

# CS681: Advanced Topics in Computational Biology

Week 5 Lecture 1

Can Alkan

EA224

[calkan@cs.bilkent.edu.tr](mailto:calkan@cs.bilkent.edu.tr)

<http://www.cs.bilkent.edu.tr/~calkan/teaching/cs681/>

---

# SNP discovery with NGS data

- ❑ SNP: single nucleotide polymorphism
    - ❑ Change of one nucleotide to another with respect to the reference genome
    - ❑ 3-4.5 million SNPs per person
    - ❑ Database: dbSNP <http://www.ncbi.nlm.nih.gov/projects/SNP/>
  - ❑ Input: sequence data and reference genome
  - ❑ Output: set of SNPs and their genotypes (homozygous/heterozygous)
  - ❑ Often there are errors, filtering required
  - ❑ SNP discovery algorithms are based on statistical analysis
  - ❑ Non-unique mappings are often discarded since they have low MAPQ values
-

# Resequencing-based SNP discovery

genome reference sequence



Read mapping



Read alignment

```
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CGATGTCTTTTCAGGG*TCTCC*ATAAAGAT
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CGATGTCTTTTCAGGG*TCTCC*ATAAAGAT
tt*act*gt*aatggaatactcatgaagtgttaagggctcaaaagaagcctccggcctt
gTT*ACT*GtcGTTGT*AA*TACTCC*AA*cgatgtCTTTTCAGGG*tctcc*ataAAGat
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CGATGTCTTTTCAGGG*TCTCC*ATAAAGAT
tgt*act*ga*gttgC*aa*tactCc*aa*cgATGtctttcaGGG*TCTcc*aTAAAGAT
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CGATGTCTTTTCAGGG*TCTCC*ATAAAGAT
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CAATGTCTTTTCAGGG*TCTCC*ATAAAGAT
gtt*act*gtcgttgt*aa*stactcc*aa*cgatgtCTTtcaggg*Tctcc*ataAagat
GTT*aa*t*kgTCGTTGT*AA*TACTCC*AA*CAATGTCTTTTCAGGG*TCTCC*ATAAAGAT
gtt*act*gtcgttgt*aa*stactcc*aa*cgatgtCTTtcaggg*Tctcc*ataaagat
gtt*act*gtcgttgt*aa*stactcc*aa*cgatgtCTTtcaggg*Tctcc*ataaagat
GTT*Act*gtcgtTGT*aa*tacTcc*aa*cgatgtCTTtcAGgg*tctCC*ATAAAGAT
gtt*act*gtcgttgt*AA*TAcTcc*AA*CGATgtctttcaggg*TCTcc*aTAAAGAT
gtt*acC*stogttGt*aa*ta*ctcc*aa*cgatGtct*gtcaggG*Tctcc*ataaagat
Gt*act*cg*scgttgt*aa*tacTcc*aa*cgatgtct*atCAGGG*TCTCC*ATAAAGAT
gtt*at*gtcgttgt*aa*stactcc*aa*cgatgtCTTtcaggg*stctcc*ataaagat
Gtt*act*gtcgtTGT*aa*ta*ctcc*aa*cgatGTCTTTTCAGGG*TCTCC*ATAAAGAT
GTT*ACT*GTCGTTGT*AA*TACTCC*AA*CGATGTCTTTTCAGGG*TCTCC*ATAAAGAT
gtt*acT*kgTCgttgt*AA*TACTcC*aa*CGATgtctttcaggg*tctcC*ATAaaGAT
```

Paralog identification

```
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
acaatcaggagccggaagcataAAGtgntaaaGctggGGTgcctaaTGAGTGagctaactc
tttGtgagtagaca*gattacaattc atttttaa* t taaag* tt t aaat aatac
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
TTTGTGAGTAGACA*GattACAattCTAttTTAAATATAaag*ttTataaaaTaaATAc
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
TTTGTGAGTAGACA*GATTACAATTCTATTTTAAATATAAAG*TTTATAAAATAAATAC
```

SNP detection + inspection

# Goal

- Given aligned short reads to a reference genome, is a read position a SNP, PSV or error?

TCTCCTCTTCCAGTGGCGAC**G**GAAC      SNP?  
CTCCTCTTCCAGTGGCGAC**A**GAACG  
CTCTTCCAGTGGCGAC**G**GAACGACC      Sequence  
CTTCCAGTGGCGAC**G**GAACGACCC      error?  
CCAGTGGCGAC**T**GAACGACCCTGGA  
CAGTGGCGAC**A**GAACGACCCTGGAG

