

DTU



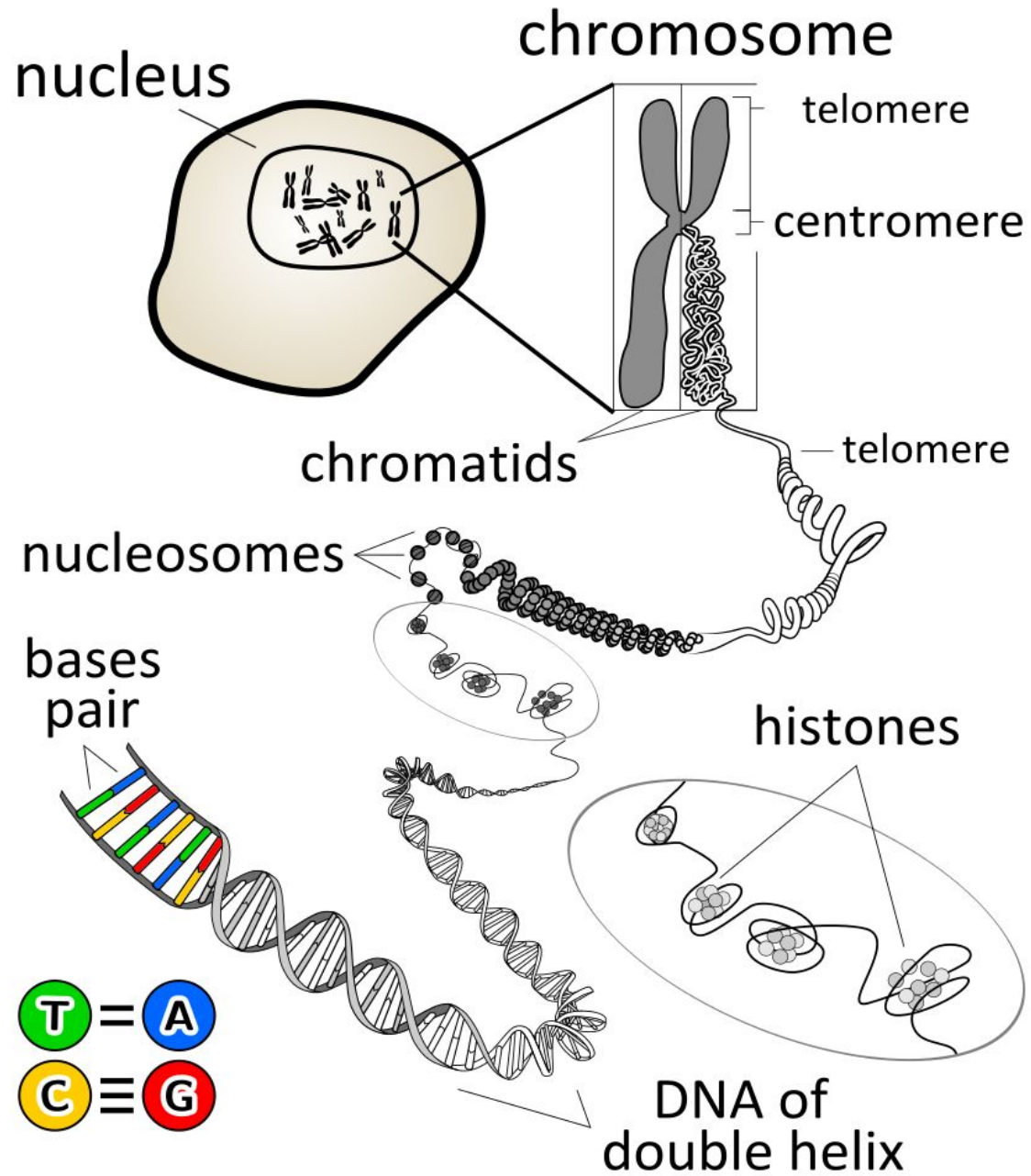


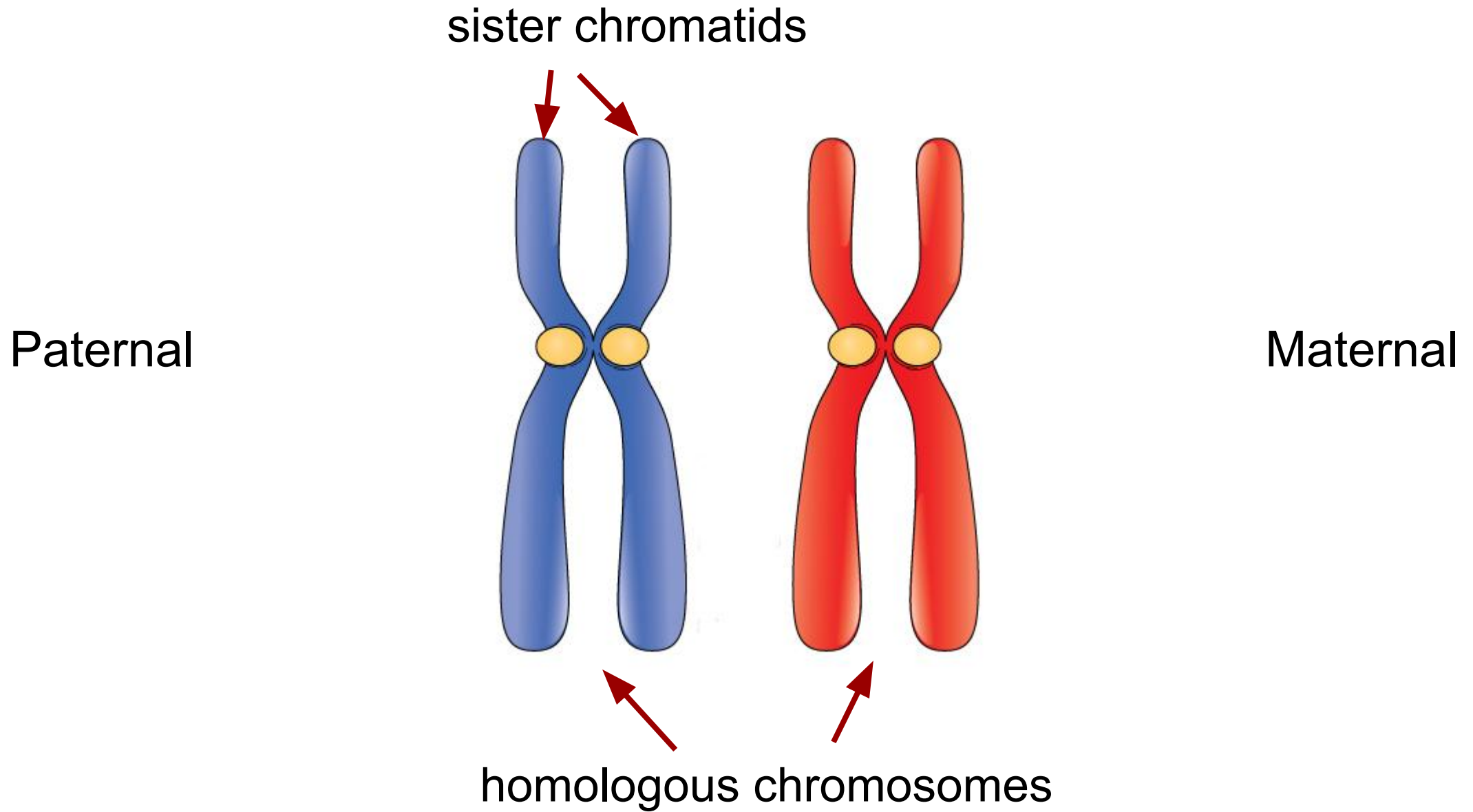
**DTU Health Technology
Bioinformatics**

Alignment post-processing and variant calling

*Gabriel Renaud
Associate Professor
Section of Bioinformatics
Technical University of Denmark
gabriel.reno@gmail.com*

A brief reminder





Heterozygosity

TACAAATAT
TACAGATAT

M:



P:

Heterozygous sites

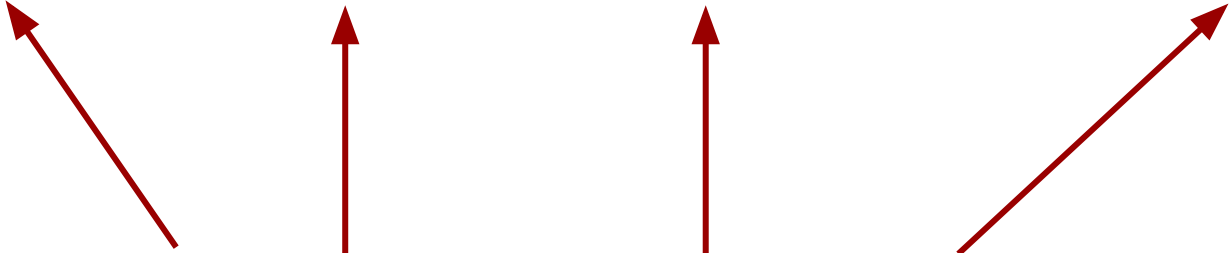
Heterozygosity

M:

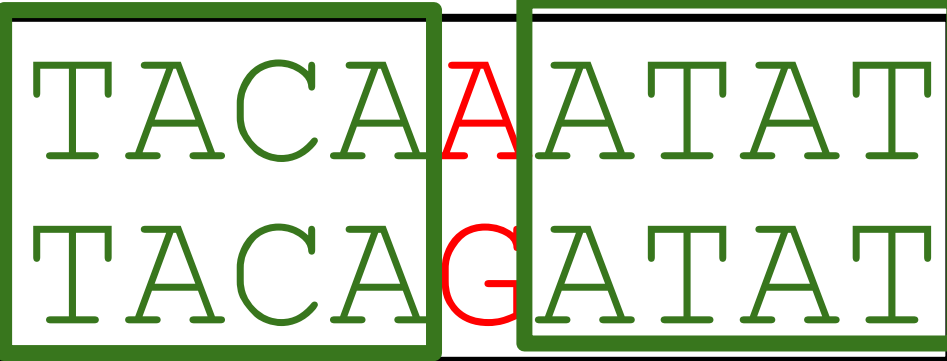


P:

segregating sites



Heterozygosity

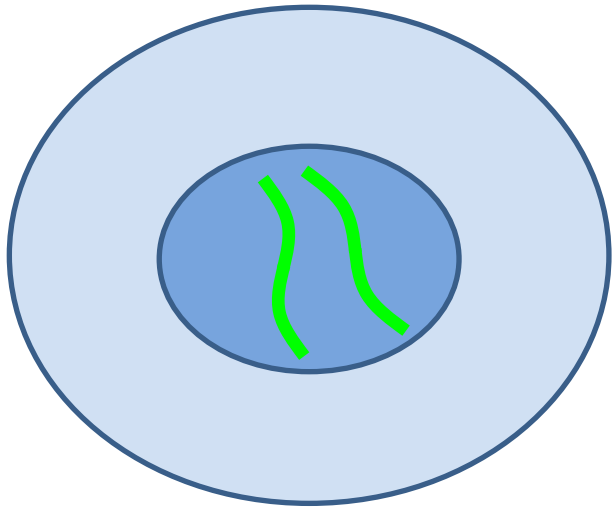


M:



P:

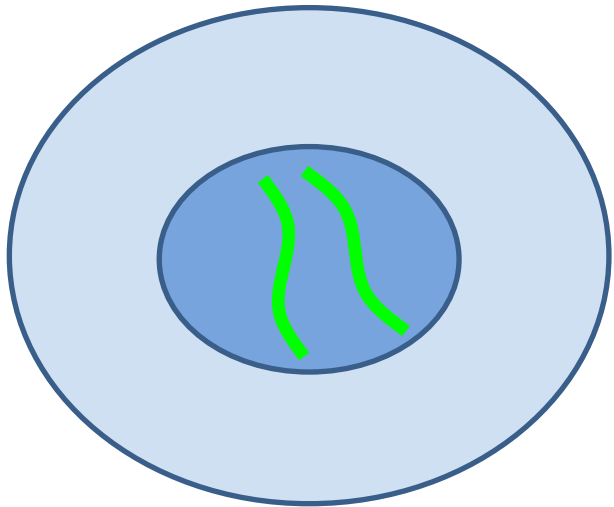
Homozygous sites



ind#A

M: **TACAAAATAT**

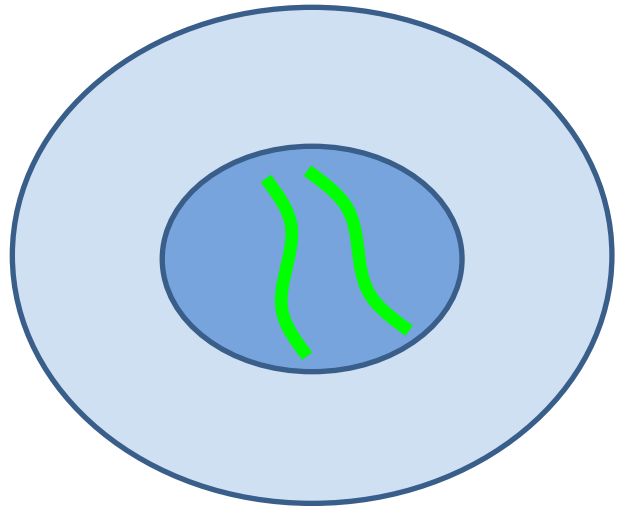
P: **TACAGATAT**



ind#B

M: **TACAGATCT**

P: **TACAGATCT**

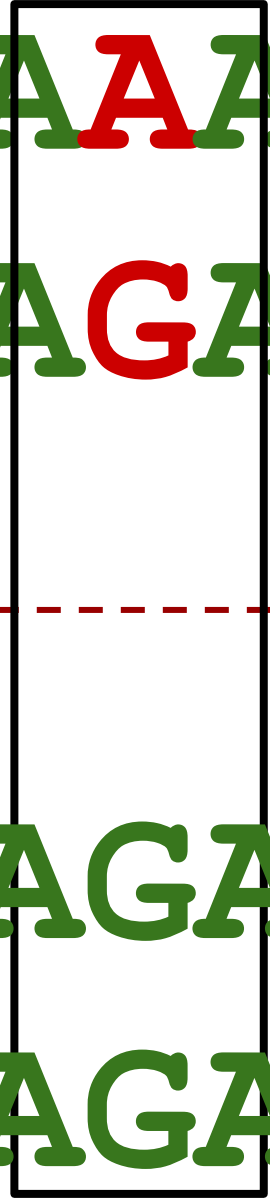


ind#A

M: TACAAATAT

P: TACAGATAT

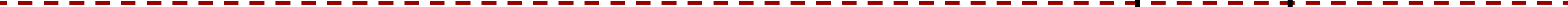
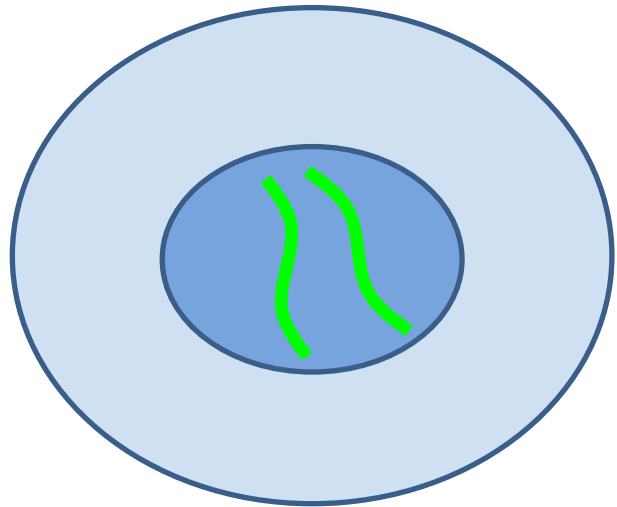
Heterozygosity

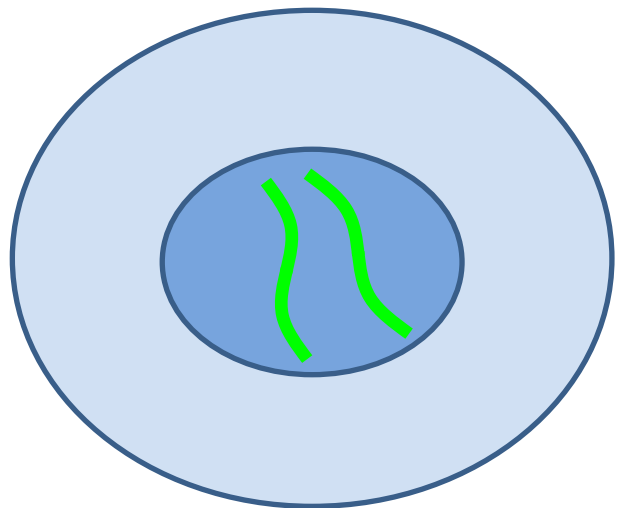


ind#B

M: TACAGATCT

P: TACAGATCT



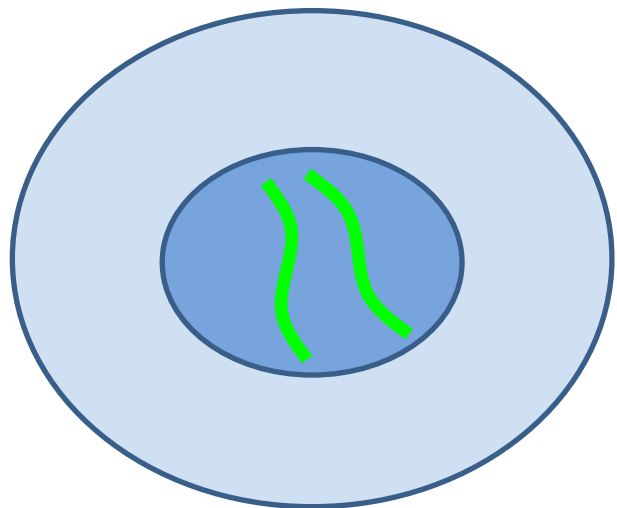


ind#A

Homozygous variant

M: TACAAATAT

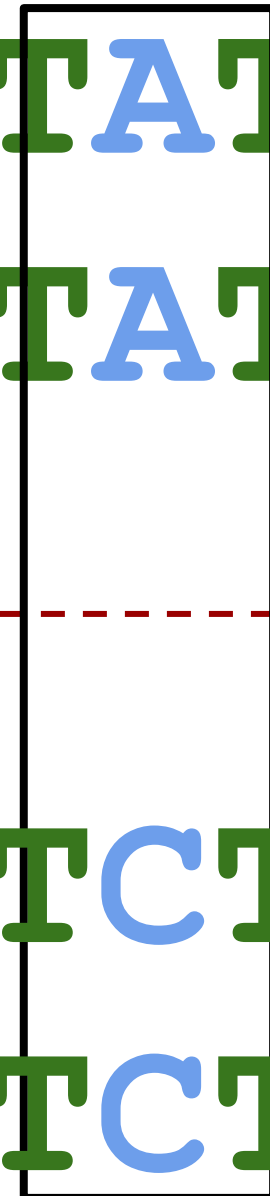
P: TACAGATAT

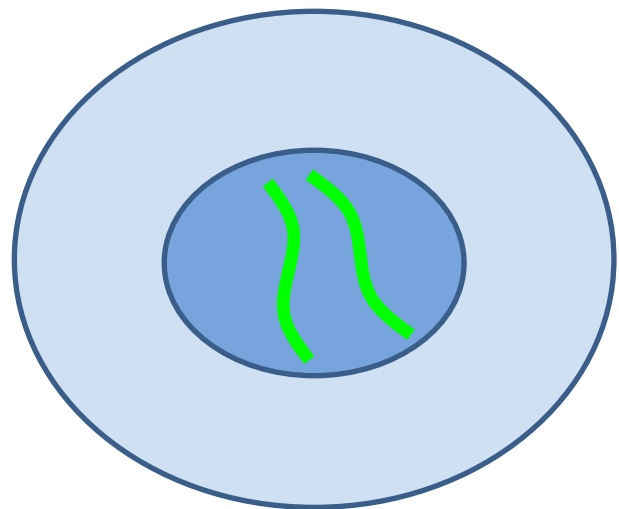


ind#B

M: TACAGATCT

P: TACAGATCT





ind#A

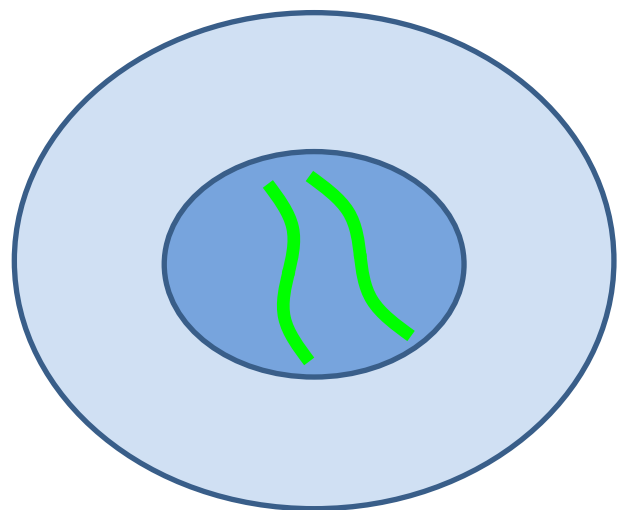
M:

TACAAATAT

P:

TACAGATAT

Homozygous invariant



ind#B

M:

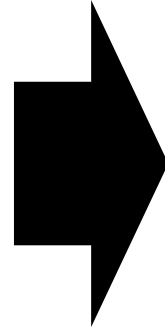
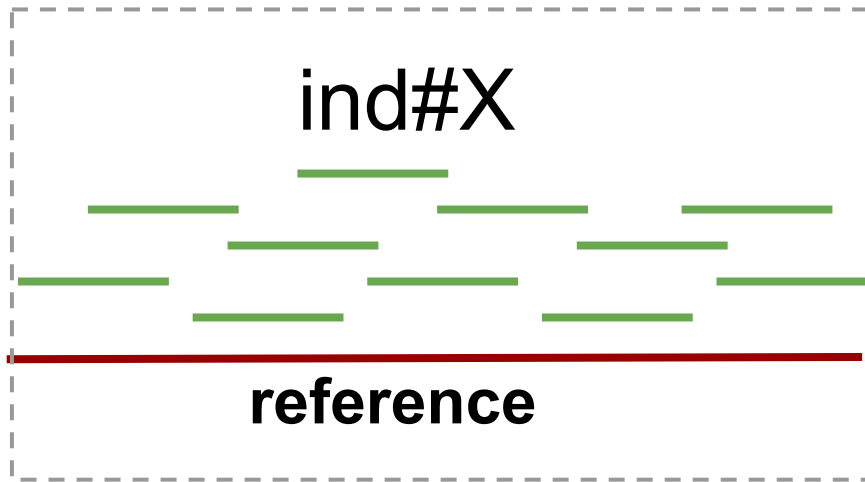
TACAGATCT

P:

TACAGATCT

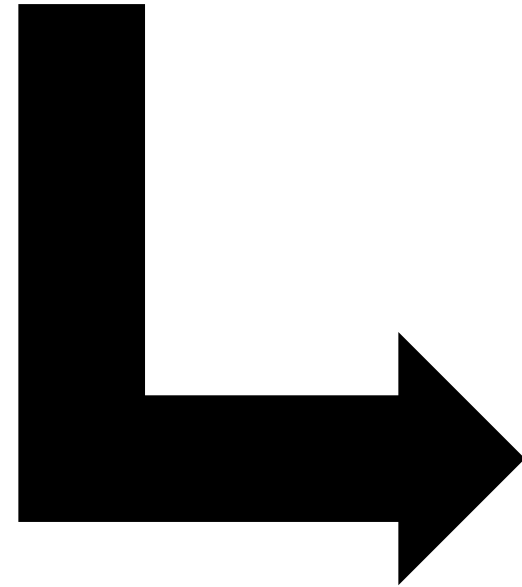


Genotyping



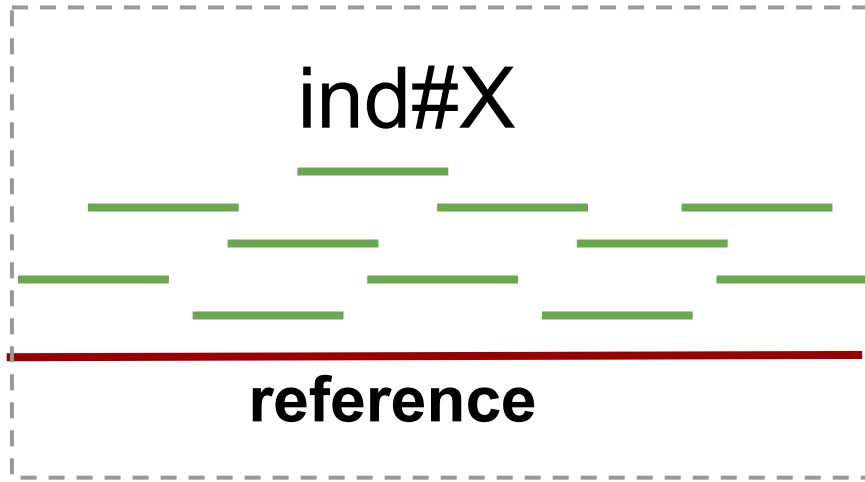
TACAAATAT
TACAGATAT

Which of the 10 possible genotypes is the most likely?

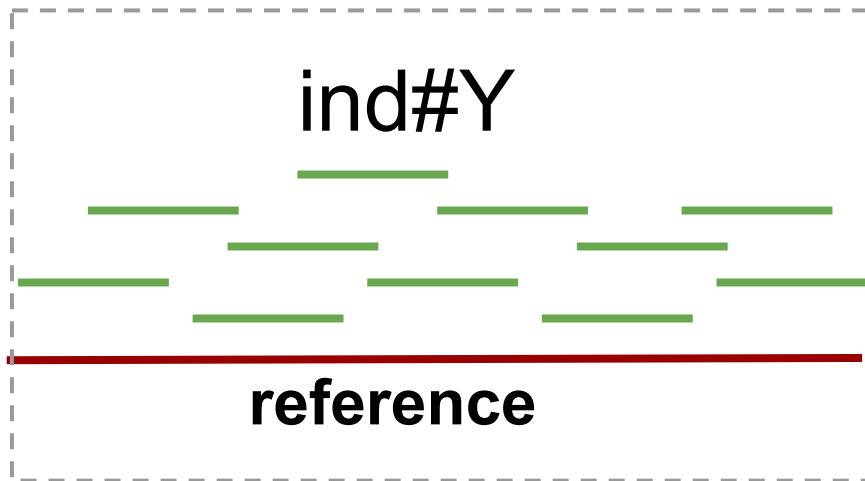


AA
AC
AG
AT
CC
CG
CT
GG
GT
TT

Joint Genotyping



TACAAATAT
TACAGATAT

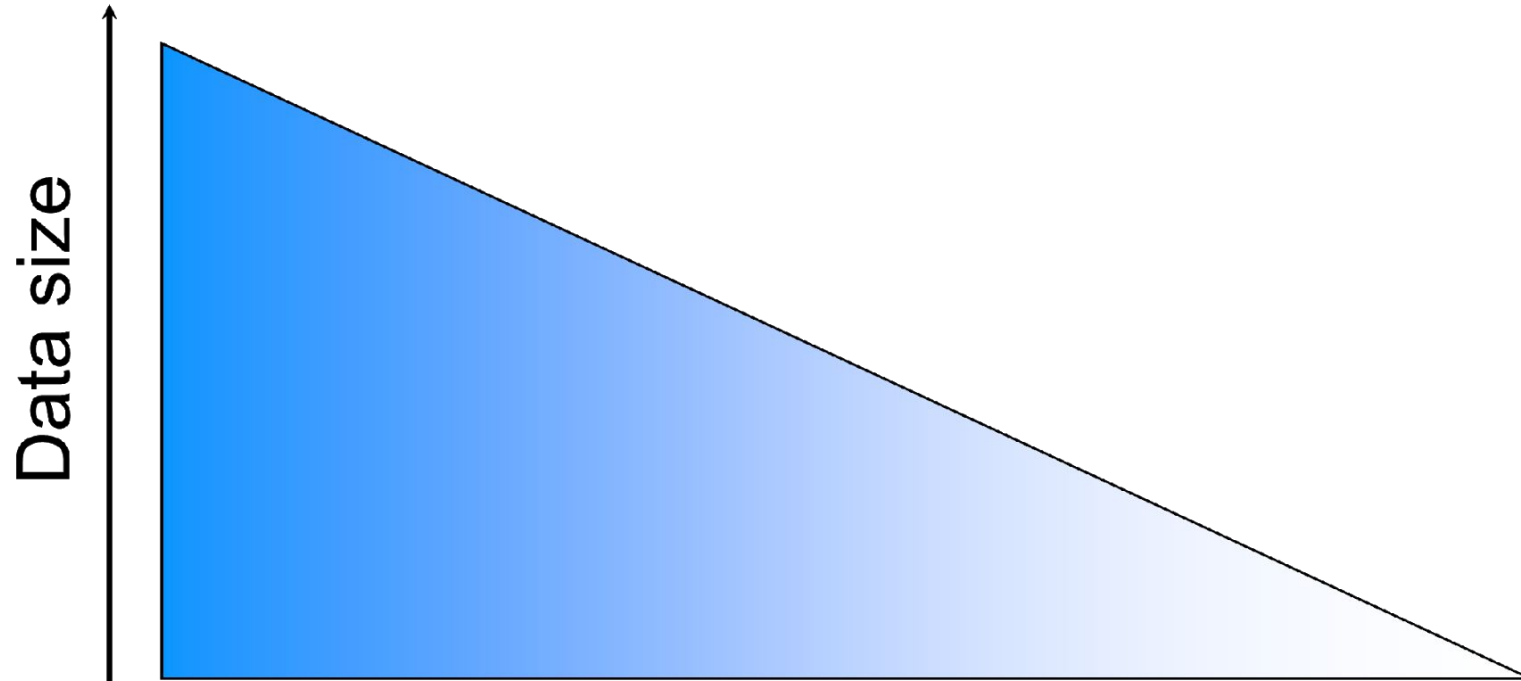


TACAGATCT
TACAGATCT

Menu

- Introduction
- Alignment post-processing
- From aligned reads to genomic variation
- Variant effect

Generalized NGS analysis



Question

Raw reads

Pre-processing

Assembly:
Alignment /
de novo

Application specific:
Variant calling,
count matrix, ...

Compare samples /
methods

Answer?

Brief probability reminder

Events:

E = I pick a random human and that person is Danish

$P(E)$ = pop. of Denmark / pop. world



Brief probability reminder

Events:

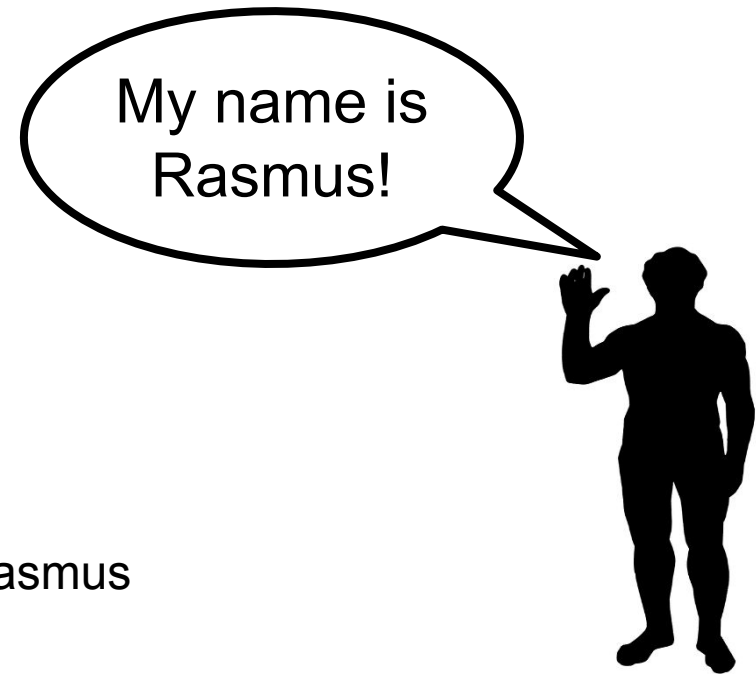
E = I pick a random human and that person is Danish

S = I pick a random human and that person's name is Rasmus

$P(E)$ = pop. of Denmark / pop. world

$P(S)$ = # of Rasmuses / pop. world

$P(E|S)$ = # of Rasmuses in Denmark / # of Rasmuses in the world



$$P(E|S) \gg P(E)$$

What is genotyping?

Genotyping is determining which genotype maximizes:

$$P(G|D)$$

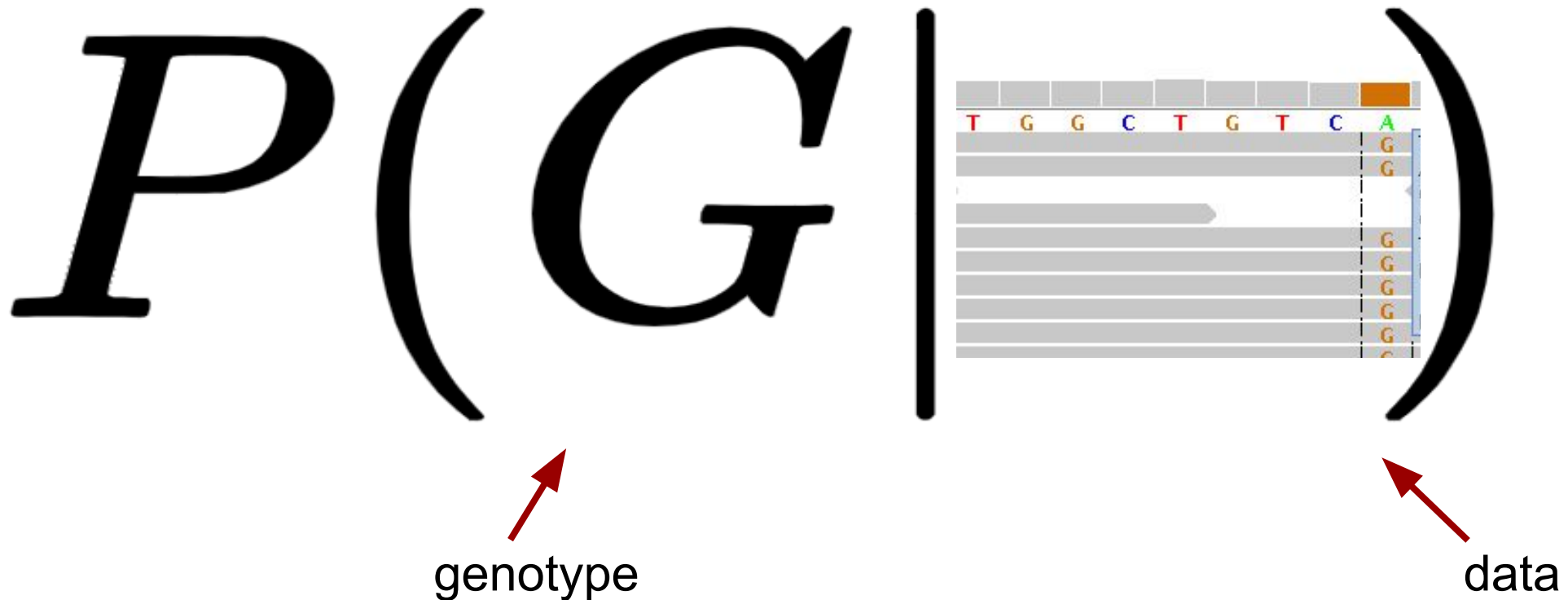
genotype

data

The diagram shows the mathematical expression $P(G|D)$ in a large, bold, black serif font. Below the letter 'G' is a red arrow pointing upwards towards it, with the word 'genotype' written in black text below the arrow. Similarly, below the letter 'D' is a red arrow pointing upwards towards it, with the word 'data' written in black text below the arrow.

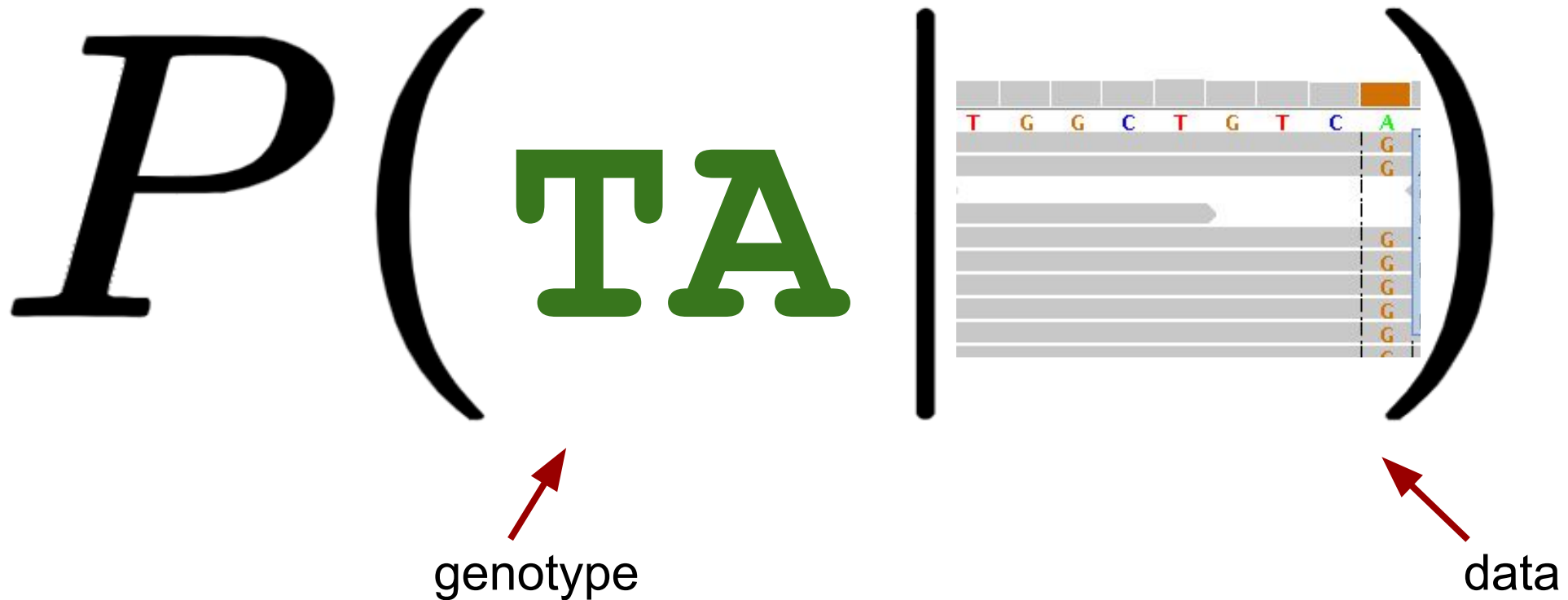
What is genotyping?

