# Variants and Filtering

RCN Bioinformatics Mentoring 2021

## OUTLINE

**Lecture**

- How do we get SNP and genotype calls?
- What is in a VCF file?
- Why and how do we filter?

**Tutorial**

- SNP calling – review code & Questions
- Filtering – Do it yourself.

# GENOMIC VARIATION



Single Nucleotide Variant
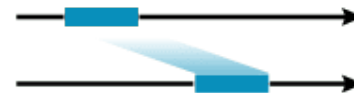
Deletion

Insertion

Tandem Duplication

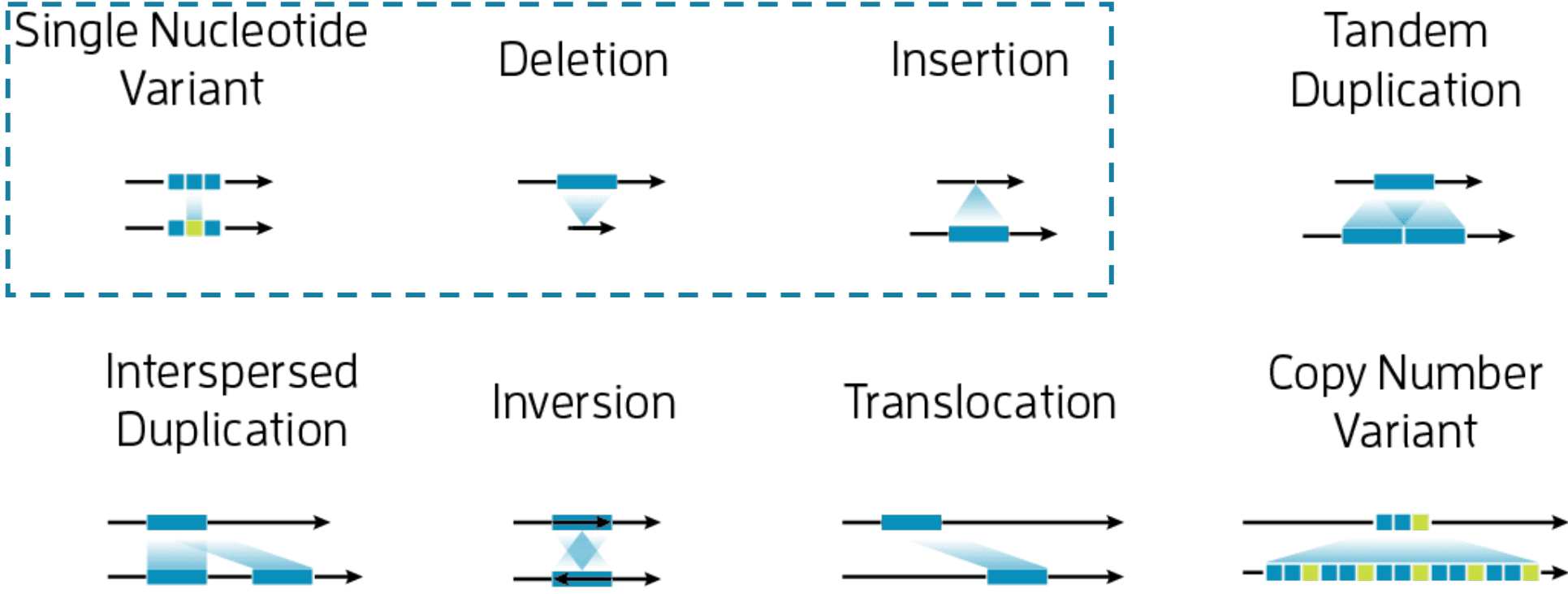Interspersed Duplication

Inversion

Translocation

Copy Number Variant

**Types of Variants**

# GENOMIC VARIATION



Types of Variants

# SNP AND GENOTYPE CALLING

- *SNP calling* finds the alleles.

- *Genotype calling* finds the genotype of the individual

# SNP AND GENOTYPE CALLING

- *SNP calling* finds the alleles.