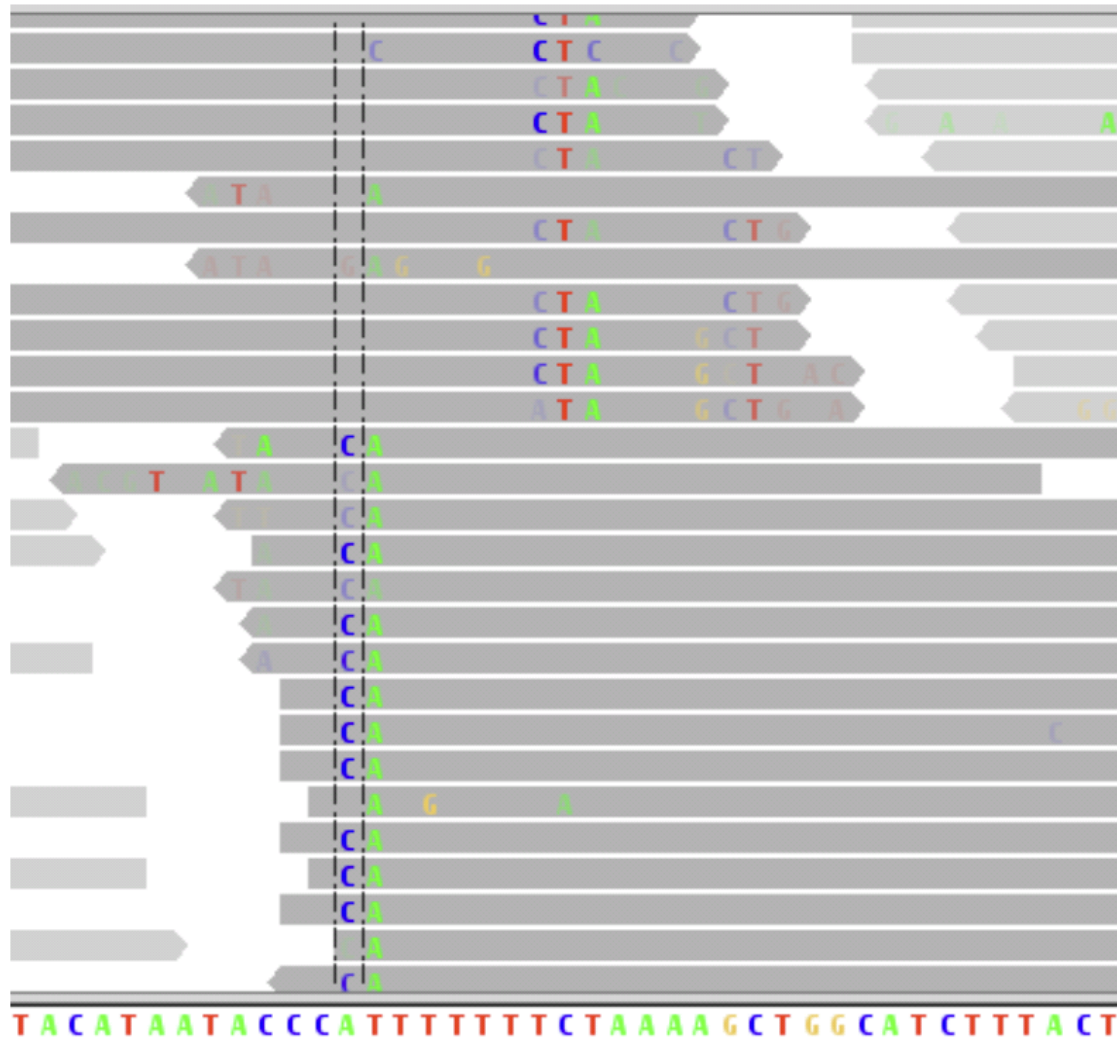
Reference    TCTCCTCTTCCAGTGGCGAC**G**GAACGACCCTGGAGCCAAGT

---

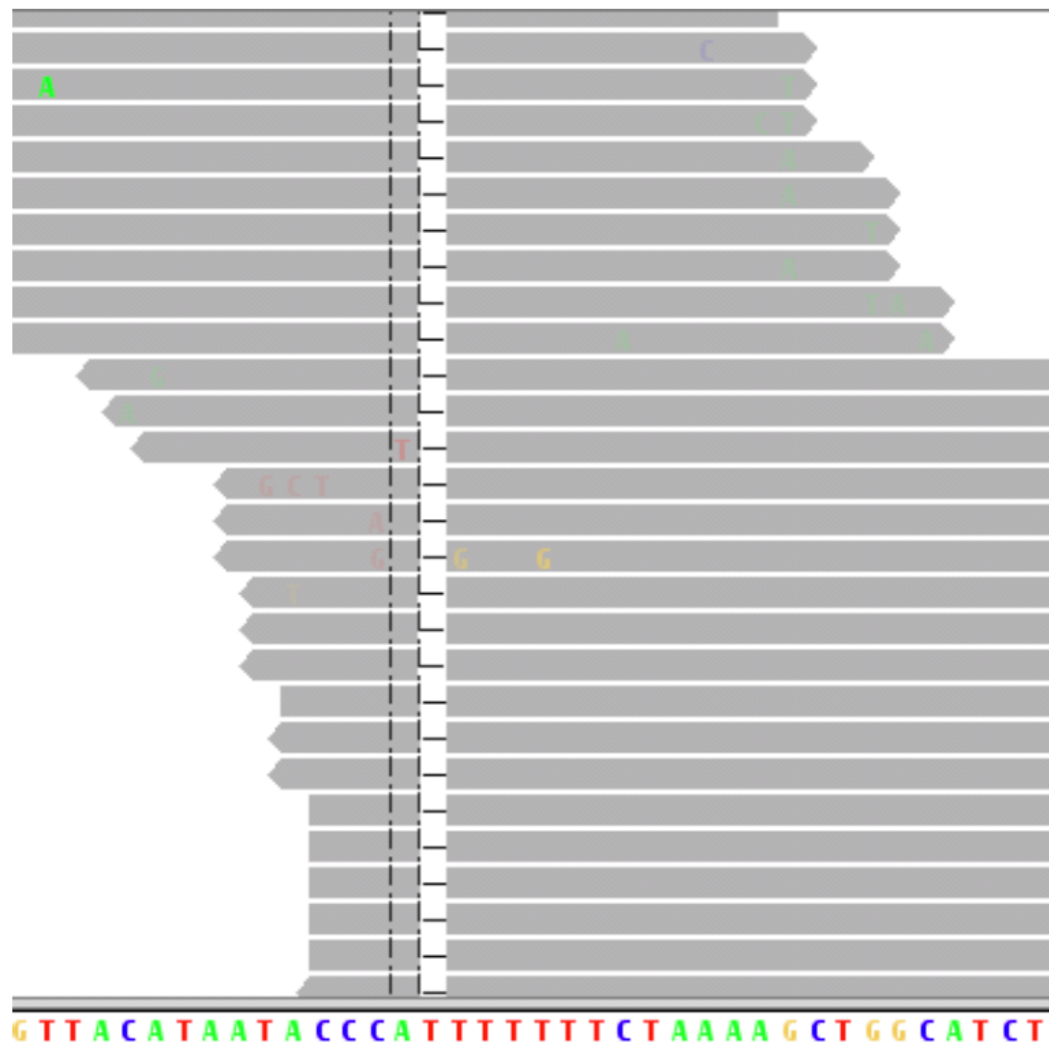
# Challenges

- Sequencing errors
  - Paralogous sequence variants (PSVs) due to repeats and duplications
  - Misalignments
    - Indels vs SNPs, there might be more than one optimal trace path in the DP table
    - Short tandem repeats
    - Need to generate multiple sequence alignments (MSA) to correct
-

