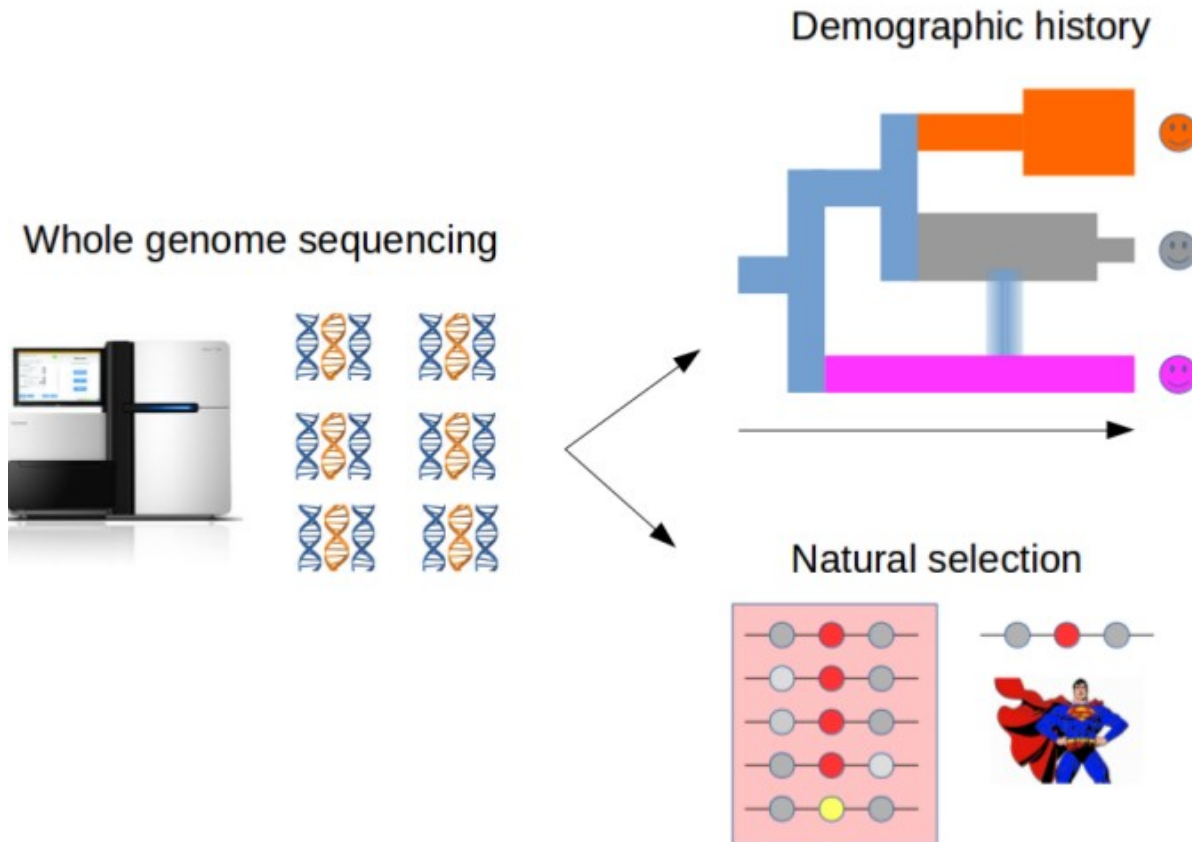


Detecting natural selection from genomic data (part 2)

Matteo Fumagalli

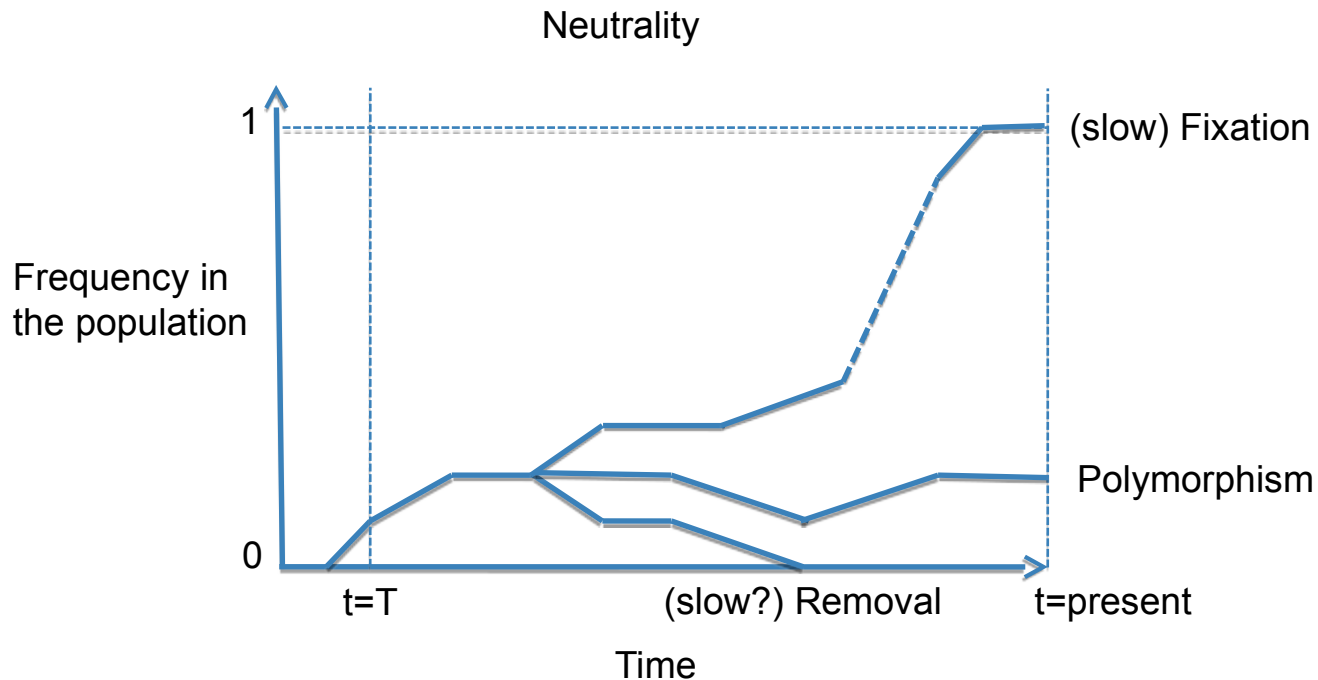
Motivation



Natural selection

Heritable traits that increase the fitness of the become more common.

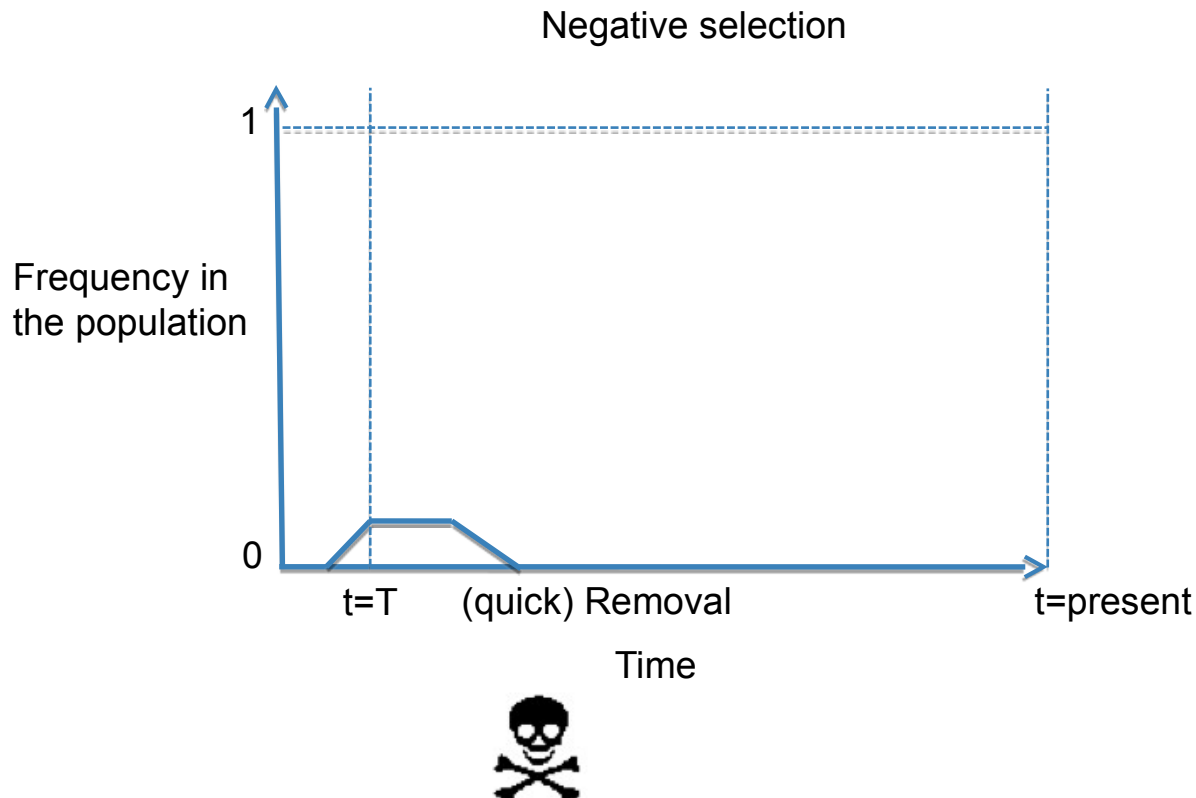
- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier



Natural selection

Heritable traits that increase the fitness of the become more common.

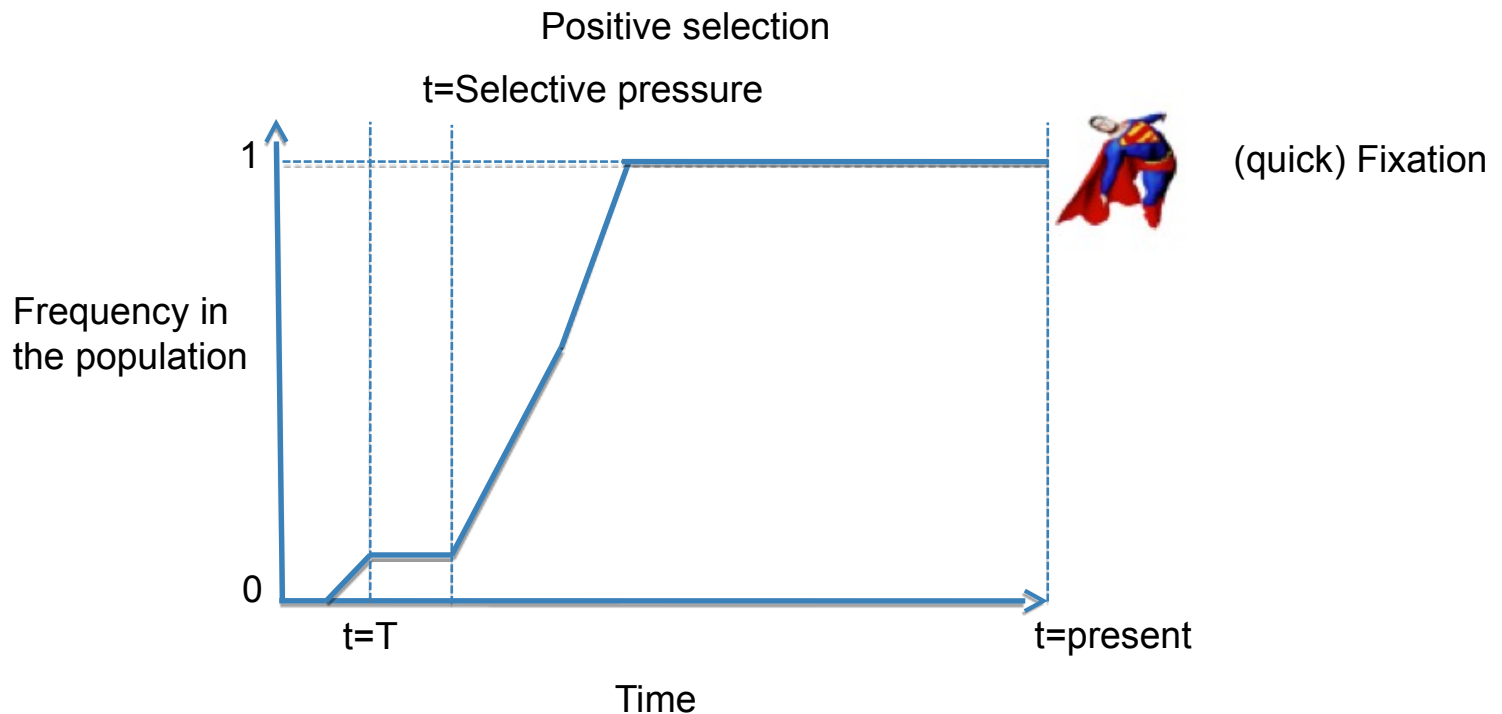
- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier



Natural selection

Heritable traits that increase the fitness of the become more common.

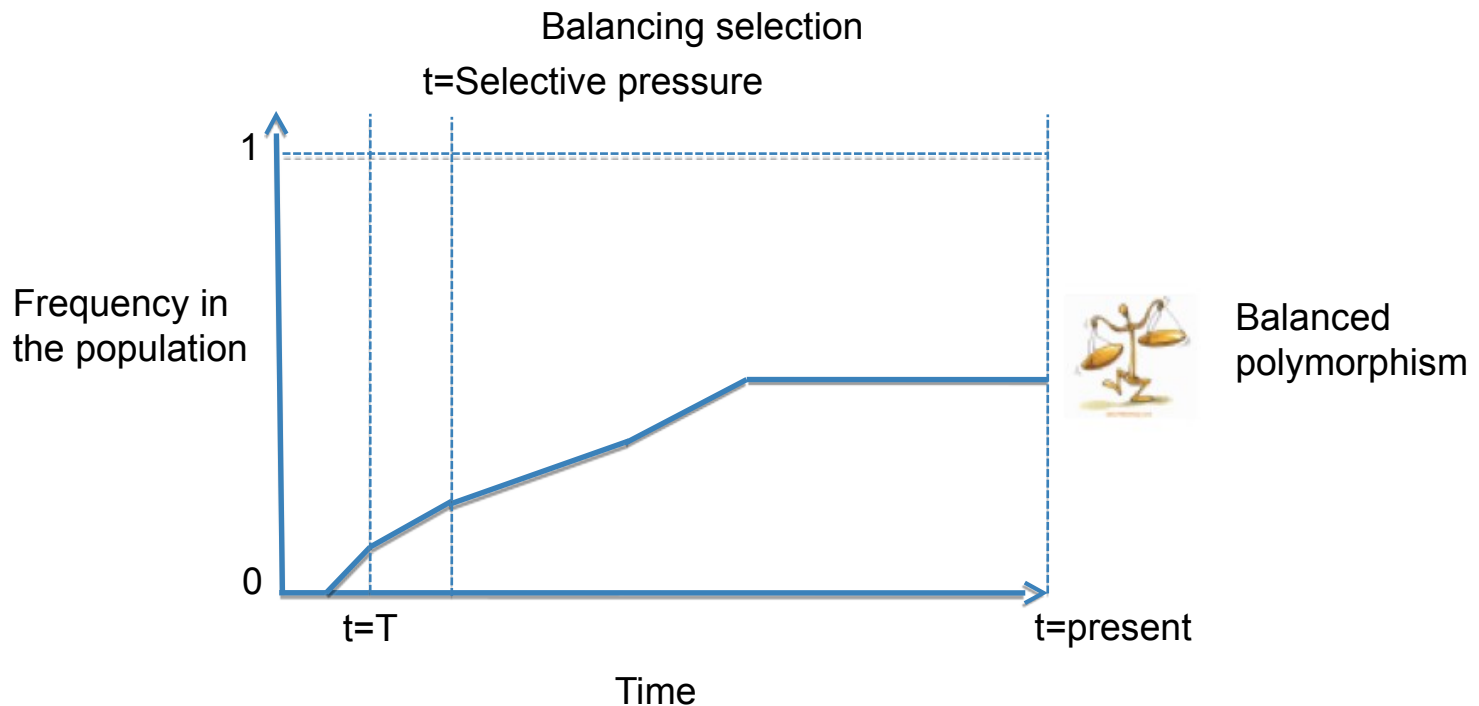
- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier



Natural selection

Heritable traits that increase the fitness of the become more common.

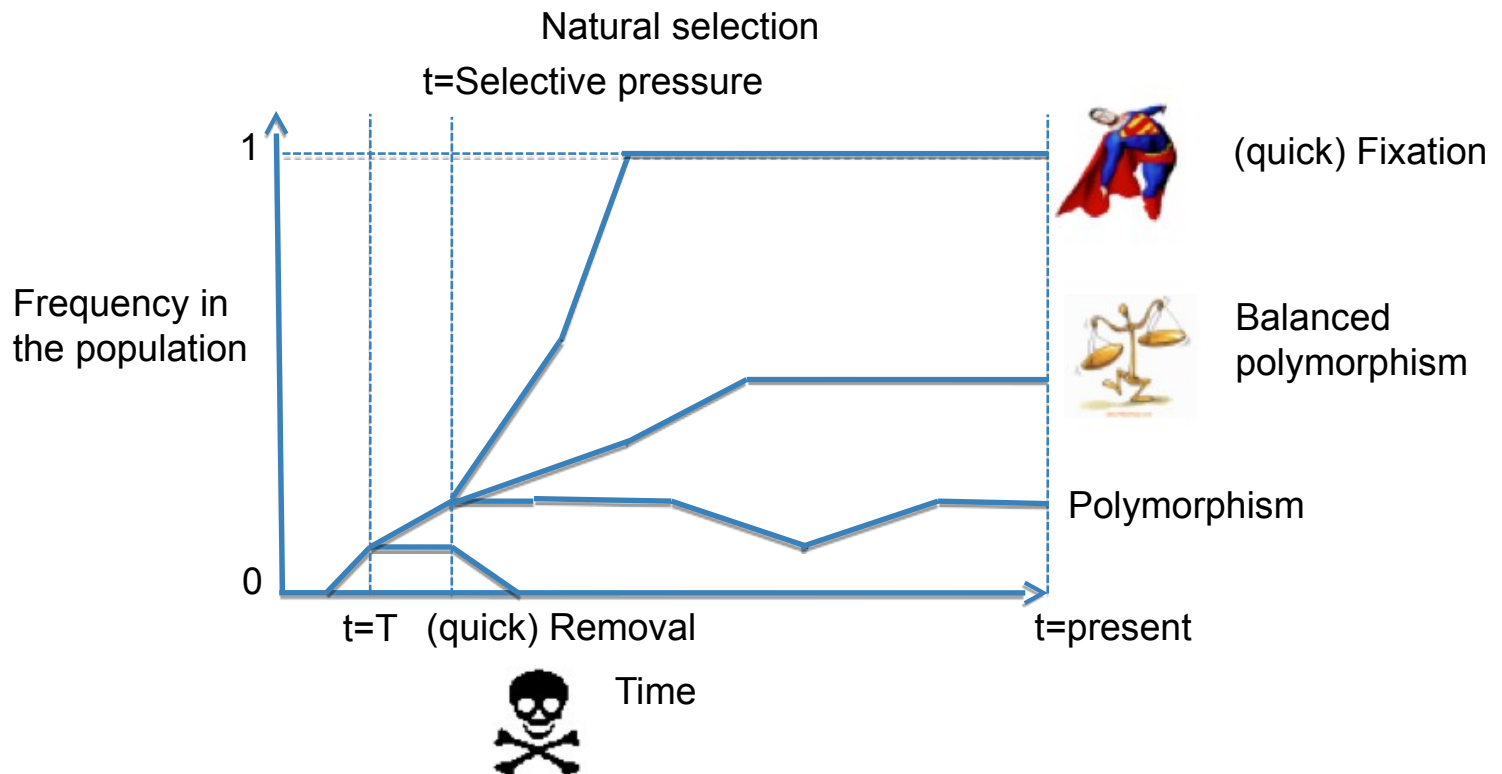
- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier



Natural selection

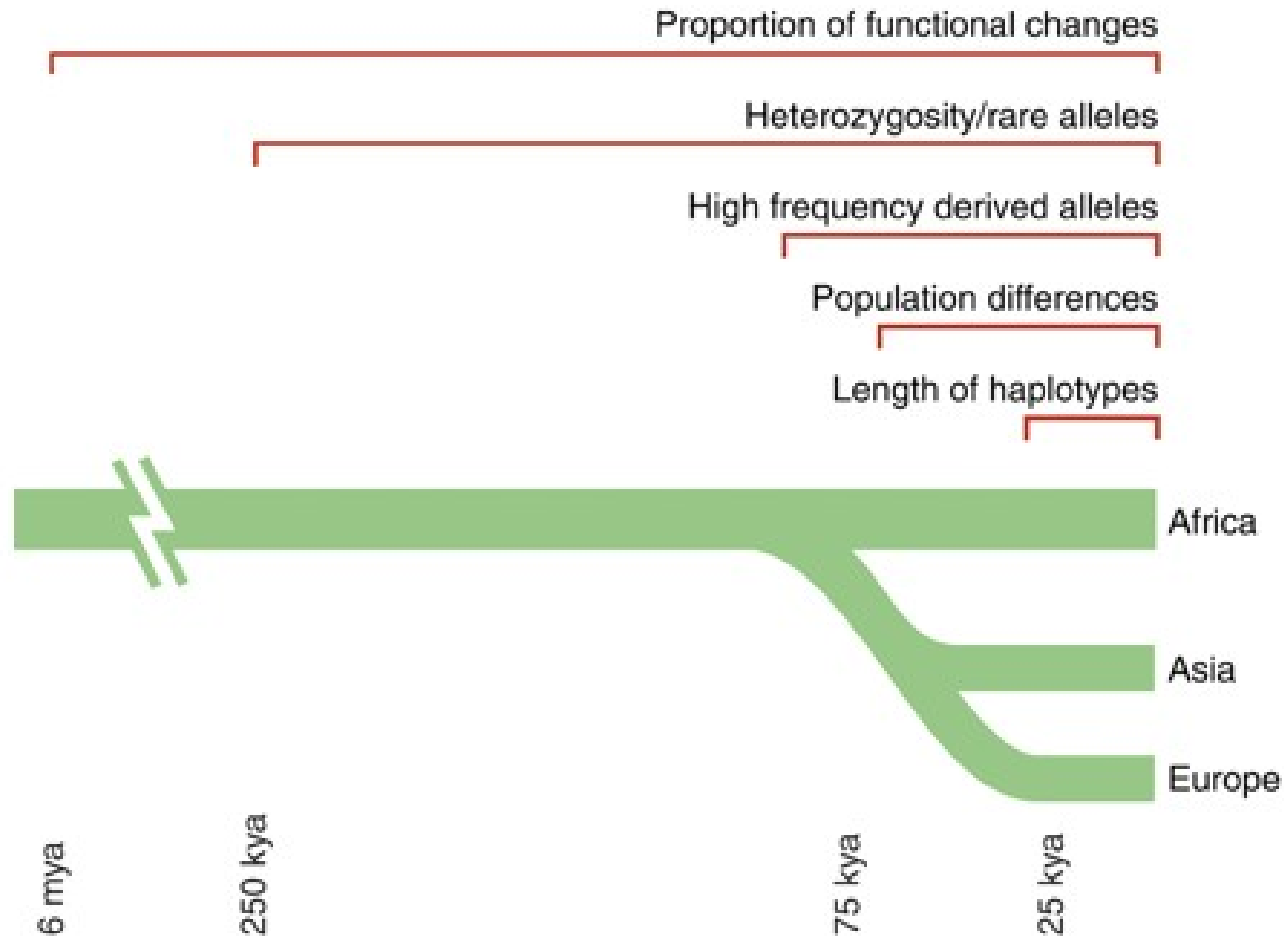
Heritable traits that increase the fitness of the become more common.

- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier

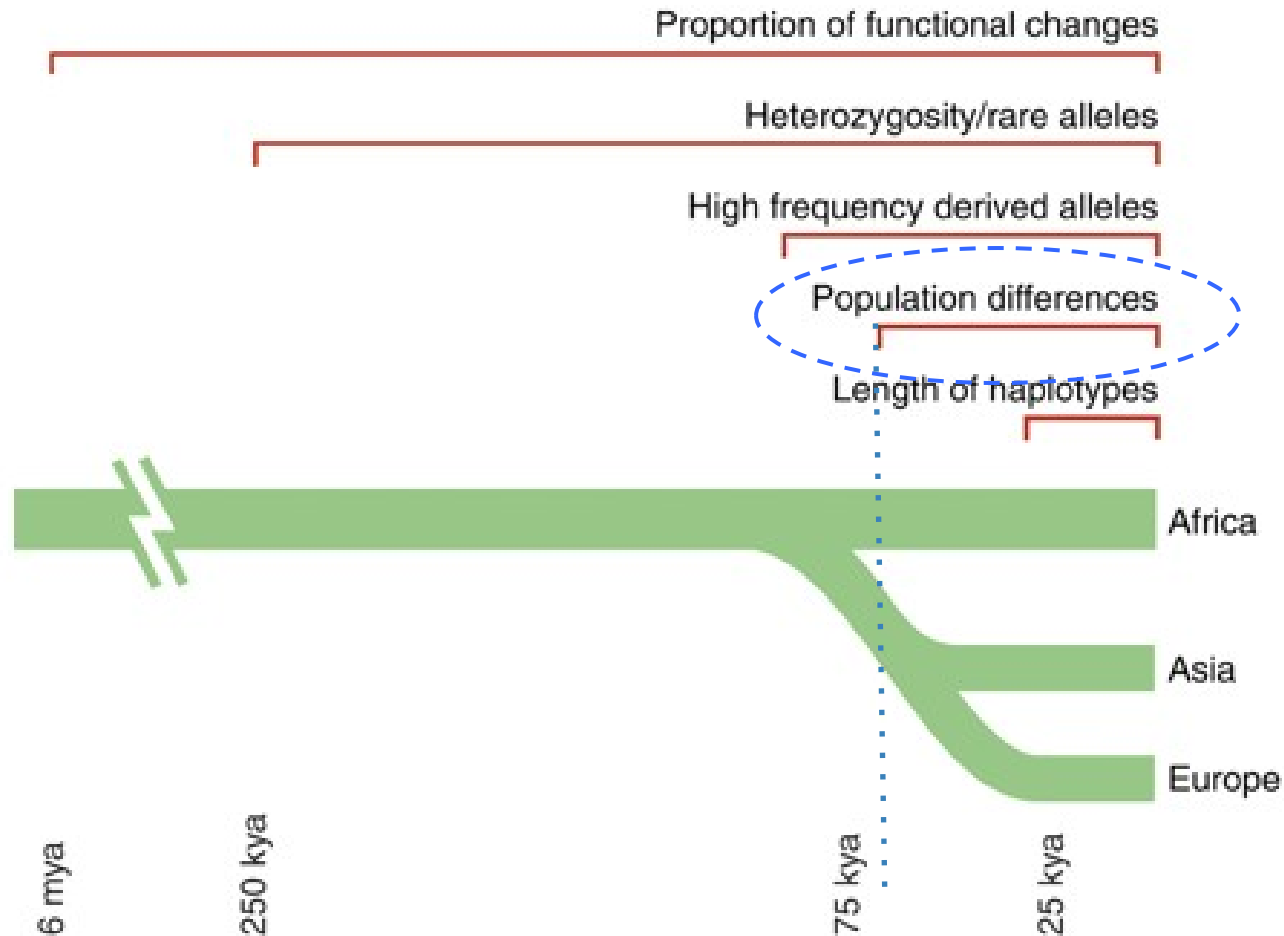


- 2) Sites targeted by natural selection are likely to harbour **functionality**

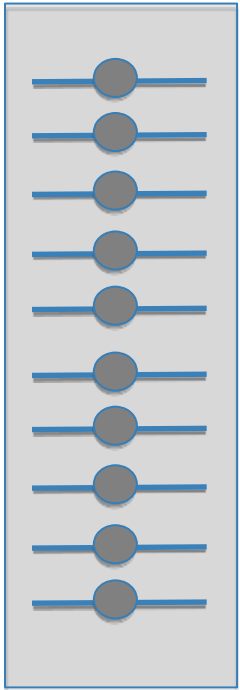
Methods to infer recent selection



Methods to infer recent selection

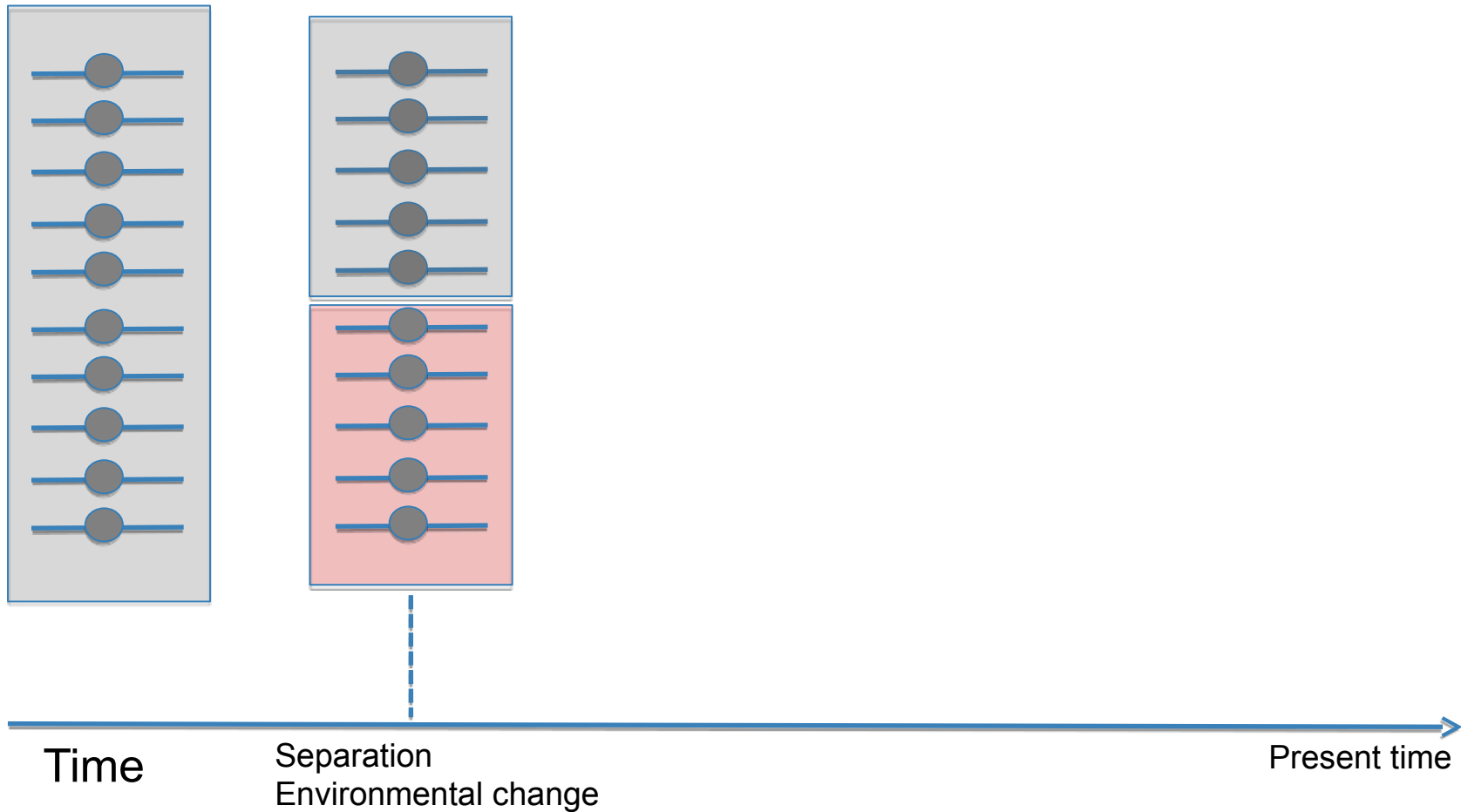


Allele frequency differentiation

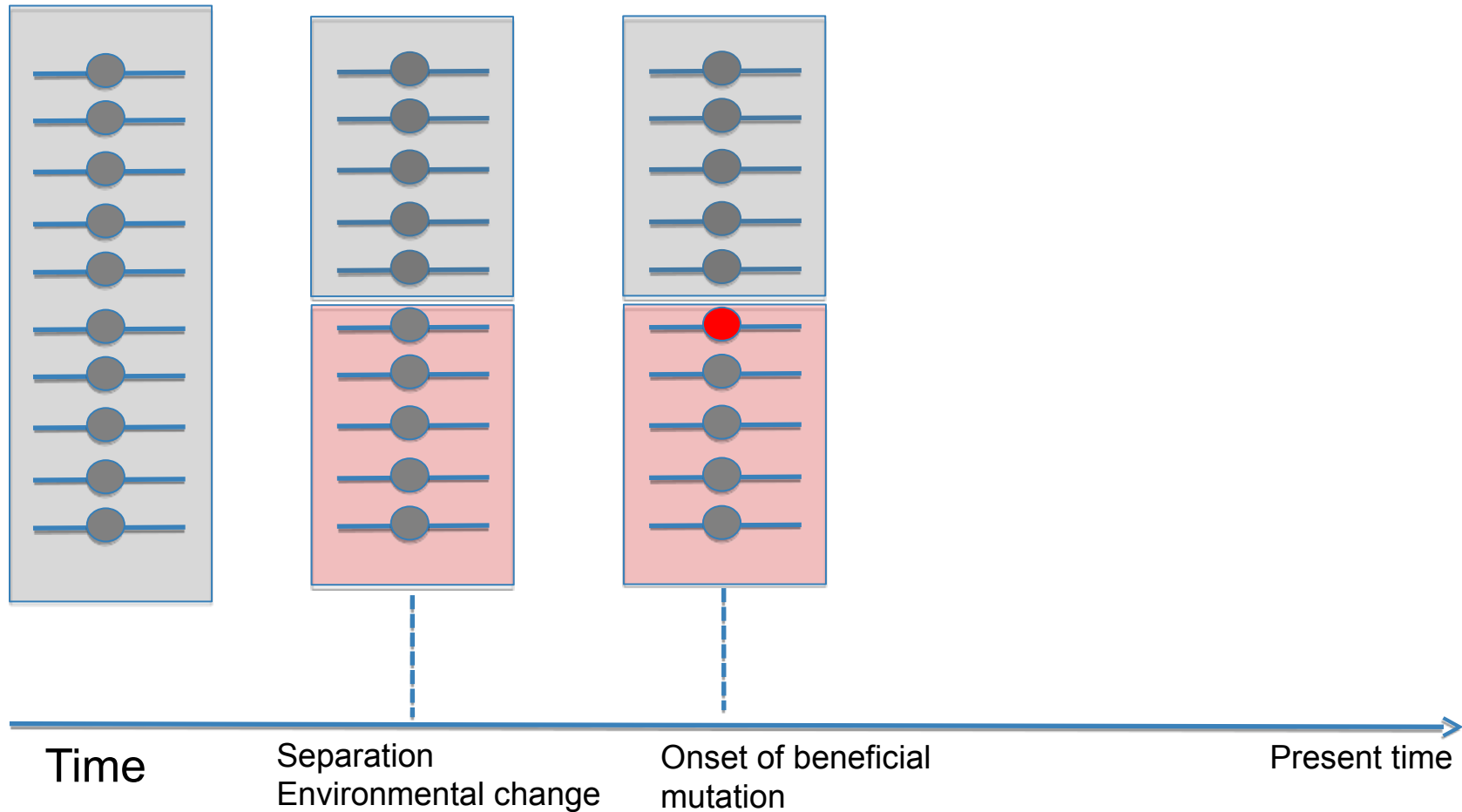


Time Present time

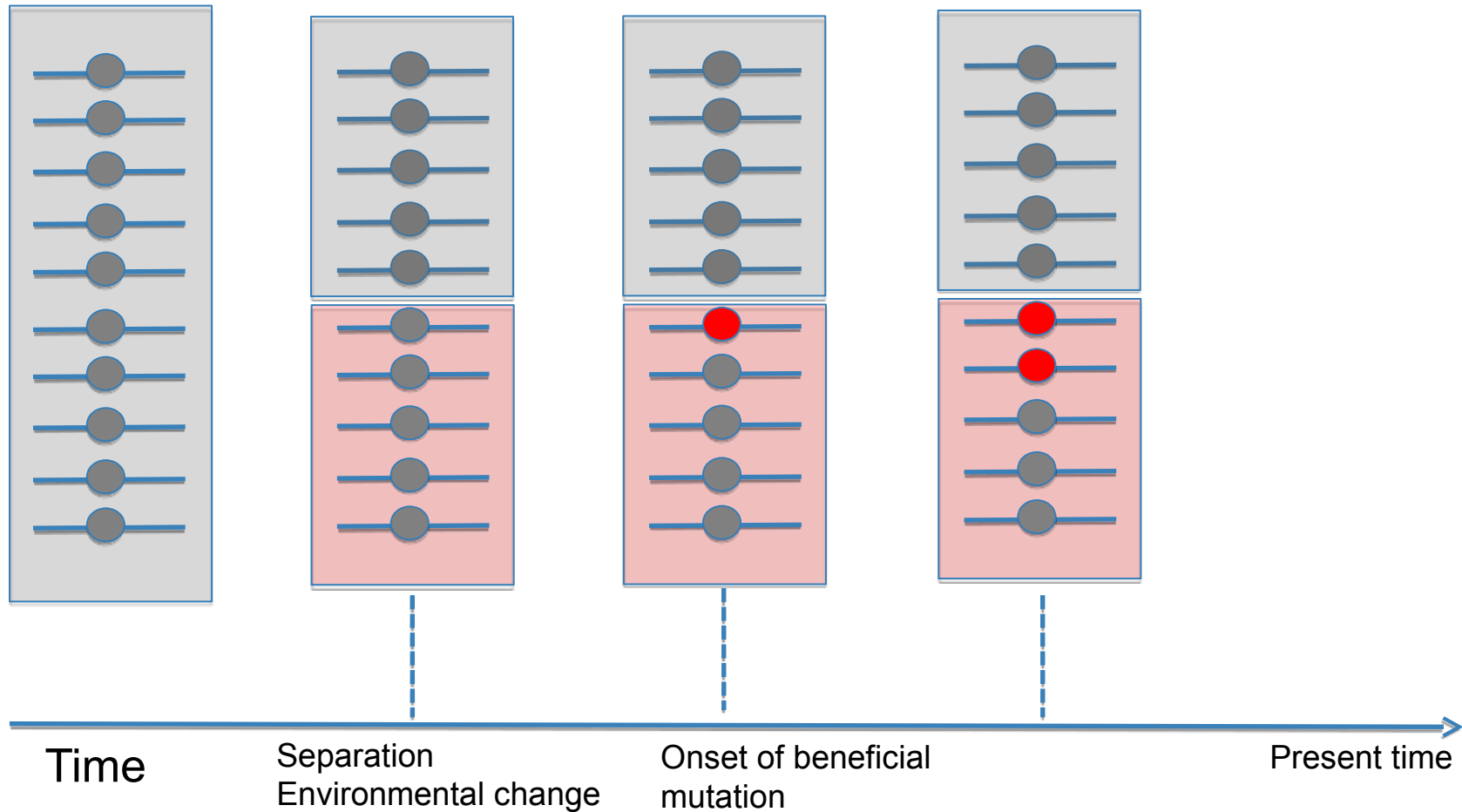
Allele frequency differentiation



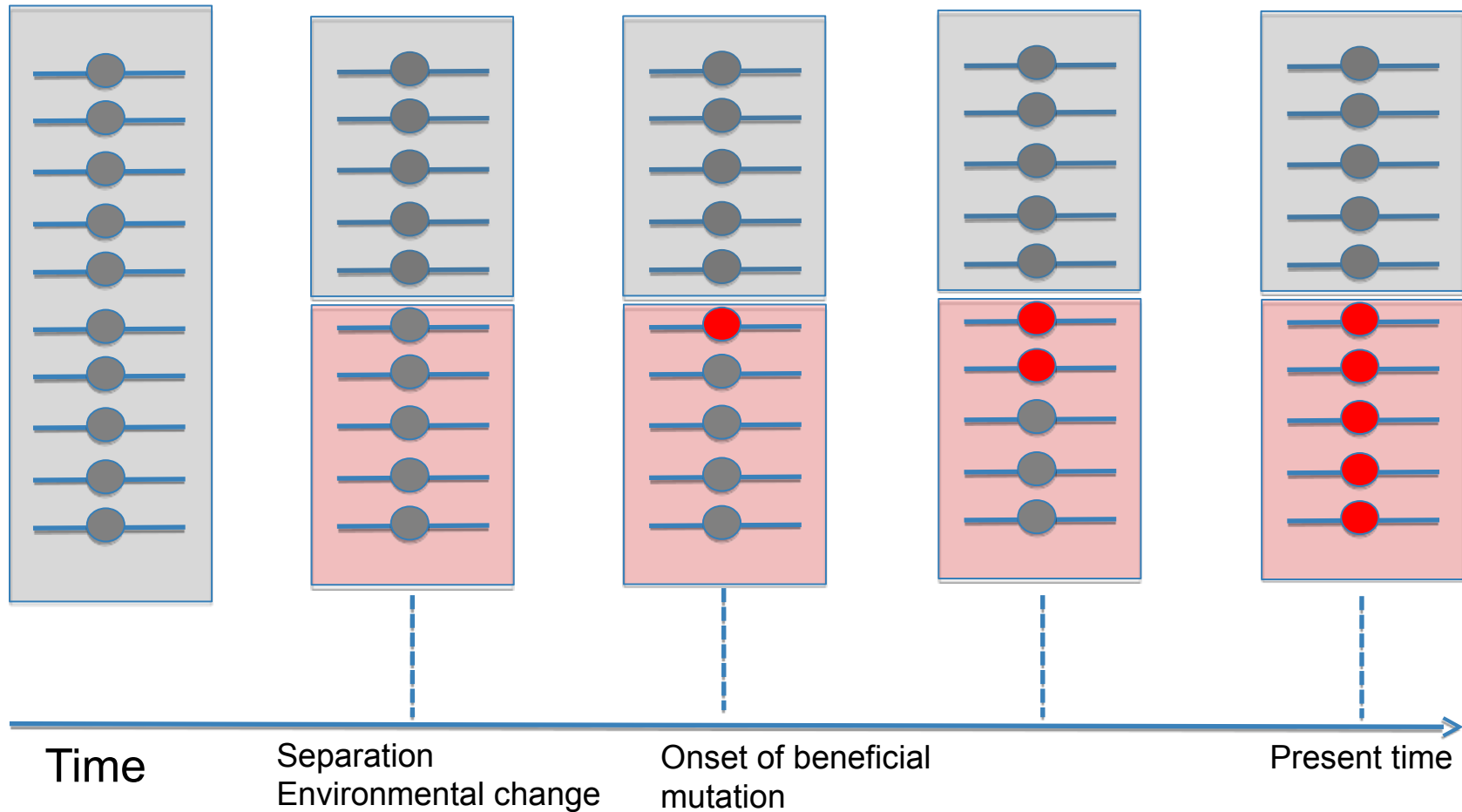
Allele frequency differentiation



Allele frequency differentiation



Allele frequency differentiation



$$F_{ST}$$

Common measure for quantifying population subdivision.

