

# *Inference of natural selection from NGS data*

at the intra-species level

Matteo Fumagalli

# Intended Learning Outcomes

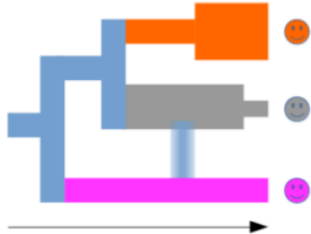
At the end of this session you will be able to:

- list commonly used methods to detect selection
  - calculate various summary statistics
  - understand main confounding factors to neutrality tests
  - assess statistical significance of tests
-

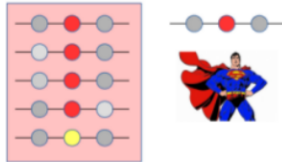
## Whole genome sequencing



## Demographic history



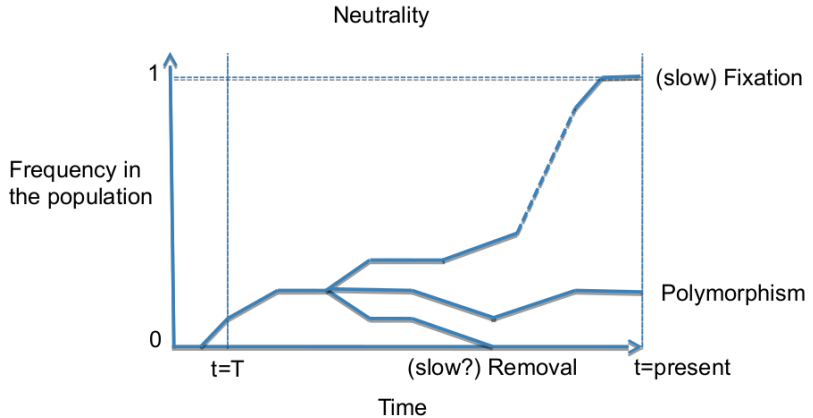
## Natural selection



# Natural selection

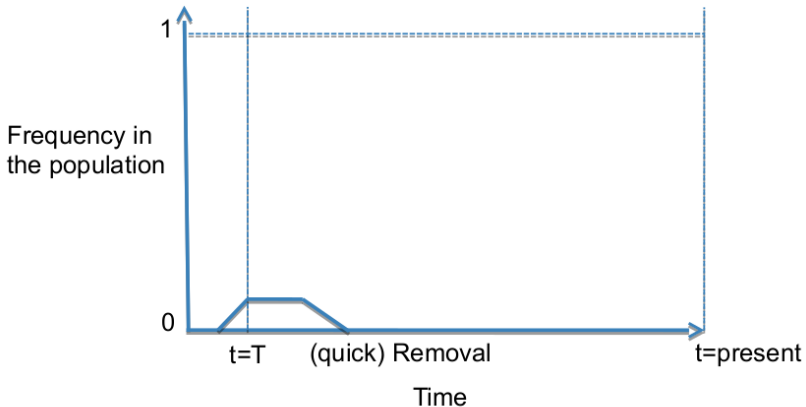
- Heritable traits that increase the fitness of the become more common.
- Sites targeted by natural selection are likely to harbour functionality
- Mutations arise randomly and evolve according to their effect on the fitness of the carrier.

