# Variant Detection and Genotype Calling Using GATK

**Part 2**

Dr Amanda Chamberlain and Dr Christy Vander Jagt
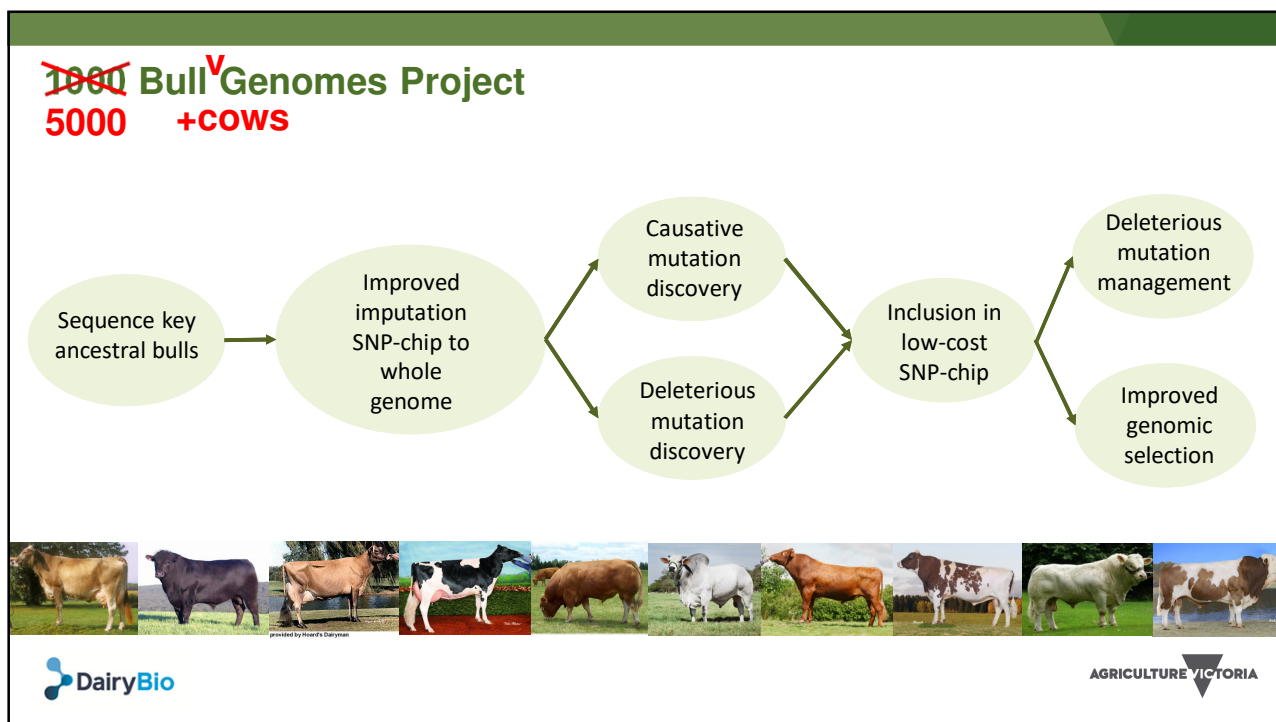
AGRICULTURE VICTORIA

1

## Overview

- The 1000 Bull Genomes Project

- Variant discovery (GATK 3.8)
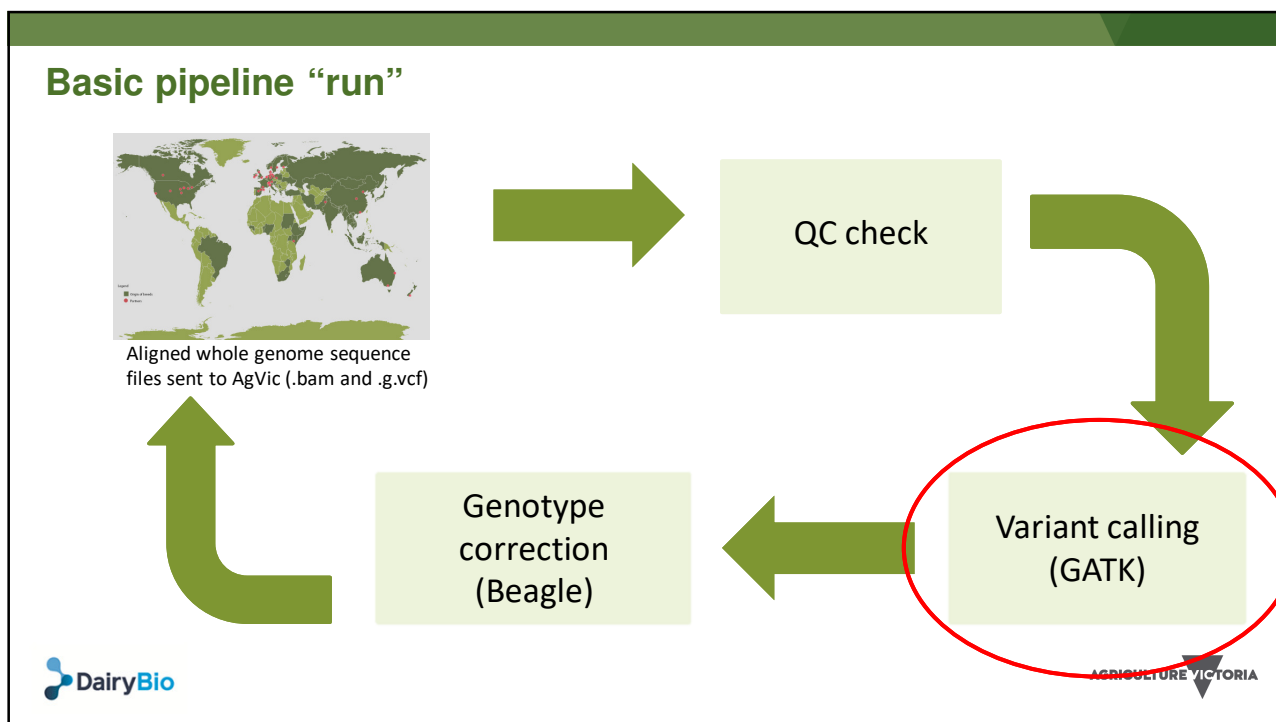
- Variant evaluation

