

# Variant calling, phasing, and imputation

Nik Baya  
Oxford Statistical Genomics Summer School  
June 19 2023

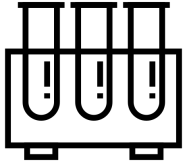


wellcome  
centre  
human  
genetics

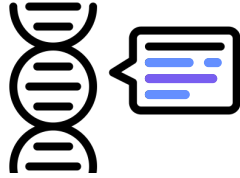


# Overview

Sequencing

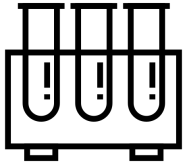


Calling

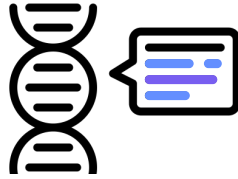


# Overview

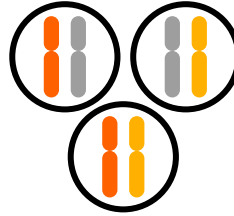
Sequencing



Calling

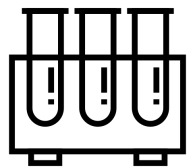


Phasing

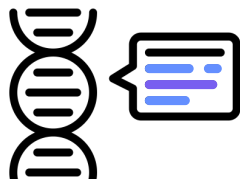


# Overview

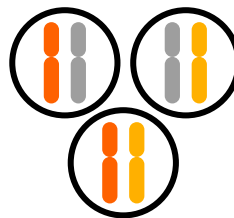
Sequencing



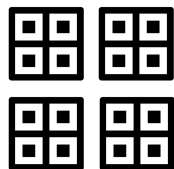
Calling



Phasing



Genotype array

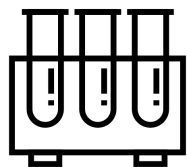


Phasing

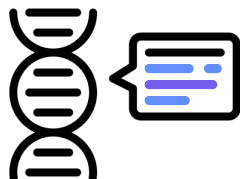


# Overview

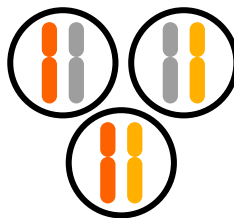
Sequencing



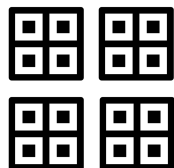
Calling



Phasing



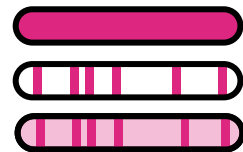
Genotype array



Phasing



Imputation

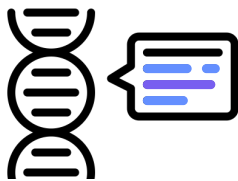


# Overview

Sequencing



Calling



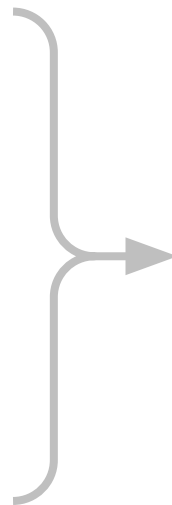
Phasing



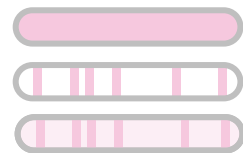
Genotype array



Phasing



Imputation



Reference sequence

C	A	C	T	C	G	G	G	G	C	G	T	C	C	G	T	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Aligned reads

A	G	T	C	A		
A	C	T	C	G		
G	G	G	G	C	G	
G	G	G	C	G	C	
G	G	G	C	G	C	
G	G	C	G	C	C	C
G	C	G	C	C	C	G
C	G	T	G	C		
C	G	T	G	C		

Reference sequence

C	A	C	T	C	G	G	G	G	C	G	T	C	C	G	T	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Aligned reads

A	G	T	C	A		
A	C	T	C	G		
G	G	G	G	C	G	
G	G	G	C	G	C	
G	G	G	C	G	C	
G	G	C	G	C	C	C
G	C	G	C	C	C	G
C	G	T	G	C		
C	G	T	G	C		

## Variant calling

*Is there evidence of a variant at this site in this dataset?*

# Variant calling

$$Pr(\text{variable site} | \text{Data})$$

*Probability that this site is  
variable, given the data*

# Variant calling

## Bayes' Theorem

$$Pr(A|B) = \frac{Pr(B|A)Pr(A)}{Pr(B)}$$

$Pr(\text{variable site} | \text{Data})$

*Probability that this site is variable, given the data*

## Bayes' Theorem

$$Pr(A|B) = \frac{Pr(B|A)Pr(A)}{Pr(B)}$$

# Variant calling

$$Pr(\text{variable site}|\text{Data}) = \frac{Pr(\text{Data}|\text{variable site})Pr(\text{variable site})}{Pr(\text{Data})}$$

*Probability that this site is variable, given the data*

# Variant calling

$$Pr(\text{variable site} | \text{Data}) = \frac{Pr(\text{Data} | \text{variable site}) Pr(\text{variable site})}{Pr(\text{Data})}$$

$$\text{Probability that this site is variable, given the data} = \frac{\text{prior} \prod_i L_x^i}{(1 - \text{prior}) \prod_i L_{ref}^i + \text{prior} \prod_i L_x^i}$$

# Variant calling

$$Pr(\text{variable site} | \text{Data}) = \frac{Pr(\text{Data} | \text{variable site}) Pr(\text{variable site})}{Pr(\text{Data})}$$

$$\text{Probability that this site is variable, given the data} = \frac{\text{prior} \prod_i L_x^i}{(1 - \text{prior}) \prod_i L_{ref}^i + \text{prior} \prod_i L_x^i}$$

“is variable”

“is not variable”

# Variant calling

$$Pr(\text{variable site} | \text{Data}) = \frac{Pr(\text{Data} | \text{variable site}) Pr(\text{variable site})}{Pr(\text{Data})}$$

$$\text{Probability that this site is variable, given the data} = \frac{\text{prior} \prod_i L_x^i}{(1 - \text{prior}) \prod_i L_{ref}^i + \text{prior} \prod_i L_x^i}$$

Prior probability that a site is variable  
(i.e. site is a variant)

