

Variant calling

with Illumina whole genome shotgun sequence data

Qi Sun, Robert Bukowski

Bioinformatics Facility, Institute of Biotechnology

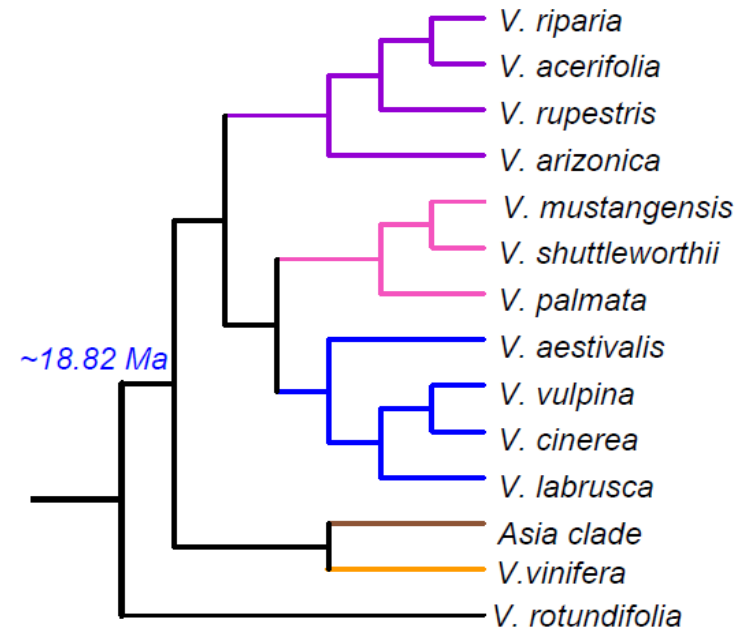
brc_bioinformatics@cornell.edu

How to sequence genomes of every grape
vine (or every cat) in the world?

- Cost effectively

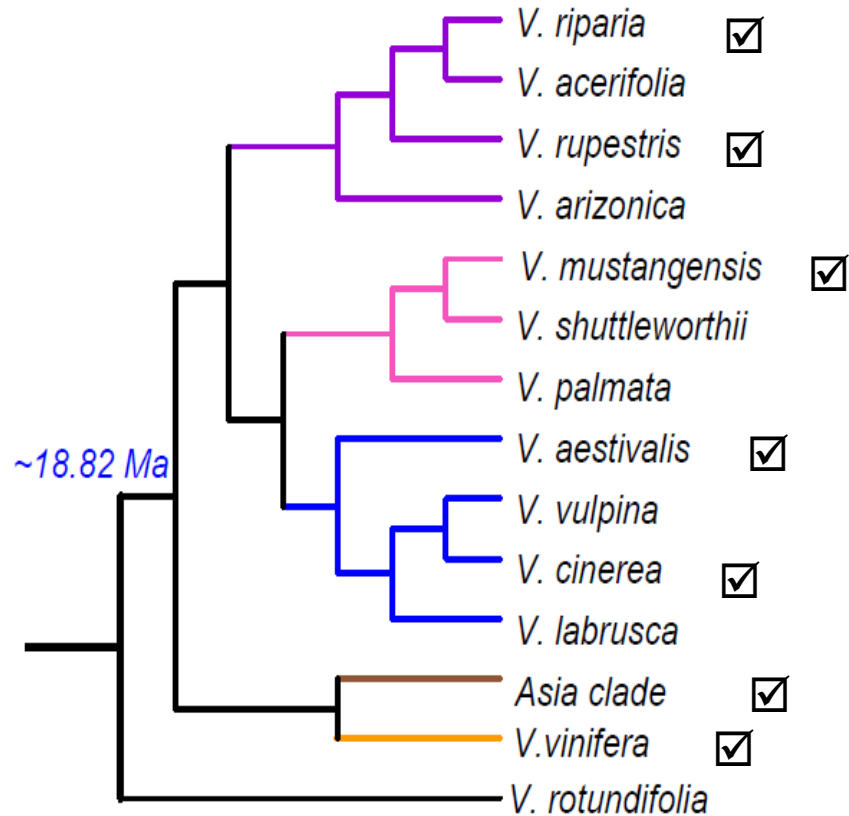
How to represent genomes of every grape
vine (or every cat) in the world?

- Variant matrix (vcf) vs Pan-genome graph



De novo genome assembly

- to capture genome structure variation



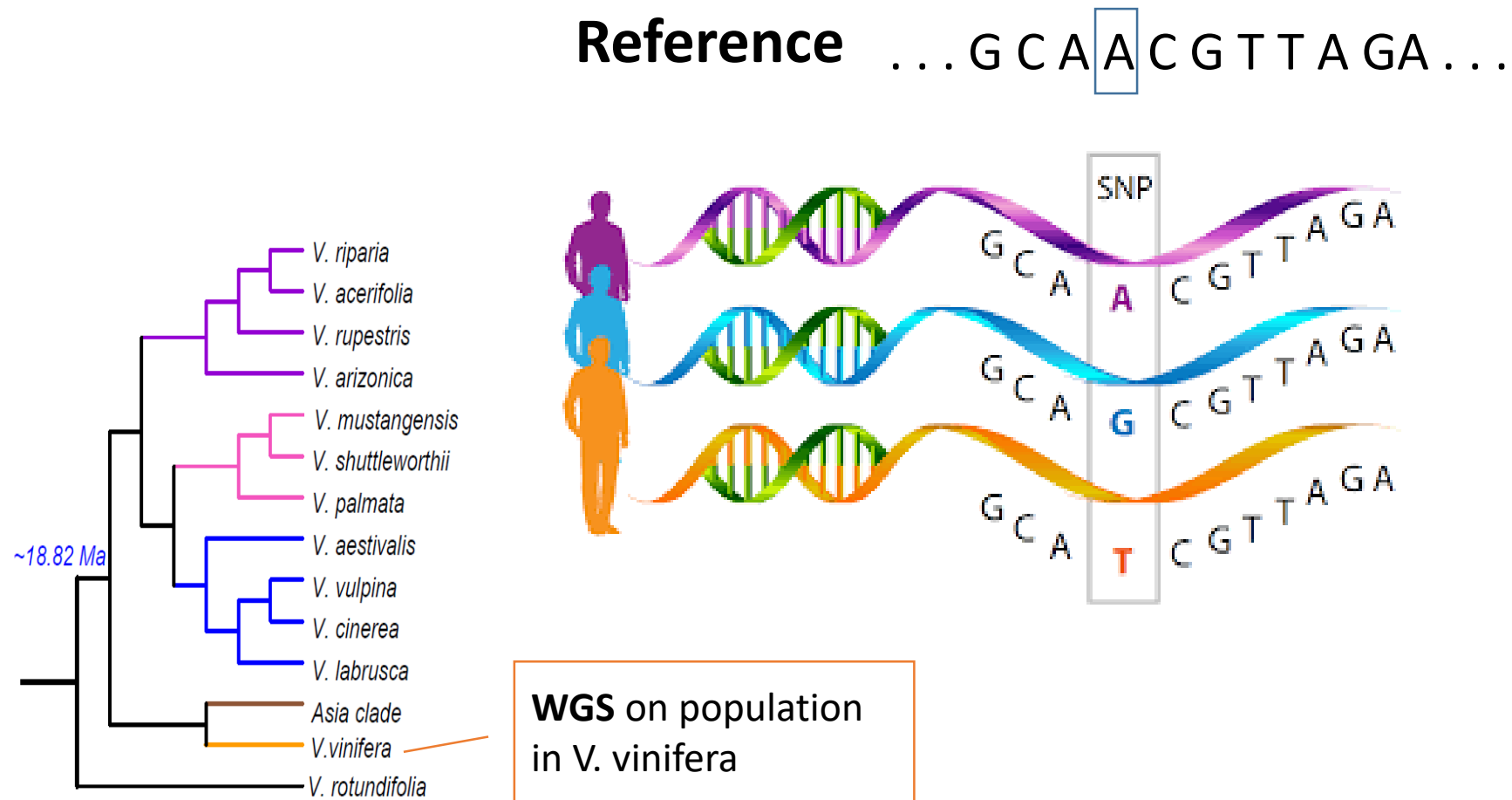
Select a panel of individuals to represent all genetic diversity in grapes;

- All gene space;
- All chromosomal structural variations, e.g. large insertions/deletions and translocations

Sequence and assemble using the best technologies we have today, e.g. PacBio, Nanopore, BioNano, Hi-C, et al.

Whole Genome Shotgun (WGS), a.k.a. re-sequencing

- to capture SNPs and short INDELs

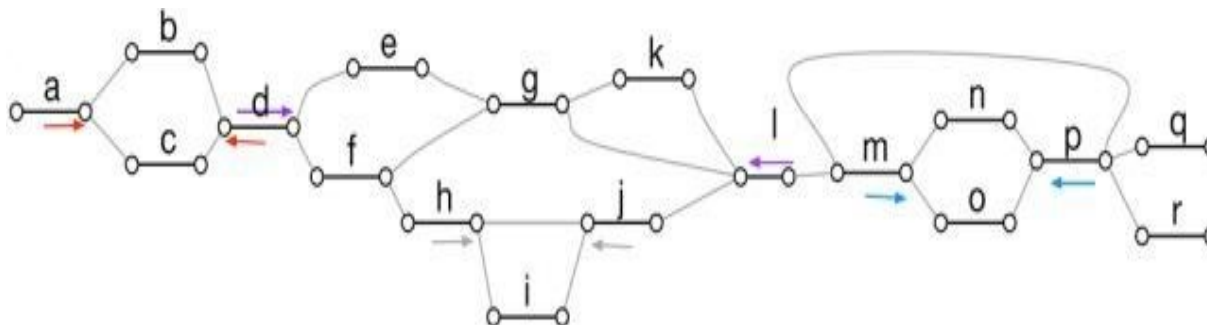


Representation of genomes of a population

SNP & INDELs: variant matrix (VCF file)

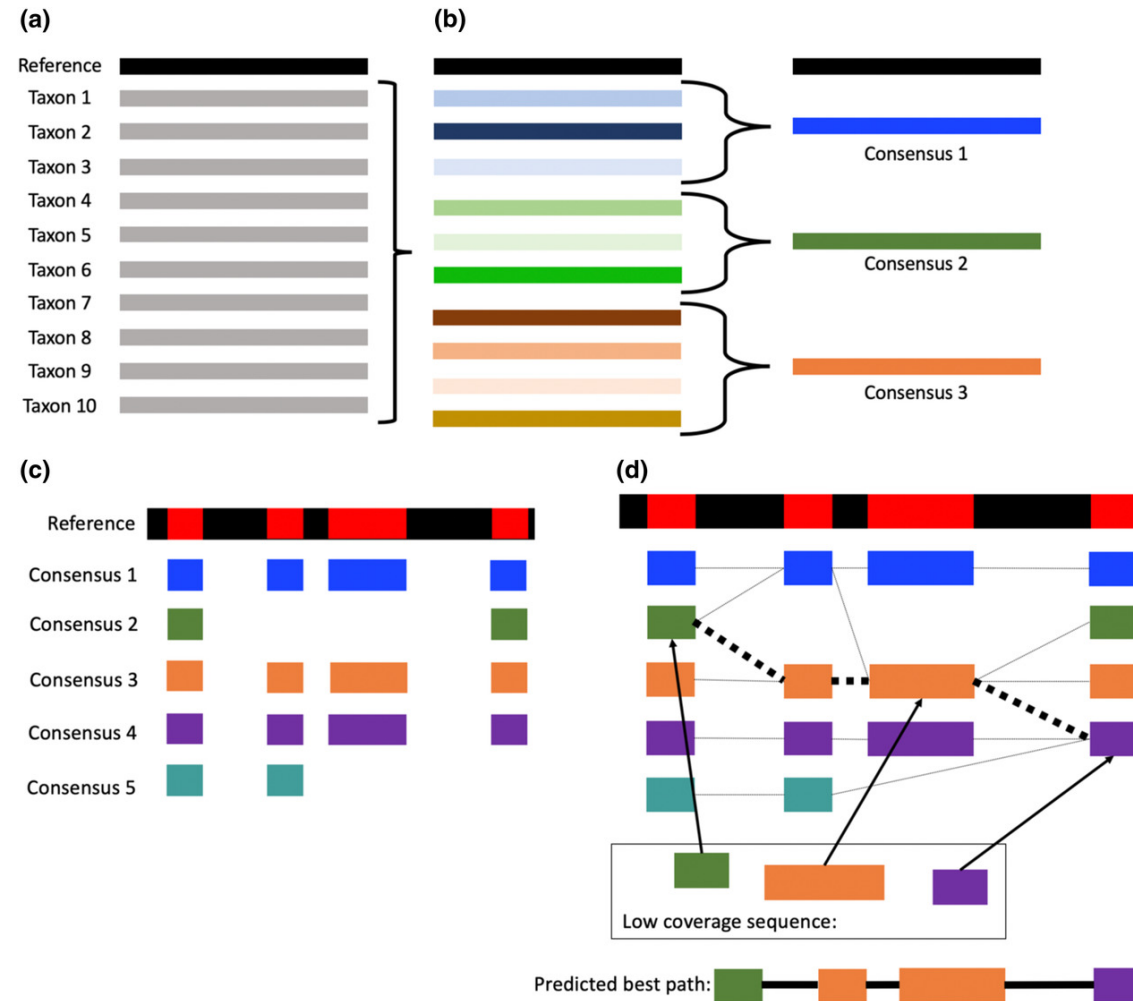
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	ZW155	ZW177
chr2R	2926	.	C	A	345.03	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:4,9:13:80:216,0,80	0/0:6,0:6:18:0,18,166
chr2R	9862	.	TA	T	180.73	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:0,5:5:15:97,15,0	1/1:0,4:4:12:80,12,0
chr2R	10834	.	A	ACTG	173.04	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/0:14,0:14:33:0,33,495	0/1:6,3:9:99:105,0,315

Major structure variation: pan genome graph



Paten et al. *Genome Res.* 2017, May; **27(5)**: 665-676.

