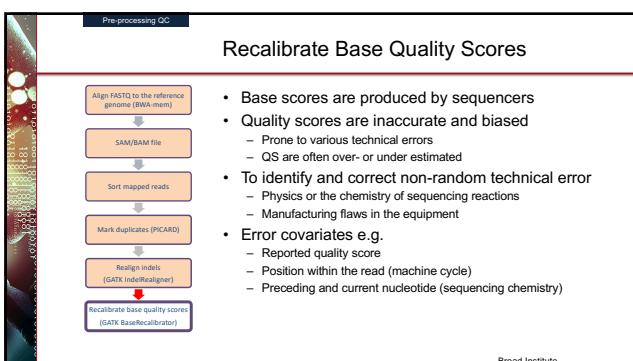
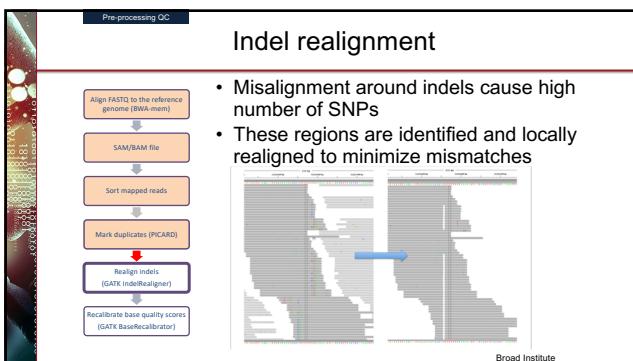
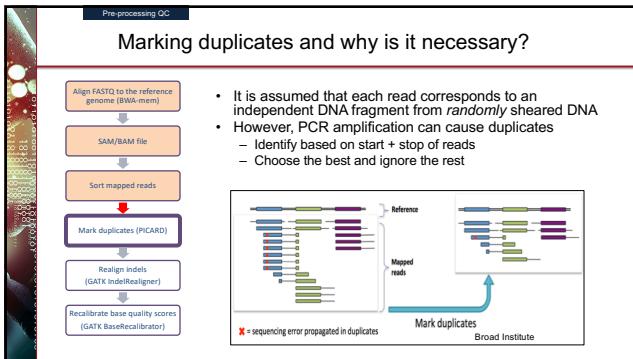
<https://www.broadinstitute.org/gatk/guide/best-practices>

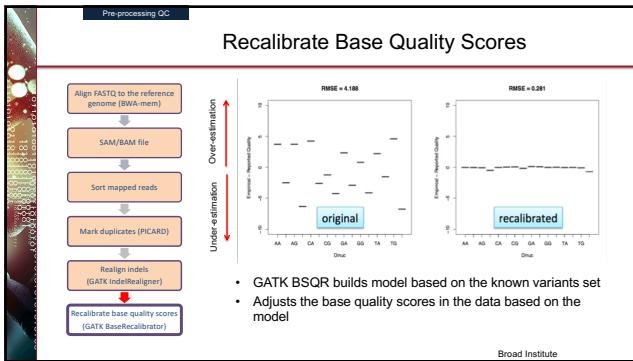
Why BWA+ GATK Haplotype caller?

- Widely accepted as the “conventional” way of processing next-gen data
- Well assessed
- Well documented
- Software is supported
- Community support for troubleshooting or information

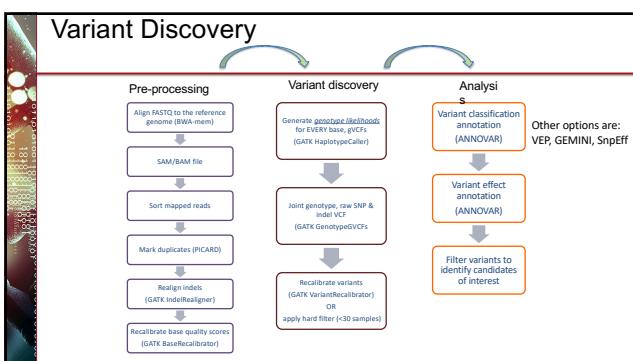
Data processing







Now we are ready to call variants!



