



Universitair Medisch Centrum Utrecht

# Mapping sequence reads & Calling variants

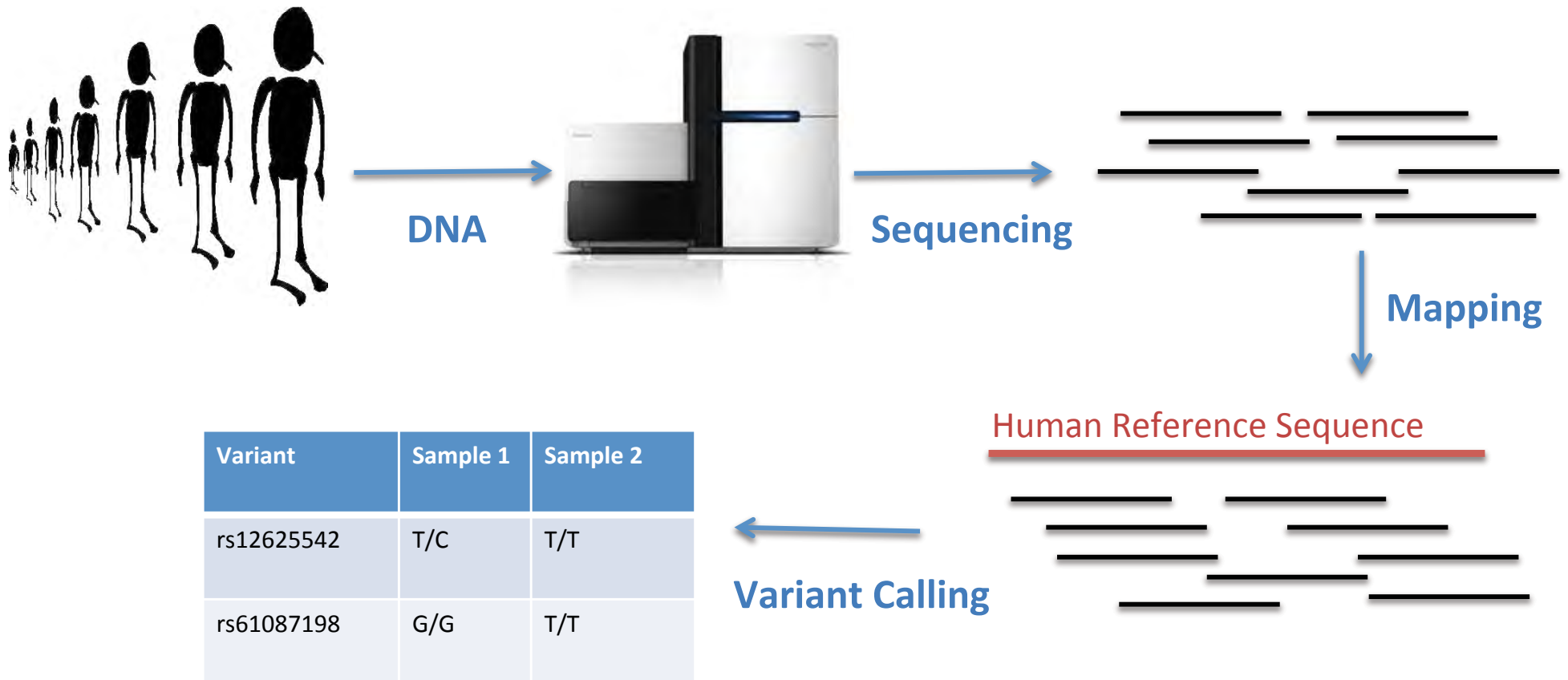
---

Laurent Francioli

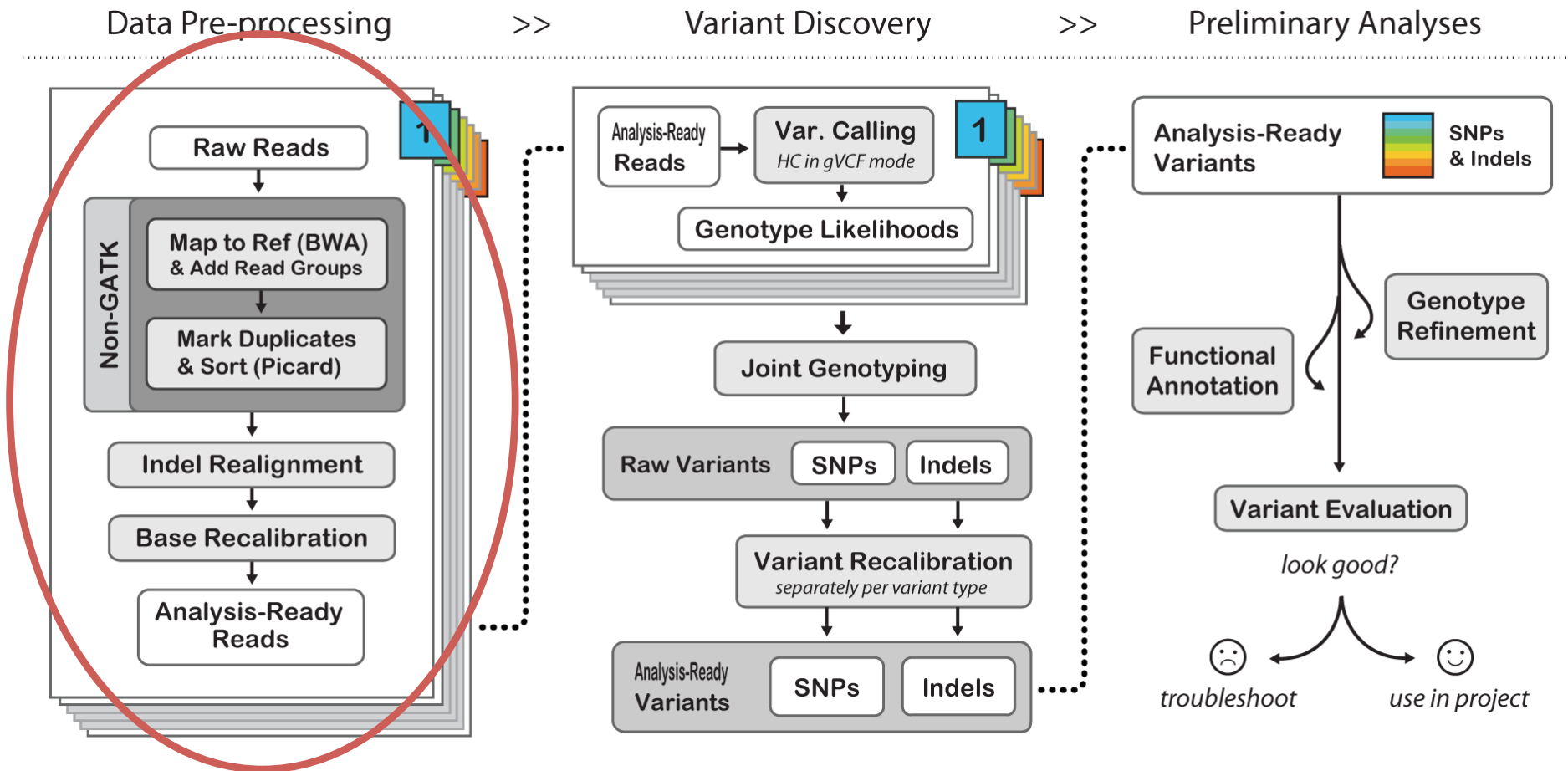
2014-10-28

[l.francioli@umcutrecht.nl](mailto:l.francioli@umcutrecht.nl)

# Next Generation Sequencing



# Data processing pipeline



From the Genome Analysis Toolkit: <http://www.broadinstitute.org/gatk/guide/best-practices>

# Alignment



# Mapping to reference

- Burrow-Wheeler Algorithm (BWA)
- Reads are scored for possible alignments against reference data
- Mapping scores are computed using mismatches and gap penalties

## Human Reference Sequence

GTGCCAGGACCAGATCGTGCCAACGGACAGGTGGTAAGGAAGGAG

## Sequence reads

GTGCCAAGGA

CAACGGACAG

AGGTGGTAAG

...

# Mapping to reference

- Burrow-Wheeler Algorithm (BWA)
- Reads are scored for possible alignments against reference data
- Mapping scores are computed using mismatches and gap penalties

## Human Reference Sequence

GTGCCAGGACCAGATCGTGCCAACGGACAGGTGGTAAGGAAGGAG  
GTGCCAAGGA

GTGCCAAGGA

# Mapping to reference

- Burrow-Wheeler Algorithm (BWA)
- Reads are scored for possible alignments against reference data
- Mapping scores are computed using mismatches and gap penalties

## Human Reference Sequence

GTGCCA - GGACCAGATCGTGCCAACGGACAGGTGGTAAGGAAGGAG

GTGCCA - AGGA

GTGCCAAGGA

GTGCCAAGGA

GTGCCAA - GGA

# Mapping to reference

- Burrow-Wheeler Algorithm (BWA)
- Reads are scored for possible alignments against reference data
- Mapping scores are computed using mismatches and gap penalties

## Human Reference Sequence

...GTGCCAGGACCAAGATCGTGCCAACGGACAGGTGGTAAGGAAGGA...