# Need to realign



# After MSA



# Indel scatter

Even when read mapper detects indels in individual reads successfully, they can be scattered around (due to additional mismatches in the read)

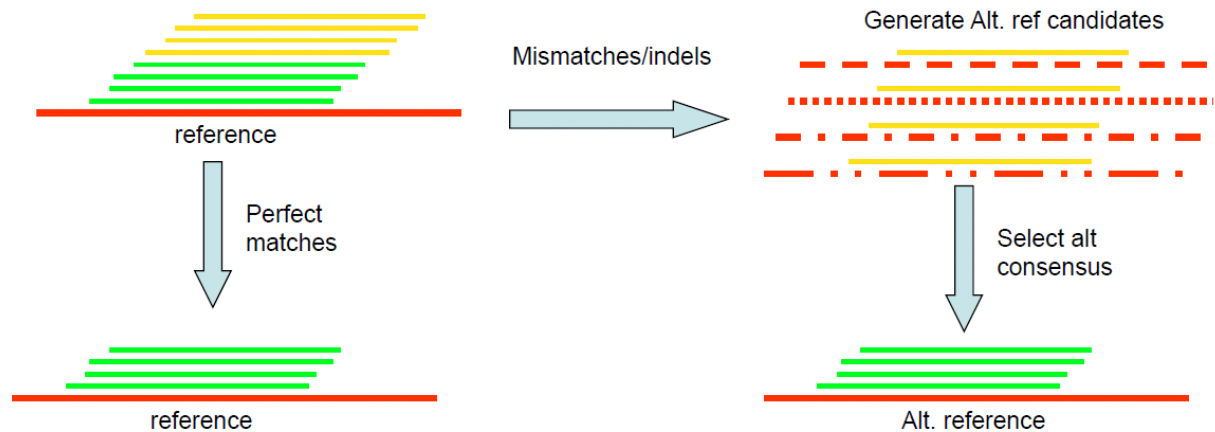
```
TAAATAATGGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT++++AGGGT++++GCACTCTCTGCTTCATAAATGGGTCTCTTGCCGCAAAAAAATCTGTTTGCTCCTCCAGATTCATCAA
<- TGGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGG
<- TGGAAATTTATTTCTCAAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGG
GGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGG
-> GGAAATTTATTTCAAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGG
-> CAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGCGCACTCTCTGCTTCATAAATGGGTCTCTTG
-> ATTTCTCAGAGTACTGGAAGCTGGGACTCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGT*****AGGGTGC
<- GTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCT
<- AATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTC
-> ATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTCTTGCCGCA
<- GTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTCTTGCCGCAAAAAAATCTGTTTGCTCCTCCAG
```

```
TAAATAATGGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT++++AGGGTGCACCTCTCTGCTTCATAAATGGGTCTCTTGCCGCAAAAAAATCTGTTTGCTCCTCCAGATTCATCAA
<- TGGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGG
<- TGGAAATTTATTTCTCAAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGG
<- GGAAATTTATTTCTCAGAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGG
-> GGAAATTTATTTCAAGTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGG
-> CAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGCGCACTCTCTGCTTCATAAATGGGTCTCTTG
-> ATTTCTCAGAGTACTGGAAGCTGGGACTCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGT*****AGGGTGC
<- GTACTGGAAGCTGGGAATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCT
<- AATCCAAGATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTC
-> ATCAAAATGCCAGCAGATTCTAAGTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTCTTGCCGCA
<- GTCTGGT*****AGGGTAGGGT*****GCACCTCTCTGCTTCATAAATGGGTCTCTTGCCGCAAAAAAATCTGTTTGCTCCTCCAG
```



# MSA for resequencing

- We have the reference and (approximate) placement
- Departures from the reference are small
- Generate alt reference as suggested by *each* non-matching read (Smith-Waterman)
- Test each non-matching read against each alt reference candidate
- Select alt reference consensus: best “home” for all non-matching reads
- Why is it MSA: look for improvement in *overall* placement score (sum across reads)
- Optimizations and constrains:
  - Expect two alleles
  - Expect a single indel
  - Downsample in regions of very deep coverage
  - Alignment has an indel: use that indel as an alt. ref candidate



---

# SNP callers

- Genome Analysis Tool Kit (GATK; Broad Inst.)
  - Samtools (Sanger Centre)
  - PolyBayes (Boston College)
  - SOAPsnp (BGI)
  - VARiD (U. Toronto)
  - ....
-