- GATK versus SAMtools

Acknowledgement: A lot of the figures/graphs in this presentation come directly from the GATK documentation. You can find this documentation plus more at https://gatk.broadinstitute.org/hc/en-us
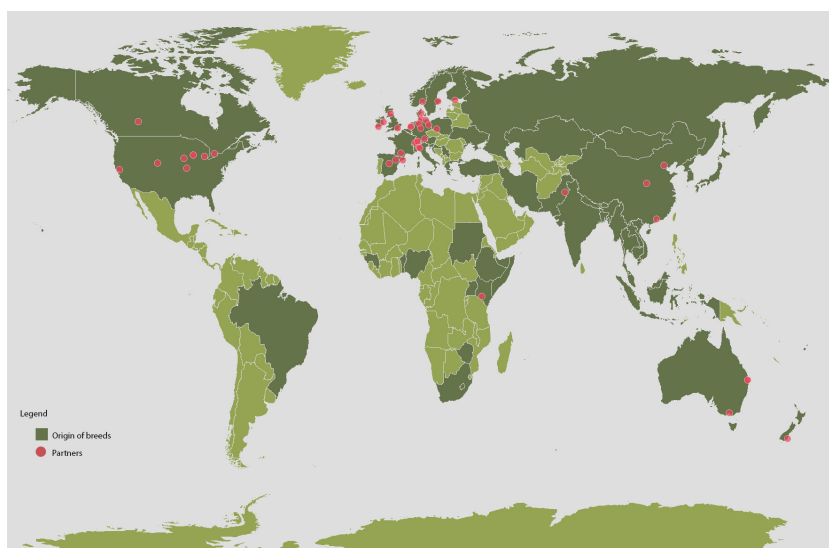
DairyBio

AGRICULTURE VICTORIA

2

1000 5000 Bull Genomes Project
+cows

Sequence key ancestral bulls → Improved imputation SNP-chip to whole genome → Causative mutation discovery / Deleterious mutation discovery → Inclusion in low-cost SNP-chip → Deleterious mutation management / Improved genomic selection

3



**Basic pipeline "run"**

Aligned whole genome sequence files sent to AgVic (.bam and .g.vcf) → QC check → Variant calling (GATK) → Genotype correction (Beagle)