- *Genotype calling* finds the genotype of
  the individual

Reference CCGTTAGAGTTACAATTCGA
Read 2        TTAGAGTAACAA
Read 3     CCGTTAGAGTTA
Read 4              TTACAATTCGA
Read 5         GAGTAACAA
Read 6       TTAGAGTAACAAT

- Alternate allele A: 3 reads

- Genotype: T/A (2/3)

- Unphased genotype

# SNP AND GENOTYPE CALLING

- *SNP calling* finds the alleles.

- *Genotype calling* finds the genotype of
  the individual



Reference CCGTTAGAGTTACAATTCGA
TTAGAGTAACAA
CCGTTAGAGTTA
TTACAATTCGA
CCGTTAGAGTTA
TTACAATTCG

- Alternate allele A: 1 reads

- Genotype: T/A (1/4)

# SNP AND GENOTYPE CALLING

- *SNP calling* finds the alleles.

- *Genotype calling* finds the genotype of the individual



Reference CCGTTAGAGTTACAATTCGA
TTAGAGTAACAA
CCGTTAGAGTTA
TTACAATTCGA
CCGTTAGAGTTA
TTACAATTCG

- Alternate allele A: 1 reads

- Genotype: T/A (1/4)

- How confident are you in that A? Do you really think it's a heterozygote?

# SNP AND GENOTYPE CALLING

- Confidence in alternate allele depends on:
  - Map quality of base in individual reads
  - Number of reads supporting the base


- Confidence in genotype depends on:
  - Map quality of base in individual reads
  - Number of reads supporting each allele in an individual.
  - Ratio of reads supporting 1 allele vs the other.

Reference CCGTTAGAGTTACAATTCGA
TTAGAGTAACAA
CCGTTAGAGTTA
TTACAATTCGA
CCGTTAGAGTTA
TTACAATTCG

# GENOTYPE LIKELIHOODS

- Describes the probability of a genotype given the data

$$L(G = \{A_1, A_2\}|D) \propto Pr(D|G = A_1, A_2), \qquad A_1, A_2 \in \{A, C, G, T\}.$$

  - e.g. $$Pr(D|G = \{A_1, A_2\}) = \prod_{i=1}^{M} Pr(b_i|G = \{A_1, A_2\}) = \prod_{i=1}^{M} (\frac{1}{2} Pr(b_i|A_1) + \frac{1}{2} Pr(b_i|A_2))$$

    Where M is sequencing depth and bi is the base at a given read.

- Different software use different genotype likelihood methods

- Phred scaled -10 log10(Genotype Likelihood) (PL in VCF)

- Take Home message: the quality of your genotypes are going to depend on your sequencing depth!

# VCF FORMAT

- Representation of SNV, and Indels for one or more individuals

- Includes a lot of information on quality, number of reads per site and per individual

- Individual genotype information

## Basic structure of a VCF file

# VCF HEADER

- Contains information describing

    - How file was created

    - Record contents metadata

- Commands

- Reference genome and contig names and lengths.

- Preceded by '##'