Genotyping is determining which genotype maximizes:



What is genotyping?

Genotyping is determining which genotype maximizes:



What is genotyping?

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

What is genotyping?

prior: what is the probability of the genotype to begin with?

likelihood: What is the probability of seeing the data given the genotype?

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

What is genotyping?

prior: what is the probability of the genotype to begin with?

likelihood: What is the probability of seeing the data given the genotype?

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

evidence: What is the probability of generating that data to begin with?

$$P(D) = \sum_{G \in \mathbb{G}} P(G)P(D|G)$$

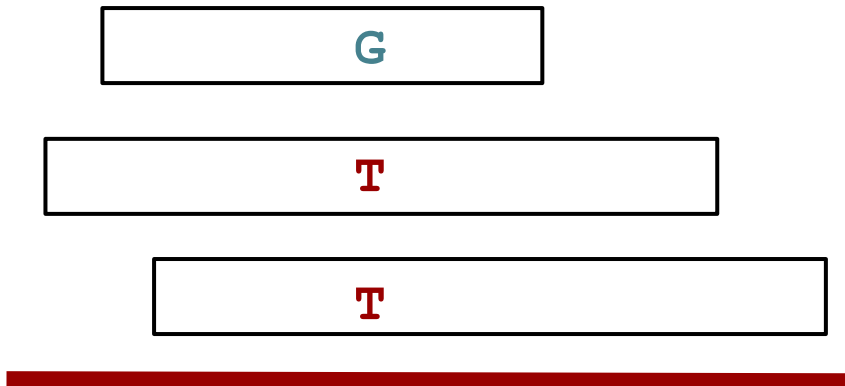
The likelihood

$$P(D|G) = \prod_{b \in \text{READS}} P(b|G)$$

i.e. each reads is an independent observation

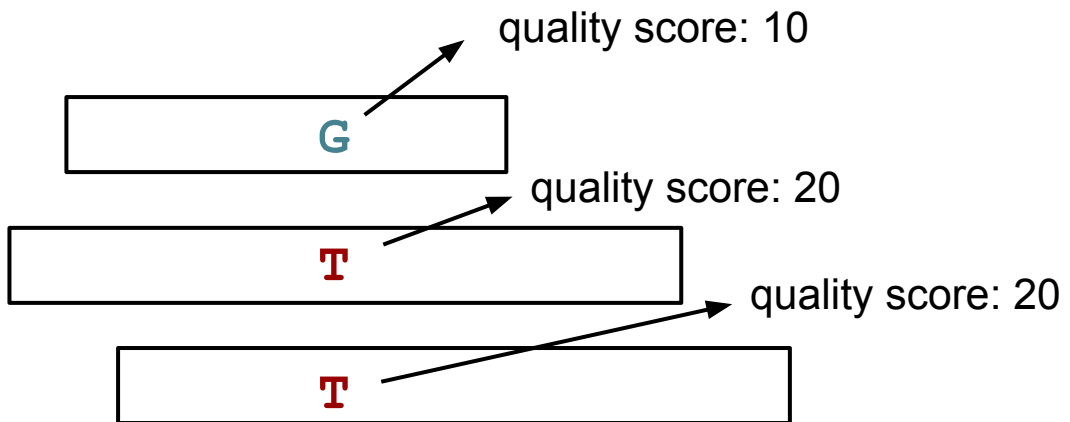
The likelihood $P(D|G)$

Toy example:



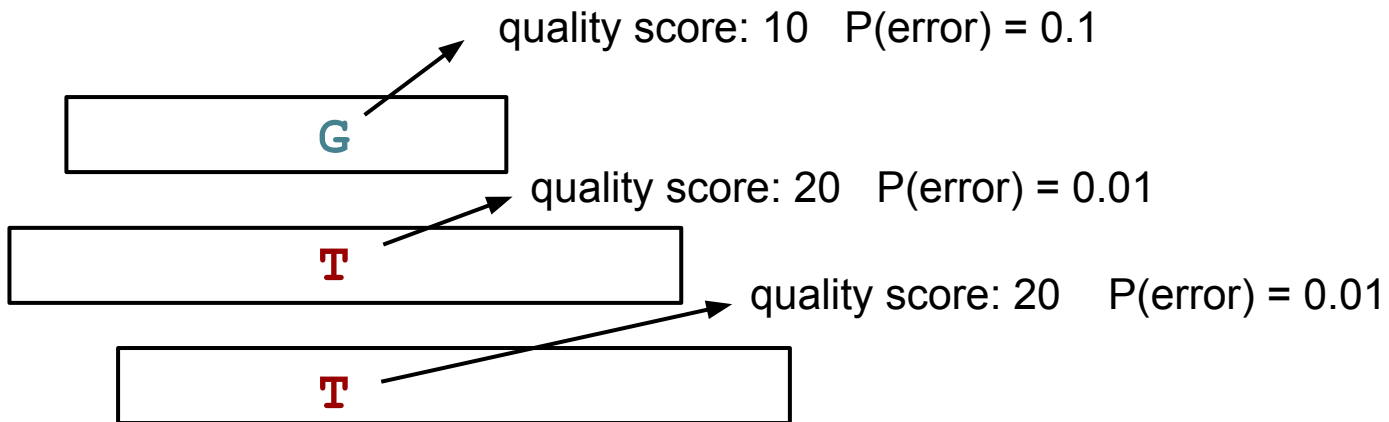
The likelihood $P(D|G)$

Toy example:



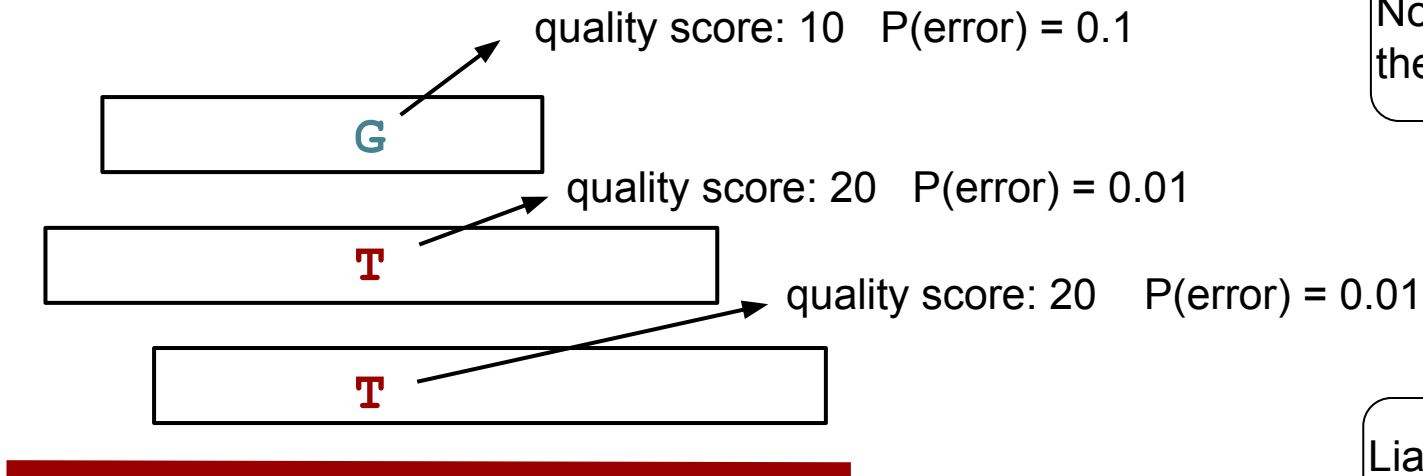
The likelihood $P(D|G)$

Toy example:



The likelihood $P(D|G)$

Toy example:



The 2 Ts are sequencing errors!
The genotype is GG



Nope! They are all correct and
the genotype is GT



Liar! The G is a sequencing error!
TT is the genotype



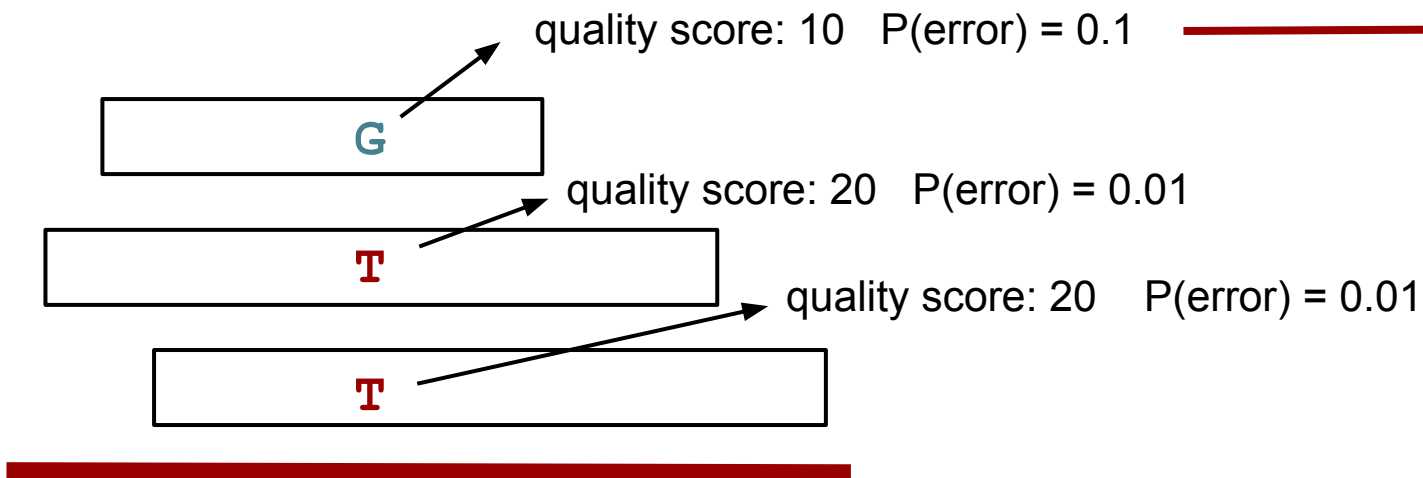
The likelihood $P(D|G)$

Error model

Toy example:

What I think is the base

probability of the data given the base

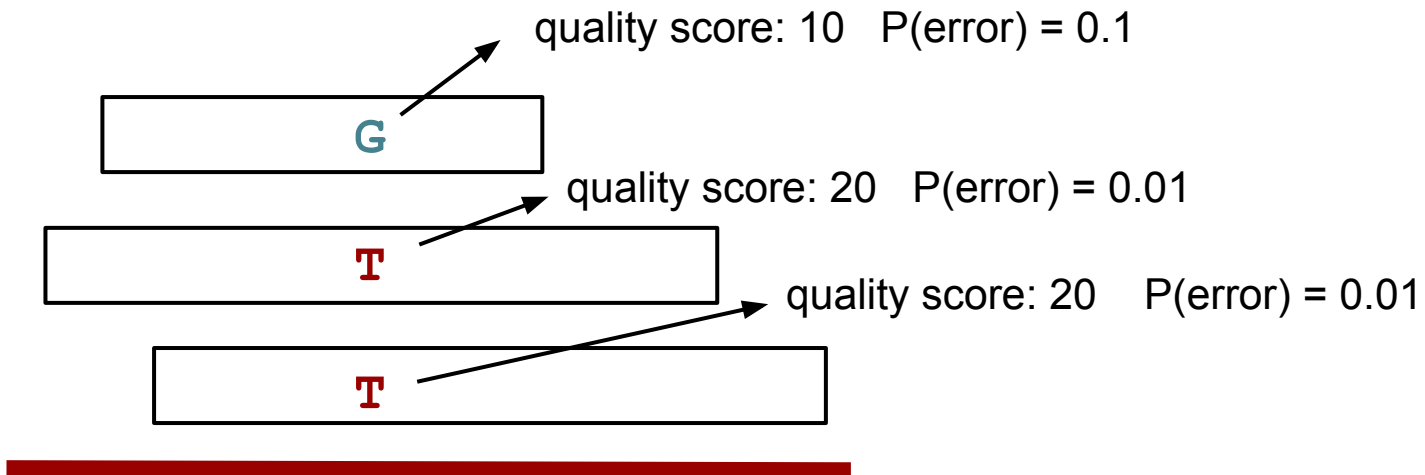


$P(\text{G} \text{A}) =$	$0.1 \frac{1}{3}$
$P(\text{G} \text{C}) =$	$0.1 \frac{1}{3}$
$P(\text{G} \text{G}) =$	0.9
$P(\text{G} \text{T}) =$	$0.1 \frac{1}{3}$

Let's evaluate 3 possible genotypes:

The likelihood $P(D|G)$

Toy example:



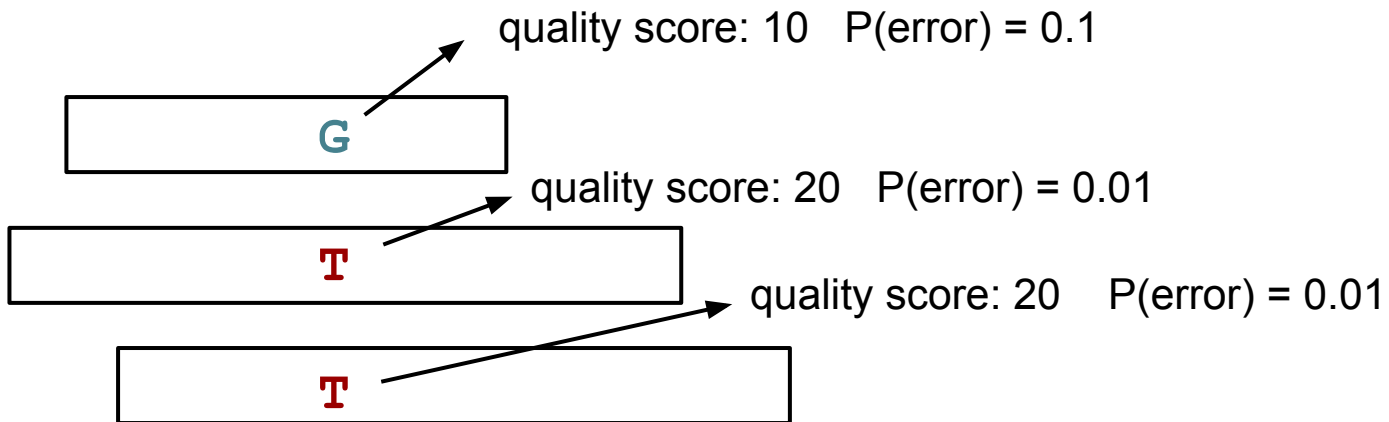
GG

GT

TT

$$P(D|GG)$$

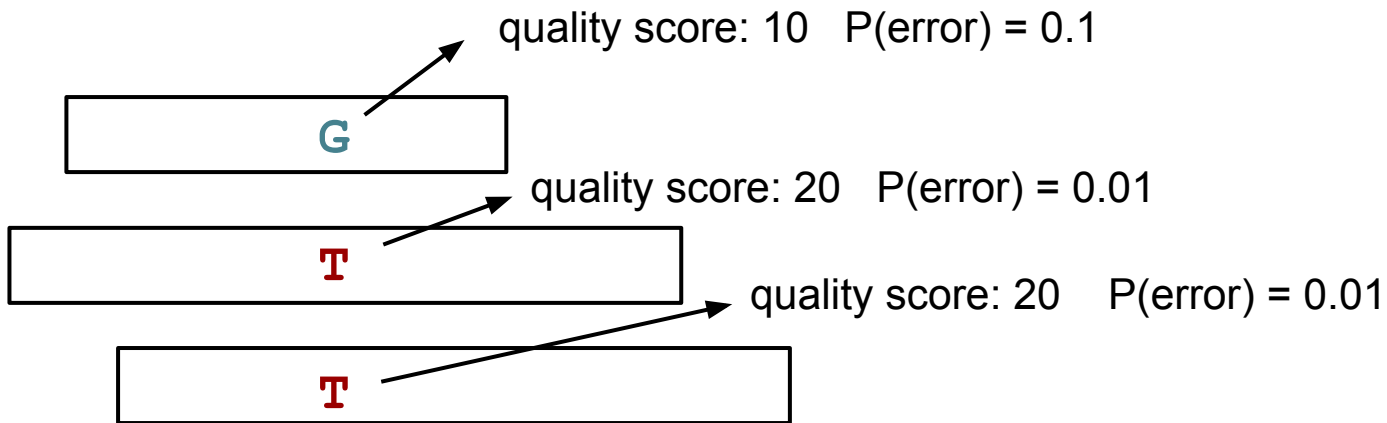
The likelihood $P(D|G)$



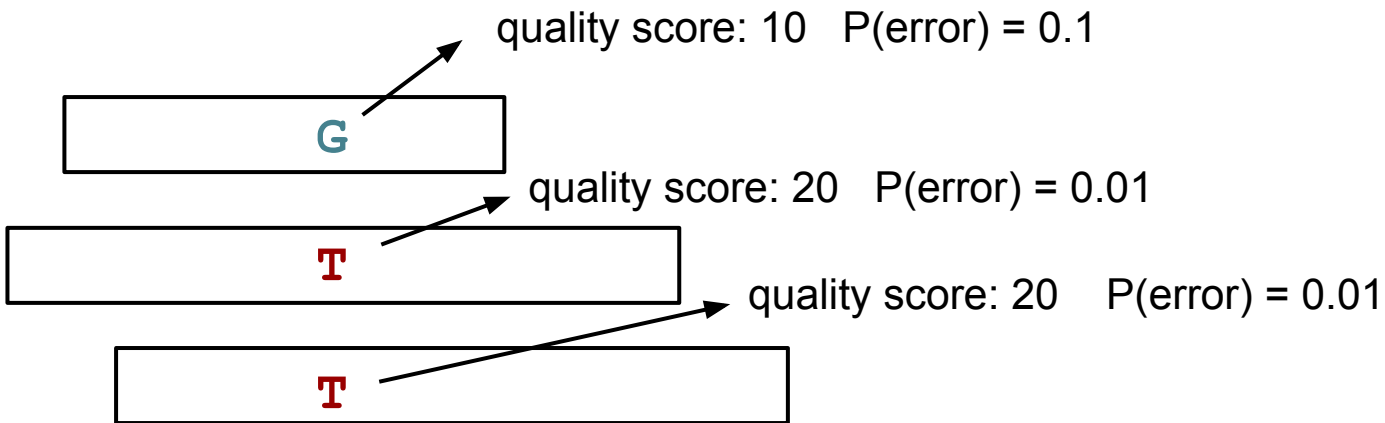
The likelihood $P(D|G)$

$$P(D|GG)$$

$$\frac{1}{2} G \quad \frac{1}{2} G$$



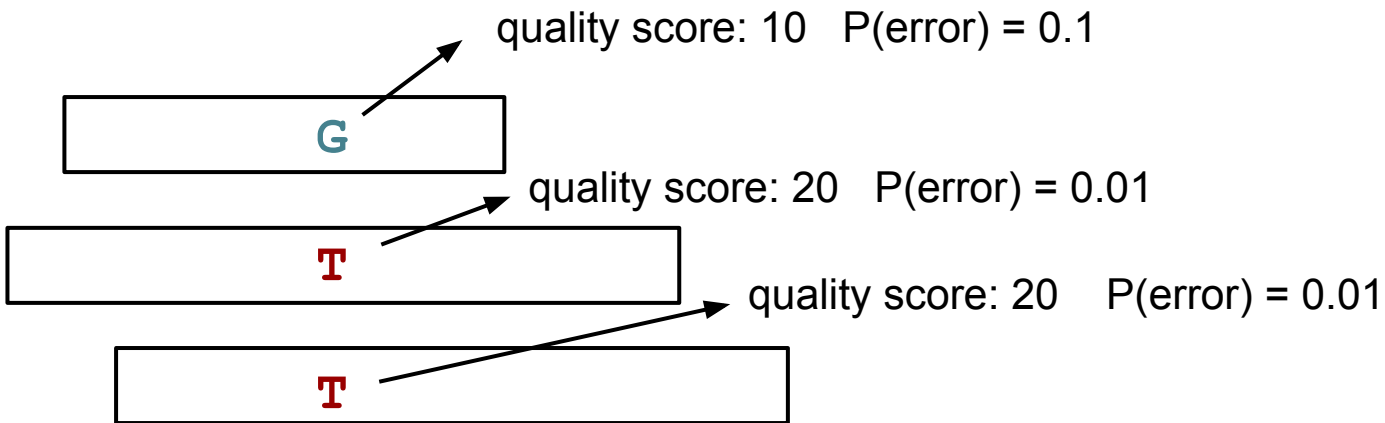
The likelihood $P(D|G)$



$$P(D|GG)$$

$\frac{1}{2}$	G		$\frac{1}{2}$	G	
	✓	0.9		✓	0.9

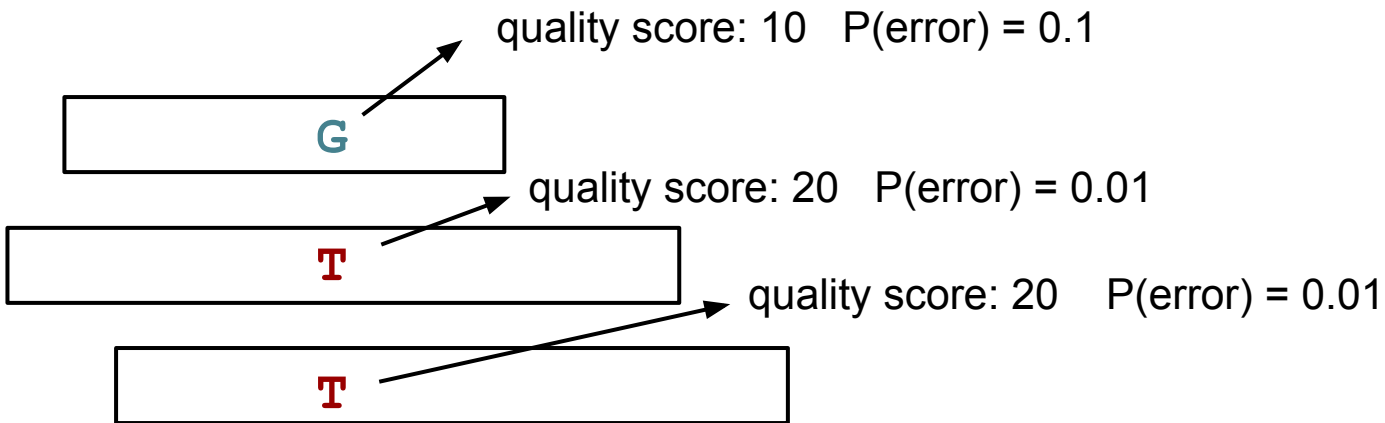
The likelihood $P(D|G)$



$$P(D|GG)$$

$\frac{1}{2}$ G		$\frac{1}{2}$ G	
✓	0.9	✓	0.9
✗	$\frac{0.01}{3}$	✗	$\frac{0.01}{3}$

The likelihood $P(D|G)$

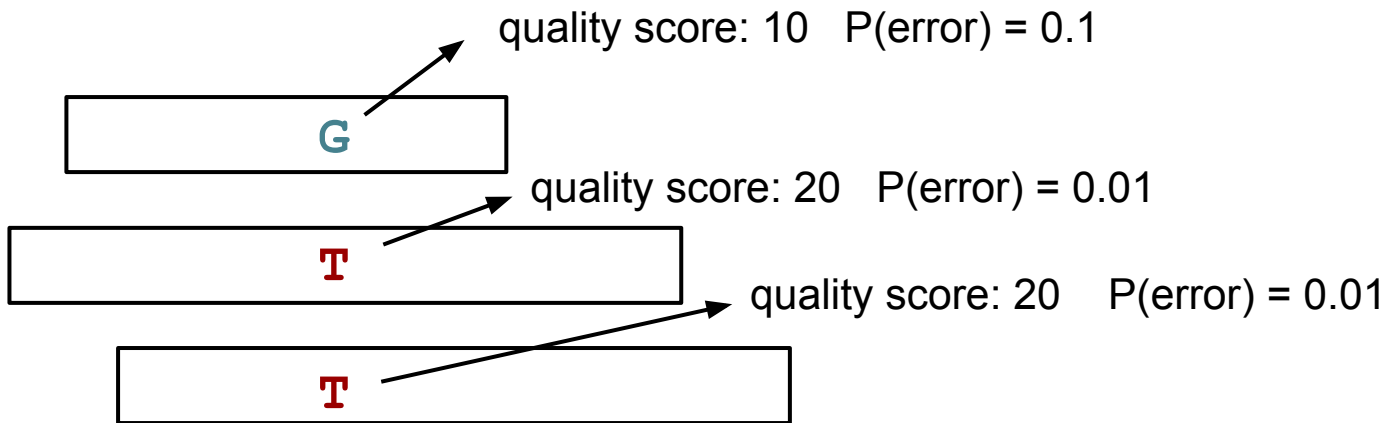


$$P(D|GG)$$

$\frac{1}{2} G$		$\frac{1}{2} G$	
✓	0.9	✓	0.9
✗	$\frac{0.01}{3}$	✗	$\frac{0.01}{3}$
✗	$\frac{0.01}{3}$	✗	$\frac{0.01}{3}$

The likelihood $P(D|G)$

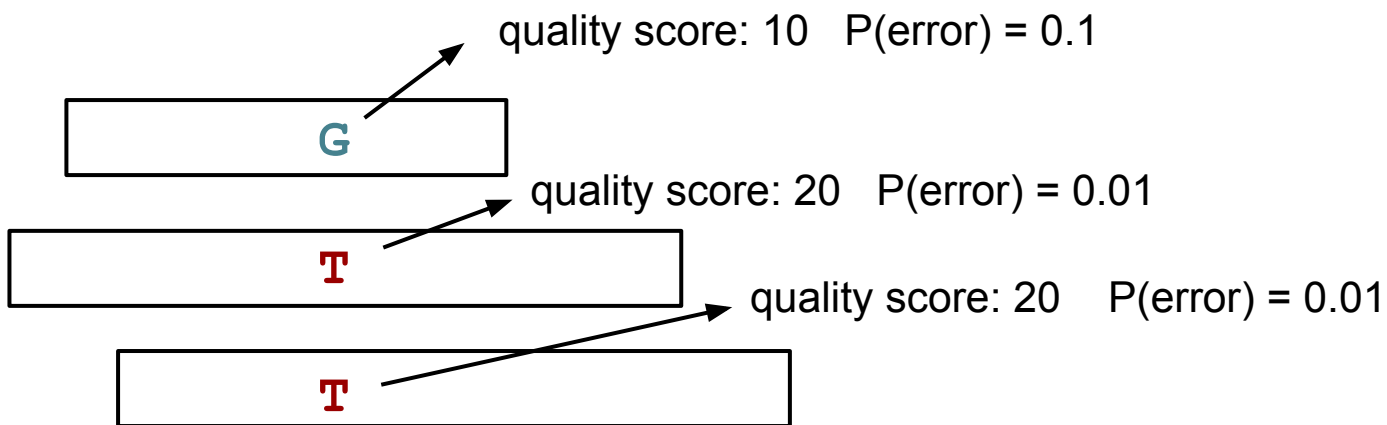
$$P(D|GG)$$



$\frac{1}{2}$ G		$\frac{1}{2}$ G	
✓	0.9	✓	0.9
✗	$\frac{0.01}{3}$	✗	$\frac{0.01}{3}$
✗	$\frac{0.01}{3}$	✗	$\frac{0.01}{3}$

$$\left(\frac{1}{2}0.9 + \frac{1}{2}0.9\right)\left(\frac{1}{2}\frac{0.01}{3} + \frac{1}{2}\frac{0.01}{3}\right)\left(\frac{1}{2}\frac{0.01}{3} + \frac{1}{2}\frac{0.01}{3}\right) = 0.000001$$

The likelihood $P(D|G)$

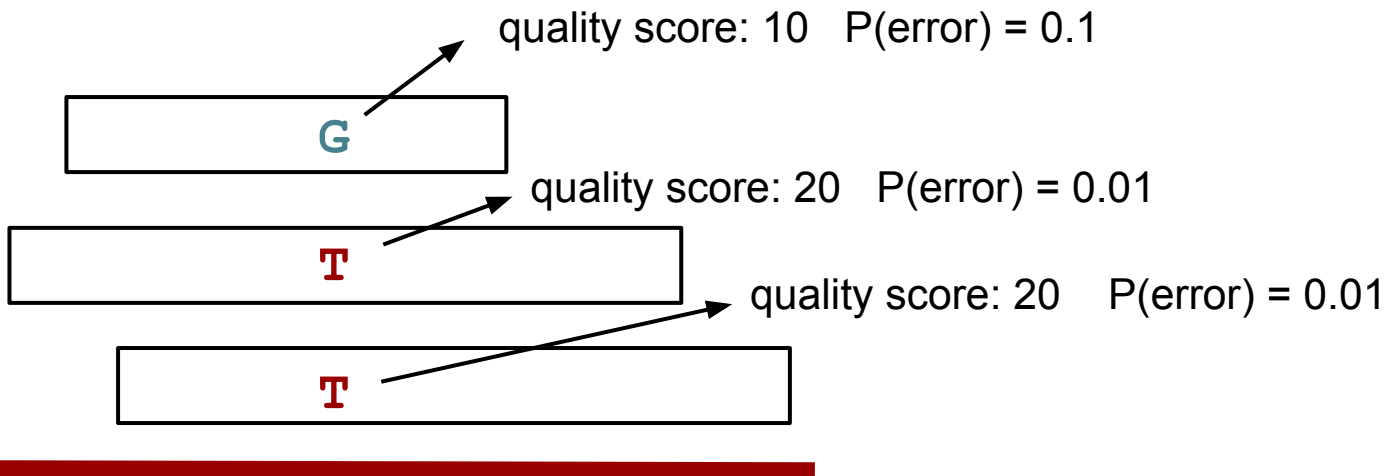


$$P(D|G^T)$$

$\frac{1}{2}$ G		$\frac{1}{2}$ T	
✓	0.9	✗	$\frac{0.1}{3}$
✗	$\frac{0.01}{3}$	✓	0.99
✗	$\frac{0.01}{3}$	✓	0.99

$$\left(\frac{1}{2}0.9 + \frac{1}{2}\frac{0.1}{3}\right)\left(\frac{1}{2}\frac{0.01}{3} + \frac{1}{2}0.99\right)\left(\frac{1}{2}\frac{0.01}{3} + \frac{1}{2}0.99\right) = 0.1151163$$

The likelihood $P(D|G)$

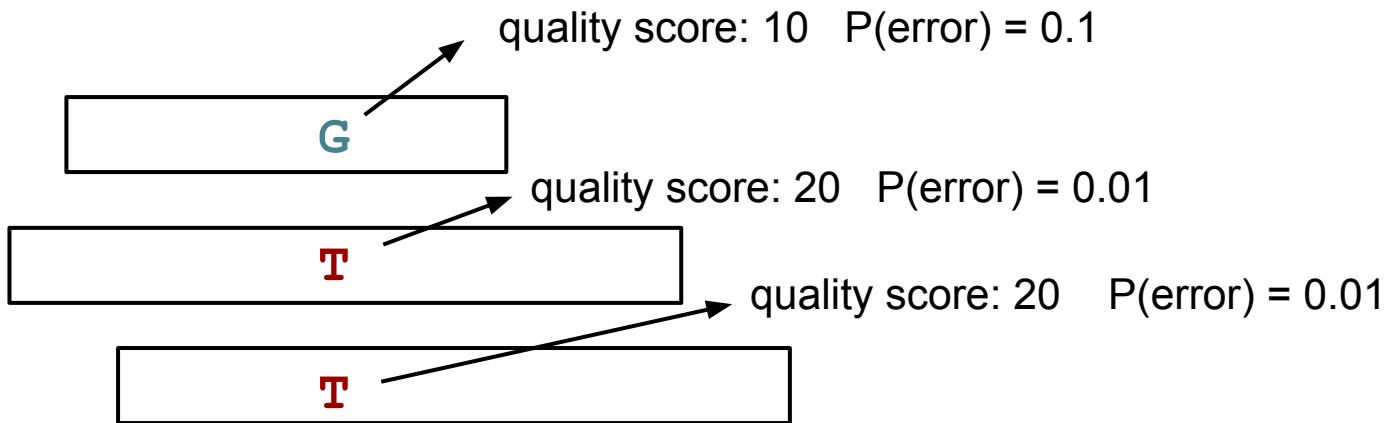


$$P(D|\mathbf{TT})$$

$\frac{1}{2}$	T	$\frac{1}{2}$	T
X	$\frac{0.1}{3}$	X	$\frac{0.1}{3}$
✓	0.99	✓	0.99
✓	0.99	✓	0.99

$$\left(\frac{1}{2} \frac{0.1}{3} + \frac{1}{2} \frac{0.1}{3}\right) \left(\frac{1}{2} 0.99 + \frac{1}{2} 0.99\right) \left(\frac{1}{2} 0.99 + \frac{1}{2} 0.99\right) = 0.03267$$

The likelihood $P(D|G)$



$$P(D|GG) = 0.00001$$

$$P(D|GT) = 0.11511$$

$$P(D|TT) = 0.0327$$

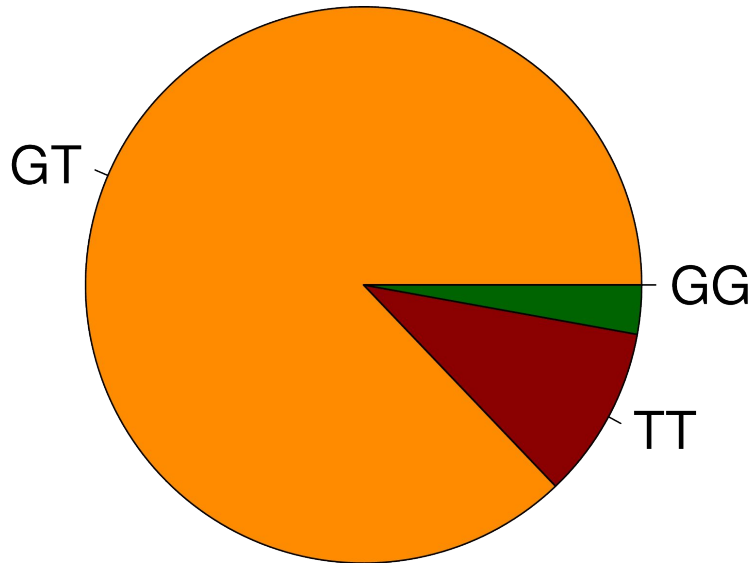
A likelihood in itself
is not meaningful,
you need to
compare it to other
models

The likelihood $P(D|G)$

$$P(D|\text{GG}) = 0.00001$$

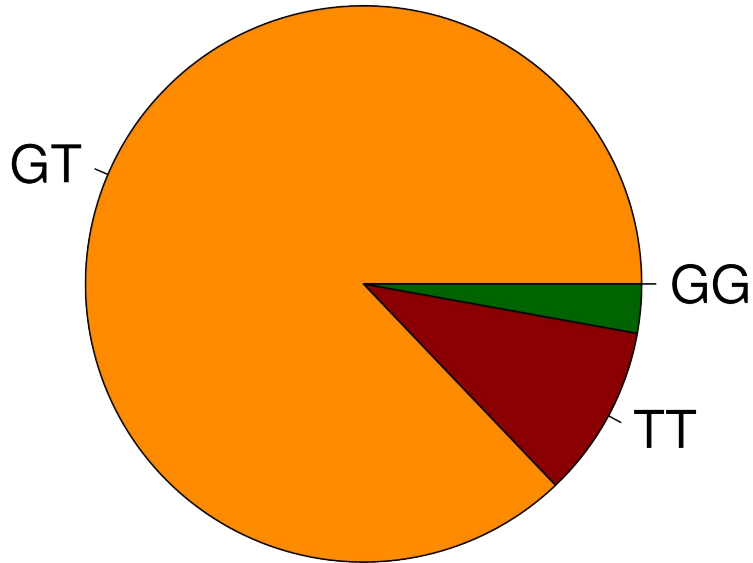
$$P(D|\text{GT}) = 0.11511$$

$$P(D|\text{TT}) = 0.0327$$



$$P(D) = P(\text{GG})P(D|\text{GG}) + P(\text{GT})P(D|\text{GT}) + P(\text{TT})P(D|\text{TT})$$