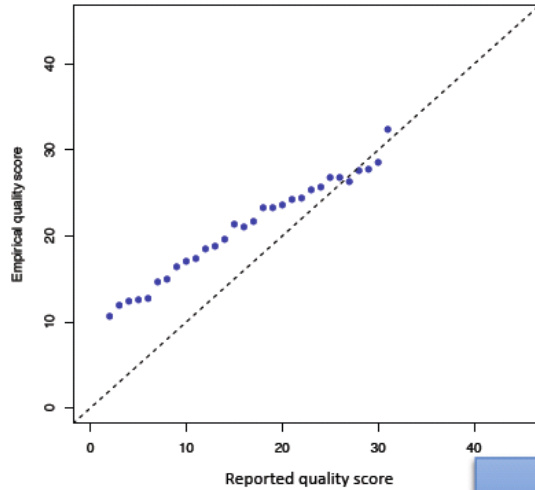
---

# Base quality recalibration

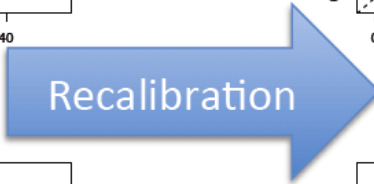
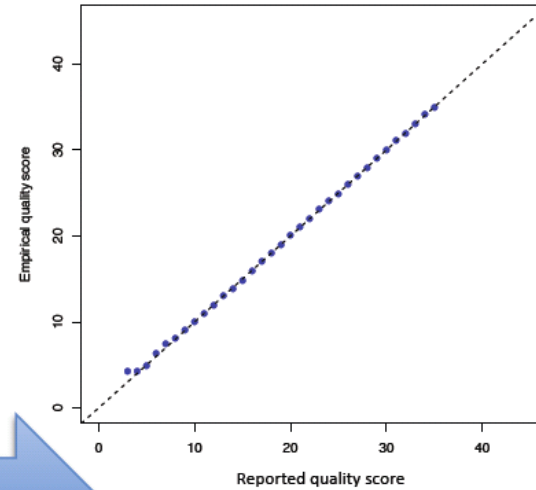
- The quality values determined by sequencers are not optimal
  - There might be sequencing errors with high quality score; or correct basecalls with low quality score
  - Base quality recalibration: after mapping correct for base qualities using:
    - Known systematic errors
    - Reference alleles
    - Real variants (dbSNP, microarray results, etc.)
  - Most sequencing platforms come with recalibration tools
  - In addition, GATK & Picard have recalibration built in
-