3.2Gb

## Millions of sequence reads

...GTGCCAAGGA...

...CAACGGACAG...

...AGGTGGTAAG...

...



50 – 120bp



# Coverage

- Number of reads covering a given base of the genome
- Overall coverage is a function of the total number of sequence reads
  - Target coverage is decided at experiment design

## Human Reference Sequence

GTGCCAGGACCAGATCGTGCCAACGGACAGGTGGTAAGGAAGGAG

GTGCCAGGAC

CCAGGATC

GATCAGATCG

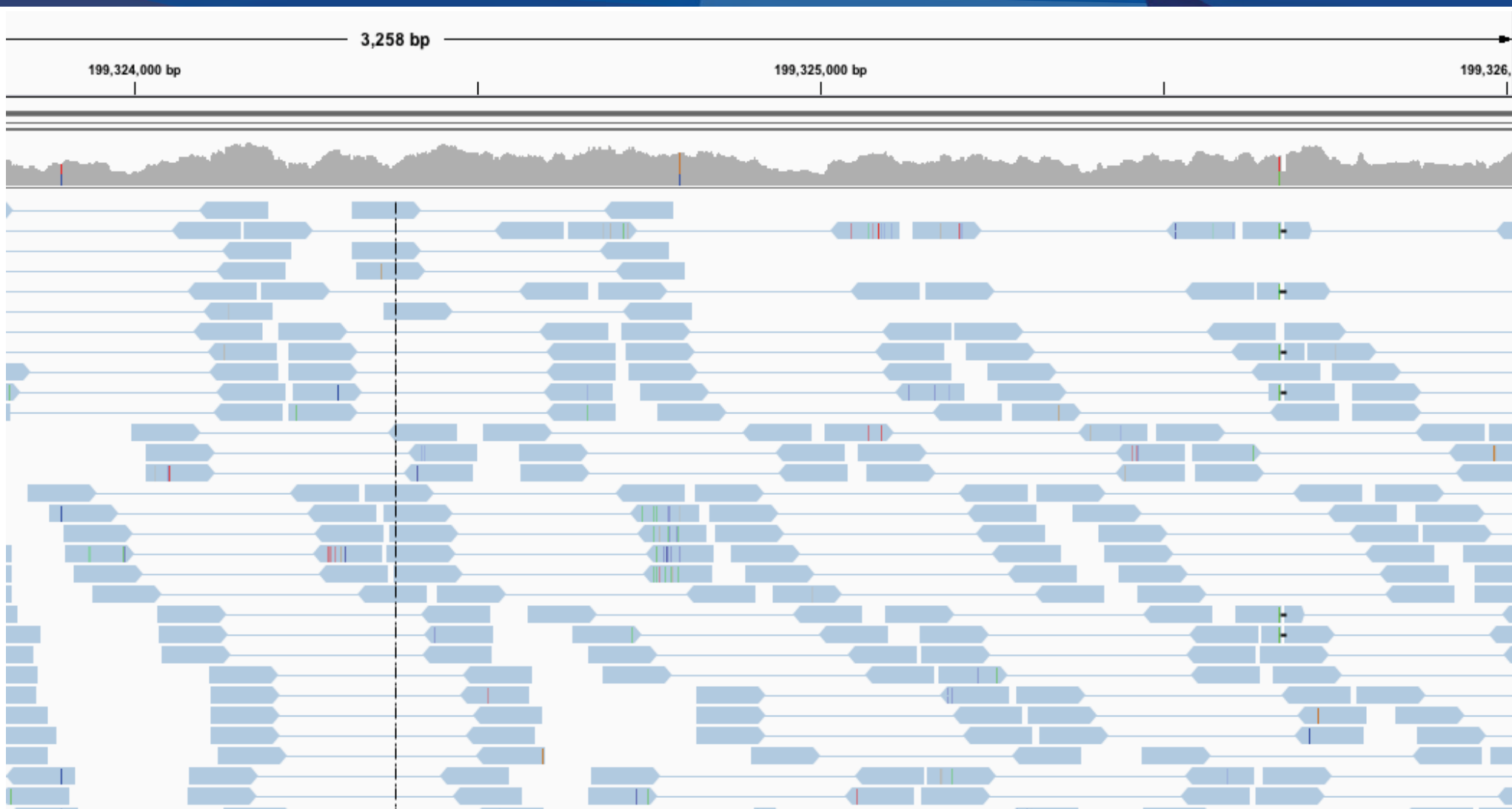
ACCAGATCGT

GTGCCAA - GGA

CAACGGACAG

GGAAGGAG

# Real NGS aligned data in IGV



# Marking Duplicates

✘ = sequencing error propagated in duplicates



FP variant call  
(bad)

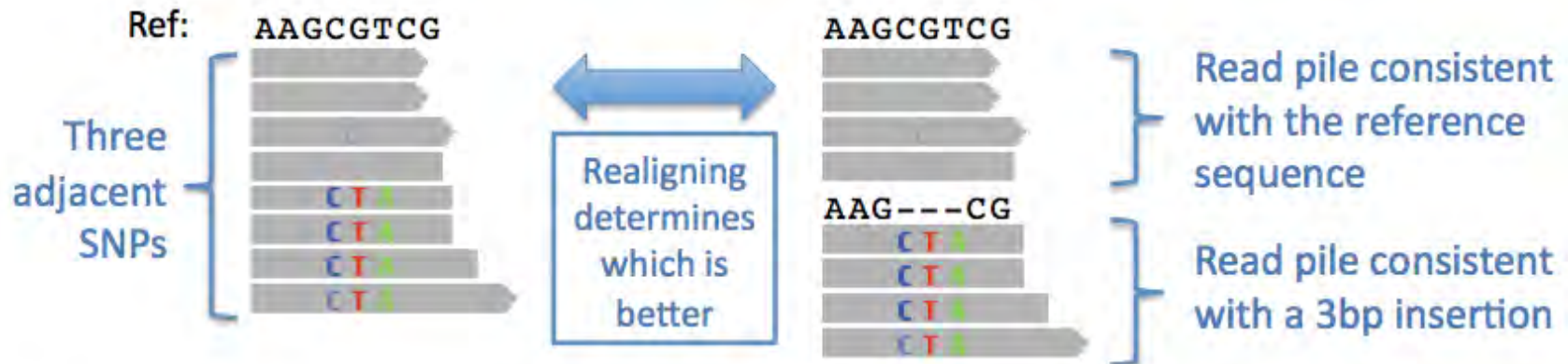
After marking duplicates, the GATK will only see :