# Prior probability of site being a variant

$$\text{prior}(N) = 4N_e\mu \sum_{k=1}^{2N-1} 1/k$$

$N$  : Size of sequenced population

$N_e$  : Effective population size

$\mu$  : Mutation rate

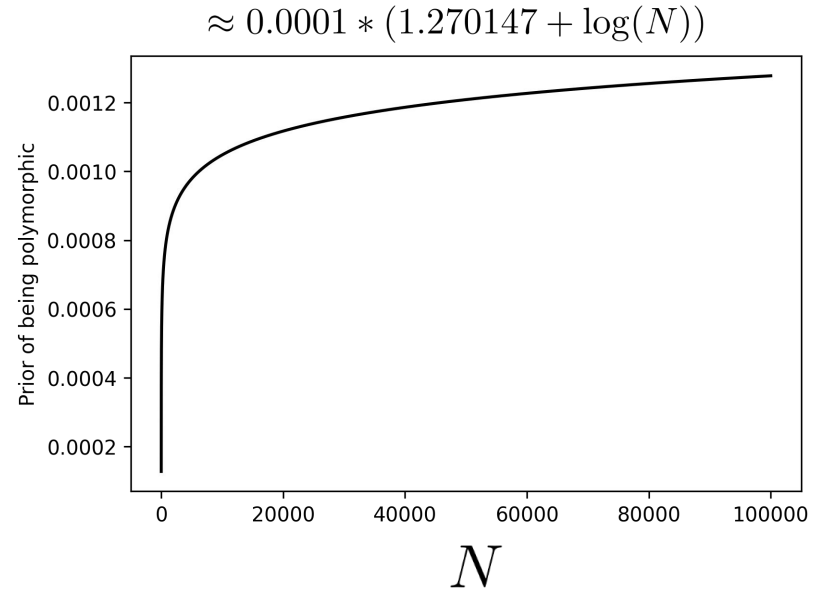
# Prior probability of site being a variant

$$\text{prior}(N) = 4N_e\mu \sum_{k=1}^{2N-1} 1/k$$

$N$  : Size of sequenced population

$N_e$  : Effective population size

$\mu$  : Mutation rate



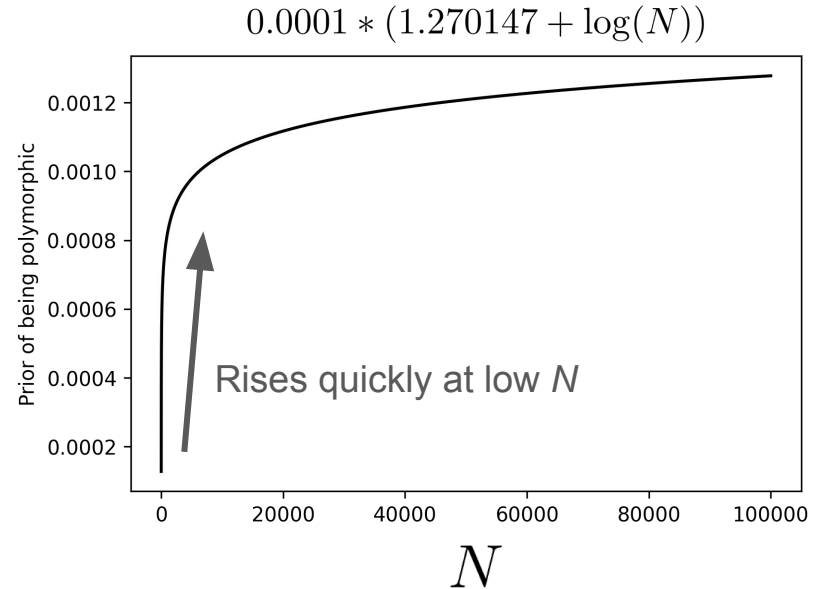
# Prior probability of site being a variant

$$prior(N) = 4N_e\mu \sum_{k=1}^{2N-1} 1/k$$

$N$  : Size of sequenced population

$N_e$  : Effective population size

$\mu$  : Mutation rate



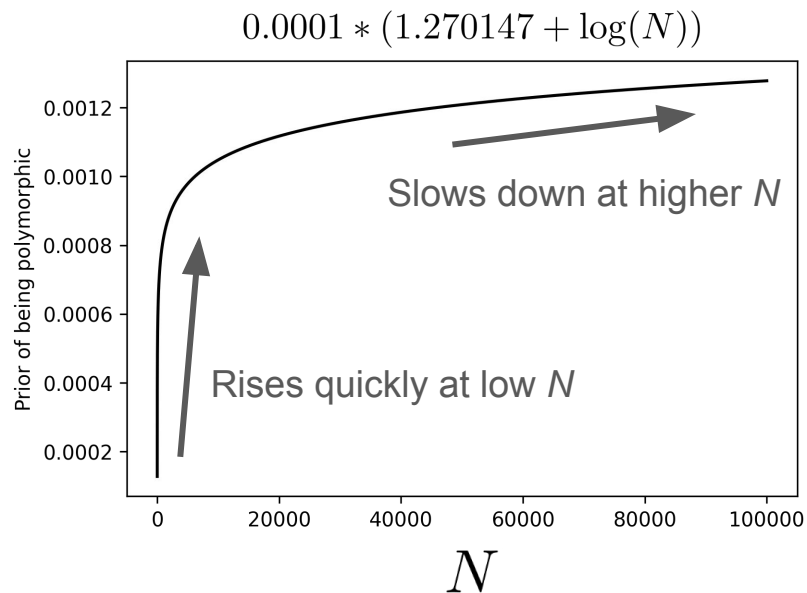
# Prior probability of site being a variant

$$\text{prior}(N) = 4N_e\mu \sum_{k=1}^{2N-1} 1/k$$

$N$  : Size of sequenced population

$N_e$  : Effective population size

$\mu$  : Mutation rate



# Variant calling

$$Pr(\text{variable site} | \text{Data}) = \frac{Pr(\text{Data} | \text{variable site}) Pr(\text{variable site})}{Pr(\text{Data})}$$

$$\text{Probability that this site is variable, given the data} = \frac{\text{prior} \prod_i L_x^i}{(1 - \text{prior}) \prod_i L_{ref}^i + \text{prior} \prod_i L_x^i}$$

# Variant calling

$$Pr(\text{variable site} | \text{Data}) = \frac{Pr(\text{Data} | \text{variable site}) Pr(\text{variable site})}{Pr(\text{Data})}$$

$$\text{Probability that this site is variable, given the data} = \frac{\text{prior} \prod_i L_x^i \quad \text{Likelihood being variable}}{(1 - \text{prior}) \prod_i L_{ref}^i + \text{prior} \prod_i L_x^i}$$

Likelihood of being invariant

# Allele likelihood

Likelihood of site being **variable**:

$$\mathcal{L}_x^i = 2 * [G_i(g = 1)f_x(1 - f_x)] + G_i(g = 2)f_x^2$$

Likelihood of being **invariant**:

$$\mathcal{L}_{ref}^i = G_i(g = 0)(1 - f_x)^2$$

# Allele likelihood

$f_x$  : Frequency of allele x

Likelihood of site being **variable**:

$$\mathcal{L}_x^i = 2 * [G_i(g = 1)f_x(1 - f_x)] + G_i(g = 2)f_x^2$$

Likelihood of being **invariant**:

$$\mathcal{L}_{ref}^i = G_i(g = 0)(1 - f_x)^2$$

# Allele likelihood

Likelihood of site being **variable**:

$$\mathcal{L}_x^i = 2 * \boxed{G_i(g = 1)} f_x (1 - f_x) + \boxed{G_i(g = 2)} f_x^2$$