# FROM THE TUTORIAL…

```
#fileformat=VCFv4.2
#ALT=<ID=NON_REF,Description="Represents any possible alternative allele not already represented at this location by REF and ALT">
#FILTER=<ID=LowQual,Description="Low quality">
#FORMAT=<ID=AD,Number=R,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">
#FORMAT=<ID=MIN_DP,Number=1,Type=Integer,Description="Minimum DP observed within the GVCF block">
#FORMAT=<ID=PGT,Number=1,Type=String,Description="Physical phasing haplotype information, describing how the alternate alleles are phased in relation to one another; will always be heterozygous and is not in
#FORMAT=<ID=PID,Number=1,Type=String,Description="Physical phasing ID information, where each unique ID within a given sample (but not across samples) connects records within a phasing group">
#FORMAT=<ID=PL,Number=G,Type=Integer,Description="Normalized, Phred-scaled likelihoods for genotypes as defined in the VCF specification">
#FORMAT=<ID=PS,Number=1,Type=Integer,Description="Phasing set (typically the position of the first variant in the set)">
#FORMAT=<ID=SB,Number=4,Type=Integer,Description="Per-sample component statistics which comprise the Fisher's Exact Test to detect strand bias.">
#GATKCommandLine=<ID=HaplotypeCaller,CommandLine="HaplotypeCaller --emit-ref-confidence GVCF --output /home/ngsclass/Bolton/ZF/alignments/tmp_gvcf/ERR1013164_1_subs.g.vcf.gz --input ERR1013164_1_subs.sorted.
#GVCFBlock0-1=minGQ=0(inclusive),maxGQ=1(exclusive)
#GVCFBlock1-2=minGQ=1(inclusive),maxGQ=2(exclusive)
#GVCFBlock10-11=minGQ=10(inclusive),maxGQ=11(exclusive)
#GVCFBlock11-12=minGQ=11(inclusive),maxGQ=12(exclusive)
#GVCFBlock12-13=minGQ=12(inclusive),maxGQ=13(exclusive)
#GVCFBlock13-14=minGQ=13(inclusive),maxGQ=14(exclusive)
#GVCFBlock14-15=minGQ=14(inclusive),maxGQ=15(exclusive)
#GVCFBlock15-16=minGQ=15(inclusive),maxGQ=16(exclusive)
#GVCFBlock16-17=minGQ=16(inclusive),maxGQ=17(exclusive)
#GVCFBlock17-18=minGQ=17(inclusive),maxGQ=18(exclusive)
#GVCFBlock18-19=minGQ=18(inclusive),maxGQ=19(exclusive)
#GVCFBlock19-20=minGQ=19(inclusive),maxGQ=20(exclusive)
#GVCFBlock2-3=minGQ=2(inclusive),maxGQ=3(exclusive)
#GVCFBlock20-21=minGQ=20(inclusive),maxGQ=21(exclusive)
#GVCFBlock21-22=minGQ=21(inclusive),maxGQ=22(exclusive)
```

# VCF RECORDS

- Each row is a site

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | Sample1… |
|--------|-----|----|-----|-----|------|--------|------|--------|----------|
| 20 | 10001298 | . | T | A | 884.77 | . | Clipped for brevity | GT:AD:DP:GQ:PL | 1/1:0,30:30:89:913,89,0 |
| 20 | 10001436 | . | A | AAGGCT | 884.77 | . | … | GT:AD:DP:GQ:PL | 1/1:0,28:28:84:1260,84,0 |
| 20 | 10004769 | . | TAAAACTATGC | T | 622.73 | . | … | GT:AD:DP:GQ:PL | 0/1:18,17:35:99:660,0,704 |

# A ZF VCF FILE

# A ZF VCF FILE

# RECORD INFO

- Exact definitions included in the header!

- QUAL = Phred scaled (-10*log10(p)) probability that there is a polymorphism at this site.

- FILTER= names of any filters applied

- INFO= site level annotations, semi-colon separated. See header for more info.

