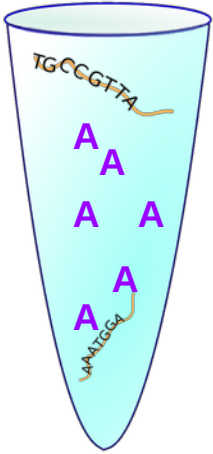# Genotype likelihoods, allele frequencies, and SNP calling from NGS data
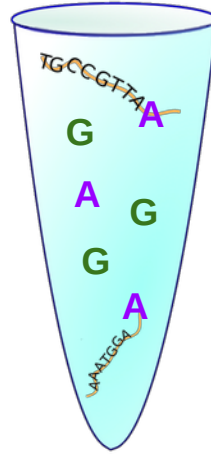
Tyler Linderoth
Physalia lcWGS course 2025

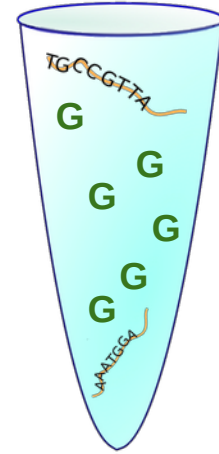# Why Probabilistic Methods?

# Why Probabilistic Methods?

The library for an individual homozygous for the **A** allele will consist only of **A**s.
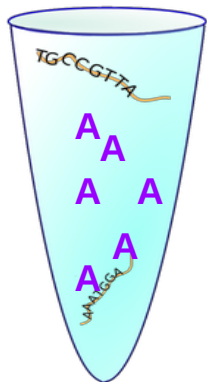
The library for a heterozygous individual at a site contains both **A**s and **G**s.

The library for an individual homozygous for the **G** allele will consist only of **G**s.

# Why Probabilistic Methods?

Sequence to average depth of 4x.

$$\text{Depth} \sim \text{Poisson}(\lambda = 4)$$
$$\text{E}[\text{depth}] = \lambda$$
$$\text{Var}[\text{depth}] = \lambda$$

Expect 4 reads, all with As at this ref position.

A
vs
G

Sequencing (sampling) the two different alleles is just like flipping a coin.

$$\#\,\text{A alleles} \sim \text{Binomial}(n\,reads, p = 0.5)$$
$$\text{E}[\text{A depth}] = np = 0.5\,n$$
$$\text{Var}[\text{A depth}] = np(1-p) = 0.25\,n$$

Expect 2 A alleles and 2 G alleles

# Why Probabilistic Methods?



Sequence to average depth of 4x.

$$\text{Depth} \sim \text{Poisson}(\lambda=4)$$
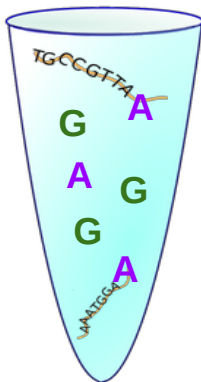$$\text{E}[\text{depth}] = \lambda$$
$$\text{Var}[\text{depth}] = \lambda$$

Expect 4 reads, all with As at this ref position.

Sequencing (sampling) the two different alleles is just like flipping a coin.

$$\#\,\text{A alleles} \sim \text{Binomial}(n\,reads, p=0.5)$$
$$\text{E}[\text{A depth}]=np=0.5n$$
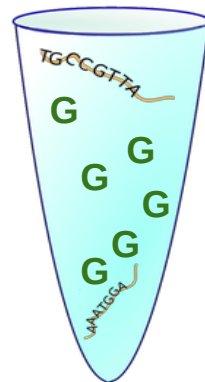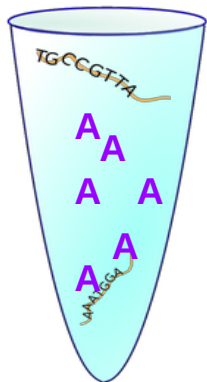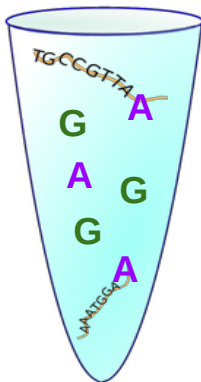$$\text{Var}[\text{A depth}]=np(1-p)=0.25n$$

vs

Expect 2 A alleles and 2 G alleles

# Why Probabilistic Methods?



Sequence to average depth of 4x.

$$\text{Depth} \sim \text{Poisson}(\lambda = 4)$$
$$\text{E}[\text{depth}] = \lambda$$
$$\text{Var}[\text{depth}] = \lambda$$

Sequencing error. Rates of ~0.1% for some Illumina platforms.

