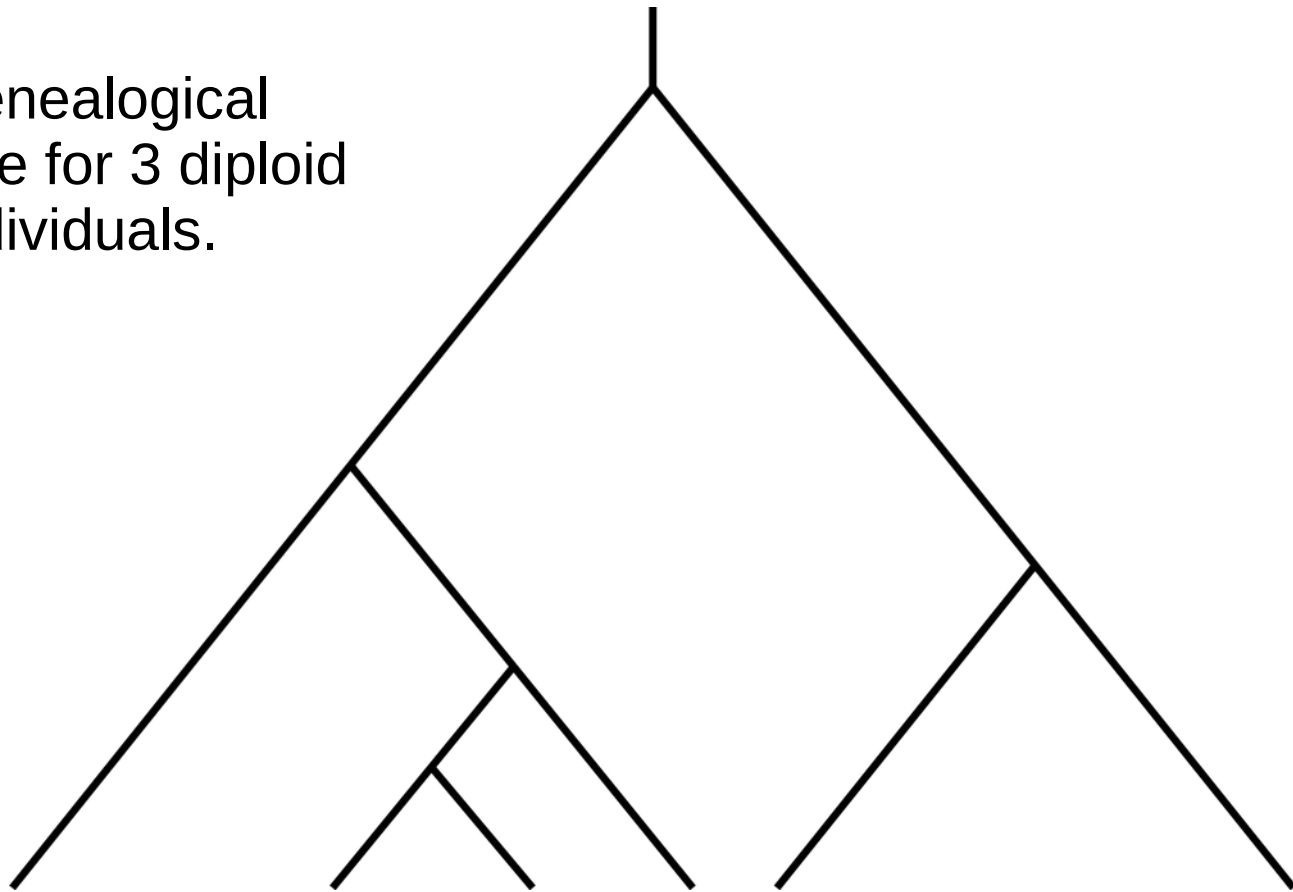# Site Frequency Spectrum (SFS)

Tyler Linderoth
Physalia lcWGS course 2025

Genealogical
tree for 3 diploid
individuals.

