# Joint variant calling

GVCF-based workflow using
GenomicsDB and GenotypeGVCFs

BROAD INSTITUTE

https://software.broadinstitute.org/gatk/
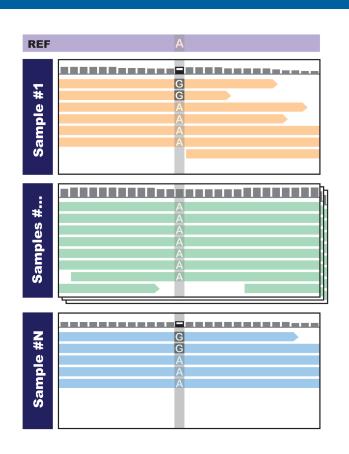
gatk

- Single genome in isolation: almost never useful

- Family or population data
  add valuable information

  – rarity of variants
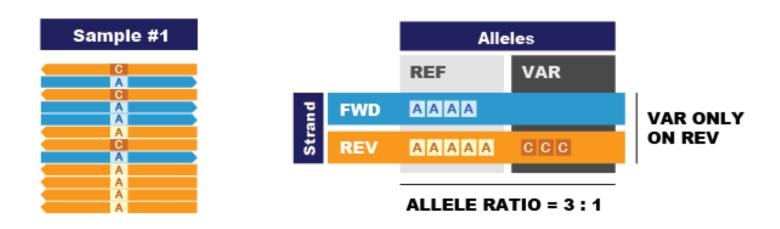
  – *de novo* mutations

  – ethnic background



Individual callsets

Underpowered

Joint callset

Superpowered!

# Discovery is empowered at difficult sites



- 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**

# Joint analysis helps resolve bias issues



**Single sample showing strand and allelic biases – would you call it?**
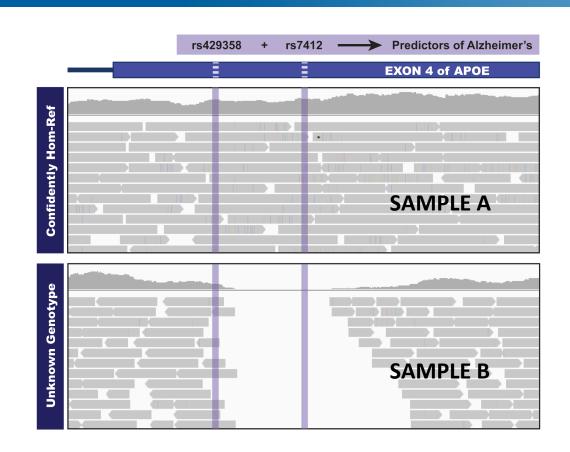
# Joint analysis helps resolve bias issues



**Decision process using evidence from multiple samples to filter out sites showing systematic biases**
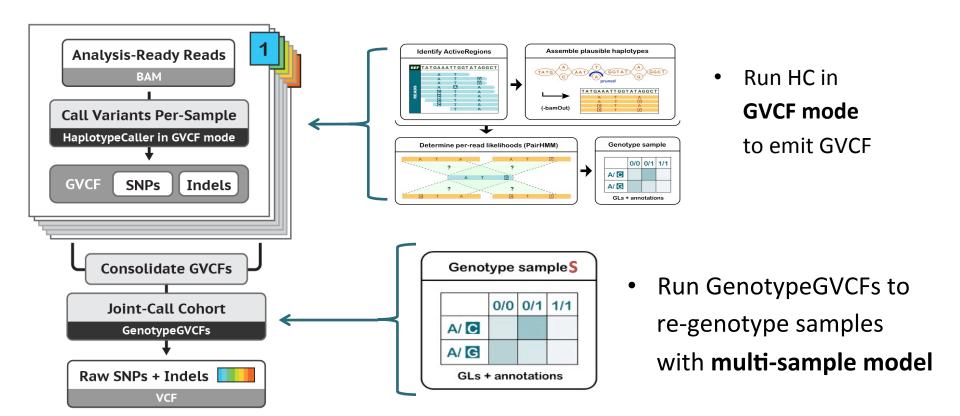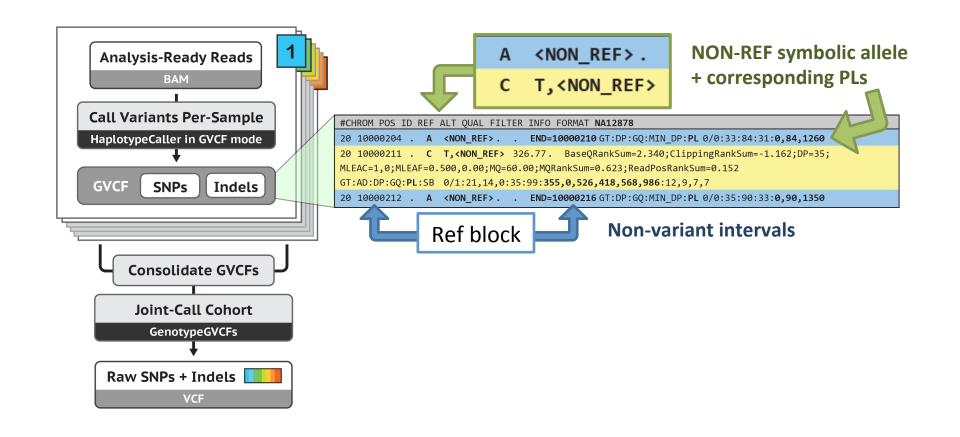
- **Analyzed individually:**
  - No call for either sample
  - Very different reasons!