A
vs
G

Sequencing (sampling) the two different alleles is just like flipping a coin.

$$\#\,\text{A alleles} \sim \text{Binomial}(n\,reads, p = 0.5)$$
$$\text{E}[\text{A depth}] = np = 0.5\,n$$
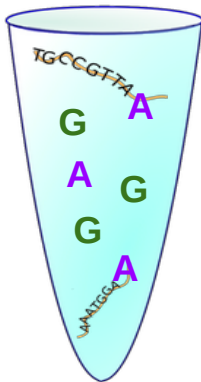$$\text{Var}[\text{A depth}] = np(1-p) = 0.25\,n$$
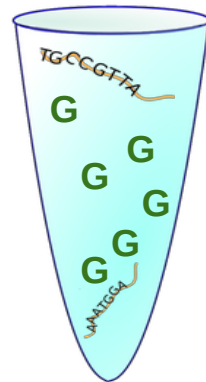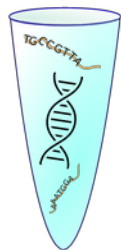
Expect 2 A alleles and 2 G alleles

# SAMtools mpileup representation of sequencing data for two individuals



Individual 1 data

Individual 2 data

```
6    .....,        DEGEGG    9    ,.,,.,,.,       DABGIIIIII
6    .....,        DEGEGG    9    ,.,,.,,.,       DABGIIIIII
6    .....,        DEGEGG    9    ,.,,.,,.,       DABGIIIIII
6    .....,        DEGEGG    10   ,.,,.,,.,^].    DABGIIIIIE
6    .....,        DEGEGG    10   ,.,,.,,.,.      DABGIIIIIII
6    .....,        DEGEGG    10   ,.,,.,,.,.      DABGIIIIIII
6    .....,        DEGEGG    10   ,.,,.,,.,.      DABGIIIIIII
7    .....,^].     DEGEGGE   10   ,.,,.,,.,.      DABGIIIIIII
7    .....,.       DEGEGGG   10   ,.,,.,,.,.      DABGIIIIIII
8    .....,.^],    DEGEGGGB  10   ,.,,.,,.,.      DABGIIIIIII
8    .....,.,      DEGEGGGB  10   ,.,,.,,.,.      DABGIIIIIII
8    .....,.,      DEGEGGGB  10   ,.,,.,,.,.      DABGIIIIIII
8    .....,.,      DEGEGGGB  10   ,.,,.,,.,.      DABGIIIIIII
9    .....,.,^],   DEGEGGGBE 10   ,.,,.,,.,.      DABGIIIIIII
9    .G...gG,,     DEGEGGGBG 10   ,.,,.,,.,.      D3BGIIIIIII
9    .....,.,,     DEGEGGGBG 10   ,.,,.,,.,.      D3BGIIIIIII
```
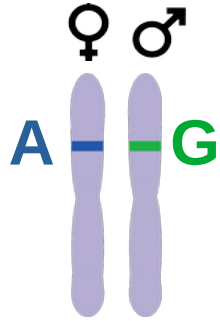
Each row is a different site in the reference genome

# reads    base IDs    base qualities

# Example sequencing data for one individual at Chr1:472



Maternally and paternally inherited chromosome 1 of a diploid individual.

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | AAG | DEG |

# Example sequencing data for one individual at Chr1:472



Maternally and paternally inherited chromosome 1 of a diploid individual.

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | AAG | DEG |

ASCII character "D" = decimal value 68
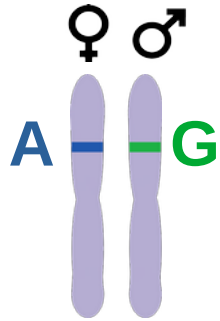
68 − 33 = base quality of 35

$$P(error) = \epsilon = 10^{\frac{-35}{10}} = 0.00032$$

# Example sequencing data for one individual at Chr1:472



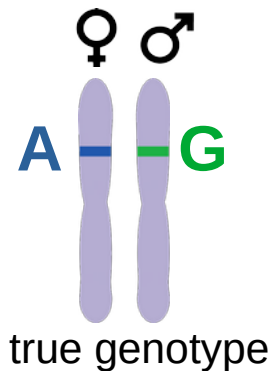| Scaffold | position | read bases | base qualities |
| --- | --- | --- | --- |
| Chr1 | 472 | AAG | DEG |

Maternally and paternally inherited chromosome 1 of a diploid individual.

This individual could have any of the following 10 genotypes (we can only see the sequencing data):

AA, AC, AG, AT, CG, CC, CT, GG, GT, TT

How do we figure out which genotype they are most likely to have based on the observed sequence data?