past

time

present

Genealogical tree for 3 diploid individuals; $2n = 6$ chromosomes.
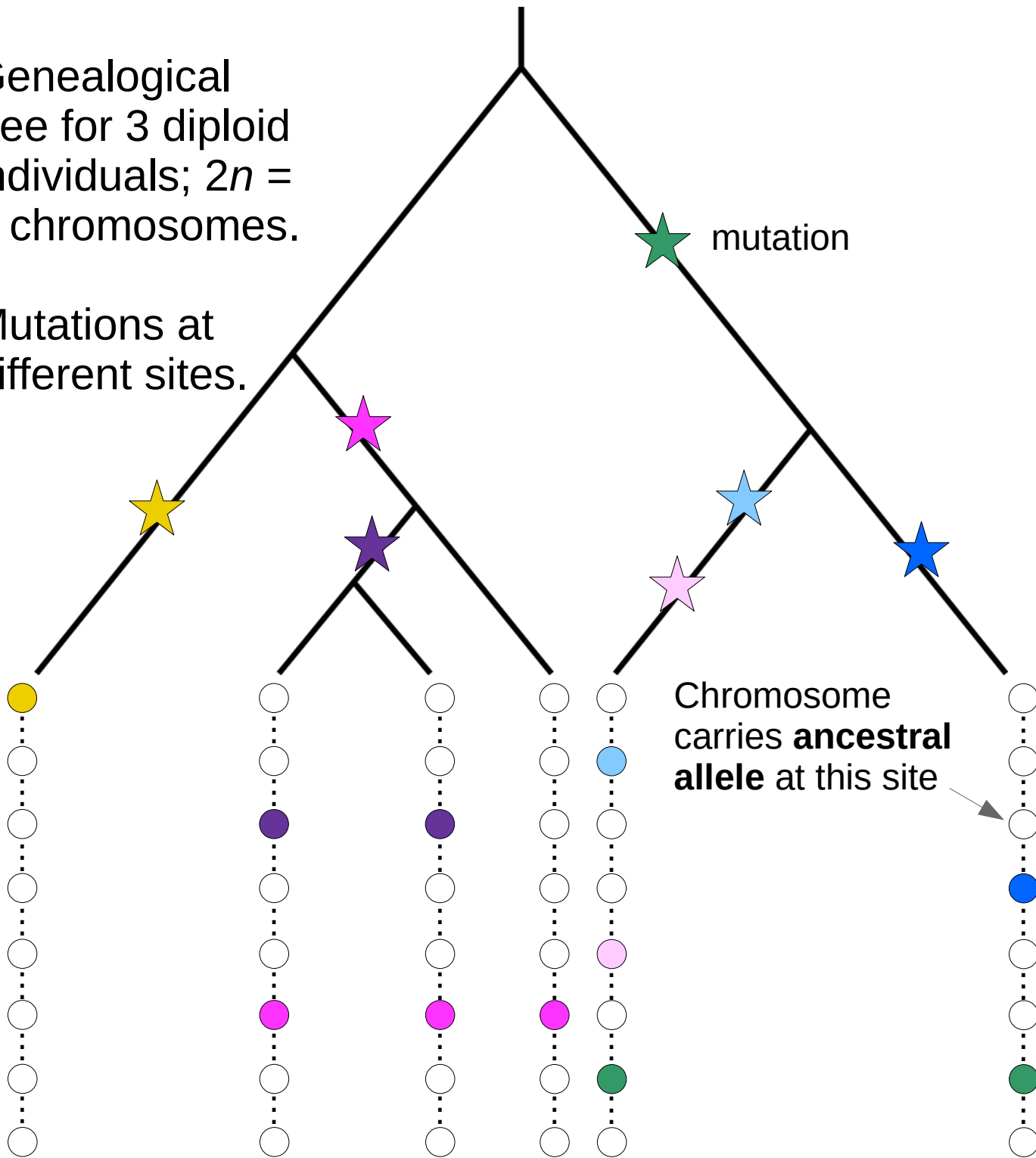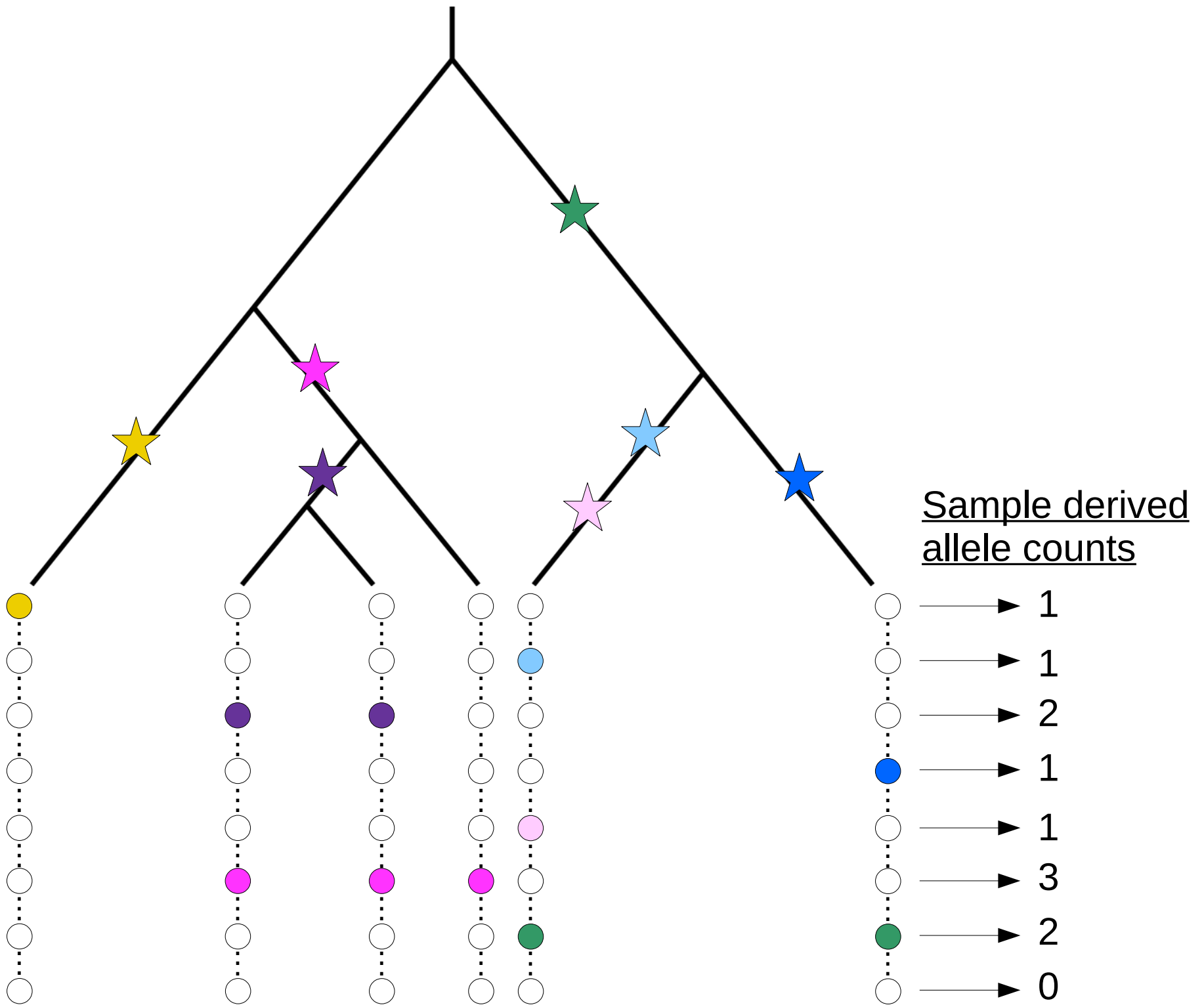
Mutations at different sites.