4

## Project Partners



DairyBio

AGRICULTURE VICTORIA

5
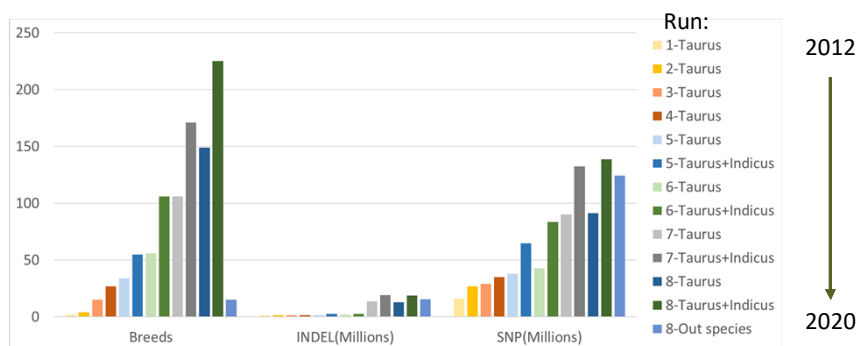
## Run8 – results

**Taurus**
- 140+ breeds
- 4,109 animals
- 12.9mil INDEL
- 91.4mil SNP

**Taurus-Indicus**
- 220+ breeds
- 4,931 animals
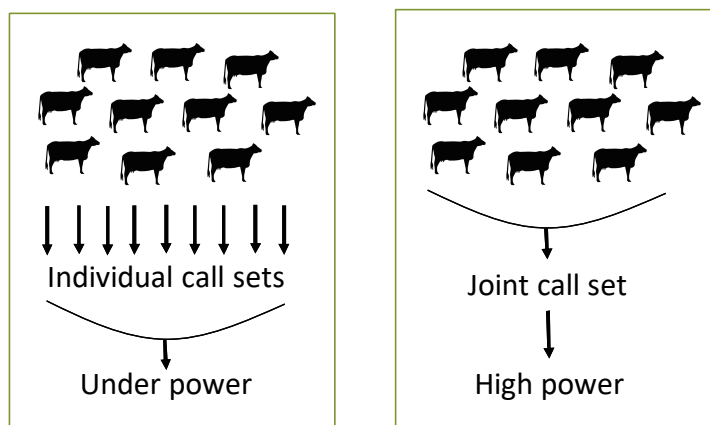- 18.9mil INDEL
- 138.9mil SNP

**Out-species**
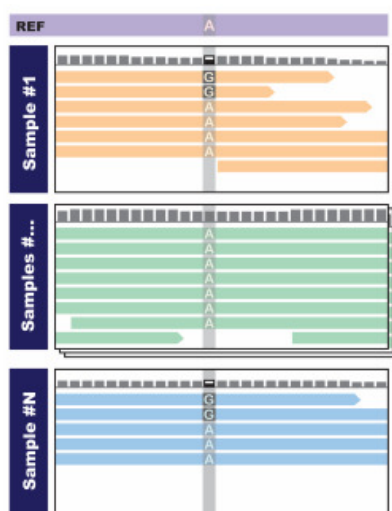- 15 breeds
- 327 animals
- 15.4mil INDEL
- 124.4mil SNP



DairyBio  Dairy Australia

GARDINER FOUNDATION   AGRICULTURE VICTORIA

6

## The power of joint calling

- Single genome not very useful
- Population data adds valuable information



Individual call sets

Under power

Joint call set

High power
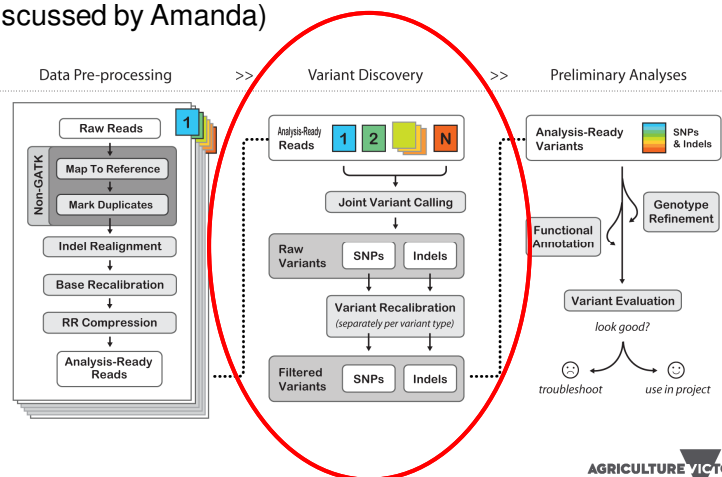
DairyBio

AGRICULTURE VICTORIA

7



- Sample #1 or Sample #N alone:
  - **weak evidence for variant**
  - **may miss calling the variant**

- Both samples seen together:
  - **unlikely to be artifact**
  - **call the variant more confidently**
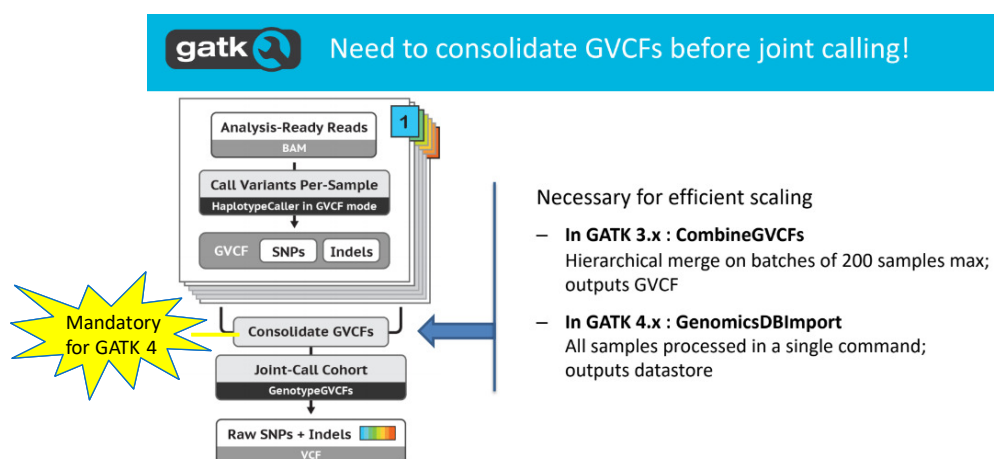
DairyBio

AGRICULTURE VICTORIA

8

# Variant calling with GATK

- GATK best practises:
  1. GVCF generation (already discussed by Amanda)
  2. Combine GVCFs (GATK4)
  3. Joint variant calling
  4. Variant recalibration
  5. Variant filtering



9

# Consolidate GVCFs prior to joint calling GATK 4



gatk — Need to consolidate GVCFs before joint calling!