# Allele frequency trajectory

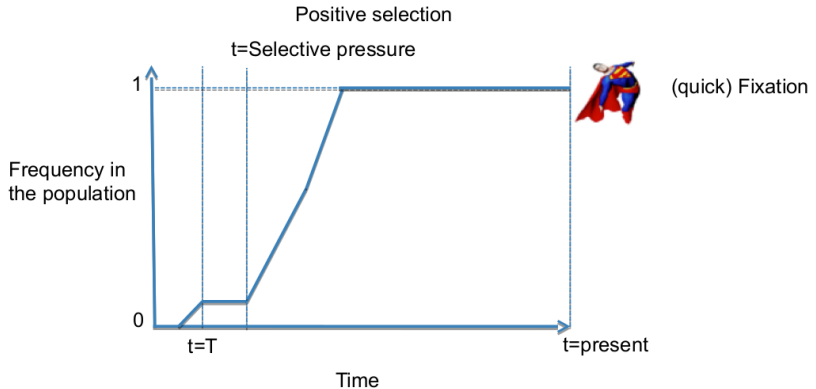


# Allele frequency trajectory

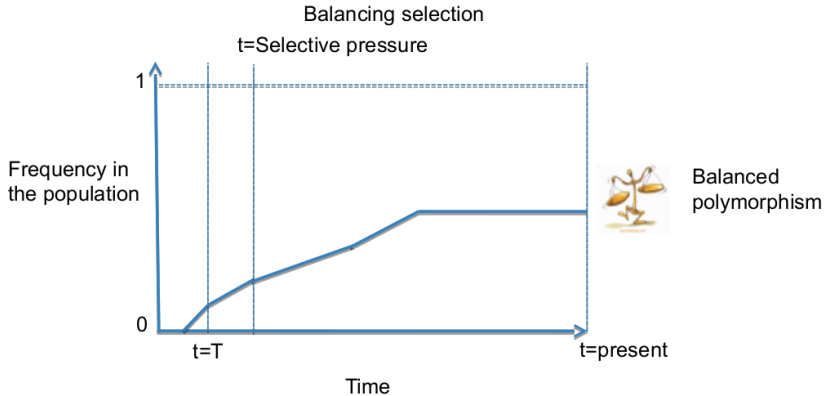
Negative selection



# Allele frequency trajectory

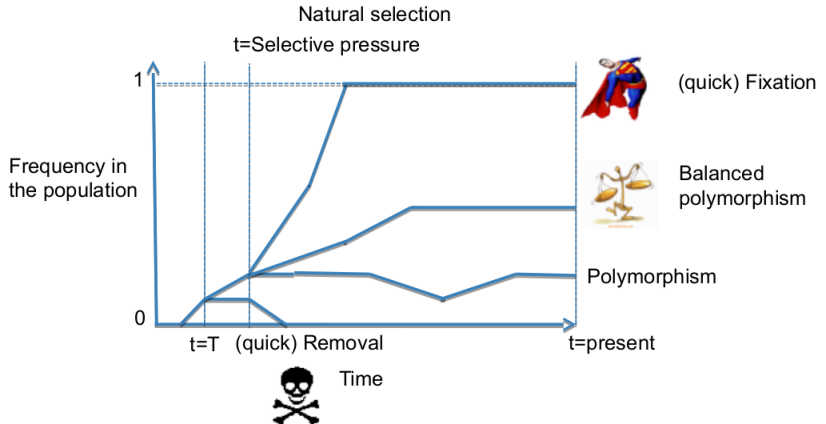


# Allele frequency trajectory





# Allele frequency trajectory



# Allele frequency trajectory - summary

## Effect of selection on alleles:

- Neutral/weak: removed, polymorphic or fixed
- Strong negative: removed or polymorphic
- Strong positive: removed, polymorphic or fixed
- Balancing: removed, polymorphic or fixed

What is "strong" selection? It depends on the effective population size.

Thus, allele frequency is (almost always) not enough to determine selection.

(slide from Anders)

# Testing for natural selection

If the simple observation of allele frequencies is not enough, what else can we do to detect signals of natural selection?

# Testing for natural selection

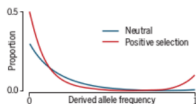
If the simple observation of allele frequencies is not enough, what else can we do to detect signals of natural selection?

- use information from the surrounding genomic region
- use information from multiple species/populations
- perform selection experiments
- use external information: candidate genes/biological knowledge, functional categories, association to phenotypes

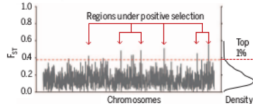
(slide from Anders)

# Common methods to detect selection

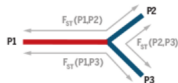
i) Change in allele frequency spectrum



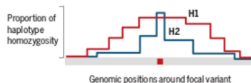
ii) Change in  $F_{ST}$  along genome



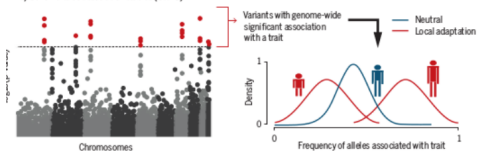
iii) Locus-specific branch length (LSBL)



iv) Extended haplotype homozygosity (EHH)

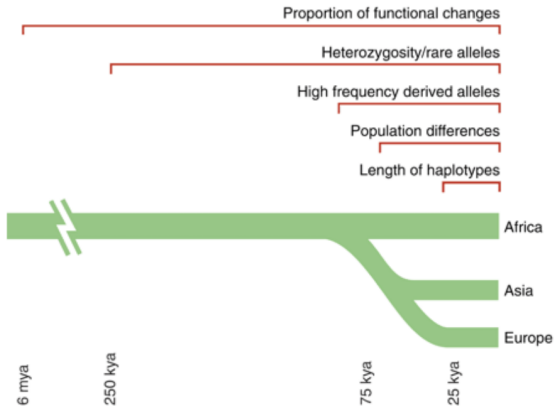


v) Genome-wide association studies (GWAS)