$$F_{ST} = H_B / (H_W + H_B)$$

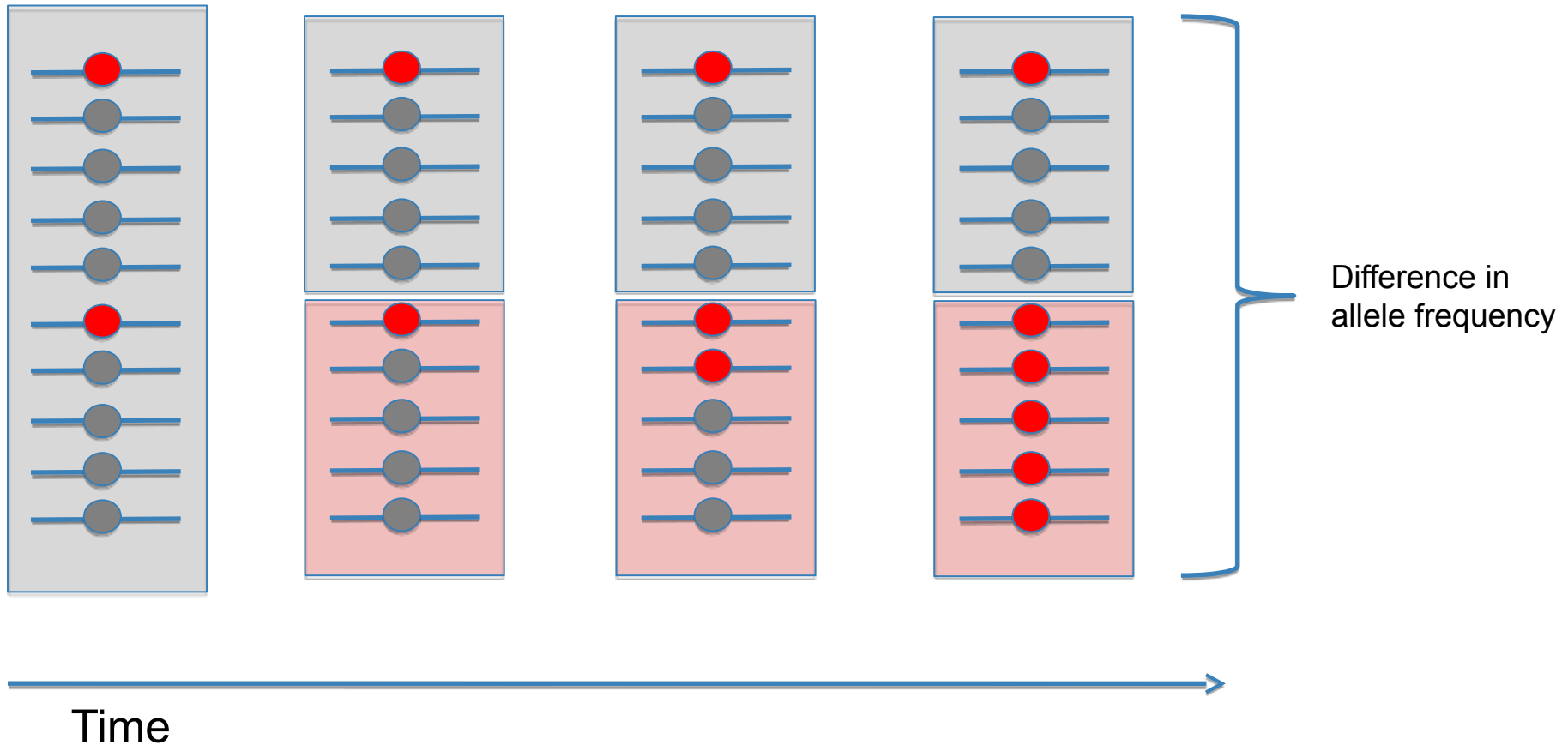
H_B : between populations

H_W : average within populations

- if $H_W \ll H_B$ then $F_{ST} \sim 1$
- if $H_B = 0$ then $F_{ST} = 0$

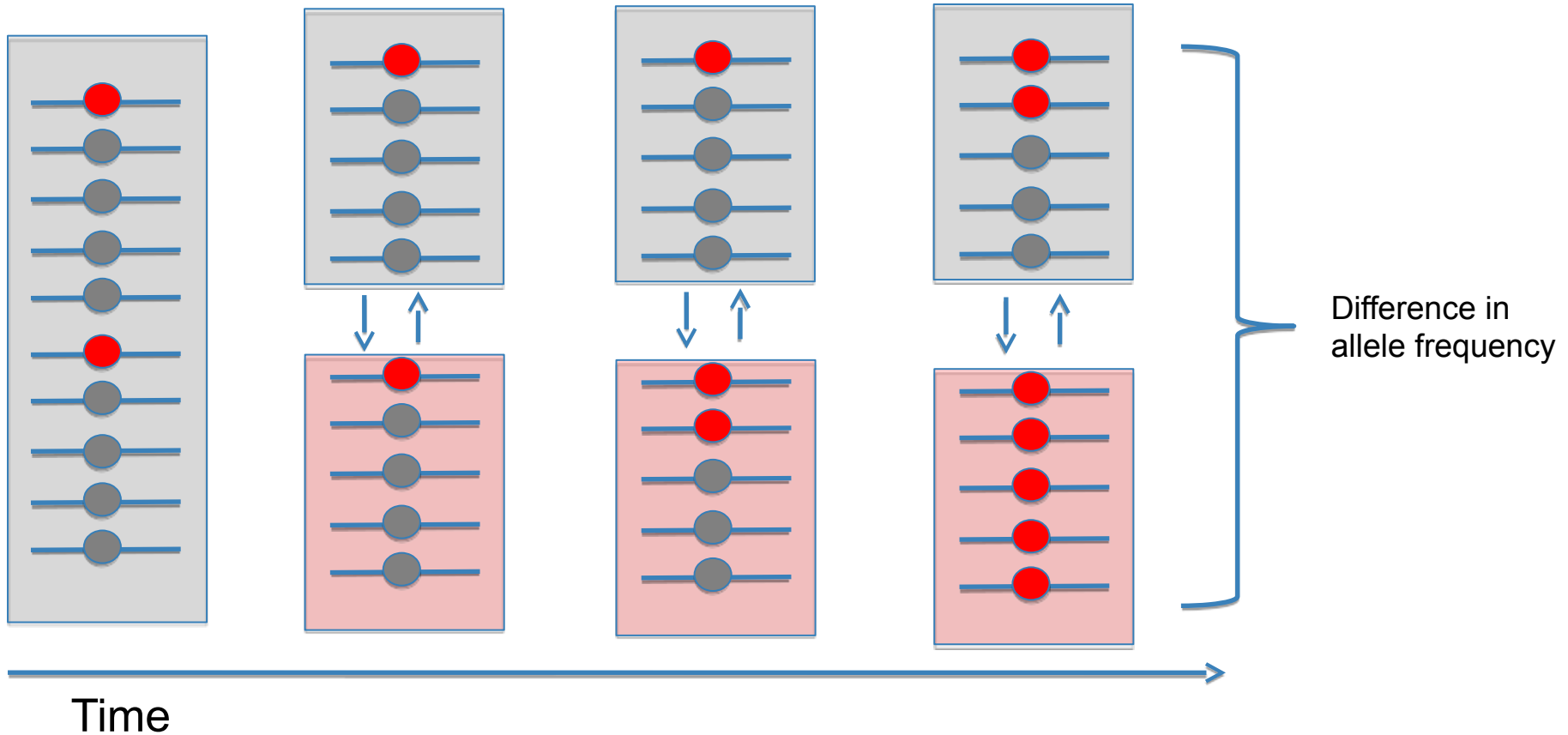
Allele frequency differentiation

From standing variation



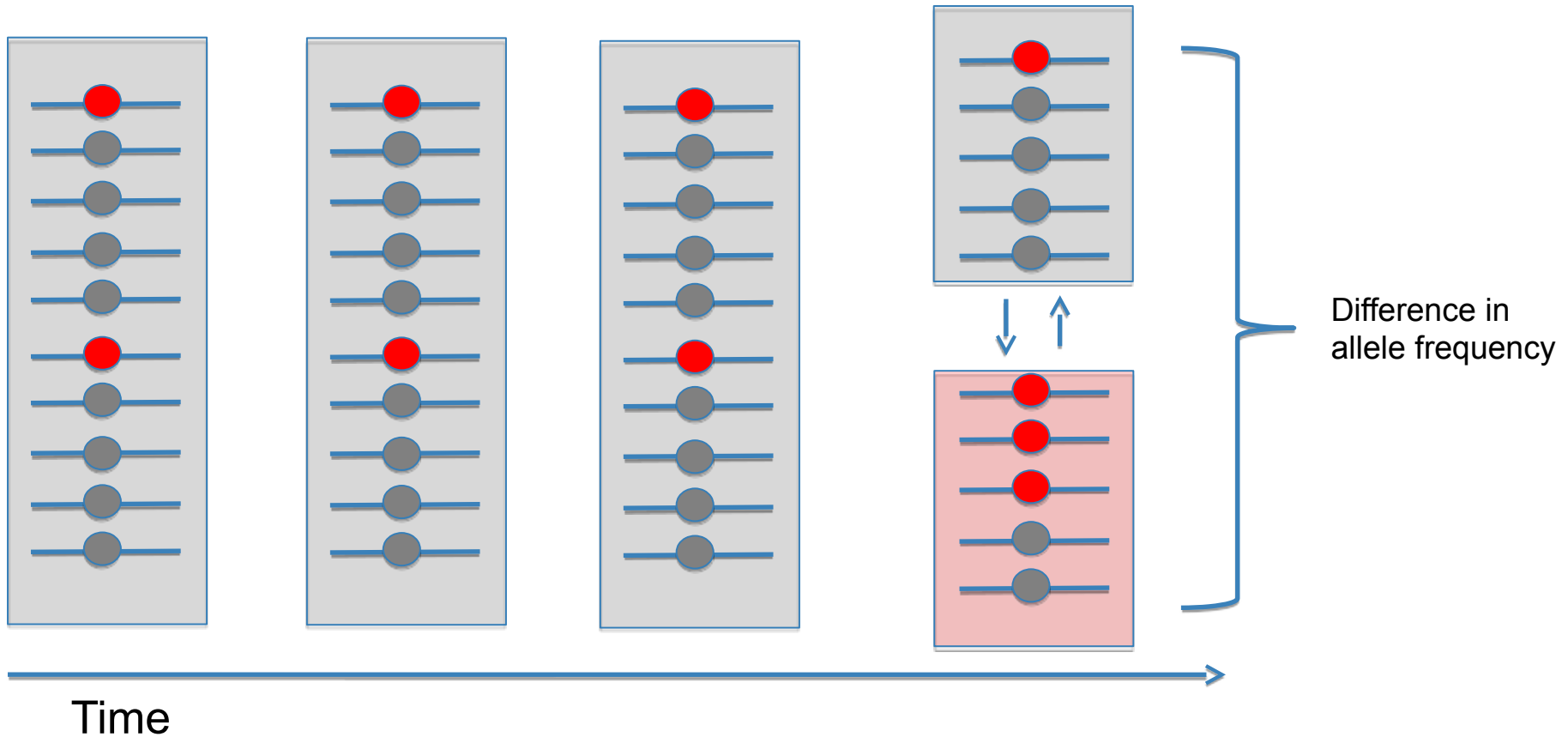
Allele frequency differentiation

With migration



Allele frequency differentiation

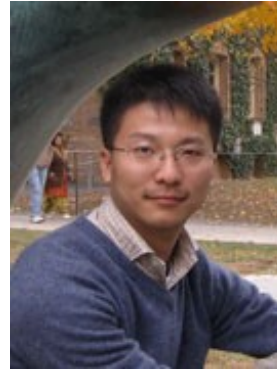
With recent divergence



Population genetic differentiation



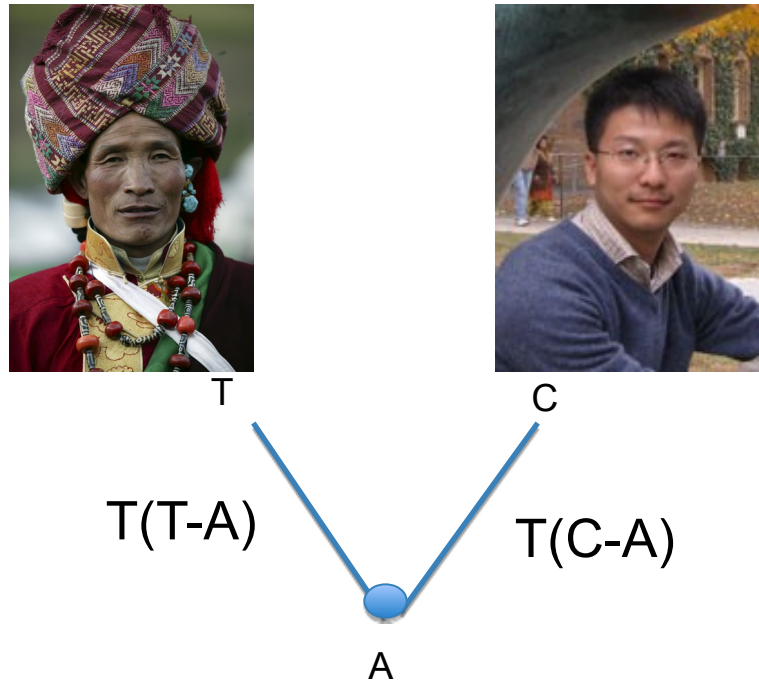
T



C

$$F_{ST}(T-C)$$

Population genetic differentiation



$$F_{ST}(T-C) \sim T(T-A-C)$$

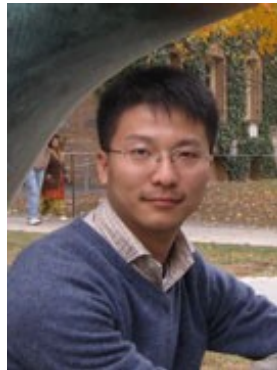
Population genetic differentiation

$$F_{ST}(T-C) \sim T(T-A-C)$$



T

$T(T-A)$



C

$T(C-A)$

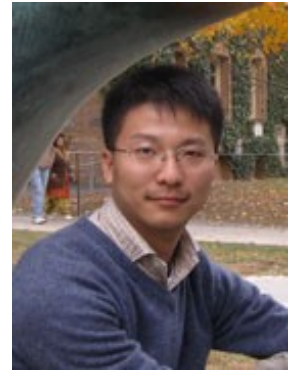
A

?



T

$T(T-A)$



C

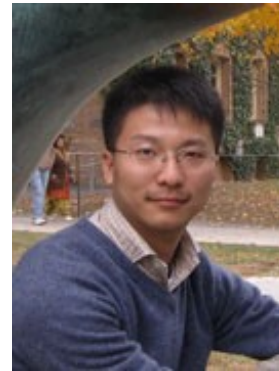
$T(C-A)$

A

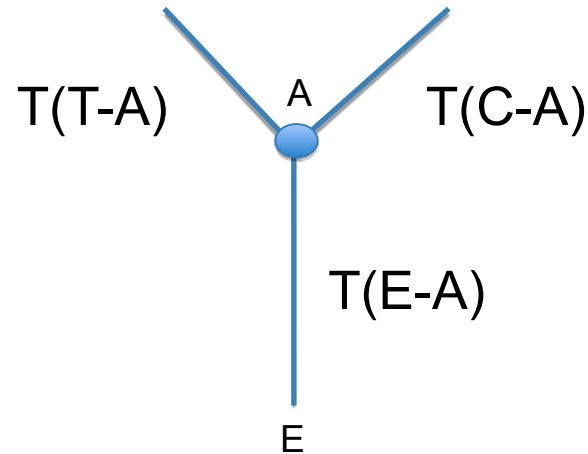
Population genetic differentiation



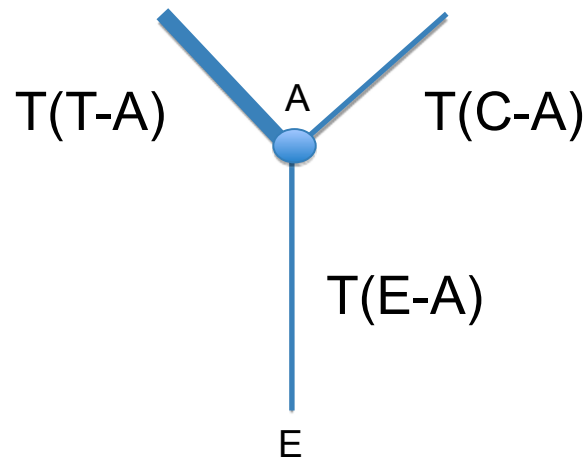
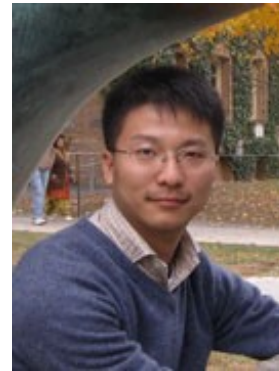
T



C



Population genetic differentiation

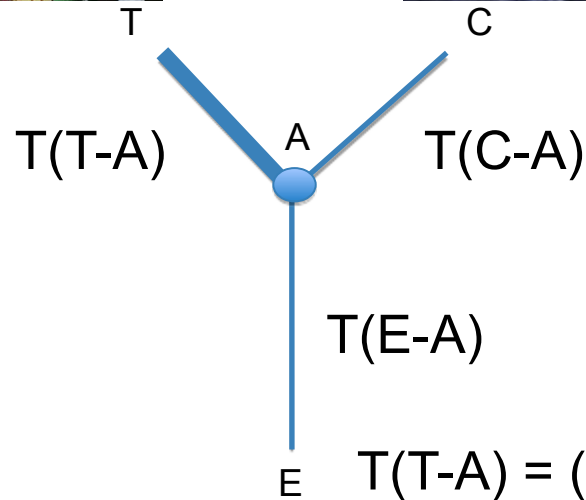
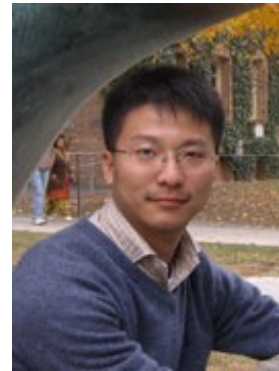


$$T(T-A-C) = -\log(1 - F_{ST}(T-C))$$

T(T-A)?



Population genetic differentiation

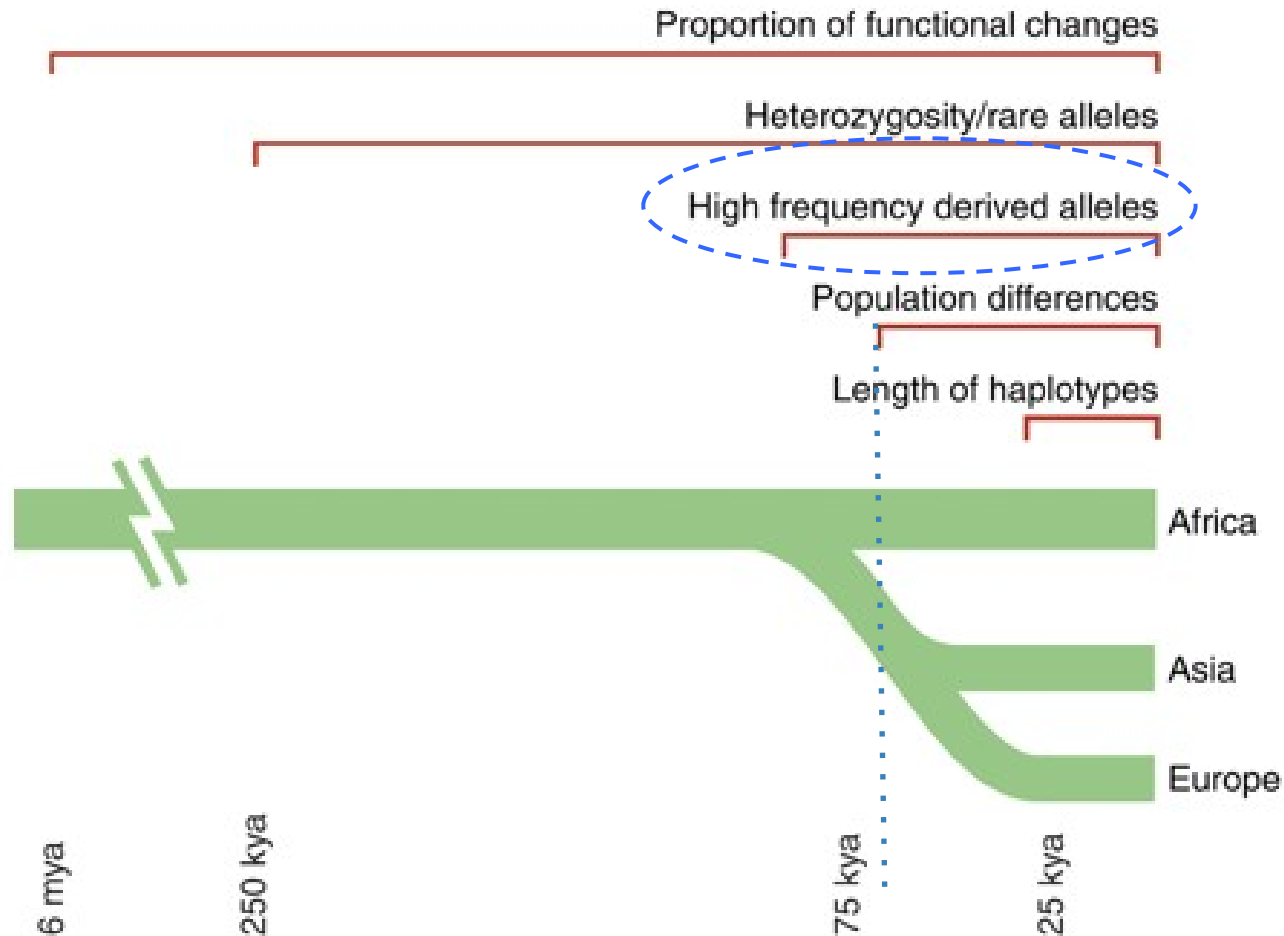


$$T(T-A-C) = -\log(1 - F_{ST}(T-C))$$

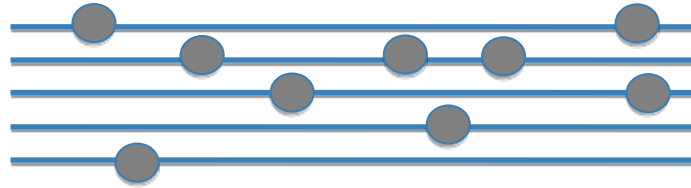
$$T(T-A) = (T(T-E) + T(C-T) - T(C-E)) / 2$$



Methods to infer selection

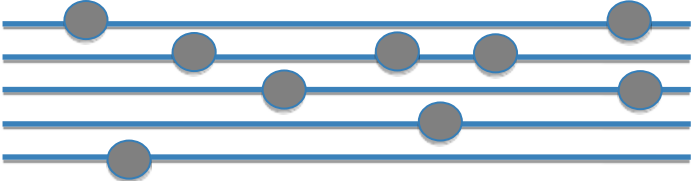
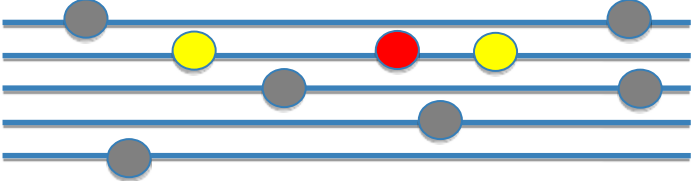


Positive selection: effect on haplotypes

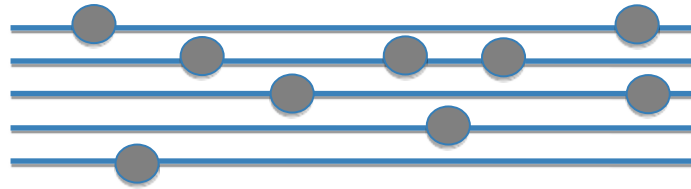


$t < T_{\text{sel}}$

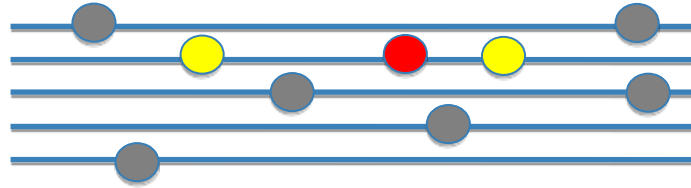
Positive selection: effect on haplotypes

 $t < T_{sel}$ 
$$t = T_{sel}$$

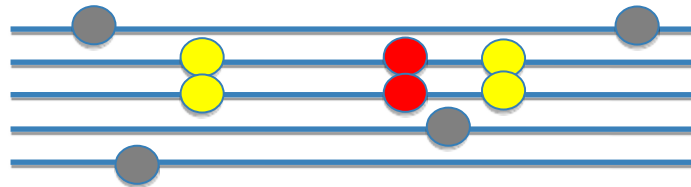
Positive selection: effect on haplotypes



$t < T_{\text{sel}}$

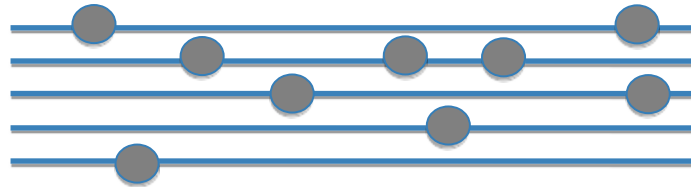


$t = T_{\text{sel}}$

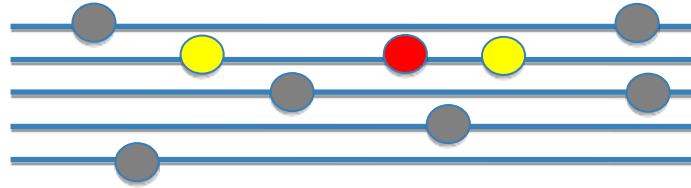


$t > T_{\text{sel}}$

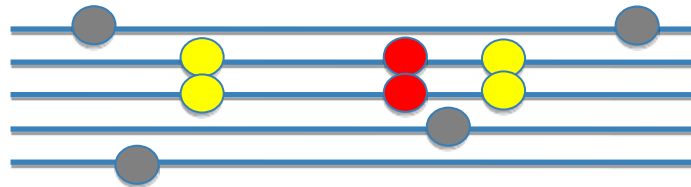
Positive selection: effect on haplotypes



$t < T_{\text{sel}}$



$t = T_{\text{sel}}$



$t > T_{\text{sel}}$



$t \gg T_{\text{sel}}$

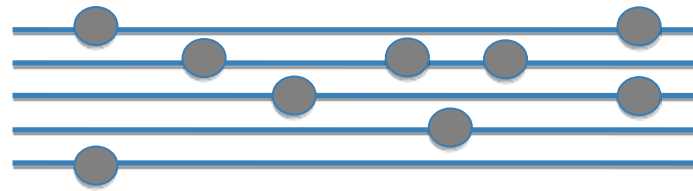
Selective sweep



Genetic hitch-hiking

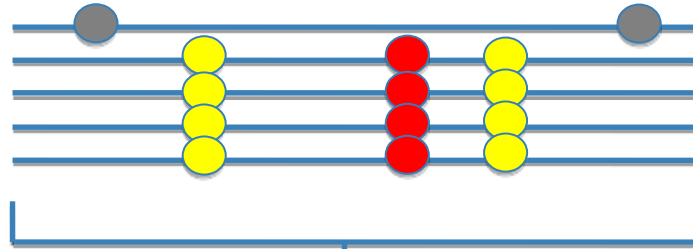


Positive selection



$t < T_{\text{sel}}$

...

