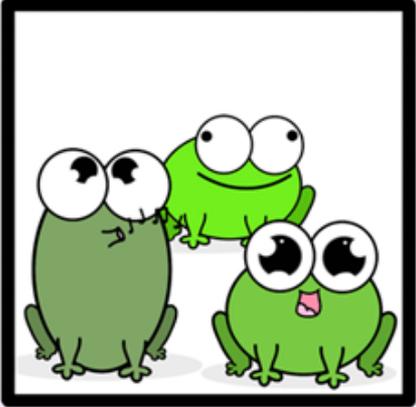


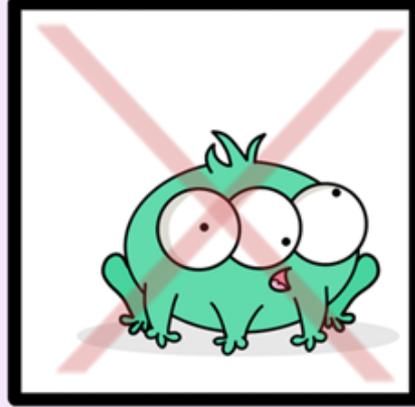
A Lapalissade model: if everything is stable, nothing changes

Assumptions of Hardy-Weinberg Equilibrium

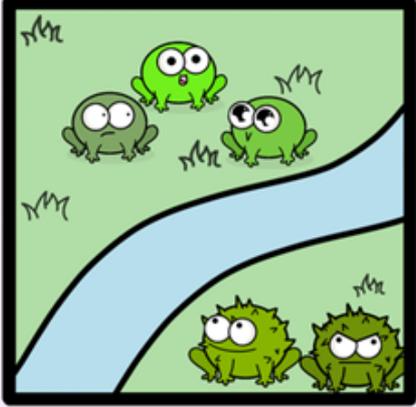
1. No selection



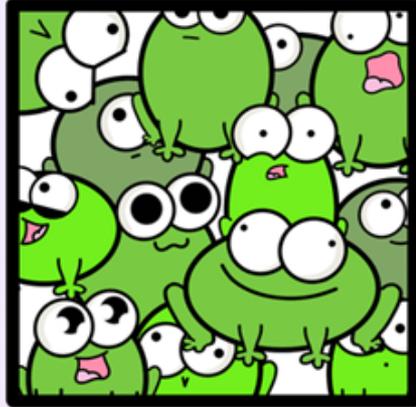
2. No Mutation



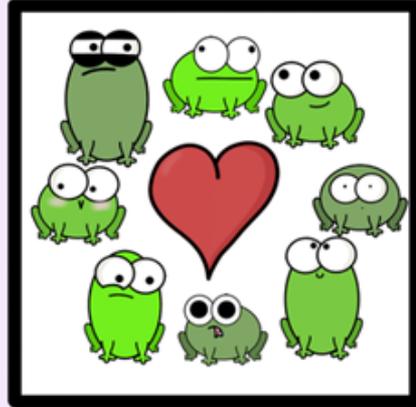
3. No Migration



4. Large Population

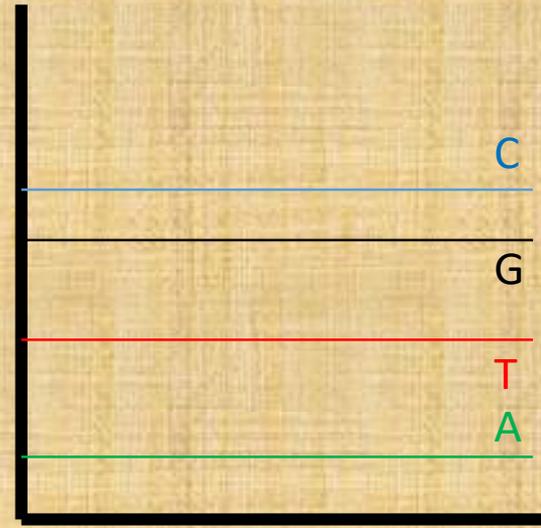


5. Random Mating



@AmoebaSisters

Allele frequency in a pop



What if???



The neutral theory paradigm

Most variation at the molecular level does not affect fitness (not due to selection)

Genetic variation is best explained by stochastic processes

Detecting selection = testing allele frequencies against the null expectation (neutral evolution)

Plan

Demography and other factors influencing genetic variation

Bottlenecks, migration, mating system etc

N_e and efficacy of selection

LD/Recombination rate

Testing for demographic confounding factors

Nucleotide diversity, heterozygosity

Testing for demographic events: Tajima's D

Estimating the "genetic load" and efficacy of selection: SFS and DFE

Detecting natural selection vs neutrality (and demography) in the genome

types of selection

selective sweeps

The MK test

DFE and SFS

Fst scans

Adaptomics methods involving additional aspects than just allele frequencies

Genetic vs phenotypic association: QTLs and GWAS

Environmental associations

Dramatic demographic events shaping genetic variation

Bottle neck

Founder effect

Population size

Self-incompatible



C. grandiflora

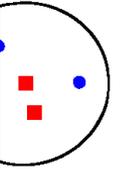
Loss of self-incompatibility
= strong bottleneck



C. rubella



C. orientalis



Reduced genetic diversity

Increased homozygosity

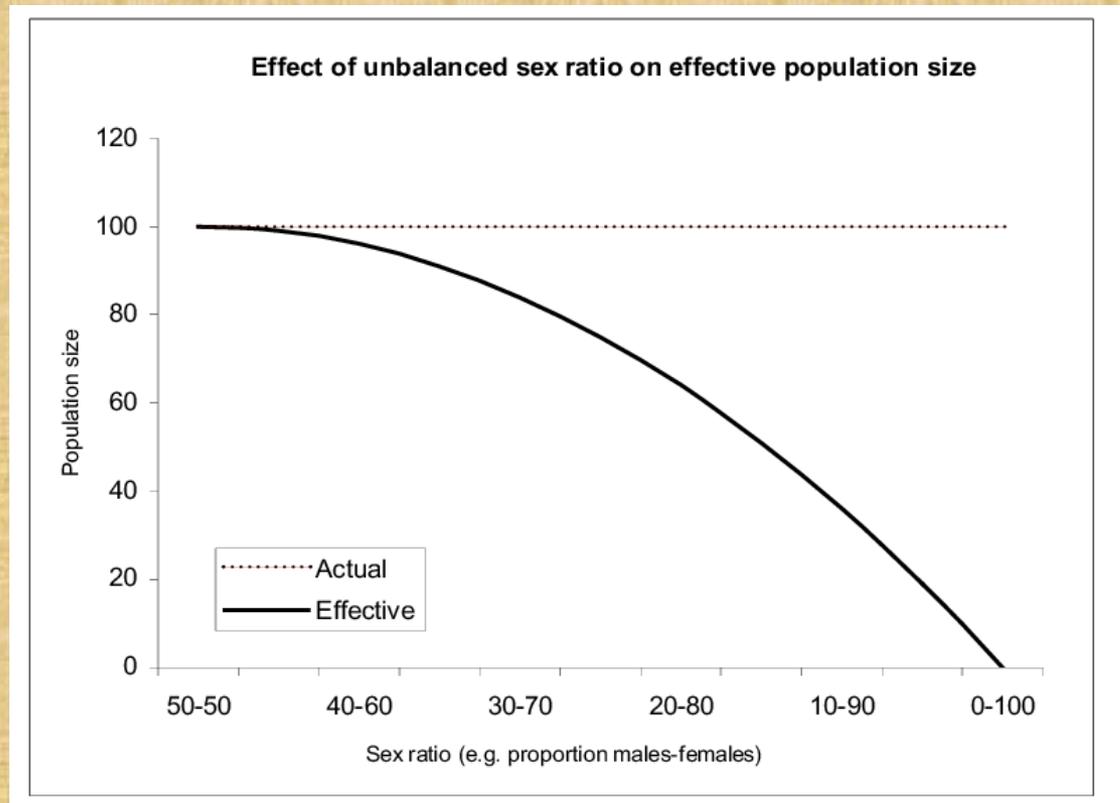
Decreased frequency of rare alleles

Effective population size, a critical factor in pop genetics

The effective population size (N_e)

= the size of an idealised population , which would give rise to the same rate of inbreeding and the same rate of change in allele frequencies

N_e of humans?