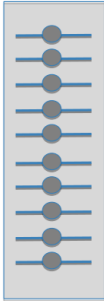
# Detect recent selection

within species / using shared variation



Sabeti et al. 2006 Science

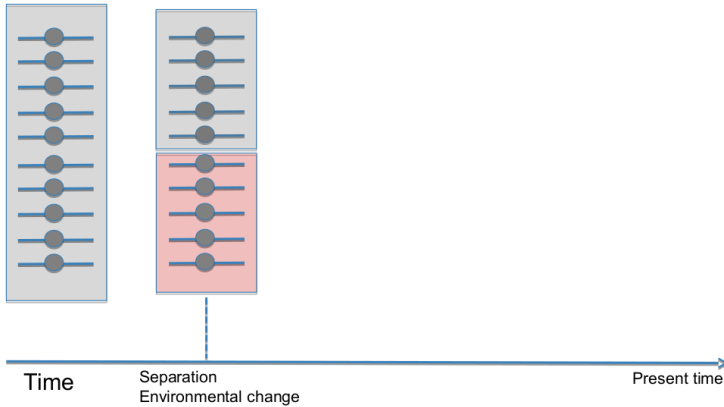
# Allele frequency differentiation



Time

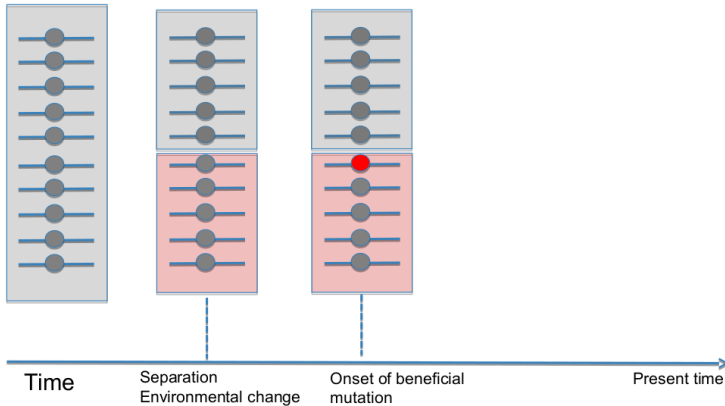
Present time

# Allele frequency differentiation

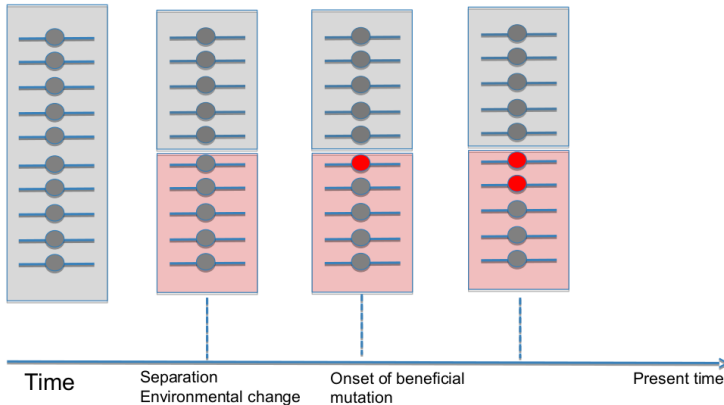




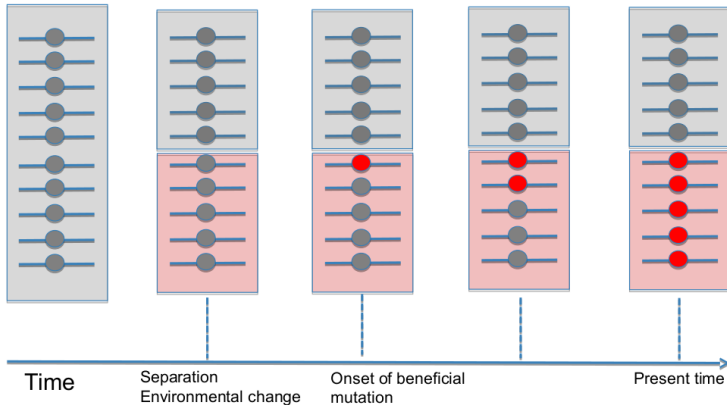
# Allele frequency differentiation



# 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

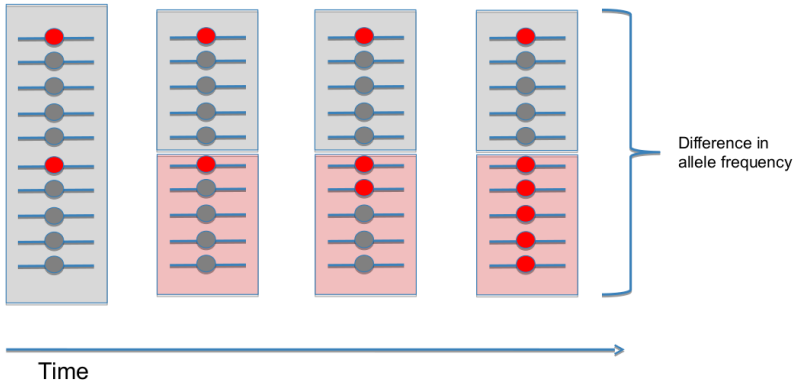
$H_W$ : average within populations

➤ if  $H_W \ll H_B$  then  $F_{ST} \sim 1$

➤ if  $H_B = 0$  then  $F_{ST} = 0$

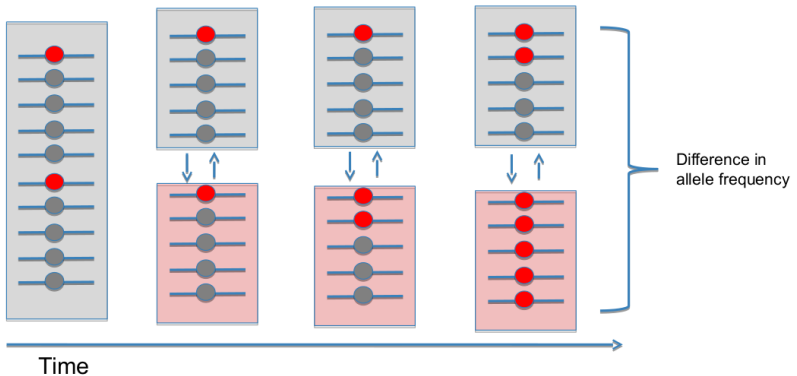
# Allele frequency differentiation

From standing variation

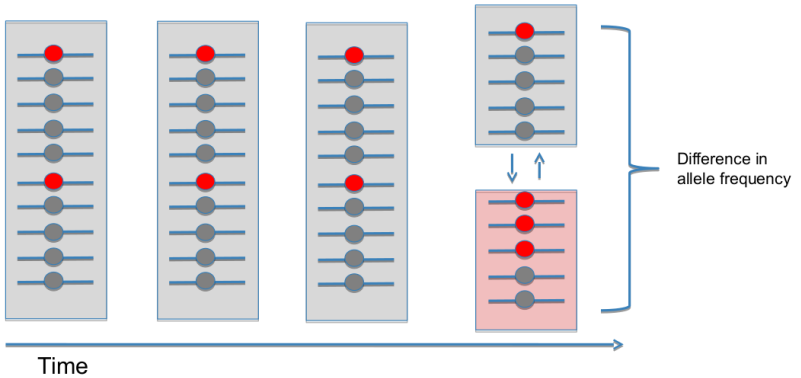


# Allele frequency differentiation

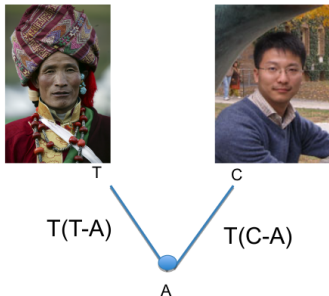
With migration



With recent divergence



# Population genetic differentiation



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



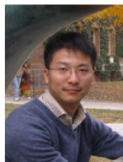
# Population genetic differentiation

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



T

T(T-A)



C

T(C-A)

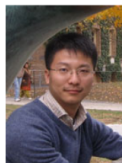
A

?



T

T(T-A)

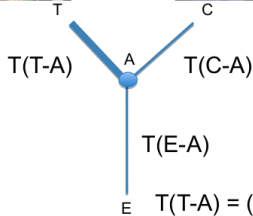


C

T(C-A)

A

# Population genetic differentiation

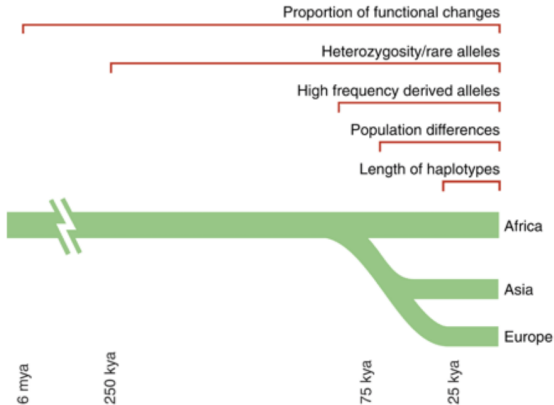


$$T(T-A-C) = -\log(1 - F_{ST}(T-C))$$



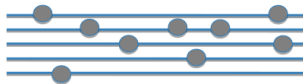
# Detect recent selection

within species / using shared variation



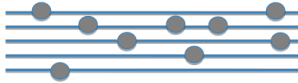
Sabeti et al. 2006 Science

## Positive selection: effect on haplotypes

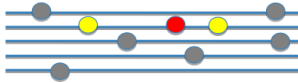