# Base quality recalibration

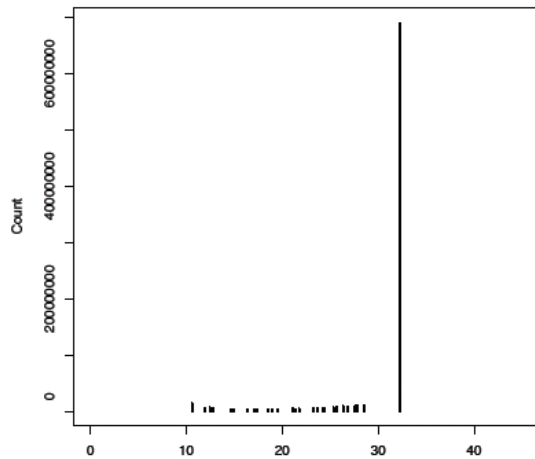
Reported vs. empirical quality scores



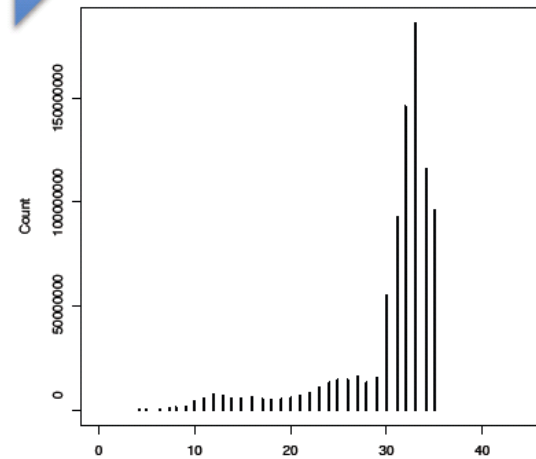
Reported vs. empirical quality scores



Reported quality score histogram

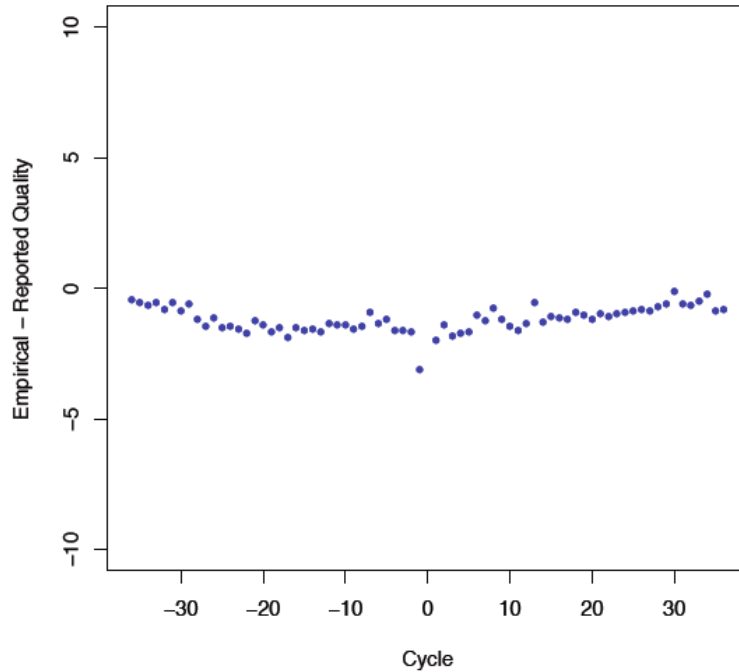


Reported quality score histogram



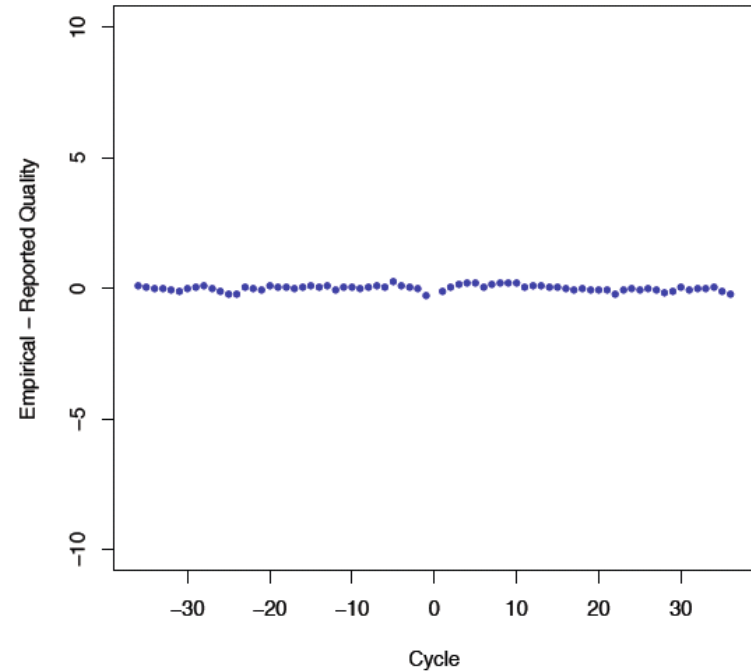
# Recalibration by machine cycle

RMSE = 1.275



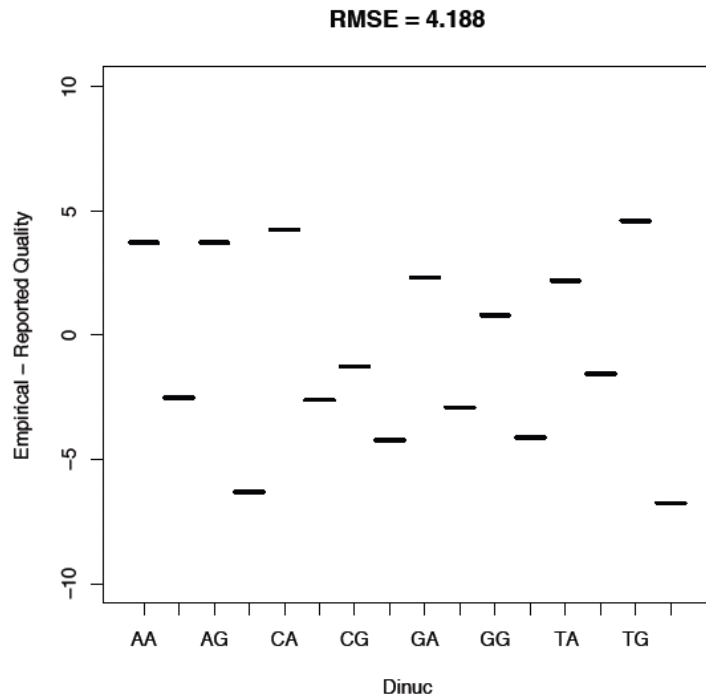
Before Recalibration

RMSE = 0.105

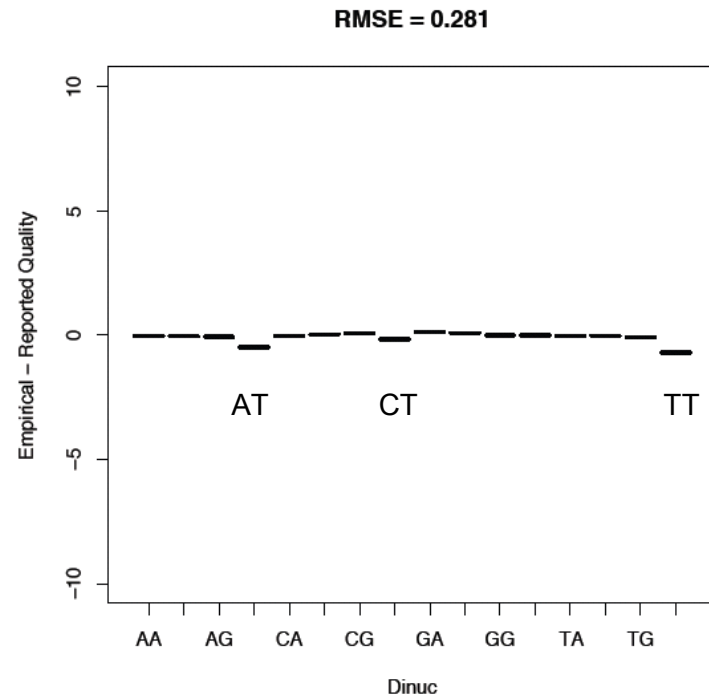


After Recalibration

# Recalibration by dinucleotide



Before Recalibration



After Recalibration

# PolyBayes

```

TCTGACCAATCTAAAATACCTGTGATTAA
TCTGACCAATCTAACAATACCTGTGATTAA
TCTGACCAATCTAACAATACCTGTGATTAA
TCTGACCAATCTAAAATACCTGTGATTAA
tctgaccaatctaacaataacctgtgattaa
    
```

**Bayesian posterior probability i.e. the SNP score**

**Base call + Base quality**

**Polymorphism rate (prior)**

$$P(SNP) = \sum_{\text{all variable } S} \frac{P(S_1 | R_1) \dots P(S_N | R_N) P_{Prior}(S_1, \dots, S_N)}{P_{Prior}(S_1) \dots P_{Prior}(S_N)}$$

$$\sum_{S_{i_1} \in [A,C,G,T]} \dots \sum_{S_{i_N} \in [A,C,G,T]} \frac{P(S_{i_1} | R_{i_1}) \dots P(S_{i_N} | R_{i_N}) P_{Prior}(S_{i_1}, \dots, S_{i_N})}{P_{Prior}(S_{i_1}) \dots P_{Prior}(S_{i_N})}$$

# Base quality values for SNP calling

$$P(SNP) = \sum_{\text{all variable } S} \frac{P(S_1 | R_1) \cdots P(S_N | R_N) \cdot P_{\text{Prior}}(S_1, \dots, S_N)}{\sum_{S_1 \in [A,C,G,T]} \cdots \sum_{S_N \in [A,C,G,T]} \frac{P(S_{i_1} | R_{i_1}) \cdots P(S_{i_N} | R_{i_1})}{P_{\text{Prior}}(S_{i_1}) \cdots P_{\text{Prior}}(S_{i_N})} \cdot P_{\text{Prior}}(S_{i_1}, \dots, S_{i_N})}$$