GATK uses a bayesian model

Bayesian model

$$\Pr\{G|D\} = \frac{\Pr\{G\} \Pr\{D|G\}}{\sum_i \Pr\{G_i\} \Pr\{D|G_i\}}$$

$$\Pr\{D|G\} = \prod_j \left(\frac{\Pr\{D_j|H_1\}}{2} + \frac{\Pr\{D_j|H_2\}}{2} \right)$$

where $G = H_1H_2$

$\Pr\{D|H\}$ is the haploid likelihood function

- It determines possible SNP and indel alleles
- Computes, for each sample, for each genotype, likelihoods of data given genotypes
- Computes the allele frequency distribution to determine most likely allele count, and emit a variant call if determined
- When it reports a variant, it assigns a genotype to each sample

http://www.broadinstitute.org/gatk/events/2038/GATKwh0-BP-5-Variant_calling.pdf

Variant Discovery

Generate genotype likelihoods in each sample (gVCF)

Generate genotype likelihoods for EVERY base, gVCFs (GATK HaplotypeCaller)

Joint genotype, raw SNP & indel VCF (GATK GenotypeGVCFs)

Recalibrate variants (GATK VariantRecalibrator) OR apply hard filter (<30 samples)

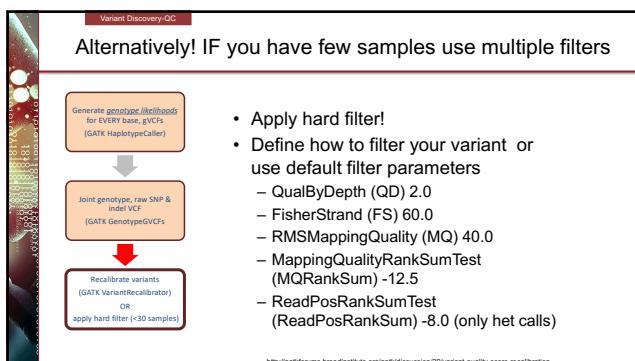
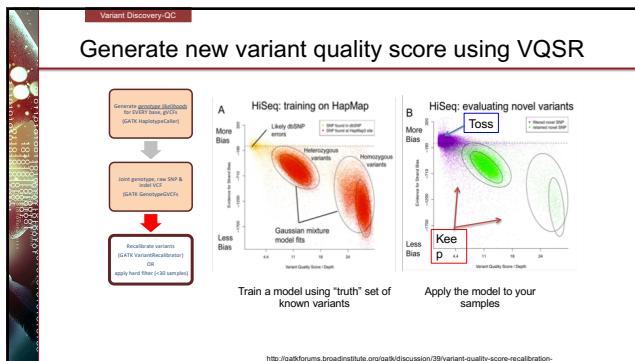
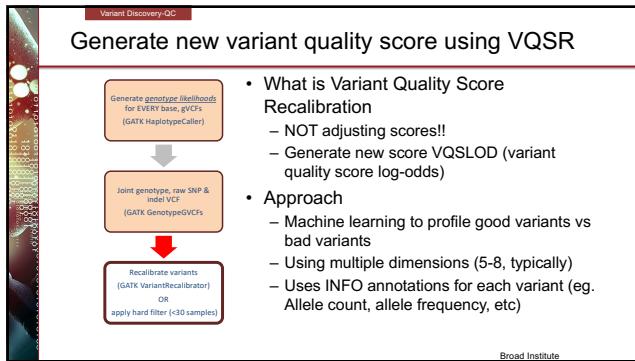
Broad Institute