$t < T_{\text{sel}}$

## Positive selection: effect on haplotypes

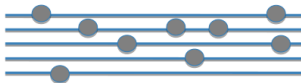


$t < T_{\text{sel}}$

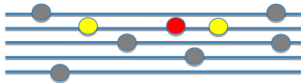


$t = T_{\text{sel}}$

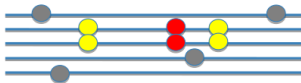
## Positive selection: effect on haplotypes



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

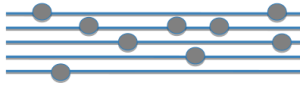


$t = T_{\text{sel}}$

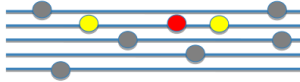


$t > T_{\text{sel}}$

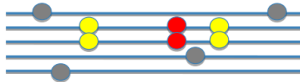
# Positive selection: effect on haplotypes



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



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



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



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

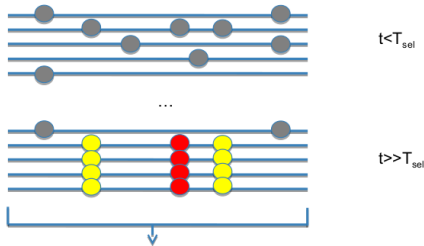
Selective sweep



Genetic hitch-hiking



# Positive selection



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

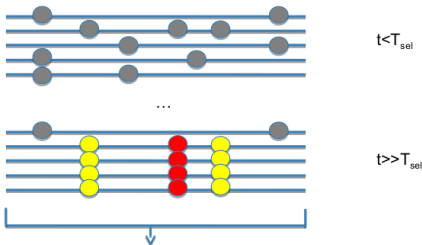
Nucleotide diversity index: Watterson's Theta  
with  $K$  SNPs and  $n$  chromosomes

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

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



# Positive selection

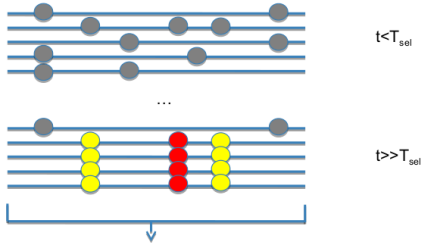


- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants

Nucleotide diversity index: average pairwise nucleotide differences ( $\pi$ ) with  $k_{ij}$  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



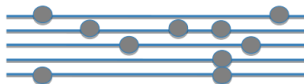
- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants (Pi)

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

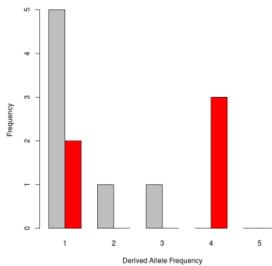


$t < T_{sel}$

...