Necessary for efficient scaling

- **In GATK 3.x : CombineGVCFs**
  Hierarchical merge on batches of 200 samples max; outputs GVCF

- **In GATK 4.x : GenomicsDBImport**
  All samples processed in a single command; outputs datastore

Source:

10

## Consolidate GVCFs prior to joint calling – GATK 3.8

- Not 100% necessary
- Recommended to run CombineGVCFs tool on batches of 200 samples
- Command:

```
gatk CombineGVCFs \
    -R reference.fasta \
    -V sample1.g.vcf \
    -V sample2.g.vcf \
    -O combined.g.vcf
```

- More information: https://gatk.broadinstitute.org/hc/en-us/articles/360037053272-CombineGVCFs

**DairyBio**

AGRICULTURE VICTORIA

11

## Consolidate GVCFs prior to joint calling – GATK 4

- Compulsory for GATK 4
- Can use CombineGVCFs tool but best practises recommends GenomicsDBImport
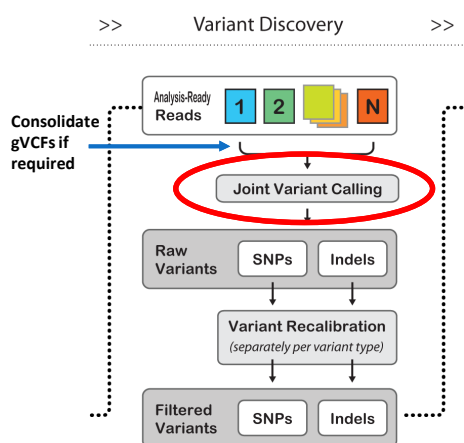- Command:

```
gatk GenomicsDBImport \
    -R reference.fasta \
    -V sample1.g.vcf \
    -V sample2.g.vcf \
    -L chr20,chr21 \
    --genomicsdb-workspace-path gvcfs_db
```

- More information: https://gatk.broadinstitute.org/hc/en-us/articles/360036883491-GenomicsDBImport

**DairyBio**

AGRICULTURE VICTORIA

12

## Joint Variant Calling – GenotypeGVCFs



GenotypeGVCFs can take multiple GVCF files in GATK 3.8 (multiple –V variants), but only a single file in GATK 4

GATK 3.8 command:

```
java –jar GenomeAnalysisTK.jar
    –T GenotypeGVCFs \
    –R human.fasta \
    –V sample1.g.vcf \
    –V sample2.g.vcf \
    –V sampleN.g.vcf \
    –o output.vcf
```

GATK 4 commands:

```
gatk GenotypeGVCFs \
    –R reference.fasta \
    –V variants.g.vcf \
    –O final_variants.vcf
```

```
gatk GenotypeGVCFs \
    –R reference.fasta \
    –V gendb://gvcfs_db \
    –O final_variants.vcf
```
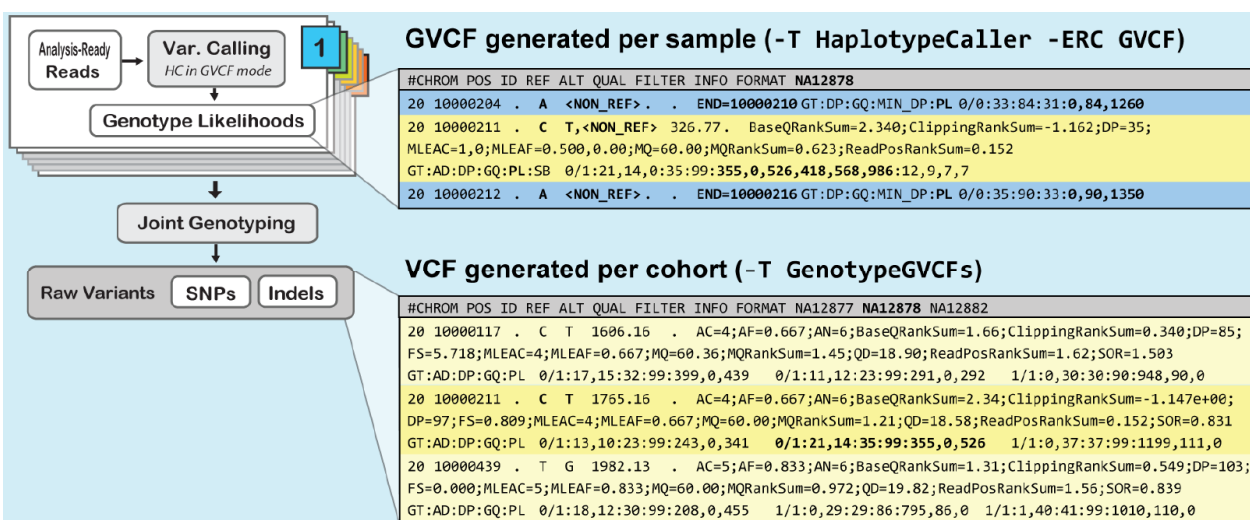
More information: https://gatk.broadinstitute.org/hc/en-us/articles/360037594731-GenotypeGVCFs
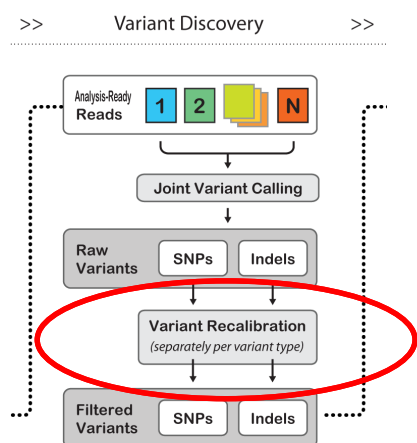
13