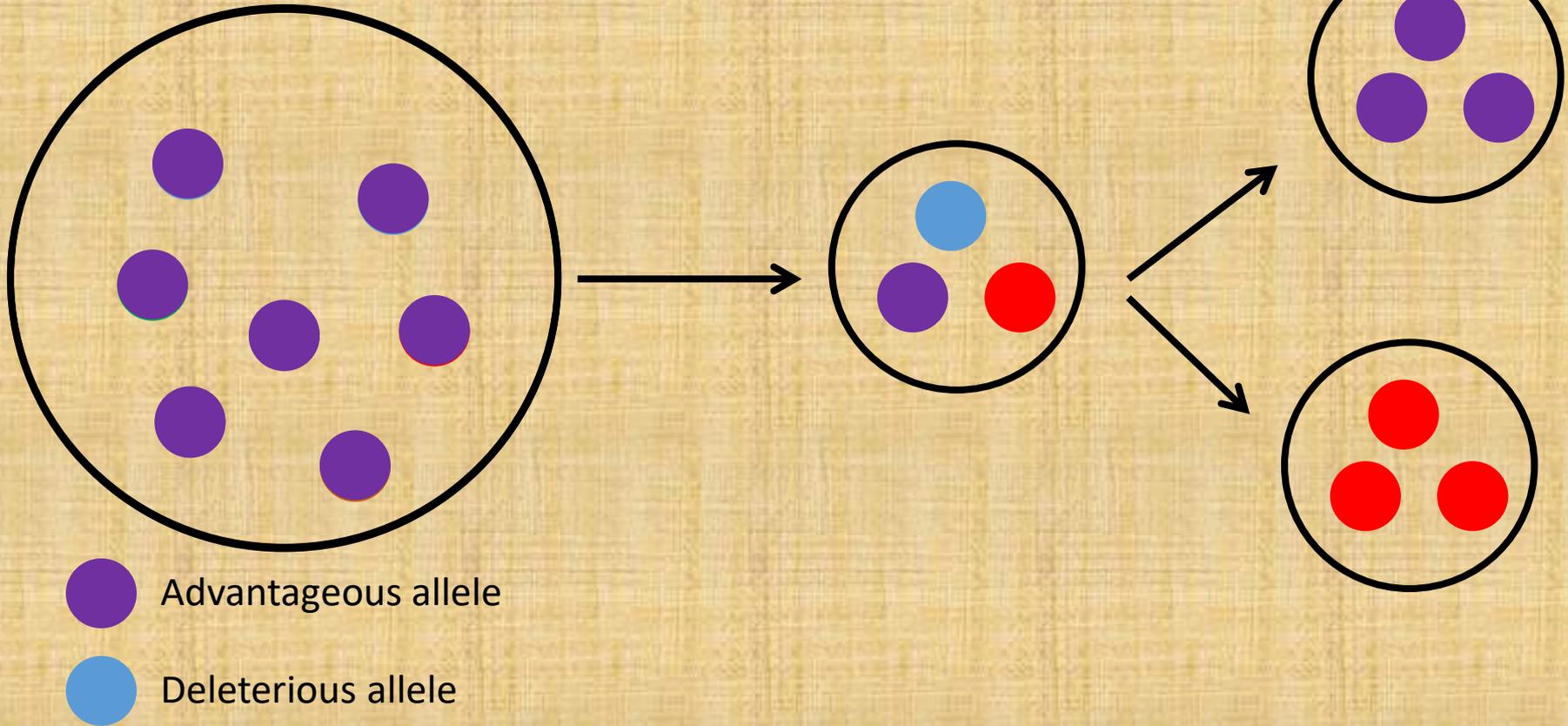
Effective population size and efficacy of selection

+++++ selection

----- stochasticity (drift)

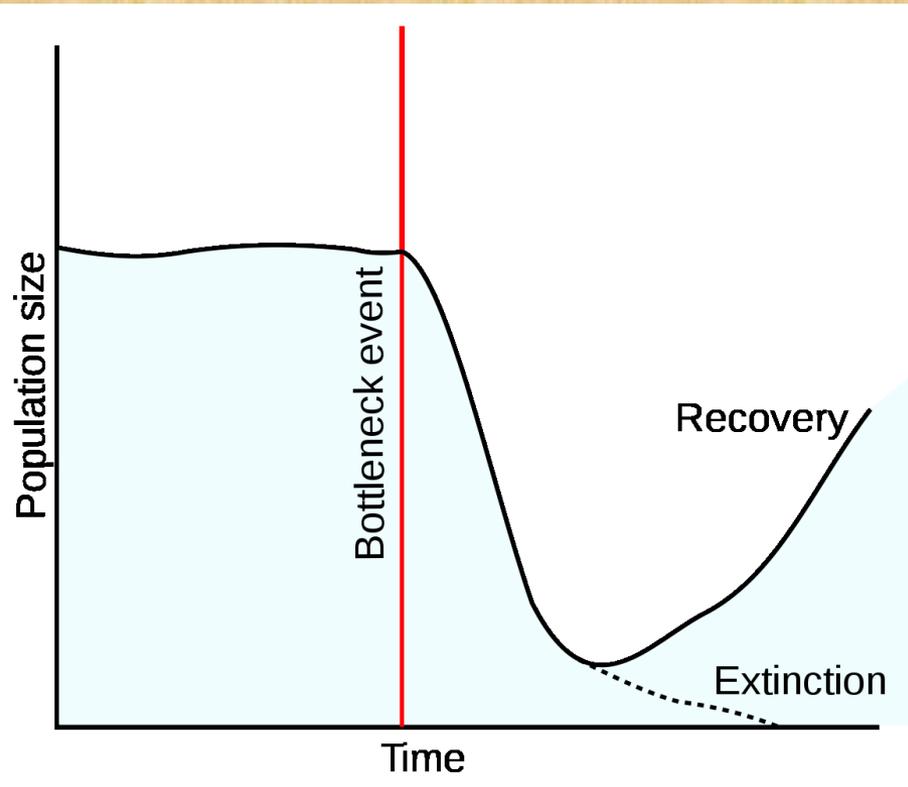
----- selection

+++++ stochasticity (drift)