Genotype likelihoods



Likelihood of being **invariant**:

$$\mathcal{L}_{ref}^i = \boxed{G_i(g = 0)} (1 - f_x)^2$$

# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m - g)\epsilon_j + g(1 - \epsilon_j) \right] \prod_{j=l+1}^k \left[ (m - g)(1 - \epsilon_j) + g\epsilon_j \right]$$

# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m-g)\epsilon_j + g(1-\epsilon_j) \right] \prod_{j=l+1}^k \left[ (m-g)(1-\epsilon_j) + g\epsilon_j \right]$$

E[ # of true reference reads ]      E[ # of true alternate reads ]

# Genotype likelihood

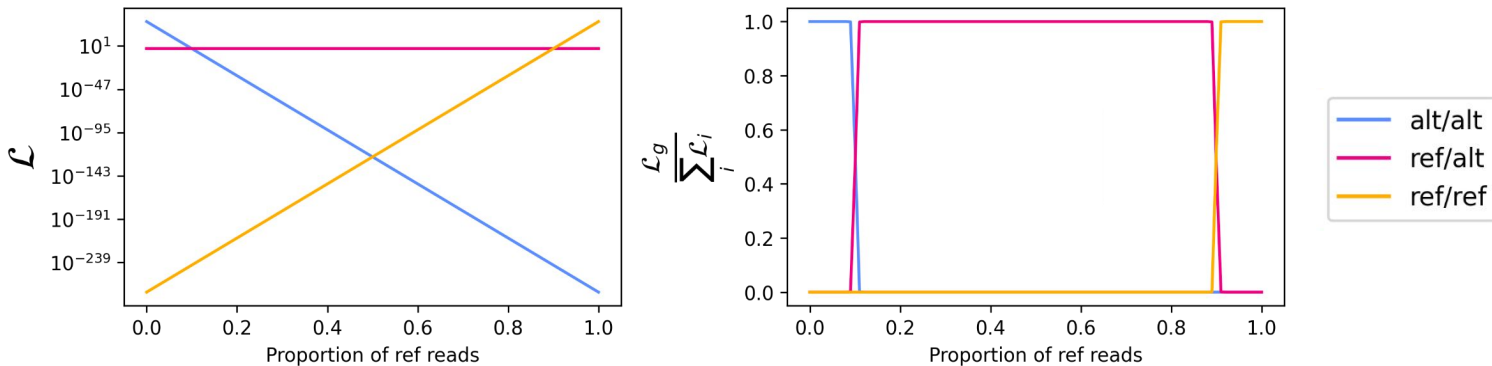
$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m-g)\epsilon_j + g(1-\epsilon_j) \right] \prod_{j=l+1}^k \left[ (m-g)(1-\epsilon_j) + g\epsilon_j \right]$$

Example: Varying proportion of reference reads

# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m-g)\epsilon_j + g(1-\epsilon_j) \right] \prod_{j=l+1}^k \left[ (m-g)(1-\epsilon_j) + g\epsilon_j \right]$$

Example: Varying proportion of reference reads



# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m-g)\epsilon_j + g(1-\epsilon_j) \right] \prod_{j=l+1}^k \left[ (m-g)(1-\epsilon_j) + g\epsilon_j \right]$$

Example: Increasing error rate

[ REF ]  
[ REF ]  
[ REF ]  
[ ALT ]

3 ref reads, 1 alt read

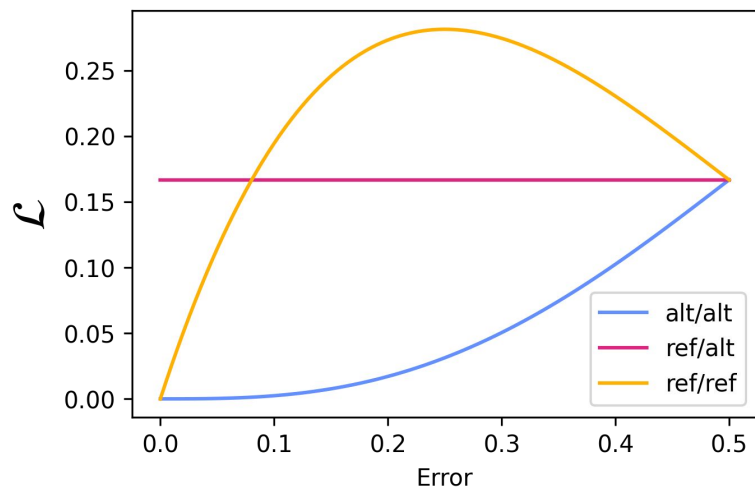
# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m - g)\epsilon_j + g(1 - \epsilon_j) \right] \prod_{j=l+1}^k \left[ (m - g)(1 - \epsilon_j) + g\epsilon_j \right]$$

Example: Increasing error rate

[ REF ]  
[ REF ]  
[ REF ]  
[ ALT ]

3 ref reads, 1 alt read



error rate constant across reads ( $\text{Var}(e_j)=0$ )

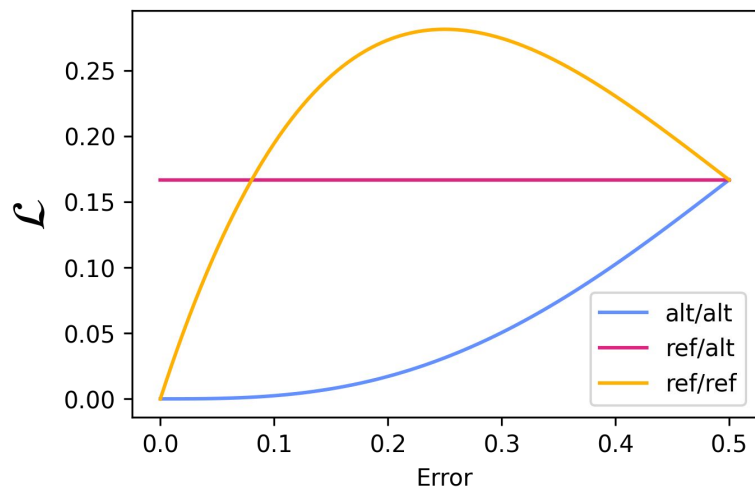
# Genotype likelihood

$$\mathcal{L}(g) = \frac{1}{m^k} \prod_{j=1}^l \left[ (m-g)\epsilon_j + g(1-\epsilon_j) \right] \prod_{j=l+1}^k \left[ (m-g)(1-\epsilon_j) + g\epsilon_j \right]$$

Example: Increasing error rate

[ REF ]  
[ REF ]  
[ REF ]  
[ ALT ]

3 ref reads, 1 alt read



error  $\rightarrow$  0.5  
(i.e. random)

all likelihoods converge

# Summary: Variant calling

# Summary: Variant calling

*Probability that this site is variable, given the data*



Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

# Summary: Variant calling

*Probability that this site is variable, given the data*



**QUAL score**

(or site quality score)

Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

#CHROM	POS	ID	REF	ALT	QUAL
chr19	48693829	.	C	T	3357.45
chr19	48693956	.	A	G	6806.98
chr19	48693992	.	G	A	1802.53
chr19	48694122	.	C	G	125.986

# Summary: Variant calling

*Probability that this site is variable, given the data*



**QUAL score**

(or site quality score)

Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

Usually expressed on the **phred scale**

#CHROM	POS	ID	REF	ALT	QUAL
chr19	48693829	.	C	T	3357.45
chr19	48693956	.	A	G	6806.98
chr19	48693992	.	G	A	1802.53
chr19	48694122	.	C	G	125.986

# Summary: Variant calling

*Probability that this site is variable, given the data*



**QUAL score**

(or site quality score)

Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

Usually expressed on the **phred scale**

$$Q = -10 \cdot \log_{10}(\text{Probability of error})$$

#CHROM	POS	ID	REF	ALT	QUAL
chr19	48693829	.	C	T	3357.45
chr19	48693956	.	A	G	6806.98
chr19	48693992	.	G	A	1802.53
chr19	48694122	.	C	G	125.986

# Summary: Variant calling

*Probability that this site is variable, given the data*



**QUAL score**

(or site quality score)

#CHROM	POS	ID	REF	ALT	QUAL
chr19	48693829	.	C	T	3357.45
chr19	48693956	.	A	G	6806.98
chr19	48693992	.	G	A	1802.53
chr19	48694122	.	C	G	125.986

Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

Usually expressed on the **phred scale**

$$Q = -10 \cdot \log_{10}(\text{Probability of error})$$

Error = 1

Q = 0

Error = 1/100

Q = 20

Error = 1/1000

Q = 30

# Summary: Variant calling

*Probability that this site is variable, given the data*



**QUAL score**

(or site quality score)

Depends on:

- Prior probability of variant
- Allele likelihood
  - Allele frequency
  - Genotype likelihoods

Usually expressed on the **phred scale**

$$Q = -10 \cdot \log_{10}(\text{Probability of error})$$

higher score = higher quality

#CHROM	POS	ID	REF	ALT	QUAL
chr19	48693829	.	C	T	3357.45
chr19	48693956	.	A	G	6806.98
chr19	48693992	.	G	A	1802.53
chr19	48694122	.	C	G	125.986

Reference sequence

C	A	C	T	C	G	G	G	G	C	G	T	C	C	G	T	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Aligned reads

A	G	T	C	A		
A	C	T	C	G		
G	G	G	G	C	G	
G	G	G	C	G	C	
G	G	G	C	G	C	
G	G	C	G	C	C	C
G	C	G	C	C	C	G
C	G	T	G	C		
C	G	T	G	C		

## Variant calling

*Is there evidence of a variant at this site in this dataset?*

# Genotype calling

*What are the genotypes for each individual?*

G/A?

Reference sequence

C A C T C G G G G C G T C C G T G C

Aligned reads

A	G	T	C	A																
A	C	T	C	G																
				G	G	G	G	C	G											
					G	G	G	C	G	C										
					G	G	G	C	G	C										
						G	G	C	G	C	C	C								
							G	C	G	C	C	C	G							
														C	G	T	G	C		
														C	G	T	G	C		

## Variant calling

*Is there evidence of a variant at this site in this dataset?*

# Genotype calling

# Genotype calling

$$\Pr(g|Data_i) = \frac{\Pr(Data_i|g)\Pr(g|f_x)}{\Pr(Data_i)}$$

*Probability of genotype  $g$ ,  
given data for sample  $i$*

# Genotype calling

$$\Pr(g|Data_i) = \frac{\Pr(Data_i|g)\Pr(g|f_x)}{\Pr(Data_i)}$$

*Probability of genotype  $g$ ,  
given data for sample  $i$*

$$= \frac{G_i(g)\Pr(g|f_x)}{\sum_{g'} G_i(g')\Pr(g'|f_x)}$$

# Genotype calling

Genotype posterior (**GP**)

$$Pr(g|Data_i) = \frac{Pr(Data_i|g)Pr(g|f_x)}{Pr(Data_i)}$$

Probability of genotype  $g$ ,  
given data for sample  $i$

$$= \frac{G_i(g)Pr(g|f_x)}{\sum_{g'} G_i(g')Pr(g'|f_x)}$$

Genotype priors

Genotype likelihoods  
(**GL**, or **PL** if phred-scaled)

# Genotype calling

Genotype posterior (GP)

$$Pr(g|Data_i) = \frac{Pr(Data_i|g)Pr(g|f_x)}{Pr(Data_i)}$$

Probability of genotype  $g$ ,  
given data for sample  $i$

$$= \frac{G_i(g)Pr(g|f_x)}{\sum_{g'} G_i(g')Pr(g'|f_x)}$$

Genotype priors

Genotype likelihoods  
(GL, or PL if phred-scaled)

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	FORMAT	HG02461	HG02462
chr19	48693829	.	C	T	3357.45	.	GT:PL:GP	0/0:0,9,75:0.810234,0.189766,2.21701e-08	1/1:85,9,0:2.87613e-09,0.213018,0.786982
chr19	48693956	.	A	G	6806.98	.	GT:PL:GP	0/0:0,99,255:1,3.1458e-11,4.93632e-28	0/0:0,105,255:1,7.9019e-12,4.93632e-28
chr19	48693992	.	G	A	1802.53	.	GT:PL:GP	0/0:0,96,255:1,1.43086e-11,2.56527e-29	0/0:0,105,255:1,1.80134e-12,2.56527e-29

# Genotype calling

Choose genotype with **highest probability**

Genotype posterior (**GP**)

$$Pr(g|Data_i) = \frac{Pr(Data_i|g)Pr(g|f_x)}{Pr(Data_i)}$$

Probability of genotype  $g$ ,  
given data for sample  $i$

$$= \frac{G_i(g)Pr(g|f_x)}{\sum_{g'} G_i(g')Pr(g'|f_x)}$$

Genotype priors

Genotype likelihoods  
(**GL**, or **PL** if phred-scaled)

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	FORMAT	HG02461	HG02462
chr19	48693829	.	C	T	3357.45	.	GT:PL:GP	0/0:0,9,75:0.810234,0.189766,2.21701e-08	1/1:85,9,0:2.87613e-09,0.213018,0.786982
chr19	48693956	.	A	G	6806.98	.	GT:PL:GP	0/0:0,99,255:1,3.1458e-11,4.93632e-28	0/0:0,105,255:1,7.9019e-12,4.93632e-28
chr19	48693992	.	G	A	1802.53	.	GT:PL:GP	0/0:0,96,255:1,1.43086e-11,2.56527e-29	0/0:0,105,255:1,1.80134e-12,2.56527e-29

## A few notes:

- Types of variants
- Single-sample vs. joint calling
- Types of errors

# Variant types

# Variant types

Single-nucleotide polymorphism (SNP)