Variant Discovery

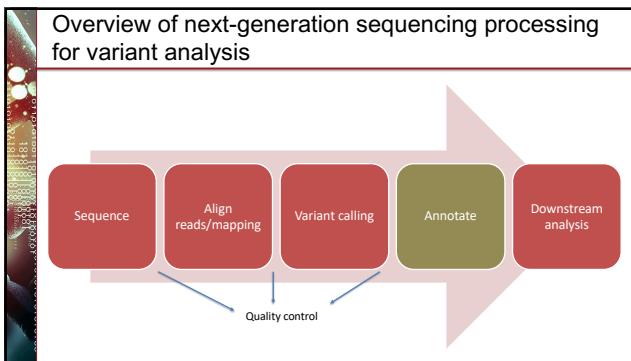
What is a joint genotyping?

- If we analyze Sample 1 or Sample N alone, we are not confident that the variant is real
- If we see both samples, we are more confident that there is real variation at this site in this cohort

Broad Institute



Variant Discovery-Final VCF		Final VCF																																																			
Header																																																					
<pre>#fileformat=VCFv4.1 #FILTER=ID=lowQual,Description="Low quality" #FORMAT<id=ad,number=1,type=integer,description="allelic alleles="" alt="" and="" depths="" for="" in="" listed"="" order="" ref="" the=""> #FORMAT<ID=DP,Number=1,Type=Integer,Description="Approximate read depth (reads with MQ<255 or with bad mates are filtered out)" #FORMAT<ID=GT,Number=1,Type=String,Description="Genotype Quality" #FORMAT<ID=GT,Number=1,Type=String,Description="Genotype" #FORMAT<ID=GL,Number=1,Type=String,Description="Normalized, Phred-scaled likelihoods for genotypes as defined in the VCF specification" #RGATCommandLine=HiMapotypeCaller<ID=HiMapotypeCaller,Version=3.4-3-gd1ac12,date="Mon May 18 23:36:4 : #INFO<ID=AC,Number=A,Type=Integer,Description="Allele count in genotypes, for each ALT allele, in the same order as listed" #INFO<ID=AF,Number=A,Type=Float,Description="Allele Frequency, for each ALT allele, in the same order as listed" #INFO<ID=AN,Number=1,Type=Integer,Description="Total number of alleles in called genotypes" #contig=chr1,length=49230623,assembly=b37 #reference=file:human_genome_b37.fasta</id=ad,number=1,type=integer,description="allelic></pre>																																																					
Body																																																					
<table border="1"> <thead> <tr> <th>CHROM</th><th>POS</th><th>ID</th><th>REF</th><th>ALT</th><th>QUAL</th><th>FILTER</th><th>INFO</th><th>FORMAT</th><th>NAL2878</th></tr> </thead> <tbody> <tr> <td>1</td><td>873762</td><td></td><td>G</td><td>5231</td><td>78</td><td>PASS</td><td>[ANNOTATIONS]</td><td>GT:AD:DP:GQ:PL</td><td>0/1:171:151:282:99:255,0,255</td></tr> <tr> <td>1</td><td>873762</td><td>r15220697</td><td>T</td><td>C</td><td>66</td><td>PASS</td><td>[ANNOTATIONS]</td><td>GT:AD:DP:GQ:PL</td><td>1/1:171:151:250:99:255,0</td></tr> <tr> <td>1</td><td>899282</td><td>r2545451</td><td>C</td><td>T</td><td>21</td><td>PASS</td><td>[ANNOTATIONS]</td><td>GT:AD:DP:GQ:PL</td><td>1/1:151:151:181:99:25</td></tr> <tr> <td>1</td><td>974165</td><td>r9442391</td><td>T</td><td>C</td><td>29.84</td><td>lowQual</td><td>[ANNOTATIONS]</td><td>GT:AD:DP:GQ:PL</td><td>0/1:114,4:14:61:61,61,255</td></tr> </tbody> </table>				CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NAL2878	1	873762		G	5231	78	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:171:151:282:99:255,0,255	1	873762	r15220697	T	C	66	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:171:151:250:99:255,0	1	899282	r2545451	C	T	21	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:151:151:181:99:25	1	974165	r9442391	T	C	29.84	lowQual	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:114,4:14:61:61,61,255
CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NAL2878																																												
1	873762		G	5231	78	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:171:151:282:99:255,0,255																																												
1	873762	r15220697	T	C	66	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:171:151:250:99:255,0																																												
1	899282	r2545451	C	T	21	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:151:151:181:99:25																																												
1	974165	r9442391	T	C	29.84	lowQual	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:114,4:14:61:61,61,255																																												



```

graph TD
    A[Variant classification annotation  
(e.g. ANNOVAR, SnpEff or GEMINI)] --> B[Region-based Annotation]
    B --> C[Gene Annotation]
    B --> D[Filter-based annotation]
    C --> E[Conserved genomic elements,  
transcription factor binding site, cytogenetic  
band, segmental duplications, GWAS...]
    D --> F[1000 GP, dbSNP, ESP, EXAC, non-synonymous variants annotation  
(SIFT/Polyphen2/MutationTaster/LRT/FATHM  
M/CADD-), ClinVar...]
    E --> G[Filter variants to identify candidates of interest]
    F --> G
  