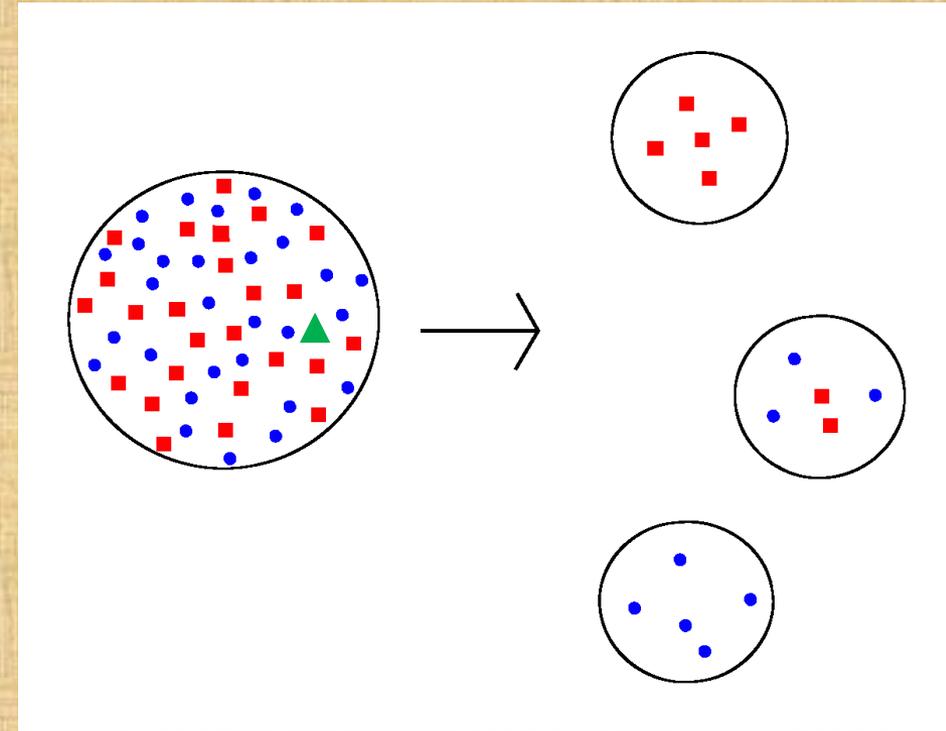


Dramatic demographic events shaping genetic variation

Bottle neck



Founder effect

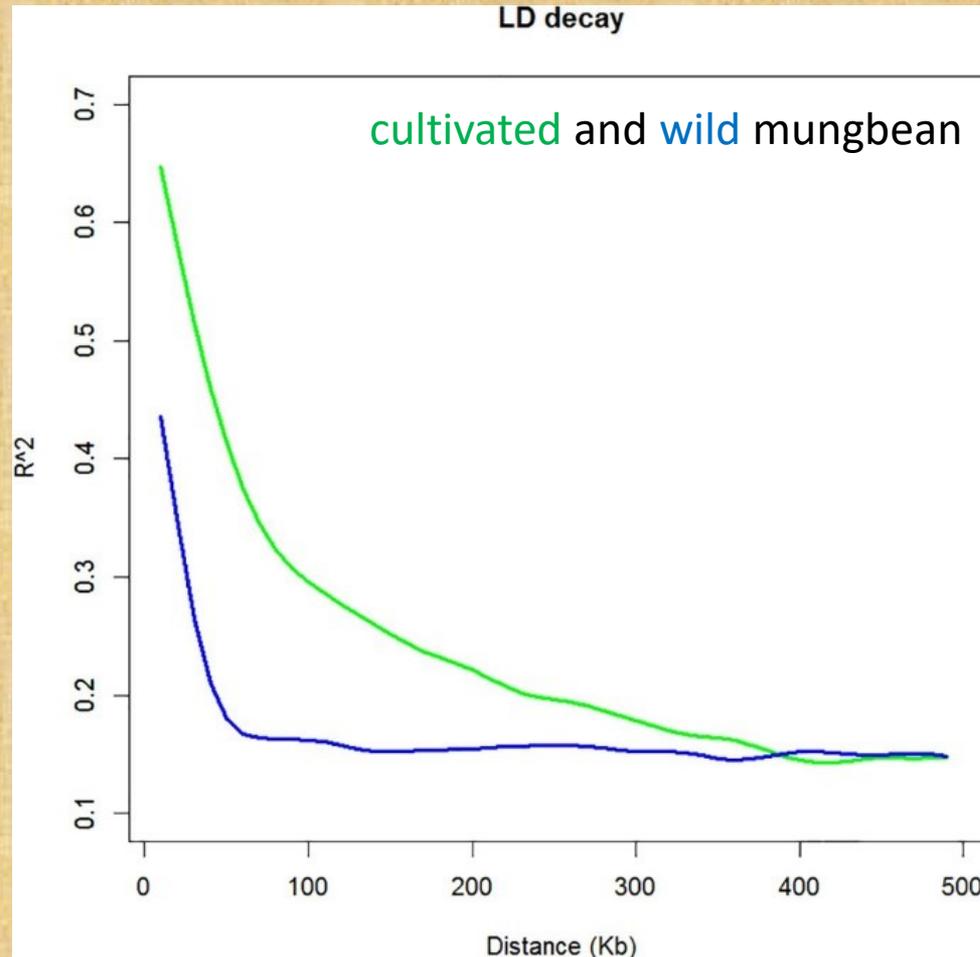


Reduced genetic diversity

Increased homozygosity

Decreased frequency of rare alleles

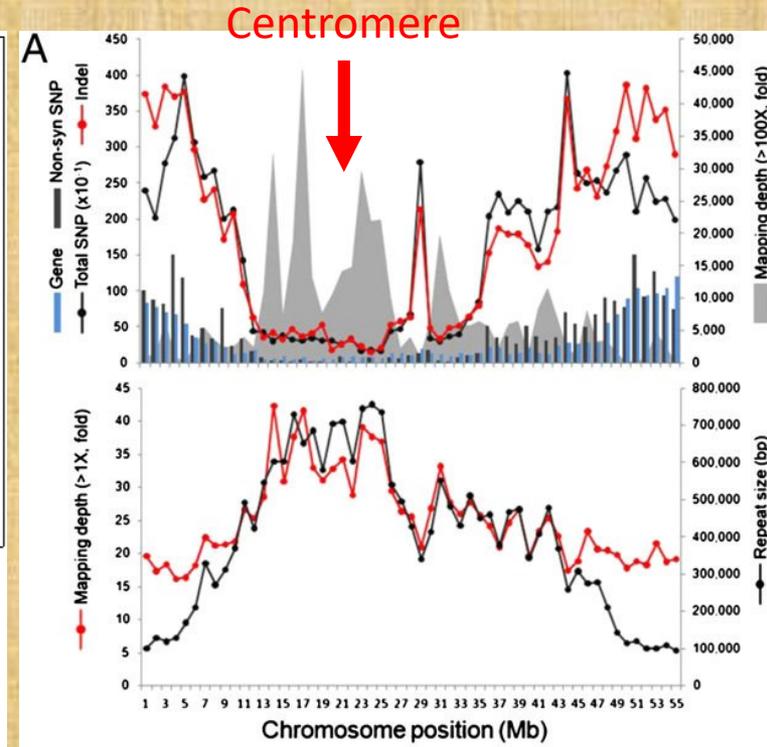
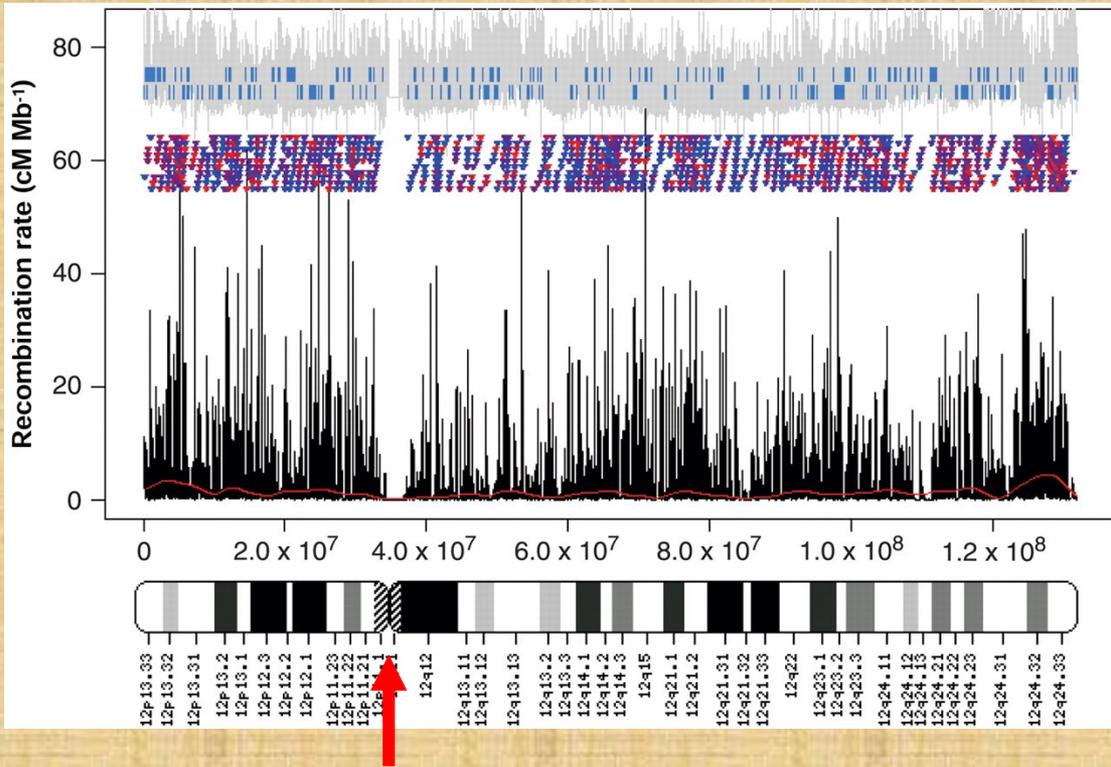
Linkage disequilibrium: a marker for effective recombination



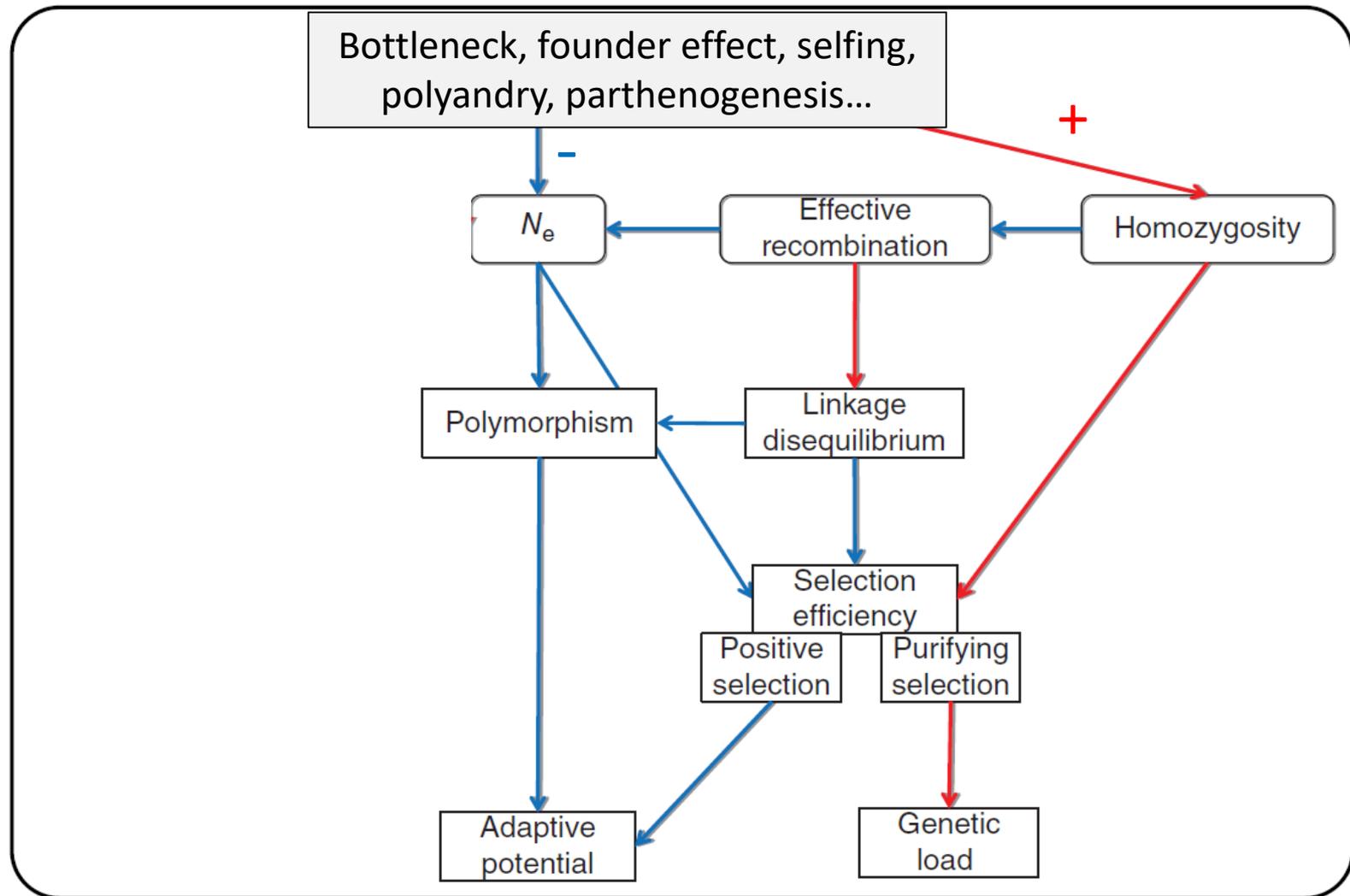
“Effective” recombination depends on:

- the actual recombination process
- genetic diversity

Recombination rate varies across genome; Influences genetic diversity



Demographic events play a major role in population genetics



Plan

Demography and other factors influencing genetic variation

Bottlenecks, migration, mating system etc

N_e and efficacy of selection

LD/Recombination rate

Testing for demographic confounding factors

Nucleotide diversity, heterozygosity

Testing for demographic events: Tajima's D

Estimating the "genetic load" and efficacy of selection: SFS and DFE

Detecting natural selection vs neutrality (and demography) in the genome

types of selection

selective sweeps

The MK test

DFE and SFS

Fst scans

Adaptomics methods involving additional aspects than just allele frequencies

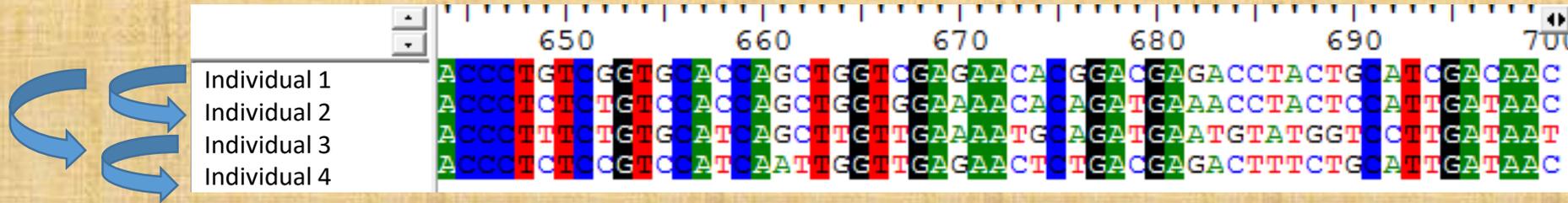
Genetic vs phenotypic association: QTLs and GWAS

Environmental associations

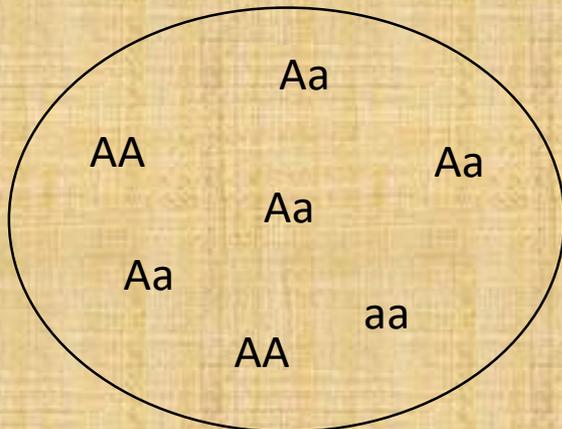
Nucleotide diversity, heterozygosity, can indicate demographic processes going on

Nucleotide diversity

Π , average pairwise differences



Heterozygosity



$$p = f(A); q = f(a)$$

$$H_{observed} = f(Aa) = ?$$

$$H_{expected} = 1 - \sum_{i=1}^n (f_i)^2 \text{ for } n \text{ alleles}$$

$$H_e = 1 - (p^2 + q^2) = ?$$

Tajima's D

$D \approx$ observed genetic variation – expected genetic variation (neutrality)

$D =$ average pairwise difference among individuals – expected number of variable sites for a given number of individuals

$$D = \frac{\sum_{i < j} d_{ij}}{n(n-1)/2} - \frac{S}{\sum_{i=1}^{n-1} 1/i}$$

Sum of pairwise differences \rightarrow $\sum_{i < j} d_{ij}$

No of individuals \rightarrow $n(n-1)/2$

No of variable sites \rightarrow S

	C	A	T	T	G	G
$n = 4$	C	A	T	T	G	G
	C	A	T	T	G	G
$S = 2$	C	A	T	T	A	G

$$D = [(2+2+2)/(4*3/2)] - [2/(1/1+1/2+1/3)]$$

$$D = 1 - 1.09 = -0.09 \approx 0$$

Observed variation = expected
 \rightarrow Population at equilibrium

$D \approx$ observed genetic variation – expected genetic variation (neutrality)

D = average pairwise difference among individuals – expected number of variable sites for a given number of individuals

$$D = \frac{\sum_{i < j} d_{ij}}{n(n-1)/2} - \frac{S}{\sum_{i=1}^{n-1} 1/i}$$

Sum of pairwise differences \rightarrow $\sum_{i < j} d_{ij}$

No of individuals \rightarrow $n(n-1)/2$

No of variable sites \rightarrow S

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

$D > 0$

More variation than expected considering the no of individuals

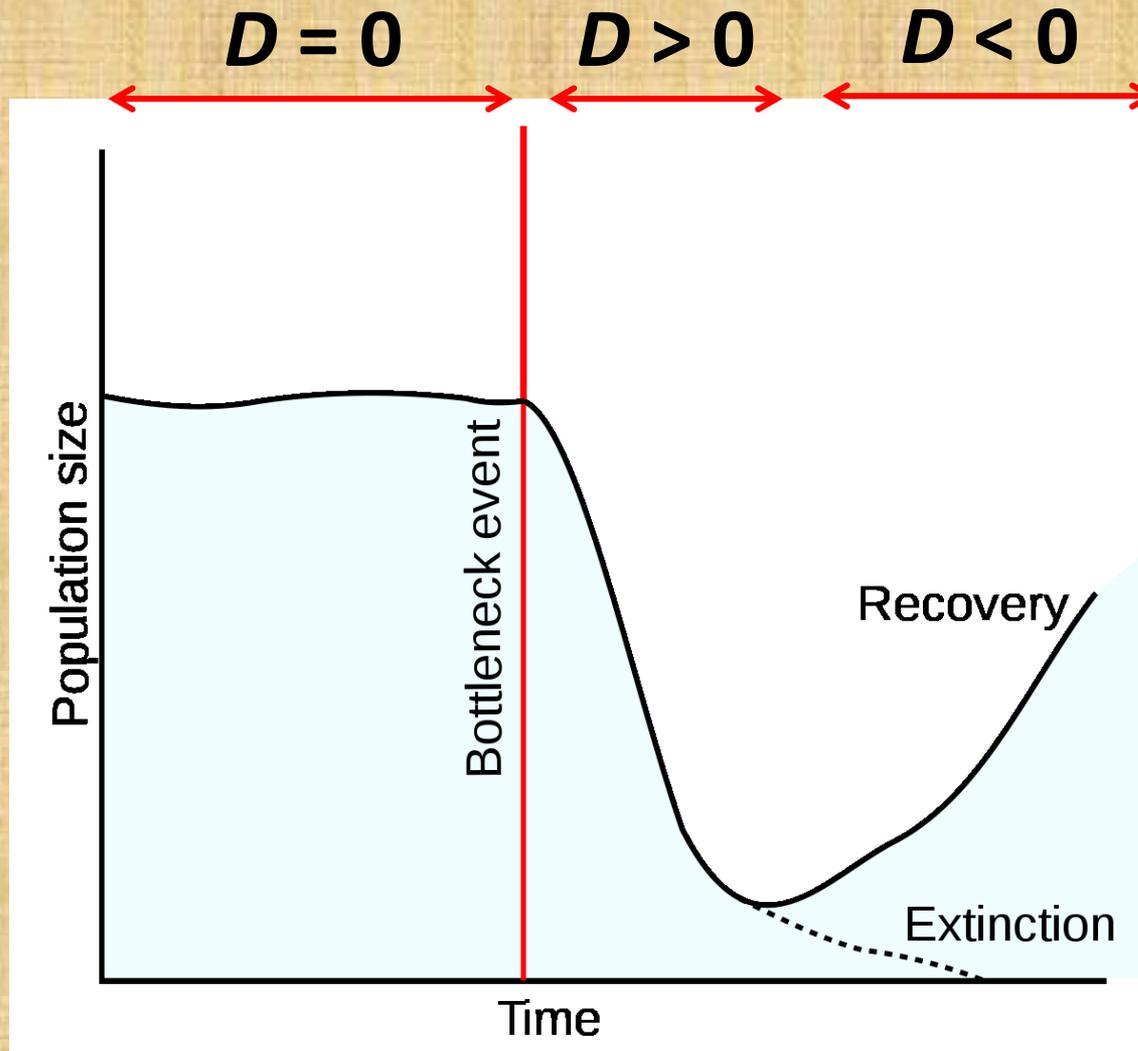
T	G	G	C	C	T
T	G	G	C	T	A
T	G	G	C	T	A
T	G	G	C	T	A
A	G	G	T	T	A