Insertion-deletion (indel)

Structural variant (SV)

Copy number variant (CNV)

Reference

Observed

A	A	G	C	T
A	A	T	C	T

# Variant types

Single-nucleotide polymorphism (SNP)

Insertion-deletion (indel)

Structural variant (SV)

Copy number variant (CNV)

Insertion

A	A	G	C	T
A	A	TCGA	C	T

Deletion

A	A	GATT	C	T
A	A	T	C	T

# Variant types

Single-nucleotide polymorphism (SNP)

Insertion-deletion (indel)

Structural variant (SV)

Copy number variant (CNV)

A	A	G	C	T
A	A	TC...CC	C	T



>50

base pairs

# Variant types

Single-nucleotide polymorphism (SNP)

Insertion-deletion (indel)

Structural variant (SV)

Copy number variant (CNV)

Reference



Gain



Loss

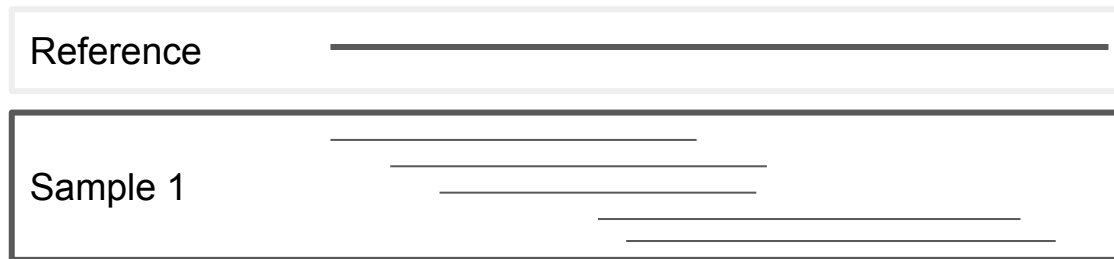


# Calling types

# Calling types

Single-sample calling

Joint calling



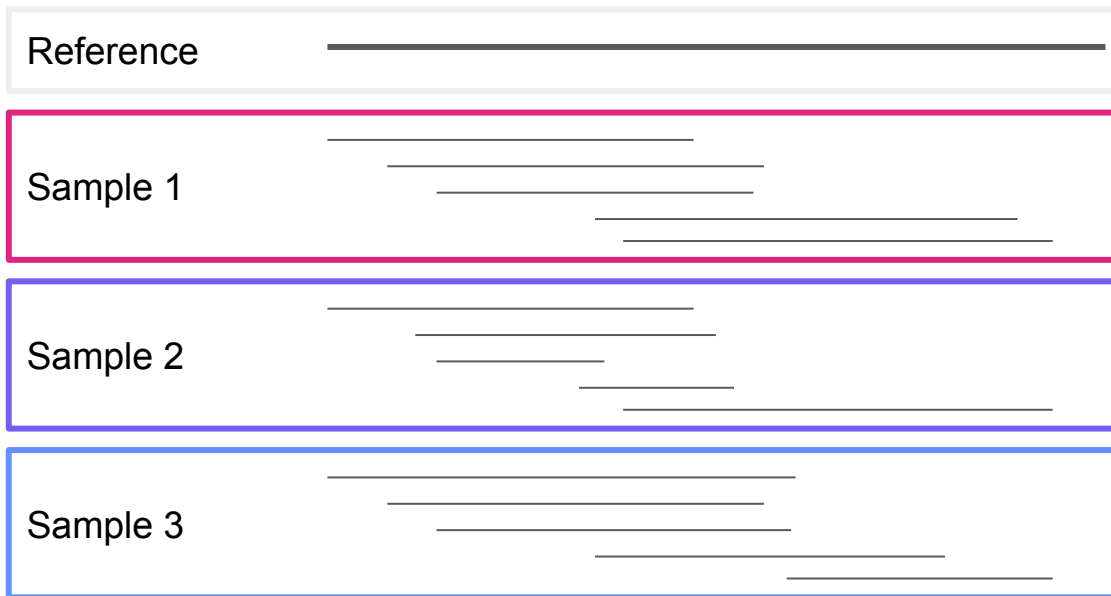
# Calling types

**Pros:** Better at identifying rare variants

**Cons:** More computationally expensive

Single-sample calling

Joint calling



# Error types

# Error types

Almost all variant callers rely on probabilistic models

(i.e. truly at random)

These models make assumptions – the biggest of which being that errors are **independent**

# Error types

Almost all variant callers rely on probabilistic models

(i.e. truly at random)

These models make assumptions – the biggest of which being that errors are **independent**

Example: Independent errors

If there's a 0.01 chance of seeing G when truth is a C, the chance of seeing it 5 times is  $0.01^5$

# Error types

Almost all variant callers rely on probabilistic models

(i.e. truly at random)

These models make assumptions – the biggest of which being that errors are **independent**

Example: Independent errors

If there's a 0.01 chance of seeing G when truth is a C, the chance of seeing it 5 times is  $0.01^5$

**But in practice, systematic (non-random) errors exist!**

# Example: Context-specific errors

Truth    AGGGCAA

Sequencing    AGGGCAA

AGGG**G**AA

AGGGCAA

AGGG**G**AA

AGGG**G**AA

“GGG” motif

# Example: Context-specific errors

Truth AGGGCAA

Sequencing AGGGCAA  
AGGG**G**AA  
AGGGCAA  
AGGG**G**AA  
AGGG**G**AA

“GGG” motif

Solution: Strand bias test

Compare enrichment of alternate alleles in forward vs. reverse strands.

# Example: Context-specific errors

Truth AGGGCAA

Sequencing AGGGCAA  
AGGG**G**AA  
AGGGCAA  
AGGG**G**AA  
AGGG**G**AA

“GGG” motif

Solution: Strand bias test

Compare enrichment of alternate alleles in forward vs. reverse strands.

Solution: Read position bias test

Correct for bias of sequencing errors typically occurring at the ends of reads.

# Example: Systematic mapping errors

Reference TCCACATGAACCACGTCTTGAAATC

Truth TCCACATGAACCACGT---GAAATC

Alignment

ACCACGT---G

CCACGT---GA

CACGT---GAA

CGT---GAAAT

This is the **true mapping**  
(CTT deletion)



# Example: Systematic mapping errors

Reference TCCACATGAACCACGTCTTCAAATC

Truth TCCACATGAACCACGT---GAAATC

Alignment ACCACGT---G

CCACGTGA

CACGTGAA

CGT---GAAAT

Solution:

Mapping quality bias test

Test whether non-reference reads have lower quality.

However, the alignment algorithm may lead to **incorrect mapping**.

This looks like a new SNP but the SNP doesn't exist!

# Example: Systematic mapping errors

Reference TCCACATGAACCACGTCTTCAAATC

Truth TCCACATGAACCACGT---GAAATC

Alignment ACCACGT---G

CCACGTGA

CACGTGAA

CGT---GAAAT

However, the alignment algorithm may lead to **incorrect mapping**.