# Genotype likelihoods



| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | AAG | DEG |

true genotype

$$P\big(\text{Data}\big|\text{Genotype}=bh\big) = L\big(\text{Genotype}\,bh\big) = \text{likelihood of genotype}\,bh$$

$$b,h \in \{\text{A,C,G,T}\}$$

Possible genotypes: AA, AC, AG, AT, CG, CC, CT, GG, GT, TT

# What is the likelihood that the individual's genotype is AC?

Possible genotypes: AA, **AC**, AG, AT, CG, CC, CT, GG, GT, TT

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | AAG | DEG |

true genotype

P(observed *i*-th read | A allele)   P(observed *i*-th read | C allele)

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)=\boxed{\phantom{XX}}\times\boxed{\phantom{XX}}+\boxed{\phantom{XX}}\times\boxed{\phantom{XX}}$$

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

Assumed true genotype

# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)   P(observed read A | C allele)

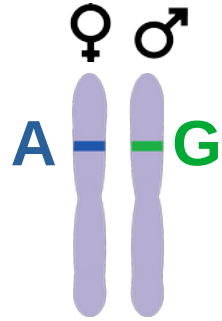$$P\left(\text{Data}\,|\,\text{Genotype}=\text{AC}\right)=\;\boxed{\phantom{xx}}\times\boxed{\phantom{xx}}\;+\;\boxed{\phantom{xx}}\times\boxed{\phantom{xx}}\times\boxed{\phantom{xx}}$$

Assumed true genotype

**Probability of sampling maternal chromosome (A allele)**

Probability of sampling paternal chromosome (C allele)

# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
| --- | --- | --- | --- |
| Chr1 | 472 | AAG | DEG |

P(observed read A | A allele)

P(observed read A | C allele)

$$P\left(\text{Data} \mid \text{Genotype} = \text{AC}\right) = \boxed{\dfrac{1}{2}} \times \boxed{\phantom{X}} + \boxed{\phantom{X}} \times \boxed{\phantom{X}} \times \boxed{\phantom{X}}$$

Assumed true genotype

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

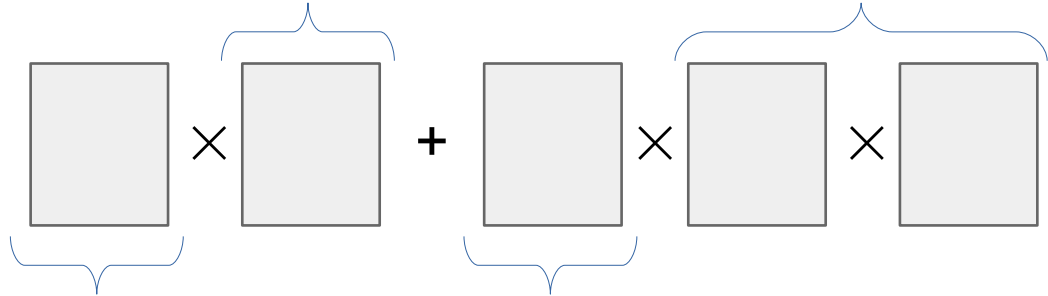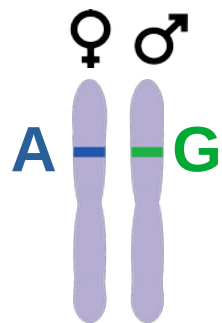# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)    P(observed read A | C allele)

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)= \boxed{\dfrac{1}{2}} \times \boxed{\phantom{xx}} + \boxed{\phantom{xx}} \times \boxed{\phantom{xx}} \times \boxed{\phantom{xx}}$$

Assumed true genotype

Probability of sampling maternal chromosome (A allele)
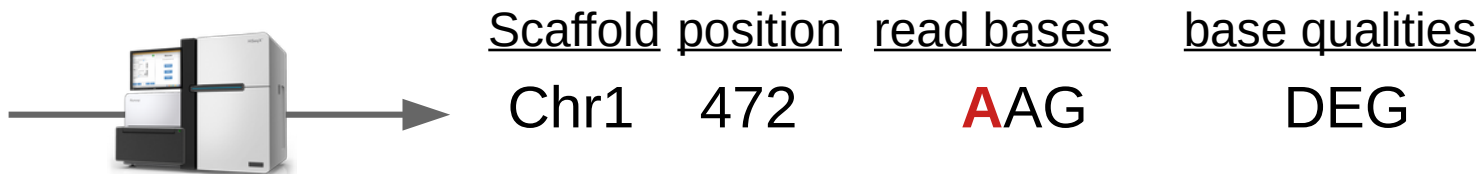
Probability of sampling paternal chromosome (C allele)

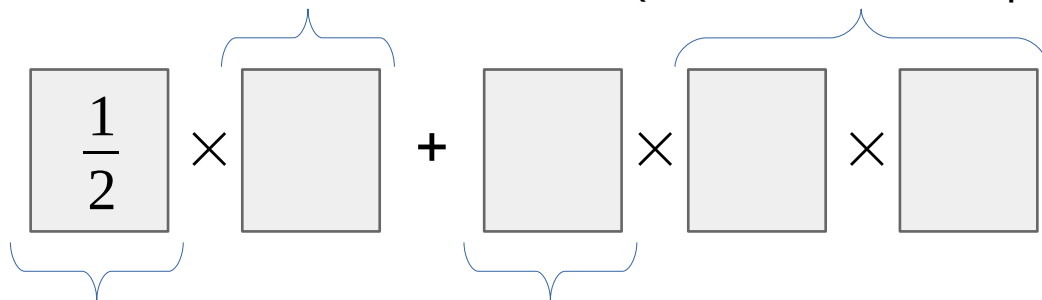# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)    P(observed read A | C allele)

$$P\left(\text{Data}\,\middle|\,\text{Genotype}=\text{AC}\right)= \boxed{\frac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\phantom{x}} \times \boxed{\phantom{x}} \times \boxed{\phantom{x}}$$