Practical Haplotype Graph



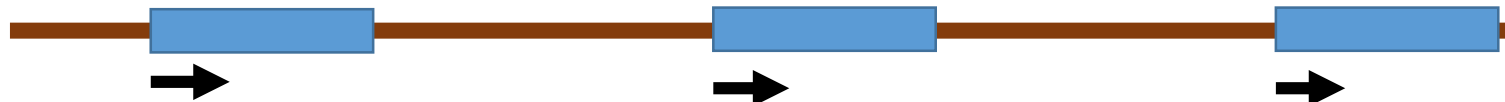
Sarah Jensen
Ed Buckler

SNP array, targeted or skim sequencing

(Not covered in this workshop)

Technologies

- SNP array
- GBS / RAD
- Amplicon
- SKIM

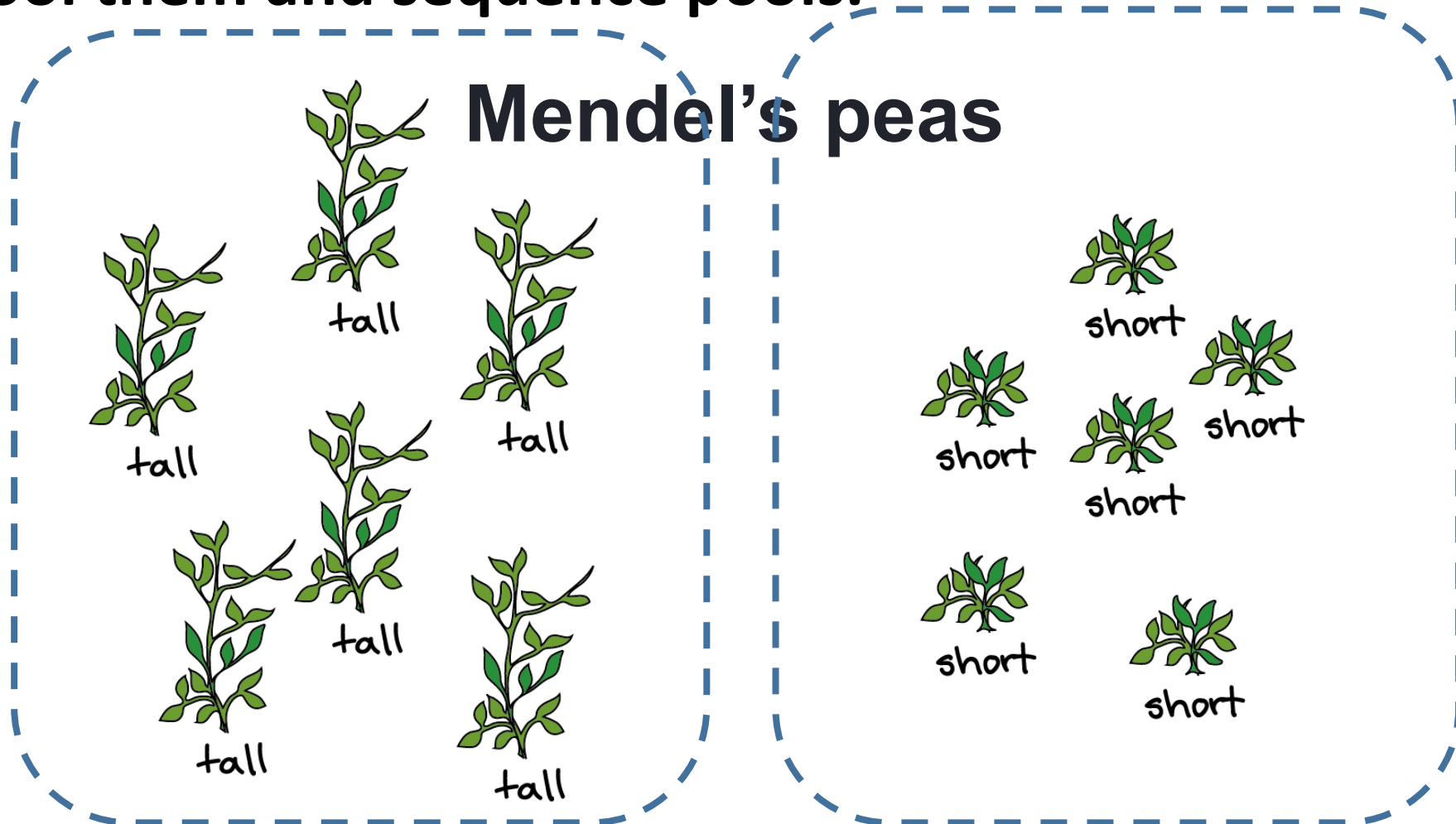


Capture major haplotype alleles. Need to be imputed.

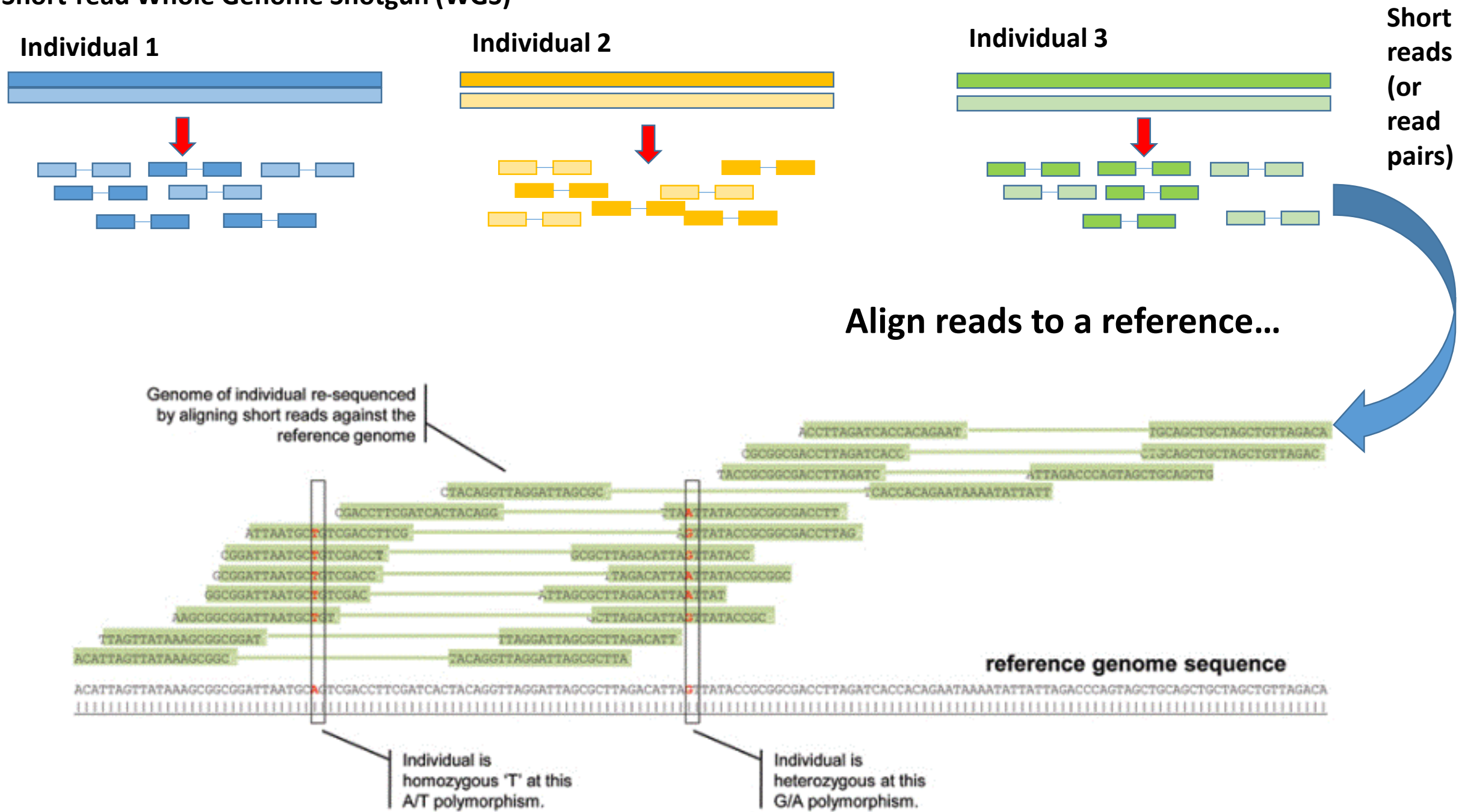
Alternatively, if you do not have the budget to sequence each individual genome?

(Bulked segregant analysis)

Pool them and sequence pools.



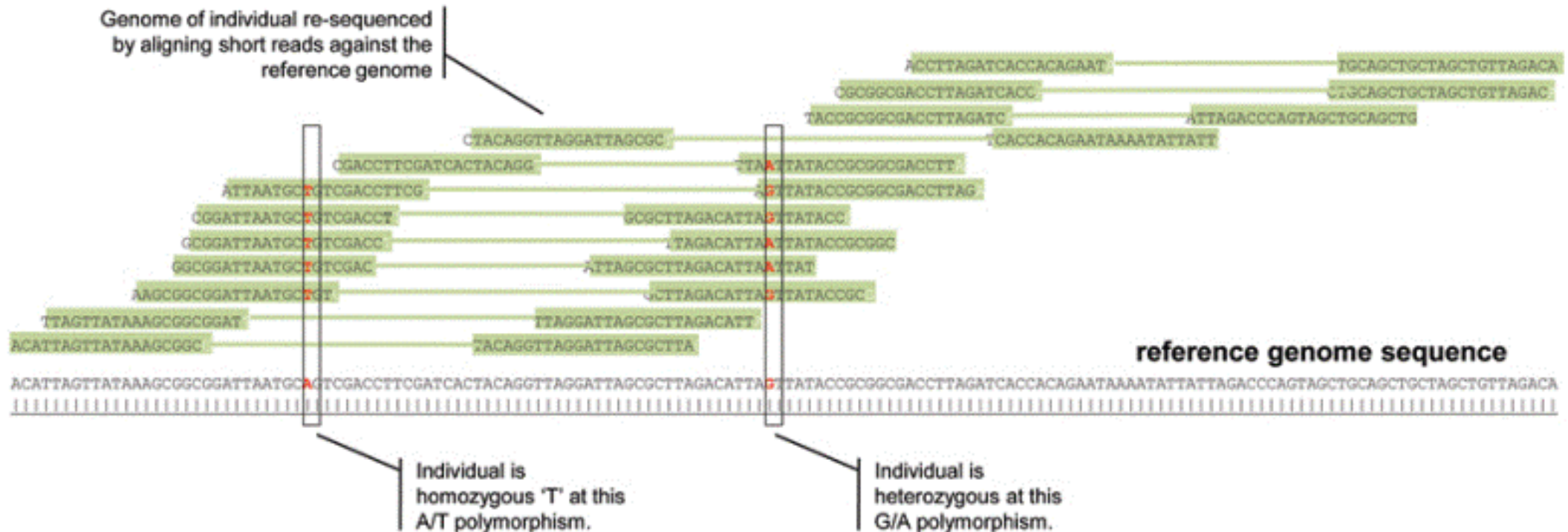
Short-read Whole Genome Shotgun (WGS)



Two major sources of errors:

- Alignment errors (reads are aligned to paralogous regions on the reference);
- Sequencing errors

Unless you are re-sequencing human genomes, majority of errors are alignment errors



Expected output: table of genotypes at variant sites


Variant site chr and position	Indiv1	Indiv2	Indiv3	...
site1	AA	AA	AC	...
site2	GT	missing	TT	...
...
siteN	CC	CC	AA	...

Table above is very schematic. In reality, genotypes are recorded in **VCF** format (**V**ariant **C**all **F**ormat)

Additional information about variants is also produced and recorded in VCF (such as call quality info)


More about VCF – coming soon

Commonly used tool: GATK from Broad Institute

 [User Guide](#) [Tool Index](#) [Blog](#) [Forum](#) [DRAGEN-GATK](#) [Events](#) [Download GATK4](#) [Sign in](#)

Genome Analysis Toolkit

Variant Discovery in High-Throughput Sequencing Data



Sequencing


READS

gatk best practices™

VARIANTS


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](#)

Find answers to your questions. Stay up to date on the latest topics. Ask questions and help others.