... and thus be more likely to make the right call

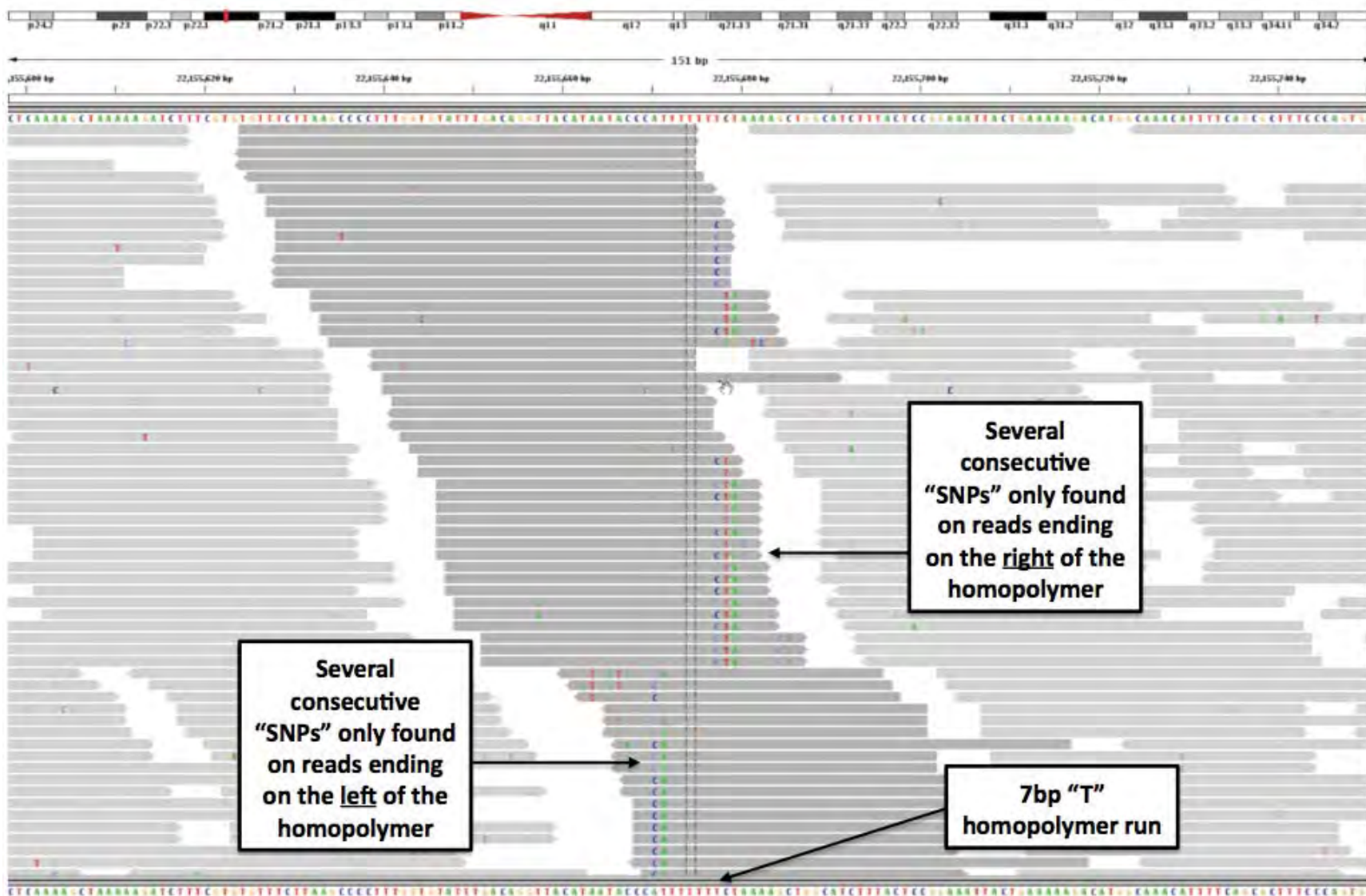
# Indel realignment

1. Find the best alternate consensus sequence that, together with the reference, best fits the reads in a pile (maximum of 1 indel)



2. The score for an alternate consensus is the total sum of the quality scores of mismatching bases
3. If the score of the best alternate consensus is sufficiently better than the original alignments (using a LOD score), then we accept the proposed realignment of the reads

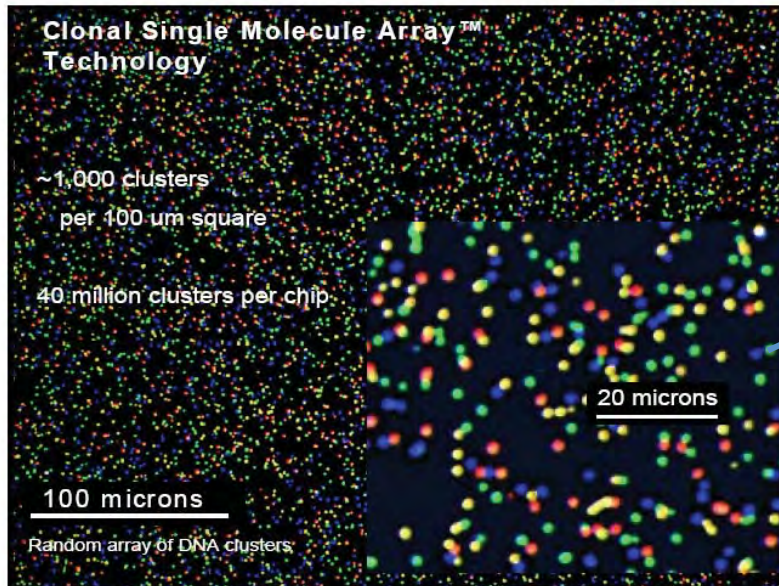
# Indel realignment - example



# Indel realignment - example



# Base qualities



ACTGCCAGGT**N**TCAGTACA

**Phred value** =  $-10 \log_{10}(E)$

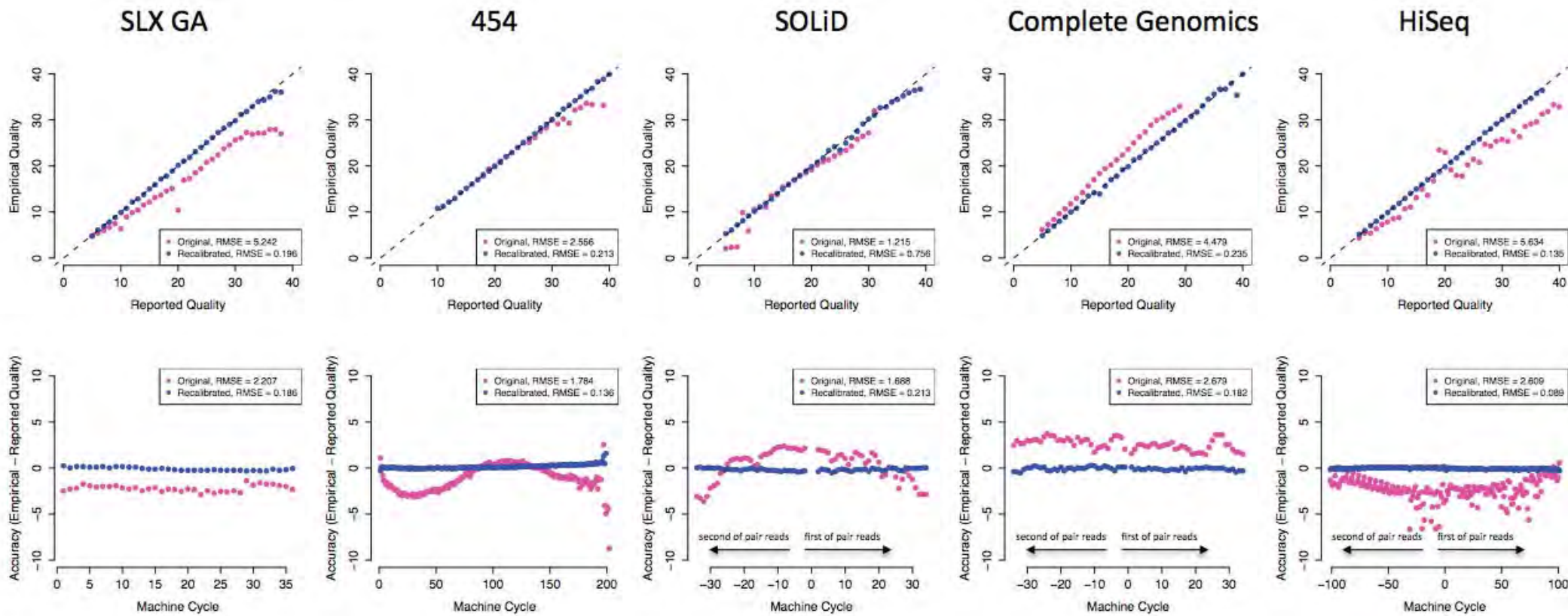
Quality 10 = 10% error (90% confidence)

Quality 20 = 1% error (99% confidence)

Quality 30 = 0.1% error (99.9% confidence)

# Base quality recalibration

- Downstream models use base qualities as input
- Base qualities have systematic machine-specific biases