## Joint Variant Calls



14

## Variant call set filtering



- Variant calling algorithms are very permissive by design
- How to filter?
  - Hard filtering
    - Multiple threshold values
    - Binary choice: pass or fail
  - Variant "recalibration"
    - Machine learning
    - Annotation profile of "good" vs "bad" variants
    - Multiple annotations
- Trade-off between sensitivity and specificity

More information: https://gatk.broadinstitute.org/hc/en-us/articles/360035531112--How-to-Filter-variants-either-with-VQSR-or-by-hard-filtering

15

## 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    0/1:12:99    0/1:11:89    0/1:28:37
```

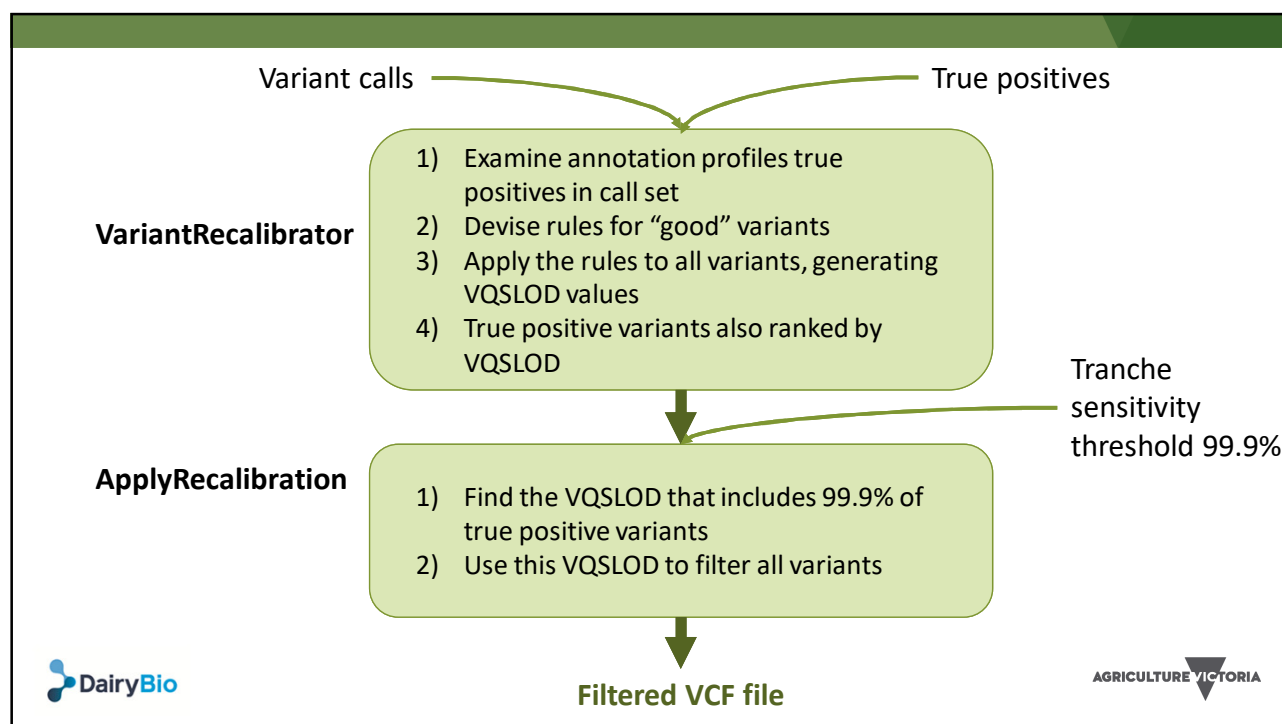| | | | |
|---|---|---|---|
| AC | No. chromosomes carrying alt allele | MLEAF | Max likelihood AF |
| AN | Total no. of chromosomes | MQ | RMS MAPQ of all reads |
| AF | Allele frequency | MQ0 | No. of MAPQ 0 reads at locus |
| DP | Depth of coverage | QD | QUAL score over depth |
| MLEAC | Max likelihood AC | | |

INFO field

DairyBio

AGRICULTURE VICTORIA

16

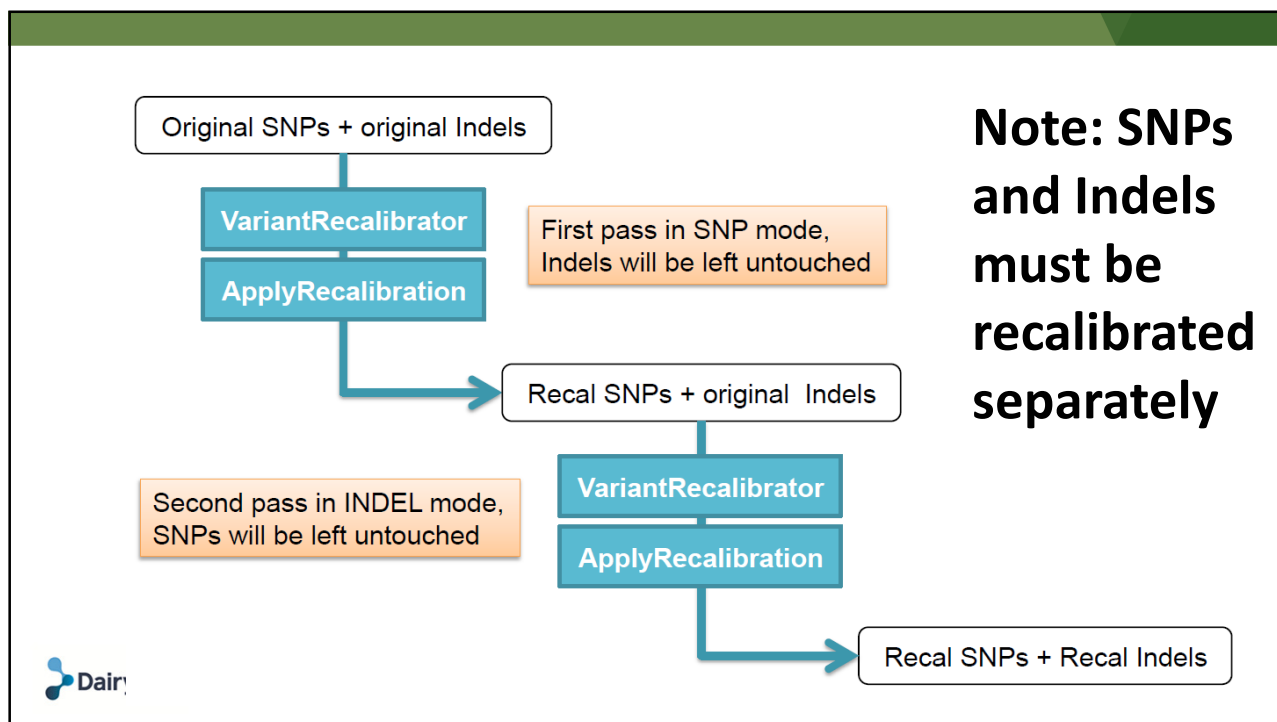## Best practises: Variant Quality Score Recalibration (VQSR)

- "Sophisticated filtering technique applied on the variant callset that uses machine learning to model the technical profile of variants in a training set and uses that to filter out probable artifacts from the callset."

- Two step process:
    1. Variant recalibration (VariantRecalibrator)