Assumed true genotype

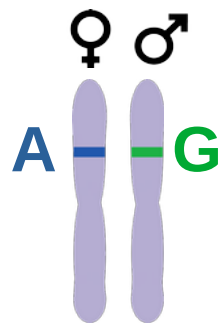Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

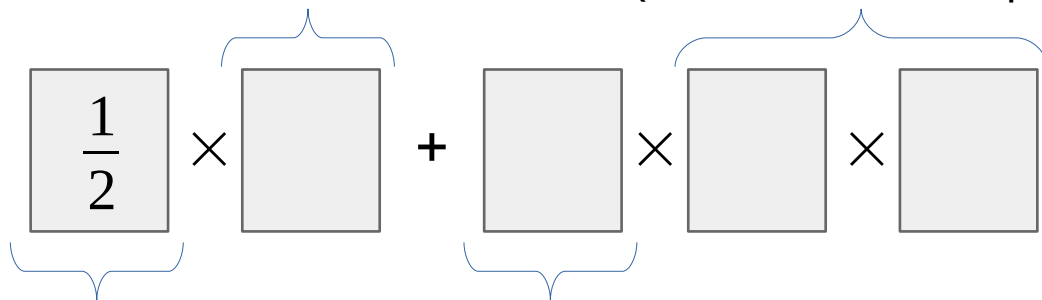# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)

P(observed read A | C allele)

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)= \boxed{\frac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\phantom{x}} \times \boxed{\phantom{x}} \times \boxed{\phantom{x}}$$

Assumed true genotype
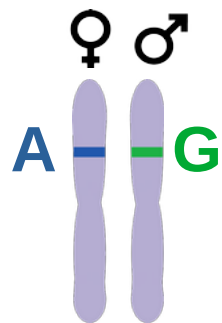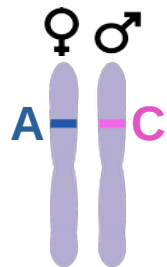
Probability of sampling maternal chromosome (A allele)

**Probability of sampling paternal chromosome (C allele)**

# What is the likelihood that the individual's genotype is AC?



Scaffold  position  read bases  base qualities

Chr1    472    **A**AG    DEG

true genotype

P(observed read A | A allele)          P(observed read A | C allele)

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)= \boxed{\frac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\frac{1}{2}} \times \boxed{\phantom{xx}} \times \boxed{\phantom{xx}}$$
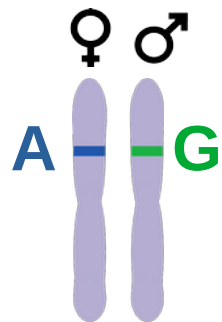
Assumed true genotype

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

# What is the likelihood that the individual's genotype is AC?



true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)

**P(observed read A | C allele)**

$$P\left(\text{Data}\big|\text{Genotype}=\text{AC}\right)= \boxed{\dfrac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\dfrac{1}{2}} \times \boxed{\phantom{xx}} \times \boxed{\phantom{xx}}$$
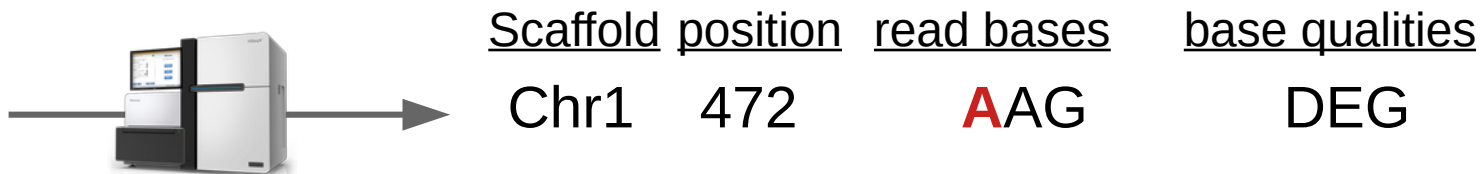
Assumed true genotype

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

# What is the likelihood that the individual's genotype is AC?
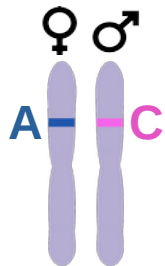


true genotype

| Scaffold | position | read bases | base qualities |
|----------|----------|------------|----------------|
| Chr1 | 472 | **A**AG | DEG |

P(observed read A | A allele)

P(observed read A | C allele)