mutation

past

time

present

Chromosome carries **ancestral allele** at this site

Chromosome carries **derived allele** at this site

Sample derived
allele counts

1
1
2
1
1
3
2
0

number of sites
with derived allele

sample derived
allele frequency

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | 4 | 2 | 1 | 0 | 0 | 0 |

Sample derived
allele counts

1
1
2
1
1
3
2
0
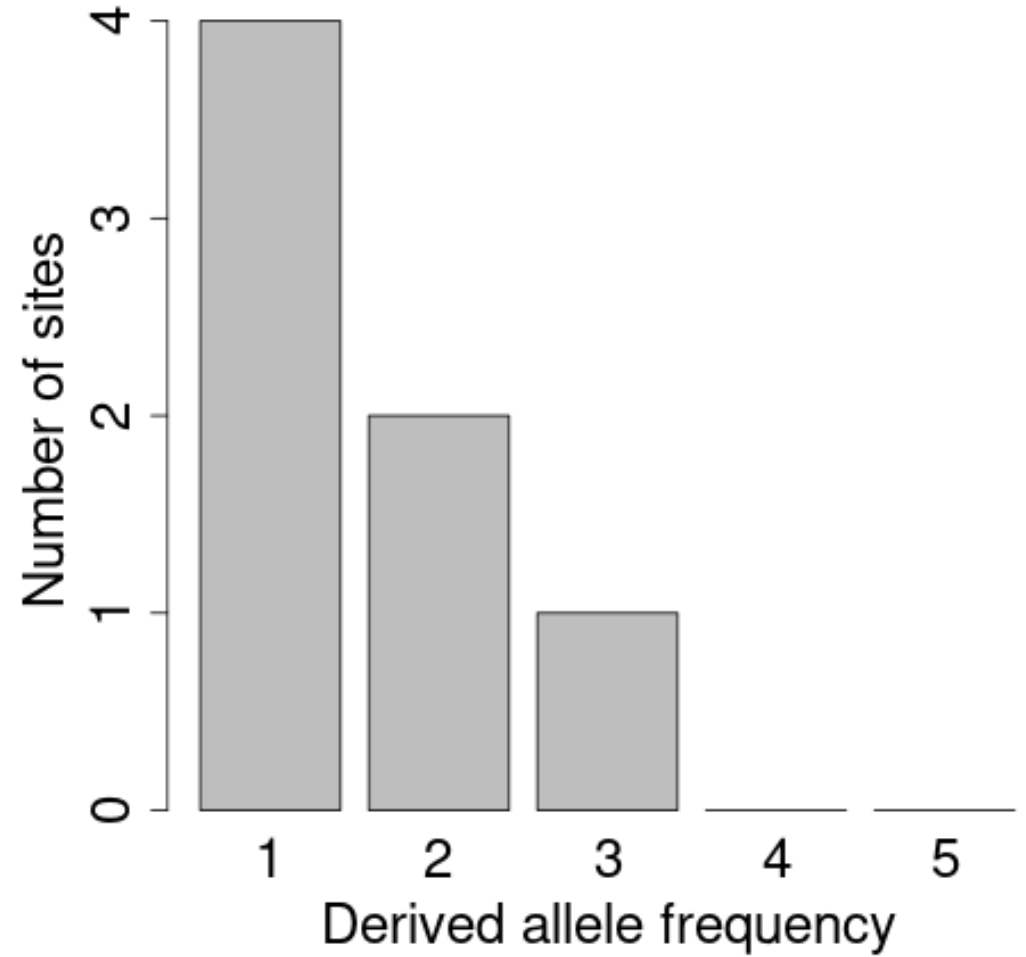
number of sites with derived allele

sample derived allele frequency

# Neutral, constant-size population SFS

Proportion of sites in each allele frequency class proportional to 1/*k,* where *k* is the derived allele frequency. In other words,

$$\mathrm{P}\big(\text{site has } x \text{ derived alleles}\big) \propto \frac{1}{k}$$



Proportion of sites vs. Derived allele frequency

Lineages coalesce at rate 1/*N. N* is the effective population size.

population expansion

More chance for mutation on long external branches

# Lineages coalesce at rate 1/*N*. *N* is the effective population size.



population bottleneck

Less chance for mutation on short external branches

Effects of selection on the SFS

Nielson 2005

# Let's play the guess the bioinformatic problem game

SFS for a sample of 4 diploid individuals
Case 1

Proportion of sites

Derived allele frequency

# SFS for a sample of 4 diploid individuals
# Case 1



- Mapping problems from duplications or repetitive sequence.

Proportion of sites

Derived allele frequency

SFS for a sample of 4 diploid individuals
Case 2

Proportion of sites

Derived allele frequency

# SFS for a sample of 4 diploid individuals
## Case 2



Could be many things:

- Insufficient map and/or base quality filtering.
- Overlapping read pairs.
- HWE genotype prior in the presence of inbreeding.
- etc.

Proportion of sites

Derived allele frequency

1   2   3   4   5   6   7

# SFS for a sample of 4 diploid individuals
## Case 3

# SFS for a sample of 4 diploid individuals
## Case 3

- Misidentification of ancestral and derived allelic state.



Proportion of sites

Derived allele frequency

1   2   3   4   5   6   7

# SFS for a sample of 4 diploid individuals
# Case 4

# SFS for a sample of 4 diploid individuals
## Case 4



- Intense inbreeding.

Proportion of sites

Derived allele frequency

1  2  3  4  5  6  7

Top-left panel:
- Mapping problems from paralogs or repetitive sequence.

(x-axis: Derived allele frequency; y-axis: Proportion of sites)

Top-right panel:
Could be many things:
- Insufficient map and/or base quality filtering.
- Overlapping read pairs.
- HWE genotype prior in the presence of inbreeding.
- etc.