TTGAT **C** CCTGT  
 TTGAT **T** CCTGT

TGAAA **g** gGAATT  
 TGAAA **t** GAATT

- base quality values help us decide if mismatches are true polymorphisms or sequencing errors
- **accurate base qualities are crucial**, especially in lower coverage



# Priors for specific scenarios

$$P(SNP) = \sum_{\text{all variable } S} \frac{P(S_1 | R_1) \cdots P(S_N | R_N) \cdot P_{\text{Prior}}(S_1, \dots, S_N)}{\sum_{S_{i_1} \in [A,C,G,T]} \cdots \sum_{S_{i_N} \in [A,C,G,T]} \frac{P(S_{i_1} | R_1)}{P_{\text{Prior}}(S_{i_1})} \cdots \frac{P(S_{i_N} | R_1)}{P_{\text{Prior}}(S_{i_N})} \cdot P_{\text{Prior}}(S_{i_1}, \dots, S_{i_N})}$$

strain 1

AACGTTAGCAT  
AACGTTAGCAT  
AACGTTAGCAT

strain 2

AACGTT**C**GCAT  
AACGTT**C**GCAT

strain 3

AACGTTAGCAT  
AACGTTAGCAT  
AACGTTAGCAT

individual 1

AACGTTAGCAT  
AACGTTAGCAT  
AACGTT**C**GCAT  
AACGTT**C**GCAT

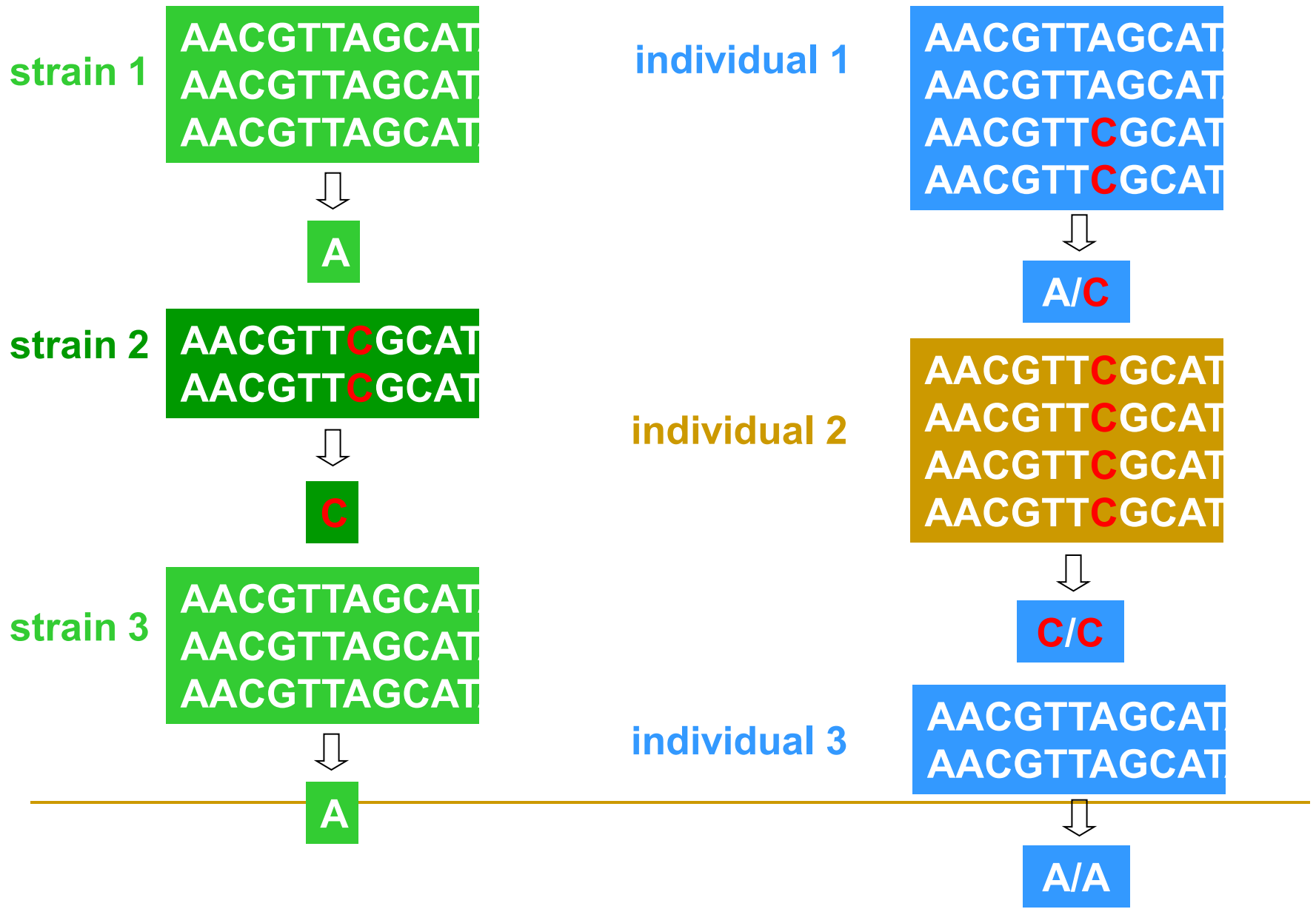
individual 2

AACGTT**C**GCAT  
AACGTT**C**GCAT  
AACGTT**C**GCAT  
AACGTT**C**GCAT

individual 3

AACGTTAGCAT  
AACGTTAGCAT

# Consensus sequence generation (genotyping)



# SOAPsnp

## ■ Bayesian model

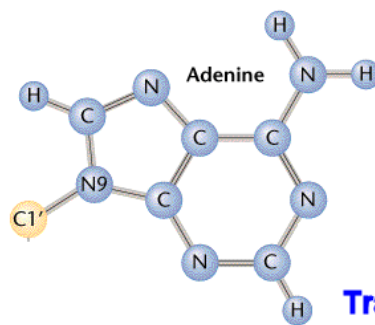
- $T_i$ : genotype
- $D$ : data at a locus
- $S$ : total number of genotypes