- FORMAT= names and order of presentation of sample level variant information

https://gatk.broadinstitute.org/hc/en-us/articles/360035531692

# RECORD INFO

- FORMAT: GT:AD:DP:GQ:PL

- GT = genotype,

  - 0/0 homozygous for reference allele

  - 0/1 heterozygous reference/alternate

  - ... | indicates a phased genotype

- AD=Depth per allele

- DP= Depth

- GQ=Genotype Quality

- PL=Phred-scaled genotype likelihood of each genotype (0/0, 0/1, 1/1)

```
GT:AD:DP:GQ:PL    0/0:1,0:1:3:0,3,42
GT:AD:DP:GQ:PL    0/0:1,0:1:3:0,3,42
```
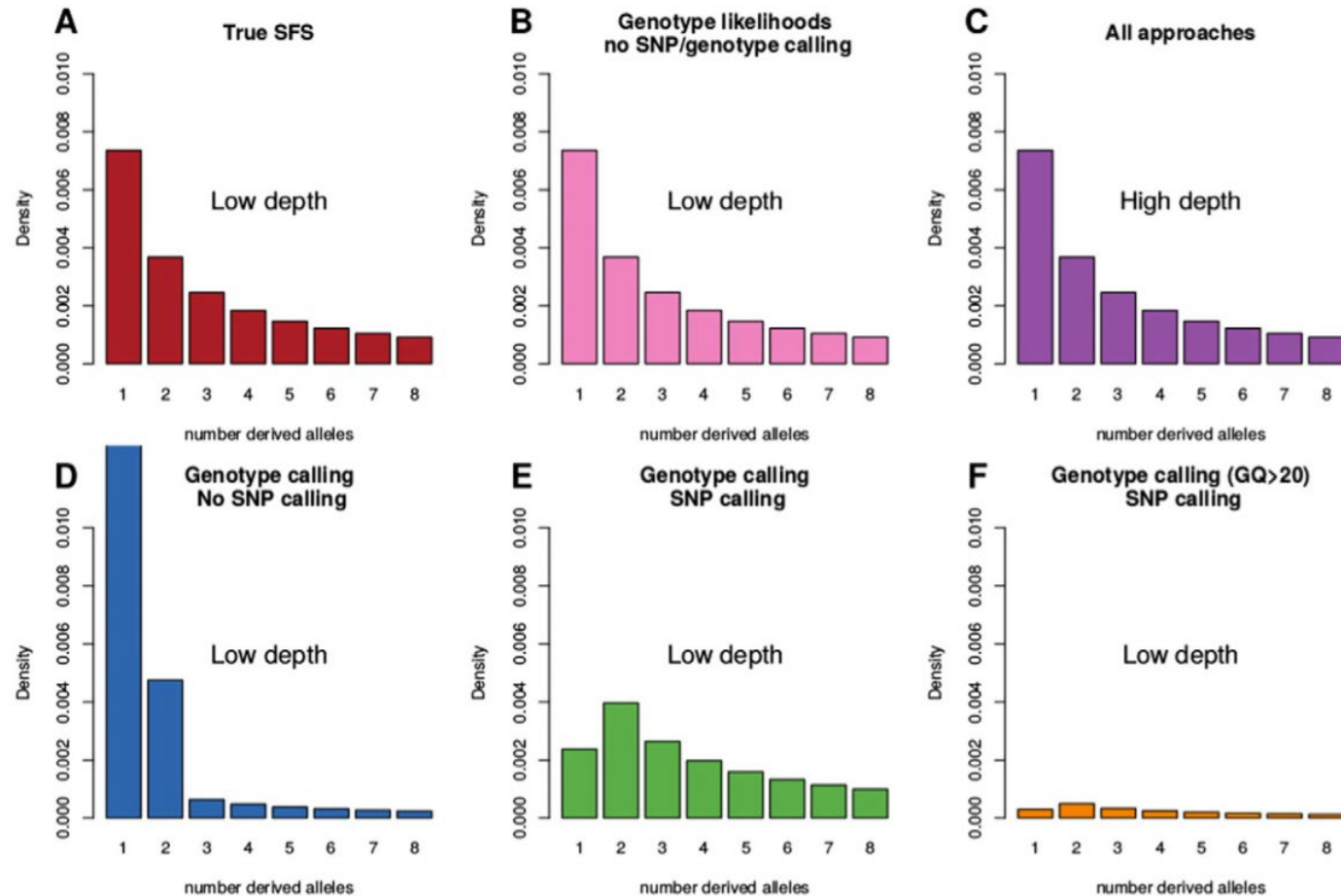
# VCF FROM SAMTOOLS MPILEUP

# SNP CALLING SOFTWARE

- samtools mpileup

- FreeBayes

- Genome Analysis Toolkit (GATK)

- Angsd (based on Genotype Likelihoods)