Getting Started

Best practices, tutorials, and other info to get you started




Technical Documentation

Algorithms, glossary, and other detailed resources




Announcements

Blog and events




Tool Index

Purpose, usage and options for each tool




Forum

Ask our team for help and report issues




GATK Showcase on Terra

Check out these fully configured workspaces




DRAGEN-GATK

Learn more about DRAGEN-GATK




Download latest version of GATK

The GATK package download includes all released GATK tools



Run on Cloud

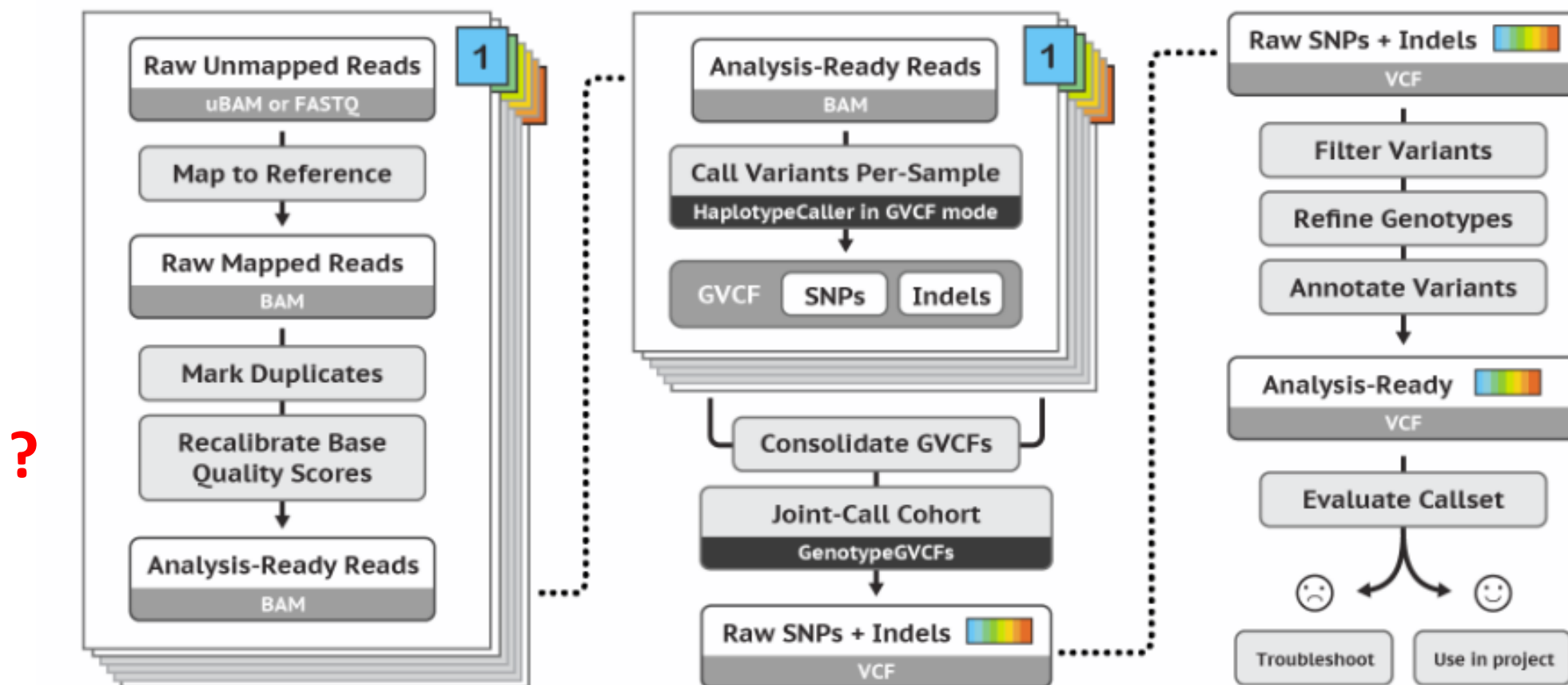


Run on HPC

“Best Practices” for DNA-Seq variant calling

Purpose

Identify germline short variants (SNPs and Indels) in one or more individuals to produce a joint callset in VCF format.



Best Practices for DNA-Seq variant calling

What are the colored tabs?

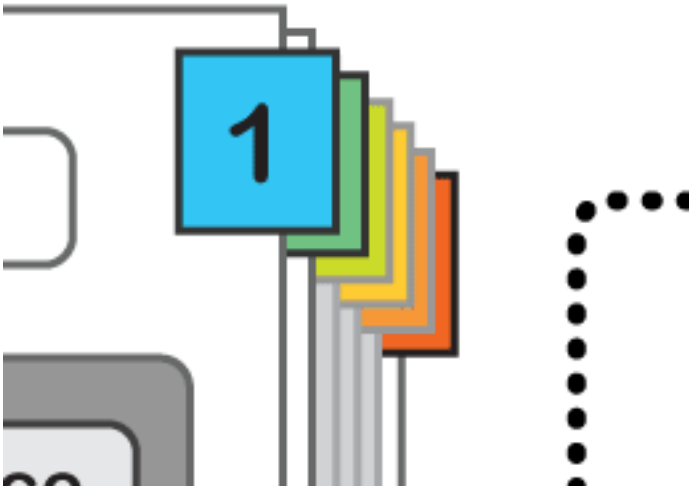
Each tab stands for a **FASTQ** file (SE case) or a **pair of FASTQ files** (PE case) with reads from **one sample, one Illumina lane, one library** (i.e. one read group)

A lane may contain reads from

- a single sample/library, OR...
- multiple samples/libraries (multiplexing)

Reads from one sample/library may initially be in

- One FASTQ file, OR.....
- Multiple FASTQ files



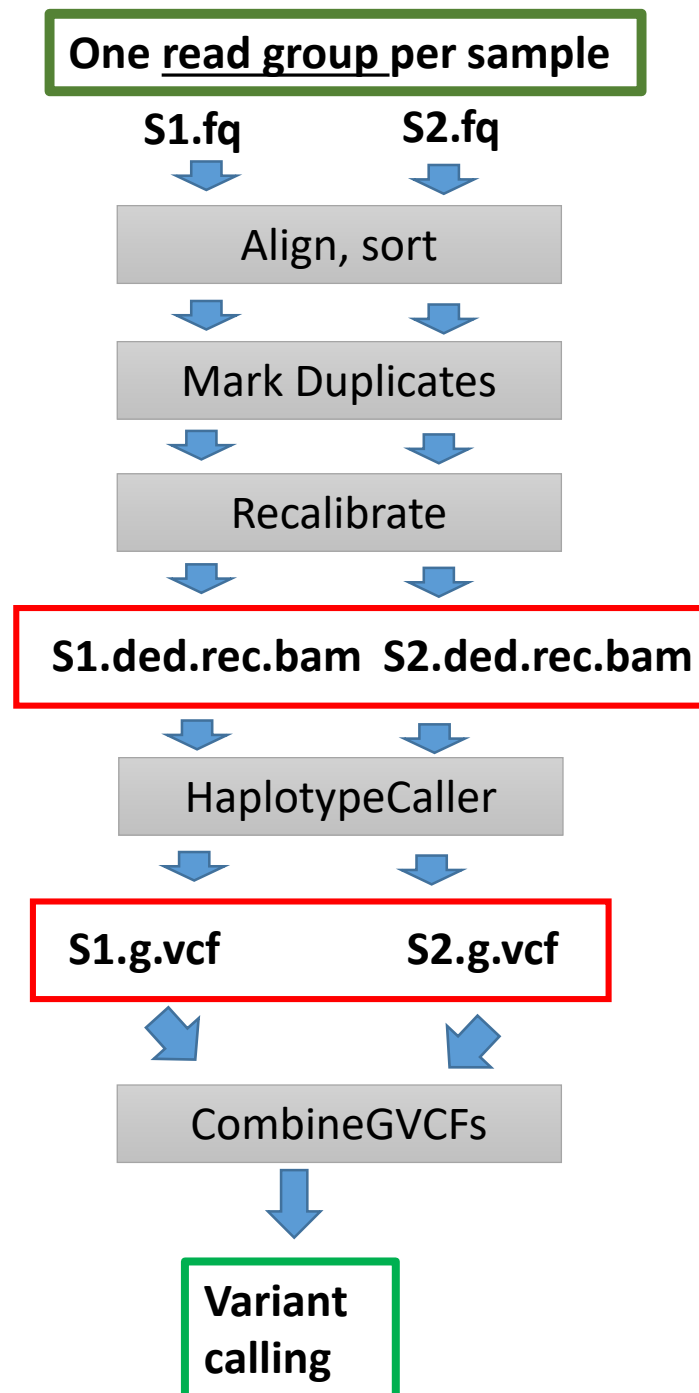
What's good for computational efficiency: process all datasets independently – in parallel

... is bad for accuracy: GATK tools prefer large datasets as possible - long compute times and loss of parallelism

- **Mark Duplicates** works best if given **all reads from a library** (sometimes scattered among files)
- **Haplotype calling** (discussed later) works best with all reads from a sample, and would be delighted to use all reads available (whole cohort)

Compromises have to be made

Pipeline in GATK



Input: paired-end (PE) reads

Paired-end case: we have two “parallel” FASTQ files – one for “left” and another for “right” end of the fragment:

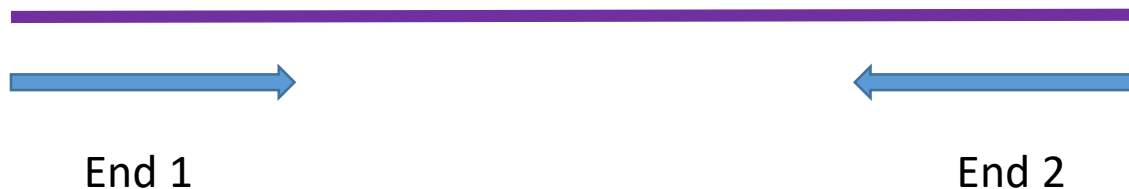
First sequence in “left” file

```
@HWI-ST896:156:D0JFYACXX:5:1101:1652:2132 1:N:0:GATCAG
ACTGCATCCTGGAAAGAATCAATGGTGGCCGGAAAGTGTTTTTCAAATACAAGAGTGACAATGTGCCCTGTTGTTT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@@CACCCCCACCCCCCCCCCCCCCCCCCCCCCCC
```

First sequence in “right” file

```
@HWI-ST896:156:D0JFYACXX:5:1101:1652:2132 2:N:0:GATCAG
CTCAAATGGTTAATTCTCAGGCTGCAAATATTCGTTTCAGGATGGAAGAACATTTTCTCAGTATTCCATCTAGCTGC
+
C<CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCACCCCACCC =
```

The two ends come from **opposite strands** of the fragment being sequenced

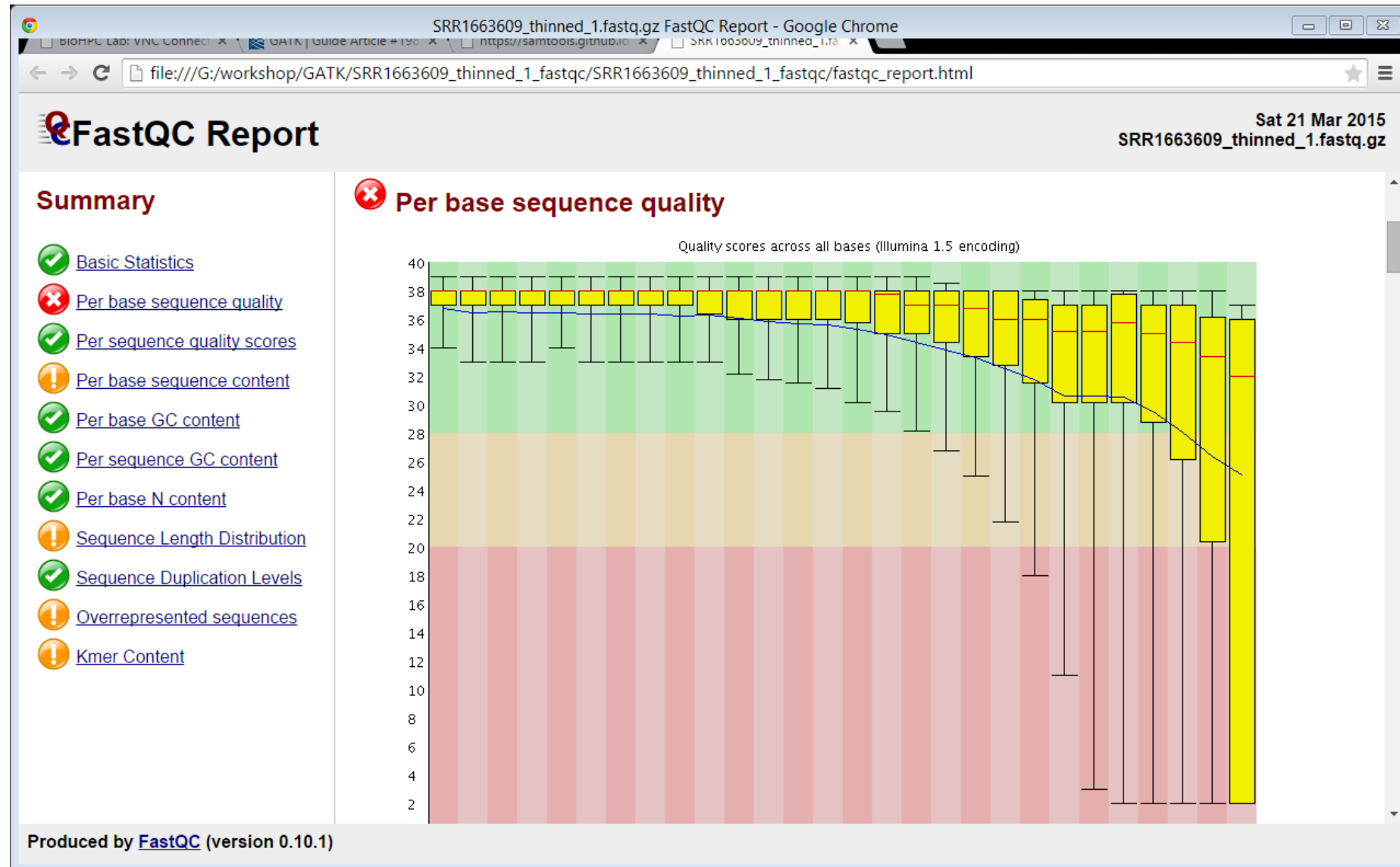


Phred base quality score

For example, “C” stands for: $67 - 33 = 34$, i.e., probability of the base (here: C) being miscalled is $10^{-3.4}$.

