# GATK Best Practices for Variant Discovery

# Introduction to
# Germline Variant Discovery
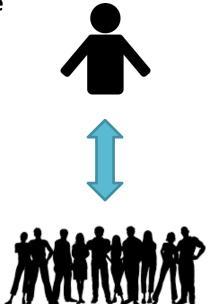
Key considerations and workflow logic
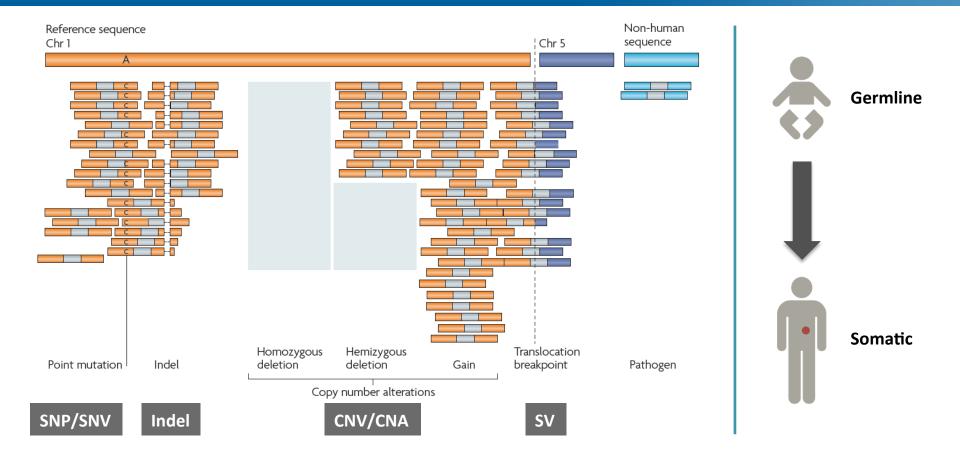
http://software.broadinstitute.org/gatk/

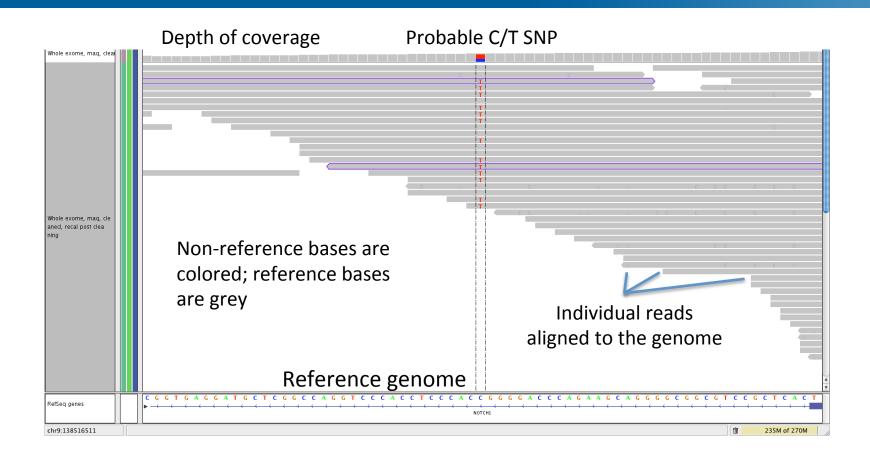# Discover **variants** relative to a reference genome

- Genetic changes in individuals **relative to a reference genome**
  - Germline (inherited)
  - Somatic (cancer)

- **Reference genome** = a standardized genomic sequence

- Human genome reference sequence
  - Previous standard: hg19 / b37
  - New standard: hg38