$D < 0$

Less variation than expected considering the no of individuals

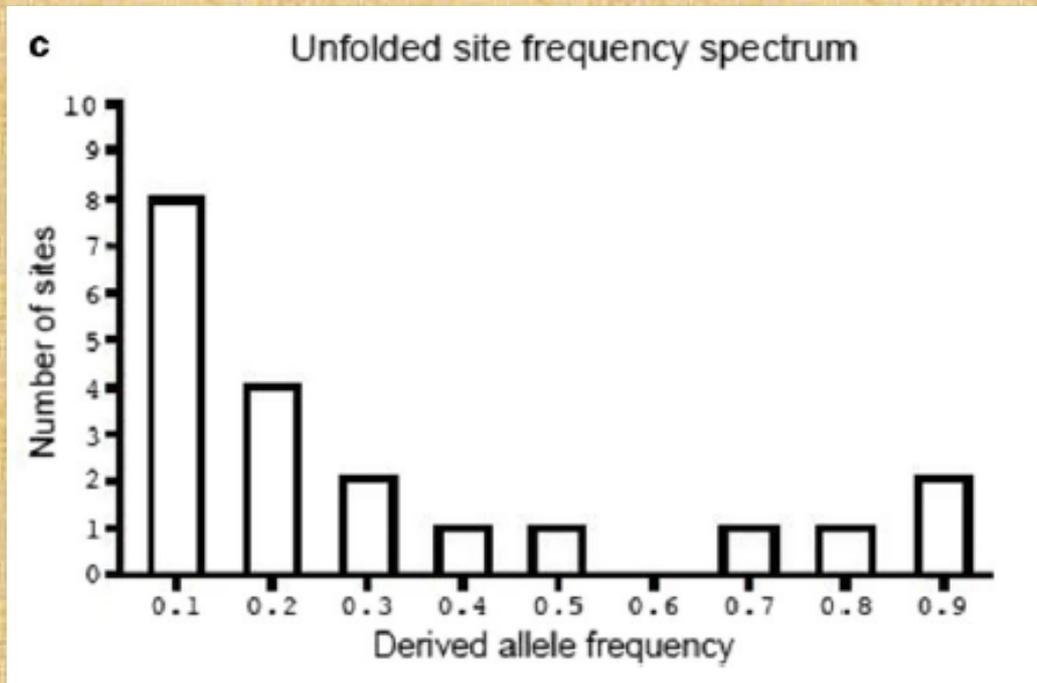
T	A	T	T	G	G	C
A	A	T	T	G	G	T
A	A	T	T	G	G	T
A	A	T	T	G	G	T
A	A	T	T	G	G	T

Tajima's D and demographic events



Estimating the Site Frequency Spectrum

outgroup sequence	CCATGATCTCCTTGAGTGGG	—	Outgroup data
sample 1	CCAAGCTCCTCTTGAGGGAG	}	Polymorphism data
sample 2	CTAAGCACCCCTGGAGGGAC		
sample 3	GTATACTCCCTCTGAGTTGG		
sample 4	CTAAGCTCTCCTTGAGGGAG		
sample 5	GTAAGATCCTCTGGAGGGAG		
sample 6	CTAAGATCCCTTCTGTGAG		
sample 7	CTGAGCTGCCTTTCAGTGAG		
sample 8	CTAAACTCCCTTGAGTTAG		
sample 9	CTAAACTCCCTTTGATTGGG		
sample 10	GTAAGCTCCCTTTGAGGGG		
	3 1 1 1 3 2 1 1 1 2 4 1 2 2 1 1 5 2 3 1	—	Minor allele count
	3 9 1 1 3 8 1 1 9 2 4 1 2 2 1 1 5 2 7 1	—	Derived allele count