…

# FILTERING AND DEPTH MATTERS
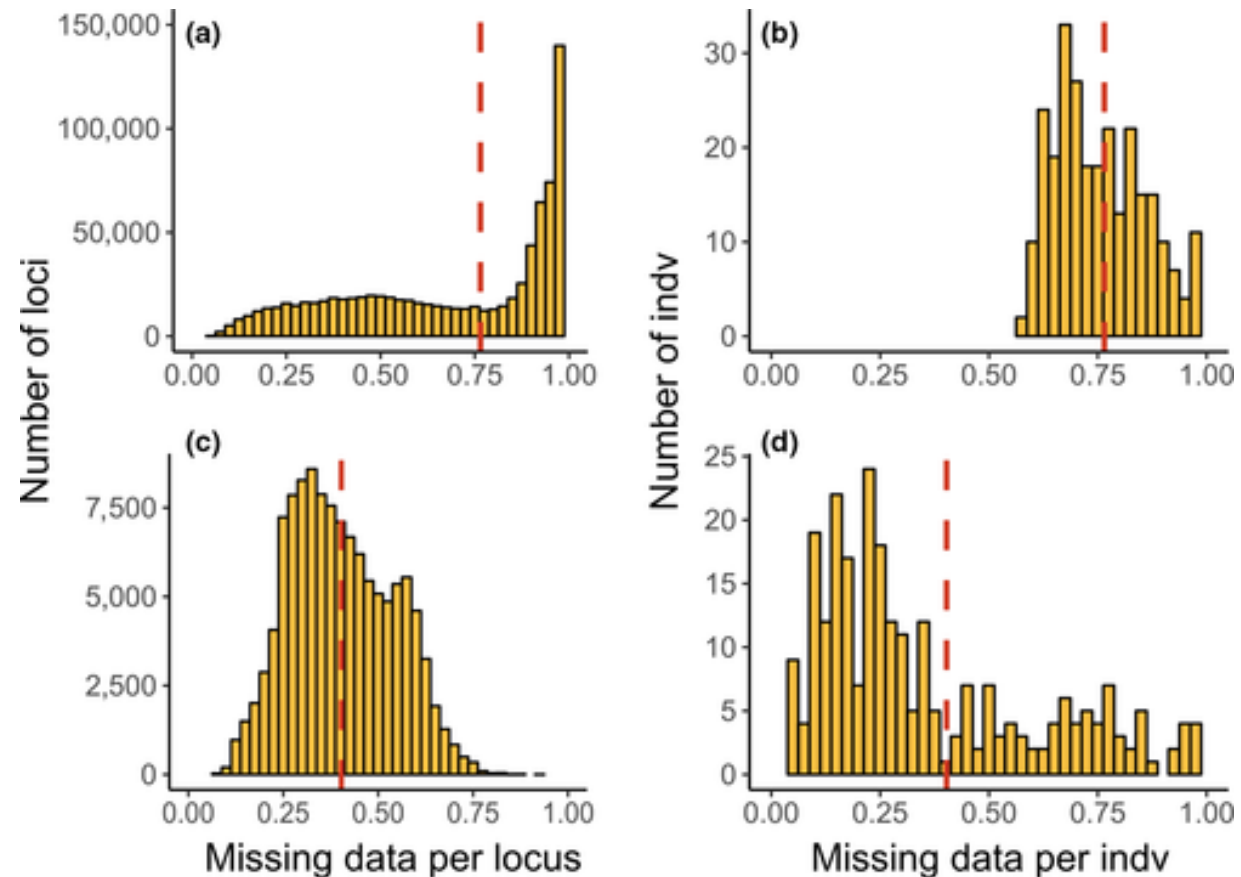
# FILTERING

- Will depend on your sampling design and your questions

- In general:

    - High quality genotypes – depth, allele balance etc

    - Low level of missing data
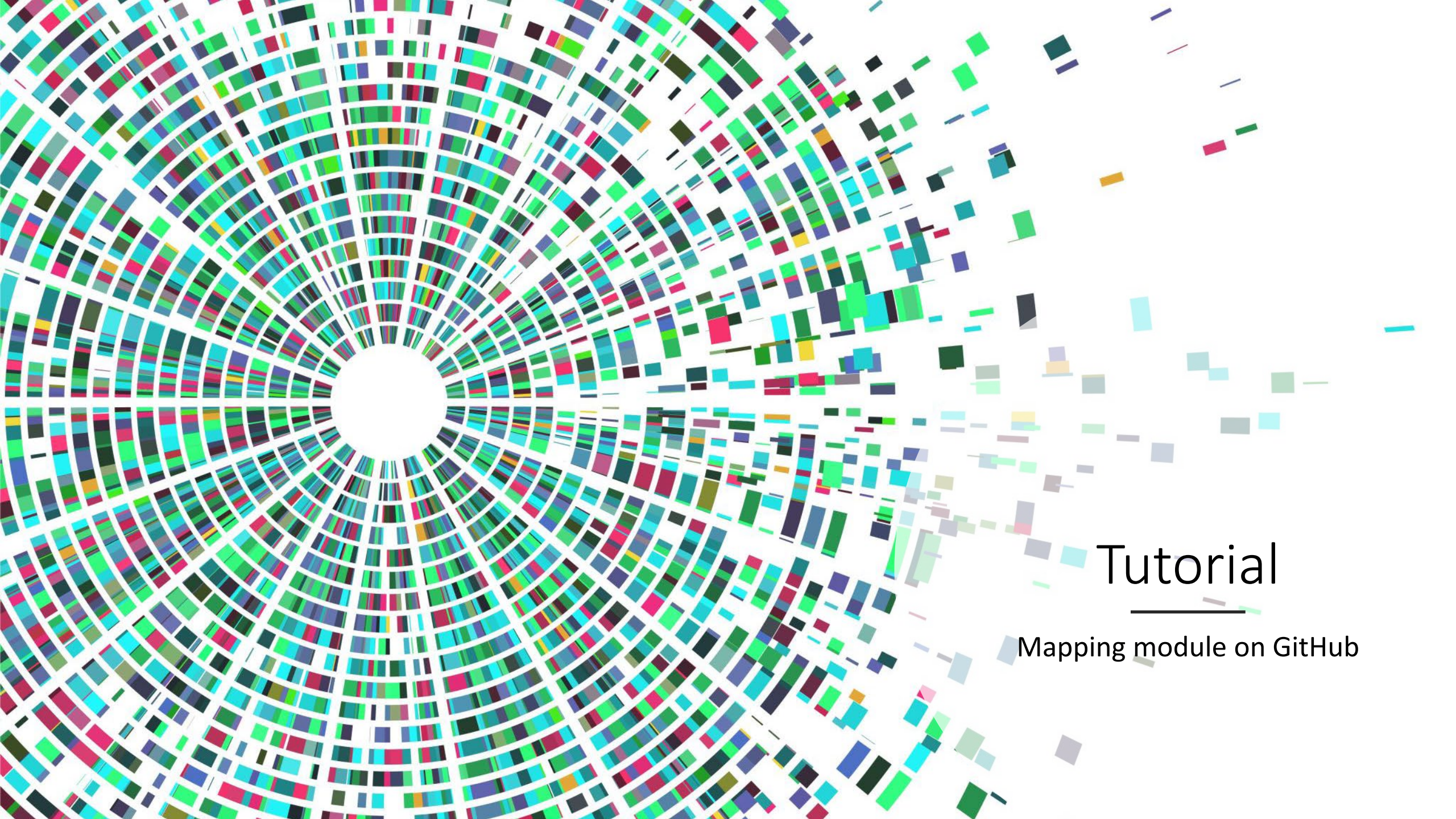
# PLOTTING THE DATA

- Before

- Apply Filter

- After

# FURTHER READING & REFERENCES

- da Fonseca et al (2016) Next-generation biology: sequencing and data analysis approaches for non-model organisms. *Marine Genomics.*

- Korneliussen et al (2014) ANGSD: Analysis of Next Generation Sequencing Data. *BMC Bioinformatics.*

- Nielsen et al (2012) SNP Calling, Genotype Calling, and Sample Allele Frequency Estimation from New-Generation Sequencing Data. *PLoS ONE*

- Nielsen et al (2011) Genotype and SNP calling from next-generation sequencing data. *Nature Reviews Genetics.*

- O'Leary et al (2018). These aren't the loci you're looking for: Principles of effective SNP filtering for molecular ecologists. *Molecular Ecology.*