This looks like a new SNP but the SNP doesn't exist!

## Solution:

Mapping quality bias test

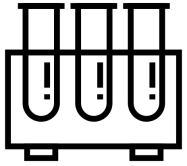
Test whether non-reference reads have lower quality.

## Solution: mapQ0

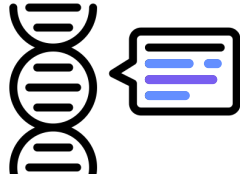
Remove variants with high counts of Q=0 reads

# Overview

Sequencing

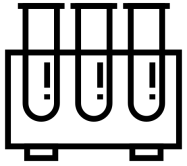


Calling

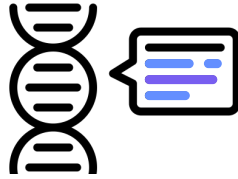


# Overview

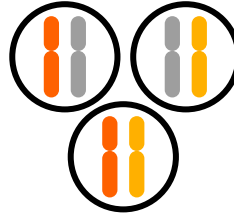
Sequencing



Calling



Phasing

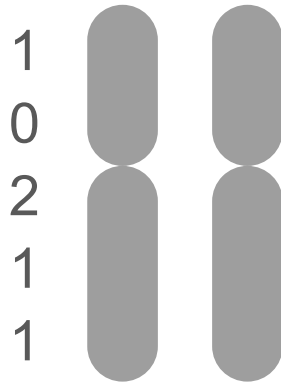


A/G  
T/T  
G/G  
A/T  
T/G

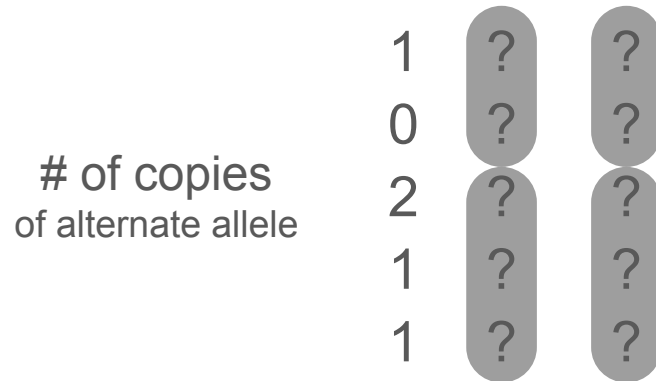


The diagram shows two vertical, rounded rectangular chromosomes. Each chromosome has five distinct horizontal bands. The bands on the left chromosome are labeled from top to bottom as A/G, T/T, G/G, A/T, and T/G. The right chromosome has an identical pattern of bands, representing a pair of homologous chromosomes.

# of copies  
of alternate allele



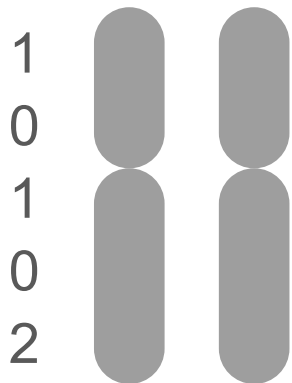
How to separate genotype into haplotypes?



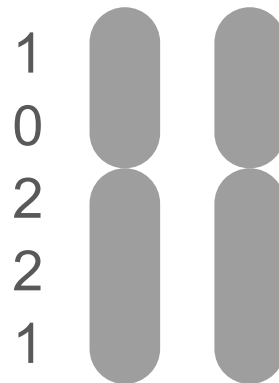
How to separate genotype into haplotypes?

	1	?	?
	0	0	0
# of copies of alternate allele	2	1	1
	1	?	?
	1	?	?

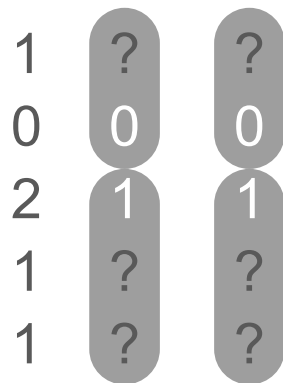
Mother



Father

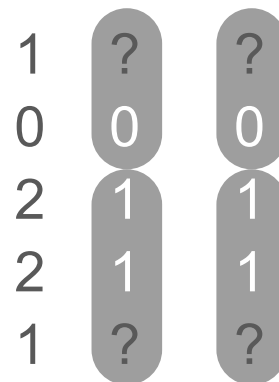
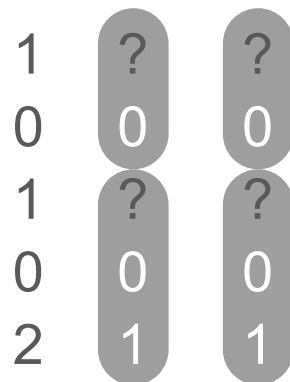