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)= \boxed{\frac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\frac{1}{2}} \times \boxed{\epsilon_1} \times \boxed{\frac{1}{3}}$$
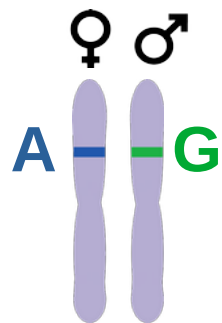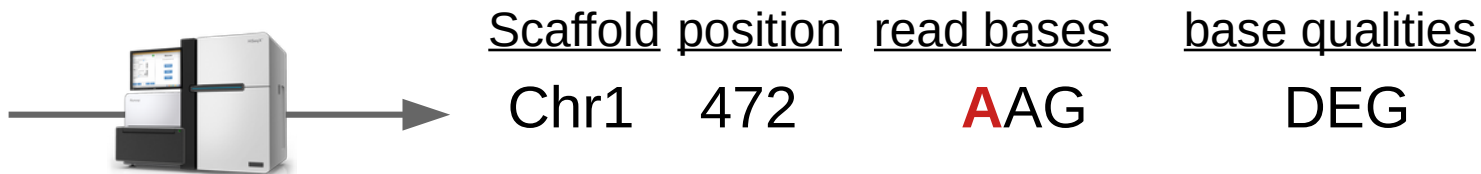
Assumed true genotype

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

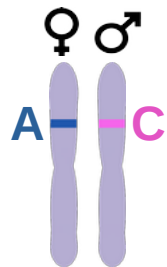# What is the likelihood that the individuals genotype is AC?



Assuming equal probability of changing to any of the possible 3 erroneous bases.

P(observed read A | A allele)    P(observed read A | C allele)

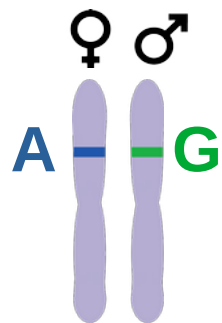$$P(\text{Data}|\text{Genotype}=\text{AC}) = \boxed{\frac{1}{2}} \times \boxed{1-\epsilon_1} + \boxed{\frac{1}{2}} \times \boxed{\epsilon_1} \times \boxed{\frac{1}{3}}$$

Probability of sampling maternal chromosome (A allele)

Probability of sampling paternal chromosome (C allele)

Assumed true genotype

true genotype

Scaffold  position    read bases    base qualities

Chr1    472    AAG    DEG

$$P\left(\text{Data}\,\middle|\,\text{Genotype}=\text{AC}\right)=\left[\frac{1}{2}\times\left(1-\epsilon_1\right)+\frac{1}{2}\times\epsilon_1\times\frac{1}{3}\right]\times$$

$$\left[\frac{1}{2}\times\left(1-\epsilon_2\right)+\frac{1}{2}\times\epsilon_2\times\frac{1}{3}\right]$$

Assumed true genotype

true genotype

Scaffold | position | read bases | base qualities
Chr1 | 472 | AA**G** | DEG

$$P\left(\text{Data}\middle|\text{Genotype}=\text{AC}\right)=\left[\frac{1}{2}\times\left(1-\epsilon_1\right)\quad\frac{1}{2}\times\epsilon_1\times\frac{1}{3}\right]\times$$

$$\left[\frac{1}{2}\times\left(1-\epsilon_2\right)+\frac{1}{2}\times\epsilon_2\times\frac{1}{3}\right]\times$$

$$\left[\frac{1}{2}\times\epsilon_3\times\frac{1}{3}+\frac{1}{2}\times\epsilon_3\times\frac{1}{3}\right]$$

**General genotype likelihood expression**

$$\mathrm{P}\big(\mathrm{Data} \mid \mathrm{Genotype}{=}bh\big)=\prod_{i=1}^{r}\left(\frac{L_b^{(i)}}{2}+\frac{L_h^{(i)}}{2}\right)$$

$$b,h \in \{\mathrm{A},\mathrm{C},\mathrm{G},\mathrm{T}\}$$

$$L_b^{(i)}=\mathrm{P}\big(\text{observed read}=x_i\big|\text{assumed true allele}=b\big)$$

$$L_h^{(i)}=\mathrm{P}\big(\text{observed read}=x_i\big|\text{assumed true allele}=h\big)$$

$$\frac{\epsilon_i}{3}\ \text{if } b,h \neq x_i$$

$$1-\epsilon_i\ \text{if } b,h = x_i$$

# Representation of genotype likelihoods in ANGSD

Chr1   472   AAG   555

ASCII character "5" = decimal value 53

53 – 33 = base quality of 20

$$\epsilon = 10^{\frac{-20}{10}} = 0.01$$

# Representation of genotype likelihoods in ANGSD



Chr1   472   AAG   555

Instead of a single genotype, AG, we have a distribution over all possible genotypes:

| Genotype | AA | AC | AG | AT | CC | CG | CT | GG | GT | TT |
|----------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Log$_{10}$ likelihood | -2.49 | -3.08 | -0.91 | -3.08 | -7.43 | -5.26 | -7.43 | -4.96 | -5.26 | -7.43 |

$$\epsilon = 10^{\frac{-20}{10}} = 0.01$$

# Representation of genotype likelihoods in ANGSD



Chr1   472   AAG   555

Instead of a single genotype, AG, we have a distribution over all possible genotypes:

| Genotype | AA | AC | **AG** | AT | CC | CG | CT | GG | GT | TT |
|---|---|---|---|---|---|---|---|---|---|---|
| **Log$_{10}$ likelihood** | -2.49 | -3.08 | **-0.91** | -3.08 | -7.43 | -5.26 | -7.43 | -4.96 | -5.26 | -7.43 |