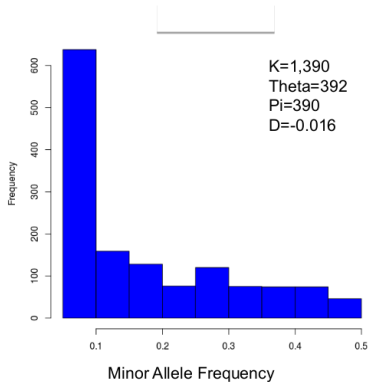


$t \gg T_{sel}$

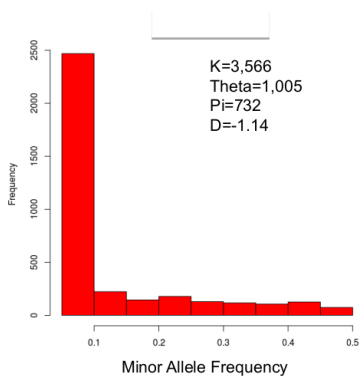


# Confounding factor

n=20; L=500kbp; no selection

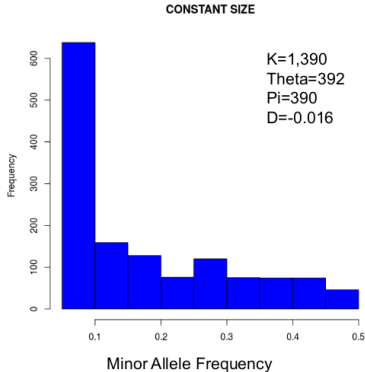


n=20; L=500kbp; no selection

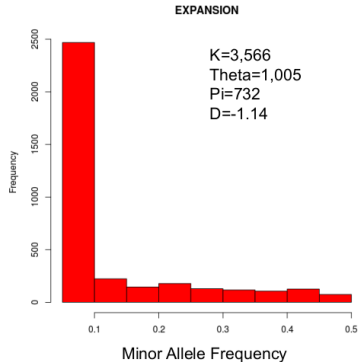


# Demography matters!

n=20; L=500kbp; no selection



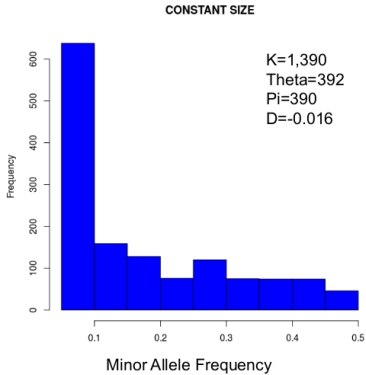
n=20; L=500kbp; no selection



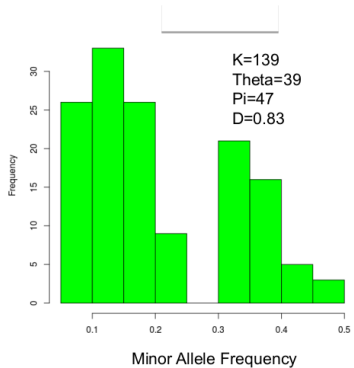
- Excess of segregating sites
- Excess of low-frequency variants
- SFS-derived summary statistics may fail to distinguish between the effects of demography and selection