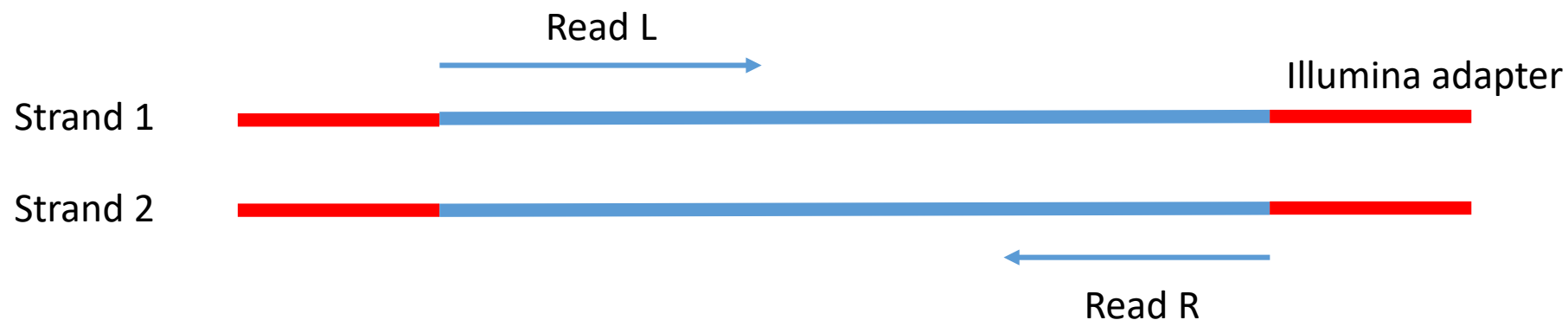
Base qualities are typically used in genotype likelihood models – they better be accurate!

Read quality assessment with fastqc

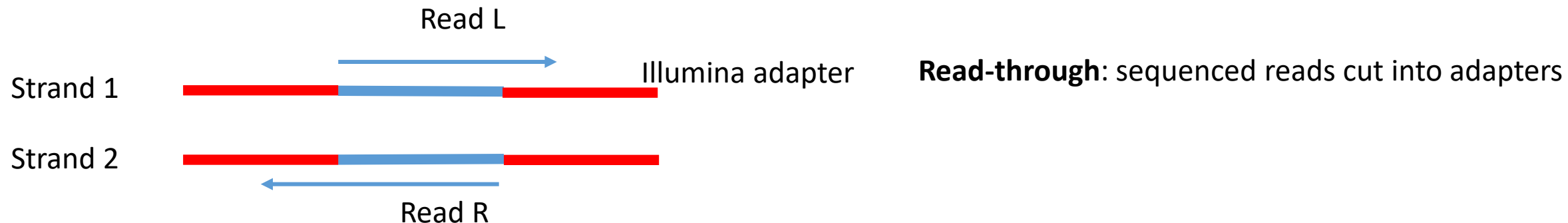


Run the command: `fastqc my_file.fastq.gz` to generate html report

Sequencing a long fragment



Sequencing a short fragment



Tool for removal of low-quality portions of reads and/or adapter sequences:

- **Trimmomatic** (<http://www.usadellab.org/cms/?page=trimmomatic>, <https://biohpc.cornell.edu/lab/userguide.aspx?a=software&i=53#c>)
- Adapter removal not that important in alignment-based methods

Alignment is fundamentally hard.....

- Genomes being re-sequenced not sufficiently similar to reference
 - Not enough reads will be mapped
 - Reads originating from parts of genome absent from reference will align somewhere anyway, leading to **false SNPs**
- Some reads cannot be mapped unambiguously in a single location (have low **Mapping Quality**)
 - if reads too short
 - reads originating from paralogs or repetitive regions
 - Having paired-end (PE) data helps
- Alignment of some reads may be ambiguous even if placement on reference correct (SNPs vs indels)
 - Need local multi-read re-alignment or local haplotype assembly (expensive!)
- Sequencing errors
 - Easier to handle and/or build into variant-calling models

Choosing a good aligner is important

Ambiguity of alignment at indel sites

Reference CTTTAGTTTCTTTT-----CTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC

CTTTAGTTTCTTTT-----GCCGCTTTCCTTCTTTCTT ←

CTTTAGTTTCTTTT-----GCCGCTTTCCTTCTTTCTT ←

Reads CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC

For these reads, aligner preferred to make a few SNPs rather than insertion

For these reads, insertion was a better choice

But we can try to shift things around a bit:

Reference CTTTAGTTTCTTTT-----CTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC

CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTT

CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTT

Reads CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC
CTTTAGTTTCTTTTGCCGCTTTCCTTCTTTCTTTTTTTTTTAAGTCTCCCTC

Aligner, like BWA, works on one read (fragment) at a time, does not see a bigger picture...)

This looks better !

Only seen after aligning all (at least some) reads!

Strategies to deal with indels

Local multiple-sequence re-alignment of all reads spanning an putative indel (GATK 3)
performed prior to variant calling
computationally expensive

Local read assembly into haplotypes (**HaplotypeCaller** in GATK 3, 4)
naturally gets rid of reads with sequencing errors
used along whole genome (not only indels)
computationally expensive
standard in modern variant calling pipelines

BWA mem – aligner of choice in GATK

- **BWA** = Burrows Wheeler Aligner (uses BW transform to compress data)
- **MEM** = Maximal Exact Match (how alignment “seeds” are chosen)
- **Performs local alignment** (rather than end-over-end)
 - Can clip ends of reads, if they do not match
 - Can split a read into pieces, mapping each separately (the best aligned piece is then the primary alignment)
- **Performs gapped alignment**
- **Utilizes PE reads** to improve mapping
- **Reports only one alignment** for each read
 - If ambiguous, one of the equivalent best locations is chosen at random
 - Ambiguously mapped reads are reported with low **Mapping Quality**
- Works well for reads 70bp to several Mbp
- Time scales linearly with the size of query sequence (at least for exact matches)
- Moderate memory requirement (few GB of RAM to hold reference genome)

Running BWA mem: align your reads

For PE reads:

```
bwa mem -M -t 4 \  
-R '@RG\tID:C6C0TANXX_2\tSM:ZW177\tLB:ZW177lib\tPL:ILLUMINA' \  
./genome_index/genome.fa \  
sample1reads_1.fastq.gz sample1reads_2.fastq.gz > sample1.sam
```

(SE version the same – just specify one read file instead of two)

What does it all mean:

- - M**: if a read is split (different parts map to different places) mark all parts other than main as “secondary alignment” (technicality, but important for GATK which ignores secondary alignments)
- –**R**: add **Read Group** description (more about it in a minute)
- –**t 4**: run of 4 CPU cores. If CPUs available, bwa mem scales well up to about 12 CPU cores.
- **./genome_index/genome.fa**: points to BWA index files (**genome.fa.***)
- Output (i.e., alignments) will be written to the file **sample1.sam**. As the name suggests, it will be in **SAM format**.

BWA mem command: define Read Group

```
-R '@RG\tID:C6C0TANXX_2\tSM:ZW177\tLB:ZW177lib\tPL:ILLUMINA'
```

What will this option do?

The **SAM/BAM file header** will contain a line (TAB-delimited) defining the group:

@RG	ID:C6C0TANXX_2	SM:ZW177	LB:ZW177lib	PL:ILLUMINA
	Unique ID of a collection of reads sequenced together, typically: Illumina lane +(barcode or sample)+library	Sample name	DNA prep Libray ID	Sequencing platform

Each alignment record will be marked with **Read Group ID** (here: **C6C0TANXX_2**), so that programs in downstream analysis know where the read is from.

Read groups, sample and library IDs are important for GATK operation!

Each **READ GROUP** contains reads from **one sample, one library, one flowcell_lane**
A library may be sequenced multiple times (on different flowcell_lanes)
Sample may be sequenced multiple times, on different lanes and from different libraries

Dad's data:

@RG	ID:FLOWCELL1.LANE1	PL:ILLUMINA	LB:LIB-DAD-1	SM:DAD	PI:200
@RG	ID:FLOWCELL1.LANE2	PL:ILLUMINA	LB:LIB-DAD-1	SM:DAD	PI:200
@RG	ID:FLOWCELL1.LANE3	PL:ILLUMINA	LB:LIB-DAD-2	SM:DAD	PI:400
@RG	ID:FLOWCELL1.LANE4	PL:ILLUMINA	LB:LIB-DAD-2	SM:DAD	PI:400

Mom's data:

@RG	ID:FLOWCELL1.LANE5	PL:ILLUMINA	LB:LIB-MOM-1	SM:MOM	PI:200
@RG	ID:FLOWCELL1.LANE6	PL:ILLUMINA	LB:LIB-MOM-1	SM:MOM	PI:200
@RG	ID:FLOWCELL1.LANE7	PL:ILLUMINA	LB:LIB-MOM-2	SM:MOM	PI:400
@RG	ID:FLOWCELL1.LANE8	PL:ILLUMINA	LB:LIB-MOM-2	SM:MOM	PI:400

Kid's data:

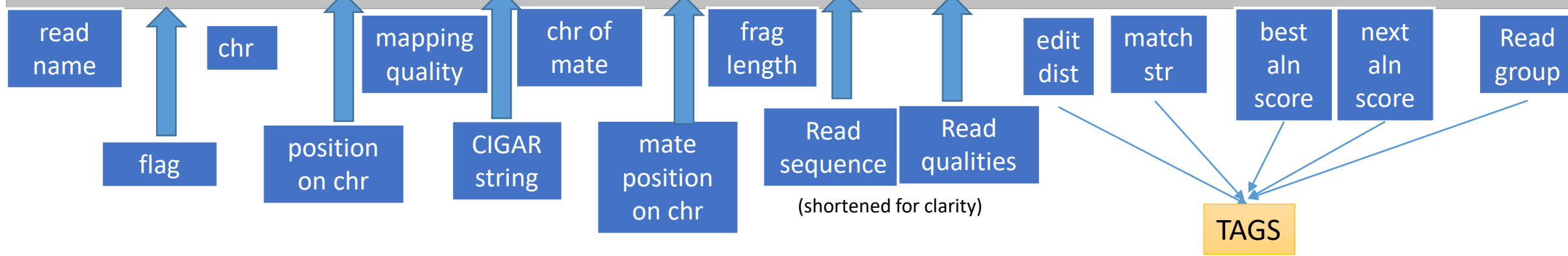
@RG	ID:FLOWCELL2.LANE1	PL:ILLUMINA	LB:LIB-KID-1	SM:KID	PI:200
@RG	ID:FLOWCELL2.LANE2	PL:ILLUMINA	LB:LIB-KID-1	SM:KID	PI:200
@RG	ID:FLOWCELL2.LANE3	PL:ILLUMINA	LB:LIB-KID-2	SM:KID	PI:400
@RG	ID:FLOWCELL2.LANE4	PL:ILLUMINA	LB:LIB-KID-2	SM:KID	PI:400

Anatomy of a SAM file

```

@SQ      SN:chr2L      LN:23011544
@SQ      SN:chr2LHet   LN:368872
@SQ      SN:chr2R      LN:21146708
@SQ      SN:chr2RHet   LN:3288761
@SQ      SN:chr3L      LN:24543557
@SQ      SN:chr3LHet   LN:2555491
@SQ      SN:chr3R      LN:27905053
@SQ      SN:chr3RHet   LN:2517507
@SQ      SN:chr4      LN:1351857
@SQ      SN:chrM      LN:19517
@SQ      SN:chrX      LN:22422827
@SQ      SN:chrXHet    LN:204112
@SQ      SN:chrYHet    LN:347038
@RG      ID:SRR1663609 SM:ZW177      LB:ZW155      PL:ILLUMINA
@PG      ID:bwa      PN:bwa      VN:0.7.8-r455      CL:bwa mem -M -t 4 -R @RG\tID:SRR1663609\tSM:ZW177\tLB:ZW155\tPL:ILLUMINA
/local_data/Drosophila_melanogaster_dm3/BWAIndex/genome.fa SRR1663609_1.fastq.gz SRR1663609_2
.fastq.gz
SRR1663609.1  97      chrX      2051224  60      6M54S      chrYHet    4586  0      GGATCGTGAT... gggfgg[gfg... NM:i:0 MD:Z:46 AS:i:46 XS:i:0 RG:Z:SRR1663609
SRR1663609.1  145     chrYHet    4586    0      100M      chrX      2051224  0      ACTTCTCTTC... BBBBbddd]c... NM:i:0 MD:Z:100 AS:i:100 XS:i:99 RG:Z:SRR1663609
SRR1663609.2  65      chr3RHet   2308288  0      100M      chrYHet    4712  0      AGAAGAGAAG... Y_b`_ccTccB... NM:i:0 MD:Z:100 AS:i:100 XS:i:100 RG:Z:SRR1663609
SRR1663609.2  129     chrYHet    4712    60      38M62S      chr3RHet   2308288  0      CTTCTCTTCT... eeeae`edee... NM:i:1 MD:Z:17T20 AS:i:33 XS:i:21 RG:Z:SRR1663609
SRR1663609.3  65      chr3RHet   2308278  0      100M      chrYHet    4649  0      AGAAGAGAAG... ffffffff... NM:i:0 MD:Z:100 AS:i:100 XS:i:100 RG:Z:SRR1663609
SRR1663609.3  129     chrYHet    4649    0      41M59S      chr3RHet   2308278  0      TCTCTTCTCT... ffffffff... NM:i:0 MD:Z:41 AS:i:41 XS:i:41 RG:Z:SRR1663609
SA:Z:chrX,5036484,-,16S41M43S,0,2;
SRR1663609.3  401     chrX      5036484  0      16H41M43H chr3RHet   2308278  0      AAAAGAAGAA... BBBBbBBBBB... NM:i:2 MD:Z:7A4G28 AS:i:31 XS:i:28 RG:Z:SRR1663609
SA:Z:chrYHet,4649,+,41M59S,0,0;
SRR1663609.4  99      chr3RHet   854491   0      100M      =      854876  485     AGAAGAAGAA... BBBBbBBBBB... NM:i:0 MD:Z:100 AS:i:100 XS:i:100 RG:Z:SRR1663609
SRR1663609.4  147     chr3RHet   854876   0      100M      =      854491  -485    GAGAAGAGAA... ffffffff... NM:i:0 MD:Z:100 AS:i:100 XS:i:100 RG:Z:SRR1663609