- Other organisms
  - Many have a fully assembled reference available
  - Many still do not -> must make one
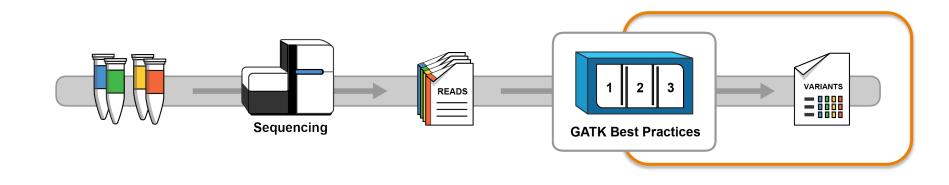
# Different types of variants



Reference sequence
Chr 1

A

Chr 5

Non-human sequence

Point mutation | Indel | Homozygous deletion | Hemizygous deletion | Gain | Translocation breakpoint | Pathogen

Copy number alterations

**SNP/SNV** | **Indel** | **CNV/CNA** | **SV**

**Germline**

**Somatic**

# This is what a good SNP looks like in a genome browser



Depth of coverage

Probable C/T SNP

Non-reference bases are colored; reference bases are grey

Individual reads aligned to the genome

Reference genome

# Short variants are reported in VCF: Variant Call Format

```
##fileformat=VCFv4.1
##reference=1000GenomesPilot-NCBI36
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
#CHROM POS     ID          REF ALT   QUAL   FILTER  INFO            FORMAT      NA00001    NA00002    NA00003
20     14370   rs6054257   G   A     29     PASS    DP=14;AF=0.5    GT:GQ:DP    0/0:48:1   1/0:48:8   1/1:43:5
20     1230237 .           T   .     47     PASS    DP=13           GT:GQ:DP    0/0:54:7   0/0:48:4   0/0:61:2
20     1234567 .           GT  G     50     PASS    DP=9            GT:GQ:DP    0/1:35:4   0/2:17:2   1/1:40:3
```

Header

Records

Format specification in
https://samtools.github.io/hts-specs/VCFv4.2.pdf

**THE WORKFLOW**
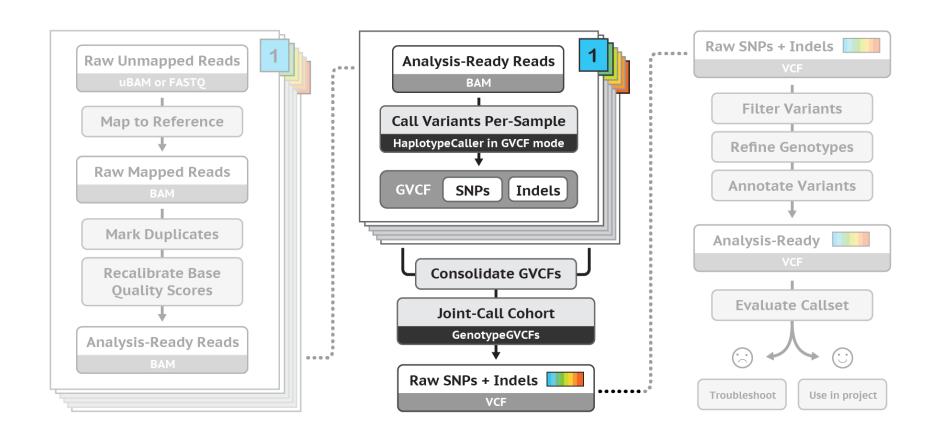
# Best Practices for Germline SNP & INDEL Discovery
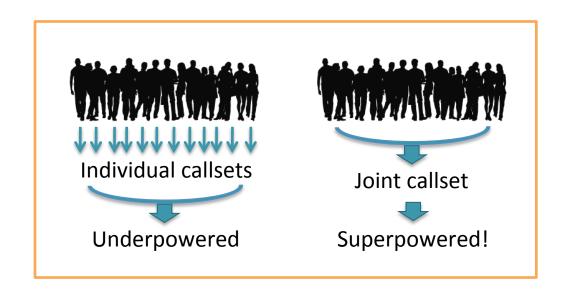
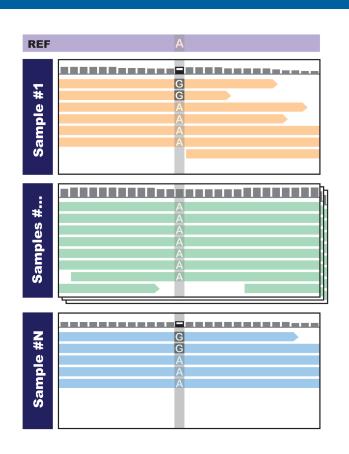# Central concept: **joint calling**