    2. Applying the recalibration (ApplyRecalibration)

More information: https://gatk.broadinstitute.org/hc/en-us/articles/360035531612-Variant-Quality-Score-Recalibration-VQSR-

17

---

Variant calls      True positives

**VariantRecalibrator**

1) Examine annotation profiles true positives in call set
2) Devise rules for "good" variants
3) Apply the rules to all variants, generating VQSLOD values
4) True positive variants also ranked by VQSLOD

Tranche sensitivity threshold 99.9%

**ApplyRecalibration**

1) Find the VQSLOD that includes 99.9% of true positive variants
2) Use this VQSLOD to filter all variants

DairyBio

AGRICULTURE VICTORIA

**Filtered VCF file**

18

19



20

## Resource datasets

- Three types of resources:
    - 1) Truth
        - Validated to a high degree of confidence
        - Representative of "true" sites (truth=true)
        - Used to train recalibration model (training=true)
        - Used to determine where to set cutoff in VQSLOD sensitivity
    - 2) Training
        - Validated to some degree of confidence
        - May contain false positives (truth=false)
        - Used to train recalibration model (training=true)
    - 3) Known
        - Not validated to a high degree of confidence (truth=false)
        - Not used to train recalibration model (training=false)
        - Only for reporting purposes, not used in any calculations
    - Prior
        - Phred-scaled estimate of data accuracy