```



Looking into a BAM file: samtools

BAM files are binary – special tool is needed to look inside

Examples:

samtools view -h myfile.bam | more
prints the file in SAM format (i.e., human-readable) to screen page by page; skip -h to omit header lines

samtools view -c myfile.bam
prints the number of records (alignments) in the file; for BWA mem it may be larger than the number of reads!

samtools view -f 4 myfile.bam
Extracts records with a given flag – here: flag 4 (unmapped); prints them to screen

Type **samtools**, or go to <http://samtools.sourceforge.net/> for more options

samtools flagstat myfile.bam

Displays basic alignment stats based on flag

```
samtools flagstat
SRR1663609.sorted.dedup.realigned.fixmate.bam
10201772 + 0 in total (QC-passed reads + QC-failed reads)
74334 + 0 secondary
0 + 0 supplementary
679571 + 0 duplicates
9685912 + 0 mapped (94.94%:-nan%)
10127438 + 0 paired in sequencing
5063719 + 0 read1
5063719 + 0 read2
8747736 + 0 properly paired (86.38%:-nan%)
9500218 + 0 with itself and mate mapped
111360 + 0 singletons (1.10%:-nan%)
252790 + 0 with mate mapped to a different chr
89859 + 0 with mate mapped to a different chr (mapQ>=5)
```

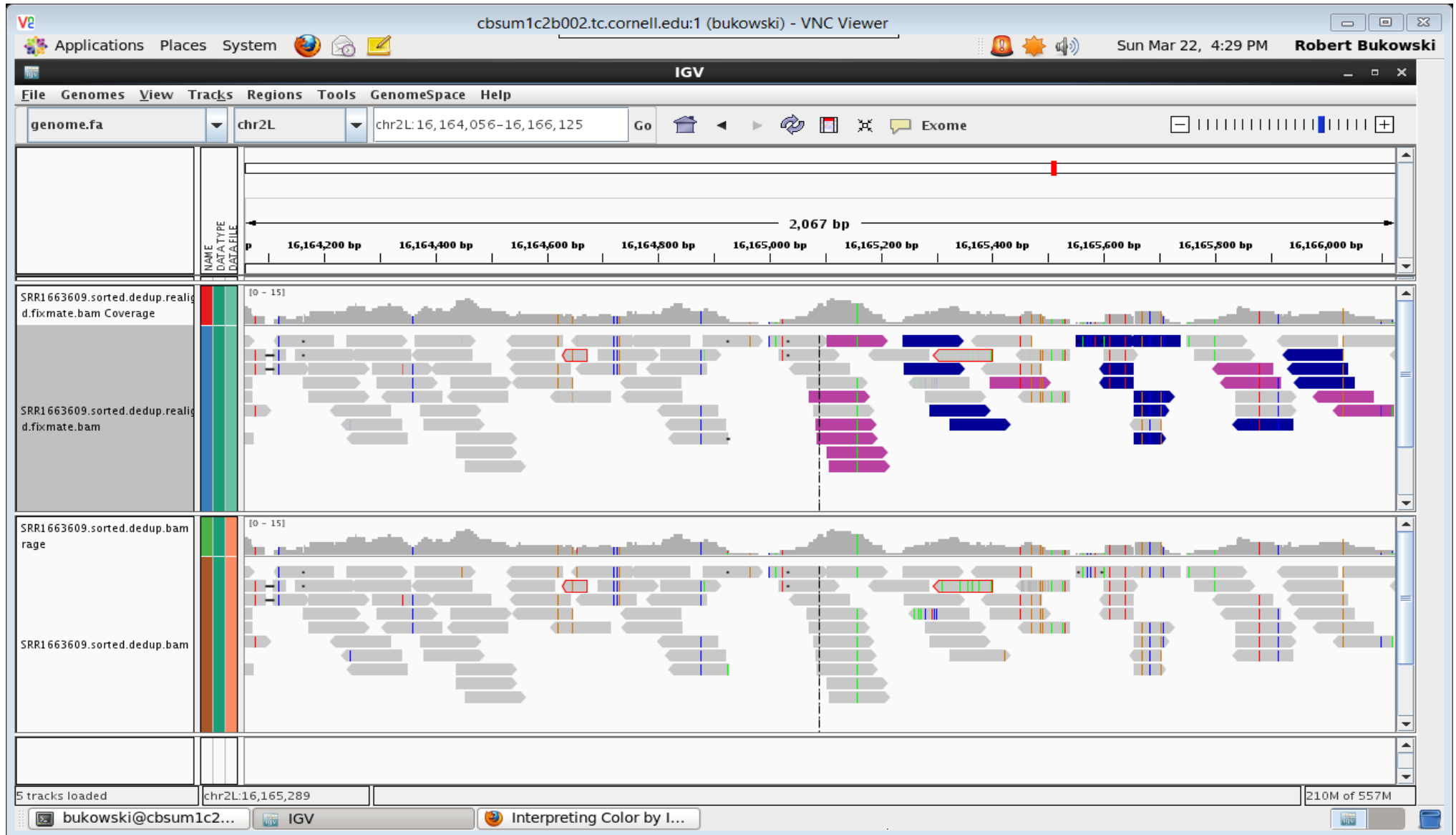
Looking into a BAM file: IGV viewer

Look at multiple
BAM files

Zoom in and out

Various color-
coding schemes

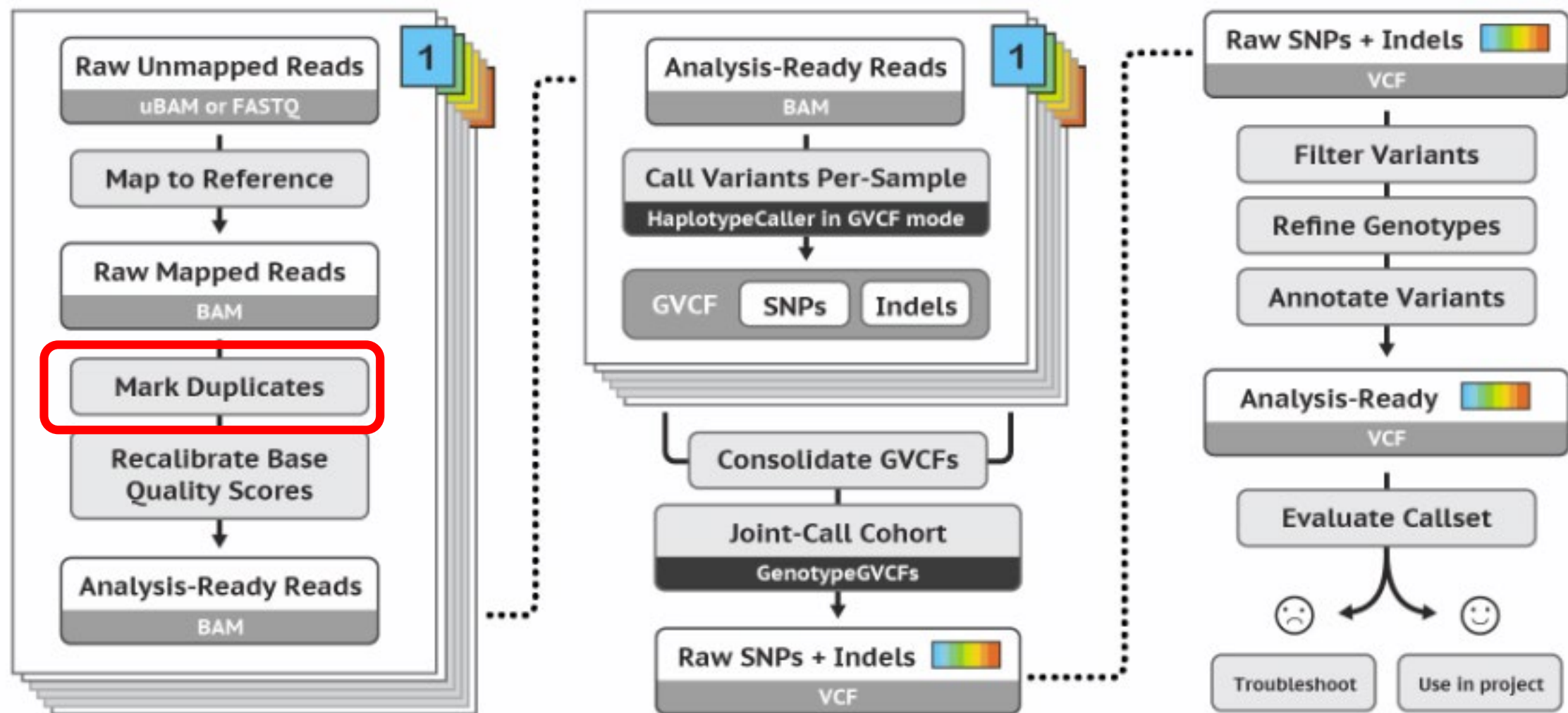
Can load
genome
annotation track



IGV is a Java program available on BioHPC machines. Can be installed on laptop, too.

<http://www.broadinstitute.org/igv/home>

“Best Practices” for DNA-Seq variant calling

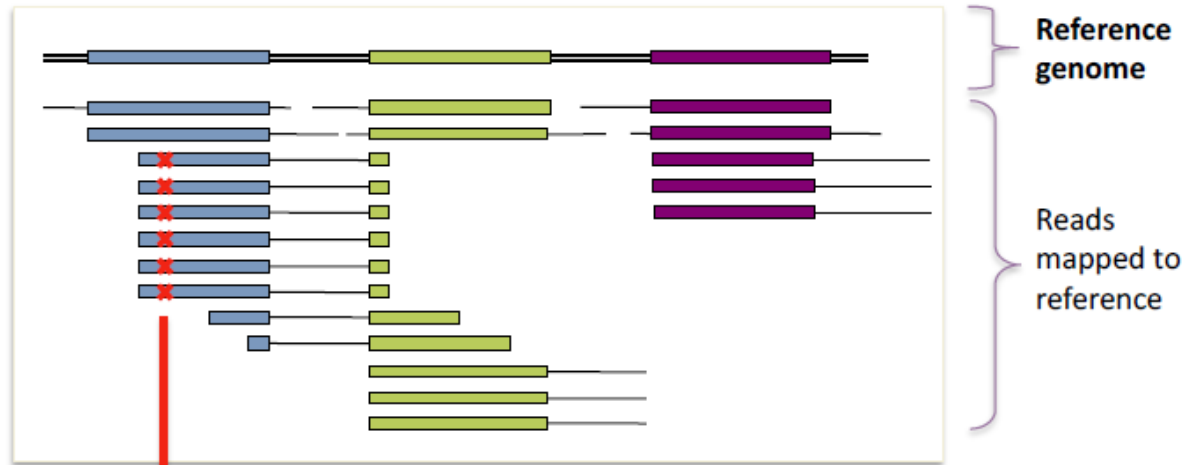


Duplicate reads (fragments)

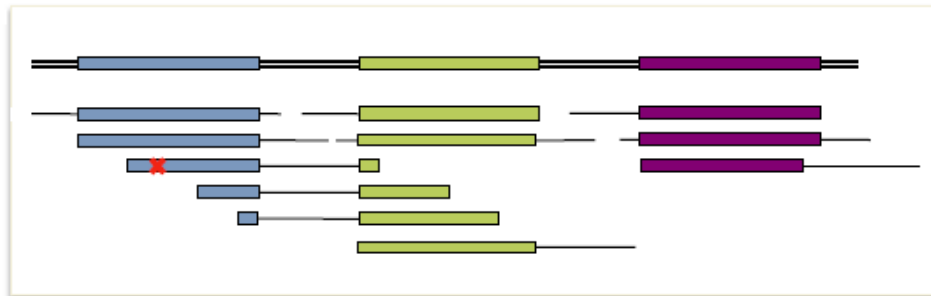
- **Optical duplicates**: (Illumina) generated when a single cluster of reads is part of two adjacent tiles' on the same slide and used to compute two read calls separately
 - Very similar in sequence (except sequencing errors).
 - Identified where the first, say, 50 bases are identical between two reads and the read's coordinates are close
- **Library duplicates (aka PCR duplicates)**: generated when the original sample is pre-amplified to such extent that initial unique targets are PCR replicated prior to library preparation and will lead to several independent spots on the Illumina slide.
 - do not have to be adjacent on the slide
 - share a very high level of sequence identity
 - align to the same place on reference
 - identified from alignment to reference