The likelihood $P(D|G)$



We will neglect the genotype prior this time

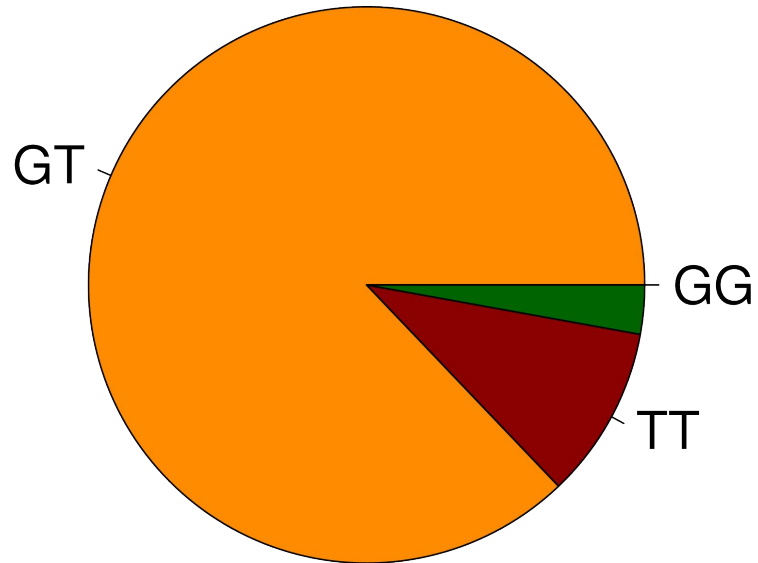
$$P(D|GG) = 0.00001$$

$$P(D|GT) = 0.11511$$

$$P(D|TT) = 0.0327$$

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

The likelihood $P(D|G)$

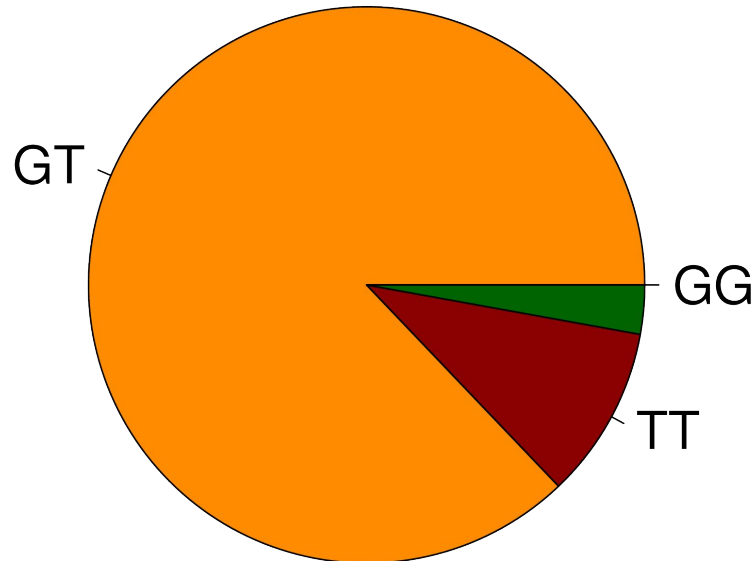


$$P(\mathbf{GG}|D) = 6.7e-05$$

$$P(\mathbf{GT}|D) = 0.77888$$

$$P(\mathbf{TT}|D) = 0.22104$$

The likelihood $P(D|G)$



$$P(\mathbf{GG}|D) = 6.7 \times 10^{-5}$$

$$P(\mathbf{GT}|D) = 0.77888$$

$$P(\mathbf{TT}|D) = 0.22104$$

Important point: More coverage \longrightarrow More multiplications \longrightarrow The relative difference between models become larger

The likelihood $P(D|G)$

	PHRED	PHRED-scaled
$P(\mathbf{GG} D) = 6.7e-05$	41.70	40.60
$P(\mathbf{GT} D) = 0.77888$	1.09	0.00
$P(\mathbf{TT} D) = 0.22104$	6.56	5.47



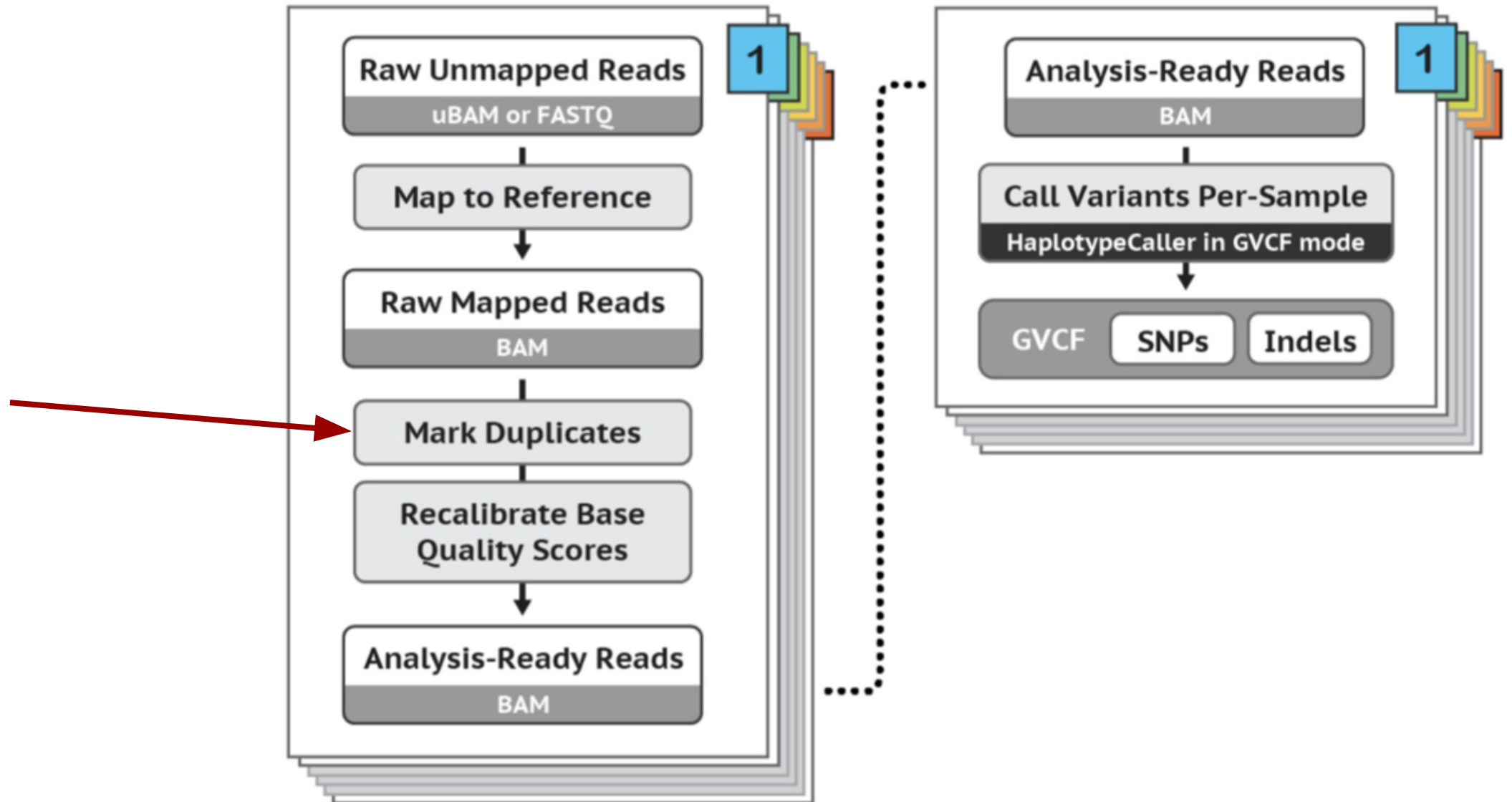
Details I did not cover

- Error model
 - Most genotypers do not simply use raw quality scores

Most common genotypers

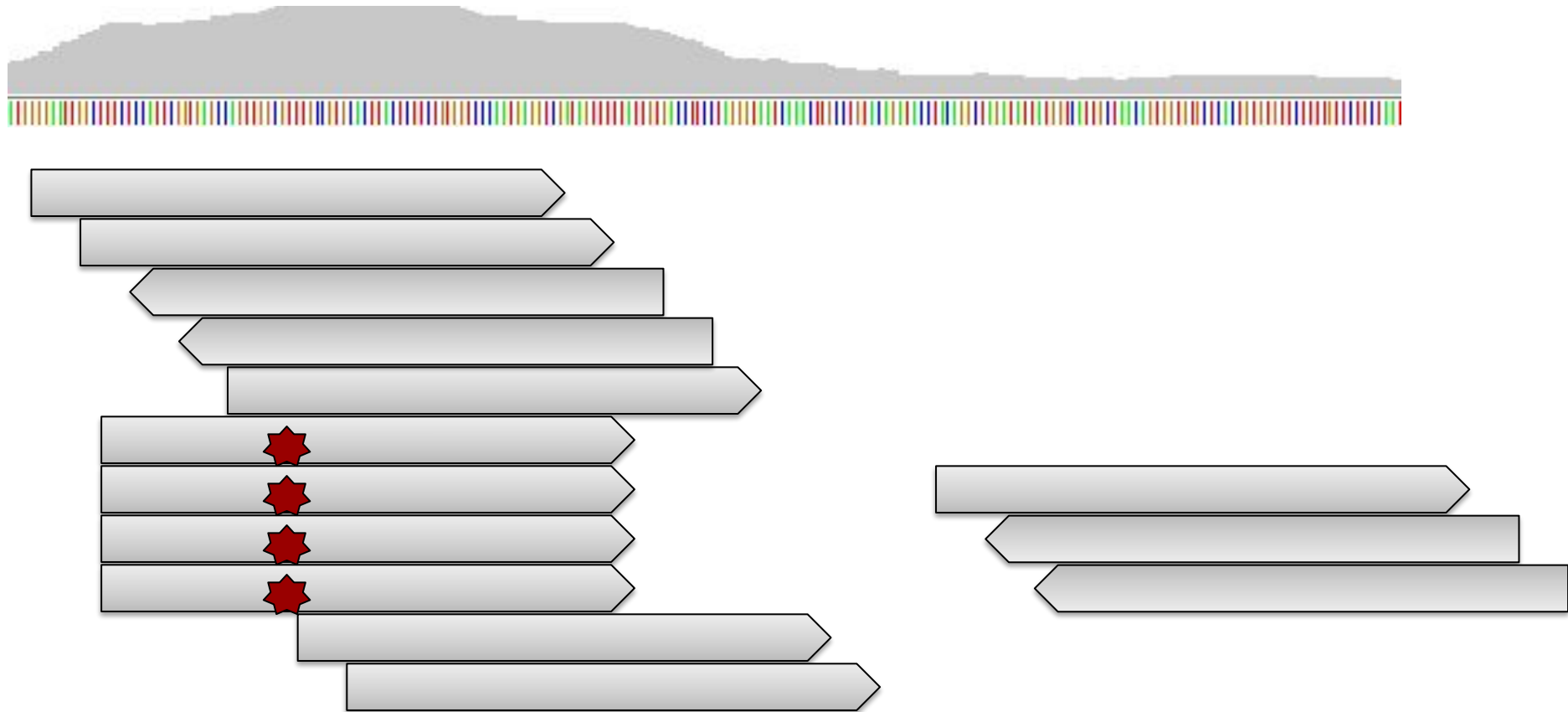
- GATK
- SAMtools/BCFtools
- graph typer
- FreeBayes


GATK's recommended workflow

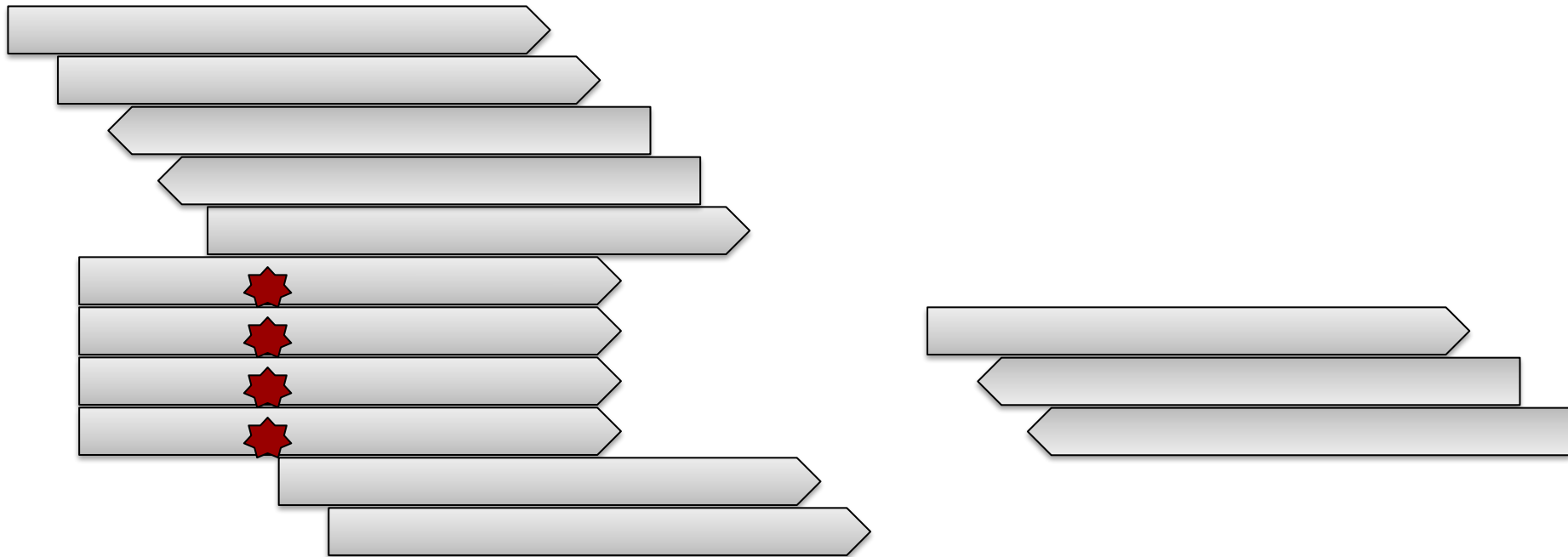
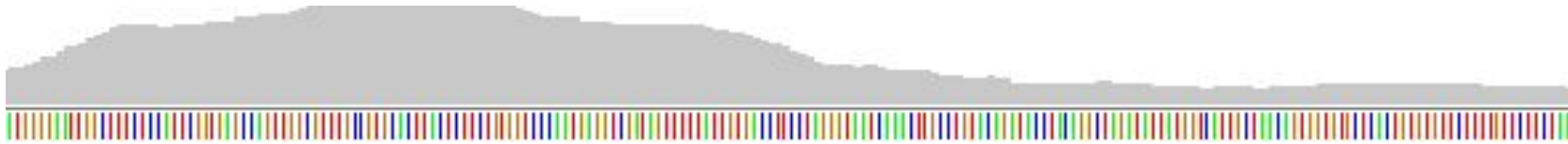



<https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels->

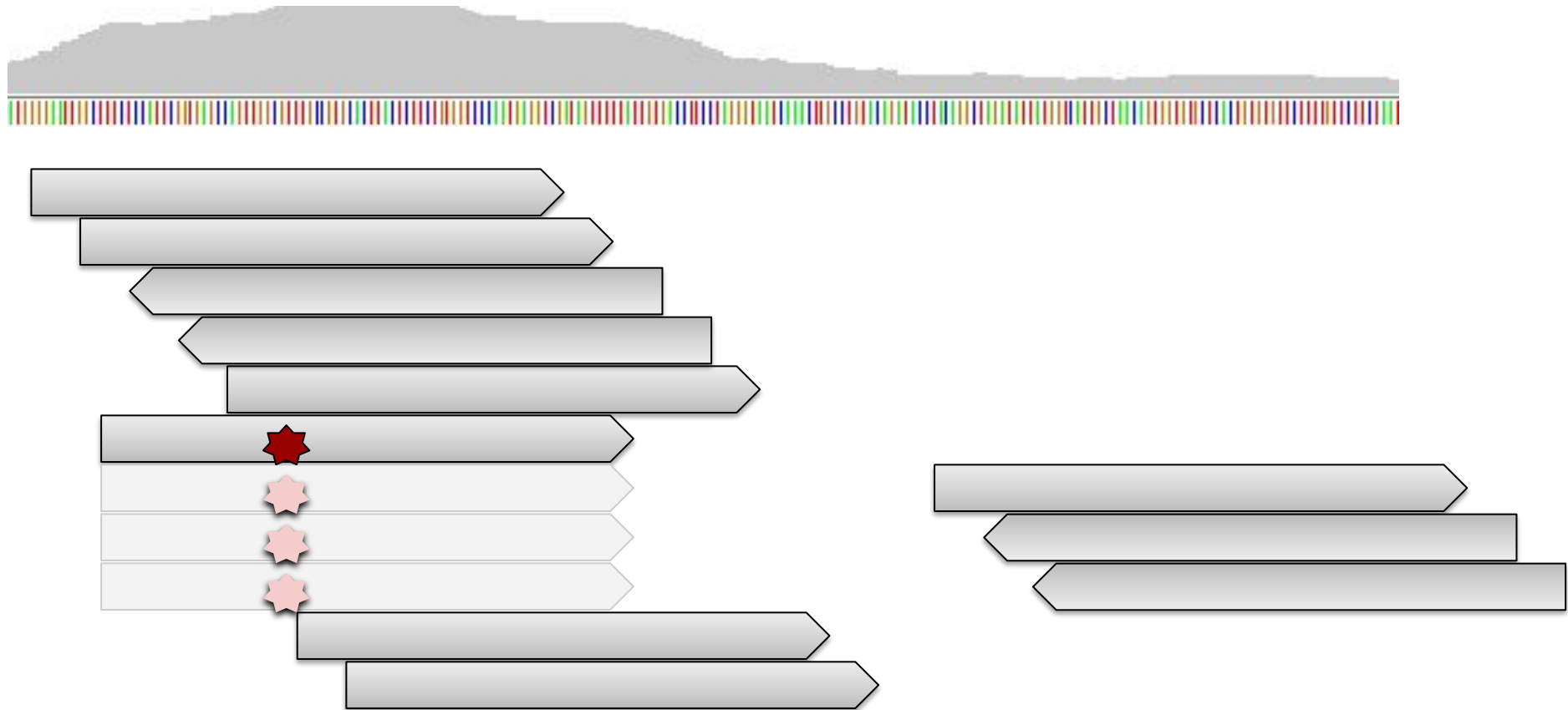
The PCR amplification step included in the majority of NGS library construction techniques can introduce duplicates in the data.



We want: remove or mark them to avoid false calls
ex: the site below is probably heterozygous (i.e. the  is the second allele)



Genotypers will ignore reads marked as duplicates
ex: the site below is probably homozygous (i.e. the  is a seq. error)



Duplicate/marking removal

Basic concepts of duplicate marking algorithm:

- Identify genomic position and strand for 5'-most bases.
- Mark reads that are duplicates of each other.
- Within a group of duplicate reads, the read with the highest sum of base quality scores is retained.

