# Callset Evaluation

Comparing statistics between
your callset and external resources

http://software.broadinstitute.org/gatk/

# Best Practices for Germline SNP & INDEL Discovery

# What do callset evaluation methods aim to determine?

| **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! |

Where are you on this spectrum?

*(not veracity of individual variant calls)*

# Key assumption: truth set is representative / comparable

- **Important to match dataset properties!**

  - Population ethnicity (European, African, etc.)
  - Sequencing / exp. design (WGS vs. WES)
  - Cohort size

  **Not easy!** You might need to use sub-cohorts (of both sets) to match all three.



http://www.nature.com/nature/journal/v526/n7571/full/nature15393.html
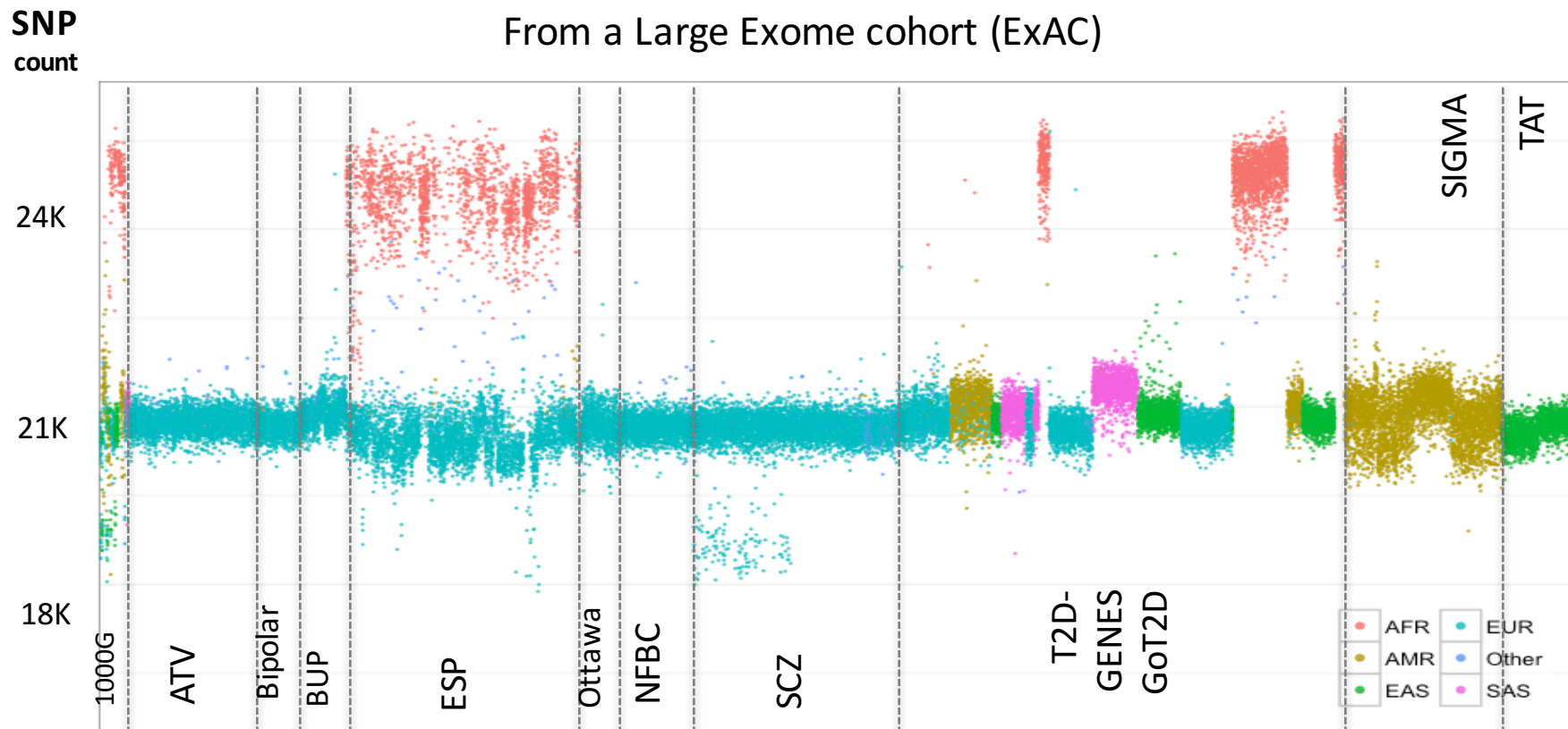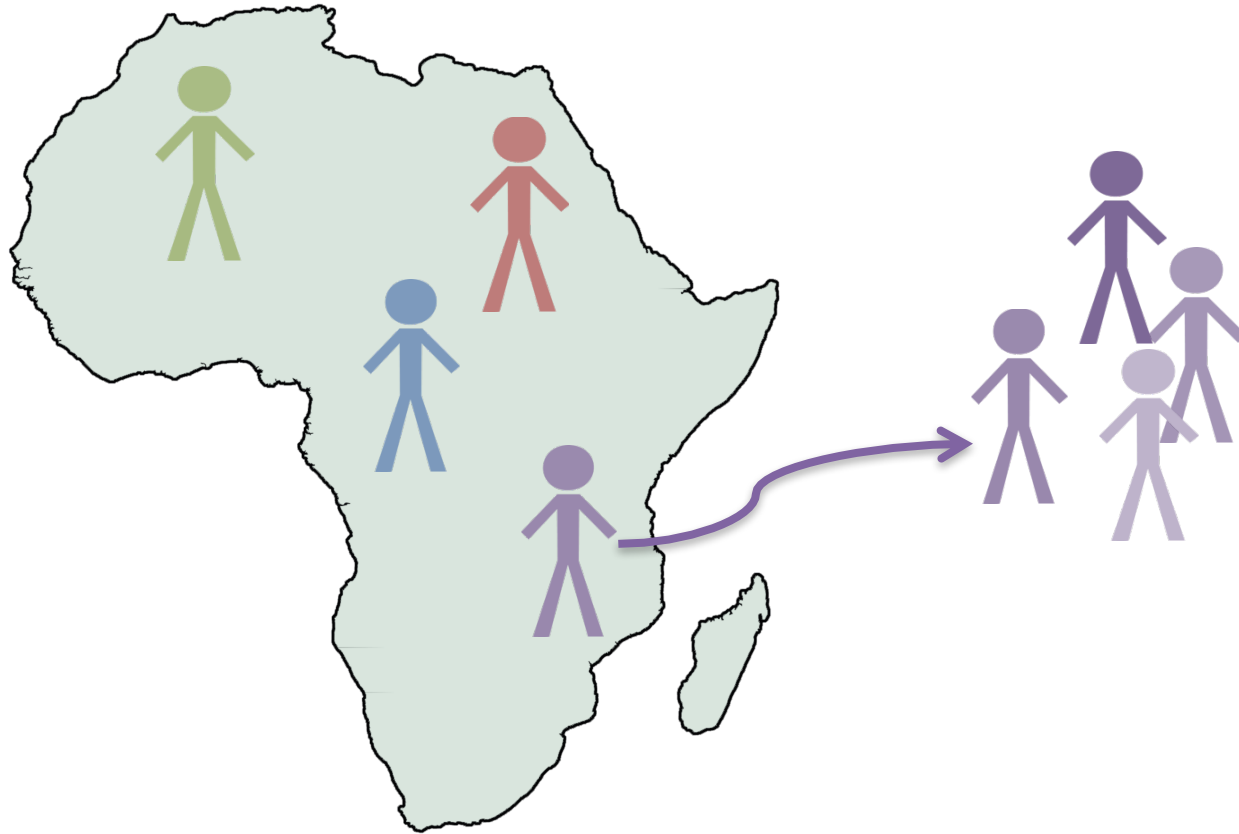
# Commonly used truth sets