# Joint analysis empowers discovery

- Single genome in isolation: almost never useful

- Family or population data
  add valuable information

  - rarity of variants

  - *de novo* mutations

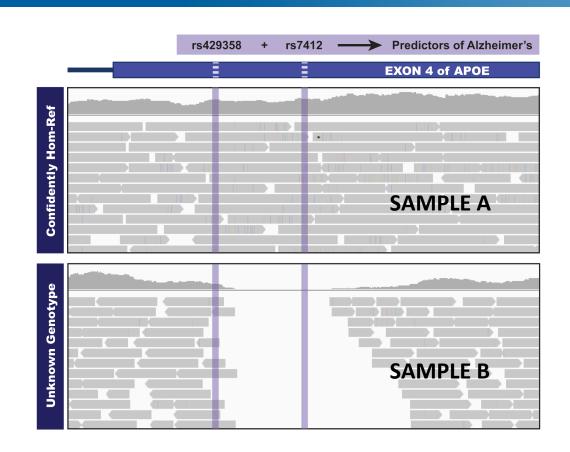  - ethnic background

# 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:
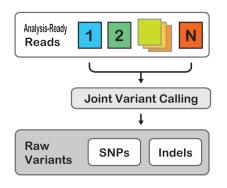  - **unlikely to be artifact**
  - **call the variant more confidently**

- **Analyzed individually:**
  - No call for either sample
  - Very different reasons!

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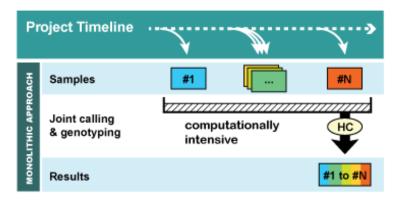  - Hom-ref call and no-call genotypes emitted

# Traditional **multi-sample calling** approach : very inefficient



**Compute requirements scale very badly with number of samples!!!**
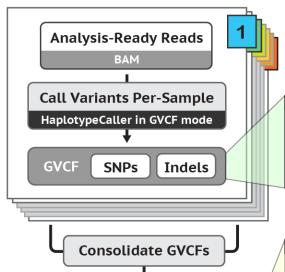
**It gives us the right answers, but…**



**Want to add new samples?**

**Got to re-run pipeline from scratch! The N+1 problem!**

# Solution: the GVCF-based joint calling workflow



Generate per-sample Genomic VCFs (GVCFs) then joint-call across all samples -> final VCF