- **In joint analysis with other samples:**
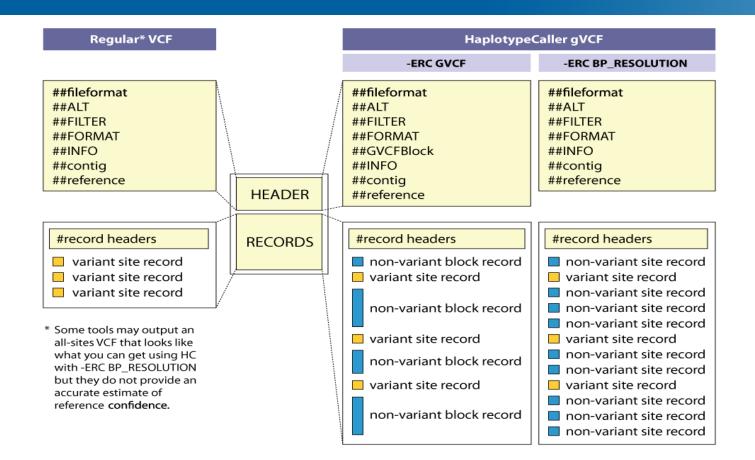  - Hom-ref call and no-call genotypes emitted



rs429358 + rs7412 → **Predictors of Alzheimer's**

**EXON 4 of APOE**

Confidently Hom-Ref

**SAMPLE A**

Unknown Genotype

**SAMPLE B**

# Joint calling implemented as a two-step process for scalability



- Run HC in **GVCF mode** to emit GVCF

- Run GenotypeGVCFs to re-genotype samples with **multi-sample model**

# GVCF intermediate contains reference confidence estimate

# GVCFs are valid VCFs with extra information

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

# Consolidating GVCFs

**With CombineGVCFs:**

```
gatk CombineGVCFs \
    -R reference.fasta \
    -V sample1.g.vcf \
    -V sample2.g.vcf \
    -O combined.g.vcf
```
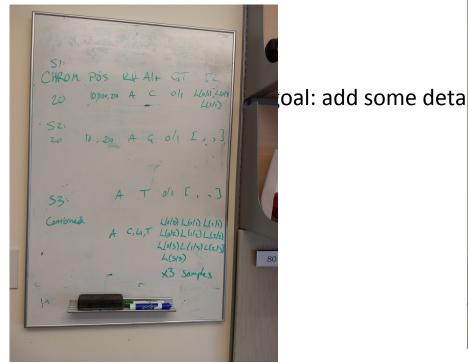
*CombineGVCFs does not scale well*

**With GenomicsDBImport:**

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

*GenomicsDBImport scales well but must be run on a single interval at a time*

goal: add some deta[...]works

Combined GVCFs have big PL arrays
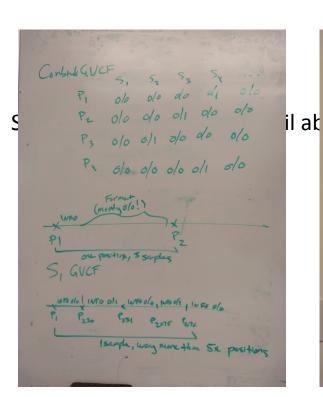at multiallelic sites
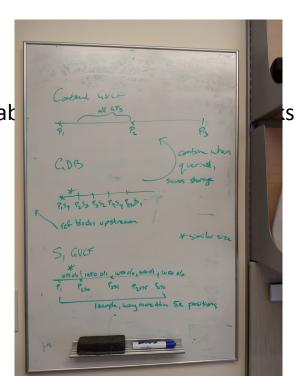
Query by sample in GDB is fast

# Storage size comparisons



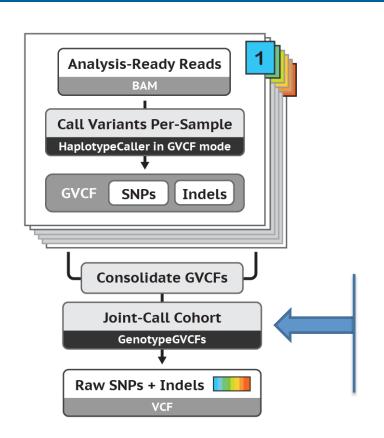Data is sparse (mostly repeated values from ref block) but gets bigger when we combine



Comparison of combined storage vs single GVCF



Theoretical size comparison

# Joint calling with GenotypeGVCFs



- **GenotypeGVCFs** can take either a **single GVCF file** (can be a merged multi-sample GVCF from CombineGVCFs) or a **GenomicsDB datastore**