```

The diagram illustrates the process of variant annotation. It starts with 'Variant classification annotation' (e.g., ANNOVAR, SnpEff or GEMINI), which branches into 'Region-based Annotation' and 'Filter-based annotation'. 'Region-based Annotation' leads to 'Gene Annotation' (refSeq Gene, mitochondrial variants, USC/Ensembl, GENCODE/CCDS). 'Filter-based annotation' leads to a list of 1000 GP, dbSNP, ESP, EXAC, non-synonymous variants annotation tools (SIFT/Polyphen2/MutationTaster/LRT/FATHM M/CADD-), and ClinVar. Both paths converge at 'Filter variants to identify candidates of interest'.

Variant Analysis...
like finding a needle in a 'deep' haystack

www.pete.co.uk

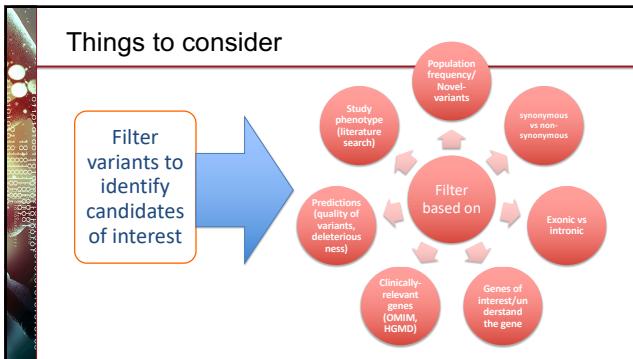
Look for evidence of variants of interest

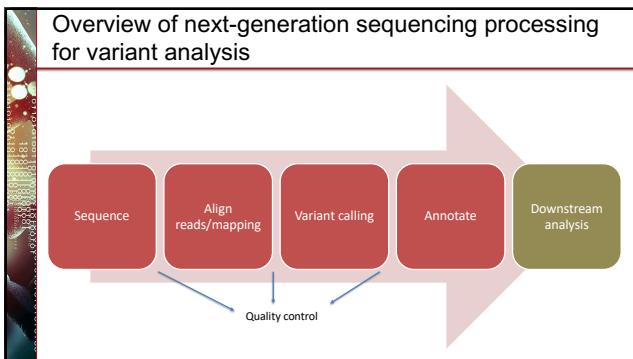
Further filters needed

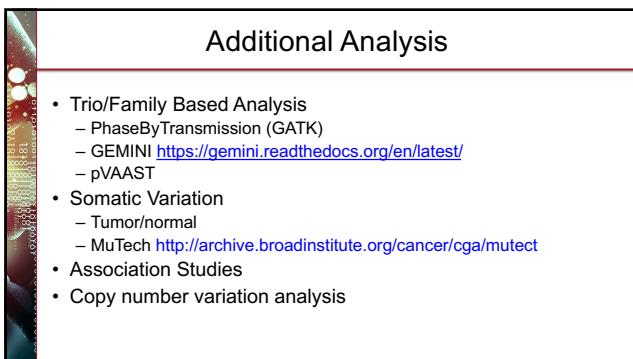
- High number of variants
- The goal is to narrow down your list of variants
- Eliminate variants that are not interesting