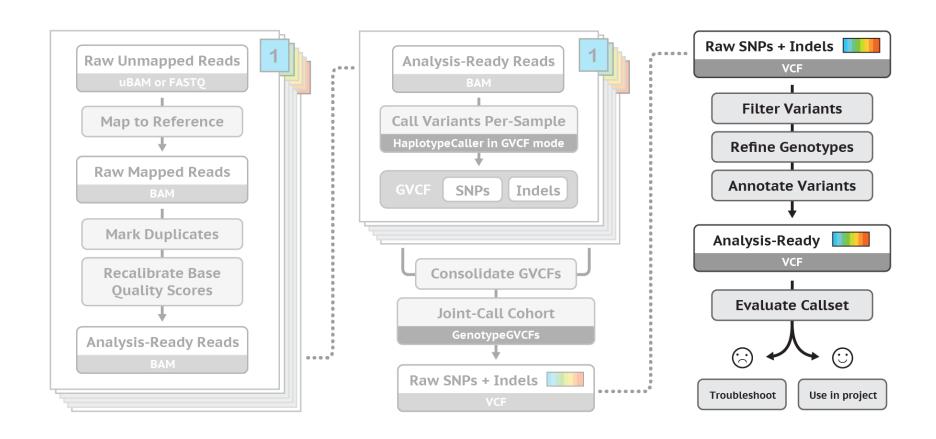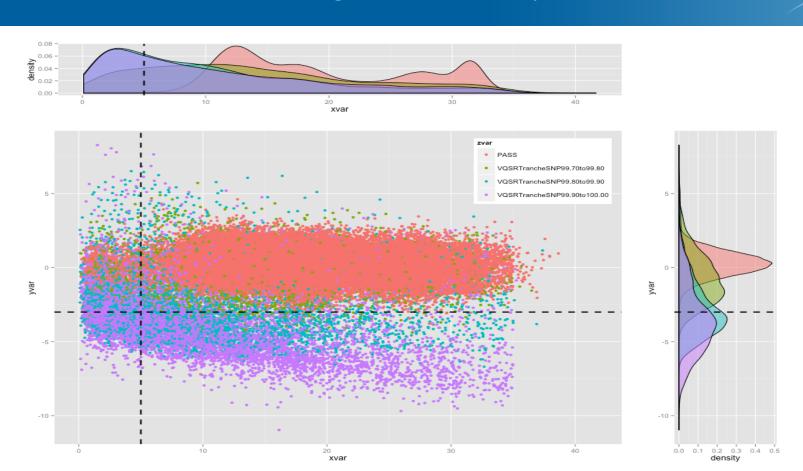# Same results as old approach - but scalable and incremental!



**Scales linearly with number of samples!**

**Want to add a new sample? Make a GVCF for that sample then re-call the cohort at will!**
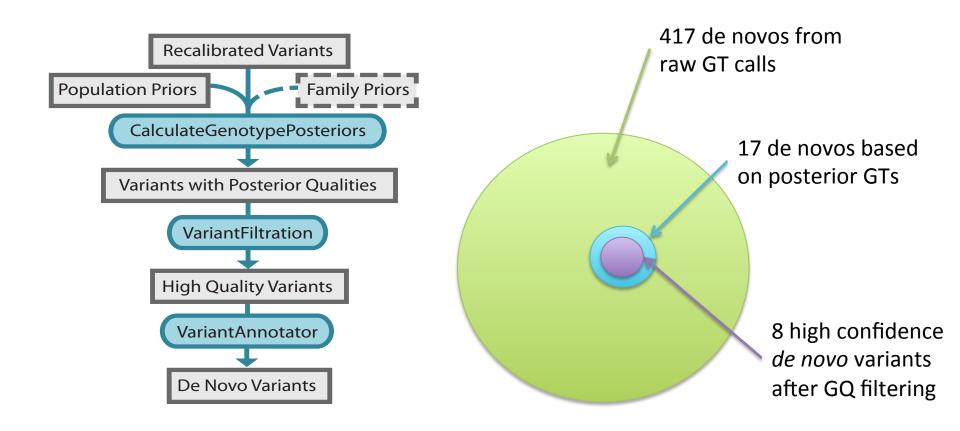
# Further refinements: filtering and more

# Variant filtering reduces false positives

# Genotype refinement improves GT quality and *de novo* calls



Recalibrated Variants

Population Priors | Family Priors

CalculateGenotypePosteriors

Variants with Posterior Qualities

VariantFiltration

High Quality Variants

VariantAnnotator

De Novo Variants

417 de novos from raw GT calls

17 de novos based on posterior GTs

8 high confidence *de novo* variants after GQ filtering

# Functional annotation predicts effects of variants

# Callset evaluation: where are you on this spectrum?

**IDEAL**

**OKAY**

**TERRIBLE**

Your variant calls perfectly match the underlying biological truth

You found many real variants and called few false positives

You didn't find any real variants and only called artifacts!

*(does not determine veracity of individual variant calls)*