- No more multiple inputs! (unlike GATK3)

# Running GenotypeGVCFs

**On a single- or multi-sample GVCF:**

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

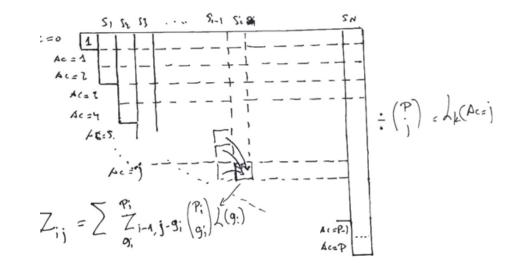**On a GenomicsDB workspace:**

```
gatk GenotypeGVCFs \
   -R reference.fasta \
   -V gendb://gvcfs_db \
   -O final_variants.vcf
```
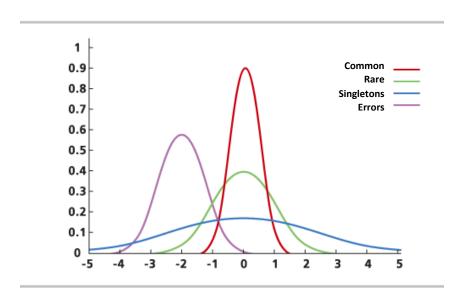
*GenotypeGVCFs cannot take multiple inputs (unlike the GATK3 version)*

- Uses human SNP heterozygosity
  1/1000 bases = Phred scale Q30
  *(can be modified)*

- QUAL > 30 means a variant is more
  likely than this base level

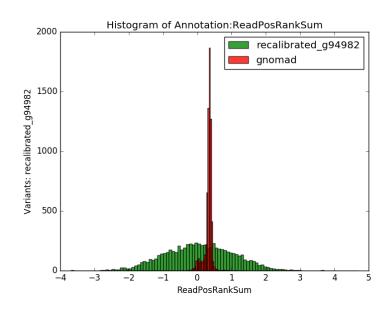- **Heuristic for the QUAL score:**
  *sum(PL[0])-30* across samples

# Combination of annotations stabilizes distributions
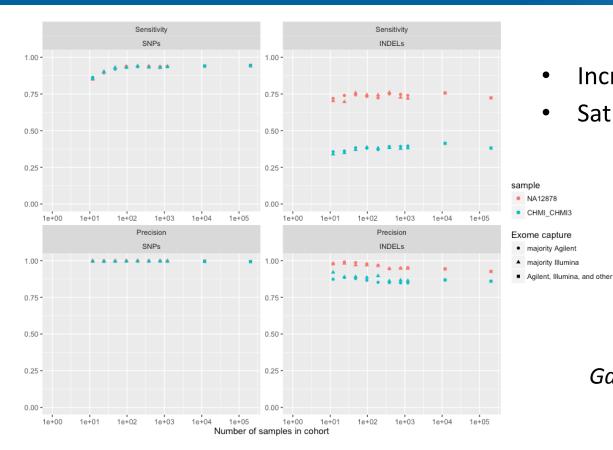


**Theoretical distribution of annotation values**

- Distinct for TP vs FP
- Tighter for common variants

**Annotation values for same variants called in :**

- Single-sample run (recalibrated_g94982, green)
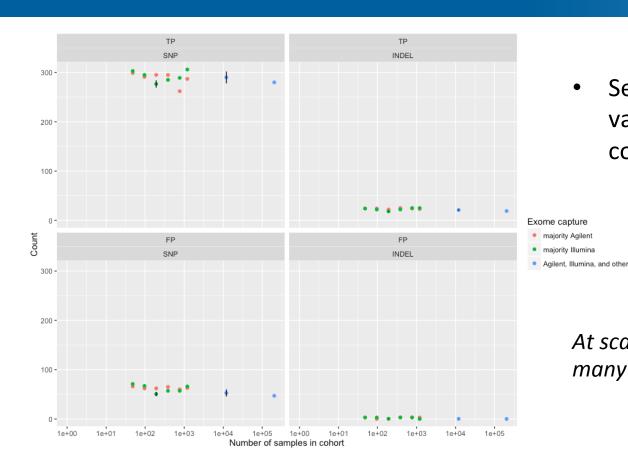- Multi-sample run (gnomad, red)

# Use of a larger cohort increases sensitivity



- Increased sensitivity
- Saturation ~600 samples

*Gauthier et al., 2016 (ASHG)*

# No loss of accuracy on singletons


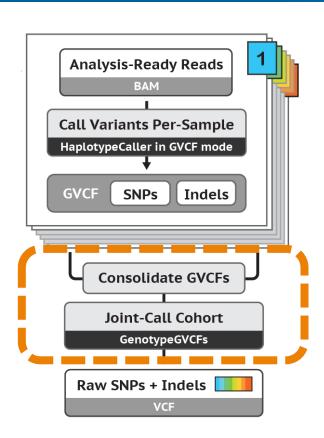
- Set of singleton truth variants compared across cohort sizes

*At scale of largest cohort, many are no longer singletons*

# Next steps: filtering and other callset refinements