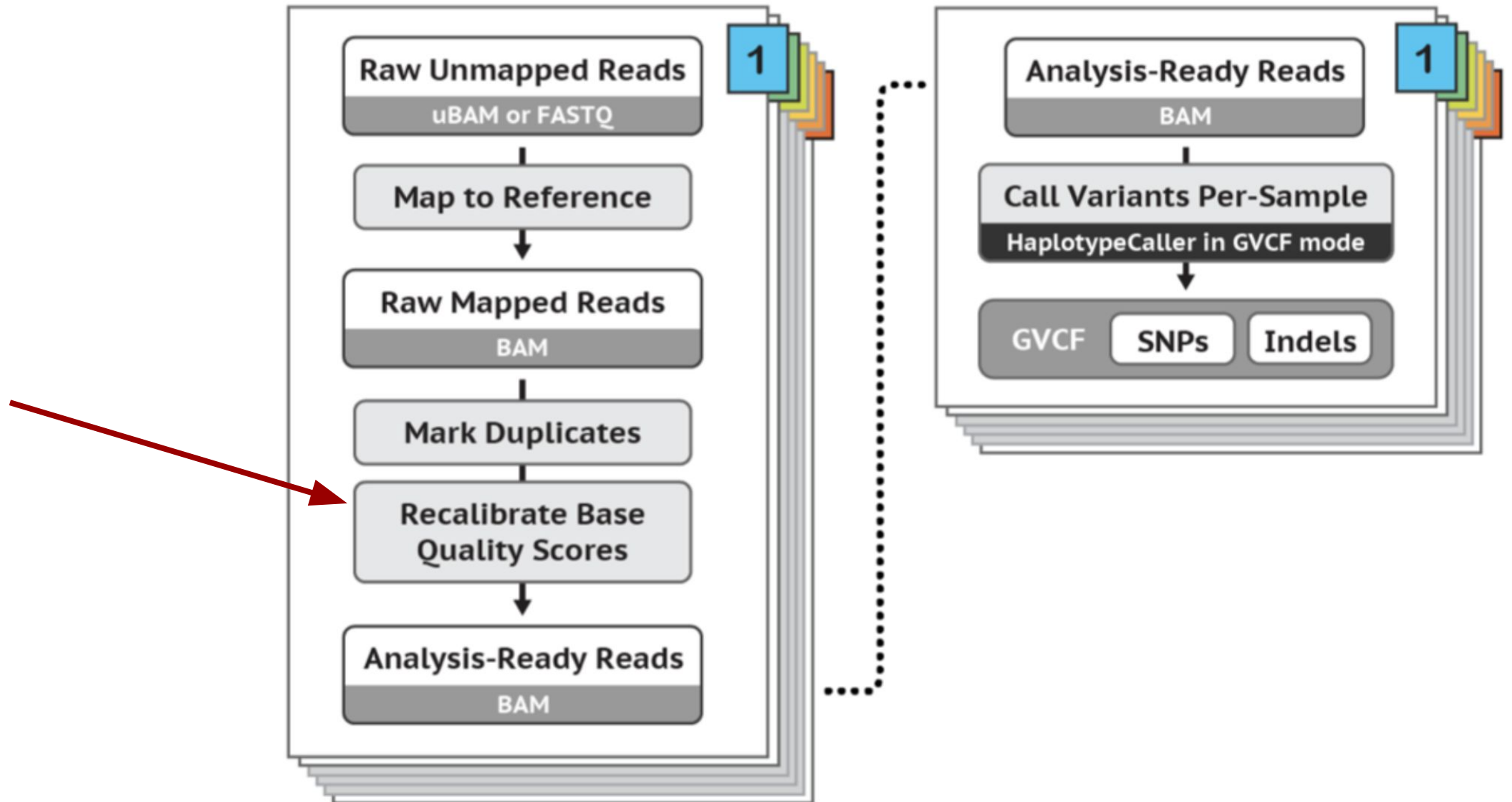
<http://picard.sourceforge.net/>

Duplicate/marking removal

Problems:

- Does not account for sequencing errors.
- Does not account for natural duplicates.
- Does not account for duplicate reads with different mapping locations.

GATK's recommended workflow



<https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels->

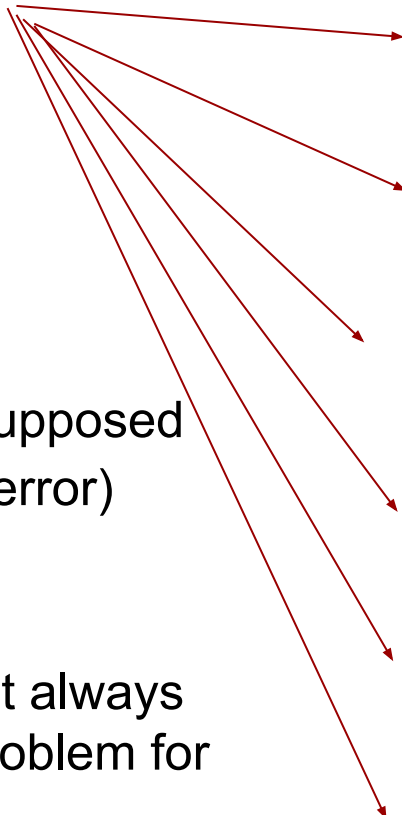
Base quality score recalibration?

- remember those?
- There are supposed to reflect $P(\text{error})$
- They are not always accurate: problem for genotyping

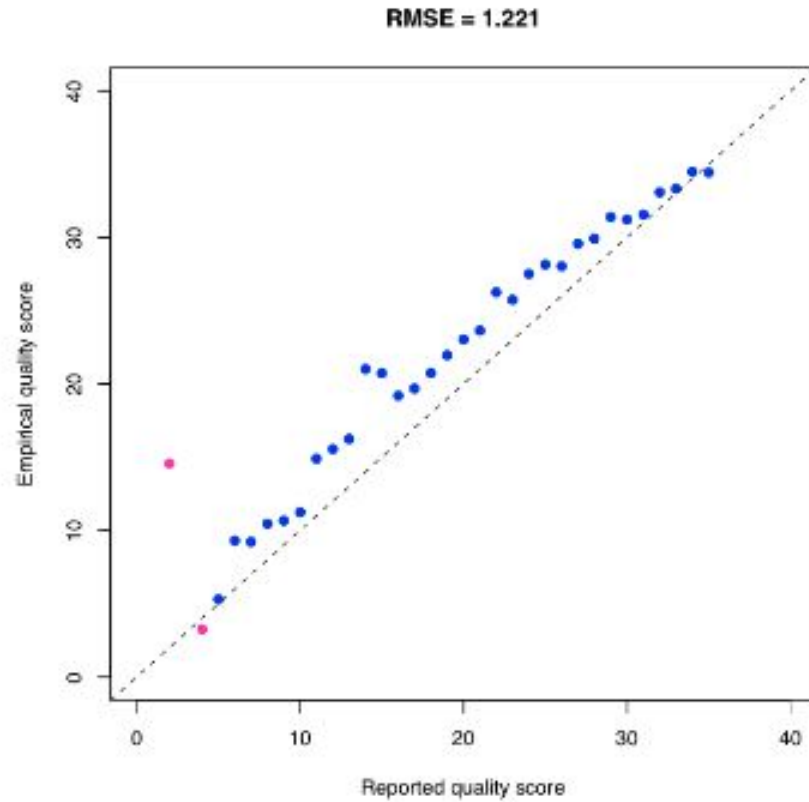
```

@A80CRCABXX:2:1:11125:1940#CCCCCCCC/1
ATCATAAAGAAAATATATTTTGCAAATGAAGTCATTAATAATTTGAGAT
+
ccccddadadd__dd\ddcddddddTdd`dd^dddcdddddcdaddd
@A80CRCABXX:2:1:12491:1939#TGTGGCTT/1
CTGACAGCATTGTTCTGTTGCAGGATTACGCTCCCTTAGATCGGAAGA
+
fffffefffffffffffffffffffffeeffffffffffffffeffdffd
@A80CRCABXX:2:1:13158:1938#AAAAAAA/1
ATGAGTGAAAAGCGTCTAATTCTCTATGCCATGCCTATTTTCTTTGTAA
+
dacdcacdeed`^c^dadd\` `bbbc`aac\^` `b` `cbc__[\bb
@A80CRCABXX:2:1:14354:1937#ACGGTTTT/1
CTCTCTTCTCTGGCTGACTGCCTGTCTCTCTCTCTCTCTCTCTCTC
+
ffffffffffffffffffffefffeefceccedeefffdffffffffffffdf
@A80CRCABXX:2:1:14546:1939#AACCGCTT/1
AGTAGTTTCAATACTCAAATTACCTCTACAGCCATGATGATAACAGCAG
+
ffffffffffffeffffffeffffffffffffffeffffffffffffffef
@A80CRCABXX:2:1:14819:1939#AACTAGAA/1
ATAGATGTTATTCAACTCCTTCAGGTTGTCTTGAAGTACTGACTCATG
+
fffffffffffffffffffffffffffffffffffffffffffffffffffff

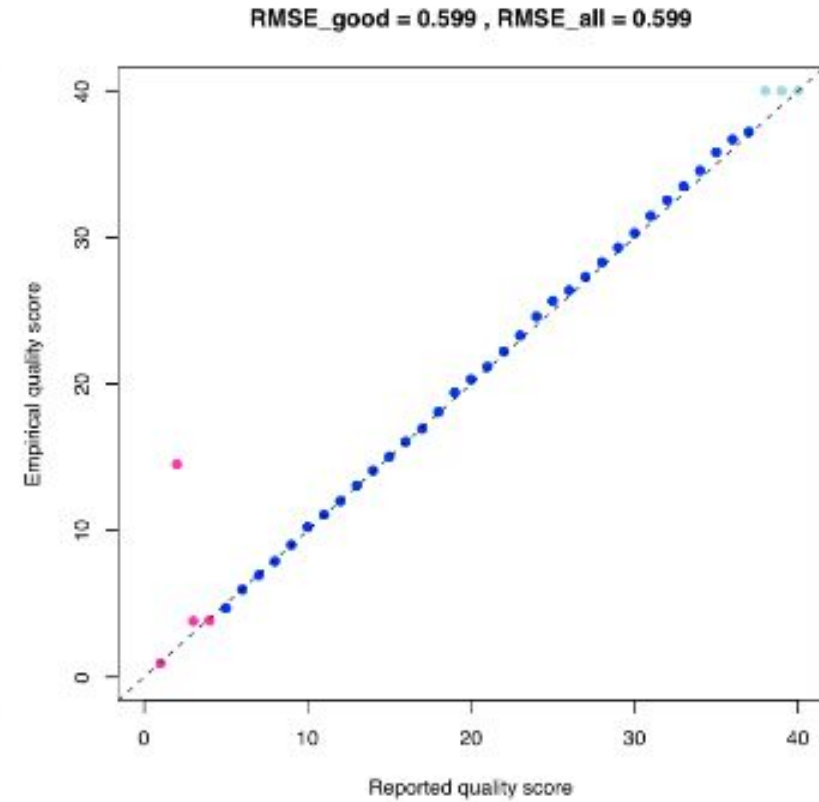
```



Reported Quality vs. Empirical Quality



Original Data



After GATK Recalibration

The Missing Diversity in Human Genetic Studies

[Giorgio Sirugo](#)  ⁶  • [Scott M. Williams](#)  ⁶  • [Sarah A. Tishkoff](#)  ⁶  • [Show footnotes](#)

DOI: <https://doi.org/10.1016/j.cell.2019.02.048> •



Check for updates

The majority of studies of genetic association with disease have been performed in Europeans. This European bias has important implications for risk prediction of diseases across global populations. In this commentary, we justify the need to study more diverse populations using both empirical examples and theoretical reasoning.

Base quality score recalibration

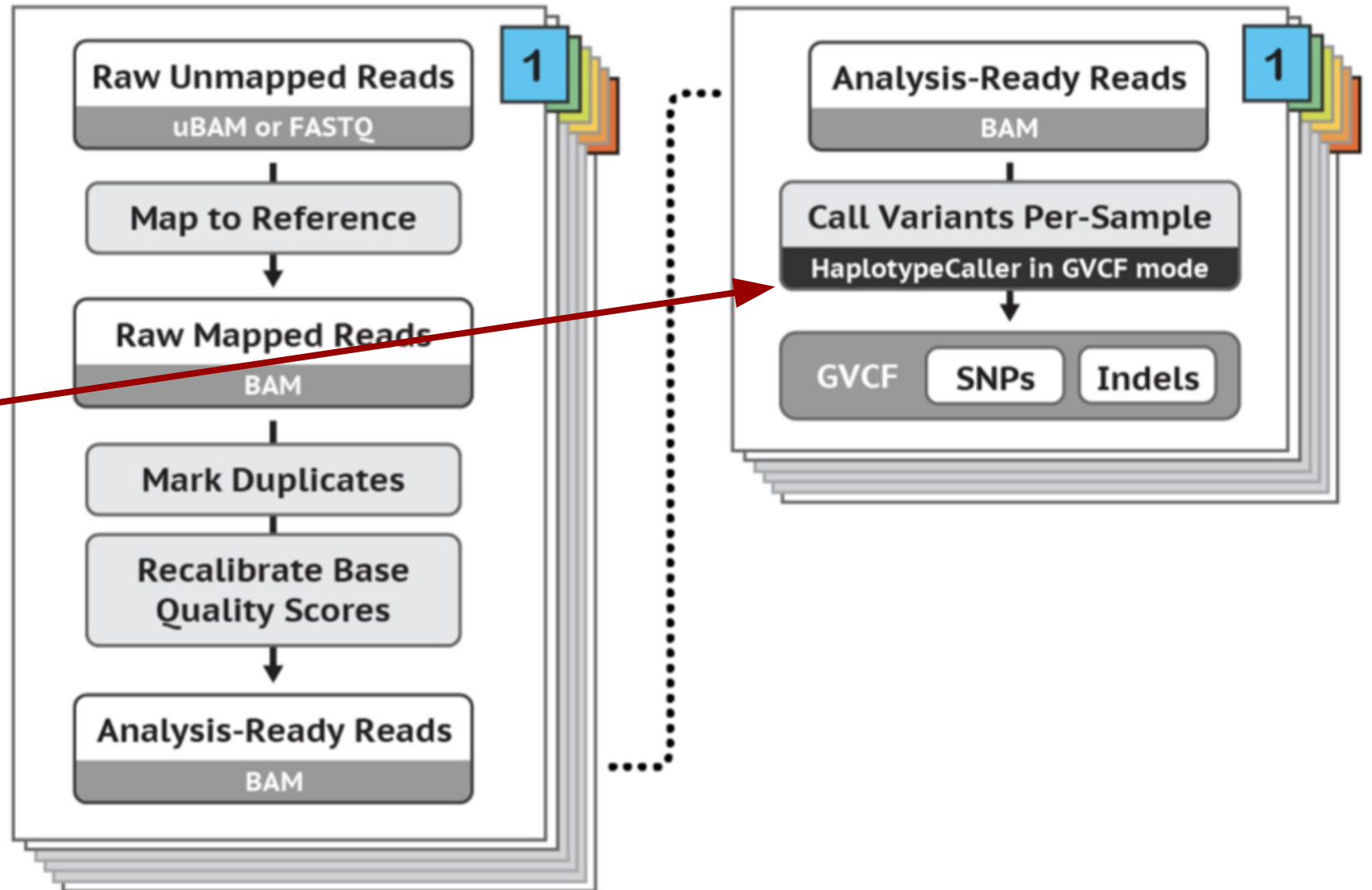
To work we need:

- East Asian or European (as in mostly West European) samples
- WGS
- Sufficient coverage

My biased opinion:

- Just don't bother

GATK's recommended workflow



We covered this before

<https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels->

Variant call format (VCF)

- Details which variants have been called
- Can be bgzip (block gzip) and indexed using tabix
- Using tabix, queries can be made like:
 - return all variants in the region chr22:323,340-361,152

Variant call format (VCF)

```

20 51391523 . A G 173.96 . AC=2;DP=5;MQ=52.03 GT:AD:DP:GQ:PL 1/1:0,5:5:15:188,15,0
20 51392469 . C T 146.14 . AC=2;DP=4;MQ=60.00 GT:AD:DP:GQ:PL 1/1:0,4:4:12:160,12,0
20 51394015 . T C 97.64 . AC=1;DP=6;MQ=60.00 GT:AD:DP:GQ:PL 0/1:3,3:6:66:105,0,66
20 51395647 . A C 89.64 . AC=1;DP=7;MQ=57.28 GT:AD:DP:GQ:PL 0/1:4,3:7:97:97,0,100
20 51397399 . C T 93.64 . AC=1;DP=7;MQ=60.00 GT:AD:DP:GQ:PL 0/1:4,3:7:99:101,0,120
20 51402308 . C T 161.64 . AC=1;DP=9;MQ=60.00 GT:AD:DP:GQ:PL 0/1:3,6:9:63:169,0,63

```

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

name of chromosome (ex: chr1, chr2 ...)

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

coordinate on chromosome

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

ID (ex: rs23534)

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

reference base

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

alternative base

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

quality field

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

Filter (ex: 'LowQual')

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

Info field ex:

AC= allele count

DP = depth

MQ = root mean square of the mapping quality

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

Format field, what do the next fields mean?

Variant call format (VCF)

```

20 51391523 . A G 173.96 . AC=2;DP=5;MQ=52.03 GT:AD:DP:GQ:PL 1/1:0,5:5:15:188,15,0
20 51392469 . C T 146.14 . AC=2;DP=4;MQ=60.00 GT:AD:DP:GQ:PL 1/1:0,4:4:12:160,12,0
20 51394015 . T C 97.64 . AC=1;DP=6;MQ=60.00 GT:AD:DP:GQ:PL 0/1:1,3:6:66:105,0,66
20 51395647 . A C 89.64 . AC=1;DP=7;MQ=57.28 GT:AD:DP:GQ:PL 0/1:4,3:7:97:97,0,100
20 51397399 . C T 93.64 . AC=1;DP=7;MQ=60.00 GT:AD:DP:GQ:PL 0/1:4,3:7:99:101,0,120
20 51402308 . C T 161.64 . AC=1;DP=9;MQ=60.00 GT:AD:DP:GQ:PL 0/1:1,6:9:63:169,0,63
  
```

Most likely genotype

Variant call format (VCF)

```

20 51391523 . A G 173.96 . AC=2;DP=5;MQ=52.03 GT:AD:DP:GQ:PL 1/1:0,5:5:15:188,15,0
20 51392469 . C T 146.14 . AC=2;DP=4;MQ=60.00 GT:AD:DP:GQ:PL 1/1:0,4:4:12:160,12,0
20 51394015 . T C 97.64 . AC=1;DP=6;MQ=60.00 GT:AD:DP:GQ:PL 0/1:3,3:6:66:105,0,66
20 51395647 . A C 89.64 . AC=1;DP=7;MQ=57.28 GT:AD:DP:GQ:PL 0/1:4,3:7:97:97,0,100
20 51397399 . C T 93.64 . AC=1;DP=7;MQ=60.00 GT:AD:DP:GQ:PL 0/1:4,3:7:99:101,0,120
20 51402308 . C T 161.64 . AC=1;DP=9;MQ=60.00 GT:AD:DP:GQ:PL 0/1:3,6:9:63:169,0,63
  
```

Allele distribution

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:65:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

Depth

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:GQ:PL	1/1:0,5:5:15:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:GQ:PL	1/1:0,4:4:12:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,3:6:66:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:GQ:PL	0/1:4,3:7:97:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:GQ:PL	0/1:4,3:7:99:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:GQ:PL	0/1:3,6:9:63:169,0,63

Genotype quality

Variant call format (VCF)

20	51391523	.	A	G	173.96	.	AC=2;DP=5;MQ=52.03	GT:AD:DP:CQ:PL	1/1:0,5:5:15	:188,15,0
20	51392469	.	C	T	146.14	.	AC=2;DP=4;MQ=60.00	GT:AD:DP:CQ:PL	1/1:0,4:4:12	:160,12,0
20	51394015	.	T	C	97.64	.	AC=1;DP=6;MQ=60.00	GT:AD:DP:CQ:PL	0/1:3,3:6:66	:105,0,66
20	51395647	.	A	C	89.64	.	AC=1;DP=7;MQ=57.28	GT:AD:DP:CQ:PL	0/1:4,3:7:97	:97,0,100
20	51397399	.	C	T	93.64	.	AC=1;DP=7;MQ=60.00	GT:AD:DP:CQ:PL	0/1:4,3:7:99	:101,0,120
20	51402308	.	C	T	161.64	.	AC=1;DP=9;MQ=60.00	GT:AD:DP:CQ:PL	0/1:3,6:9:63	:169,0,63