(x-axis: Derived allele frequency; y-axis: Proportion of sites)

Bottom-left panel:
- Misidentification of ancestral and derived allelic state.

(x-axis: Derived allele frequency; y-axis: Proportion of sites)

Bottom-right panel:
- Intense inbreeding.

(x-axis: Derived allele frequency; y-axis: Proportion of sites)

Likelihood function for the SFS

$$\text{SFS} = \boldsymbol{P} = \left( p_0, p_1, p_2, p_3, p_4, ..., p_{2n} \right)$$

$p_k$: probability that a random site in the genome has $k$ derived alleles
$n$ = diploid sample size

$$\text{SFS} = \boldsymbol{P} = \left( p_0, p_1, p_2, p_3, p_4, \ldots, p_{2n} \right)$$

Likelihood of the SFS with *known* genotypes considering a single site:

*X*: observed data
(sequencing reads)

$$\mathrm{P}(X, \boldsymbol{G} | \boldsymbol{P}) = \sum_{k=0}^{2n} \mathrm{P}(X, G | K = k) \, \mathrm{P}(K = k | \boldsymbol{P})$$

**G**: genotype vector
= $(G_1, G_2, G_3, \ldots, G_n)$

*K*: number of derived alleles in genotype vector **G**

Likelihood of the SFS with *known* genotypes considering a single site:

$$\mathrm{P}(X, \boldsymbol{G} | \boldsymbol{P}) = \sum_{k=0}^{2n} \mathrm{P}(X, \boldsymbol{G} | K=k) \, \mathrm{P}(K=k | \boldsymbol{P})$$

$$= p_k$$

$$\mathrm{P}(a,b|c) = \mathrm{P}(a|b)\mathrm{P}(b|c)$$

Recall that SFS = $\boldsymbol{P}$ = $(p_0, p_1, \ldots, p_{2n})$

$$= \sum_{k=0}^{2n} \mathrm{P}(X | \boldsymbol{G}) \mathrm{P}(\boldsymbol{G} | K=k) \, \mathrm{P}(K=k | \boldsymbol{P})$$

Likelihood of the SFS with *known* genotypes considering a single site:

$$\mathrm{P}(X,\boldsymbol{G}|\boldsymbol{P})=\sum_{k=0}^{2n}\boxed{\mathrm{P}(X,\boldsymbol{G}|K=k)}\boxed{\mathrm{P}(K=k|\boldsymbol{P})}$$

$$=p_k$$

$$\mathrm{P}(a,b|c)=\mathrm{P}(a|b)\mathrm{P}(b|c)$$

Recall that SFS = $\boldsymbol{P}$ = $(p_0, p_1, \ldots, p_{2n})$

$$=\sum_{k=0}^{2n}\boxed{\mathrm{P}(X|\boldsymbol{G})\,\mathrm{P}(\boldsymbol{G}|K=k)}\mathrm{P}(K=k|\boldsymbol{P})$$

Consider the each individual's data and genotype independently.

$$=\sum_{k=0}^{2n}\boxed{\left(\prod_{i=1}^{n}\mathrm{P}(X_i|G_i)\right)\mathrm{P}(\boldsymbol{G}|K=k)\,\mathrm{P}(K=k|\boldsymbol{P})}$$

Likelihood of the SFS with *known* genotypes considering a single site:

$$\mathrm{P}(X,\boldsymbol{G}|\boldsymbol{P})=\sum_{k=0}^{2n}\left(\prod_{i=1}^{n}\mathrm{P}(X_i|G_i)\right)p(\boldsymbol{G}|K=k)\mathrm{P}(K=k|\boldsymbol{P})$$

Probability based on the number of ways to have *k* derived alleles in **G** out of the total number of ways to to arrange *k* derived alleles among 2*n* sampled chromosomes. ⟶ **Assumes HWE**

**Example for 4 diploid individual with derived allele *A* and ancestral allele *a*.**

*k* = 4

*m* = # heterozygotes = 2

```
AA  aA  aa  aA
AA  Aa  aa  aA
AA  aA  aa  Aa
AA  Aa  aa  Aa
```
$2^m = 2^2$
= 4 combinations

$$\mathrm{P}(\boldsymbol{G}|K=k)=\frac{2^m}{\dbinom{2n}{k}}=\frac{2^2}{\dbinom{8}{4}}=\frac{4}{70}$$

$$P(X, \boldsymbol{G}|\boldsymbol{P}) = \sum_{k=0}^{2n} \left( \prod_{i=1}^{n} P(X_i|G_i) \right) p(\boldsymbol{G}|K=k) P(K=k|\boldsymbol{P})$$

Allow for unknown genotypes (consider the likelihood for all possible genotypes)

$$P(X|\boldsymbol{P}) =$$

$$\sum_{k=0}^{2n} P(K=k|\boldsymbol{P}) \sum_{G_1 \in \{0,1,2\}} \cdots \sum_{G_n \in \{0,1,2\}} P(G|K=k) \prod_{i=1}^{n} P(X_i|G_i)$$

Summing over the uncertainty of the identity of genotypes contained in **G**

Likelihood of **P** considering site $v$

$$P\left(X^v|\boldsymbol{P}\right)=$$

$$\sum_{k=0}^{2n} P\left(K=k|\boldsymbol{P}\right) \sum_{G_1\in\{0,1,2\}} \cdots \sum_{G_i\in\{0,1,2\}} P\left(G^v|K=k\right)\prod_{i=1}^{n} P\left(X_i^v|G_i^v\right)$$

**Assume sites are independent**

Likelihood of **P** considering all sites (the SFS)

$$P\left(X|\boldsymbol{P}\right)=$$

$$\prod_{v=1}^{\text{all sites}} \sum_{k=0}^{2n} P\left(K=k|\boldsymbol{P}\right) \sum_{G_1\in\{0,1,2\}} \cdots \sum_{G_d\in\{0,1,2\}} P\left(G^v|K=k\right)\prod_{d=1}^{n} P\left(X_i^v|G_i^v\right)$$