# Demography matters?

n=20; L=500kbp; no selection

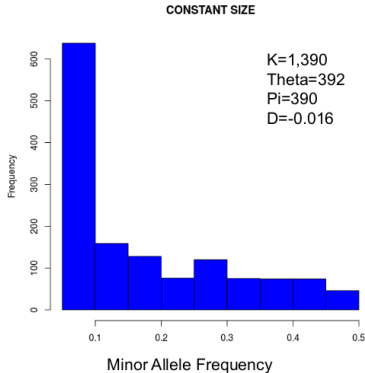


n=20; L=500kbp; no selection

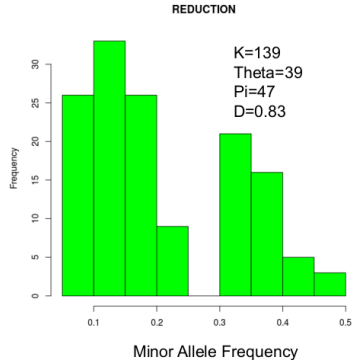


# Demography matters!

n=20; L=500kbp; no selection



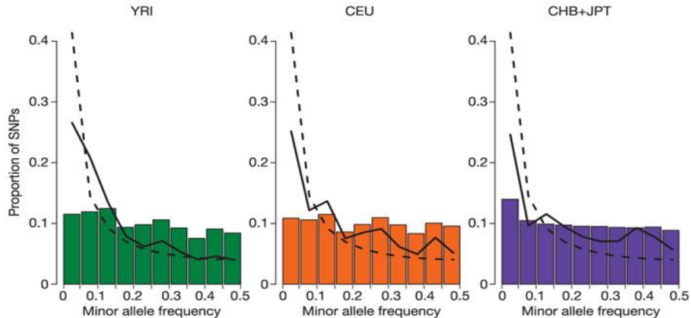
n=20; L=500kbp; no selection



- 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



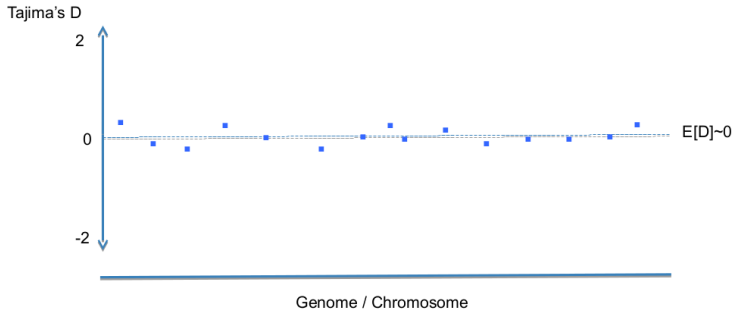
Deficiency of low-frequency variants

HapMap Consortium. Nature 2005



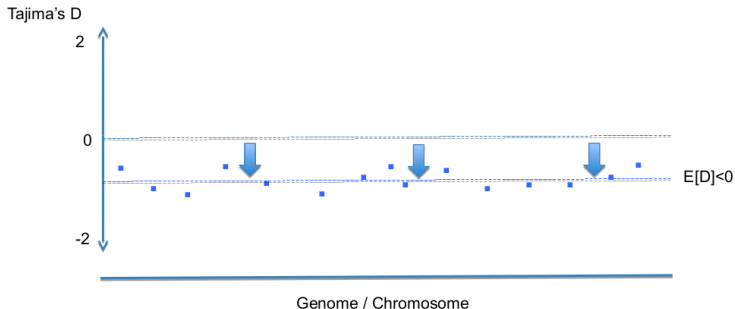
# How to take neutral confounding factors into account?

Under constant population size:



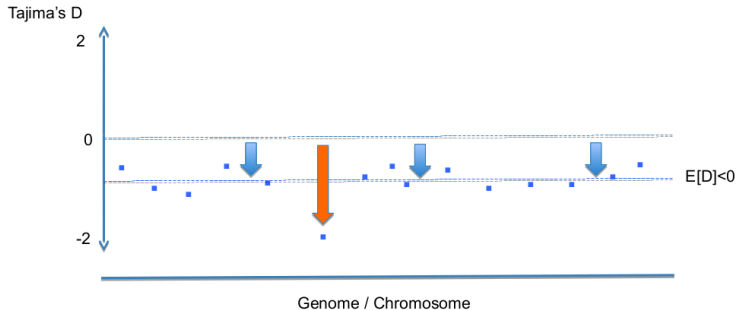
# How to take neutral confounding factors into account?

Under expanding population size:



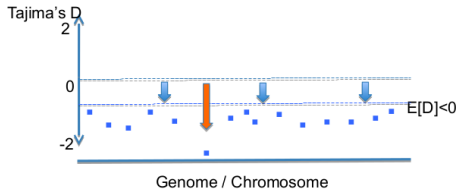
# How to take neutral confounding factors into account?

Under expanding population size and positive selection:

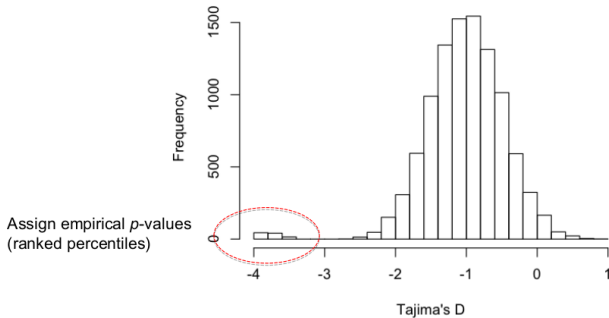


- Demography affects all loci equally, while selection changes local patterns

# Outlier approach

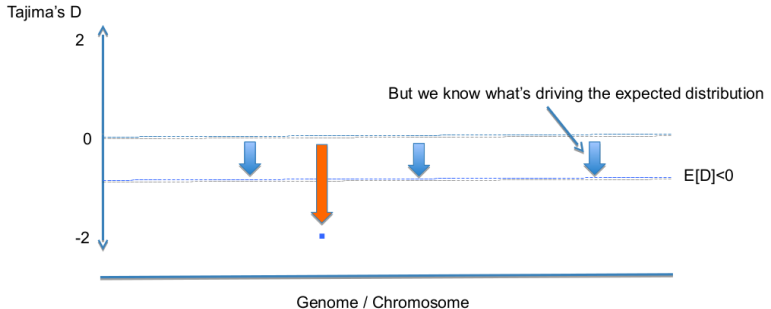


## Empirical distribution



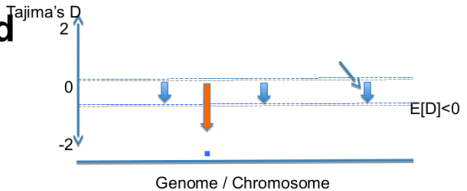
## How to take neutral confounding factors into account?