$$P(T_i | D) = \frac{P(D | T_i)P(T_i)}{\sum_{x=1}^S P(D | T_x)P(T_x)}$$

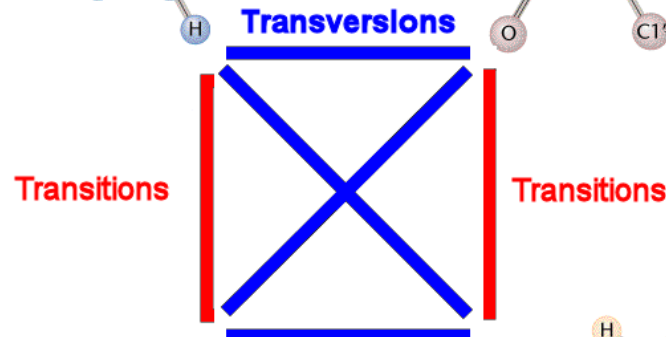
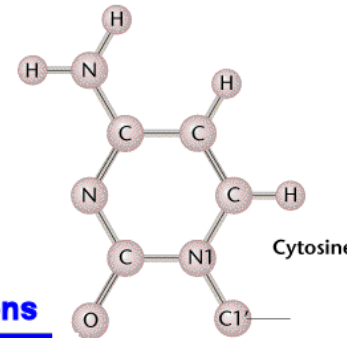
# SOAPsnp priors: Haploid

- SNP rate = 0.001. Assuming ref allele is **G**

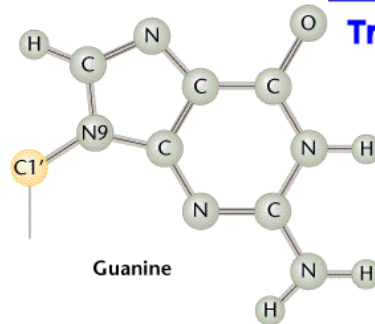
**A**  
 $4/6 \times 0.001$



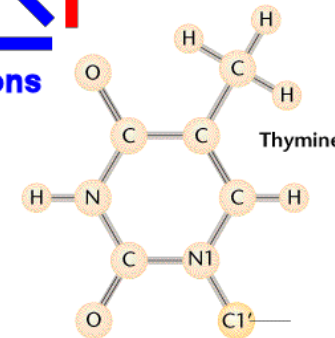
**C**  
 $1/6 \times 0.001$



**G**  
0.999



**T**  
 $1/6 \times 0.001$



Ideally;  $T_i / T_v = 2.1$

Li et al, Genome Research, 2009

# SOAPsnp priors: Diploid

- Heterozygous SNP rate = 0.001
- Homozygous SNP rate = 0.0005
- Assuming ref allele is **G**

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
A	$3.33 \times 10^{-4}$	$1.11 \times 10^{-7}$	$6.67 \times 10^{-4}$	$1.11 \times 10^{-7}$
C		$8.33 \times 10^{-5}$	$1.67 \times 10^{-4}$	$2.78 \times 10^{-8}$
G			0.9985	$1.67 \times 10^{-4}$
T				$8.33 \times 10^{-5}$

# SOAPsnp: Genotype Likelihood

$$P(D | T) = \prod_{k=1}^n P(d_k | T)$$

$$P(d_k | T)$$

$$= P(o_k, q_k, c_k | T)$$

$$= P(o_k, c_k | q_k, T) P(q_k | T)$$

**T:** genotype (GG/GA/AA)

**o:** observed allele type

**q:** quality score

**c:** cycle

TCTCCTCTTCCAGTGGCGAC**G**GAAC

CTCCTCTTCCAGTGGCGAC**A**GAACG

CTCTTCCAGTGGCGAC**G**GAACGACC

CTTCCAGTGGCGAC**G**GAACGACCC

CCAGTGGCGAC**T**GAACGACCCTGGA

CAGTGGCGAC**A**GAACGACCCTGGAG

  $d_k$ : observed allele

# GATK SNP calling

$$P(G | D) = \frac{P(G)P(D | G)}{\sum (G_i)P(D | G_i)}$$

$$P(D | G) = \prod_j \left( \frac{P(D_j | H_1)}{2} + \frac{P(D_j | H_2)}{2} \right) \text{ where } G = \{H_1, H_2\}$$

$$P(D_j | H) = P(D_j | b)$$

$$P(D_j | b) = \begin{cases} 1 - \varepsilon_j & D_j = b \\ \varepsilon_j & \text{otherwise} \end{cases}$$

# GATK genotype likelihoods

$$\begin{array}{c} \text{Likelihood for} \\ \text{the genotype} \end{array} \quad \begin{array}{c} \text{Prior for} \\ \text{the genotype} \end{array} \quad \begin{array}{c} \text{Likelihood for} \\ \text{the data} \\ \text{given genotype} \end{array} \quad \begin{array}{c} \text{Independent base model} \end{array}$$
$$L(G | D) = P(G)P(D | G) = \prod_{b \in \{good\_bases\}} P(b | G)$$

- Likelihood of data computed using pileup of bases and associated quality scores at given locus
- Only “good bases” are included: those satisfying minimum base quality, mapping read quality, pair mapping quality
- $P(b | G)$  uses platform-specific confusion matrices
- $L(G|D)$  is computed for all 10 genotypes