The SFS can tell about demographic events

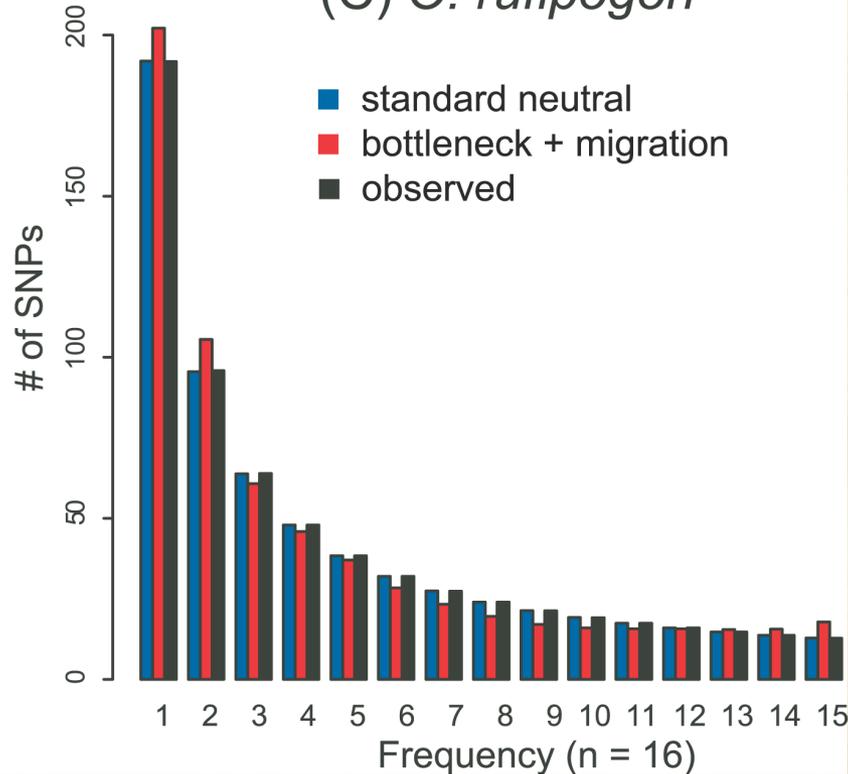


Wild rice

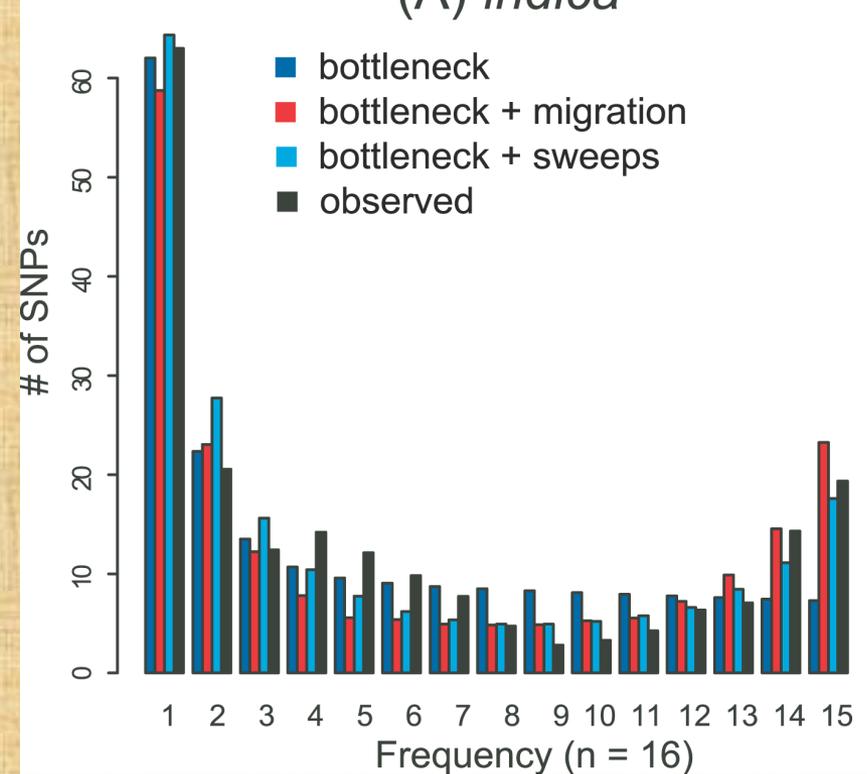


Cultivated rice

(C) *O. rufipogon*



(A) *indica*

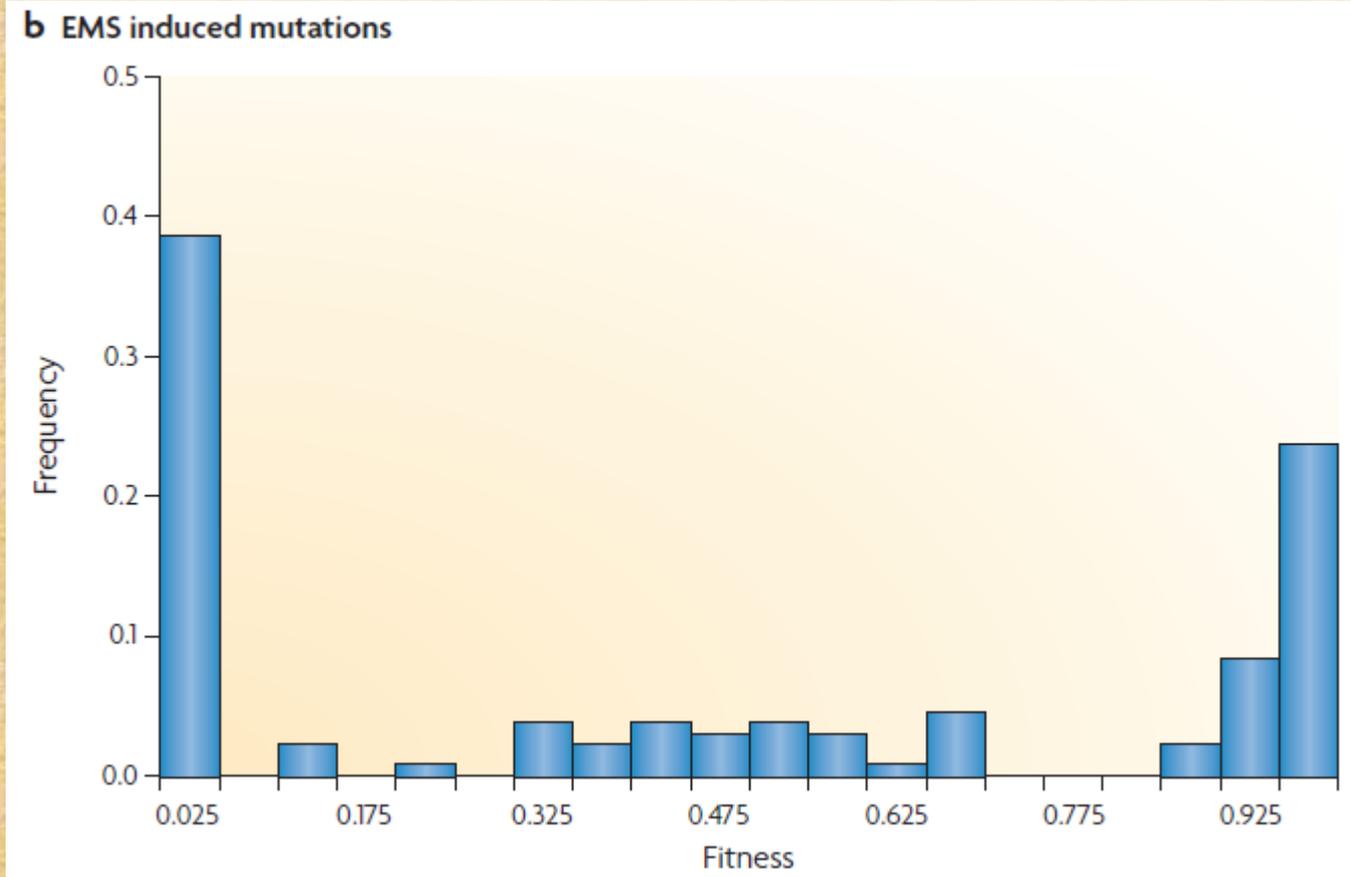


Bottleneck: few high frequency mutations



S. cerevisiae

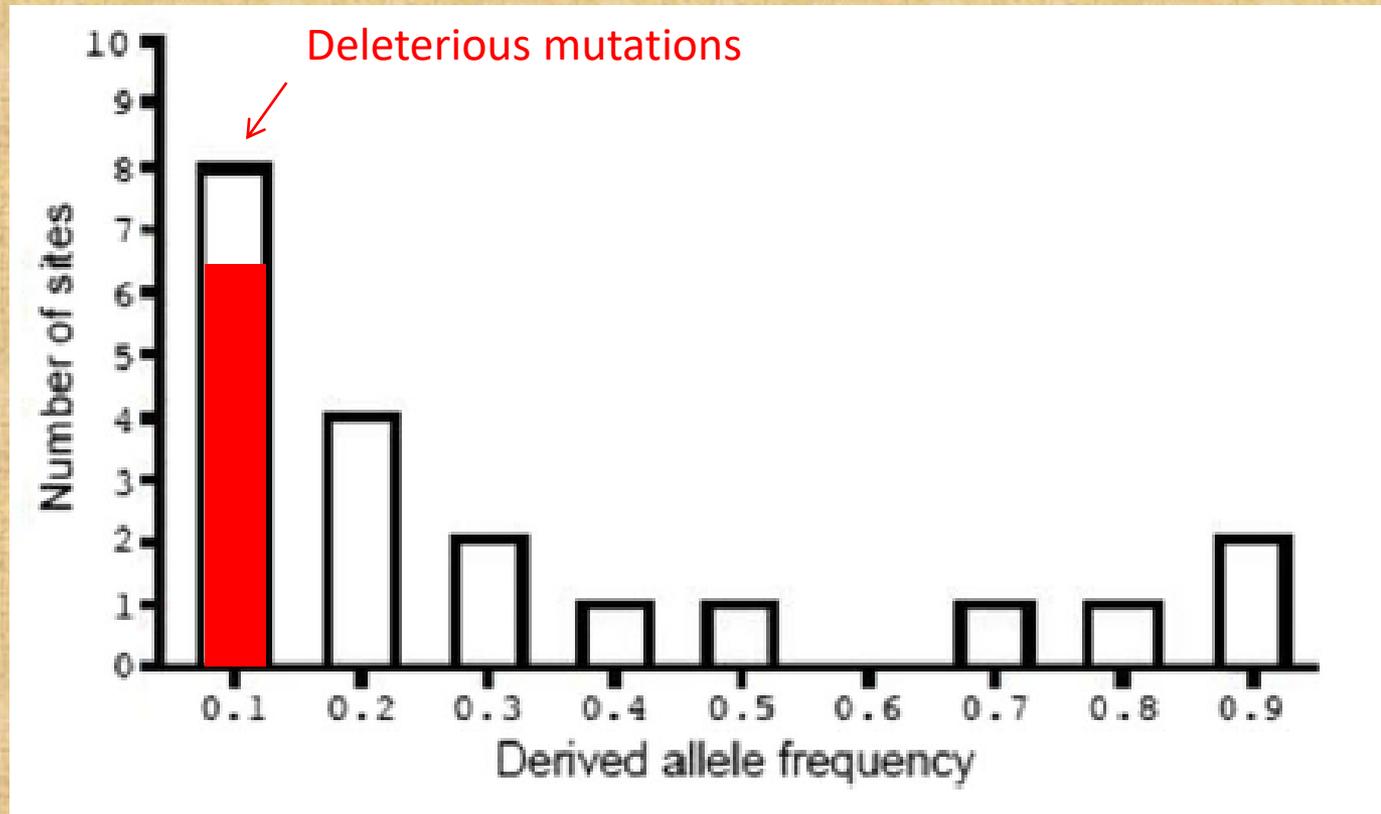
Most new mutations are deleterious



It is not always possible to determine the effect of mutations on fitness

Inference from allele frequencies in pop genomics

Estimating the Site Frequency Spectrum

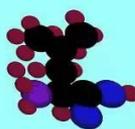


Syn vs non syn. mutations

		SECOND NUCLEOTIDE IN CODON				
		U	C	A	G	
FIRST NUCLEOTIDE IN CODON	U	UUU PHE UUC UUA LEU UUG	UCU SER UCC UCA UCG	UAU TYR UAC UAA STOP UAG	UGU CYS UGC UGA STOP UGG TRP	THIRD NUCLEOTIDE IN CODON
	C	CUU LEU CUC CUA CUG	CCU PRO CCC CCA CCG	CAU HIS CAC CAA GLN CAG	CGU ARG CGC CGA CGG	
	A	AUU ILE AUC MET AUA AUG	ACU THR ACC ACA ACG	AAU ASN AAC AAA LYS AAG	AGU SER AGC AGA ARG AGG	
	G	GUU VAL GUC GUA GUG	GCU ALA GCC GCA GCG	GAU ASP GAC GAA GLU GAG	GGU GLY GGC GGA GGG	

LEUCINE IS ENCODED BY 6 CODONS:

UUA
UUG
CUU
CUC
CUA
CUG



Synonymous: same AA encoded; non-synonymous: different AA

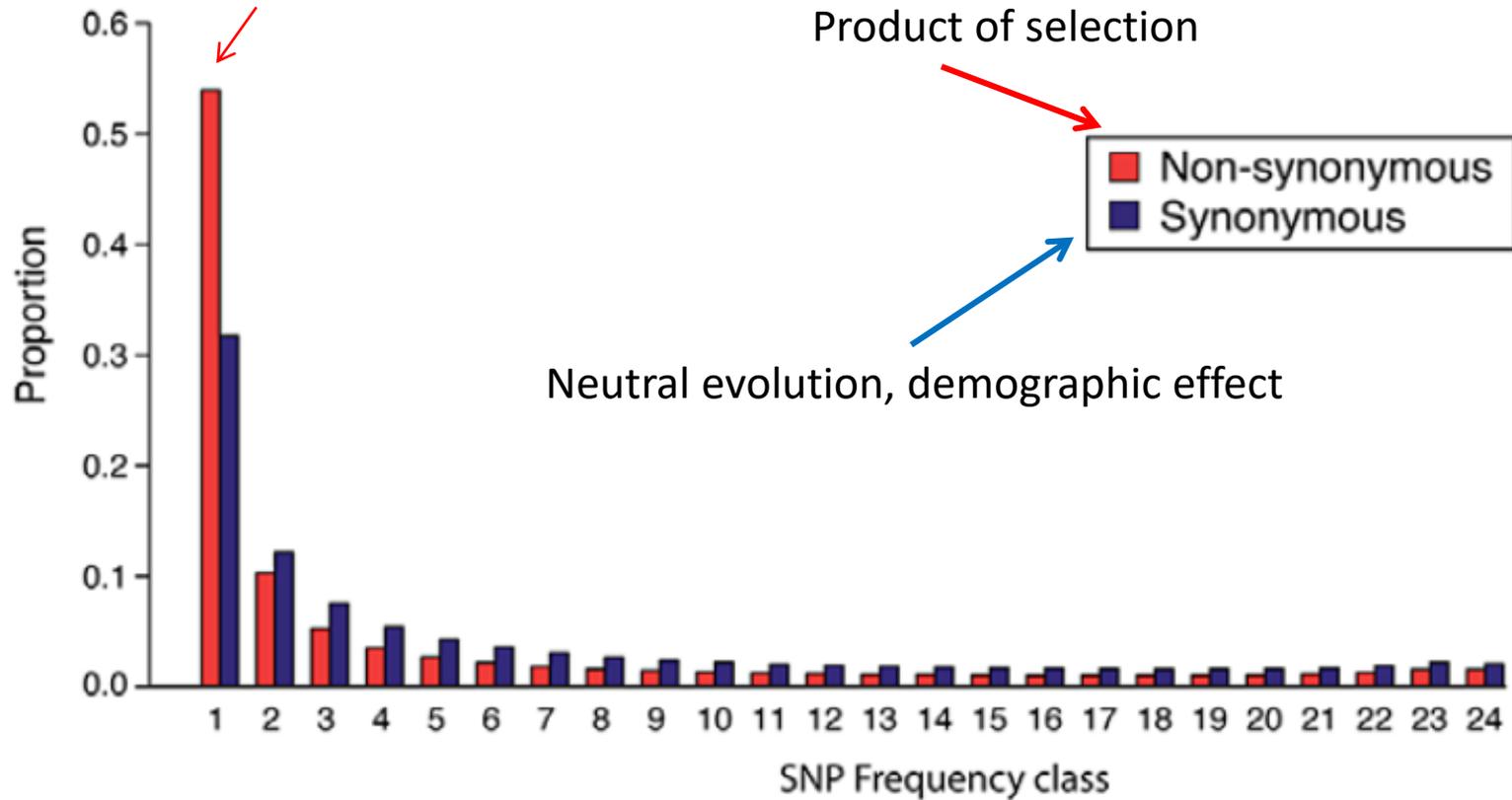
Ratio between syn and non-syn will tell about a given gene evolution:

syn = neutral mutations, due to population demography, mutation rate...

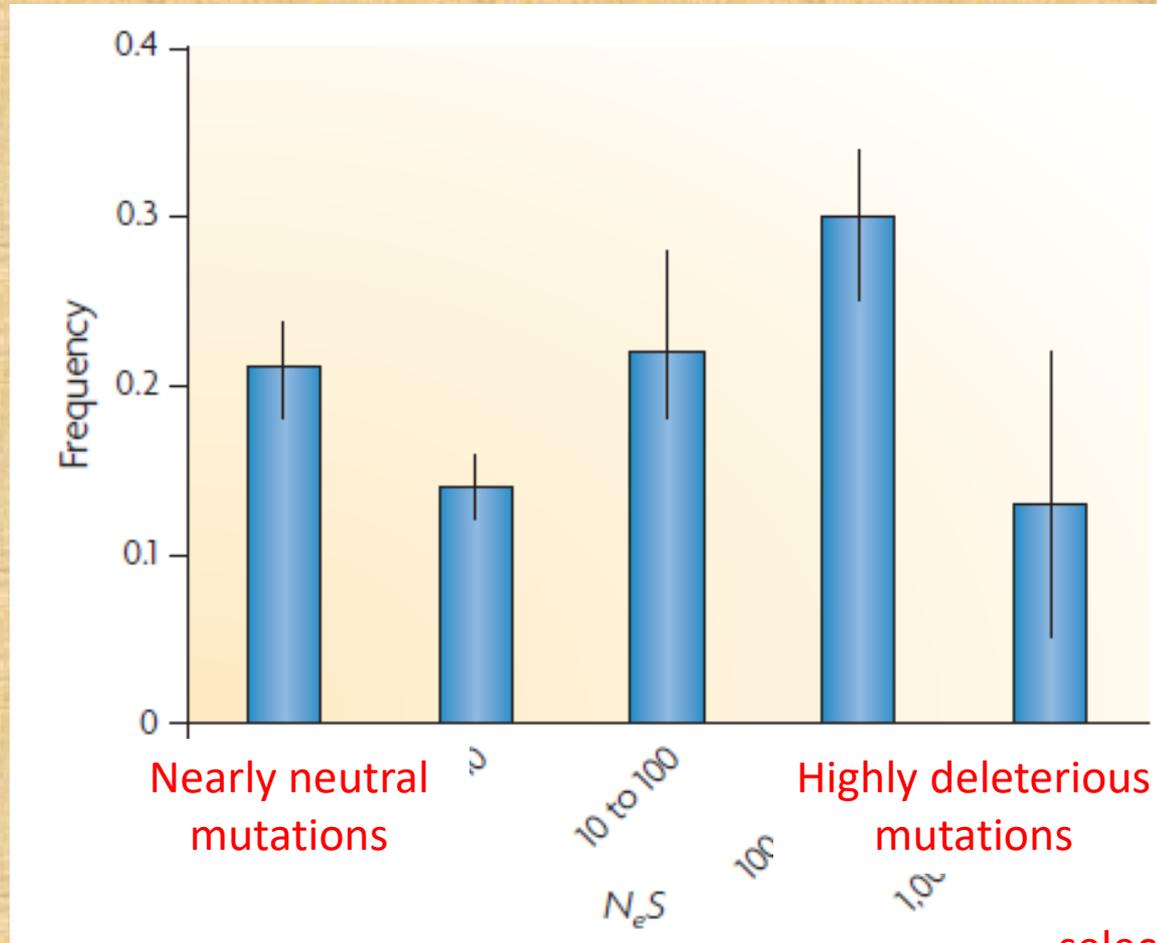
non syn = important for prot function; can be purged (if deleterious) or selected (if advantageous)

SFS to infer the frequency of deleterious mutations

Deleterious mutations



The distribution of fitness effect



selection coefficient

The strength of selection: $N_e S$

Plan

Demography and other factors influencing genetic variation

Bottlenecks, migration, mating system etc

N_e and efficacy of selection

LD/Recombination rate

Testing for demographic confounding factors

Nucleotide diversity, heterozygosity

Testing for demographic events: Tajima's D

Estimating the "genetic load" and efficacy of selection: SFS and DFE

Detecting natural selection vs neutrality (and demography) in the genome

types of selection

selective sweeps

the MK test

DFE and SFS

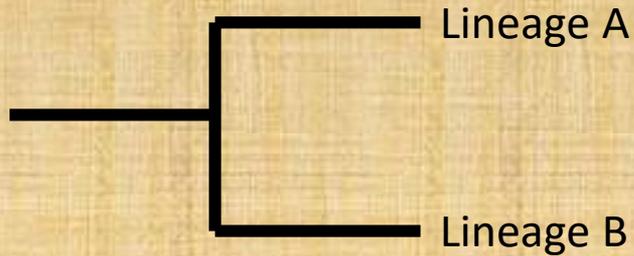
Fst scans

Adaptomics methods involving additional aspects than just allele frequencies

Genetic vs phenotypic association: QTLs and GWAS

Environmental associations

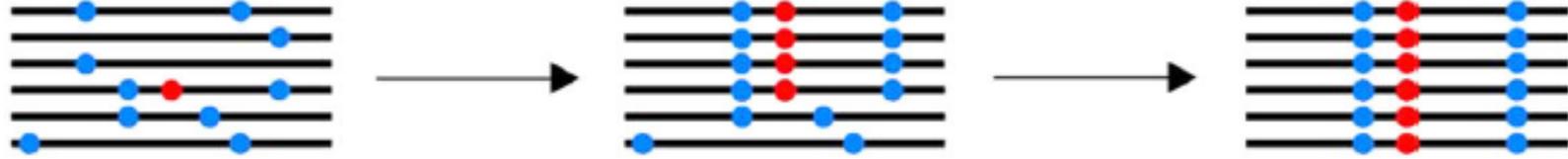
Types of selection



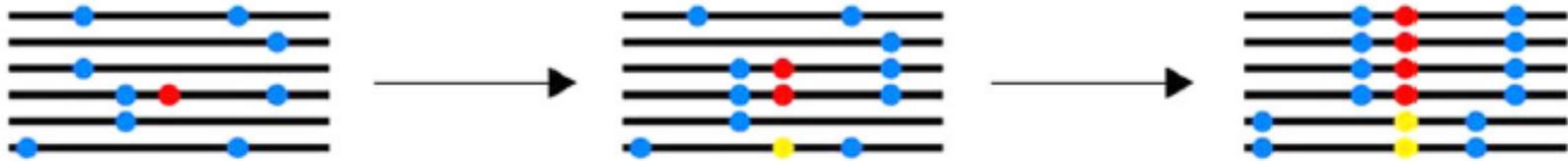
Type of selection	Consequence on intralinear variability	On interlineage divergence	On site frequency spectrum
Negative/Purifying	Reduced	Reduced	more low frequency derived alleles
Positive selection	Reduced	Increased	more high frequency derived alleles
Balancing selection	Increased	none	more intermediate frequency derived/ancestral alleles