Under expanding population size and positive selection:

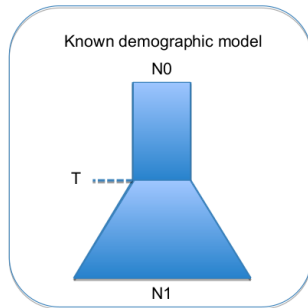
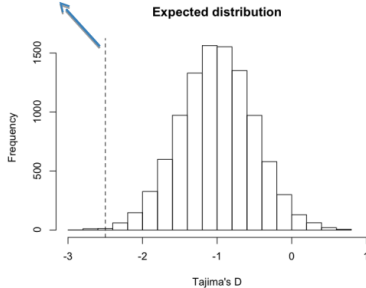


- Demography affects all loci equally, while selection changes local patterns  
What should we do if we don't have genome-wide data?

# Simulations-based approach

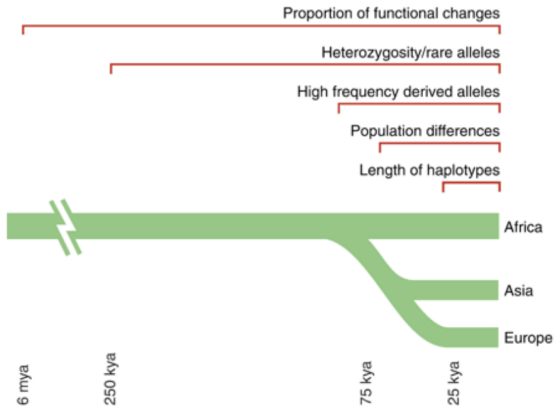


Assign  $p$ -values  
(based on ranked percentile of observed value)



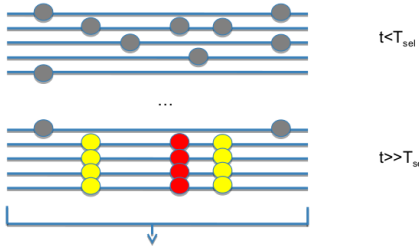
# Detect recent selection

within species / using shared variation



Sabeti et al. 2006 Science

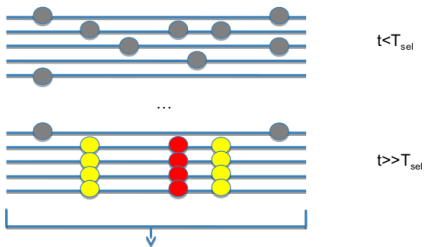
# Positive selection



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

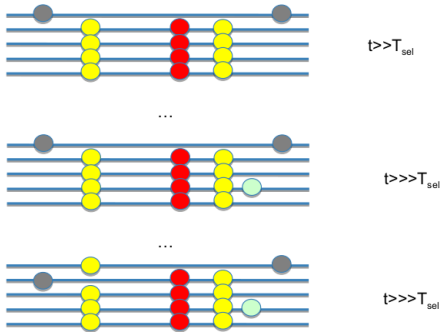


# Positive selection



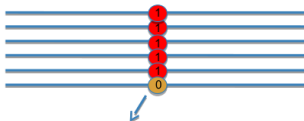
- Reduction of polymorphisms levels (Theta)
- Excess of low-frequency variants (Pi, Tajima's D, SFS)
- Extended haplotype homozygosity / Extended LD

## Extended Haplotype Homozygosity



Extended haplotype homozygosity (EHH): EHH at distance  $x$  from the core region is the probability that two randomly chosen chromosomes carry a tested core haplotype are homozygous at all SNPs for the entire interval from the core region to the distance  $x$ .

# Extended Haplotype Homozygosity

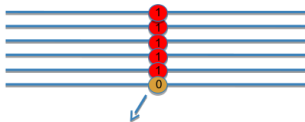


Core haplotype is 1  
(Biallelic: 0 is ancestral, 1 is derived allele)

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

Core SNP

# Extended Haplotype Homozygosity

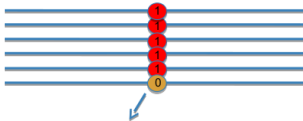


Core haplotype is 1  
(Biallelic: 0 is ancestral, 1 is derived allele)

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

Until marker  $x_i$   
(starting from  $x_0$ )

# Extended Haplotype Homozygosity

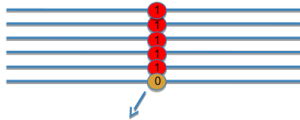


Core haplotype is 1  
(Biallelic: 0 is ancestral, 1 is derived allele)

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

Sum across all unique haplotypes  
carrying the core SNP

# Extended Haplotype Homozygosity



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

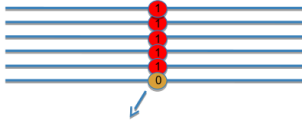
$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}}$$

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)