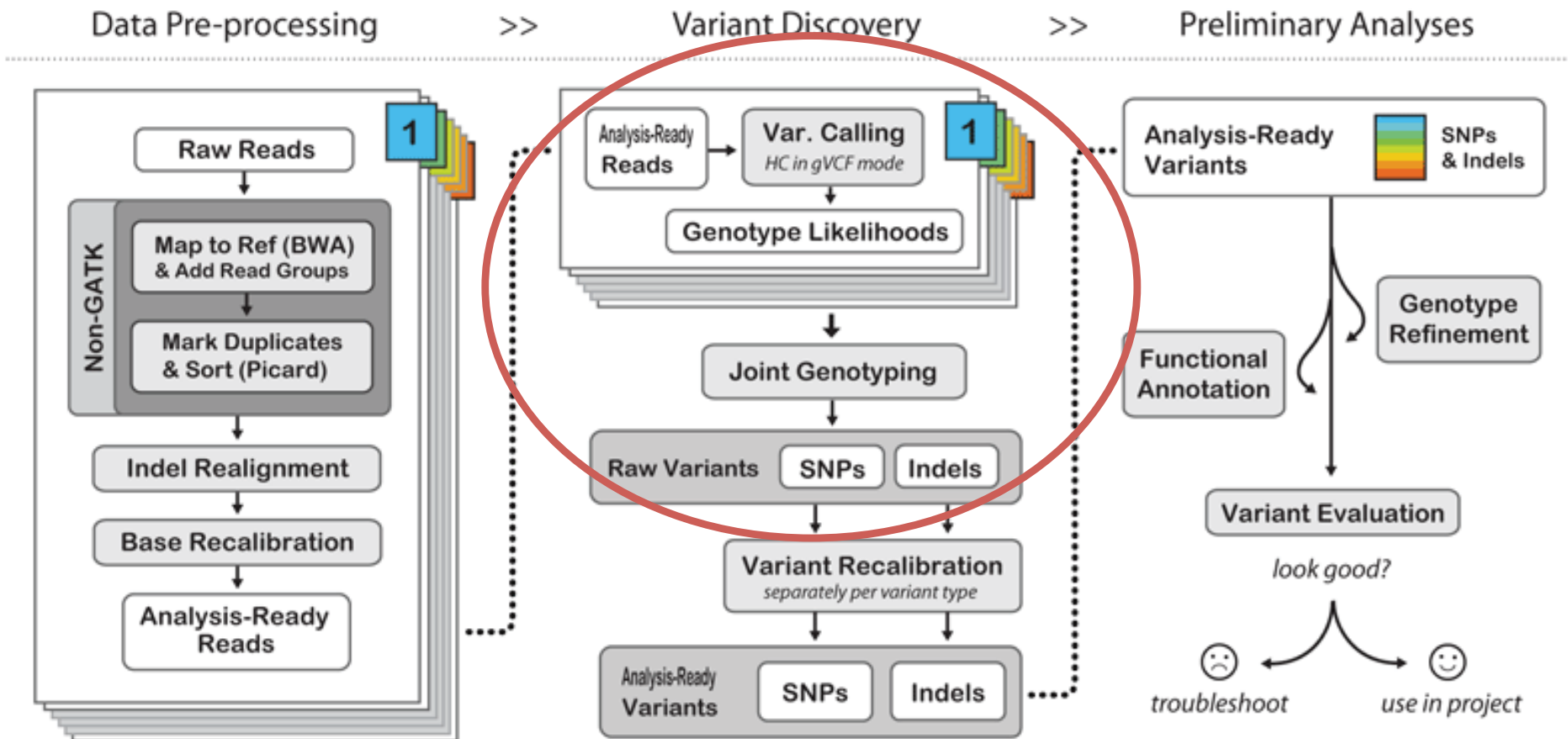




# Analysis-ready reads

- Reads are correctly placed in the genome
  - Mapping qualities model the uncertainty of mapping
- Reads are independent
  - Clonal (duplicate) reads are marked
- Base qualities are calibrated
  - Base qualities model the uncertainty of the base call

# Data processing pipeline



# SNPs and indels

- Single Nucleotide Variants (SNVs)

ATATTATTGCCAGACCAATGTCTGGAGTTATTCCCCTGT

ATATTATTACCAGACCAATGTCTCGGAGTTATTCCCCTGT

- Insertions

ATATTATTGCCAGACCAATGTCTGGAGTTATTCCCCTGT

ATATTATTGCCAGACCAGTTATGTCTGGAGTTATTCCCCTGT

- Deletions

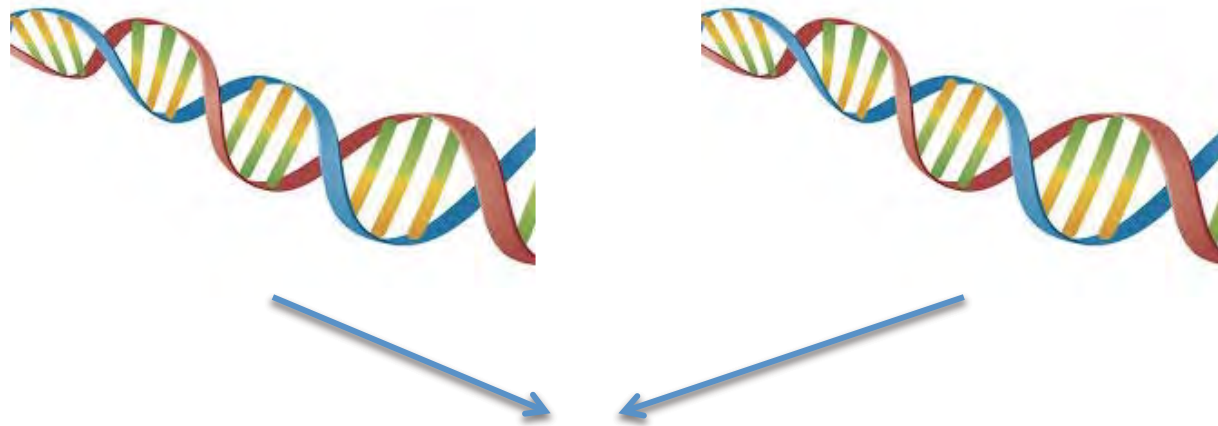
ATATTATTGCCAGACCAATGTCTGGAGTTATTCCCCTGT

ATATTATTGCCAGACCAATGGAGTTATTCCCCTGT

# NGS data is noisy



# Diploid is difficult



Chromosome pair reads are mixed in sequence data

- **Homozygous**

- All reads should show the same allele

- **Heterozygous**

- ~50% of the reads should show each allele

# SNP Calling and Genotyping

- For each base in the genome, find the sample genotype given the observed reads.
- NGS reads contain errors (typically ~1% - 0.1% bases)

## Human Reference Sequence

...GTGCCAGGACCAGATCG...

```
...GTGCCAGG
...GTGCCAGG
...GTGCCATGA
GTGCCAGGACC
  GCCAGGACCAGA
  GCCAGGACCAGAT
  CCAGGACCAGA
  CCAGGACCAGAT
    AGGACCAGATCG
    GGACCAGATCG...
```

## Homozygous reference (GG)

- 9 correct reads
- 1 error

## Heterozygous (GT)

- 9 Correct reads (allele G)
- 1 Correct read (allele T)

## Homozygous non-reference (TT)

- 1 correct reads
- 9 errors

# SNP Calling and Genotyping

- For each base in the genome, find the sample genotype given the observed reads.
- NGS reads contain errors (typically ~1% - 0.1% bases)

## Human Reference Sequence

...GTGCCAGGACCAGATCG...

```
...GTGCCAGG
...GTGCCAGG
...GTGCCATGA
GTGCCAGGACC
GCCATGACCAGA
GCCAGGACCAGAT
CCATGACCAGA
CCAGGACCAGAT
AGGACCAGATCG
TGACCAGATCG...
```

## Homozygous reference (GG)

- 6 correct reads
- 4 errors

## Heterozygous (GT)

- 6 correct reads (allele G)
- 4 correct reads (allele T)

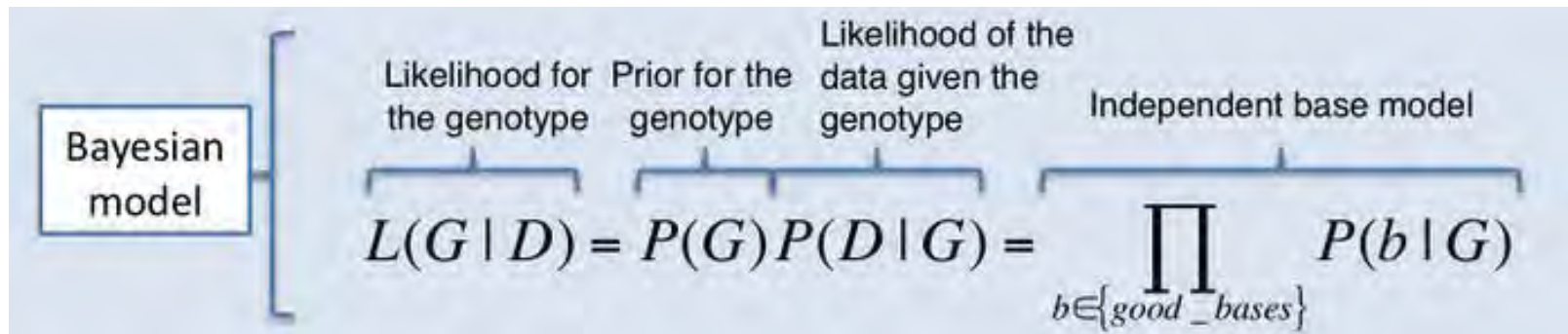
## Homozygous non-reference (TT)

- 4 correct
- 6 errors

# Genotyping Formalism (GATK)

## Bi-Allelic Model

- *A*: Reference allele
- *B*: Alternate allele



**Homozygous**

$$P(b | AA) = \begin{cases} 1 - e_b, & \text{if } b = A \\ e_b / 3, & \text{if } b \neq A \end{cases}$$