```
-resource:hapmap,known=false,training=true,truth=true,prior=15.0
hapmap_3.3.b37.sites.vcf
-resource:omni,known=false,training=true,truth=false,prior=12.0
omni2.5.b37.sites.vcf
-resource:1000G,known=false,training=true,truth=false,prior=10.0
1000G.b37.sites.vcf
-resource:dbsnp,known=true,training=false,truth=false,prior=2.0
dbsnp_137.b37.vcf
```

21

## Where to find resources

- Human genome training, truth and known resource datasets:
  https://gatk.broadinstitute.org/hc/en-us/articles/360035890811-Resource-bundle
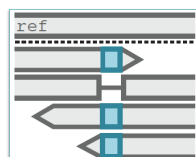
- GATK do not provide resources for non-human organisms
    - you need to have at least truth and training resource datasets (with assigned prior likelihoods)
    - GATK forum topic called "Non-Human": https://gatk.broadinstitute.org/hc/en-us/community/topics/360001496611-Non-Human

**DairyBio**

**AGRICULTURE VICTORIA**

22

## Call set evaluation

- Minimum recommended metrics for call set evaluation

**Number of Indels & SNPs**



**Indel Ratio**



**Genotype Concordance**



**TiTv Ratio**



More information: https://gatk.broadinstitute.org/hc/en-us/articles/360035531572-Evaluating-the-quality-of-a-germline-short-variant-callset

23

## Number of Indels and SNPs

- Variants = Indels + SNPs
- Useful for order-of-magnitude sanity check
- For WGS the number of variants should be ~4.4M for a single sample (human)



DairyBio

AGRICULTURE VICTORIA

24

## TiTv Ratio

- Ratio of transition to transversion SNPs
- If random, expect ratio of 0.5
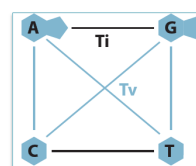- Twice as many possible transversions versus transitions
- Low TiTv ratio indicates high rate of false positives
- For WGS TiTv ratio should be 2.0-2.1 for humans



DairyBio

AGRICULTURE VICTORIA

25

## Indel Ratio

- Ratio of insertions to deletions
- Varies by type of study e.g. rare variant association vs common variant association

| Variant Association Study type | Indel Ratio |
|---|---|
| Common | ~1 |
| Rare | 0.2-0.5 |

DairyBio

AGRICULTURE VICTORIA

26

## Concordance

**SENSITIVITY vs. FALSE DISCOVERY RATE**

**Variant level concordance**

my callset
**12**

gold standard
**10**

SENSITIVITY

$$\frac{TP}{TP + FN} = \frac{7}{7 + 3} = 70\%$$

FALSE DISCOVERY RATE

$$\frac{FP}{FP + TP} = \frac{5}{5 + 7} = 42\%$$

false positives (FP)
**5**

false negatives (FN)
**3**

true positives (TP)
**7**

**Genotype concordance**

**GENOTYPE CONCORDANCE**

gold standard ☆ ☆ ★ ☆ ★ ★ ☆

my callset ☆ ★ ☆ ☆ ★ ☆ ☆

matches (4) | | | | |
1      1   1    1

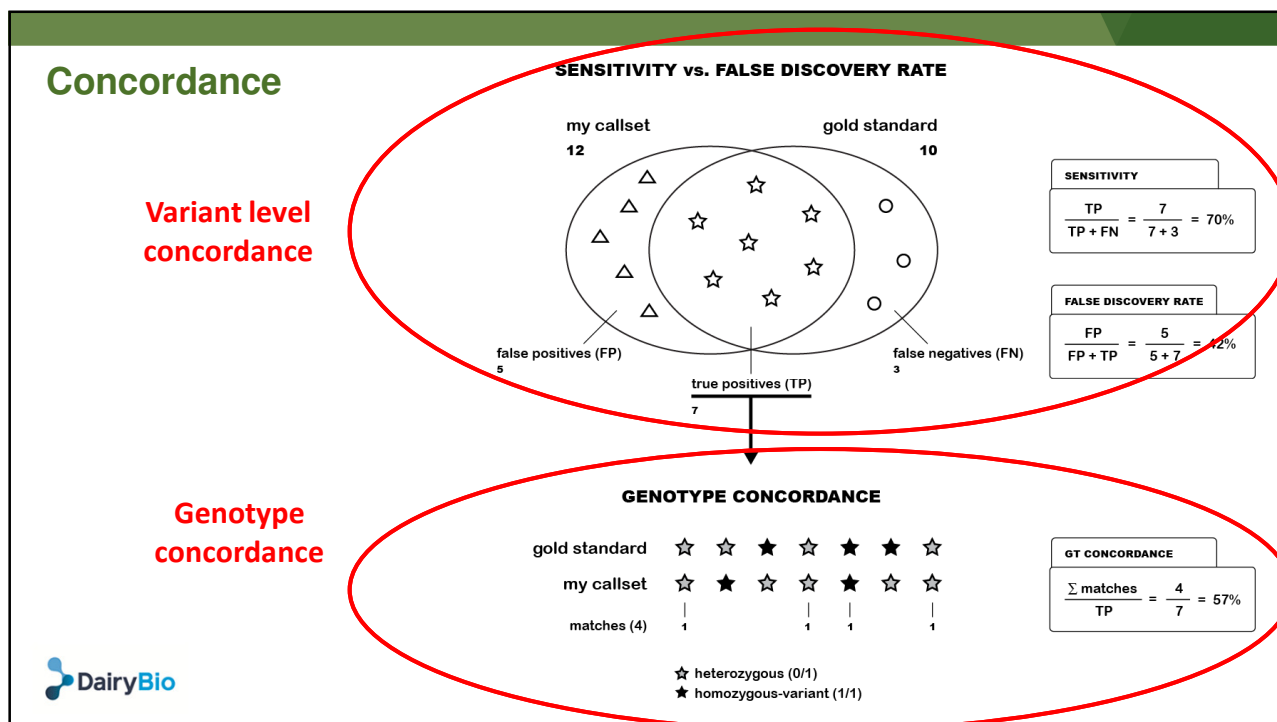GT CONCORDANCE

$$\frac{\sum matches}{TP} = \frac{4}{7} = 57\%$$

☆ heterozygous (0/1)
★ homozygous-variant (1/1)

DairyBio

27

---

| | **Variant Level Evaluation** | **Genotype Level Evaluation** |
|---|---|---|