PHRED-scaled likelihood

The likelihood $P(D|G)$



PHRED

PHRED-scaled

$$P(\mathbf{GG}|D) = 6.7e-05$$

41.70

40.60

$$P(\mathbf{GT}|D) = 0.77888$$

1.09

0.00

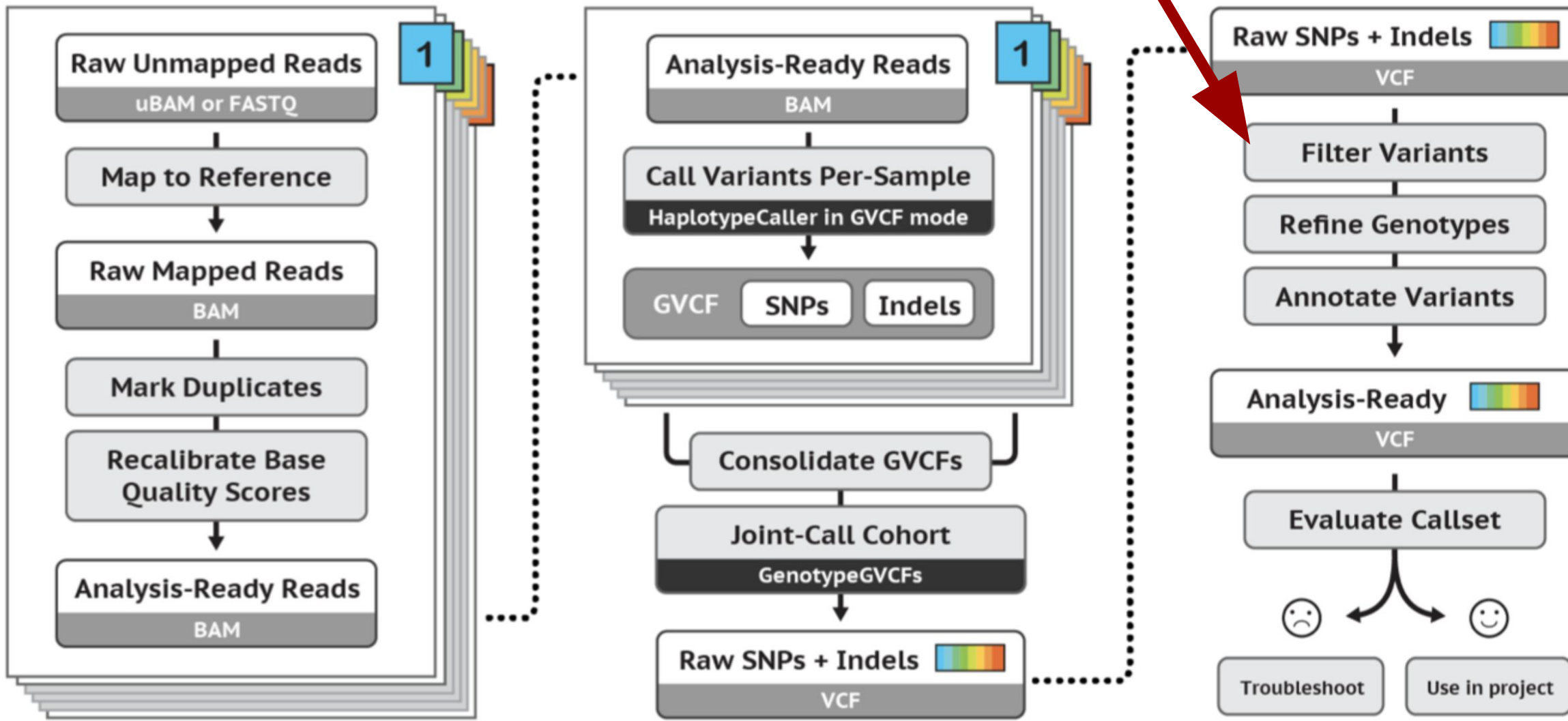
$$P(\mathbf{TT}|D) = 0.22104$$

6.56

5.47



GATK's recommended workflow

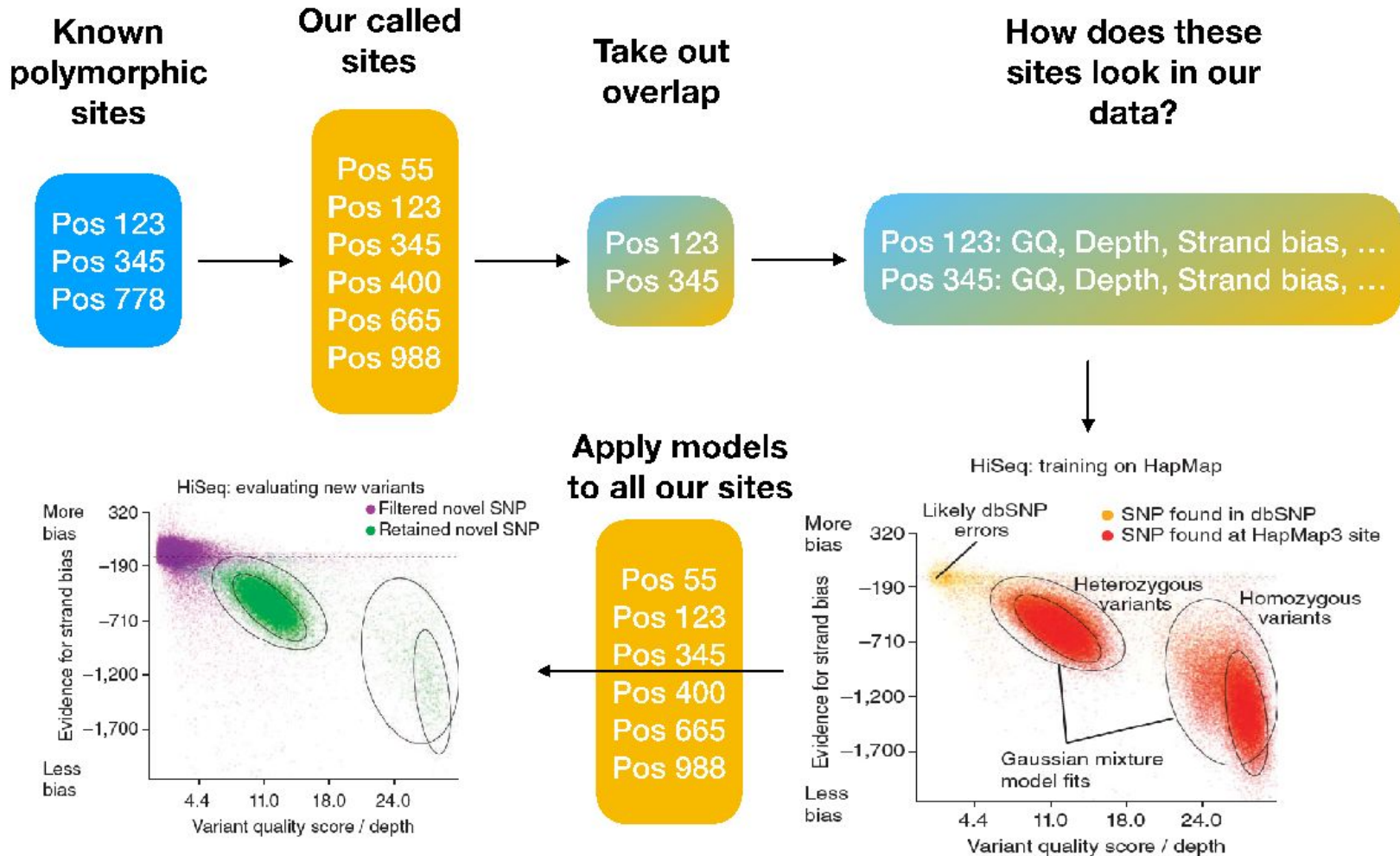


<https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels->

Variant filtration (soft)

- How do we remove false positive calls?
- Use known polymorphic sites to estimate what a real variant and a false variant “looks like”
- Learn how does the known sites (=truth set) look like in our data
- Evaluate on all our data, filter sites that look different!

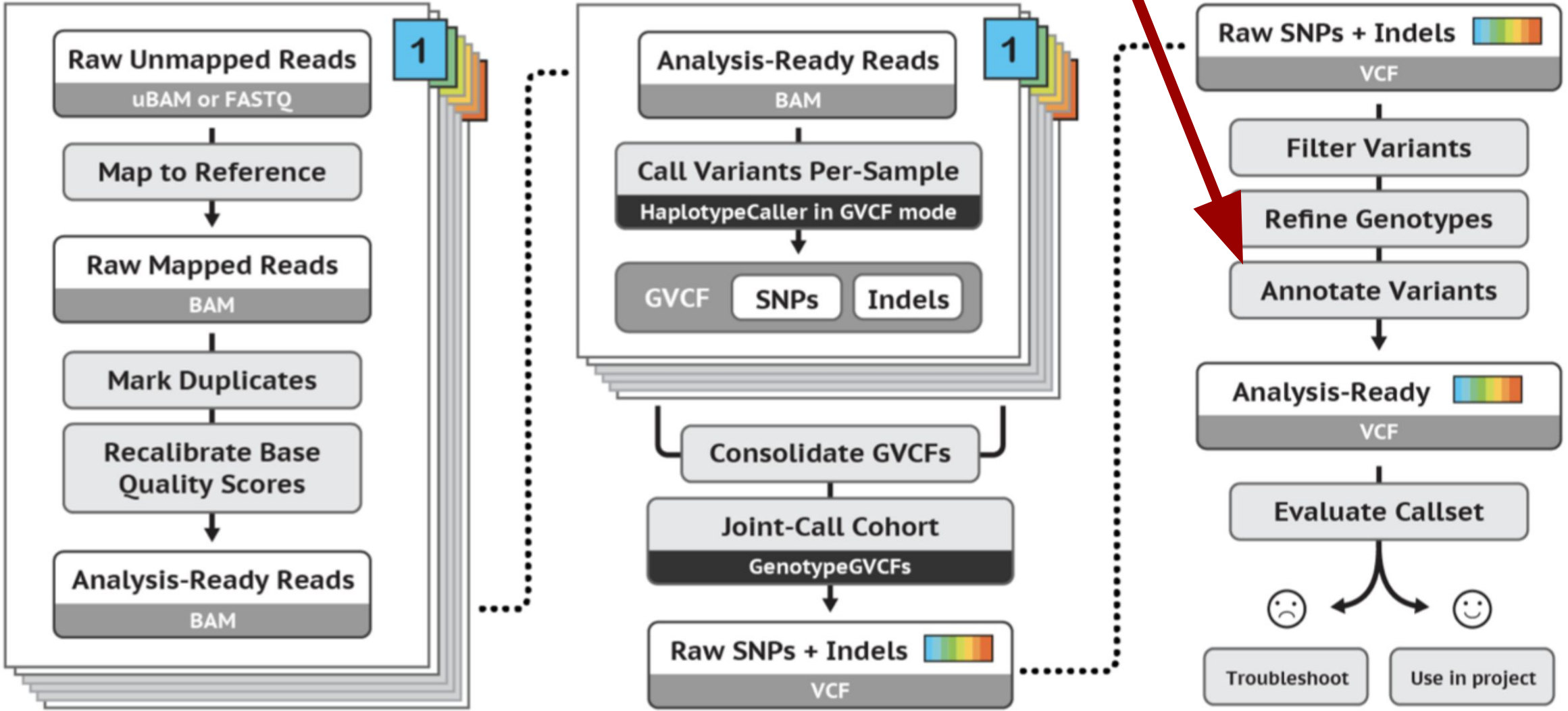
Train a predictor & Test:



Variant filtration (hard)

- Hard filtering:
 - Variant quality score /depth
 - Mapping quality
 - Mappability
 - Strand bias (the variant being seen only on the forward strand or only on the reverse strand)
 - Depth
- BCFtools can perform this
- Depends on the project at hand
- Be careful of introducing a bias in favor of certain types of variants

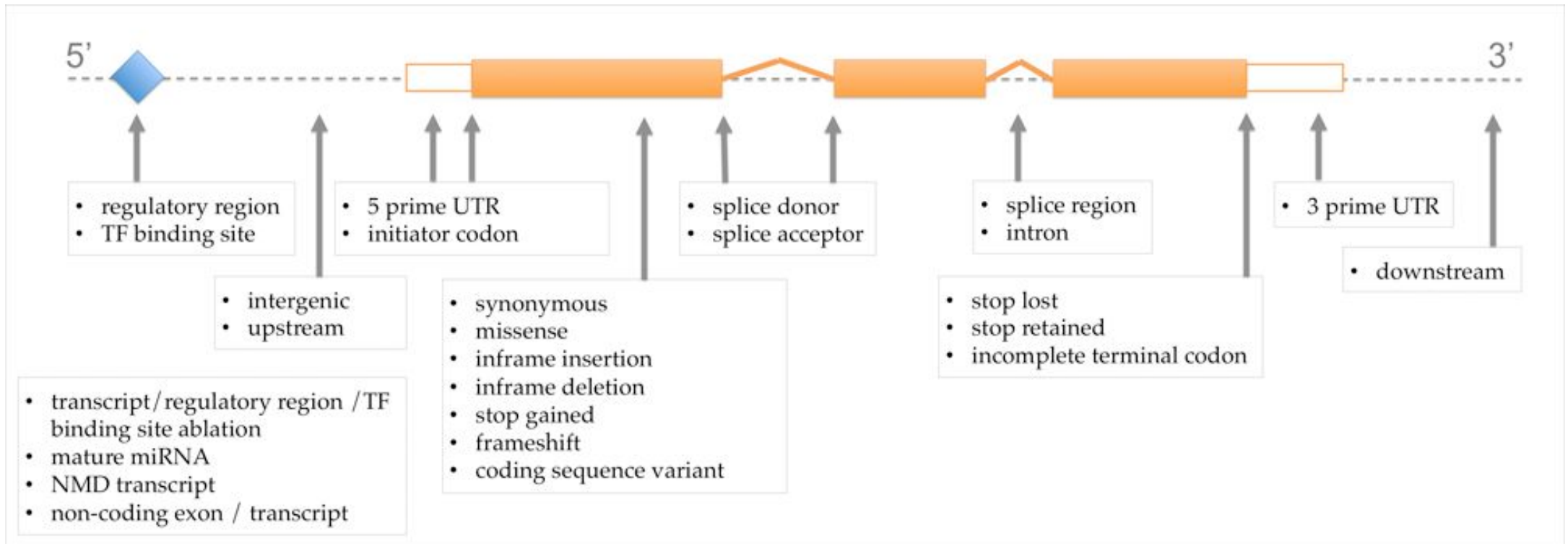
GATK's recommended workflow



<https://gatk.broadinstitute.org/hc/en-us/articles/360035535932-Germline-short-variant-discovery-SNPs-Indels->

Variant annotation

What does the SNP do?



Variant annotation

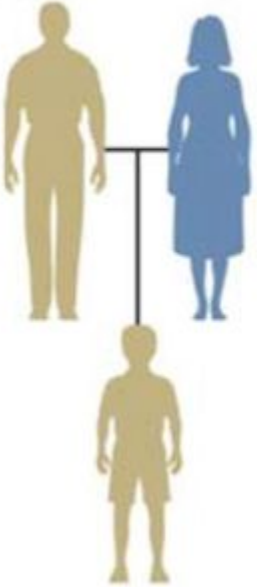
- Some example of tools:
 - Annovar
 - Ensembl Variant Effect Predictor (VEP)
 - SnpEff
- As good as annotations
- Beware of gene expression

we did not cover (in detail)...

Germline vs somatic

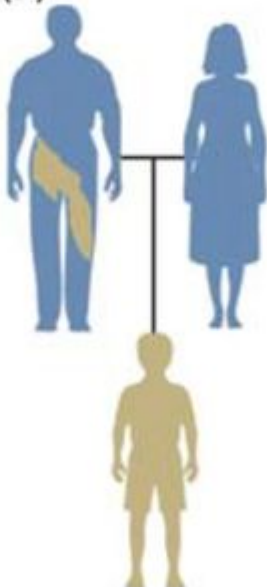
Inherited

(A)



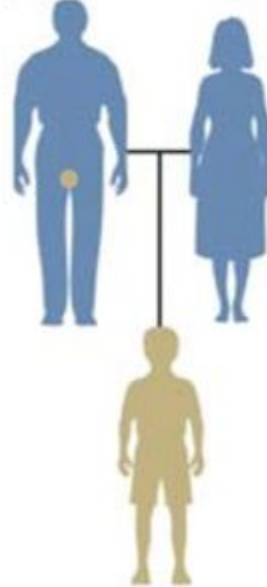
Father has mutation in all cells and transmits it on to his child. Child is heterozygous in every cell.

(B)



Father has mosaic mutation that affects germline and somatic cells. Child is heterozygous in every cell.

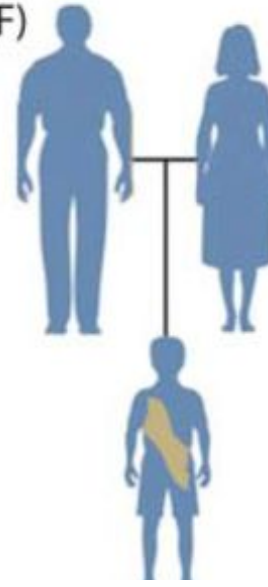
(C)



Father has germline mosaic mutation. Child is heterozygous in every cell.

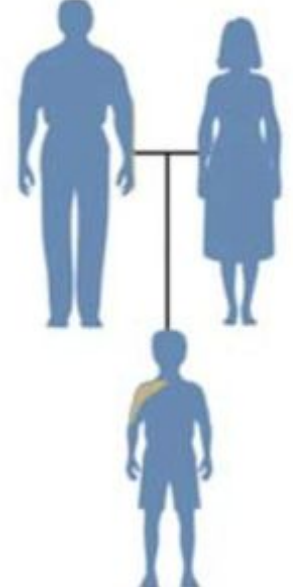
Somatic

(F)



Child has mosaic somatic mutation that occurs early in postzygotic development and is present in a percentage of his cells.

(G)

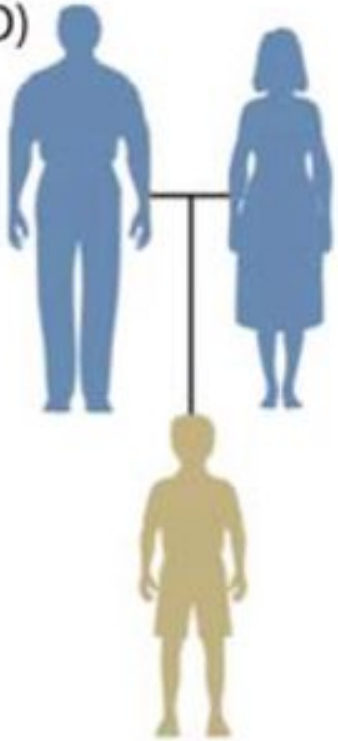


Child has mosaic mutation that occurs later in development and affects fewer cells (e.g. skin cells)

de novo

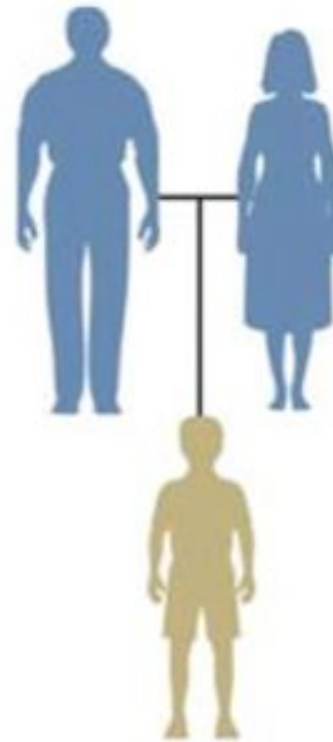
De novo

(D)



Father has mutation in a single sperm cell and transmits it to the child. Child is heterozygous in every cell.