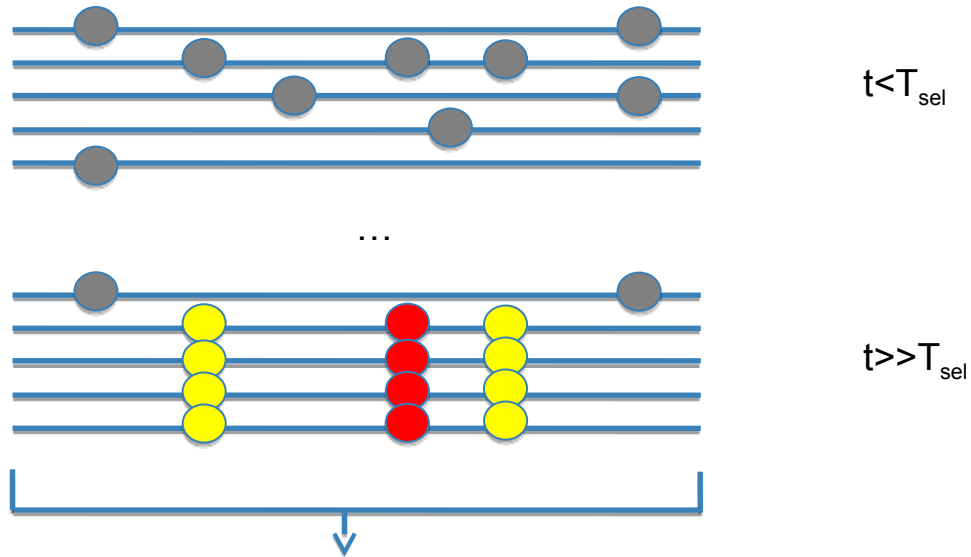


$t > T_{\text{sel}}$



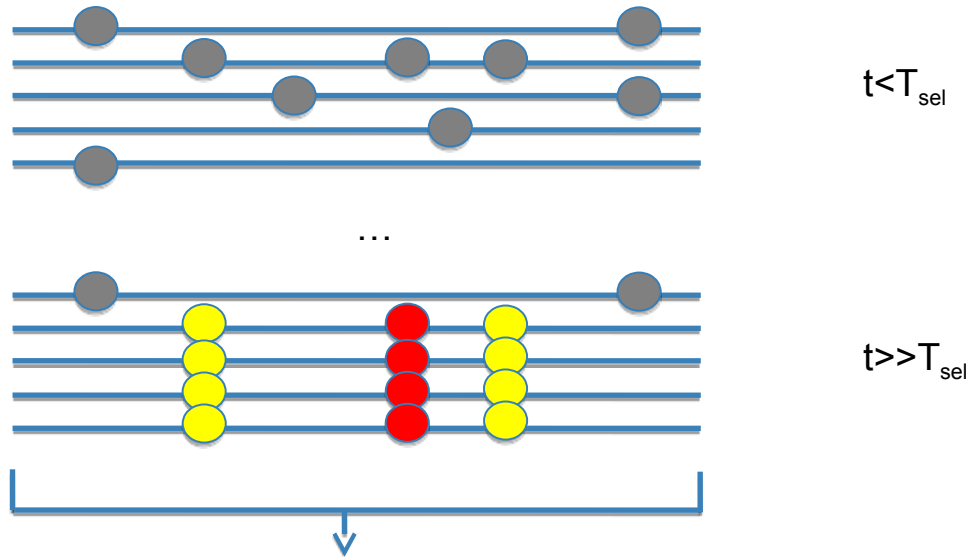
?

Positive selection



- Reduction of polymorphisms levels
(e.g. from 7 to 5 SNPs)

Positive selection



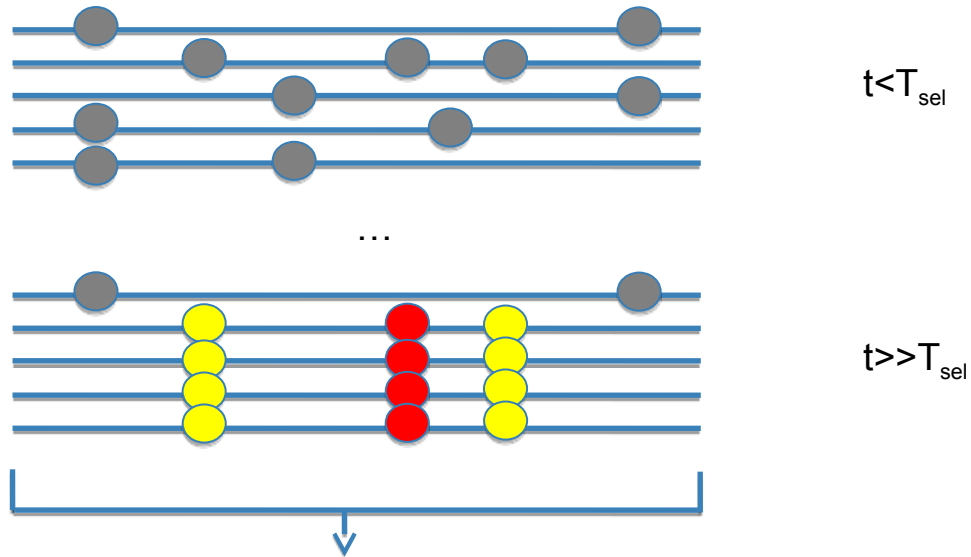
- Reduction of polymorphisms levels
(e.g. from 7 to 5 SNPs)

Nucleotide diversity index: Watterson's Theta
with K SNPs and n chromosomes

$$\theta_w = \frac{K}{a_n}$$

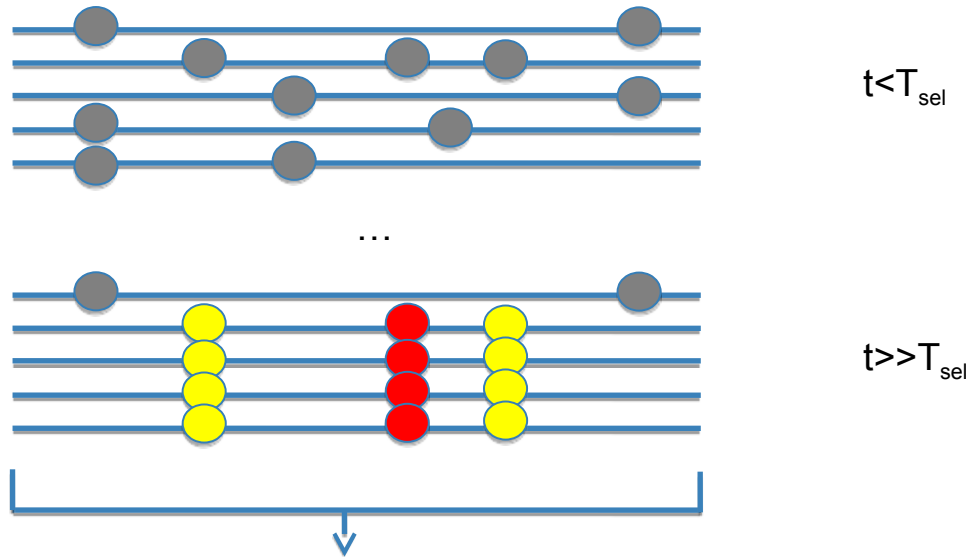
$$a_n = \sum_{i=1}^{n-1} \frac{1}{i}$$

Positive selection



- Reduction of polymorphisms levels (Theta)
- ?

Positive selection

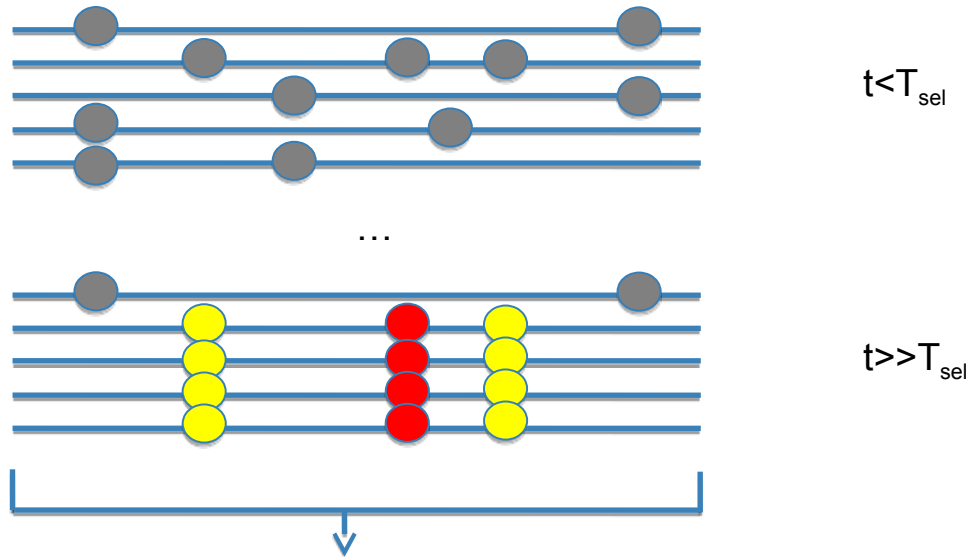


- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants

Nucleotide diversity index: average pairwise nucleotide differences (π)
with $k_{i,j}$ equal to the number of nucleotide differences between sequences i and j

$$\pi = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n k_{i,j}}{\binom{n}{2}}$$

Positive selection



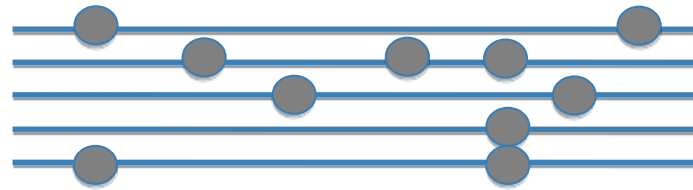
- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants (Pi)

Under neutrality, Theta and Pi are expected to be the same.
Tajima's D measures their difference.

$$D = \frac{\pi - \theta_w}{\sqrt{\hat{V}(\pi - \theta_w)}}$$

$D < 0$ is suggestive of an excess of low-frequency variants

The Site Frequency Spectrum

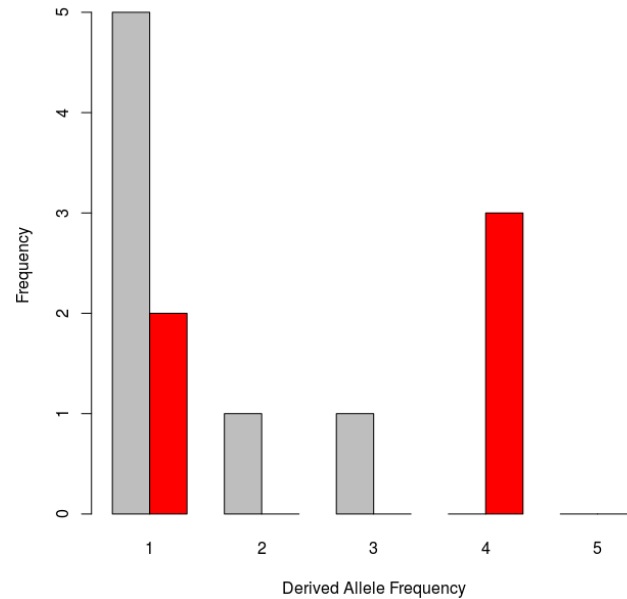


$t < T_{\text{sel}}$

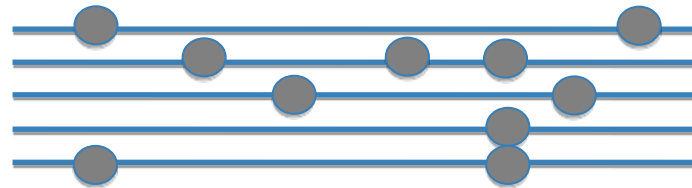
...



$t \gg T_{\text{sel}}$



The Site Frequency Spectrum

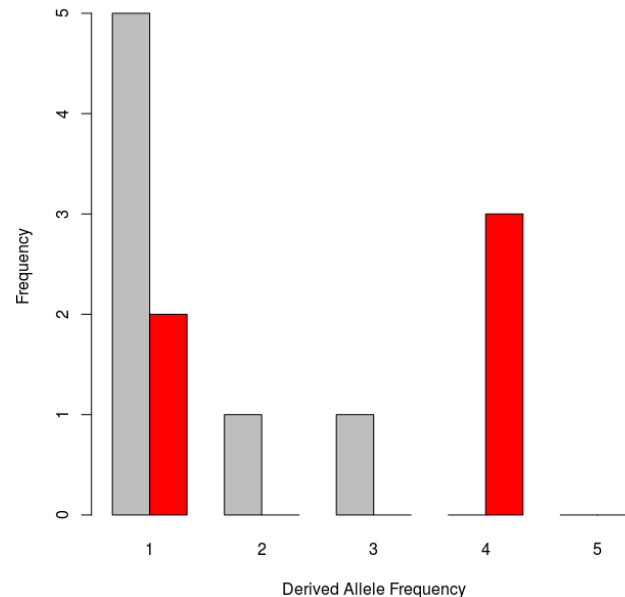


$t < T_{\text{sel}}$

...



$t \gg T_{\text{sel}}$



Tajima's D?

$$D = \frac{\pi - \theta_w}{\sqrt{\hat{V}(\pi - \theta_w)}} \quad \rightarrow \sim 0.33 = 1/3$$

$$\theta_w = \frac{K}{a_n}$$

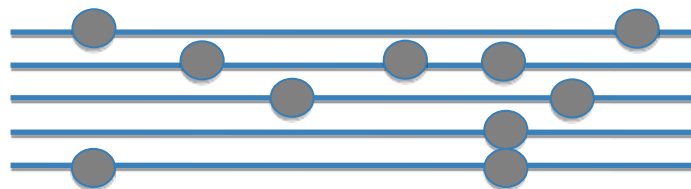
$$a_n = \sum_{i=1}^{n-1} \frac{1}{i}$$

$$\pi = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n k_{i,j}}{\binom{n}{2}}$$

= 10, the number of comparisons you need to make

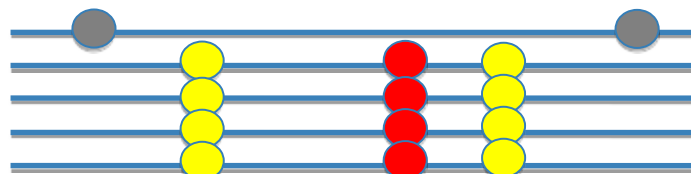
The importance of being...

The Site Frequency Spectrum

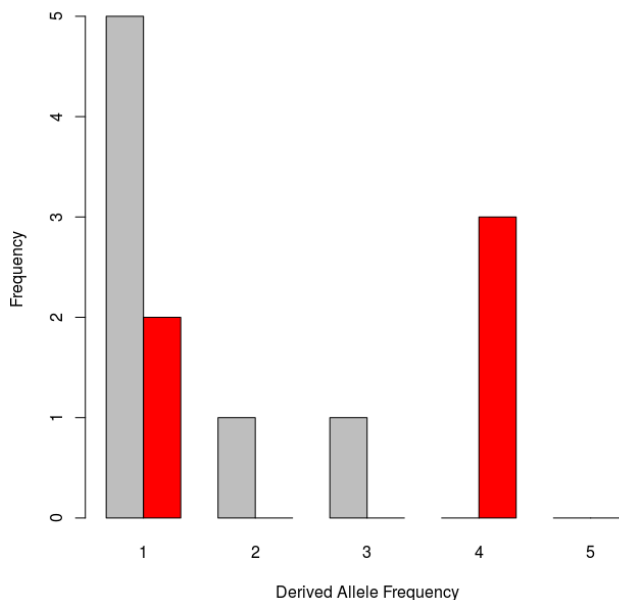


$t < T_{\text{sel}}$

...



$t \gg T_{\text{sel}}$



$K=5$

$a_n = 1/1 + 1/2 + 1/3 + 1/4 = (12+6+4+3)/12 = 25/12$

$\Theta = 5/(25/12) = 12/5$

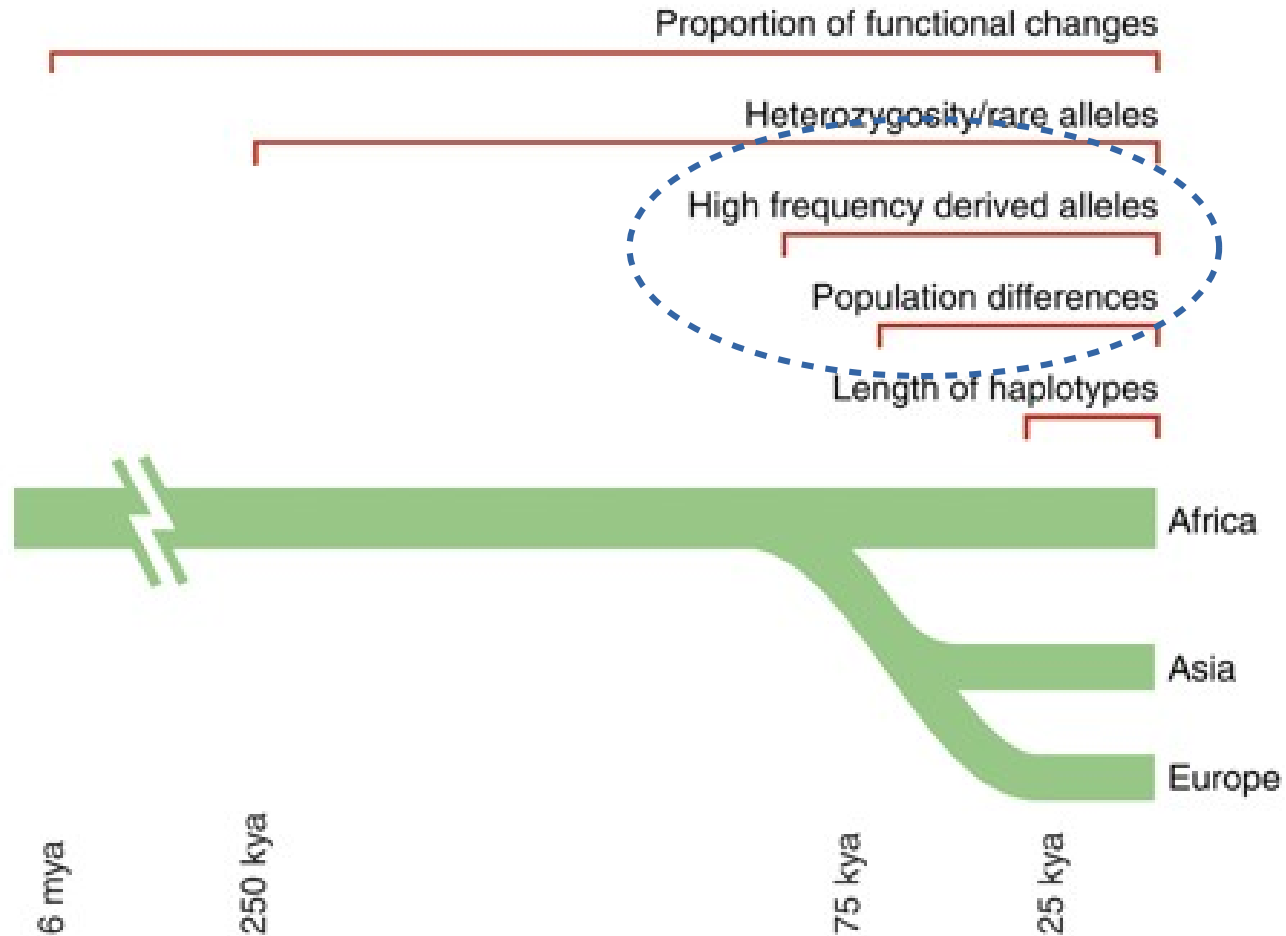
$P_i = (5+5+5+5+0+0+0+0+0+0)/10 = 20/10 = 2$