- Bhatia et al. (2013). Estimating and interpreting FST: the impact of rare variants. *Genome Research.*

- Linck and Battey (2019). Minor allele frequency thresholds strongly affect population structure inference with genomic datasets. *Moleuclar Ecology Resources.*

- Broad Institute (2021) GATK Best Practices Workflow: Germline short variant discovery. https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels-

- Salter & Faircloth (2020) Running GATK in Parallel. http://protocols.faircloth-lab.org/en/latest/protocols-computer/analysis/analysis-gatk-parallel.html
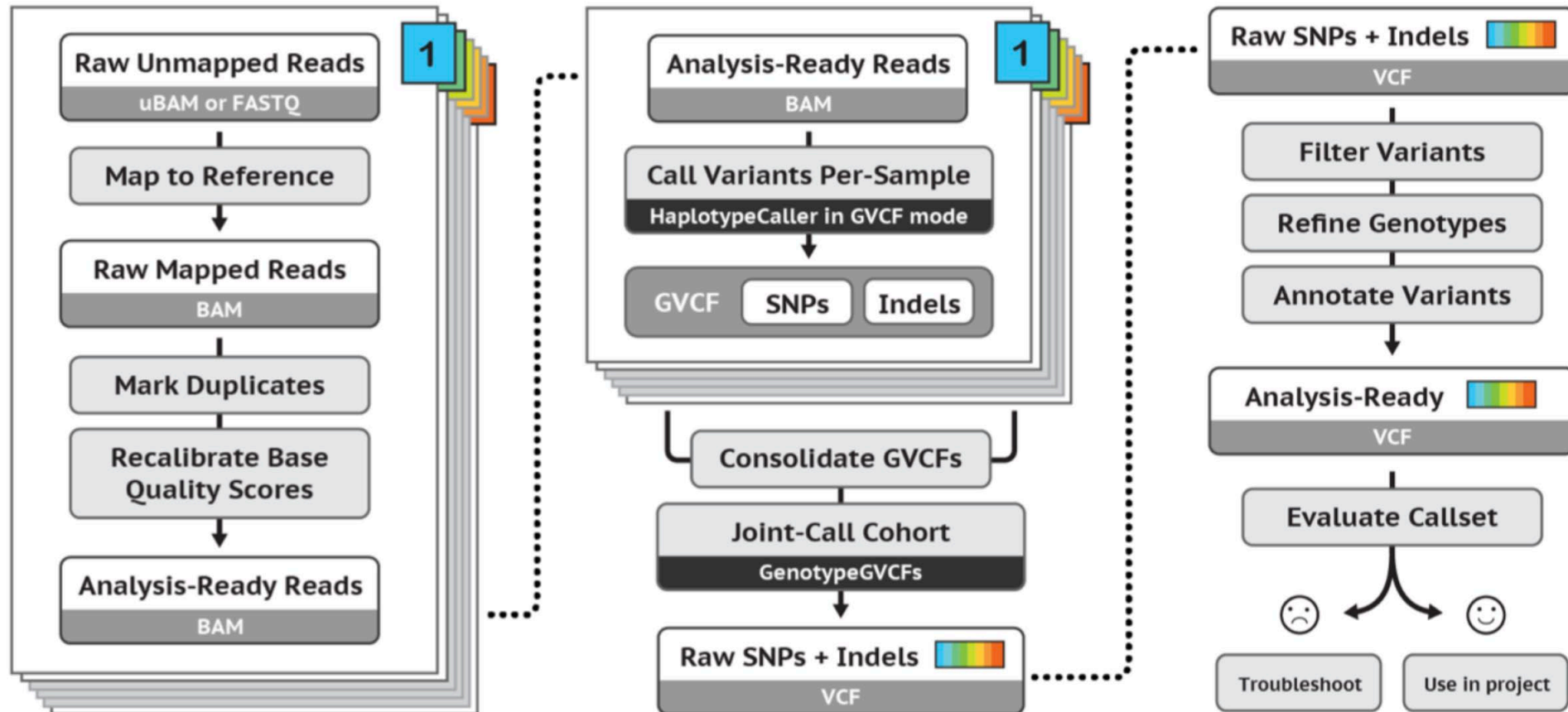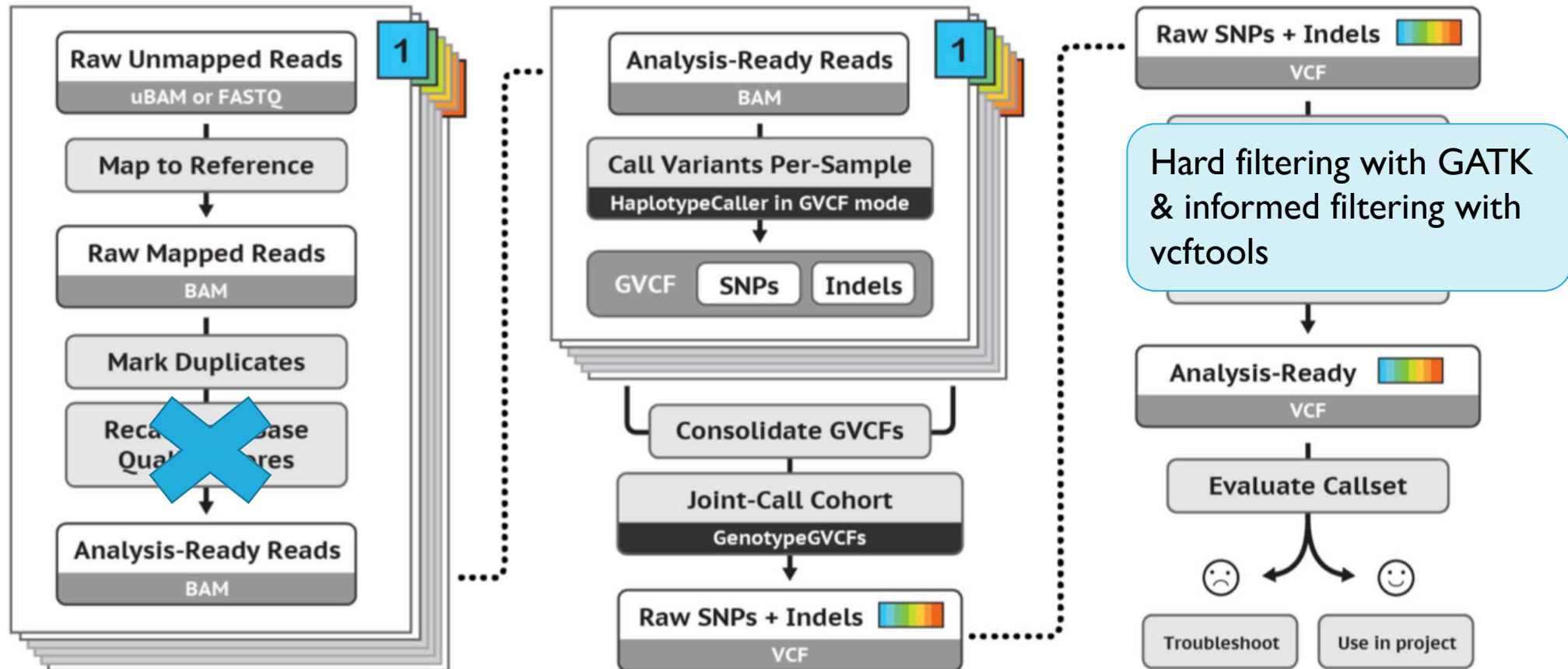
# Tutorial

Mapping module on GitHub

# GATK GERMLINE BEST PRACTICES

# GATK WITH NON-MODEL ORGANISM



Hard filtering with GATK & informed filtering with vcftools

# A NOTE ABOUT READ GROUPS

- Set at mapping step or afterwards using PicardTools

- Fastq files can be split by lane

- Enables multiple libraries, flowcells, and lanes to be run per "individual" – SM.

- GATK can take into account some of the batch effects.