Maximum likelihood estimate of the genotype

$$\epsilon = 10^{\frac{-20}{10}} = 0.01$$

# Exercise. Calculate genotype likelihoods with ANGSD.

# Estimating allele frequencies

When genotypes are known, allele frequencies can be calculated by simply counting alleles.

| | Ind1 | Ind2 | Ind3 | Ind4 | Ind5 | Ind6 | Ind7 | Ind8 | Ind9 | Ind10 | Minor allele frequency (MAF) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Site 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0.3 |
| Site 2 | 0 | 2 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0.4 |

Genotype notation
0 = zero minor alleles
1 = one minor alleles
2 = two minor alleles

# Estimating allele frequencies

When genotypes are known, allele frequencies can be calculated by simply counting alleles.

|  | Ind1 | Ind2 | Ind3 | Ind4 | Ind5 | Ind6 | Ind7 | Ind8 | Ind9 | Ind10 |  | Minor allele frequency (MAF) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Site 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | = | 0.3 |
| Site 2 | 0 | 2 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | = | 0.4 |

But how do you estimate allele frequencies when you have a distribution of genotype likelihoods?

**Maximum likelihood estimation of allele frequencies**

$$\mathrm{P}\left(\mathrm{Data}|f\right)=\prod_{i=1}^{n}\sum_{g\in\{0,1,2\}}\mathrm{P}\left(D_i|\mathrm{Genotype}_i=g\right)\mathrm{P}\left(\mathrm{Genotype}_i=g|f\right)$$

$D_i=$ sequencing data for individual $i$

$f=$ population minor allele frequency

Genotype notation

0 = zero minor alleles
1 = one minor alleles
2 = two minor alleles

# Maximum likelihood estimation of allele frequencies

$$P\left(\text{Data}|f\right) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} \underbrace{P\left(D_i|\text{Genotype}_i = g\right)} P\left(\text{Genotype}_i = g|f\right)$$

This the likelihood of genotype *g* for individual *i* calculated as shown previously.

$D_i =$ sequencing data for individual $i$

$f =$ population minor allele frequency

Genotype notation

0 = zero minor alleles
1 = one minor alleles
2 = two minor alleles

# Maximum likelihood estimation of allele frequencies

$$P(\text{Data}|f) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} P(D_i|\text{Genotype}_i = g) \, P(\text{Genotype}_i = g|f)$$

This the likelihood of genotype *g* for individual *i* calculated as shown previously.

Hardy-Weinberg frequency of genotypes.

$D_i =$ sequencing data for individual $i$

$f =$ population minor allele frequency

Genotype notation

0 = zero minor alleles
1 = one minor alleles
2 = two minor alleles

# Genotype frequencies under Hardy-Weinberg Equilibrium (HWE)

$$P(\text{Data}|f) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} P(D_i|\text{Genotype}_i = g) \, P(\text{Genotype}_i = g|f)$$

An "infinitely" large population of sexually reproducing diploid organisms segregating for alleles *A* and *a*, and for which

- Mating is random.

- Generations are nonoverlapping.

- Allele frequencies are the same in males and females.

- No migration, mutation, or selection.

# Genotype frequencies under Hardy-Weinberg Equilibrium (HWE)

$$\mathrm{P}\left(\mathrm{Data}|f\right)=\prod_{i=1}^{n}\sum_{g\in\{0,1,2\}}\mathrm{P}\left(D_i|\mathrm{Genotype}_i=g\right)\mathrm{P}\left(\mathrm{Genotype}_i=g|f\right)$$

Offspring are formed through independent draws of gametes from this population, expressed as $(f_A + f_a)^2$, which expanded yields:

$$f_{AA} = f_A^2 = \left(1-f_a\right)^2 = \mathrm{P}\left(\mathrm{Genotype}=AA|f_a\right)$$

$$f_{Aa} = 2f_A f_a = 2\left(1-f_a\right)f_a = \mathrm{P}\left(\mathrm{Genotype}=1|f_a\right)$$

$$f_{aa} = f_a^2 = \mathrm{P}\left(\mathrm{Genotype}=2|f_a\right)$$

*Aa* *Aa* *Aa* *AA* *AA* *aa* *AA* *AA* *Aa* *AA* *Aa* *aa* *Aa* *AA* *Aa*

**Maximum likelihood estimation of allele frequencies**

$$P(\text{Data}|f) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} P(D_i|\text{Genotype}_i = g) P(\text{Genotype}_i = g|f)$$

Summing over the possible genotypes accounts for genotyping uncertainty.