# of copies  
of alternate allele



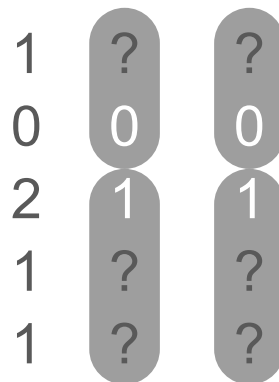
Child

Mother



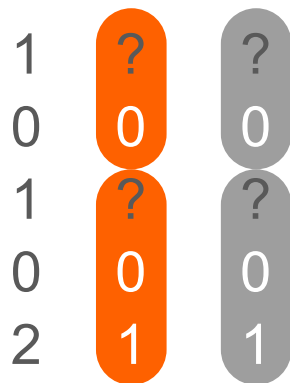
Father

# of copies  
of alternate allele

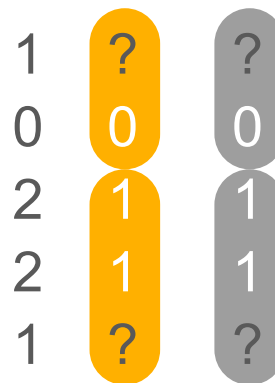


Child

Mother



Father

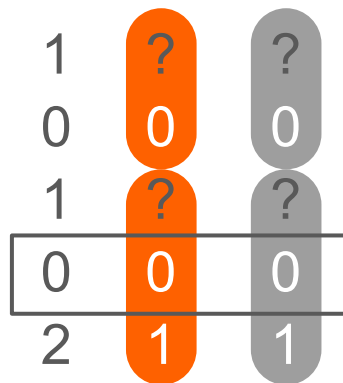


# of copies  
of alternate allele

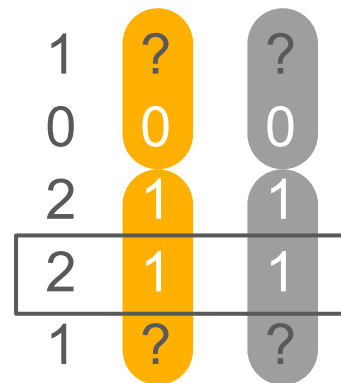


Child

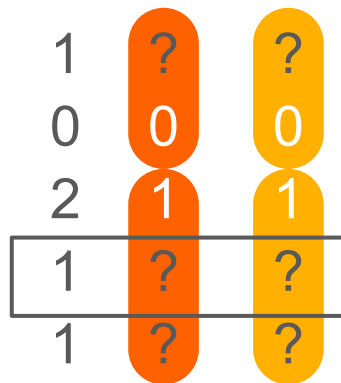
Mother



Father

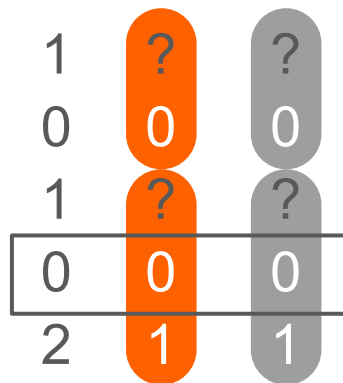


# of copies  
of alternate allele

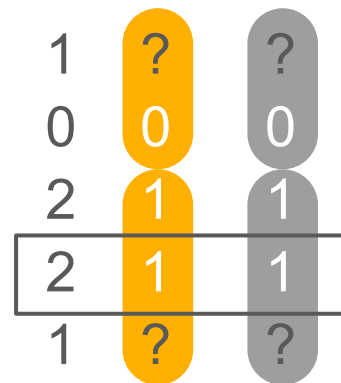


Child

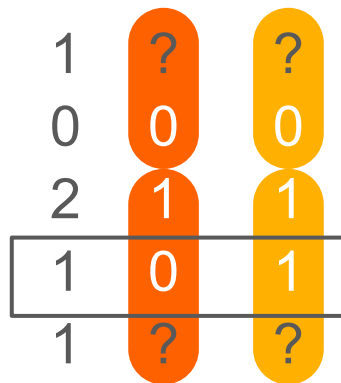
Mother



Father

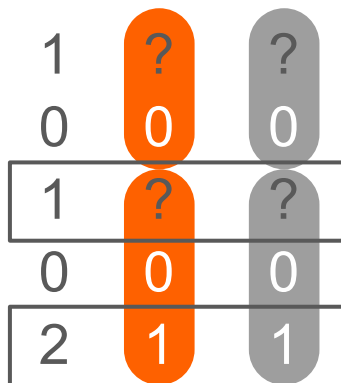


# of copies  
of alternate allele

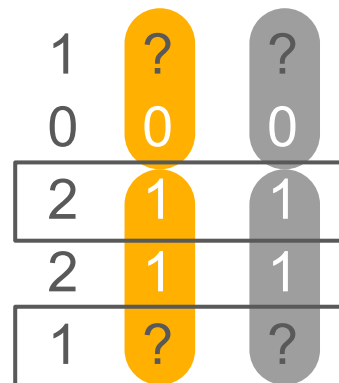


Child

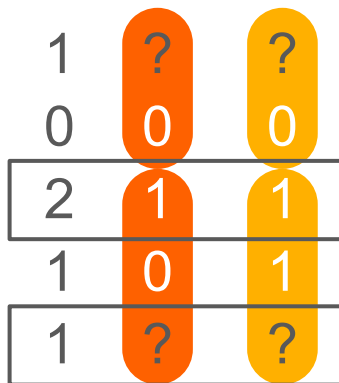
Mother



Father

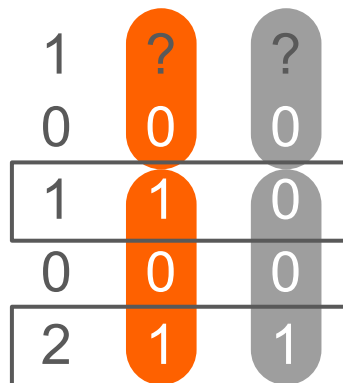


# of copies  
of alternate allele

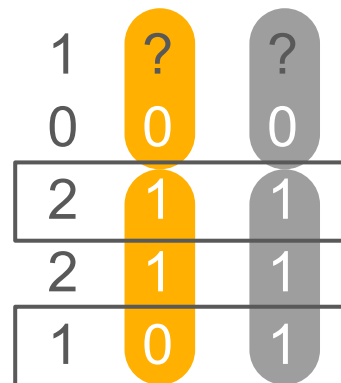


Child

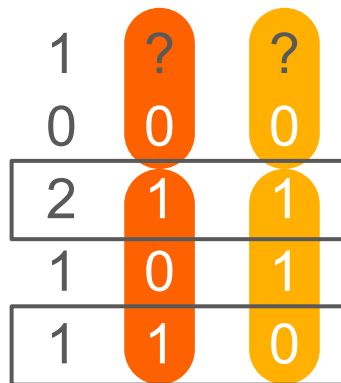
Mother



Father

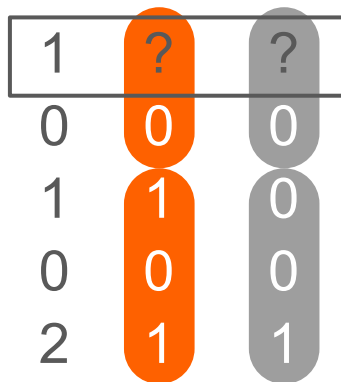


# of copies  
of alternate allele

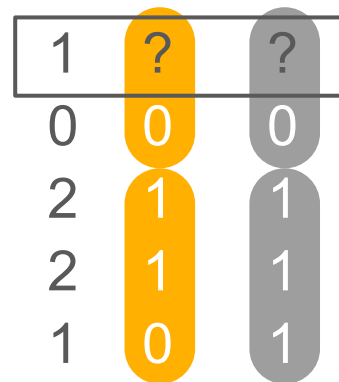


Child

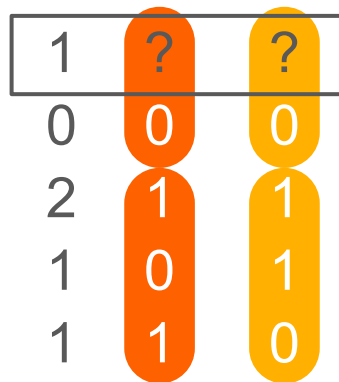
Mother



Father



# of copies  
of alternate allele



Child

Sample 1



Sample 2



Sample 3



Sample 4



Inferred  
ancestral  
haplotypes



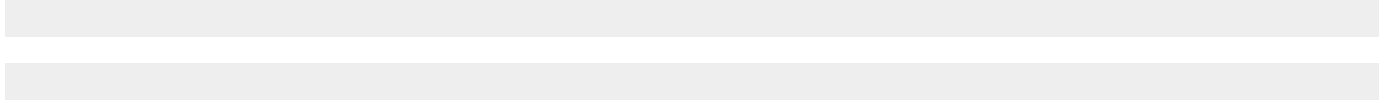
Sample 1



Sample 2



Sample 3



Sample 4



Inferred  
ancestral  
haplotypes



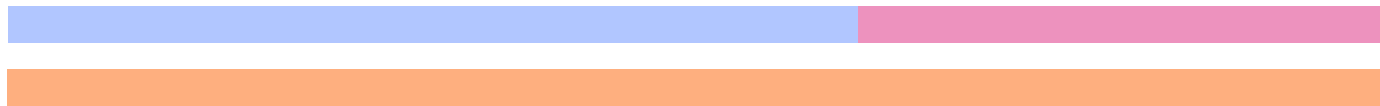
Sample 1



Sample 2



Sample 3

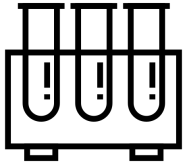


Sample 4

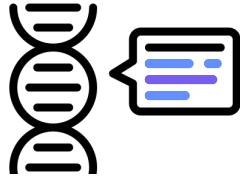


# Overview

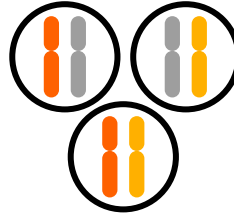
Sequencing



Calling

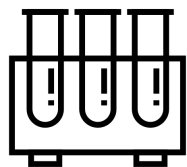


Phasing

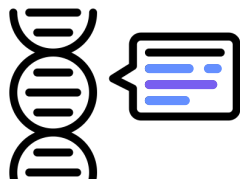


# Overview

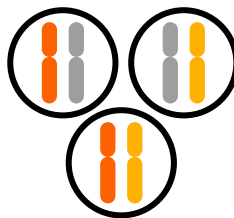
Sequencing



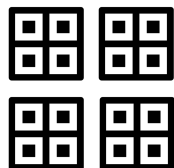
Calling



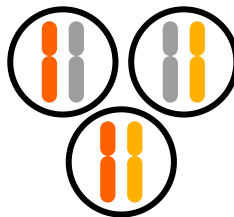
Phasing



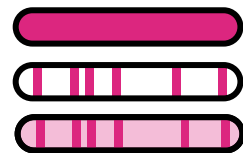
Genotype array



Phasing



Imputation



ATGGGCAAGCTAGAATTCGGCGCCAAGTAGGATCCGTATTCTATACCCGGTATGGGCAA