Why duplicates are bad for variant calling

✖ = sequencing error propagated in duplicates

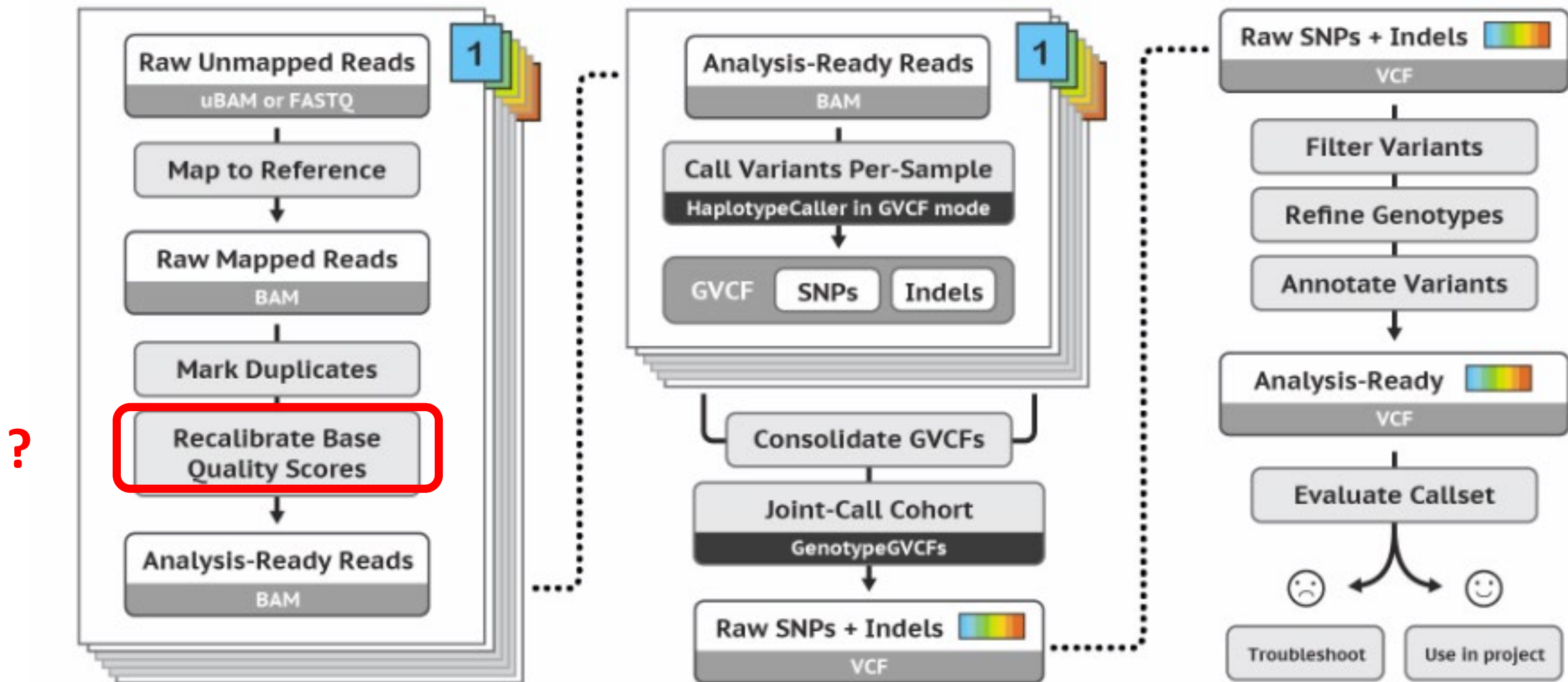


After marking duplicates, the GATK will only see :



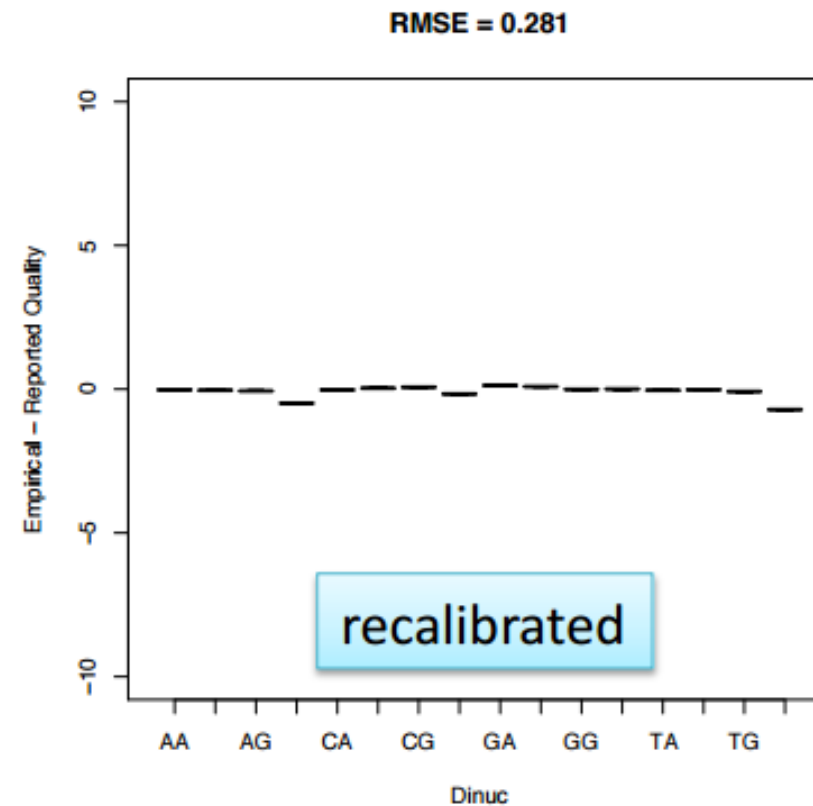
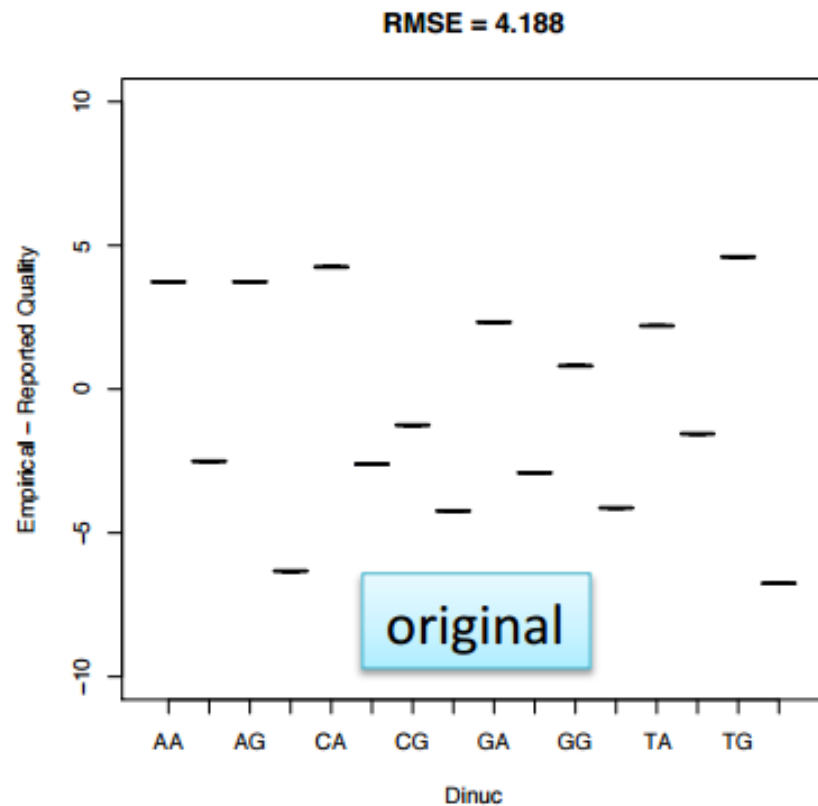
... and thus be more likely to make the right call

“Best Practices” for DNA-Seq variant calling



Base quality scores reported by a sequencer may be inaccurate and biased

Example: Bias in the qualities reported depending on nucleotide context



Base quality score recalibration: good or bad?

Implicit assumption behind recalibration: sequencing error rate higher than SNP rate

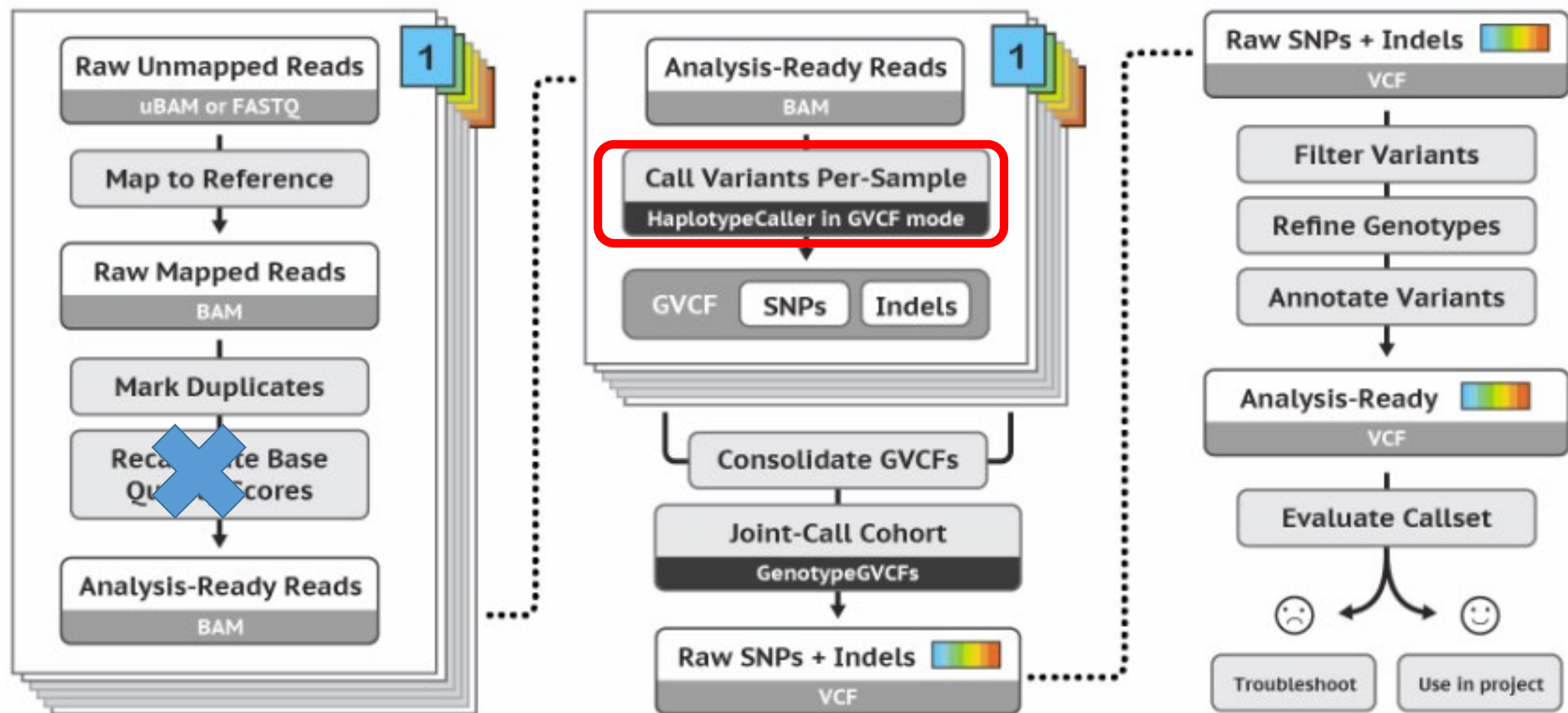
applicable only to populations with very little diversity (humans)

in most (non-human) cases, 'empirical errors' are not sequencing errors (either real variants or misalignments)

'known variants' not always available

Conclusion: do not recalibrate (unless dealing with human genomes)

“Best Practices” for DNA-Seq variant calling



SNP and Indel calling is a large-scale Bayesian modeling problem

Bayesian model

$$\Pr\{G|D\} = \frac{\Pr\{G\} \Pr\{D|G\}}{\sum_i \Pr\{G_i\} \Pr\{D|G_i\}}, \text{ [Bayes' rule]}$$

$\Pr\{G\}$ Prior of the genotype
 $\Pr\{D|G\}$ Likelihood of the genotype
 Reported as PL in our VCF example