Marginal probabilities are used to account for uncertainty in various quantities associated with low coverage sequencing. In general, for random variables *X* and *Y*

$$P(X=x) = \sum_{y} P(X=x|Y=y) P(Y=y)$$

**Maximum likelihood estimation of allele frequencies**

$$P(\text{Data}|f) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} P(D_i|\text{Genotype}_i = g) P(\text{Genotype}_i = g|f)$$

The value of *f* that maximizes the likelihood function above yields a **maximum likelihood estimate of *f*** :

$$\hat{f} = \text{argmax}_f P(Data|f)$$

# Using the ML allele frequency estimate to identify polymorphic sites

$$P\left(\text{Data}|f\right) = \prod_{i=1}^{n} \sum_{g \in \{0,1,2\}} P\left(D_i | \text{Genotype}_i = g\right) P\left(\text{Genotype}_i = g | f\right)$$

$$\hat{f} = \text{argmax}_f\, P\left(Data|f\right)$$

←—— **maximum likelihood estimate of $f$**

Probability of the sequencing data when $f = 0$, i.e., the site is monomorphic (null case).

$$\lambda = -2 \ln\left(\frac{P\left(\text{Data}|f_0\right)}{P\left(\text{Data}|\hat{f}\right)}\right) = -2\left[\ln\left(P\left(\text{Data}|f_0\right)\right) - \ln\left(P\left(\text{Data}|\hat{f}\right)\right)\right]$$

$$\lambda \sim \chi^2\left(1\,\text{degree of freedom}\right) \longrightarrow$$ Call SNPs at a given level of statistical confidence.

# Exercise. Estimate allele frequencies with ANGSD.

Can we use other information from our data to further increase our genotyping accuracy?

```
6  .....,          DEGEGG    9   ,.,,.,,.,      DABGIIIII   5  ,,.,.  >AB/A   6  .,,...  FGG
6  .....,          DEGEGG    9   ,.,,.,,.,      DABGIIIII   5  ,,.,.  >ABDA   6  .,,...  3GGBGB
6  .....,          DEGEGG    9   ,.,,.,,.,      DABGIIIII   5  ,,.,.  >ABDA   6  .,,...  3GGBGB
6  .....,          DEGEGG    10  ,.,,.,,.,^].   DABGIIIIIE  5  ,,.,.  >AB/A   6  .,,...  3GGBGB
6  .....,          DEGEGG    10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >AB/A   6  .,,...  3GGBGB
6  .....,          DEGEGG    10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABDA   6  .,,...  BGGBGB
6  .....,          DEGEGG    10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABDA   6  C,,...  5G/BGB
7  .....,^].       DEGEGGE   10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABDA   6  .,,...  5GIBGB
7  ......,.        DEGEGGG   10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABDA   6  .,,...  5GIBGB
8  ......,.^],     DEGEGGGB  10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABDA   6  .,,...  5GIBGB
8  ......,.,       DEGEGGGB  10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >AB/A   6  .,,...  DGIBGB
8  ......,.,       DEGEGGGB  10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >ABAA   6  .,,...  DGIBGB
8  ......,.,       DEGEGGGB  10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >/BAA   6  .,,...  DGIBGB
9  ......,.,^],    DEGEGGGBE 10  ,.,,.,,.,.     DABGIIIIII  5  ,,.,.  >BBAA   6  .,,...  DGIBGB
9  .G...gG,,       DEGEGGGBG 10  ,.,,.,,.,.     D3BGIIIIII  5  ,,.,C  >BBAA   6  .,,...  DGIBGB
9  ......,.,,      DEGEGGGBG 10  ,.,,.,,.,.     D3BGIIIIII  5  ,,.,.  >BBAA   6  .,,...  /GIBGB
```

Wouldn't it be awesome if you knew what the frequency of C was in the population.

# Bayesian Inference

$$P(\theta|X) = \frac{P(X|\theta)\,P(\theta)}{P(X)} = \frac{P(X|\theta)\,P(\theta)}{\sum_{\theta} P(X|\theta)\,P(\theta)}$$



$P(X|\theta)$  
$P(\theta|X)$  
Posterior

Likelihood

$P(\theta)$

Prior

density

value

$X$ = Data  
$\theta$ = Parameter

Distribution plots from Bink 2008

# Genotype posterior probabilities to improve genotyping accuracy

Using Bayes' Theorem, the posterior probability of genotype *g* is

$$\mathrm{P}\big(\mathrm{Genotype}{=}g\big|\mathrm{Data}\big) = \frac{\mathrm{P}\big(\mathrm{Data}\big|\mathrm{Genotype}{=}g\big)\,\mathrm{P}\big(\mathrm{Genotype}{=}g\big)}{\displaystyle\sum_{g\in\{0,1,2\}}\mathrm{P}\big(\mathrm{Data}\big|\mathrm{Genotype}{=}g\big)\,\mathrm{P}\big(\mathrm{Genotype}{=}g\big)}$$