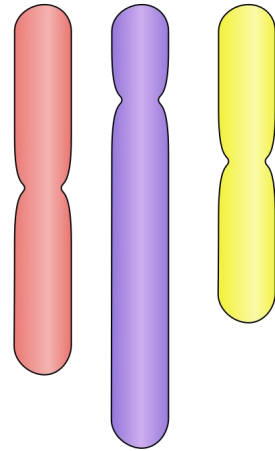
(E)



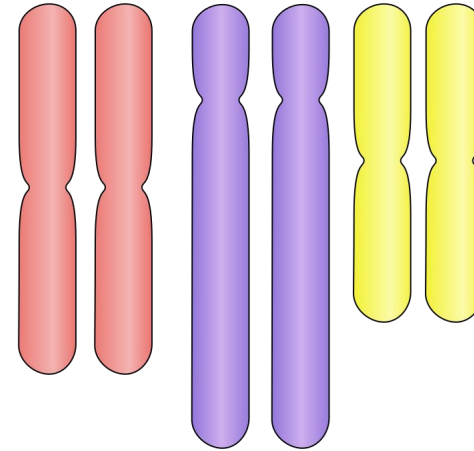
Mutation occurs in zygote within first few cell divisions. Child is heterozygous in every cell.

Polyploid

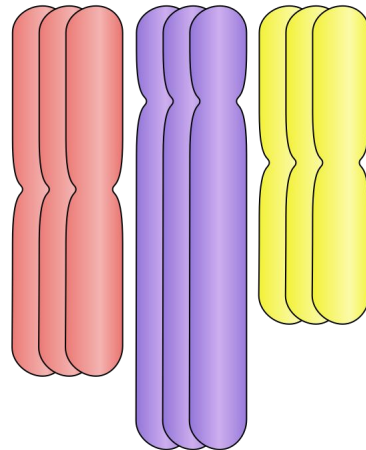
Haploid (N)



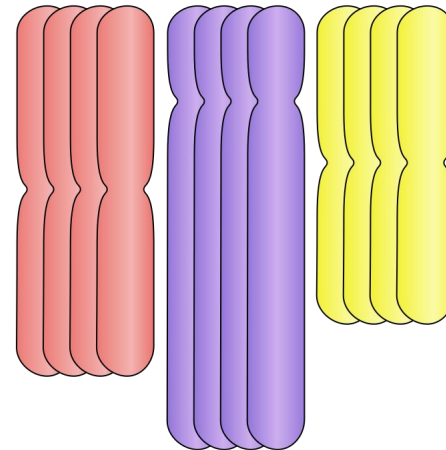
Diploid (2N)



Triploid (3N)

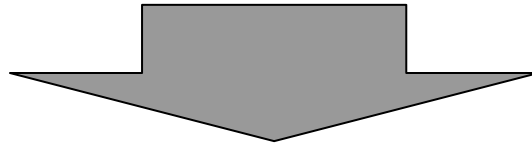


Tetraploid (4N)



Phasing

TAC^CGAAA^TCAT



TAC^CGAAA^TCAT

VS

TAC^CGAAA^CCAT

TAG^GGAAA^CCAT

TAG^GGAAA^TCAT

INDELS

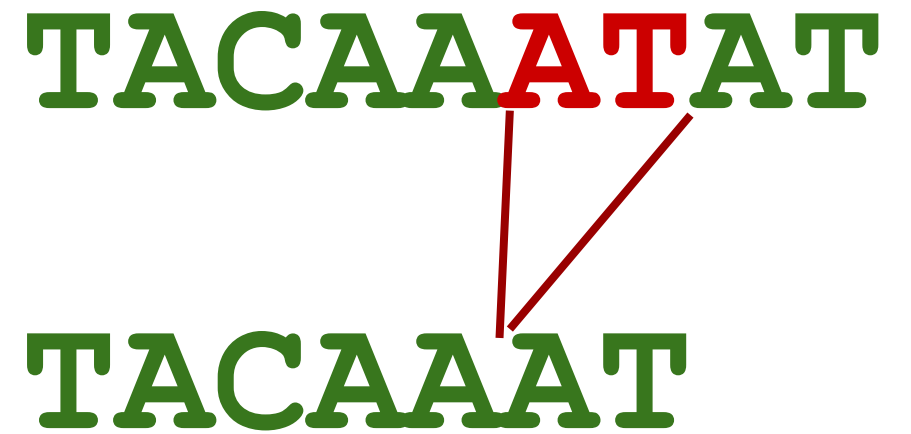
Insertions

TACAAAATAT
TACAAA**GC**TAT



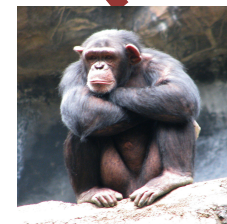
Deletion

TACAAA**AT**TAT
TACAAAAT



INDELS

Caution:



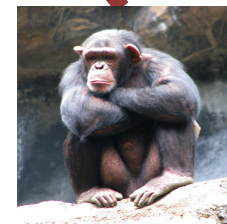
TACAAA--TAT

TACAAA**GC**TAT

GC was inserted

INDELS

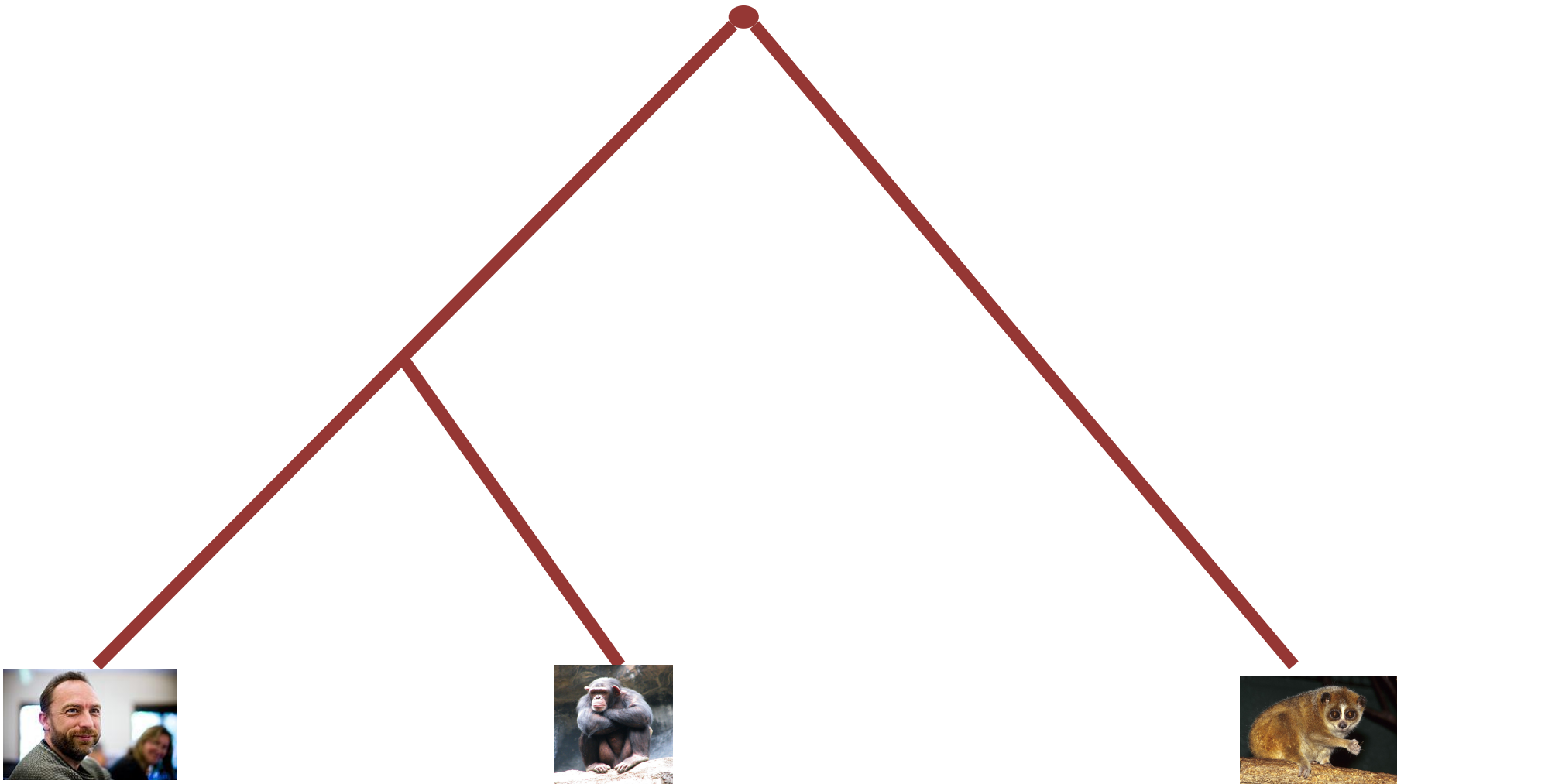
Caution:



TACAAA--TAT

TACAAA**G****C**TAT





TACAAA--TAT

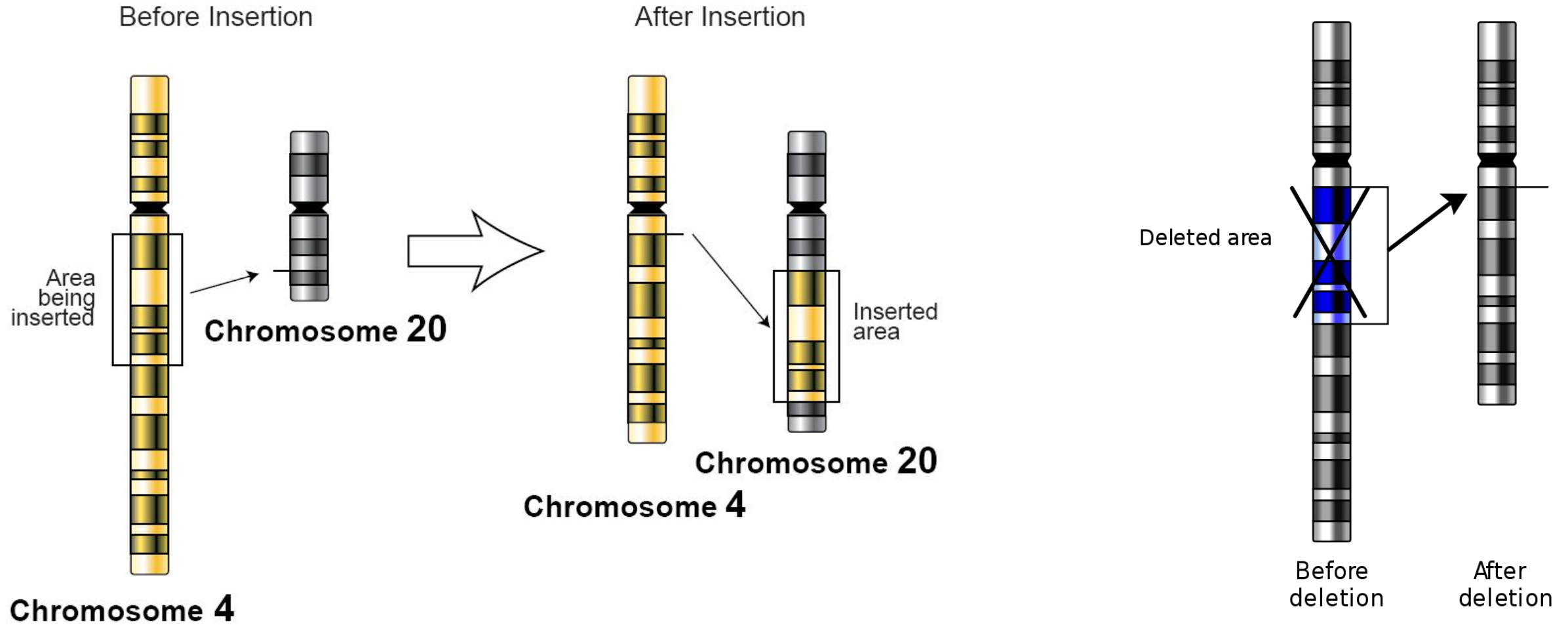
TACAAAGCTAT

TACAAAGCTAT

← GC was deleted

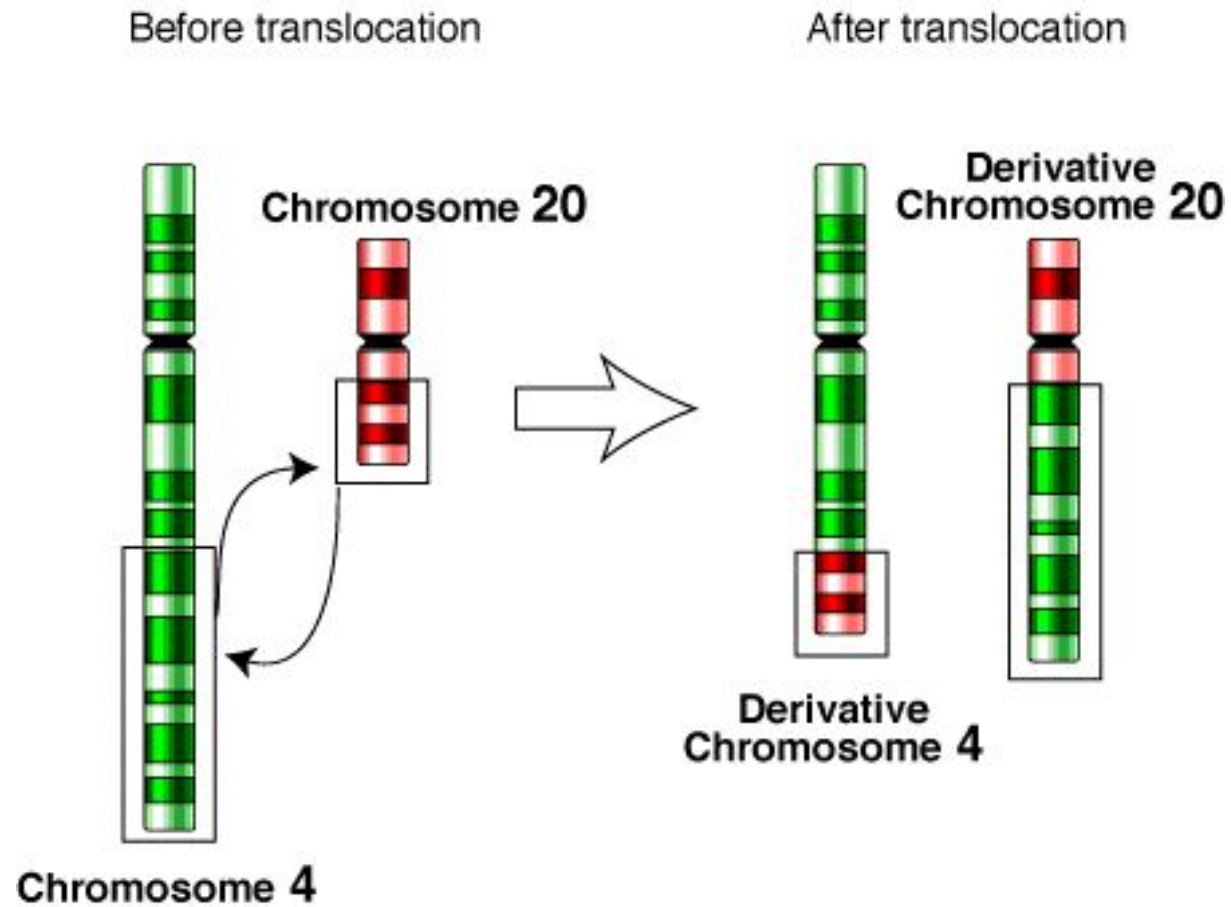
more likely, not guaranteed!

Structural variants



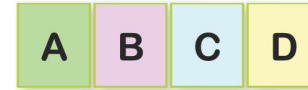
Structural variants

Translocation:



Structural variants

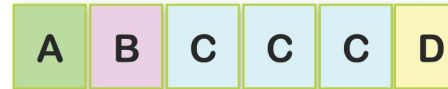
Copy number variations (CNV)



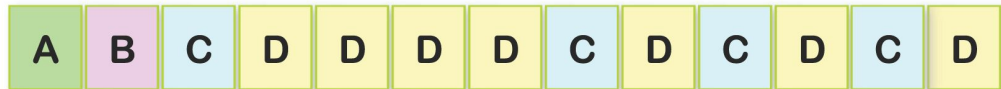
Reference



Segmental Duplication - Biallelic CNV (C)₂



Multiallelic Copy Number Variant (C)_{0-n}



Complex CNV (D)₄(CD)₃



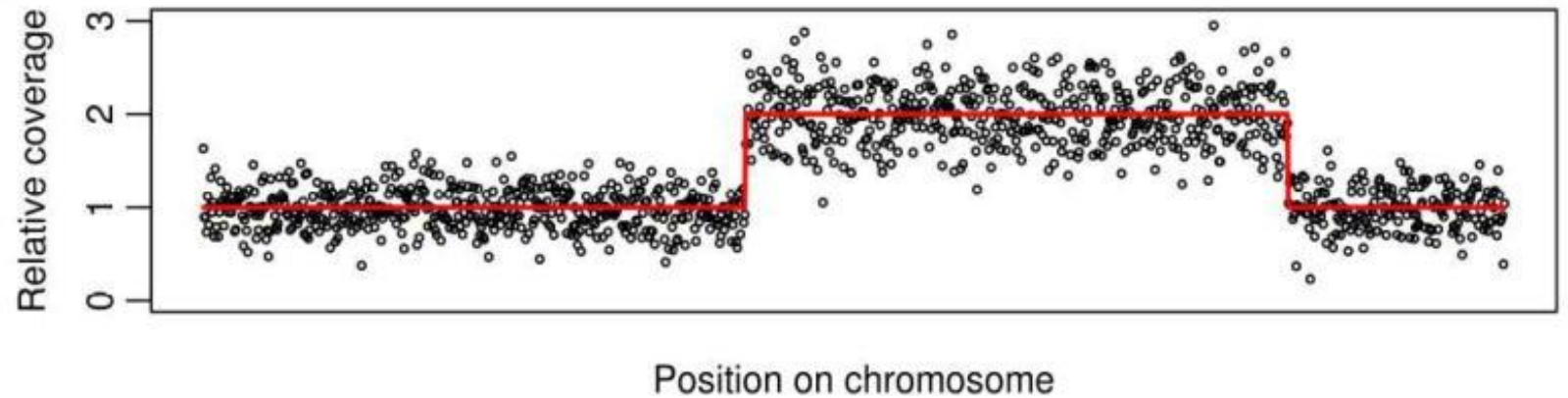
Inversion (CB)

Chromosome

Estivill, Xavier, and Lluís Armengol. "Copy Number Variants and Common Disorders: Filling the Gaps and Exploring Complexity in Genome-Wide Association Studies." *PLoS Genet* 3.10 (2007): e190.

Structural variants

Copy number variations (CNV)
effect on coverage



Weetman, David, Luc S. Djogbenou, and Eric Lucas. "Copy number variation (CNV) and insecticide resistance in mosquitoes: evolving knowledge or an evolving problem?." *Current Opinion in Insect Science* 27 (2018): 82-88.

Ethical concerns

privacy, justice, fairness etc..

Exercise time!

http://teaching.healthtech.dtu.dk/22126/index.php/Postprocess_exercise