Product over all sites

# Folding the SFS when ancestral state is unknown

1 derived allele, $2n$-1
ancestral alleles

Collapse equivalent allele frequency
classes if derived/ancestral state is
unknown

+

+

+

2n-1 derived alleles, 1
ancestral allele

Number of sites

1    2    3    4    5    6    7

Derived allele frequency

# Folding the SFS when ancestral state is unknown



Unfolded SFS

Number of sites

Derived allele frequency

Folded SFS

Number of sites

Minor allele frequency

# How to estimate posterior probabilities of allele frequencies: allele frequency likelihoods

ANGSD -doSaf 1: Table of allele frequency likelihoods

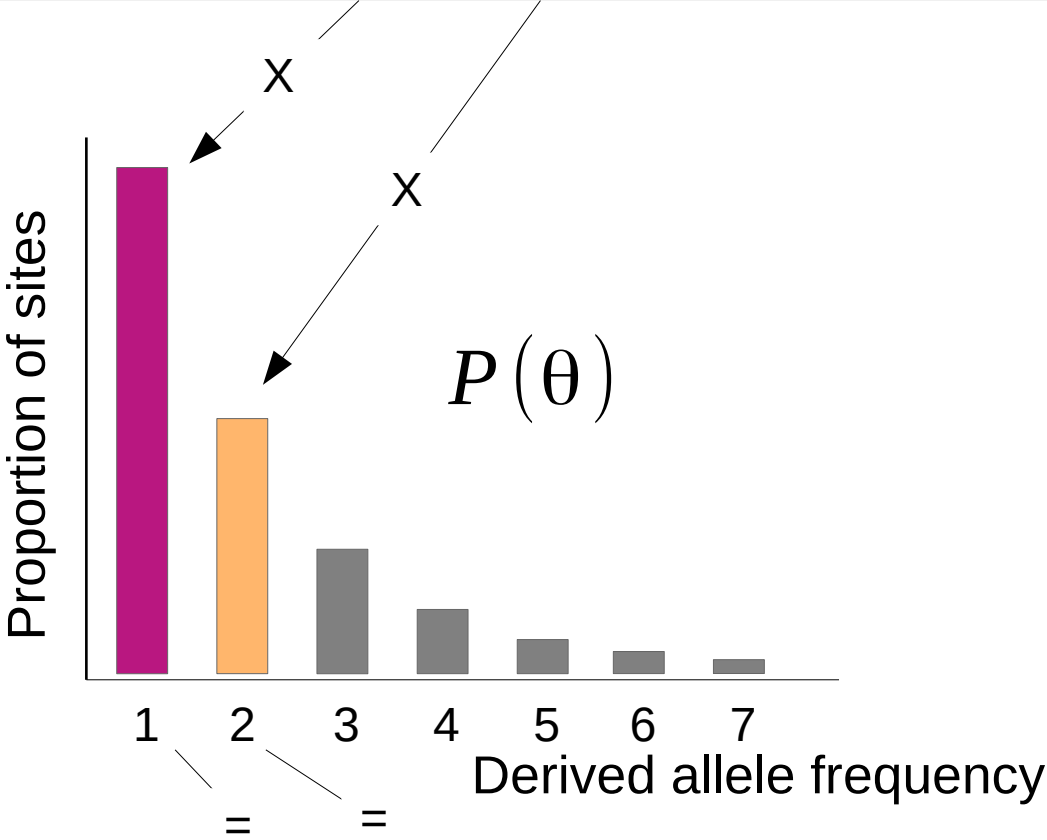| | 0 | 1 | 2 | 3 | 4 | . . . | 2n |
|---|---|---|---|---|---|---|---|
| site1 | 0.00 | -2.24 | -4.53 | -6.99 | -9.63 | | -232.69 |
| site2 | 0.00 | -2.24 | -4.53 | -6.99 | -9.63 | | -232.69 |
| site3 | -76.63 | -37.87 | -10.42 | 0.00 | -9.59 | | -467.13 |
| site4 | 0.00 | -2.24 | -5.53 | -6.99 | -9.63 | | -237.55 |
| . | | | | | | | |
| . | | | | | | | |
| . | | | | | | | |
| . | | | | | | | |
| . | | | | | | | |
| . | | | | | | | |
| sitem | 0.00 | -8.62 | -19.22 | -30.67 | -43.27 | | -626.78 |

← # derived alleles

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

# How to estimate posterior probabilities of allele frequencies: Allele frequency prior probabilities



Number of sites

Gives the probability of sampling a site with 4 derived alleles

Derived allele frequency

$$P(\theta | X) = P(X | \theta) \, P(\theta)$$

|         | 0    | 1      | 2      | 3     | 4     | . . . | 2n       | $P(X\|\theta)$ |
|---------|------|--------|--------|-------|-------|-------|----------|------|
| site1   | 0.00 | -2.24  | -4.53  | -6.99 | -9.63 |       | -232.69  |      |
| site2   | 0.00 | -2.24  | -4.53  | -6.99 | -9.63 |       | -232.69  |      |

$P(\theta)$

Proportion of sites

Derived allele frequency

|         | 0      | 1      | 2      | 3      | 4      | . . . . | 2n     | $P(\theta\|X)$ |
|---------|--------|--------|--------|--------|--------|---------|--------|------|
| site1   | 0.9892 | 0.0101 | 0.0006 | 0.0000 | 0.0000 |         | 0.0000 |      |
| site2   | 0.9892 | 0.0101 | 0.0006 | 0.0000 | 0.0000 |         | 0.0000 |      |

# Population genetic summary statistics derived from the SFS

# Estimating diversity from the SFS



number of
segregating sites

number of
sites with $j$
derived alleles

$$\hat{\theta}_W = \frac{S}{\sum_{j=1}^{2n-1} \frac{1}{j}} = \frac{\sum_{j=1}^{2n-1} \eta_j}{\sum_{j=1}^{2n-1} \frac{1}{j}}$$

Expected length
of a coalescent
tree of $n$ diploid
individuals. This
scales $S$ based
on sample size.