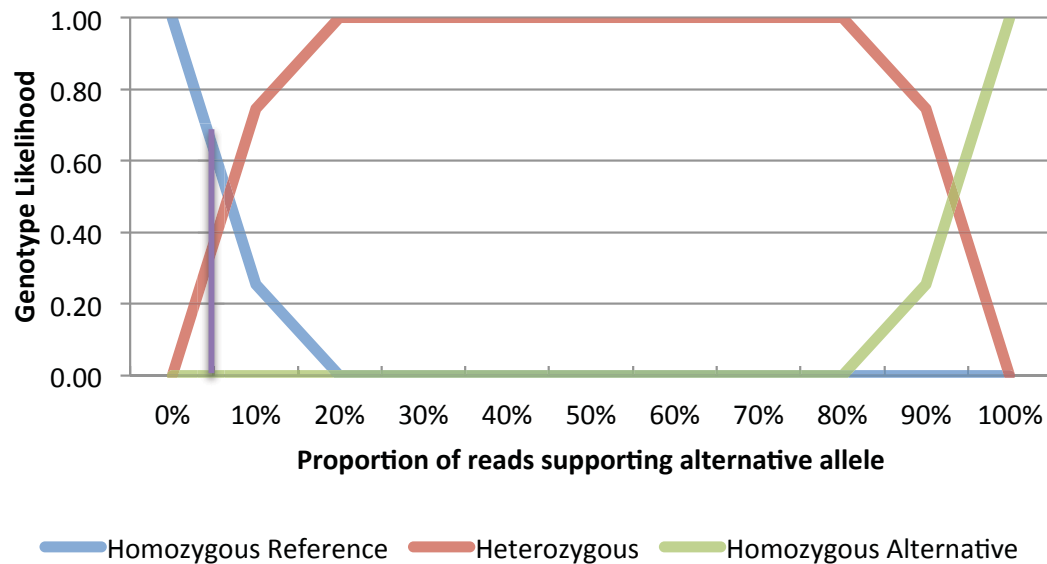
**Heterozygous**

$$P(b | AB) = \begin{cases} (1 - e_b) / 2 + e_b / 6, & \text{if } b \in \{A, B\} \\ e_b / 3, & \text{if } b \notin \{A, B\} \end{cases}$$



# Interpreting NGS genotypes

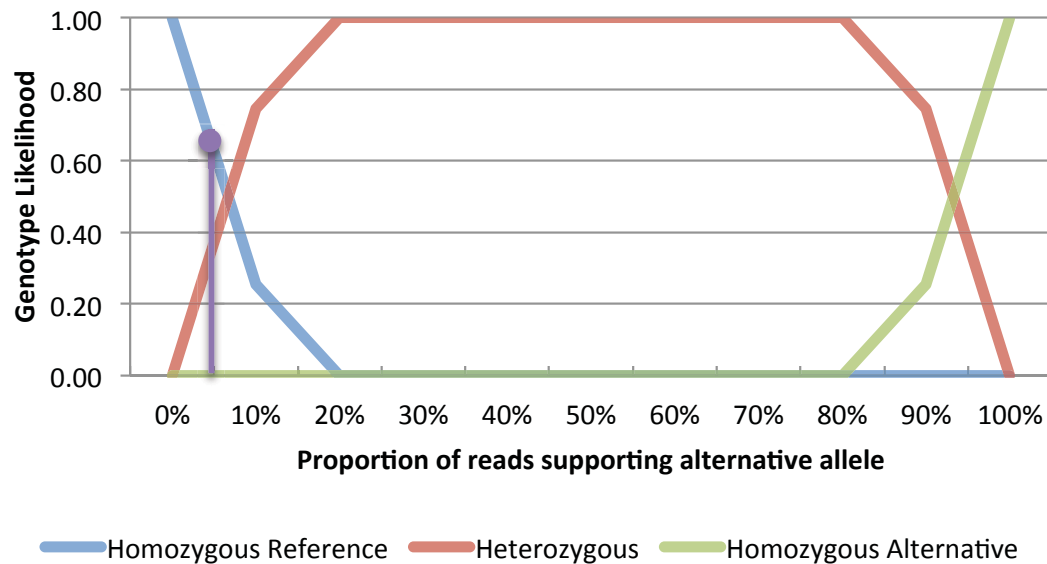
SNP Genotype confidence for a 12x coverage site



Genotyping a site given N supporting reads

# Interpreting NGS genotypes

SNP Genotype confidence for a 12x coverage site

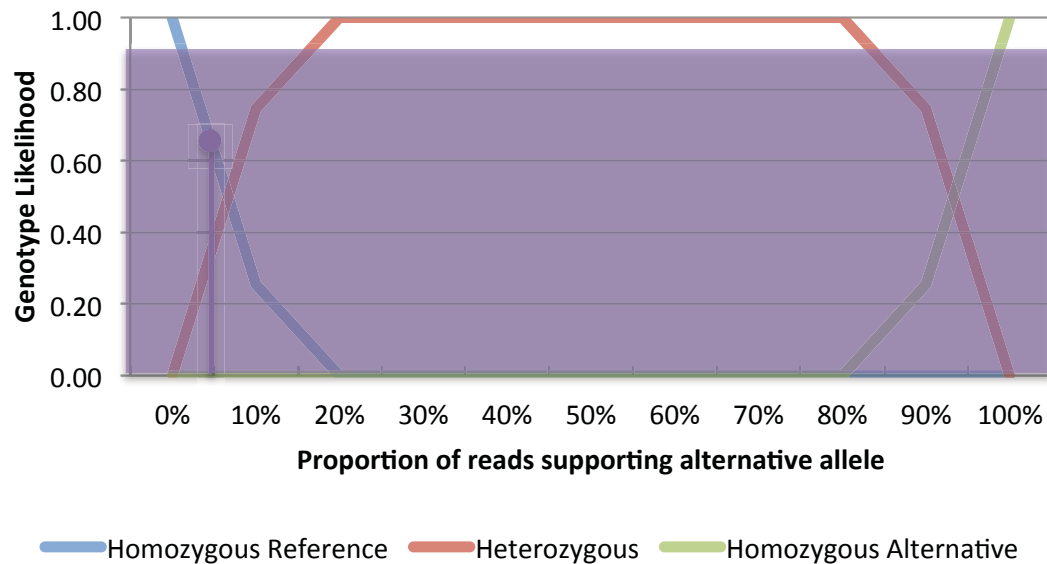


Genotyping a site given N supporting reads

- **Most likely genotype**

# Interpreting NGS genotypes

SNP Genotype confidence for a 12x coverage site

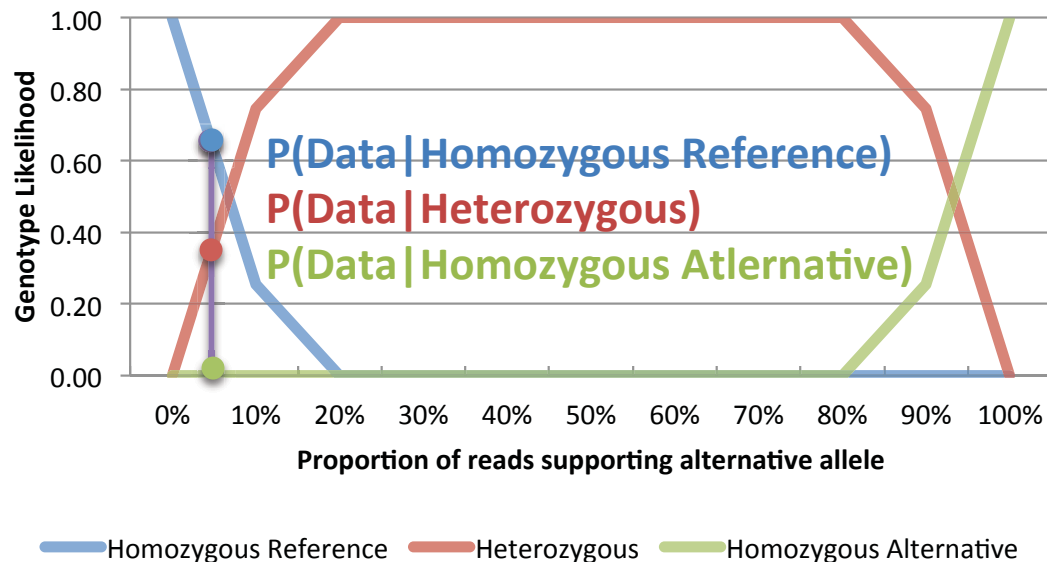


Genotyping a site given N supporting reads

- Most likely genotype
- **Most likely genotype, only if within set confidence threshold**

# Interpreting NGS genotypes

## SNP Genotype confidence for a 12x coverage site



Genotyping a site given N supporting reads

- Most likely genotype
- Most likely genotype, only if within set boundaries
- **Genotype likelihoods for all possible genotypes**