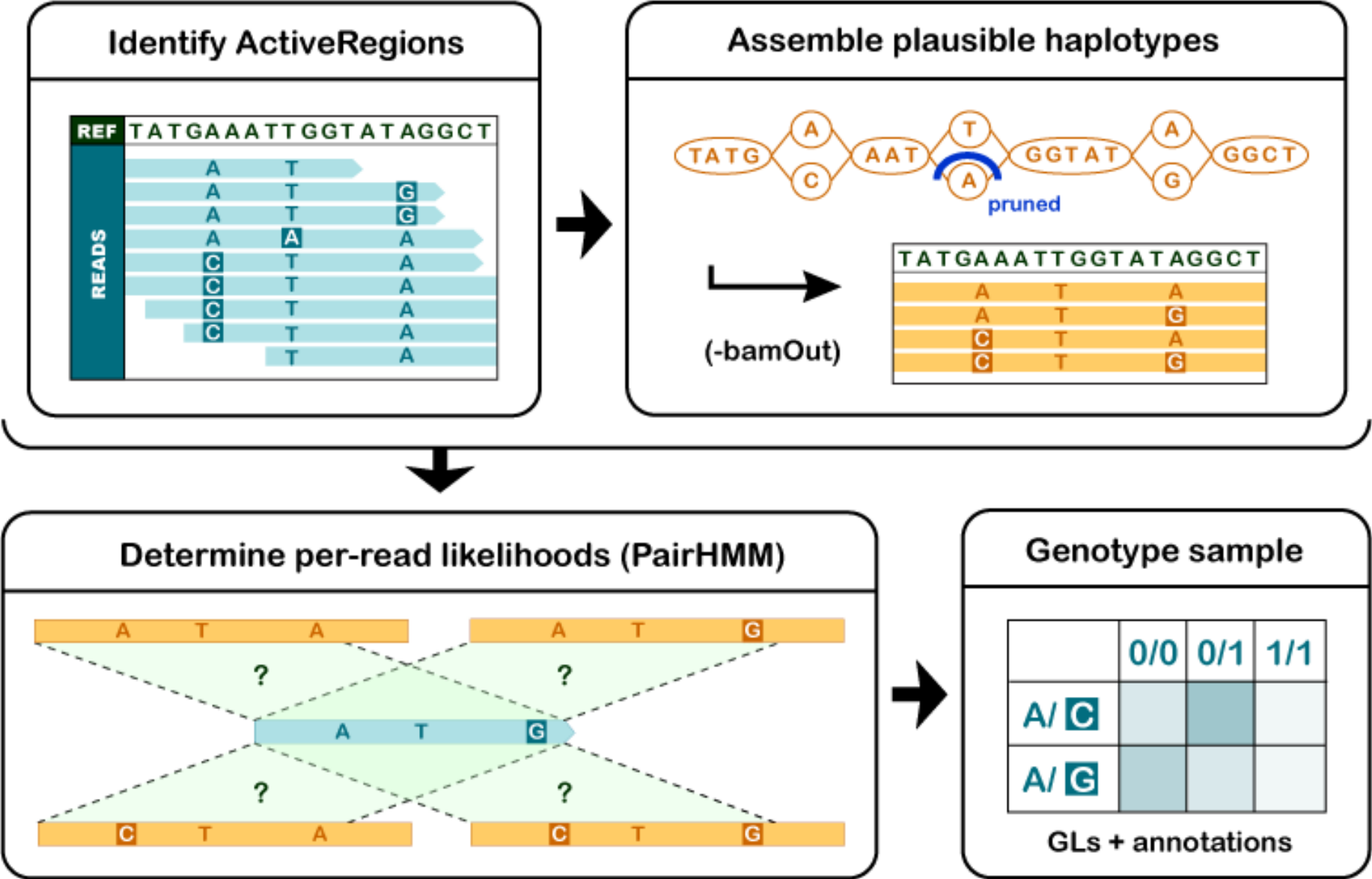
$$\Pr\{D|G\} = \prod_j \left(\frac{\Pr\{D_j|H_1\}}{2} + \frac{\Pr\{D_j|H_2\}}{2} \right) \text{ where } G = H_1 H_2$$

Diploid assumption

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

- Inference: what is the genotype G of each sample given read data D for each sample?
- Calculate via Bayes' rule the probability of each possible G
- Product expansion assumes reads are independent
- Relies on a likelihood function to estimate probability of sample data given proposed haplotype

HaplotypeCaller (in GVCF mode): extract variation from alignments for each sample



$P\{D_j | H\}$ determined from PairHMM scores of reads alignments to haplotypes (based on base qualities)

Genomic Variant Call (g.vcf) file: result of HaplotypeCaller

```
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA12878
20 10001567 . A <NON_REF> . . END=10001616 GT:DP:GQ:MIN_DP:PL 0/0:38:99:34:0,101,1114
20 10001617 . C A,<NON_REF> 493.77 .
BaseQRankSum=1.632;ClippingRankSum=0.000;DP=38;ExcessHet=3.0103;MLEAC=1,0;MLEAF=0.500,0.00;MQRankSum=0.000;RAW_MQ=136800.00;ReadPosRankSum=0.170 GT:AD:DP:GQ:PL:SB 0/1:19,19,0:38:99:522,0,480,578,538,1116:11,8,13,6
20 10001618 . T <NON_REF> . . END=10001627 GT:DP:GQ:MIN_DP:PL 0/0:39:99:37:0,105,1575
20 10001628 . G A,<NON_REF> 1223.77 .
DP=37;ExcessHet=3.0103;MLEAC=2,0;MLEAF=1.00,0.00;RAW_MQ=133200.00 GT:AD:DP:GQ:PL:SB
1/1:0,37,0:37:99:1252,111,0,1252,111,1252:0,0,21,16
20 10001629 . G <NON_REF> . . END=10001660 GT:DP:GQ:MIN_DP:PL 0/0:43:99:38:0,102,1219
20 10001661 . T C,<NON_REF> 1779.77 .
DP=42;ExcessHet=3.0103;MLEAC=2,0;MLEAF=1.00,0.00;RAW_MQ=151200.00 GT:AD:DP:GQ:PL:SB
1/1:0,42,0:42:99:1808,129,0,1808,129,1808:0,0,26,16
20 10001662 . T <NON_REF> . . END=10001669 GT:DP:GQ:MIN_DP:PL 0/0:44:99:43:0,117,1755
20 10001670 . T G,<NON_REF> 1773.77 .
DP=42;ExcessHet=3.0103;MLEAC=2,0;MLEAF=1.00,0.00;RAW_MQ=151200.00 GT:AD:DP:GQ:PL:SB
1/1:0,42,0:42:99:1802,129,0,1802,129,1802:0,0,25,17
20 10001671 . G <NON_REF> . . END=10001673 GT:DP:GQ:MIN_DP:PL 0/0:43:99:42:0,120,180
20 10001674 . A <NON_REF> . . END=10001674 GT:DP:GQ:MIN_DP:PL 0/0:42:96:42:0,96,1197
```

Positional information: chromosome, start, end (if non-variant block)

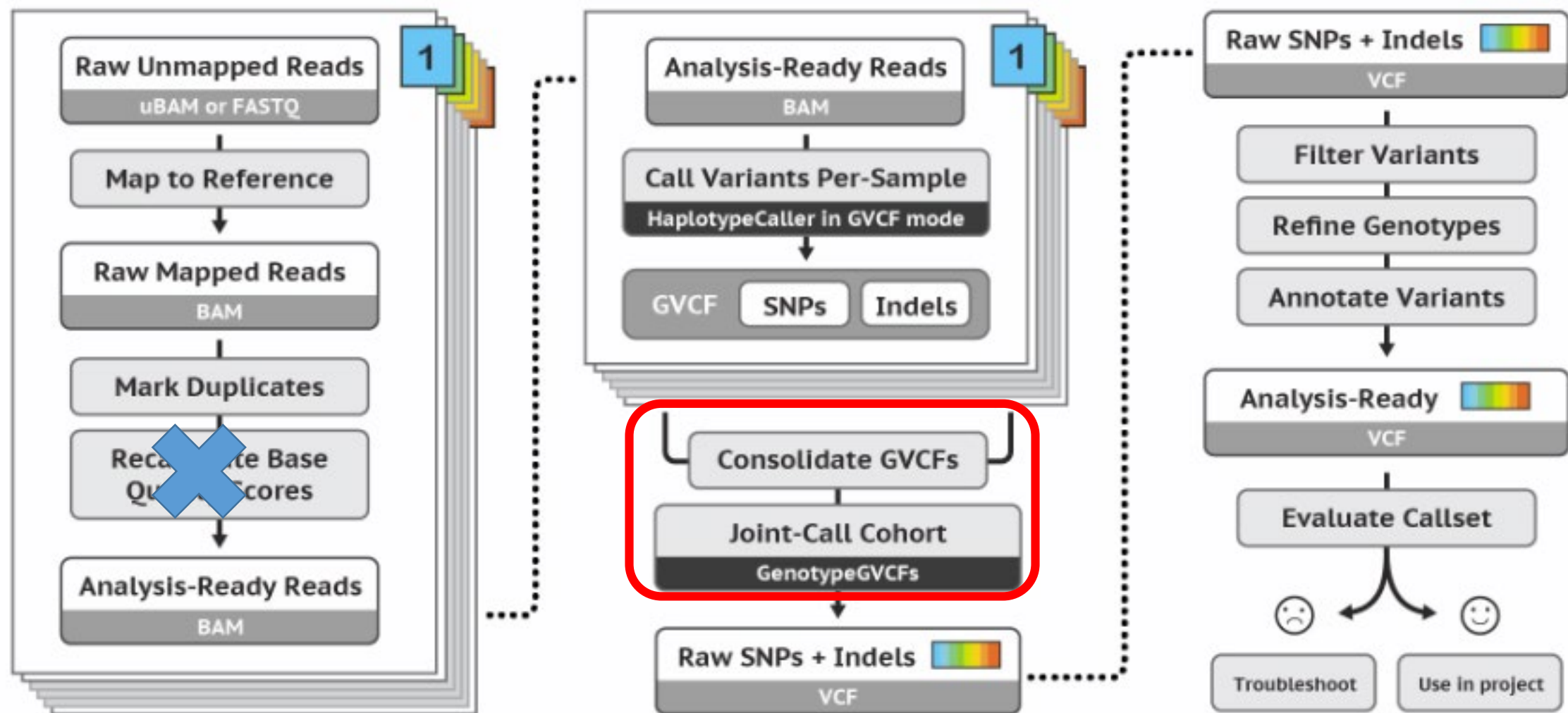
Non-reference allele info; <NON_REF> stands for any non-reference allele

Genotype (GT): may be 0/0, 0/1, 1/1, 0/2, 1/2, 2/2, ... where '0' is REF allele and 1, 2, ... are ALT alleles in order listed

Genotype likelihoods (PL): example: 0,120,180 means that 0/1 is 10^{-12} times less likely than 0/0

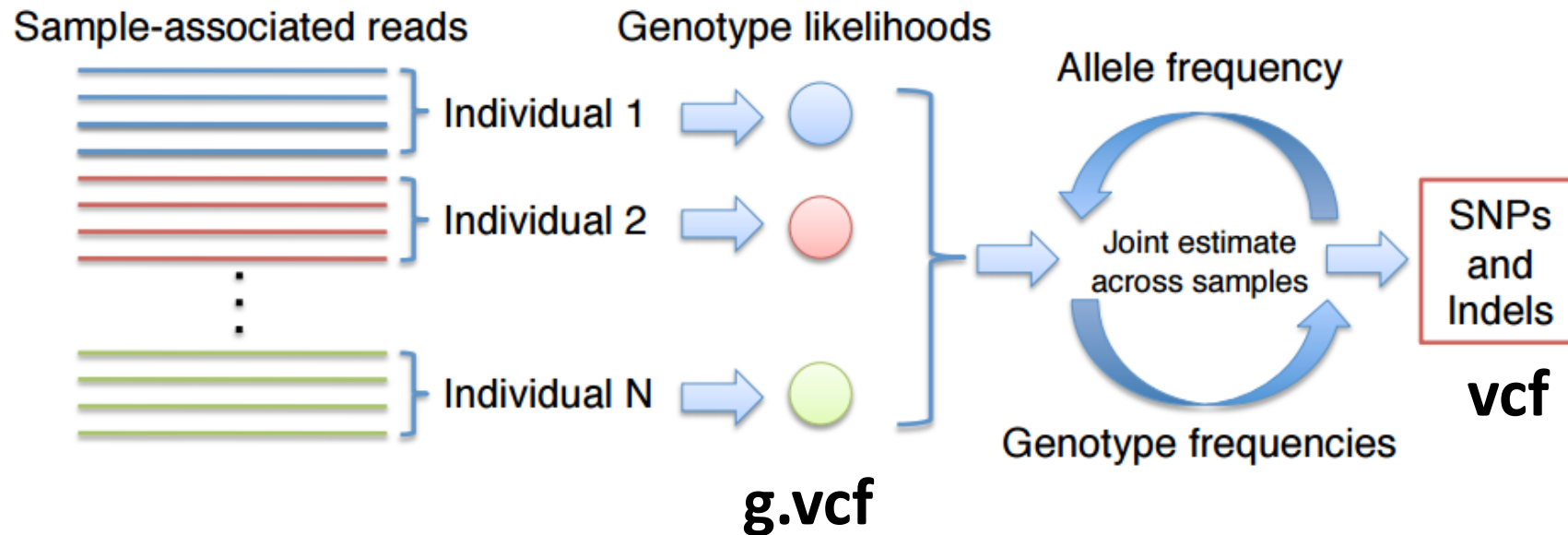
All symbols defined in the **header** of the g.vcf file (e.g., entries in INFO field for variant sites)

“Best Practices” for DNA-Seq variant calling



Variation across cohort

Multi-sample calling integrates per sample likelihoods to jointly estimate allele frequency of variation



For each site, obtain distribution of count of non-reference allele (**AC**):

$$\Pr\{AC=i \mid D\} \leftarrow \text{Per sample Genotype Likelihoods} + \text{Prior}$$

Prior: $\Pr\{AC=i\} = \text{Het}/i$ (where Het is population heterozygosity; or define your own prior)

QUAL = $-10 * \log \Pr\{AC=0 \mid D\}$ (reported in VCF file)

From BAM files to population variants

Run for each sample (on a multi-CPU machine, run a few simultaneously)

gatk HaplotypeCaller
In gVCF mode
(1 sample calls,
preferably in parallel)



sample1.g.vcf
sample2.g.vcf
...
sampleN.g.vcf



allsample.g.vcf



gatk GenotypeGVCFs
(joint variant calling)



N-sample VCF file

```
gatk HaplotypeCaller \  
  -R genome.fa \  
  -I sample1.sorted.dedup.recal.bam \  
  --minimum-mapping-quality 30 \  
  -ERC GVCF \  
  -O sample1.g.vcf
```

Slow

....Followed by combining g.vcf files

```
gatk CombineGVCFs \  
  -R genome.fa \  
  --variant sample1.g.vcf \  
  .....  
  --variant sampleN.g.vcf  
  -O allsample.g.vcf
```

Fast

....and by joint variant calling with GenotypeGVCFs

```
gatk GenotypeGVCFs \  
  -R genome.fa \  
  -V allsample.g.vcf \  
  -stand-call-conf 5 \  
  -O allsample.vcf
```

Fast

Variant Call Format (VCF)

Similar to **g.vcf**, but used to describe sites deemed **variant** across a cohort

HEADER LINES: start with “##”, describe all symbols found later on in FORMAT and ANNOTATIONS, e.g.,

```
##fileformat=VCFv4.1
##FORMAT=<ID=AD,Number=.,Type=Integer,Description="Allelic depths for the ref and alt alleles in the order listed">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Approximate read depth (reads with MQ=255 or with bad mates are filtered)">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
```

SITE RECORDS:

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	ZW155	ZW177
chr2R	2926	.	C	A	345.03	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:4,9:13:80:216,0,80	0/0:6,0:6:18:0,18,166
chr2R	9862	.	TA	T	180.73	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:0,5:5:15:97,15,0	1/1:0,4:4:12:80,12,0
chr2R	10834	.	A	ACTG	173.04	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/0:14,0:14:33:0,33,495	0/1:6,3:9:99:105,0,315

ID: some ID for the variant, if known (e.g., dbSNP)

REF, ALT: reference and alternative alleles (on forward strand of reference)

QUAL = $-10 \cdot \log(1-p)$, where p is the probability of variant being present given the read data

FILTER: whether the variant failed a filter (filters defined by the user or program processing the file)

Variant Call Format (VCF)

[HEADER LINES]										
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	ZW155	ZW177
chr2R	2926	.	C	A	345.03	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/1:4,9:13:80:216,0,80	0/0:6,0:6:18:0,18,166
chr2R	9862	.	TA	T	180.73	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	1/1:0,5:5:15:97,15,0	1/1:0,4:4:12:80,12,0
chr2R	10834	.	A	ACTG	173.04	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	0/0:14,0:14:33:0,33,495	./.