$sd(D) = 1/3$

$D = (2 - 12/5)/(1/3) = ((10-12)/5)*3 = -6/5 = -1.2$

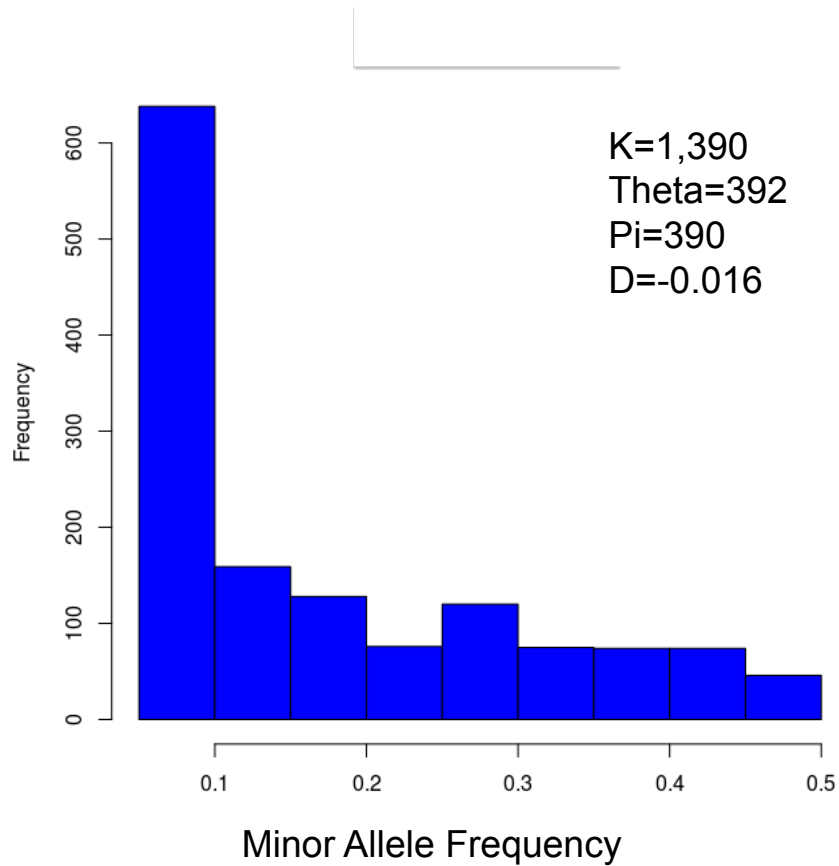
$D < 0$

Inference of positive selection



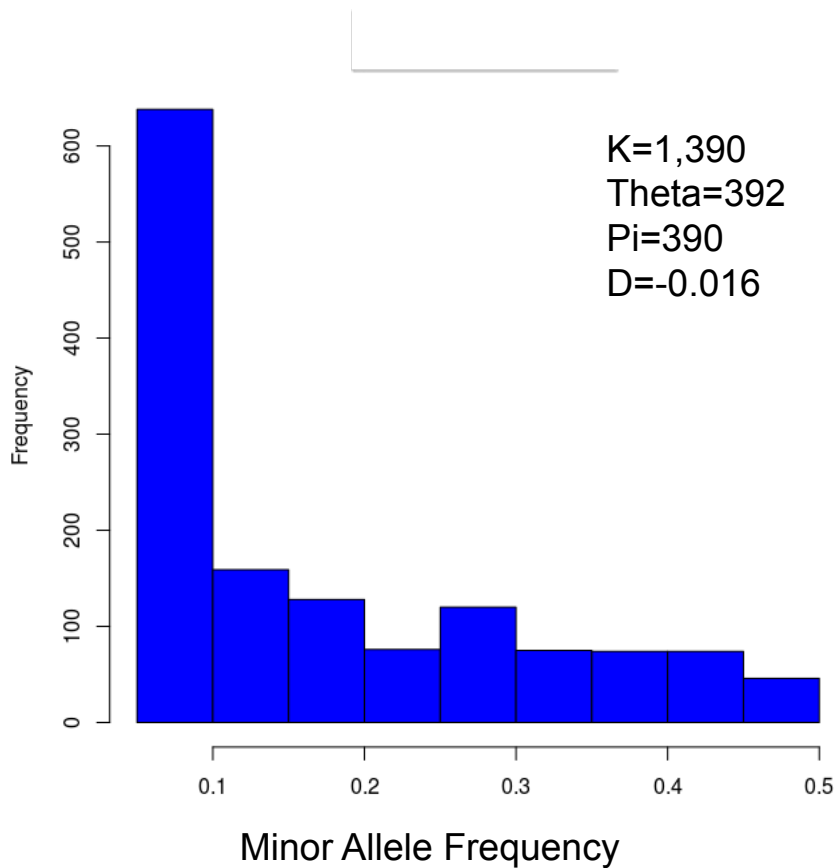
Confounding factor

n=20; L=500kbp; no selection

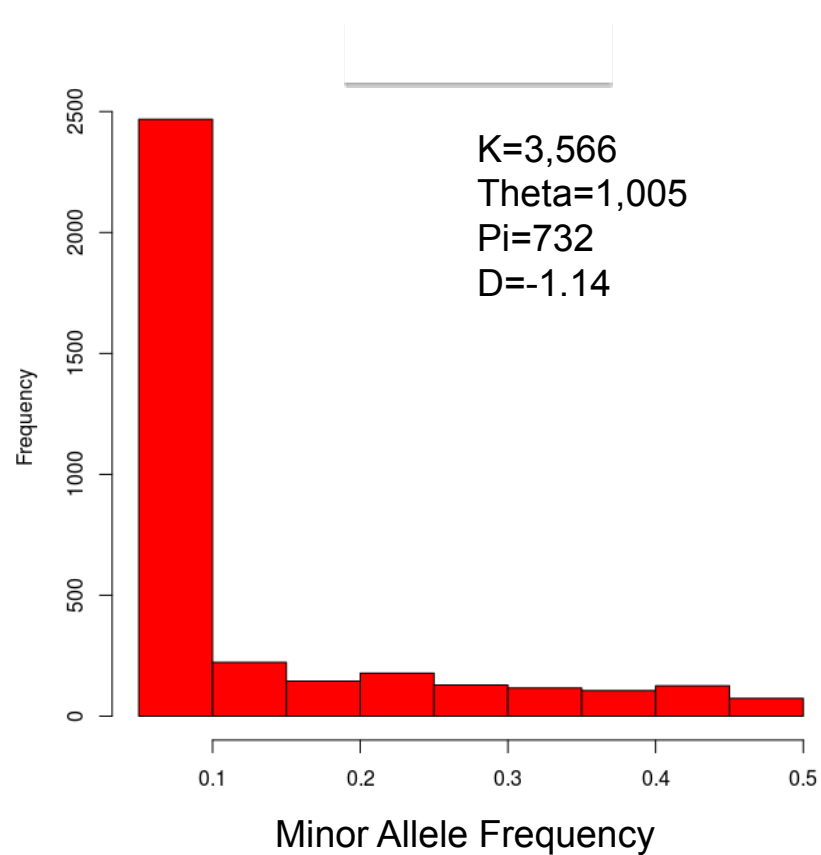


Confounding factor

n=20; L=500kbp; no selection



n=20; L=500kbp; no selection

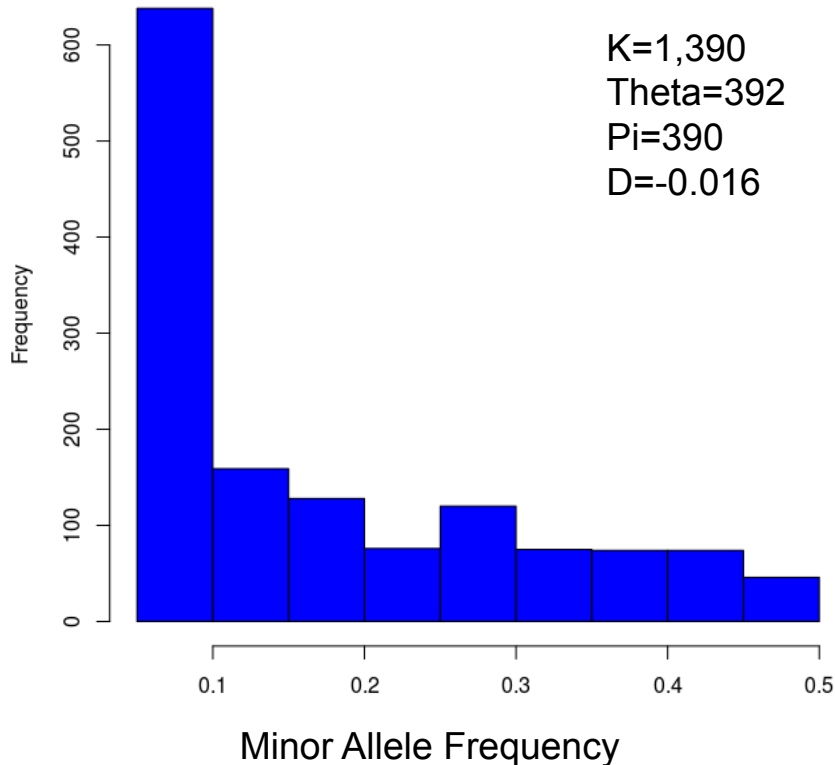


Demography matters!

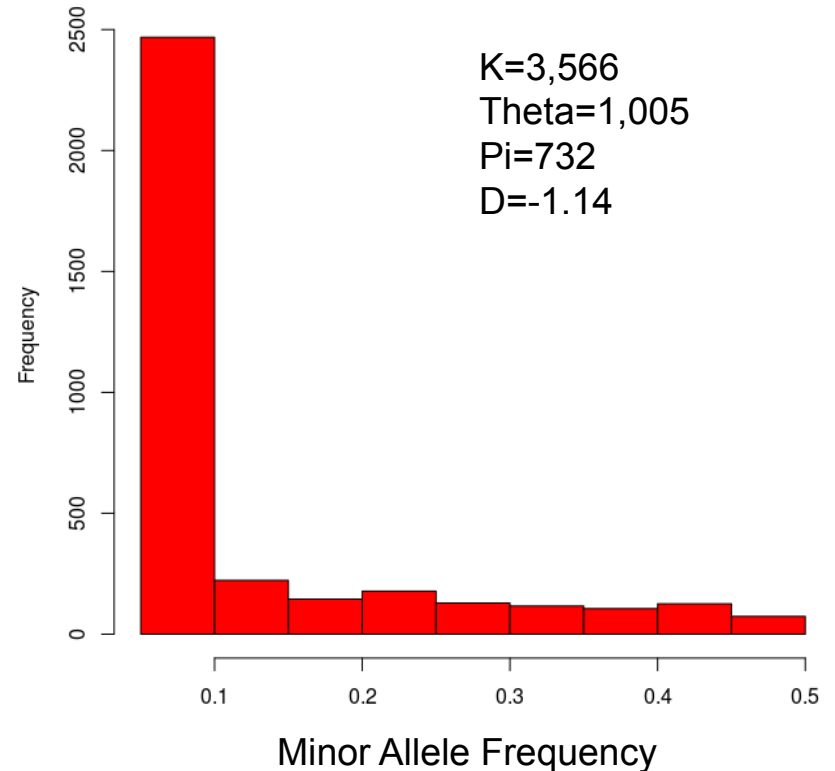
n=20; L=500kbp; no selection

n=20; L=500kbp; no selection

CONSTANT SIZE



EXPANSION

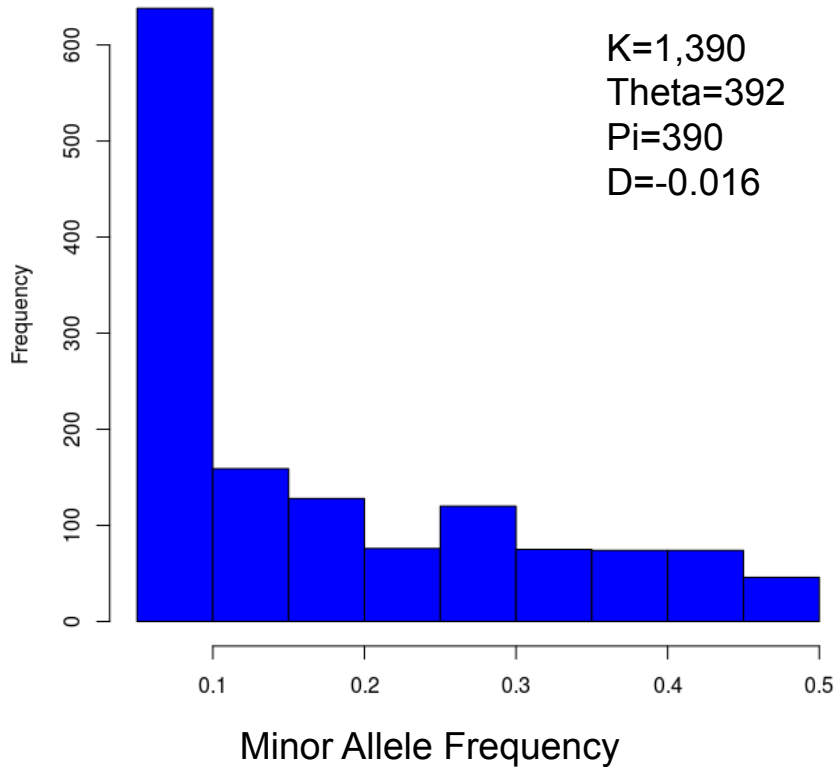


- Excess of segregating sites
- Excess of low-frequency variants
- SFS-derived summary statistics may fail to distinguish between the effects of demography and selection

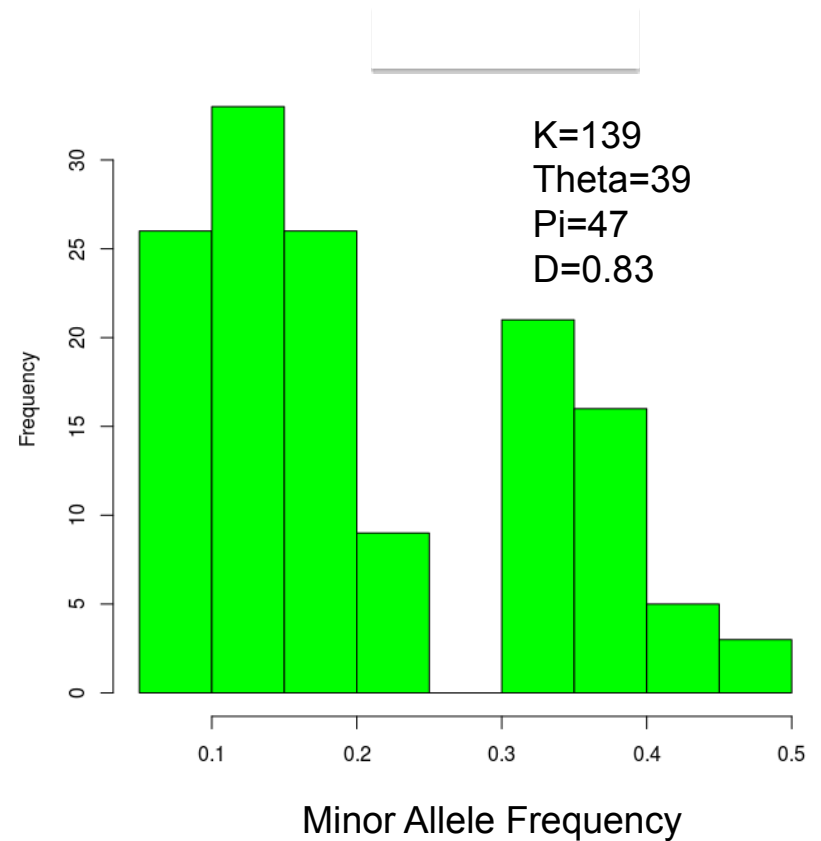
Demography matters?

n=20; L=500kbp; no selection

CONSTANT SIZE



n=20; L=500kbp; no selection

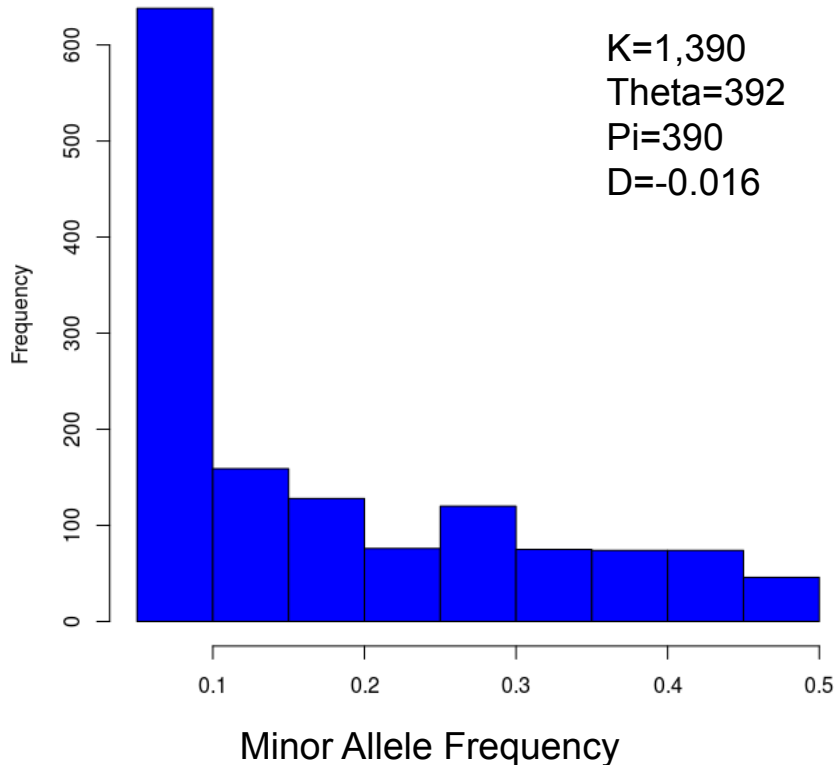


Demography matters!

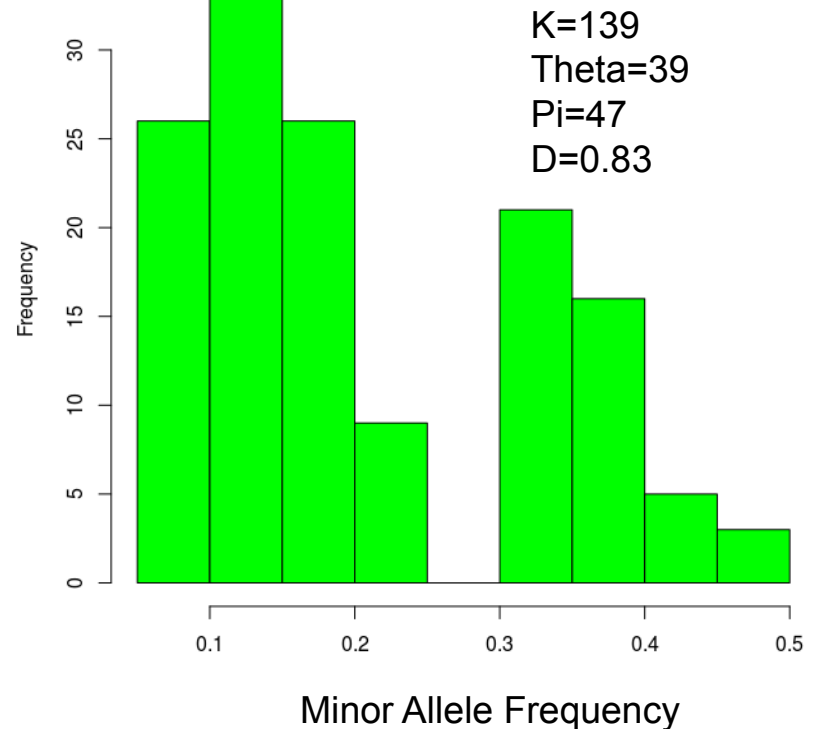
n=20; L=500kbp; no selection

n=20; L=500kbp; no selection

CONSTANT SIZE



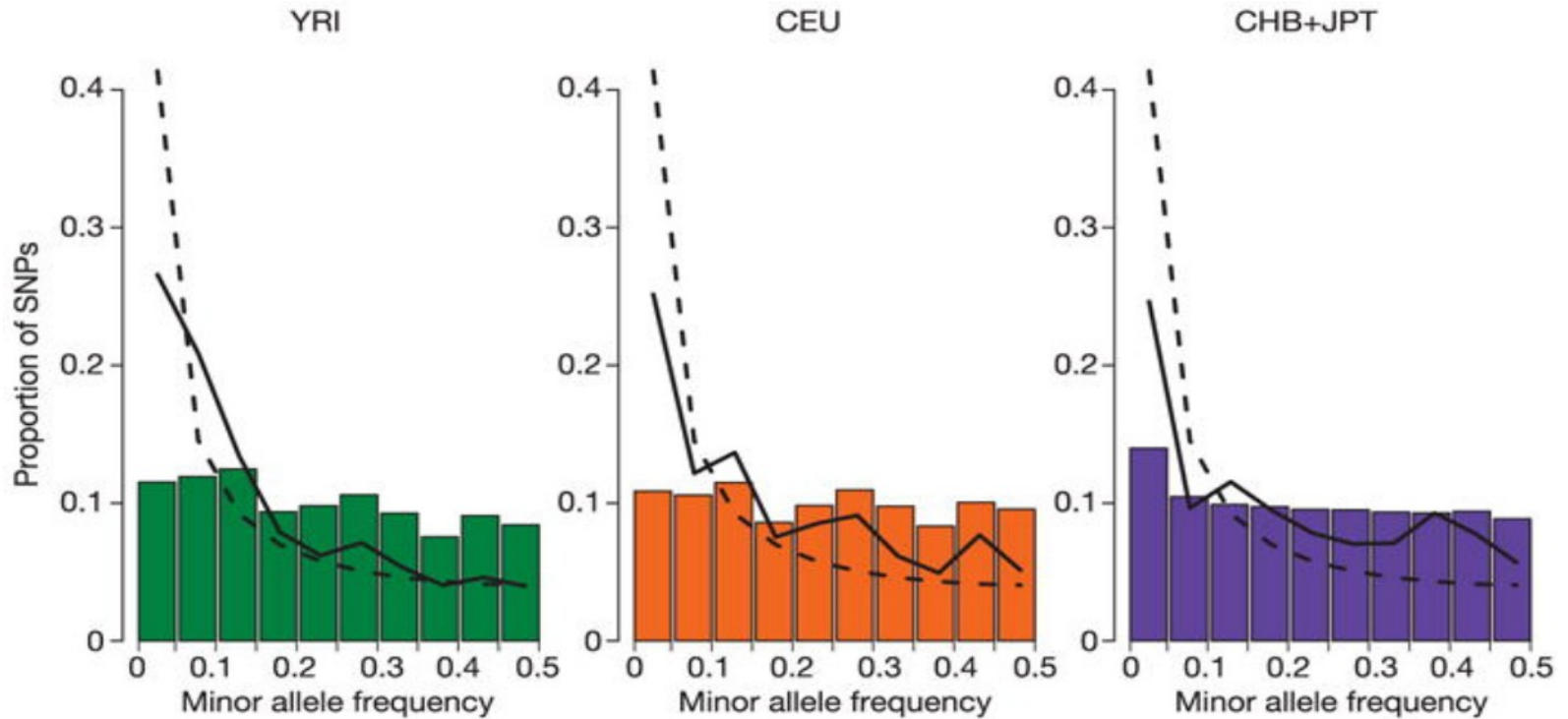
REDUCTION



- Depletion of segregating sites
- Excess of intermediate-frequency variants
- SFS-derived summary statistics may fail to distinguish between the effects of demography and selection

Experimental design matters?