ATGGGCAAGCTAGAATTCGGCTCCAACACTACGATCCGAATTCTAGACCCGGTATGGGTAA

ATGGGCAAGCTAGAGTTCGGCGCCAAGTAGGATCCGAATTCTATACCCGGTATGGGCAA

ACGGGCAAGCTAGAGTTCGGCGCCAAGTAGGATCCGAATTCTAGACCCGGTATGGGCAA

Reference panel  
haplotypes  
(sequenced)

Reference panel haplotypes (sequenced)

ATGGGCAAGCTAGAATTCGGCGCCAAGCTAGGATCCGTATTCTATACCCGGTATGGGCAA

ATGGGCAAGCTAGAATTCGGCTCCAAGCTACGATCCGAATTCTAGACCCGGTATGGGTAA

ATGGGCAAGCTAGAGTTTCGGCGCCAAGCTATGATCCGAATTCTATACCCGGTATGGGCAA

ACGGGCAAGCTAGAGTTTCGGCGCCAAGCTAGGATCCGAATTCTAGACCCGGTATGGGCAA

Genotype array (phased)

Sample 1

T	A	G	G	T	T	C
C	G	G	G	A	G	C

Sample 2

T	G	G	T	A	T	C
T	A	T	C	A	G	T

Sample 3

T	G	G	T	A	T	C
C	G	G	G	A	G	C

Sample 4

T	A	T	C	A	G	T
C	G	G	G	A	G	C

Reference panel haplotypes (sequenced)

ATGGGCAAGCTAGAAATTCGGCGCCAAGTAGGATCCGTATTCTATACCCGGTATGGGCAA

ATGGGCAAGCTAGAAATTCGGCTCCAAGTACGATCCGAATTCTAGACCCGGTATGGGTAA

ATGGGCAAGCTAGAGTTTCGGCGCCAAGTATGATCCGAATTCTATACCCGGTATGGGCAA

ACGGGCAAGCTAGAGTTTCGGCGCCAAGTAGGATCCGAATTCTAGACCCGGTATGGGCAA

Genotype array (phased)

Sample 1	T	A	G	G	T	T	C
	C	G	G	G	A	G	C
Sample 2	T	G	G	T	A	T	C
	T	A	T	C	A	G	T
Sample 3	T	G	G	T	A	T	C
	C	G	G	G	A	G	C
Sample 4	T	A	T	C	A	G	T
	C	G	G	G	A	G	C

Reference panel  
haplotypes  
(sequenced)

ATGGGCAAGCTAGAATTCGGCGCCAAC TAGGATCCGTATTCTATACCCGGTATGGGCAA  
ATGGGCAAGCTAGAATTCGGCTCCAAC TACGATCCGAATTCTAGACCCGGTATGGGTAA  
ATGGGCAAGCTAGAGTTCGGCGCCAAC TAGATCCGAATTCTATACCCGGTATGGGCAA  
ACGGGCAAGCTAGAGTTCGGCGCCAAC TAGGATCCGAATTCTAGACCCGGTATGGGCAA

Genotype array  
(phased)

Sample 1

ATGGGCAAGCTAGAATTCGGCGCCAAC TAGGATCCGTATTCTATACCCGGTATGGGCAA  
ACGGGCAAGCTAGAGTTCGGCGCCAAC TAGGATCCGAATTCTAGACCCGGTATGGGCAA

Sample 2

ATGGGCAAGCTAGAGTTCGGCGCCAAC TAGATCCGAATTCTATACCCGGTATGGGCAA  
ATGGGCAAGCTAGAATTCGGCTCCAAC TACGATCCGAATTCTAGACCCGGTATGGGTAA

Sample 3

ATGGGCAAGCTAGAGTTCGGCGCCAAC TAGATCCGAATTCTATACCCGGTATGGGCAA  
ACGGGCAAGCTAGAGTTCGGCGCCAAC TAGGATCCGAATTCTAGACCCGGTATGGGCAA

Sample 4

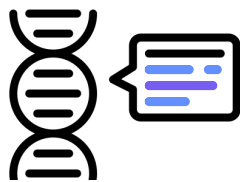
ATGGGCAAGCTAGAATTCGGCTCCAAC TACGATCCGAATTCTAGACCCGGTATGGGTAA  
ACGGGCAAGCTAGAGTTCGGCGCCAAC TAGGATCCGAATTCTAGACCCGGTATGGGCAA

# Practical

Sequencing



Calling



Phasing



Genotype array



Phasing



Imputation