Number of sites

Derived allele frequency

# Estimating diversity from the SFS



number of
sites with $j$
derived alleles

number of
segregating sites

$$\hat{\theta}_W = \frac{S}{\sum_{j=1}^{2n-1} \frac{1}{j}} = \frac{\sum_{j=1}^{2n-1} \eta_j}{\sum_{j=1}^{2n-1} \frac{1}{j}}$$

Expected length
of a coalescent
tree of $n$ diploid
individuals. This
scales $S$ based
on sample size.

Number of sites

$+$ $+$ $+$ $+$ $+$ $+$ $= S$

1    2    3    4    5    6    7

Derived allele frequency

# Estimating diversity from the SFS



$$\hat{\pi} = \frac{\sum_{j=1}^{2n-1} \overbrace{j(2n-j)}\,\eta_j}{\binom{2n}{2}}$$

Number of pairwise differences when you have $j$ alleles of one type and $2n\text{-}j$ of another

Number of pairwise allele comparisons given $2n$ alleles

Number of sites

Derived allele frequency

# Estimating diversity from the SFS



$$\hat{\pi} = \frac{\sum_{j=1}^{2n-1} j(2n-j)\eta_j}{\binom{2n}{2}}$$

Example average number of pairwise differences comparing the following set of 8 alleles for one site:

$j = 3$ A alleles

$2n$-$j = 5$ a alleles:

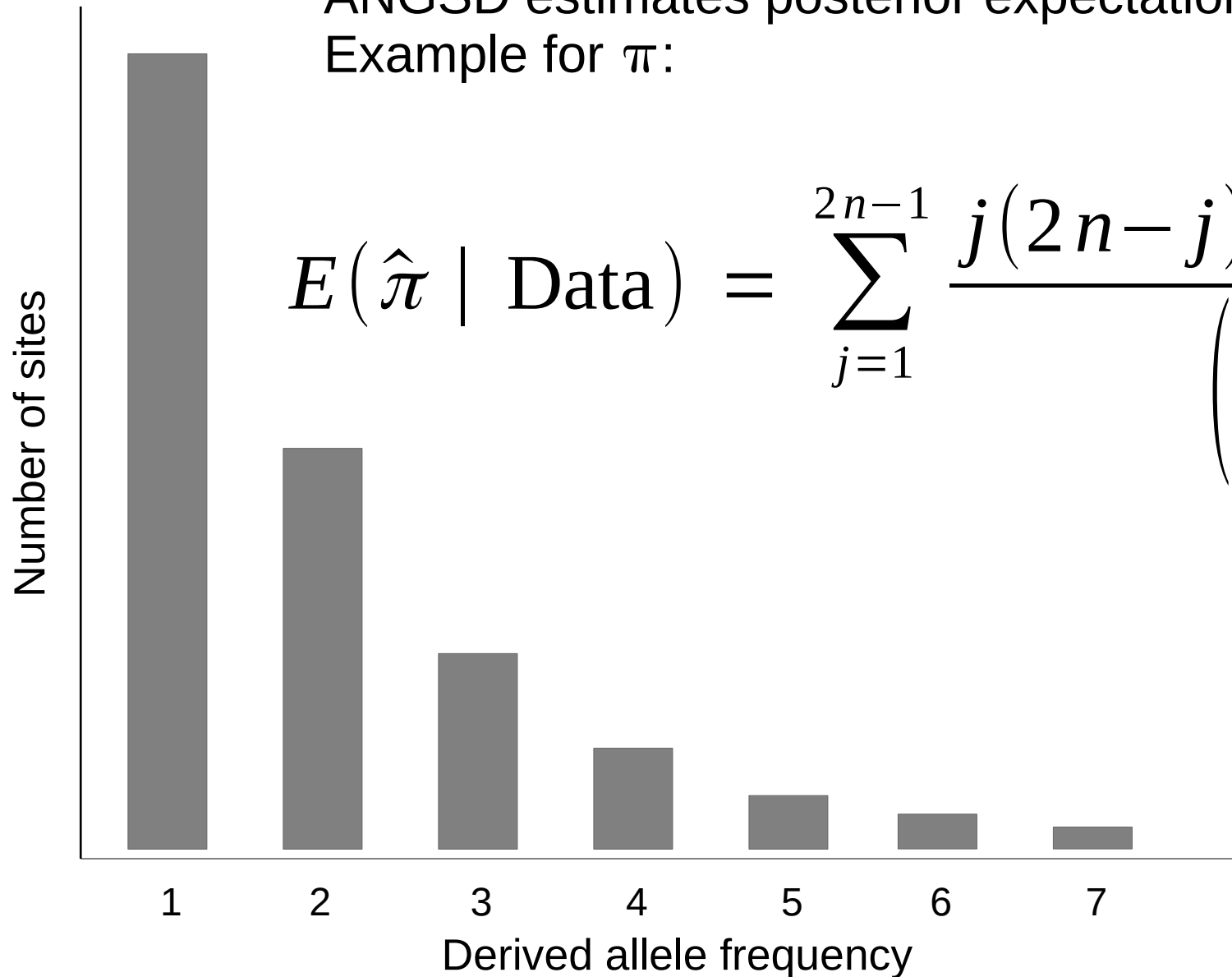A   A   A   a   a   a   a   a

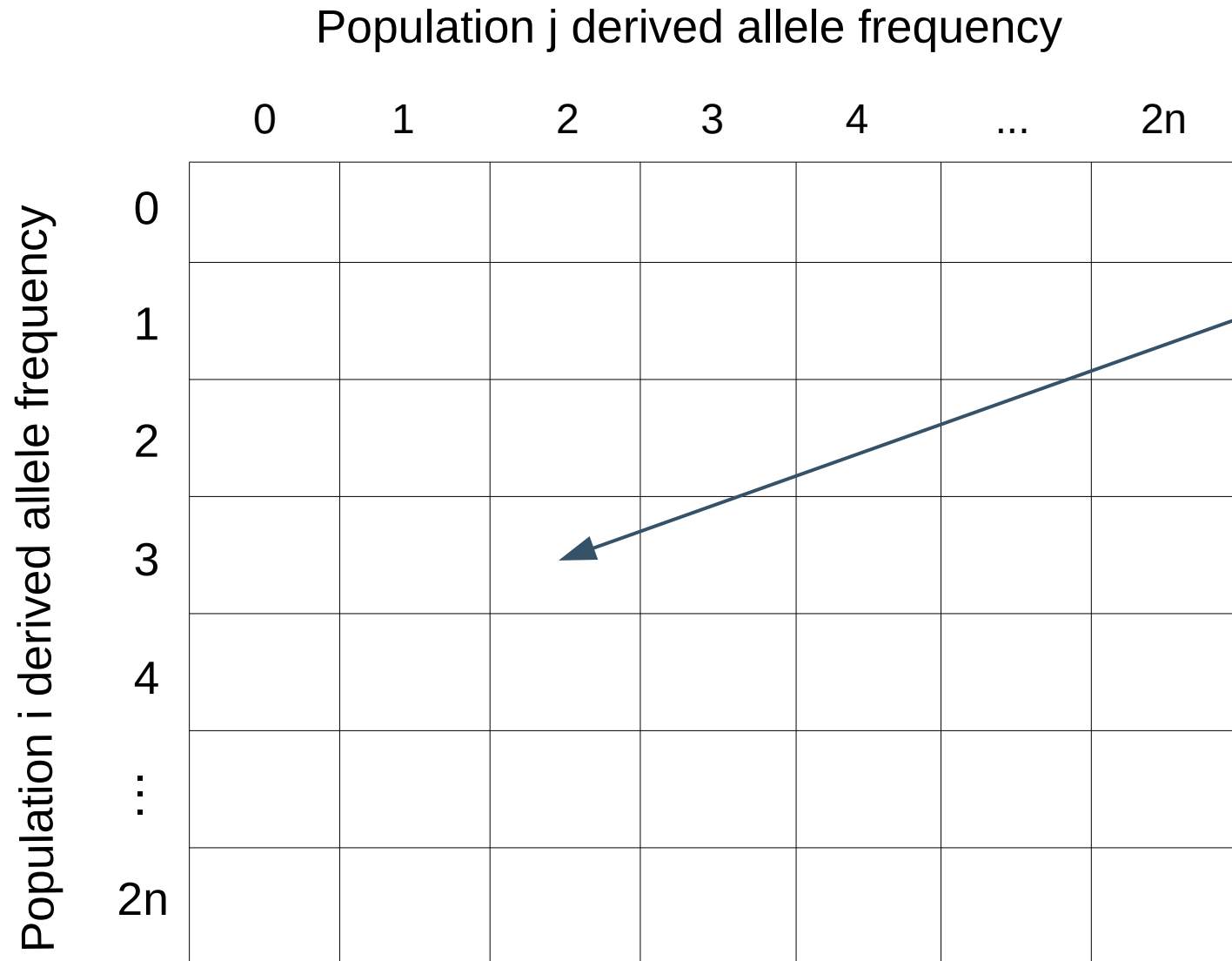$$\frac{5+5+5+0+0+0+0}{\binom{8}{2}} = \frac{15}{28}$$

Number of sites

Derived allele frequency

1     2     3     4     5     6     7

# Estimating diversity from the SFS

ANGSD estimates posterior expectations of θ.
Example for $\pi$:

$$E(\hat{\pi} \mid \text{Data}) = \sum_{j=1}^{2n-1} \frac{j(2n-j)\,\text{E}(\eta_j|\text{Data})}{\binom{2n}{2}}$$

Number of sites

Derived allele frequency

1   2   3   4   5   6   7

# 2-dimensional SFS

Population j derived allele frequency

|   | 0 | 1 | 2 | 3 | 4 | ... | 2n |
|---|---|---|---|---|---|-----|----|
| 0 |   |   |   |   |   |     |    |
| 1 |   |   |   |   |   |     |    |
| 2 |   |   |   |   |   |     |    |
| 3 |   |   |   |   |   |     |    |
| 4 |   |   |   |   |   |     |    |
| ... |  |   |   |   |   |     |    |
| 2n |  |   |   |   |   |     |    |

Population i derived allele frequency

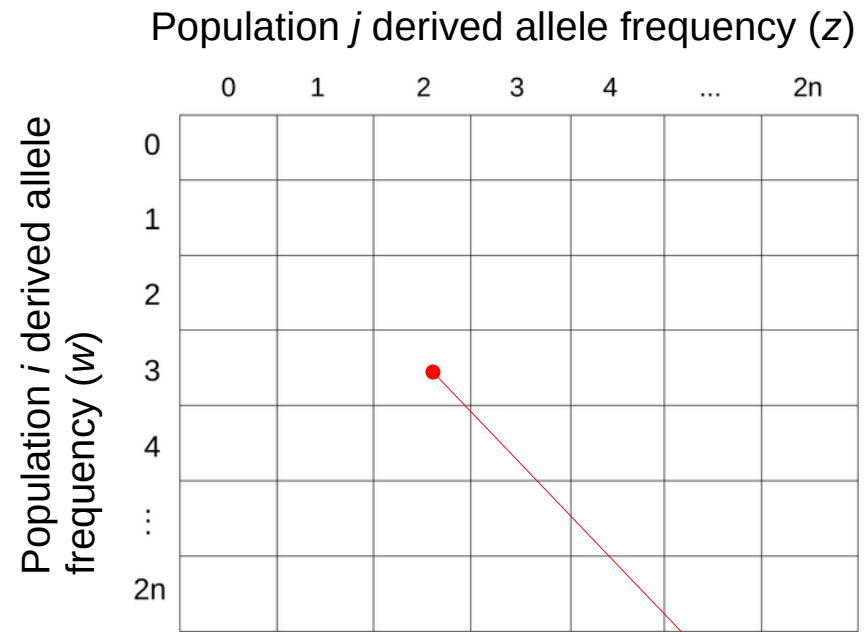Proportion of sites where:

pop *i* has 3 derived alleles

&

pop *j* has 2 derived alleles.