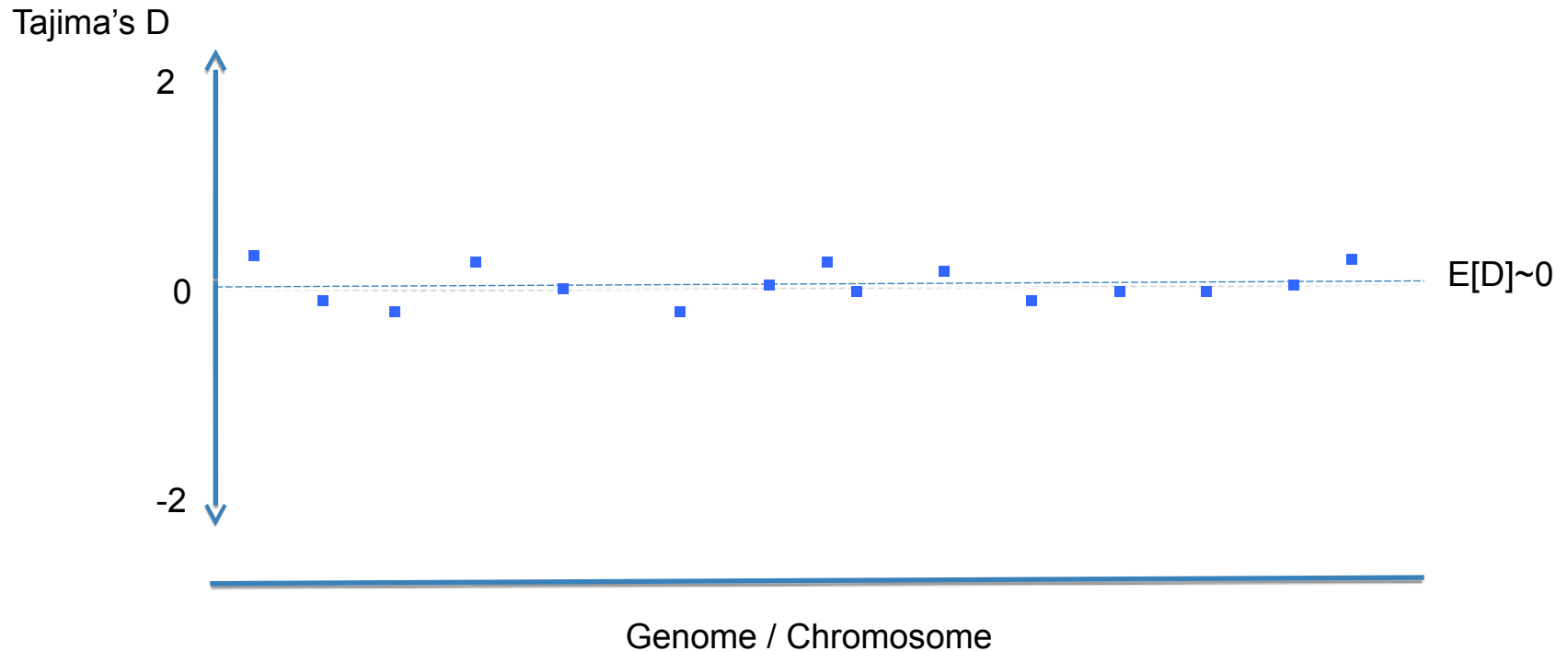
The effect of ascertainment bias



Deficiency of low-frequency variants

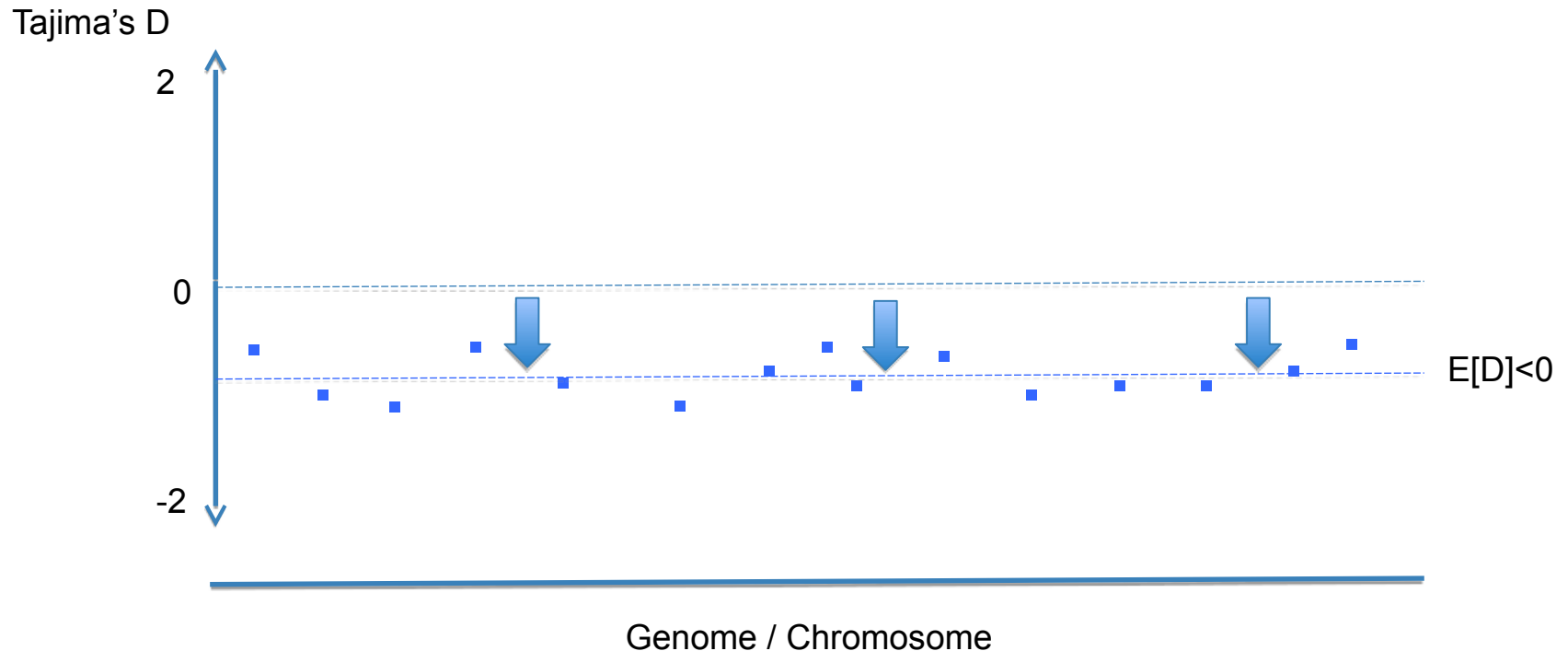
How to take neutral confounding factors into account?

Under constant population size:



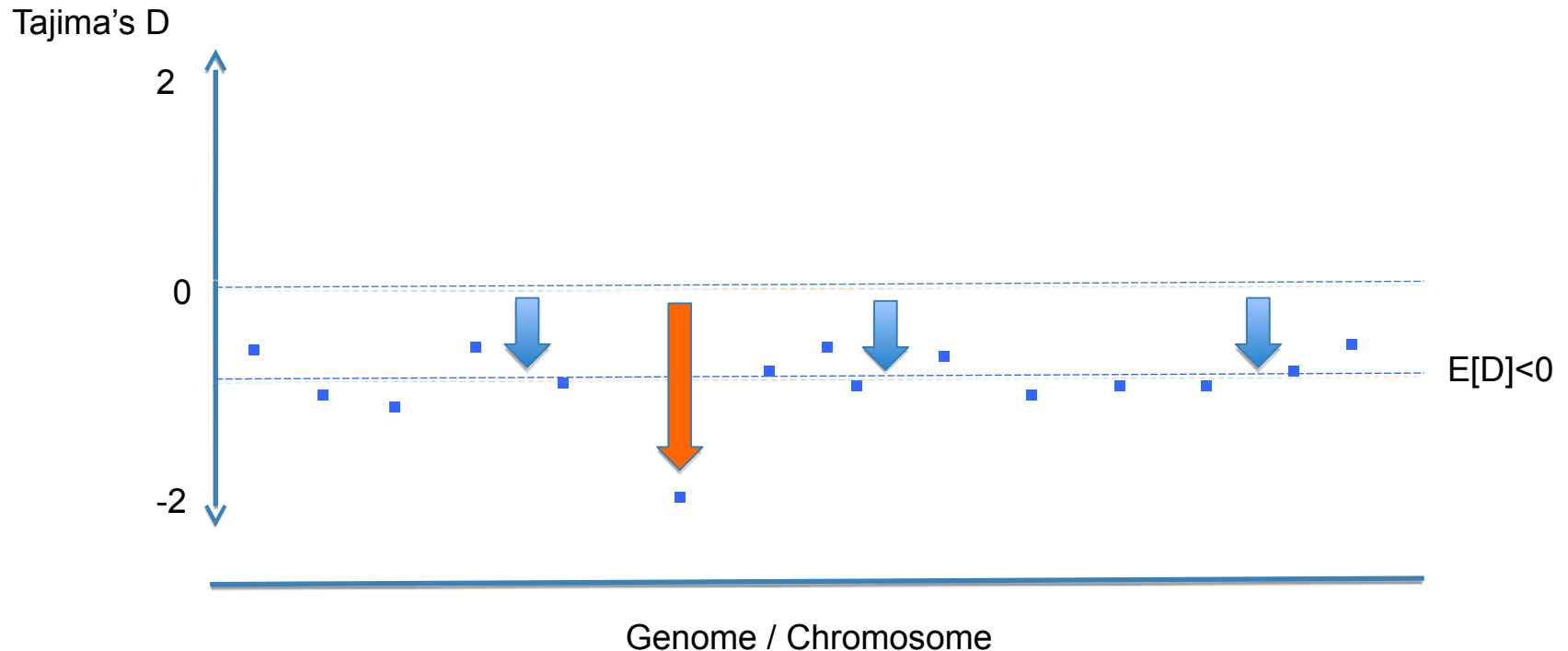
How to take neutral confounding factors into account?

Under expanding population size:



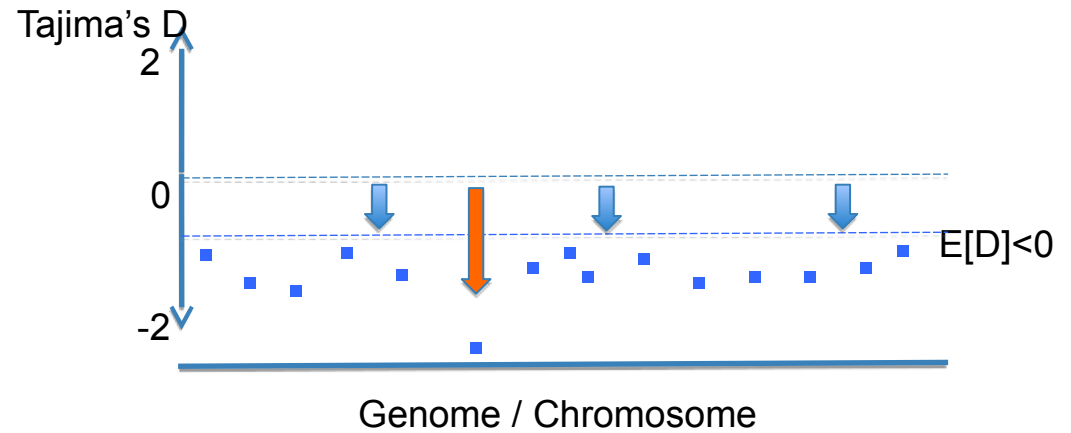
How to take neutral confounding factors into account?

Under expanding population size and positive selection:

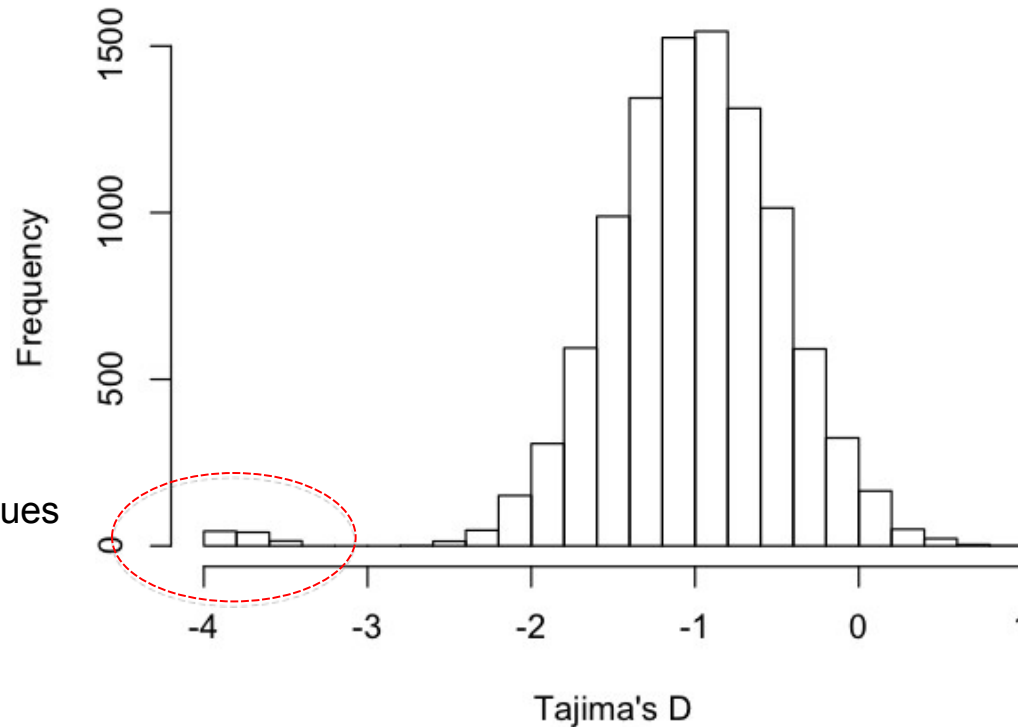


- Demography affects all loci equally, while selection changes local patterns

Outlier approach

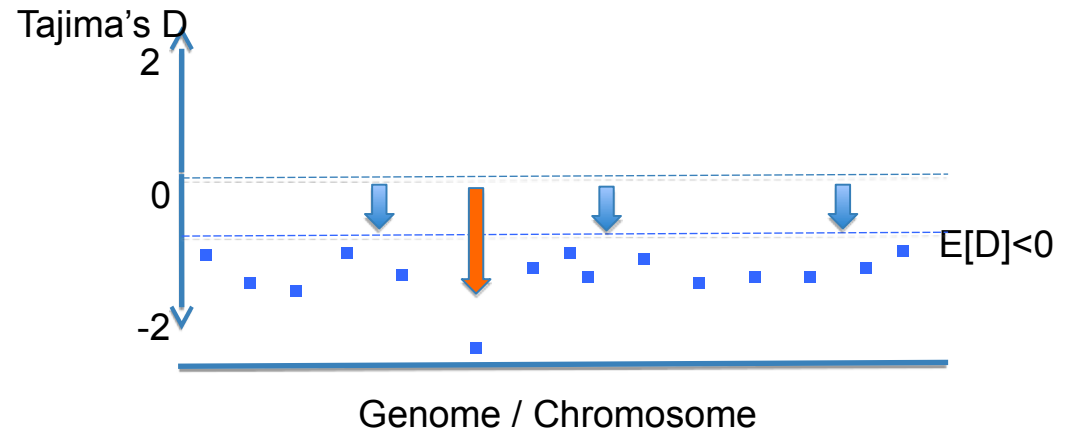


Empirical distribution

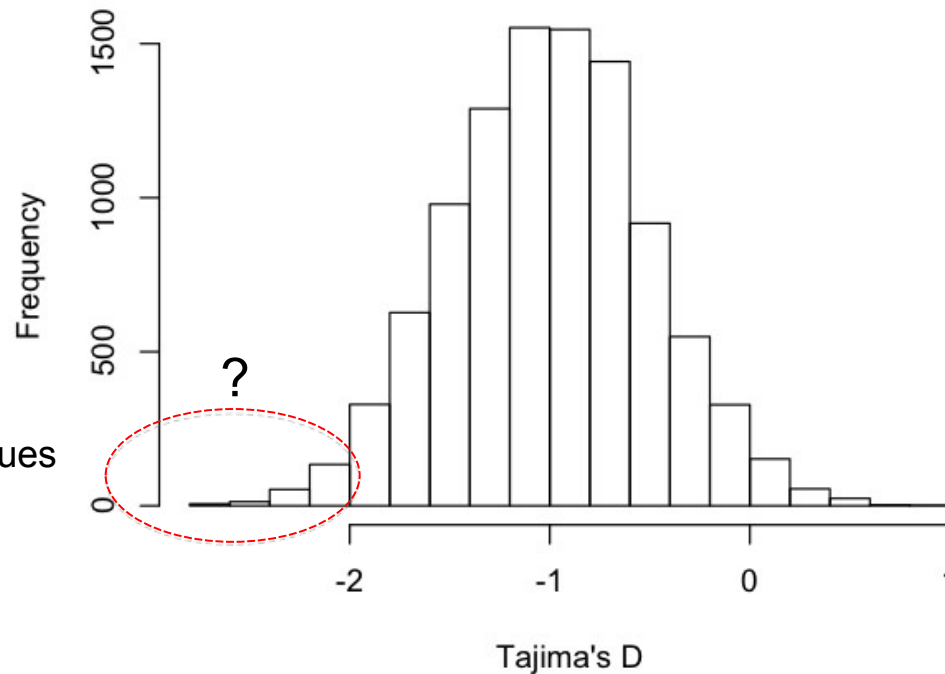


Assign empirical p -values
(ranked percentiles)

Outlier approach



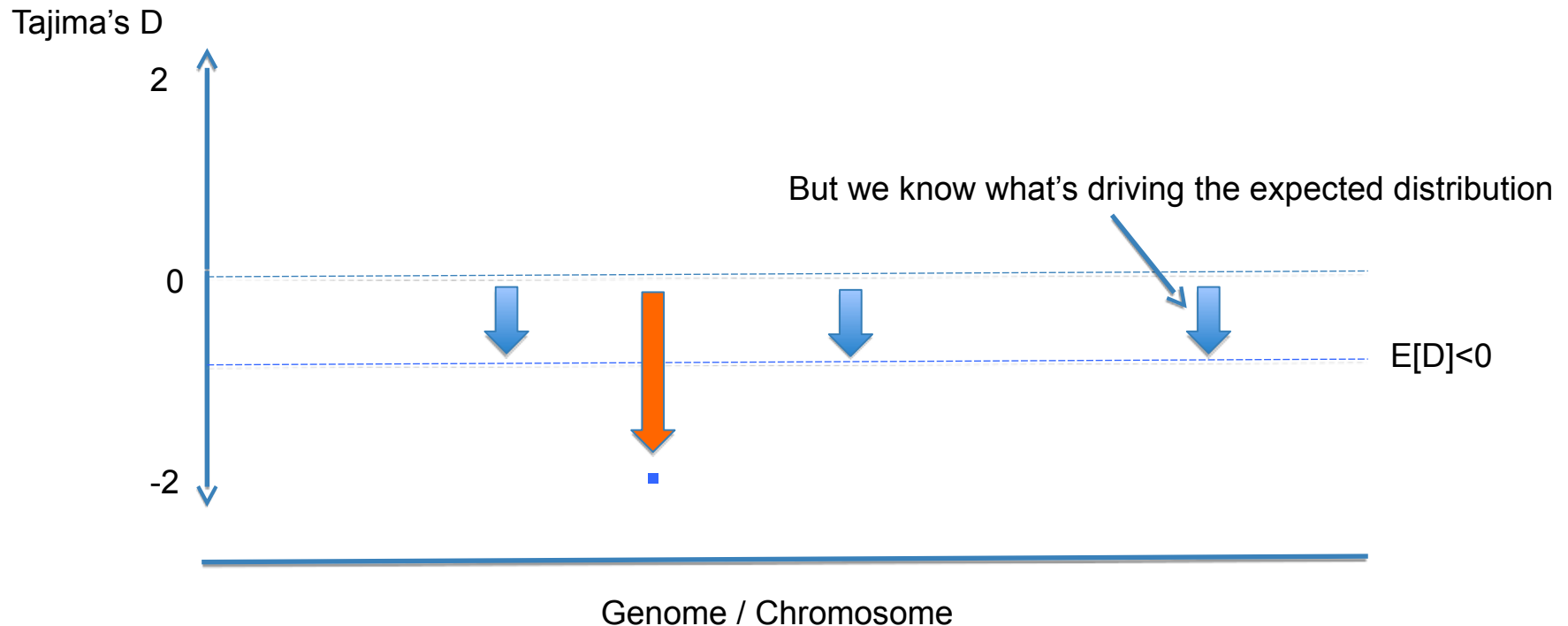
Empirical distribution



Assign empirical p -values
(ranked percentiles)

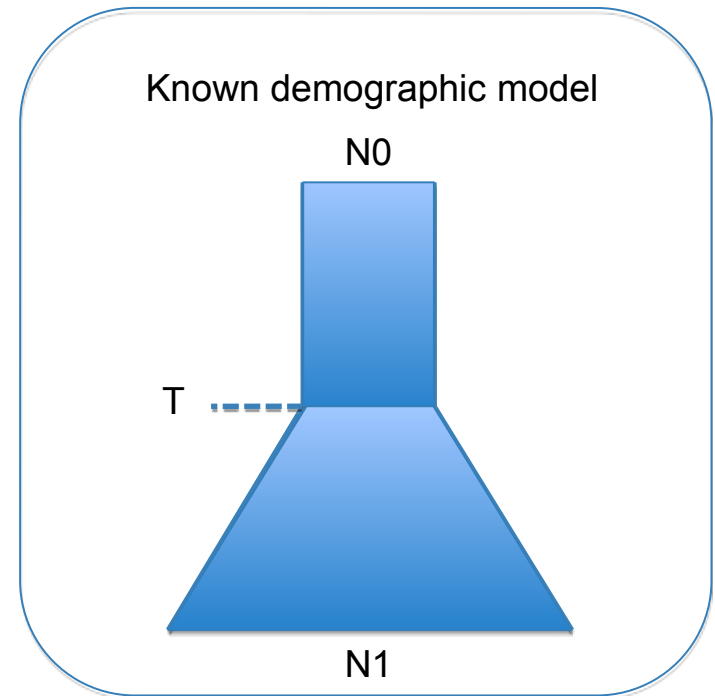
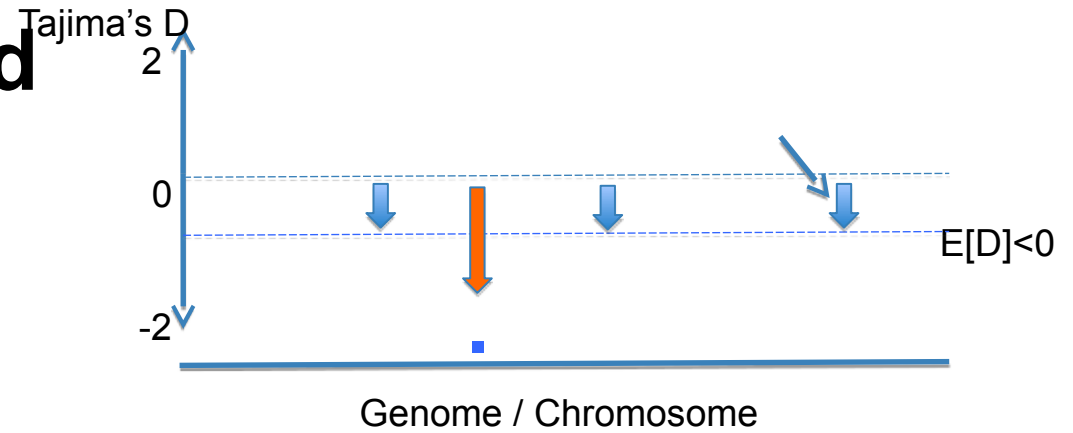
How to take neutral confounding factors into account?

Under expanding population size and positive selection:

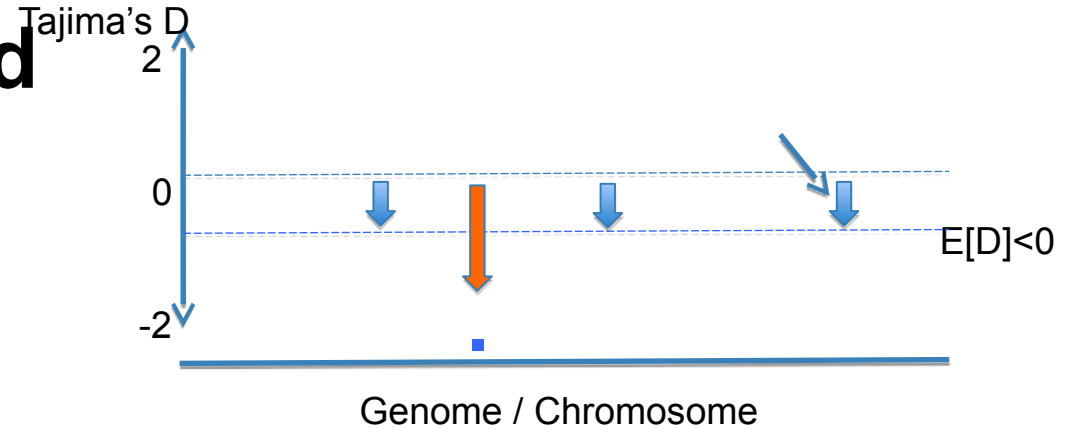


- Demography affects all loci equally, while selection changes local patterns
What should we do if we don't have genome-wide data?