| **GATK** | VariantEval<br><br>```java -jar GenomeAnalysisTK.jar \    -T VariantEval \    -R reference.b37.fasta \    -eval callset.vcf \    -D truthset.vcf \    -o results.eval.grp``` | GenotypeConcordance<br><br>```java -jar GenomeAnalysisTK.jar \    -T GenotypeConcordance \    -R reference.b37.fasta \    --comp truthset.vcf \    --eval callset.vcf \    -o results.grp``` |
| **Picard** | CollectVariantCallingMetrics<br><br>```java -jar picard.jar \    CollectVariantCallingMetrics    INPUT=callset.vcf \    DBSNP=truthset.vcf \    OUTPUT=results``` | GenotypeConcordance<br><br>```java -jar picard.jar \    GenotypeConcordance \    CALL_VCF=callset.vcf \    TRUTH_VCF=truthset.vcf \    CALL_SAMPLE=sampleName \    TRUTH_SAMPLE=sampleName \    OUTPUT=results``` |

Dai

VICTORIA

28

# Call set evaluations used for the 1000 Bull Genomes Project

- Concordance to HD chip (~1000 samples)

- Opposing homozygotes
  - Requires pedigree information

- Number of unique variants

- Level of heterozygosity

DairyBio

AGRICULTURE VICTORIA

29

# Other variant callers

- SAMtools, FreeBayes, Platypus, VarScan and more…

- Comparison between GATK and SAMtools
  - http://www.wcgalp.org/system/files/proceedings/2018/which-best-variant-caller-large-whole-genome-sequencing-datasets.pdf

DairyBio

AGRICULTURE VICTORIA

30

## SAMtools versus GATK: SNP calls

| Coverage | SAMtools_mpileup | | GATK_1000bullFilters | | GATK_ t99.9 | |
|---|---|---|---|---|---|---|
| | High | 10x | High | 10x | High | 10x |
| **Number filtered SNP** | 23,303,340 | 22,012,522 | 24,130,168 | 22,662,445 | 25,140,036 | 23,828,447 |
| **Mean Concordance** | 0.982 | 0.980 | 0.982 | 0.979 | 0.982 | 0.979 |
| **Mean unique variants** | 171.772 | 128.579 | 196.489 | 165.053 | 169.365 | 140.962 |
| **Mean Heterozygosity** | 0.172 | 0.180 | 0.169 | 0.171 | 0.175 | 0.179 |
| **Mean OppHom** | 0.0015 | 0.0020 | 0.0013 | 0.0019 | 0.0020 | 0.0030 |
| **Percent 800k SNP** | 97.02% | 97.00% | 96.77% | 94.99% | 98.92% | 98.94% |

DairyBio

AGRICULTURE VICTORIA

31

## SAMtools versus GATK: Indel calls

| Coverage | SAMtools_mpileup | | GATK_1000bullFilters | | GATK_ t99.9 | |
|---|---|---|---|---|---|---|
| | High | 10x | High | 10x | High | 10x |
| **Number filtered INDEL** | 2,022,663 | 1,956,676 | 2,476,684 | 2,274,590 | 2,319,278 | 2,573,292 |
| **Mean unique variants** | 10.639 | 8.034 | 23.504 | 17.596 | 12.555 | 14.871 |
| **Mean Heterozygosity** | 0.175 | 0.184 | 0.155 | 0.160 | 0.180 | 0.172 |
| **Mean OppHom** | 0.0034 | 0.0044 | 0.0012 | 0.0020 | 0.0022 | 0.0036 |

DairyBio

AGRICULTURE VICTORIA

32

## Acknowledgements

**Intellectual Climate Fund (ICF)**

**The Broad (creators of GATK)**

**1000 Bull Genomes Project**

**Agriculture Victoria BASC team**

DairyBio

AGRICULTURE VICTORIA

33

# Thank you

AGRICULTURE VICTORIA

34