$$F_{ST} = \frac{E(a|X)}{E(a|X) + E(b|X)}$$

Genetic variance between populations

Genetic variance within populations

Population $j$ derived allele frequency ($z$)

Population $i$ derived allele frequency ($w$)

Probability of $k$ derived alleles in pop $i$ and $z$ derived alleles in pop $j$

Likelihood of $w$ derived alleles in pop $i$ at site $v$

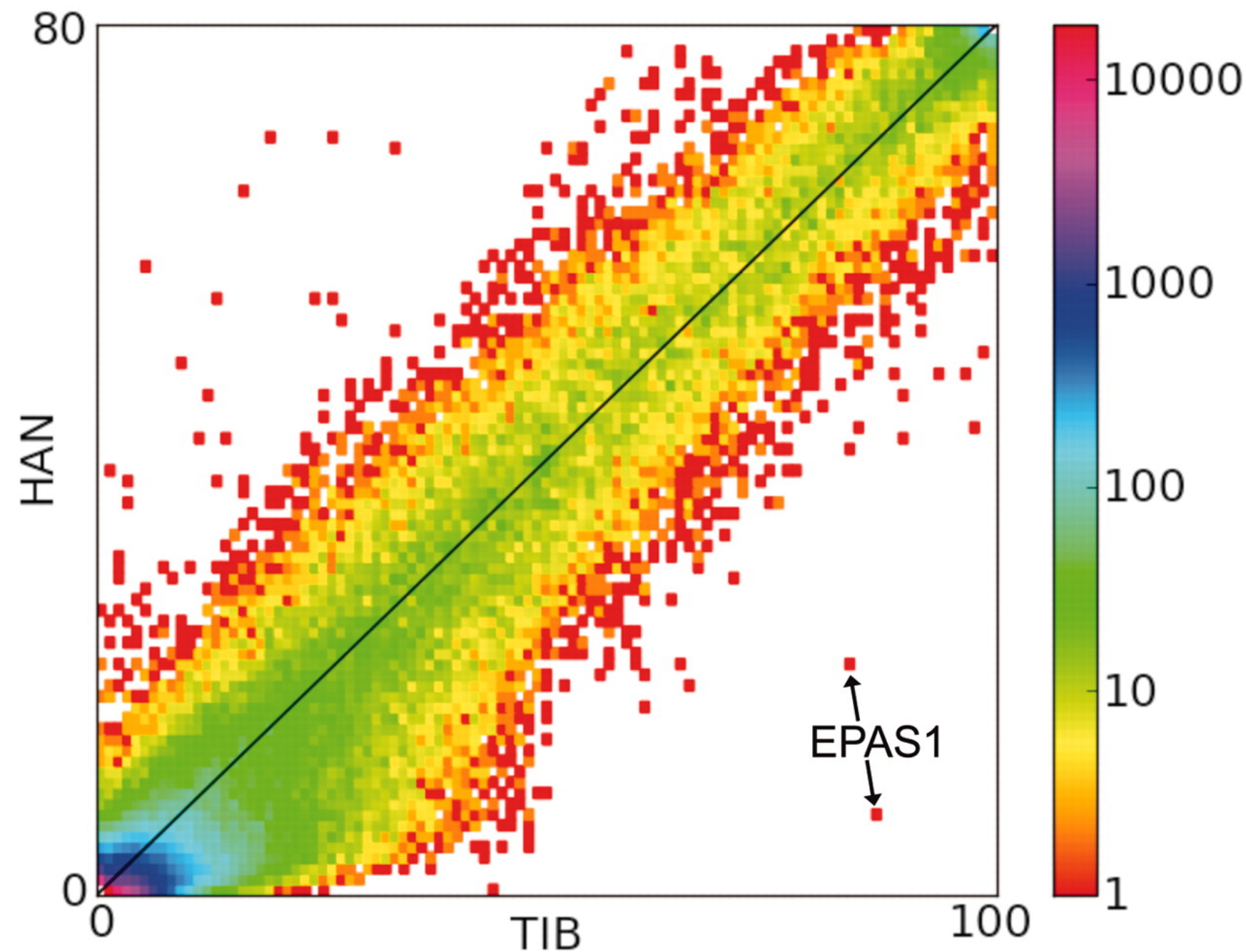Likelihood of $z$ derived alleles in pop $j$ at site $v$

$$E(a|X) = \sum_{k=0}^{2n} \sum_{z=0}^{2n} a_{pop\,i,\,pop\,j}^{w,z} \, P(X_{i,v}|K_{i,v}=w) \, P(X_{j,v}|K_{j,v}=z) \, Q_{i,j}^{w,z}$$

$$E(b|X) = \sum_{k=0}^{2n} \sum_{z=0}^{2n} b_{pop\,i,\,pop\,j}^{w,z} \, P(X_{i,v}|K_{i,v}=k) \, P(X_{j,v}|K_{j,v}=z) \, Q_{i,j}^{w,z}$$

# Identifying loci under selection



Yi *et al.* (2010)

# Exercise. Estimating the SFS and summary statistics