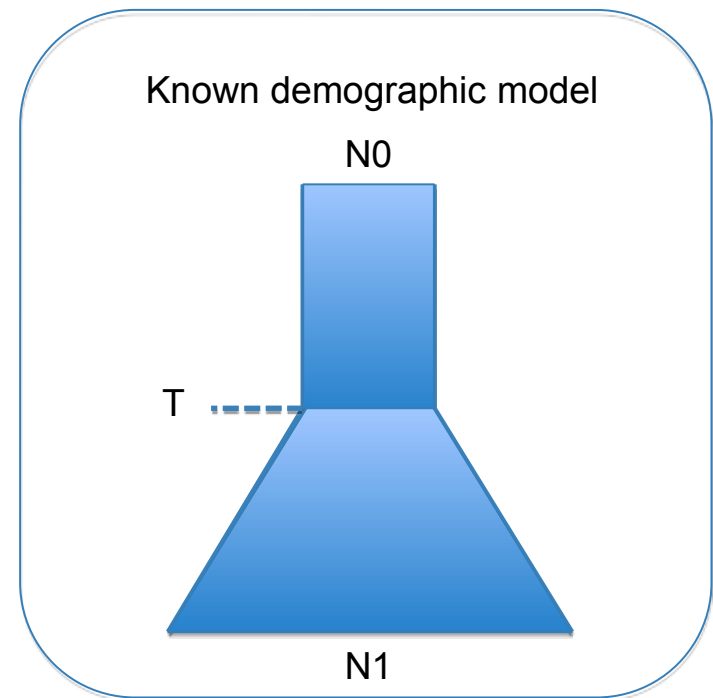
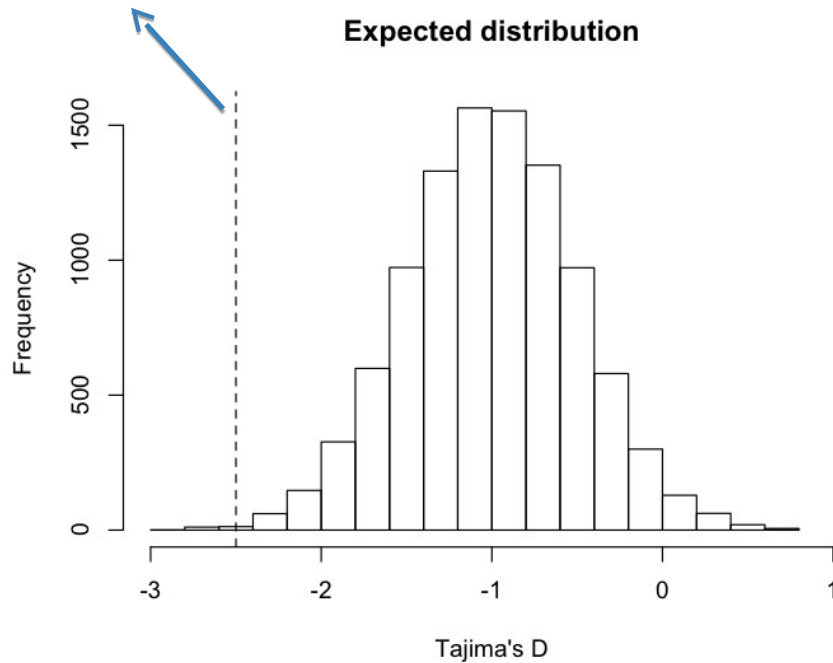
Simulations-based approach



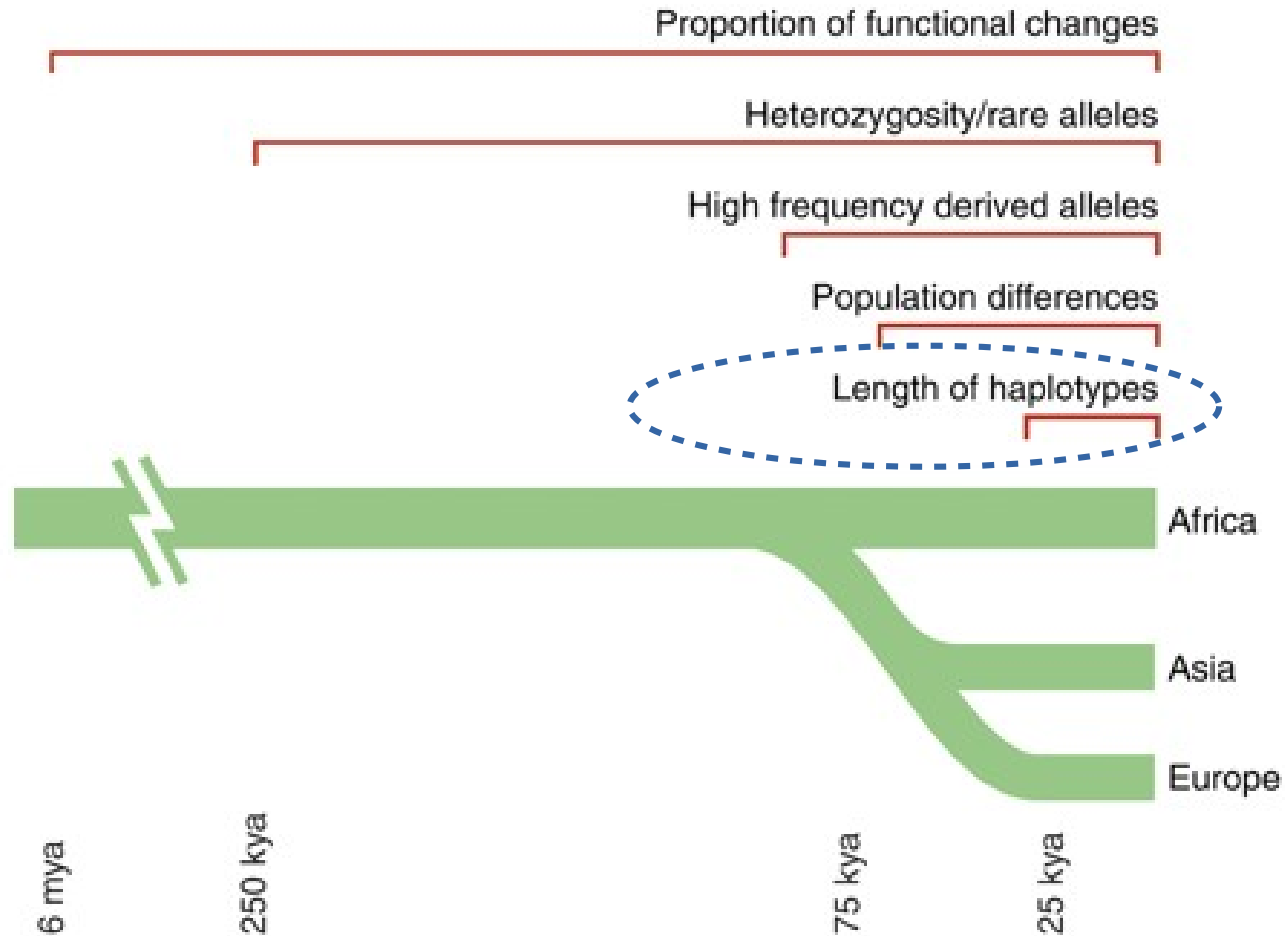
Simulations-based approach



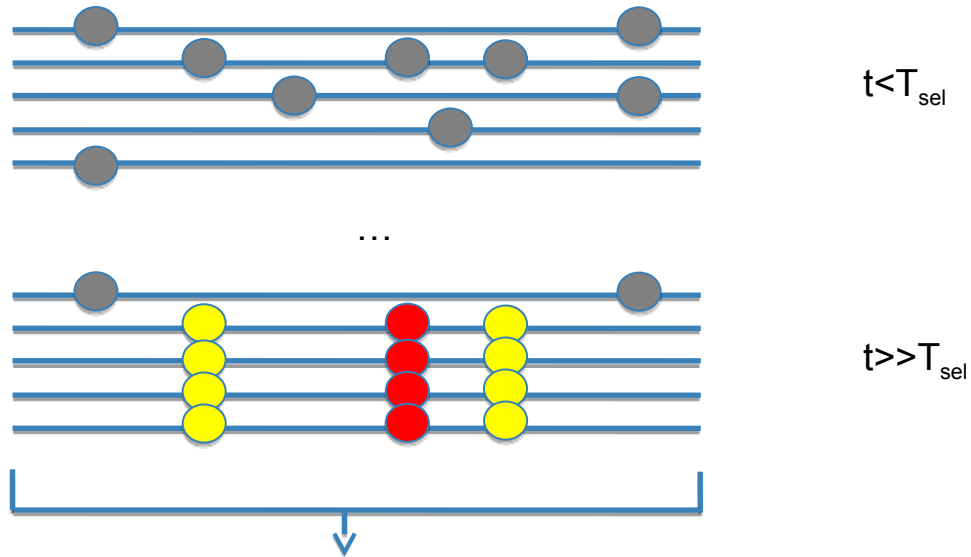
Assign p -values
(based on ranked percentile of observed value)



Inference of positive selection

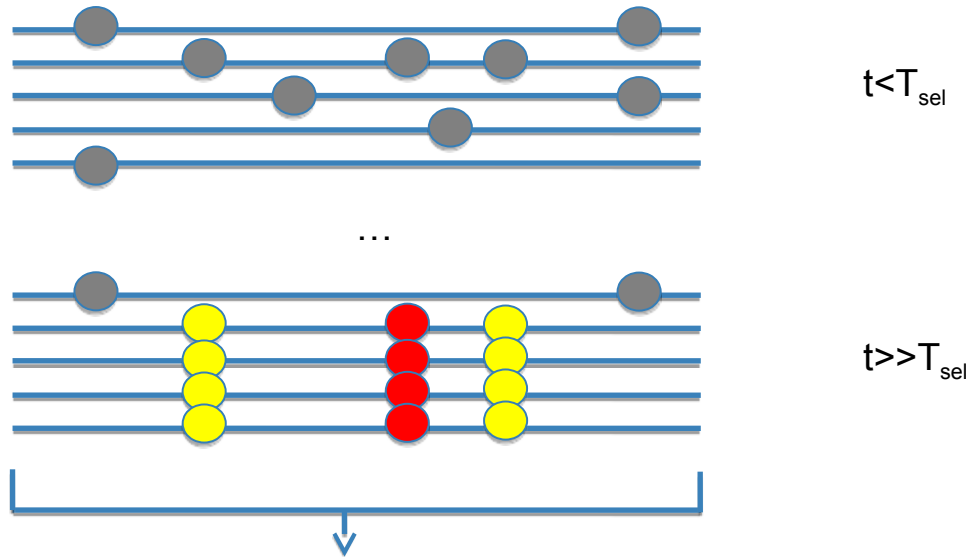


Positive selection



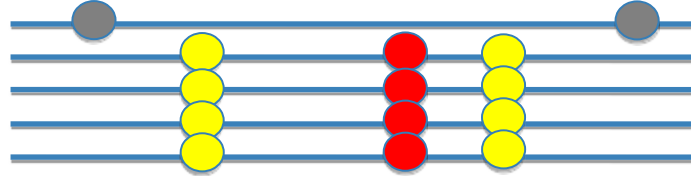
- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants (Pi, Tajima's D, SFS)
- ?

Positive selection



- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants (Pi, Tajima's D, SFS)
- Extended haplotype homozygosity / Extended LD

Extended Haplotype Homozygosity



$t \gg T_{\text{sel}}$

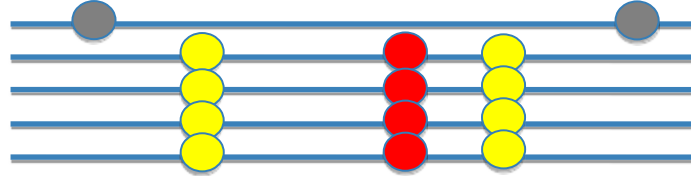
...



$t \gg \gg T_{\text{sel}}$

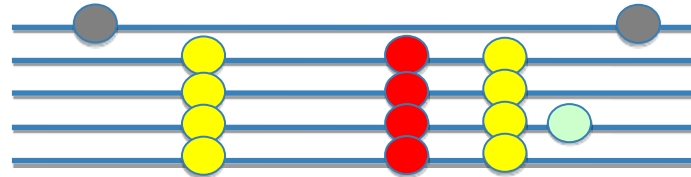
...

Extended Haplotype Homozygosity



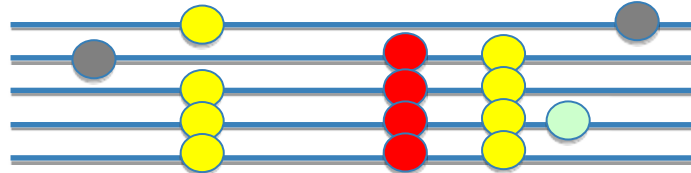
$t \gg T_{\text{sel}}$

...



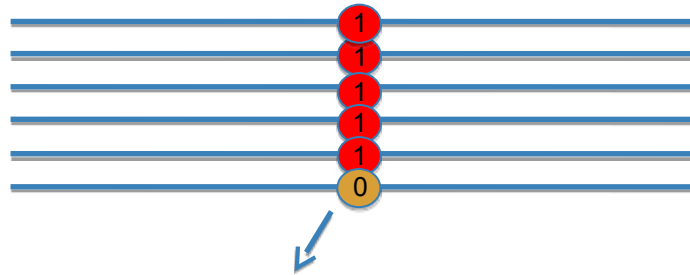
$t \gg T_{\text{sel}}$

...



$t \gg T_{\text{sel}}$

Extended Haplotype Homozygosity

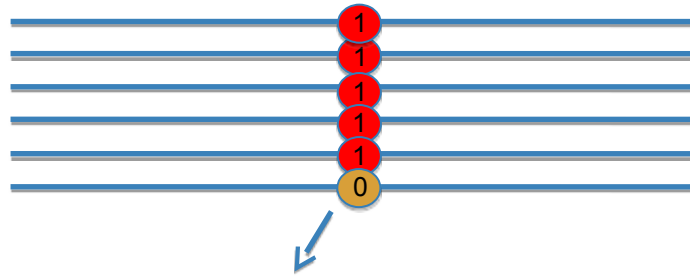


Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

Core SNP

Extended Haplotype Homozygosity

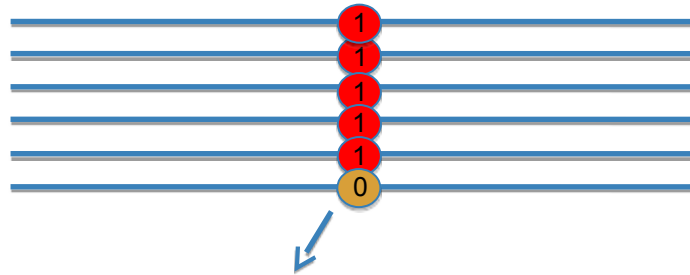


Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

Until marker x_i
(starting from x_0)

Extended Haplotype Homozygosity

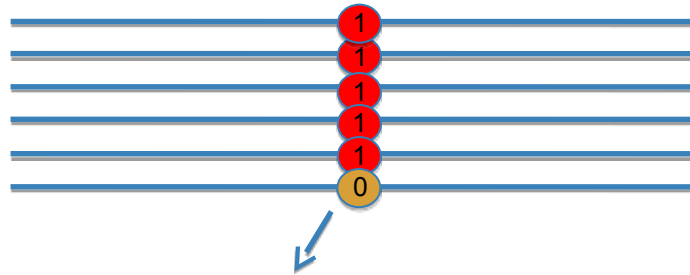


Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

Sum across all unique haplotypes
carrying the core SNP

Extended Haplotype Homozygosity



Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

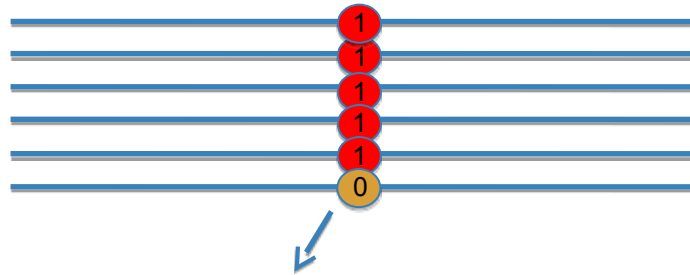
$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

Sum across all unique haplotypes carrying the core SNP

n_h is haplotype frequency of h

n_h is haplotype frequency of the core SNP

Extended Haplotype Homozygosity



Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

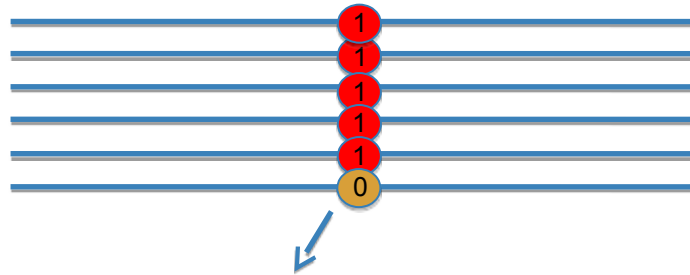
$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

n_h is haplotype frequency of h
 n_c is haplotype frequency of the core SNP

Sum across all unique haplotypes carrying the core SNP

$$EHH_c(x_i = 0) = ?$$

Extended Haplotype Homozygosity



Core haplotype is 1
(Biallelic: 0 is ancestral, 1 is derived allele)

$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

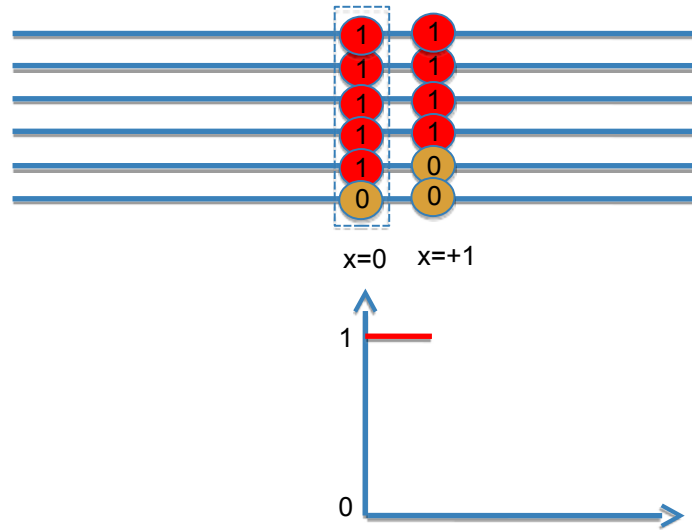
n_h is haplotype frequency of h

n_c is haplotype frequency of the core SNP

Sum across all unique haplotypes carrying the core SNP

$$EHH_c(x_i = 0) = \frac{\binom{5}{2}}{\binom{5}{2}} = 1$$

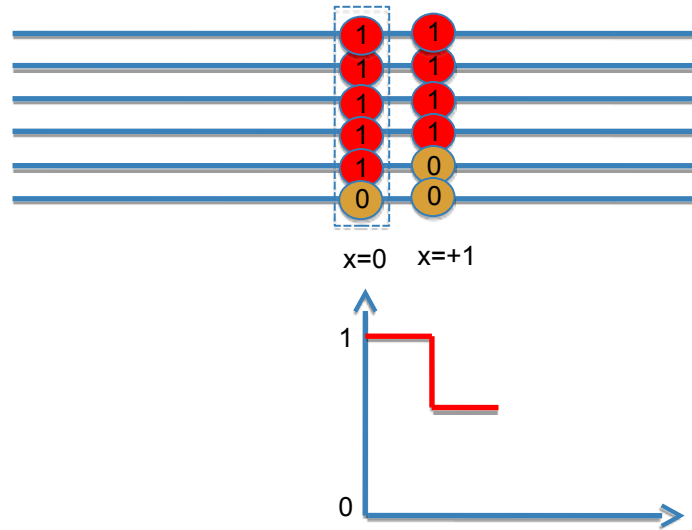
Extended Haplotype Homozygosity



$$EHH_c(x_i = +1) = ?$$

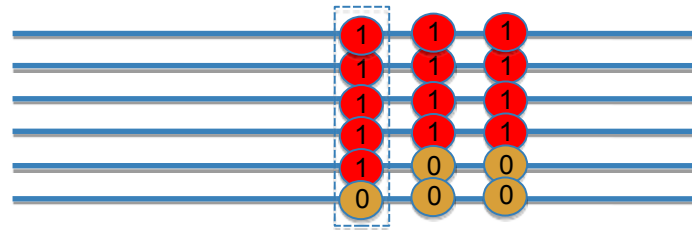
How many unique haplotypes carrying the core SNP?
What is their frequency?