17,687 SNPs
dbsSNPs
1000 genomes
In house exome
2 novel in chr10

PLoS One. 2012;7(1):e29708 54

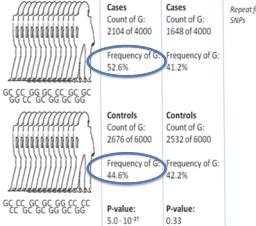




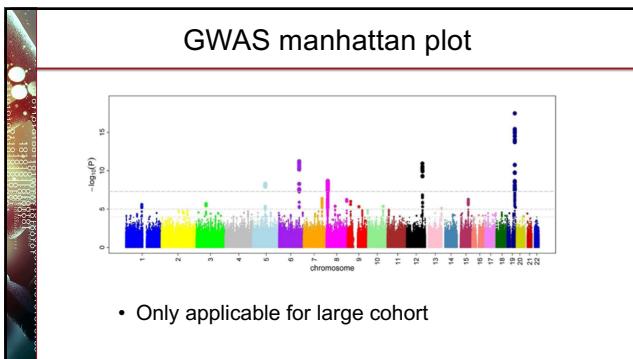


GWAS Association studies

- A typical analysis:
 - Identify SNPs where one allele is significantly more common in cases than controls
 - Hardy-Weinberg Chi Square



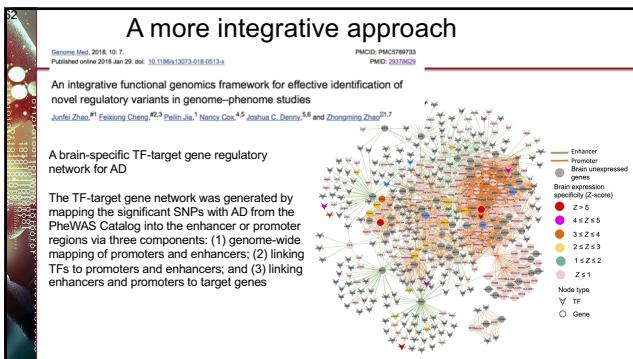
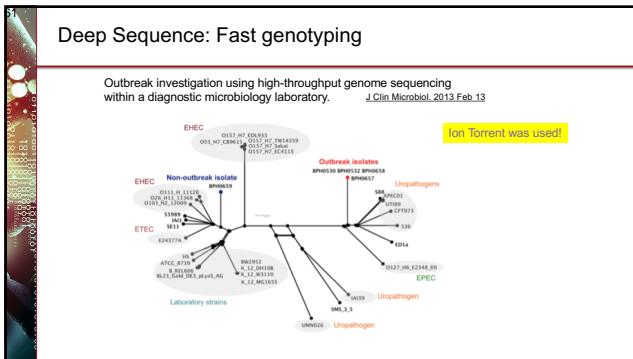
http://en.wikipedia.org/wiki/Genome-wide_association_study <http://www.ebi.ac.uk/gwas/>



CNVs from exome

- High variability of read-depth in exomes
- CNV prediction is very challenging problem
 - High false positives
 - Break-points limited due to targets
 - Long range CNVs have higher positive predictive value
- Useful when other alternatives (SNPs or aCGH) are not available
- Lots of prediction tools!! (some examples below)

Population Caller	Somatic Caller
XHMM CoNIFER EXCAVATOR ExomeDepth CONTRA	ADTEx ExomeCNV VarScan2 Control-FREEC



Let's practice