# Two approaches for calling SNPs and indels

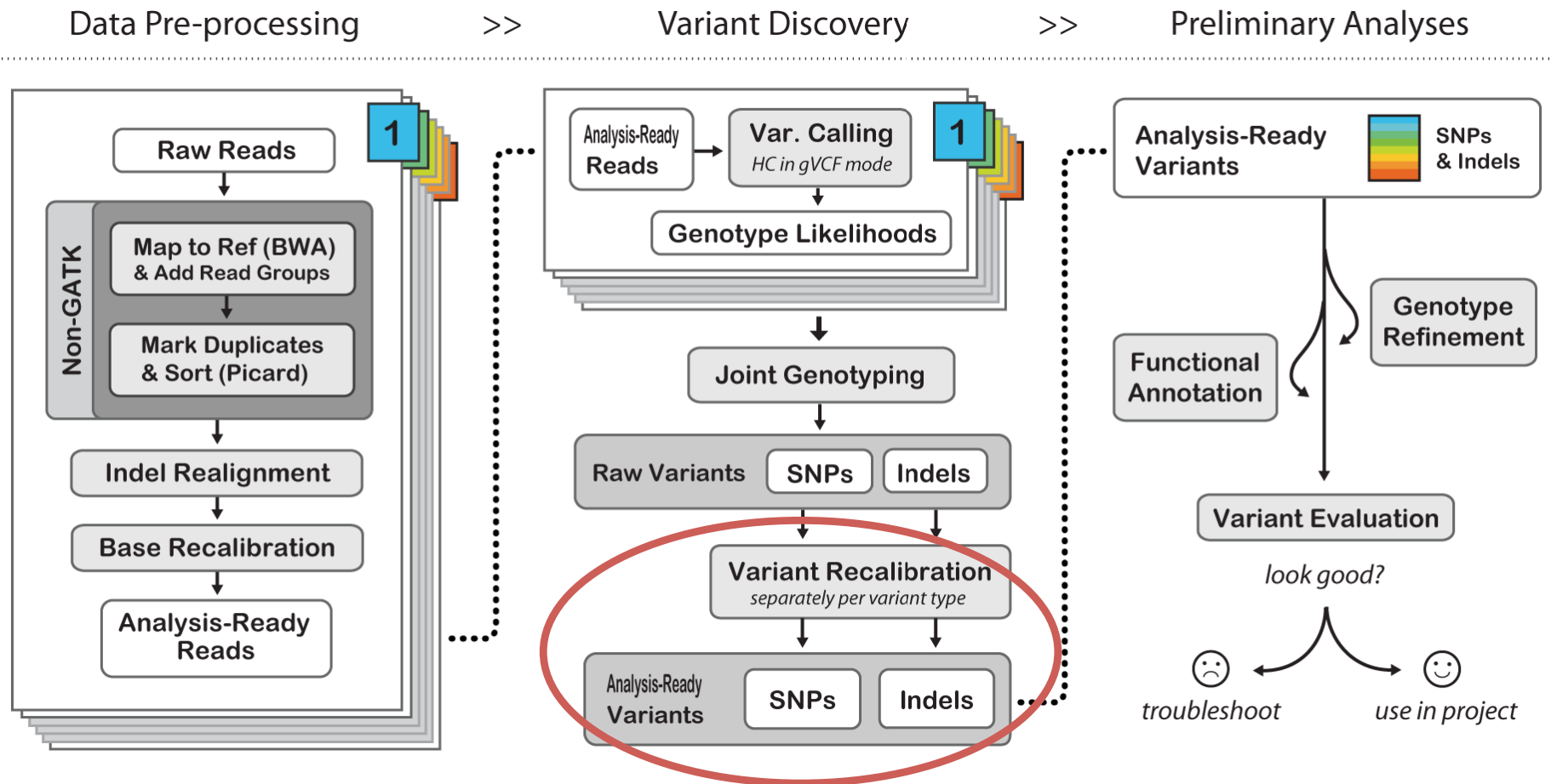
## Unified Genotyper

- Calls SNPs and INDELS **separately**
- Evaluates each base **independently**
  - Fast!

## Haplotype Caller

- Calls SNPs and INDELS **simultaneously**
- Uses all bases connected through local *de novo* assembly
  - No need for indel realignment
  - More accurate, especially for indels
  - Supports “N+1” genotyping

# Data processing pipeline



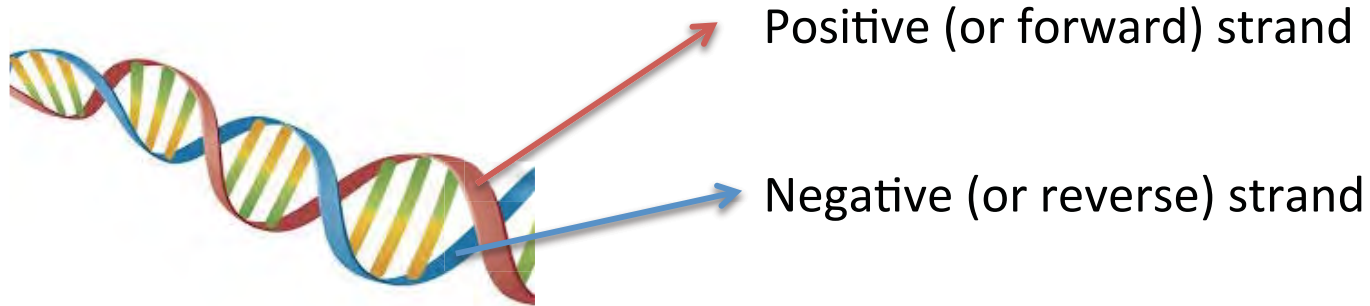
From the Genome Analysis Toolkit: <http://www.broadinstitute.org/gatk/guide/best-practices>

# Variant Annotations

- Reads
  - Strand bias
  - Quality / Depth
  - Mapping quality ranksum
  - ...
- Biology / Population
  - Hardy-Weinberg Equilibrium
  - Haplotype Score

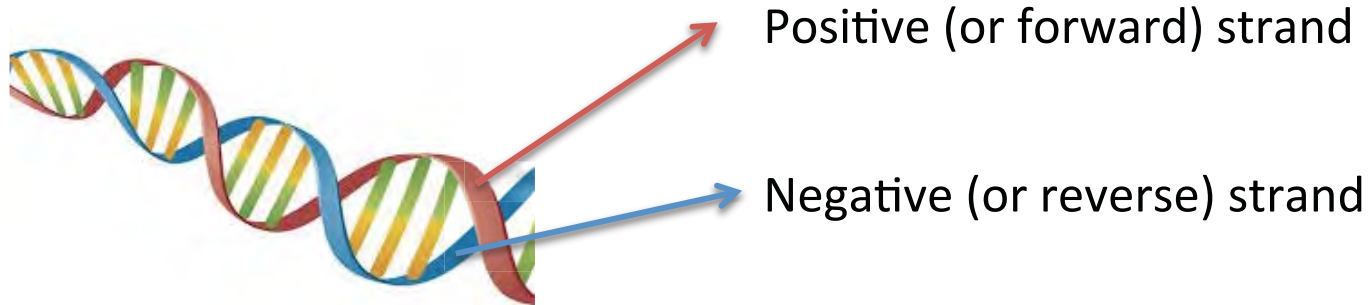
See GATK guidelines for a full list of recommended annotations to use!

# Example: strand bias





# Example: strand bias



Allele	Positive strand	Negative strand
A	30	5
C	35	30

Fisher's test p-value = 0.002

# Two ways of filtering

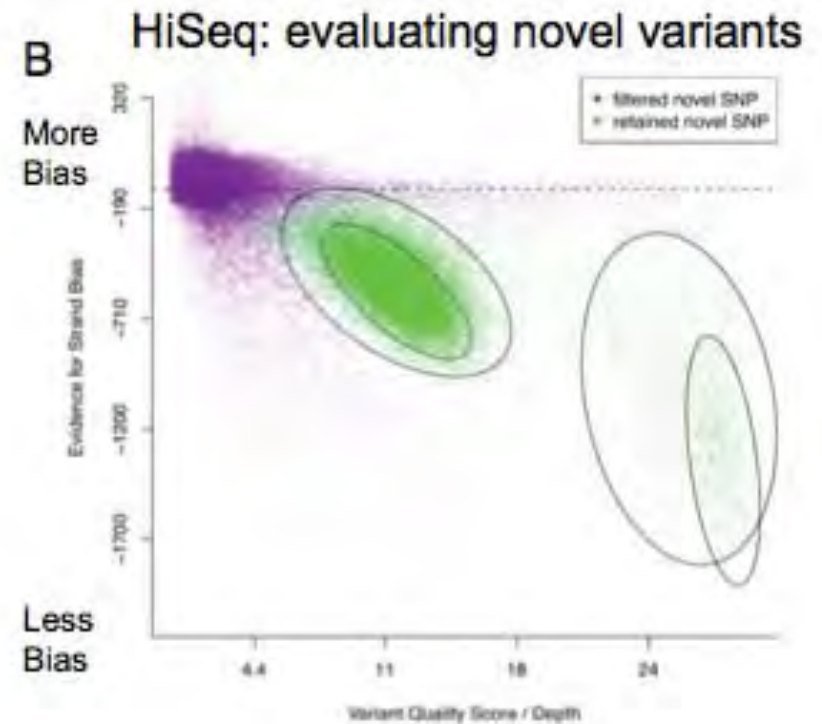
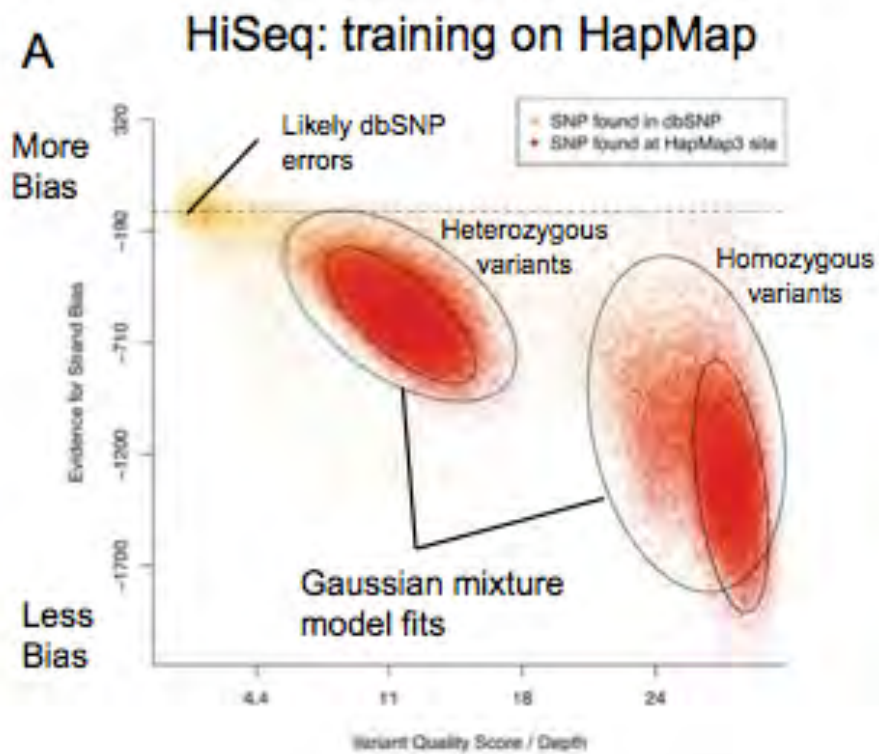
## Thresholds