- **dbSNP**
  <u>All</u> previously reported variation (lots of junk!)

- **ExAC and GnomAD**
  Extensive catalog of human variation built by aggregating results from many studies

- **HapMap**
  Highly validated common human variants

- **OMNI**
  Common variation validated by array

- **NIST's Genomes in a Bottle, or Illumina's Platinum Genomes**
  high confidence callsets from a handful of common benchmarking samples
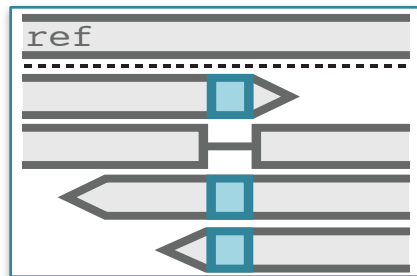
# Ethnicity affects many variant call metrics



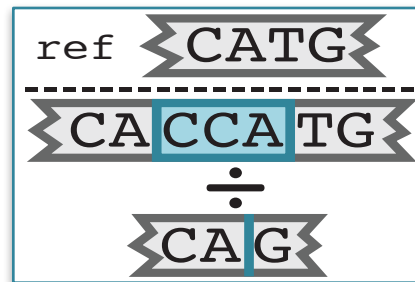From a Large Exome cohort (ExAC)

Monkol Lek, 2014

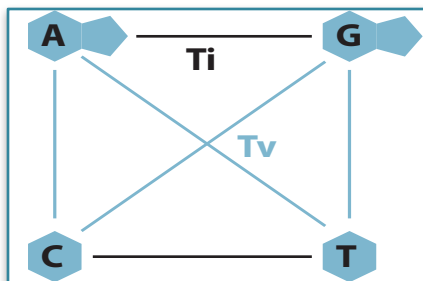# Recommended metrics for callset evaluation