---

# SNP calling artifacts

- SNP calls are generally infested with false positives
  - From systematic machine artifacts, mismapped reads, aligned indels/CNV
  - Raw SNP calls might have between 5-20% FPs among novel calls
- Separating true variation from artifacts depends very much on the particulars of one's data and project goals
  - Whole genome deep coverage data, whole genome low-pass, hybrid capture, pooled PCR are have significantly different error models

---

# Filtering

- Hard filters based on
    - Read depth (low and high coverage are suspect)
    - Allele balance
    - Mapping quality
    - Base quality
    - Number of reads with MAPQ=0 overlapping the call
    - Strand bias
    - SNP clusters in short windows
-

---

# Filtering

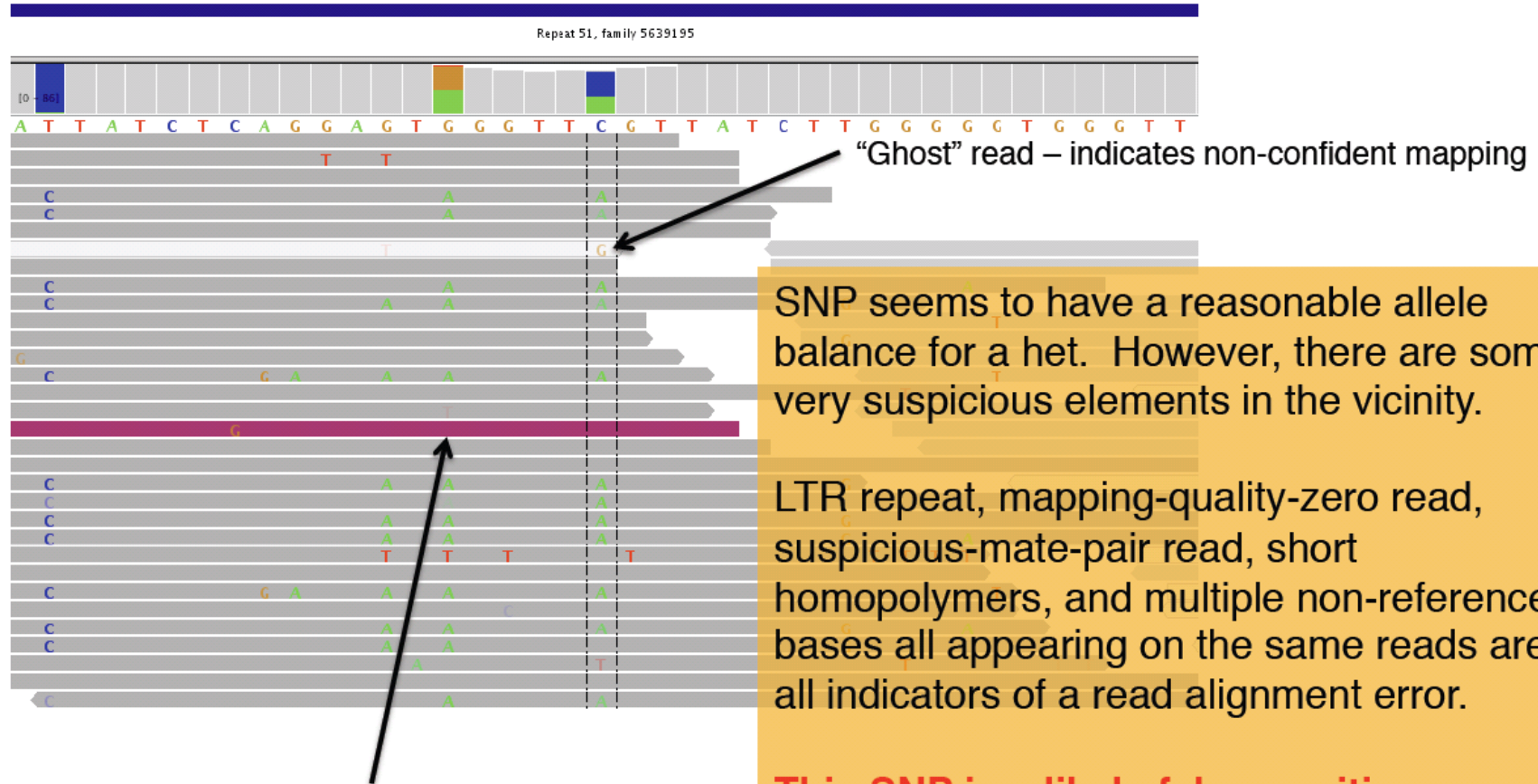
- Statistical determination of filtering parameters:
    - Training data: dbSNP, HapMap, microarray experiments, other published results
    - Based on the distribution of values over the training data adjust cut off parameters depending on the sequence context
      - VQSR: Variant Quality Score Recalibration
-

# Indicators of call set quality

- Number of variants
  - Europeans and Asians: ~3 million; Africans: ~4-4.5 million
- Transition/transversion ratio
  - Ideally  $Ti/Tv = 2.1$
- Hardy Weinberg equilibrium
  - Allele and genotype frequencies in a population remain constant
  - For alleles **A** and **a**;  $\text{freq}(\mathbf{A})=p$  and  $\text{freq}(\mathbf{a})=q$ ;  $p+q=1$
  - If a population is in equilibrium then
    - $\text{freq}(\mathbf{AA}) = p^2$
    - $\text{freq}(\mathbf{aa}) = q^2$
    - $\text{freq}(\mathbf{Aa}) = 2pq$
- Presence in databases: dbSNP, HapMap, array data
- Visualization

# Validation through visualization

NA19240, chr1:5,639,327-5,639,365



Read's mate-pair maps to another chromosome

**This SNP is a likely false-positive.**

---

# Pooled sequencing

- When sequence coverage is low, pool mapping of data from multiple samples (ideally from the same population) into a single file
  - SNP calling is more challenging
    - Allele frequencies close to error rate
    - Track which read comes from which individual
-

---

**NEXT: INDELS**

---