Extended Haplotype Homozygosity

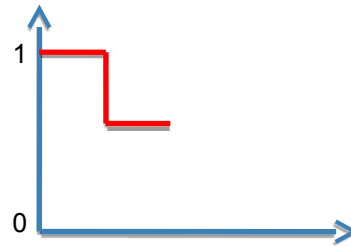


$$EHH_c(x = +1) = \frac{\binom{4}{2} + \binom{1}{2}}{\binom{5}{2}} = \frac{6 + 0}{10} = 0.60$$

Extended Haplotype Homozygosity

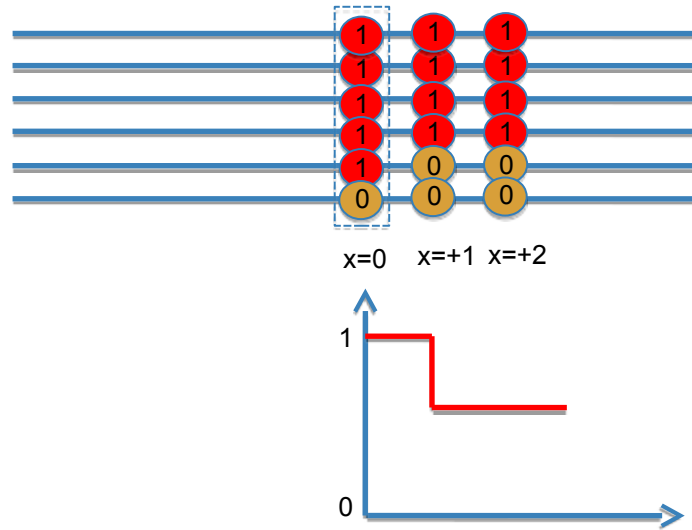


x=0 x=+1 x=+2



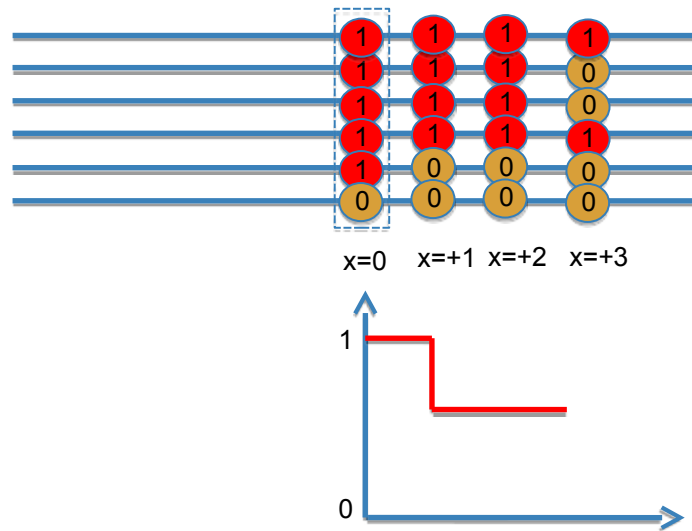
$$EHH_c(x_i = +2) = ?$$

Extended Haplotype Homozygosity



$$EHH_c(x_i = +2) = EHH_c(x_i = +1) = 0.60$$

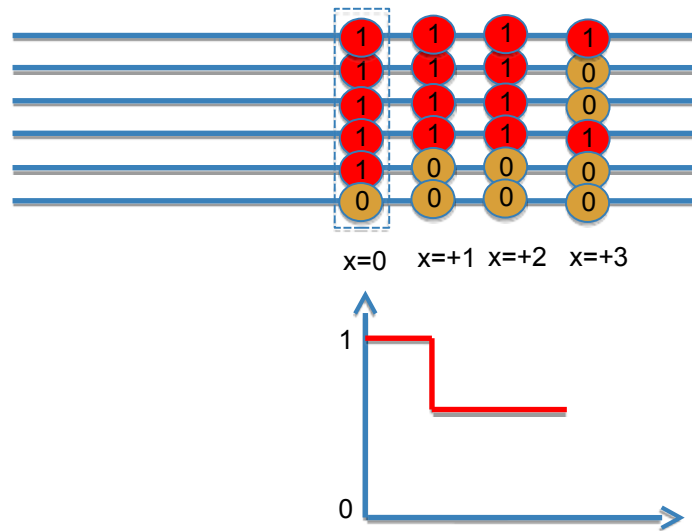
Extended Haplotype Homozygosity



$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

How many unique haplotypes carrying the core SNP?
What is their frequency?

Extended Haplotype Homozygosity



$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

How many unique haplotypes carrying the core SNP?

What is their frequency?

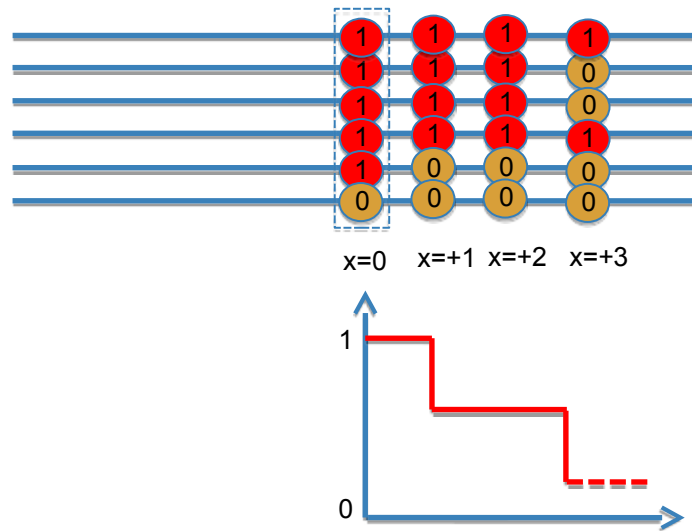
1111 with freq=2

1110 with freq=2

1000 with freq=1

$$EHH_c(x_i = +3) = ?$$

Extended Haplotype Homozygosity



$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

How many unique haplotypes carrying the core SNP?

What is their frequency?

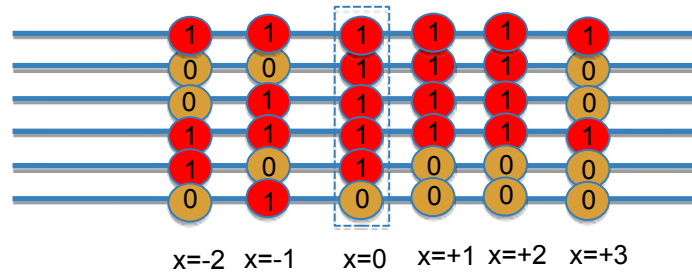
1111 with freq=2

1110 with freq=2

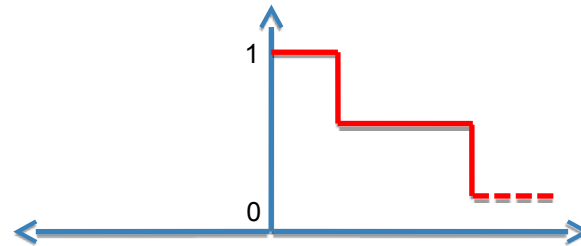
1000 with freq=1

$$EHH_c(x_i = +3) = \frac{\binom{2}{2} + \binom{2}{2} + \binom{1}{2}}{\binom{5}{2}} = \frac{1+1+0}{10} = 0.20$$

Extended Haplotype Homozygosity



$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$



n n choose 2

1 0

2 1

3 3

4 6

5 10

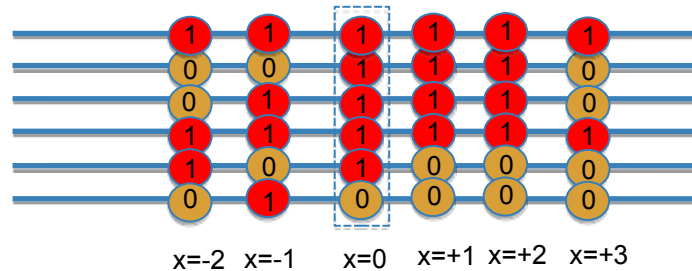
6 15

$$EHH_c(x_i = -1) = ?$$

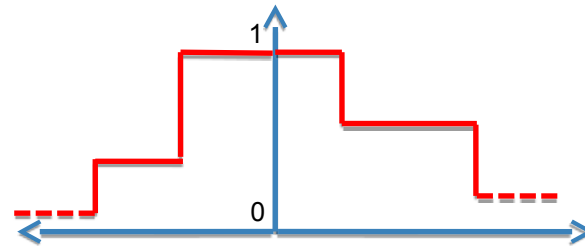
$$EHH_c(x_i = -2) = ?$$

Comment on differences (if any) between $EHH(x=+2)$ and $EHH(x=-2)$.

Extended Haplotype Homozygosity



$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$



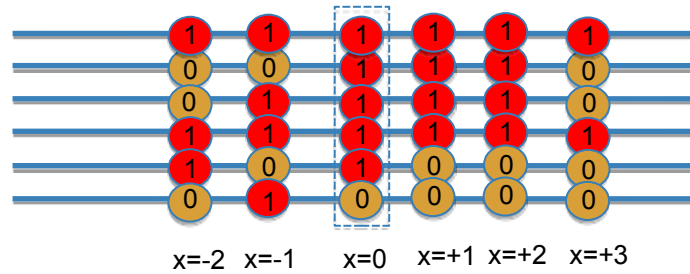
$$EHH_c(x_i = -1) = \frac{\binom{3}{2} + \binom{2}{2}}{\binom{5}{2}} = \frac{3+1}{10} = 0.4$$

$$EHH_c(x_i = -2) = \frac{\binom{2}{2} + \binom{1}{2} + \binom{1}{2}}{\binom{5}{2}} = \frac{1+0+0}{10} = 0.1$$

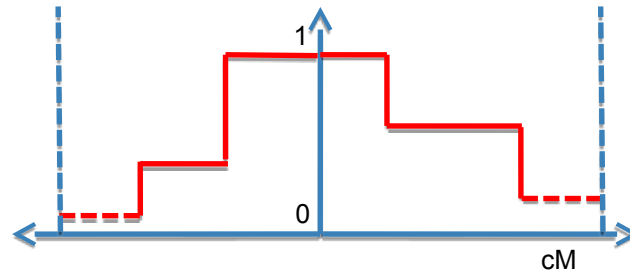
+ (1 choose 2)

Comment on differences (if any) between $EHH(x=+2)$ and $EHH(x=-2)$?

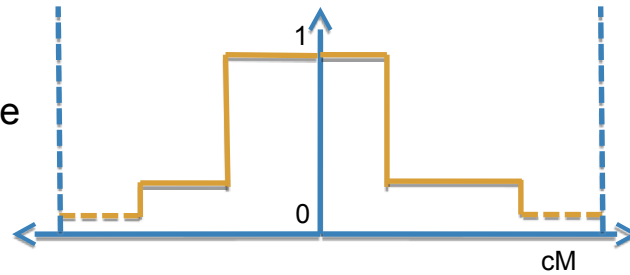
Integrated Haplotype Score



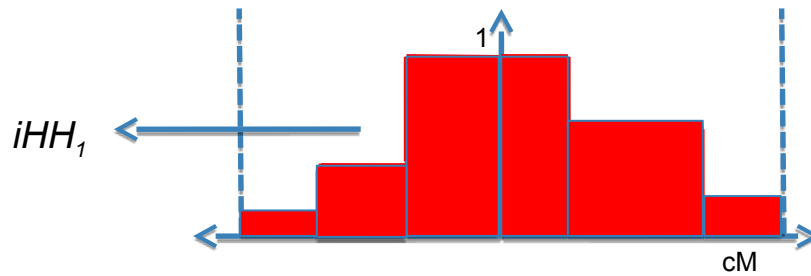
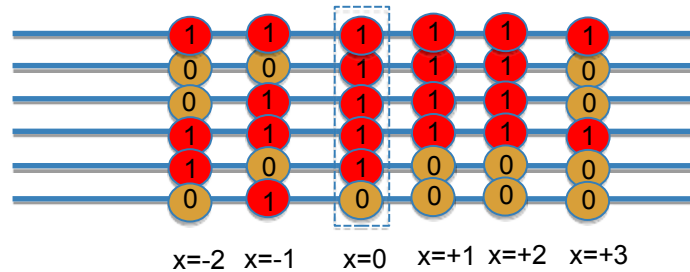
For the derived allele



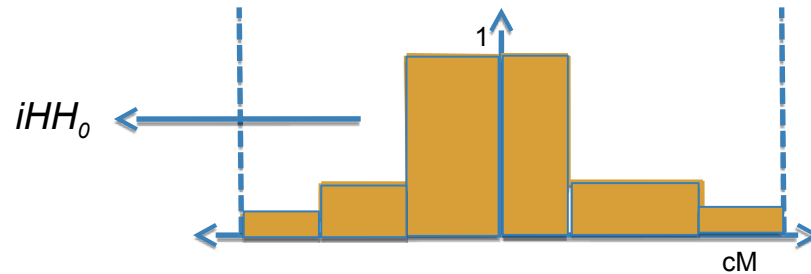
For the ancestral allele



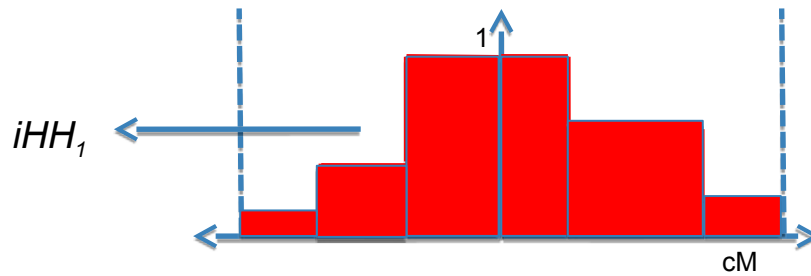
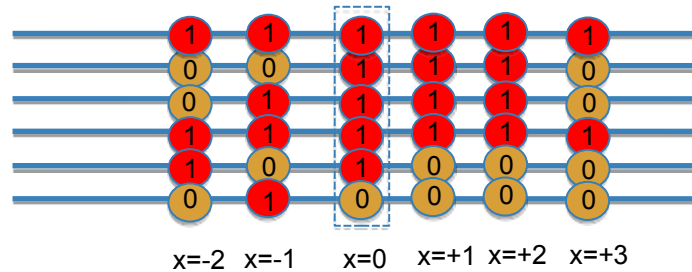
Integrated Haplotype Score



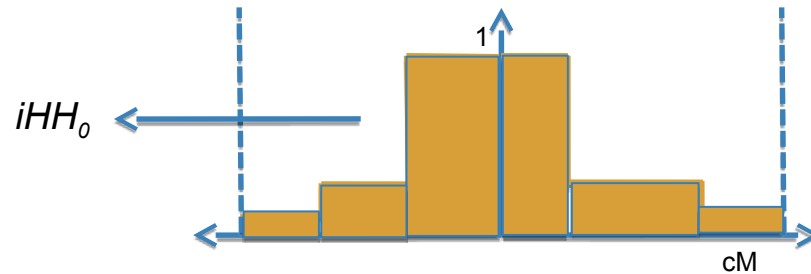
Integrated haplotype homozygosity (iHH)



Integrated Haplotype Score



Integrated haplotype homozygosity (iHH)



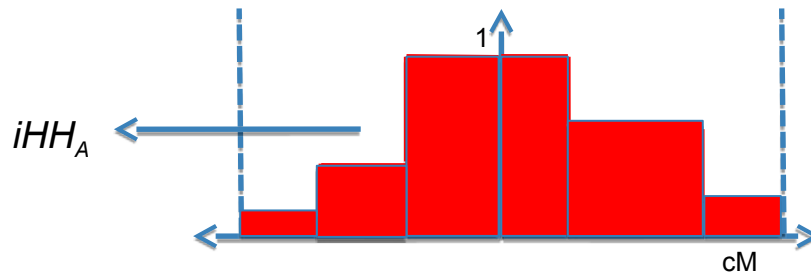
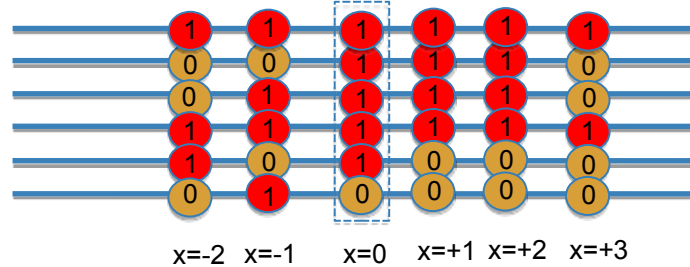
Integrated haplotype score:

$$iHs = \ln(iHH_1/iHH_0)$$

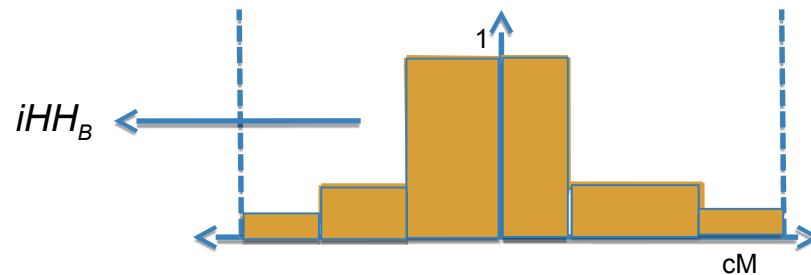


Genome-wide normalization in frequency bins
(to mean=0 and sd=1)

Cross-population Extended Haplotype Homozygosity



Integrated haplotype homozygosity (iHH)
for **populations A and B**

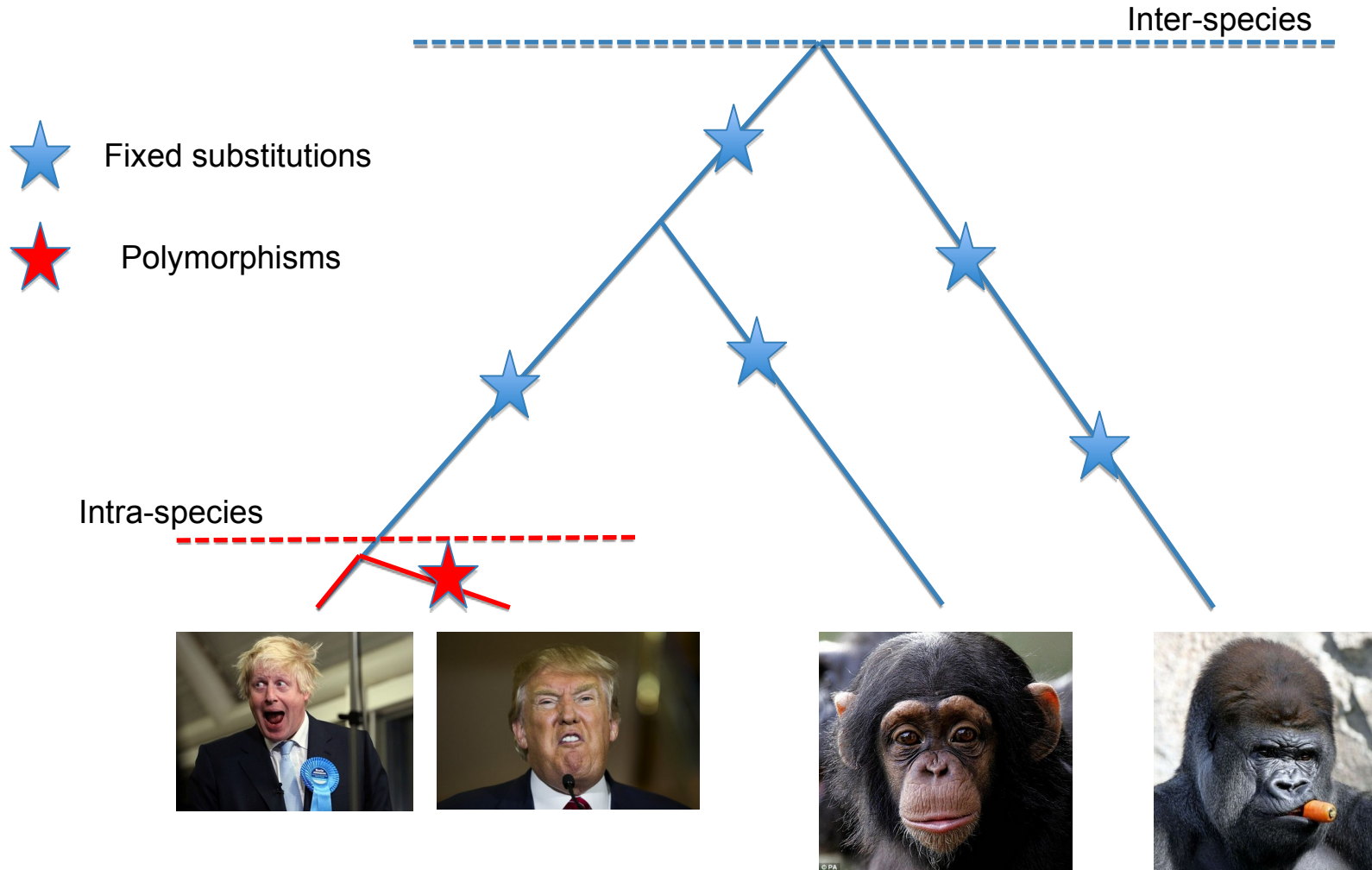


Integrated haplotype score:
 $XP-EHH = \ln(iHH_A/iHH_B)$



Genome-wide normalization in frequency bins
(to mean=0 and sd=1)

Inferring inter-species selection



State-of-the-art methods to detect natural selection

1. Composite scores (Grossman et al. 2013)

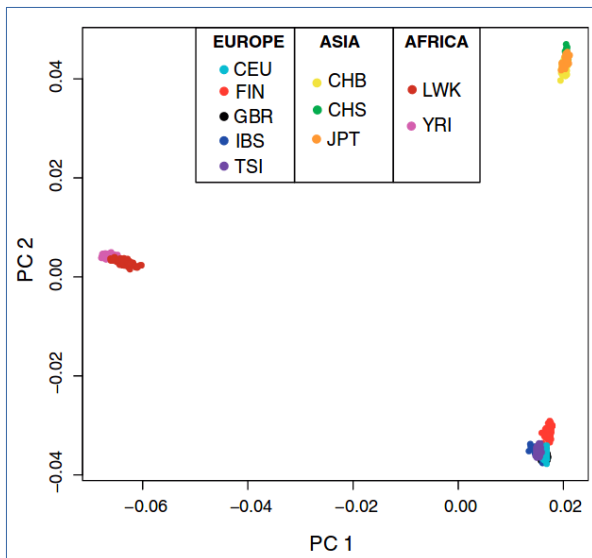
$$BF_t = \frac{P(v_t \in \text{bin}_{t,k} | \text{selected})}{P(v_t \in \text{bin}_{t,k} | \text{unselected})}$$

and defined the composite score as the product of the Bayes factor of each test:

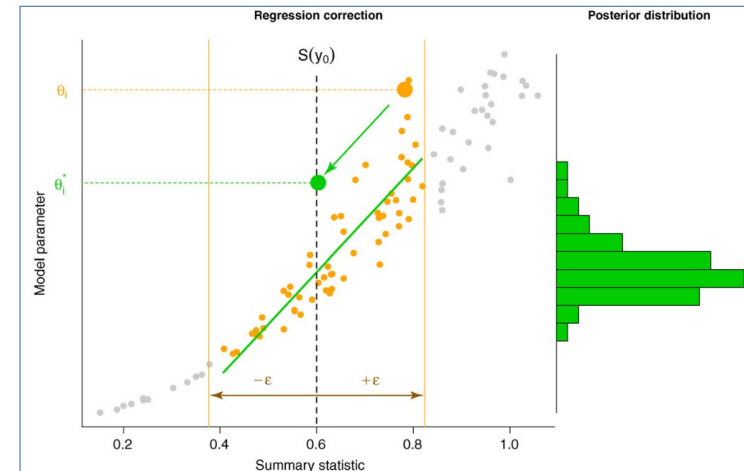
$$CMS_{GW} = \prod_{t \in \text{tests}} BF_t$$

3. Unsupervised machine learning

(PCA, Duforet-Frebourg et al. 2016)

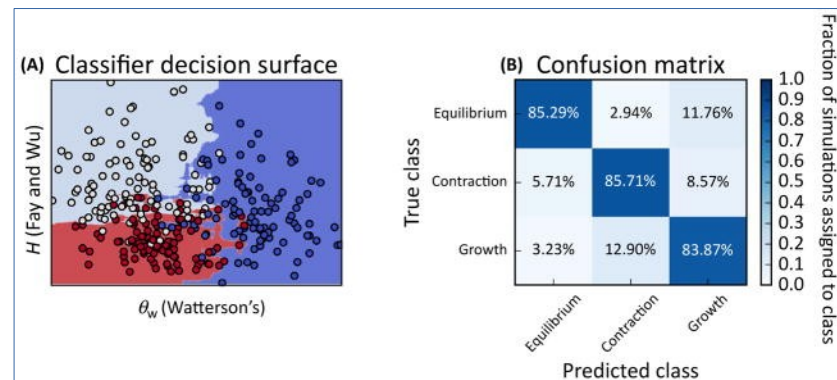


2. Simulations-based (rejection, ABC)



4. Supervised machine learning

(SVM, Schrider & Kern 2018)



How do we infer signals of natural selection from genomic data using machine learning and deep learning?

(let's move to part 3)

What are the main limitations of currently employed methods to detect selection?