## Number of Indels & SNPs



## Indel Ratio



## TiTv Ratio



## Genotype Concordance

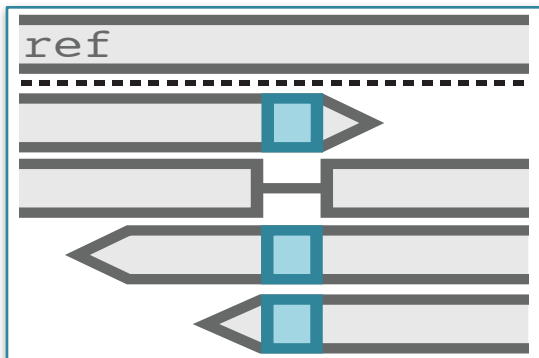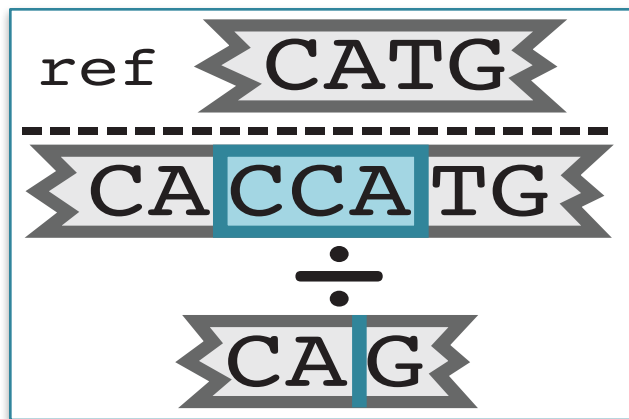| Sequencing Type | # of Variants (in 1 sample) |
|---|---|
| WGS | ~4.4 M |
| WES | ~21 k |

- Variants = Indels + SNPs
- Useful for order-of-magnitude sanity check
- Vary by size and diversity of cohort

# Indel Ratio



| Variant prevalence | Indel Ratio |
|---|---|
| Common | ~1 |
| Rare | 0.2-0.5 |

- Ratio of **in**sertions to **del**etions
- Varies by allele frequency: common ("known") vs. rare ("novel")

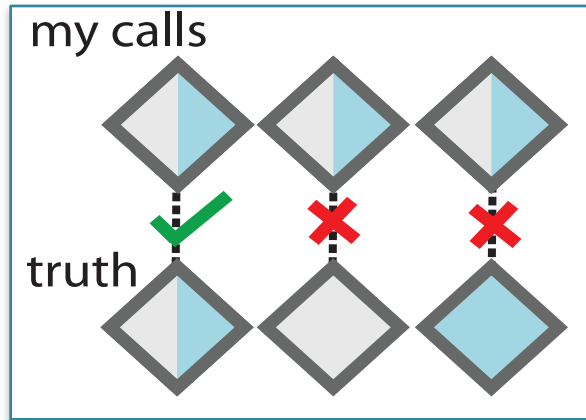# TiTv Ratio (Transitions/Transversions)



| Sequencing Type | TiTv Ratio |
|---|---|
| WGS | 2.0-2.1 |
| WES | 3.0-3.3 |

In Humans...

- Used for SNPs only

- If variation were random: expect ratio of 0.5 as there are twice as many possible transversions vs transitions!

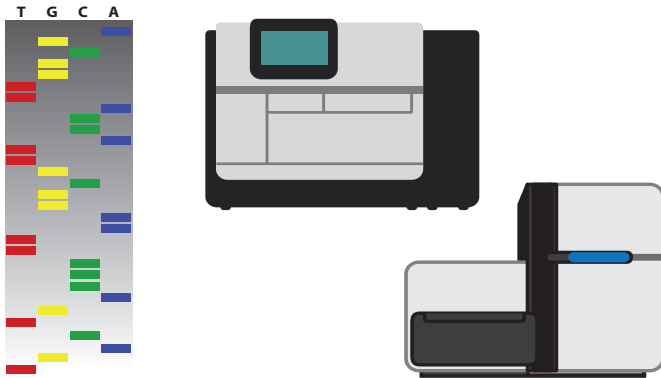- Low TiTv ratio indicates high rate of false positives

- Most appropriate truth set is genotyping chip for same sample

- % Genotype calls in callset that match calls in truth set

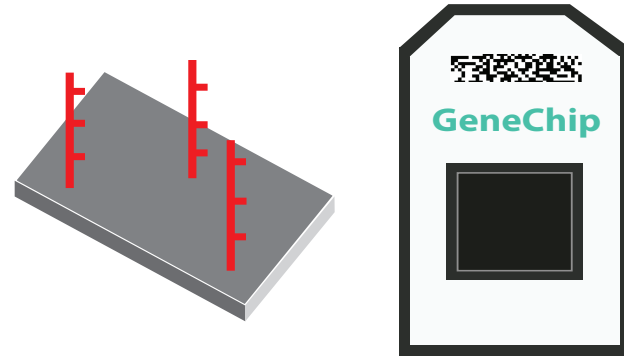- Unmatched variants considered false positives

Variant Concordance

# So how do I get these metrics?

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

- Genotype level evaluation tools equivalent—these tools will be merged in GATK4

- Variant level evaluation tools are different

# Which variant-level evaluator should I use?

## GATK VariantEval

- More detailed analysis
- More options for stratification
- Ability to compare to multiple truth sets

## Picard CVCM

- Best performance on very large callsets
- Ability to interpret no-call as confident reference in a "confidence region"
- Few options beyond the metrics discussed here

Best Practices for Germline SNP & INDEL Discovery