- Pros
  - Work for small datasets
  - No need for external resources
- Cons
  - Does not adapt to particular data
  - Manual tuning
  - Human bias

## Machine-learning

- Pros
  - Multi-dimensional and complex model can fit the data better
  - Ranking of variants using a single metric
- Cons
  - Need training data
  - Need enough data

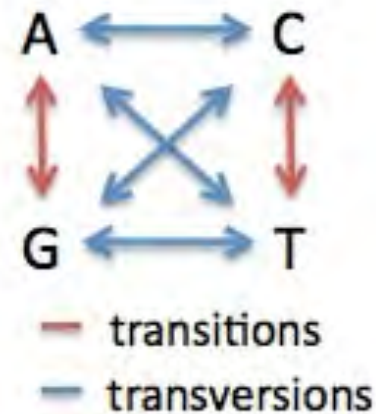
# Machine learning filtering



DePristo *et al.*, Nature Genetics 2011

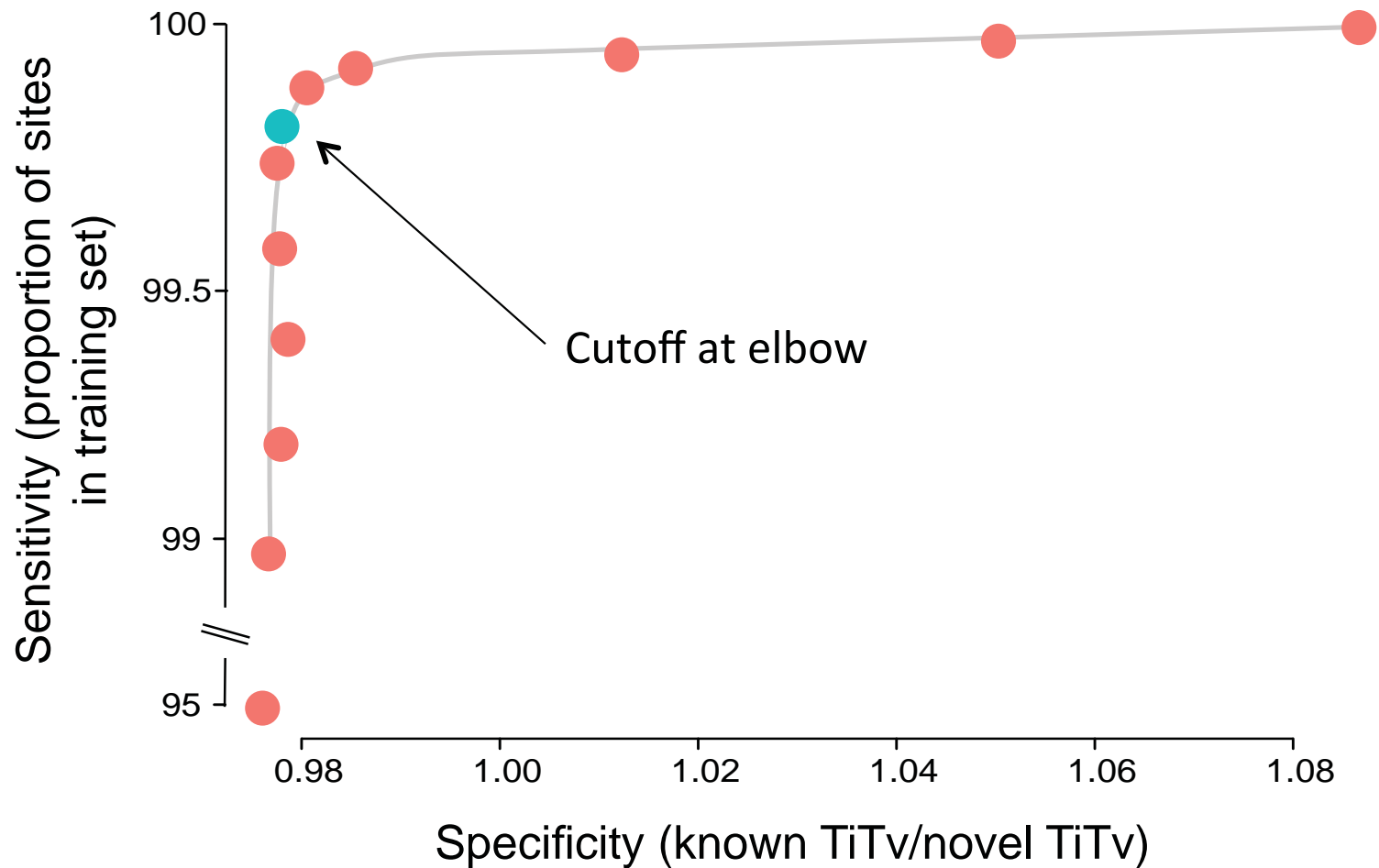
# Filtering Evaluation

- Transition / Transversion (Ti/Tv)
- Proportion of novel variants
- Heterozygosity
- Allele frequency spectrum of variants
- Genotype concordance against chip (if available)