# Genotype posterior probabilities to improve genotyping accuracy

Using Bayes' Theorem, the posterior probability of genotype $g$ is

Likelihood of genotype $g$ calculated as shown previously.

$$\mathrm{P}\left(\text{Genotype}=g \middle| \text{Data}\right) = \frac{\mathrm{P}\left(\text{Data}\middle|\text{Genotype}=g\right)\mathrm{P}\left(\text{Genotype}=g\right)}{\displaystyle\sum_{g\in\{0,1,2\}} \mathrm{P}\left(\text{Data}\middle|\text{Genotype}=g\right)\mathrm{P}\left(\text{Genotype}=g\right)}$$

# Genotype posterior probabilities to improve genotyping accuracy

Using Bayes' Theorem, the posterior probability of genotype $g$ is

Given an estimate of the population minor allele frequency, $f$, under HWE

$$P(\text{Genotype}=0|f)=(1-f)^2$$

$$P(\text{Genotype}=1|f)=2f(1-f)$$

$$P(\text{Genotype}=2|f)=f^2$$

Likelihood of genotype $g$ calculated as shown previously.

$$P(\text{Genotype}=g|\text{Data}) = \frac{P(\text{Data}|\text{Genotype}=g)\,P(\text{Genotype}=g)}{\displaystyle\sum_{g\in\{0,1,2\}} P(\text{Data}|\text{Genotype}=g)\,P(\text{Genotype}=g)}$$

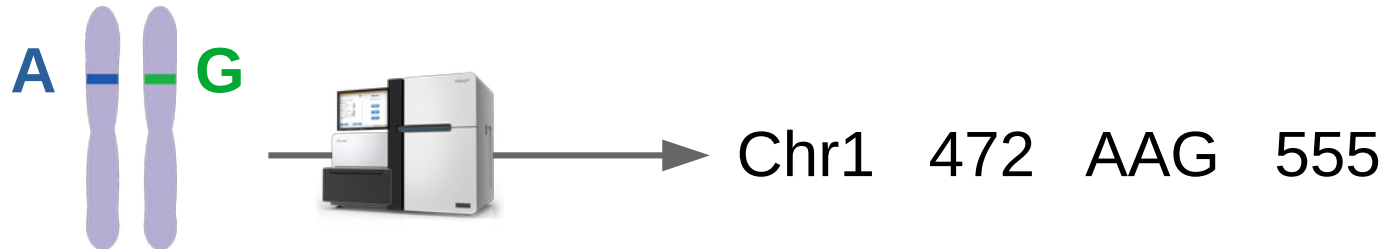# Genotype posterior probabilities to improve genotyping accuracy

Using Bayes' Theorem, the posterior probability of genotype *g* is

Likelihood of genotype *g* calculated as shown previously.

Note: factors like inbreeding can easily be incorporated into the genotype posterior probabilities by conditioning on the allele frequency and inbreeding coefficient.

$$\mathrm{P}\left(\mathrm{Genotype}=g\,|\,\mathrm{Data}\right) = \frac{\mathrm{P}\left(\mathrm{Data}\,|\,\mathrm{Genotype}=g\right)\mathrm{P}\left(\mathrm{Genotype}=g\right)}{\sum_{g\in\{0,1,2\}}\mathrm{P}\left(\mathrm{Data}\,|\,\mathrm{Genotype}=g\right)\mathrm{P}\left(\mathrm{Genotype}=g\right)}$$

# Example genotype posterior probability distribution

A — G

Chr1  472  AAG  555

Assume we estimate f(A) = 0.7, f(G) = 0.3

| Genotype | $Log_{10}$ likelihood | Prior | Posterior probability |
|---|---|---|---|
| AA | -2.49 | $P(\text{Genotype} = \text{AA}) = 0.7^2 = 0.49$ | 0.03 |
| AG | -0.91 | $P(\text{Genotype} = \text{AG}) = 2 \times 0.7 \times 0.3 = 0.42$ | 0.97 |
| GG | -4.96 | $P(\text{Genotype} = \text{GG}) = 0.3^2 = 0.09$ | 0.00 |

# Example genotype posterior probability distribution

A — G

Chr1   472   AAG   555

Assume we estimate f(A) = 0.7, f(G) = 0.3

We could call most probable genotype, AG, and have an associated degree of confidence (prob = 0.97).

| Genotype | Log$_{10}$ likelihood | Prior | Posterior probability |
|---|---|---|---|
| AA | -2.49 | $P(\text{Genotype} = \text{AA}) = 0.7^2 = 0.49$ | 0.03 |
| AG | -0.91 | $P(\text{Genotype} = \text{AG}) = 2 \times 0.7 \times 0.3 = 0.42$ | 0.97 |
| GG | -4.96 | $P(\text{Genotype} = \text{GG}) = 0.3^2 = 0.09$ | 0.00 |

Exercise. Calculate genotype posterior probabilities and call genotypes with ANGSD.