GT (genotype):

- 0/0 reference homozygote
- 0/1 reference-alternative heterozygote
- 1/1 alternative homozygote
- 0/2, 1/2, 2/2, etc. - possible if more than one alternative allele present
- ./. missing data

AD: allele depths

DP: total depth (may be different from sum of AD depths, as the latter include only reads significantly supporting alleles)

PL: genotype likelihoods (phred-scaled), normalized to the best genotype, e.g.,
$$PL(0/1) = -10 \cdot \log[\text{Prob}(\text{data} | 0/1) / \text{Prob}(\text{data} | \text{best_genotype})]$$

GQ: genotype quality – this is just PL of the second-best genotype

Variant Call Format (VCF)

[HEADER LINES]									
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	
chr2R	2926	.	C	A	345.03	PASS	[ANNOTATIONS]	GT:AD:DP:GQ:PL	ZW155
chr2R	9862	.	TA	T	180.73	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	ZW177
chr2R	10834	.	A	ACTG	173.04	.	[ANNOTATIONS]	GT:AD:DP:GQ:PL	

[ANNOTATIONS]: all kinds of quantities and flags that characterize the variant; supplied by the variant caller (different callers will do it differently)

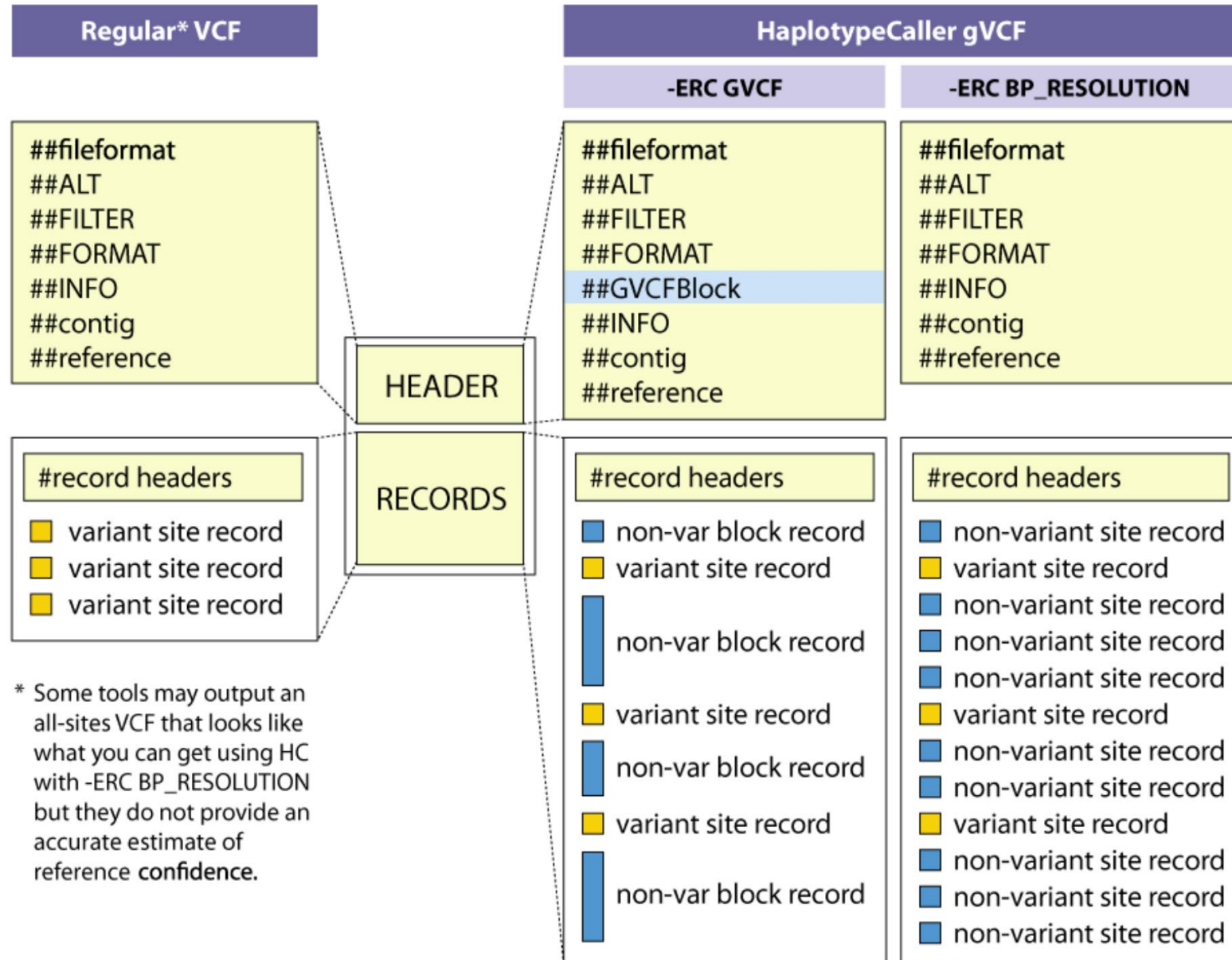
Example:

```
AC=2;AF=0.333;AN=6;DP=16;FS=0.000;GQ_MEAN=16.00;GQ_STDDEV=10.54;MLEAC=2;MLEAF=0.333;MQ=25.00;MQ0=0;NCC=1;QD=22.51;SOR=3.611
```

All ANNOTATION parameters are defined in the **HEADER LINES** on top of the file

```
...
##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">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Approximate read depth; some reads may have been filtered">
##INFO=<ID=FS,Number=1,Type=Float,Description="Phred-scaled p-value using Fisher's exact test to detect strand bias">
##INFO=<ID=GQ_MEAN,Number=1,Type=Float,Description="Mean of all GQ values">
##INFO=<ID=MQ,Number=1,Type=Float,Description="RMS Mapping Quality">
##INFO=<ID=NCC,Number=1,Type=Integer,Description="Number of no-called samples">
##INFO=<ID=QD,Number=1,Type=Float,Description="Variant Confidence/Quality by Depth">
##INFO=<ID=SOR,Number=1,Type=Float,Description="Symmetric Odds Ratio of 2x2 contingency table to detect strand bias">
...
```

VCF versus GVCF format



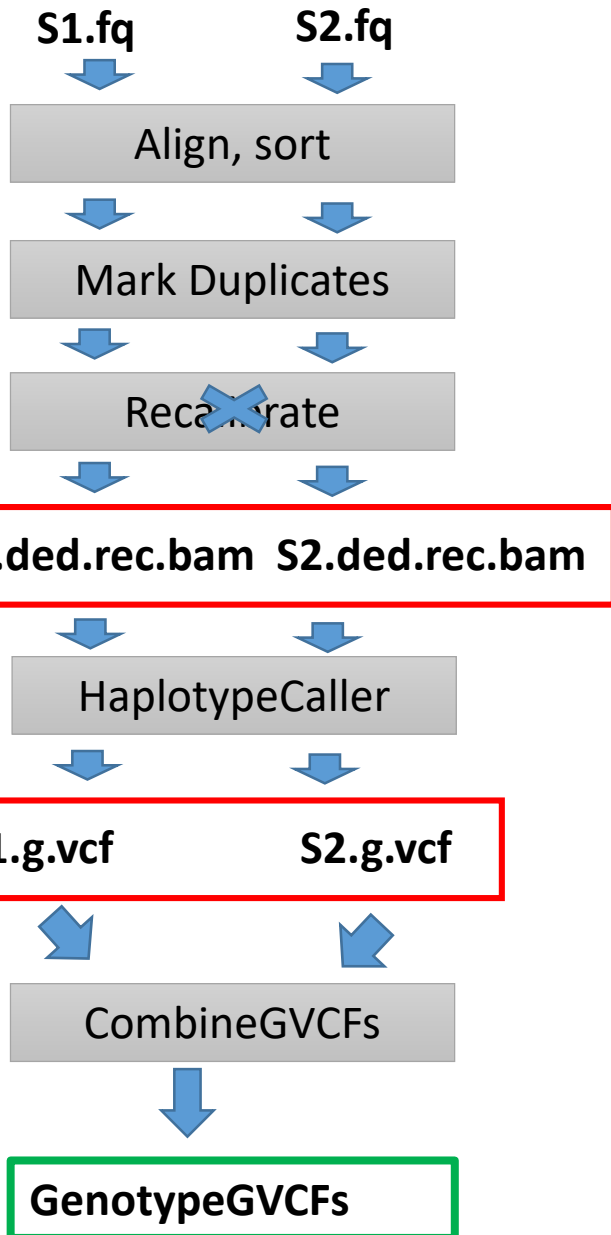
VCF, BCF, and BGZF compressed VCF

VCF (.vcf): text file

BCF (.bcf): binary version of VCF

BGZF VCF (.vcf.gz): vcf compressed with gzip, and can be indexed with tabix

Running things in parallel



Process individual datasets in parallel

Process genomic regions (e.g., chromosomes) in parallel (after alignment)

Some pipeline stages allow multithreading, i.e., processing within one dataset may be sped up by running on multiple CPUs, e.g.:

BWA alignment on 10 threads:

```
bwa mem -t 10 [other options]
```

HaplotypeCaller on 4 threads:

```
gatk HaplotypeCaller --native-pair-hmm-threads 4 [other options]
```

Most GATK4 tools have multithreaded versions (add **Spark** at end of tool name, like **HaplotypeCallerSpark**) – some still in BETA stage...

Caution:

total number of requested threads should never exceed the number of CPUs on the machine!

Using too many threads or running too many simultaneous jobs may decrease performance. Experiment!

BWA and **HaplotypeCaller** scale decently on up to 8-10 threads

Technical considerations

Set **PATH** to see the latest GATK (execute once in terminal or in the beginning of a script):

```
export PATH=/programs/gatk-4.1.4.0:$PATH
```

Use “--java-options” to control Java virtual machine, e.g., give java 8GB of RAM to work in:

```
gatk --java-options "-Xmx8g" [other options]
```

Specify scratch directory (important for most tools), e.g.,

```
gatk HaplotypeCaller --tmp-dir /workdir/$USER/tmp [other options]
```

Compressed (gzipped) versions of all FASTQ and VCF files can be used with all commands

For **GenotypeGVCFs**, use permissive variant emission threshold (you can filter bad variants later)

```
gatk GenotypeGVCFs -stand_call_conf 5 [other options]
```


Sample-by-sample or joint (cohort-level) variant calling?

“Seeing” reads from multiple samples (mapped to a region of reference genome) allows smarter decisions about which alleles are real and which are sequencing or alignment errors...

More confidence in variant calling

Multiple samples data allow calling a variant even if individual sample calls are of low quality

Joint calling is better, but....

Scales badly with the number of samples

“N+1” problem: what if one more (or a few more) individuals are added?
Need to repeat the calling! (in finite time....)

Alternatives to GATK

bcftools, samtools (Sanger Institute, Broad Institute,
<http://www.htslib.org/doc/bcftools.html>)

FreeBayes (Erik Garrison et al., <https://github.com/ekg/freebayes>)

- Haplotype-based variant detection (no re-alignment around indels needed)
- Better (than GATK's) Bayesian model, directly incorporating a number of metrics, such as read placement bias and allele balance
- **In our tests – a few times faster than GATK HaplotypeCaller**
- Still suffers from “N+1” problem

Sentieon (<http://sentieon.com>)

- Commercial version of GATK (currently equivalent to GATK 3.8)
- **10-30 times faster than GATK on most parts of the pipeline**
- Command syntax different from GATK (although functionality the same)
- Available on BioHPC Lab for \$500/year (need to recover license cost)
 - License: 1500 CPU cores of can run simultaneously (across all machines) at any time

BCFtools and VCFtools

Tools to filter vcf files

- Use **bcftools** if possible, as it uses htslib, developed by the consortium that maintains the vcf file format;
- VCFtools has many biologically meaningful filters not available in BCFtools
e.g. filter based on individual and site statistics.

Annotation of the variants

Consequence of the variants on the protein sequence
(e.g. missense, stop codon, frameshift)

Software: Ensembl Variant Effect Predictor (VEP)

Input: vcf file, gff file, genome sequence FASTA file;

Output: text file with annotation for each variant;

<https://biohpc.cornell.edu/lab/userguide.aspx?a=software&i=566#c>