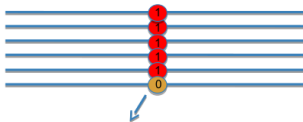
$$EHH_c(x_i) = \sum_{h \in H_c(x_i)} \left[ \frac{\binom{n_h}{2}}{\binom{n_c}{2}} \right]$$

$n_h$  is haplotype frequency of  $h$   
 $n_h$  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)

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

Sum across all unique haplotypes carrying the core SNP

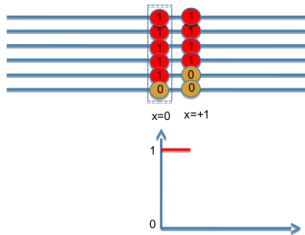
$n_h$  is haplotype frequency of  $h$

$n_c$  is haplotype frequency of the core SNP

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



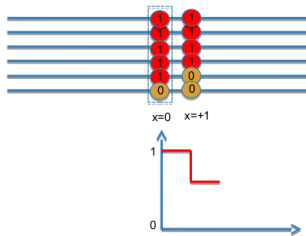
# Extended Haplotype Homozygosity



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

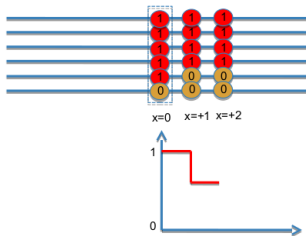
How many unique haplotypes carrying the core SNP?  
What is their frequency?

# Extended Haplotype Homozygosity



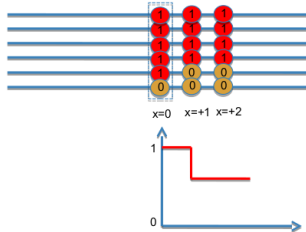
$$EHH_c(x_1 = +1) = \frac{\binom{4}{2} + \binom{1}{2}}{\binom{5}{2}} = \frac{6+0}{10} = 0.60$$

# Extended Haplotype Homozygosity



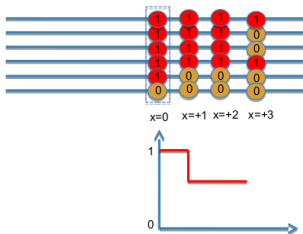
$$EHH_c(x_i = +2) = ?$$

# Extended Haplotype Homozygosity



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

# Extended Haplotype Homozygosity



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

How many unique haplotypes carrying the core SNP?

What is their frequency?

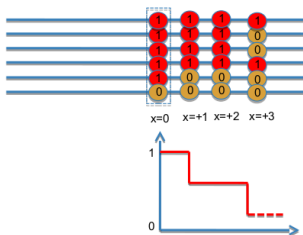
1111 with freq=2

1110 with freq=2

1000 with freq=1

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

# 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?

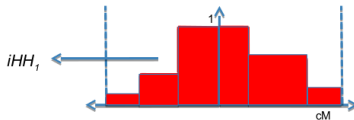
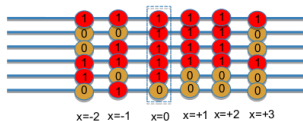
1111 with freq=2

1110 with freq=2

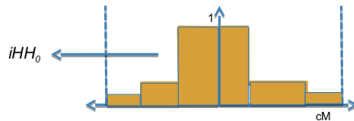
1000 with freq=1

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

# Integrated Haplotype Score



Integrated haplotype homozygosity ( $iHH$ )



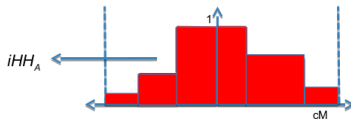
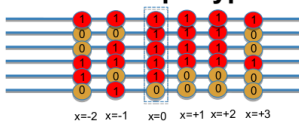
Integrated haplotype score:

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

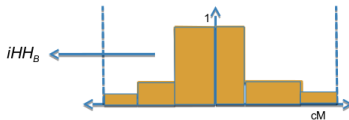


Genome-wide normalization in frequency bins  
(to mean=0 and sd=1)

# Cross-population Extended Haplotype Homozygosity



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



Integrated haplotype score:

$$XP-EHH = \ln(iHH_A/iHH_B)$$

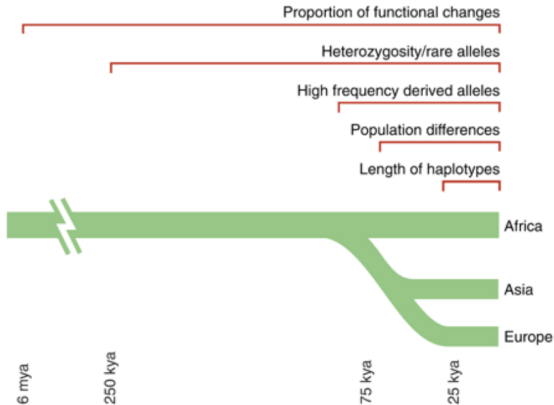


Genome-wide normalization in frequency bins  
(to mean=0 and sd=1)



# Detect recent selection

within species / using shared variation



Sabeti et al. 2006 Science

# Intended Learning Outcomes

At the end of this session you are now be able to:

- list commonly used methods to detect selection
  - calculate various summary statistics
  - understand main confounding factors to neutrality tests
  - assess statistical significance of tests
-