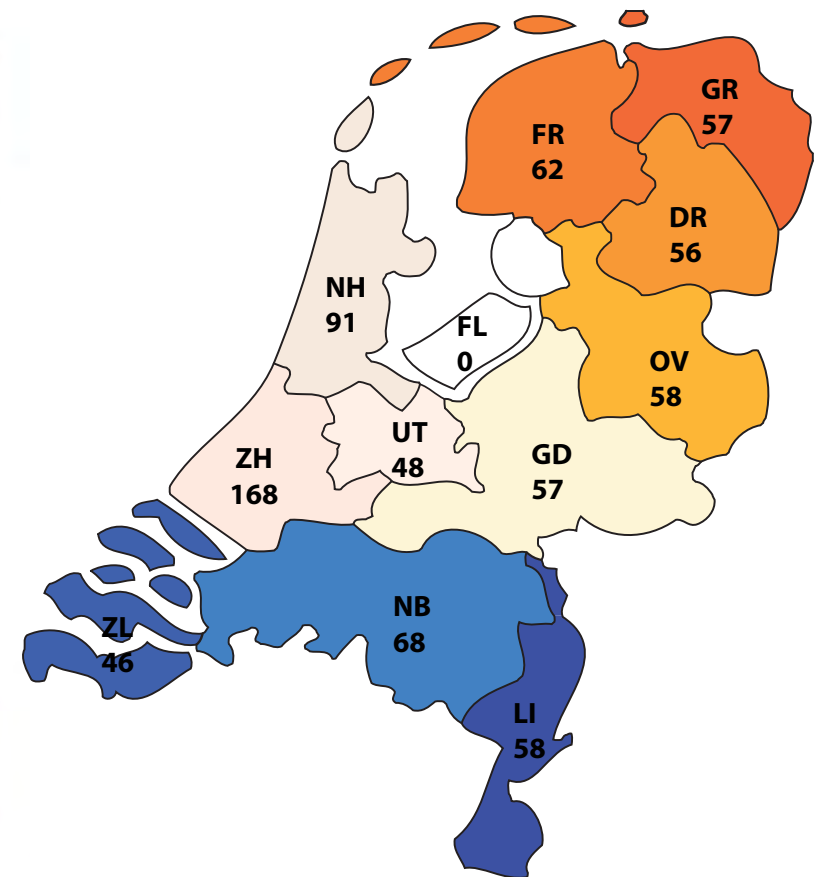
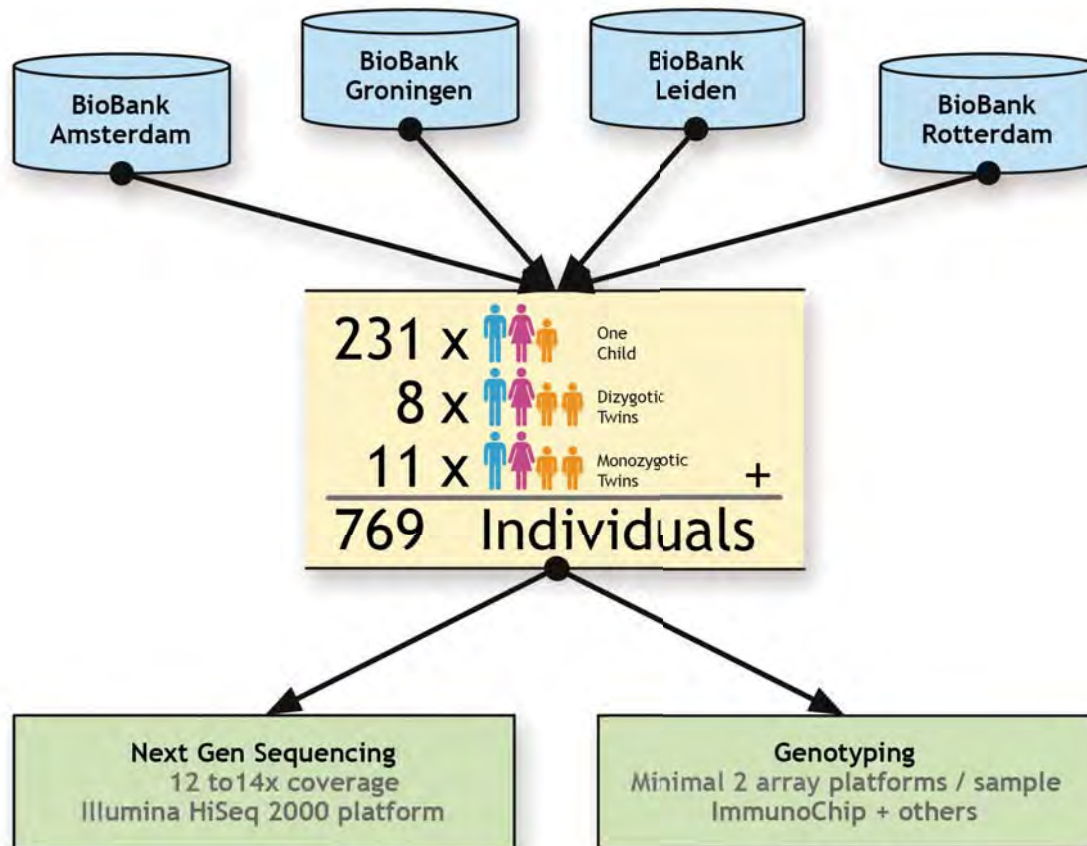


- Expected human Ti/Tv ratio
  - Whole-genome: 2.1
  - Whole-exome: 3.0
- FP SNPs should have a Ti/Tv of 0.5

# Ti/Tv-based filtering threshold



# Genome of the Netherlands (GoNL)

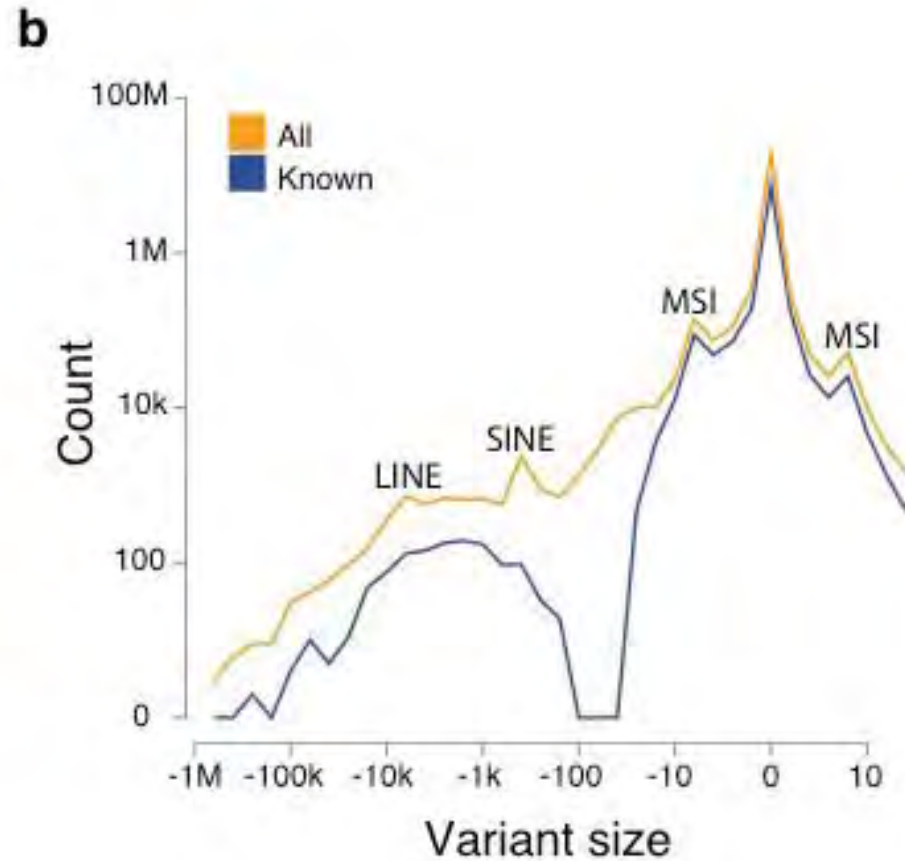
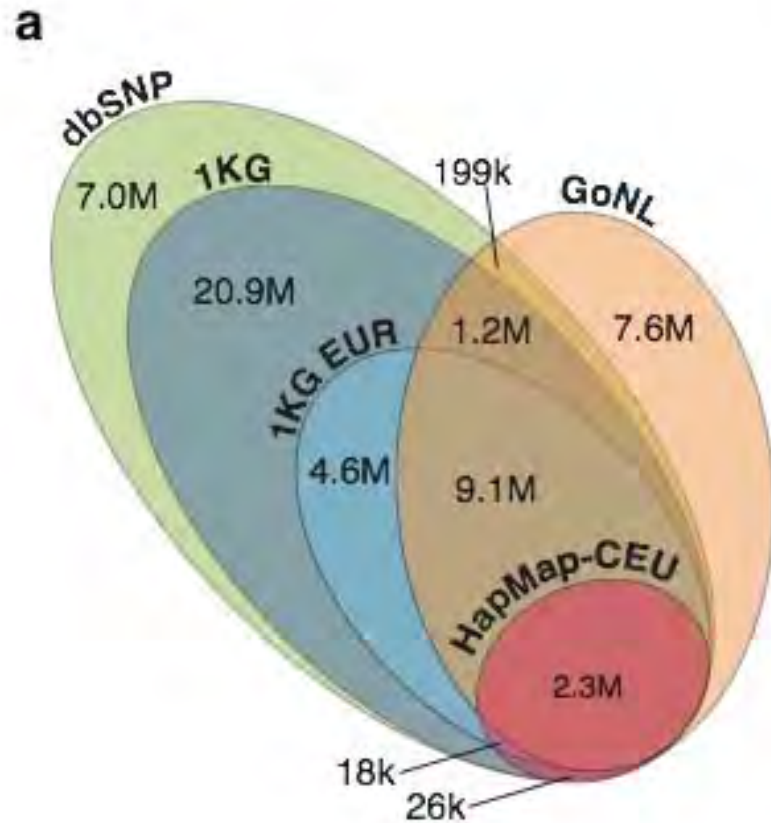


# GoNL SNPs

250 Dutch trios sequenced at 12x coverage on Illumina HiSeq

Step	SNPs (mln)	% in dbSNP 137	Ti/Tv Known	Ti/Tv Novel	#novel SNPs per sample
<b>Initial SNP Calling</b> <i>Unified Genotyper</i>	22.10	44.71	2.18	1.85	154k
<b>Variant Filtering</b> <i>Variant Quality Score Recal.</i>	<b>19.76</b>	<b>60.01</b>	<b>2.25</b>	<b>2.14</b>	<b>26k</b>

# Variation captured by GoNL





# Resources

- GATK homepage
  - <http://www.broadinstitute.org/gatk/index.php>
- GATK video workshop
  - <http://www.broadinstitute.org/gatk/guide/events?id=2038#materials>
- GoNL homepage
  - <http://www.nlgenome.nl>