A selective sweep

a Incomplete,
then complete
hard sweep

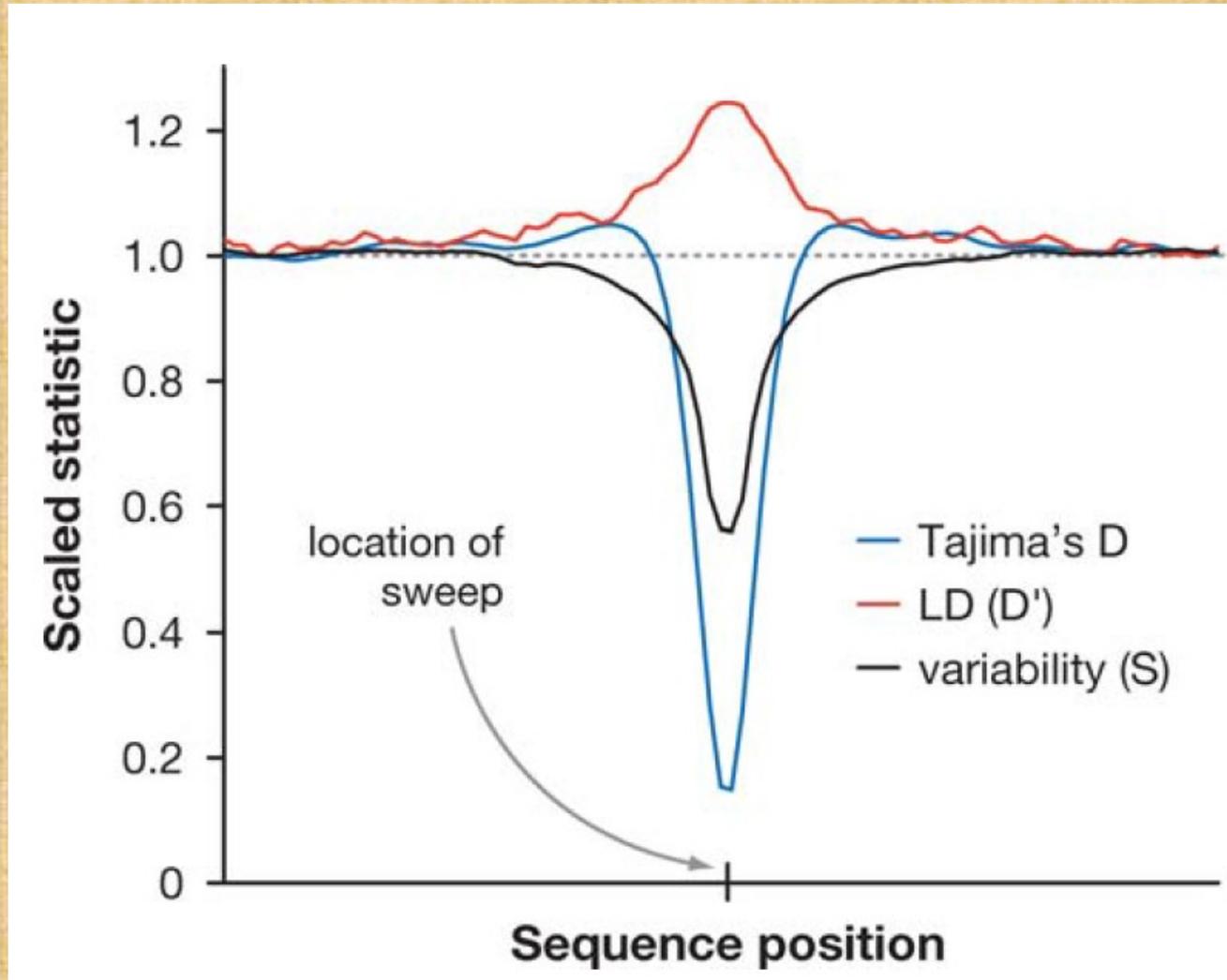


c Soft sweep,
recurrent
mutation



Because of genetic linkage, reduction of local
genetic diversity around the selected mutation

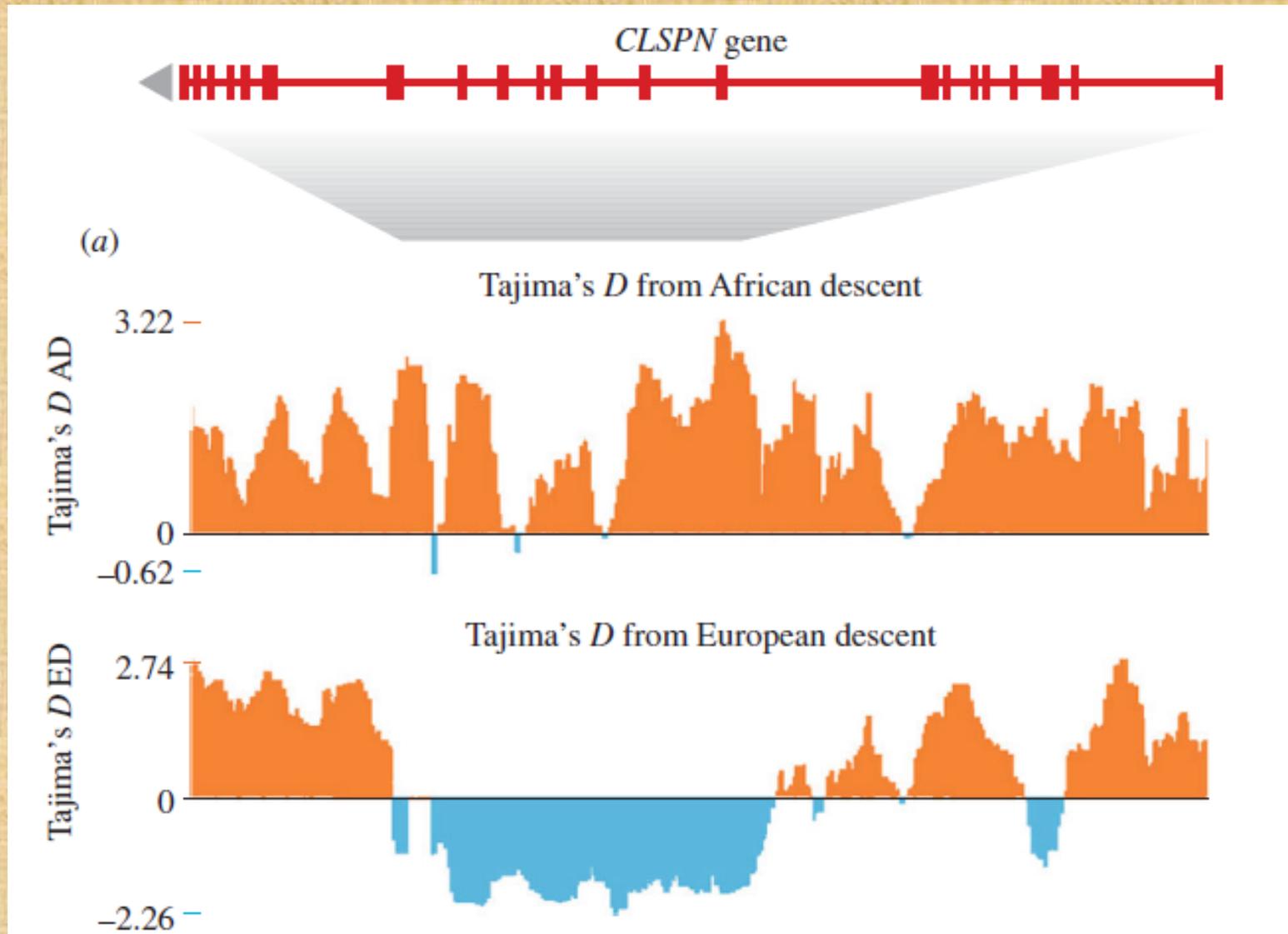
A typical selective sweep



As a general rule, loci under potential selection need to be compared to the overall genome:

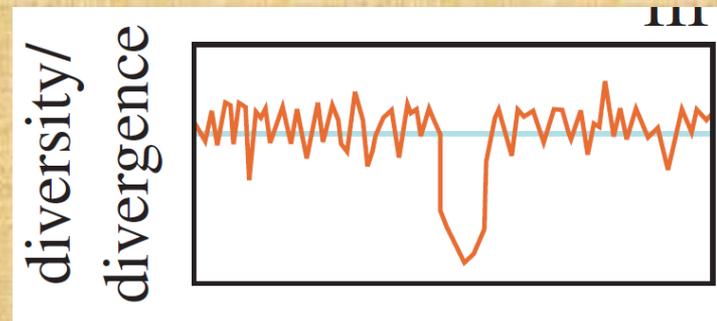
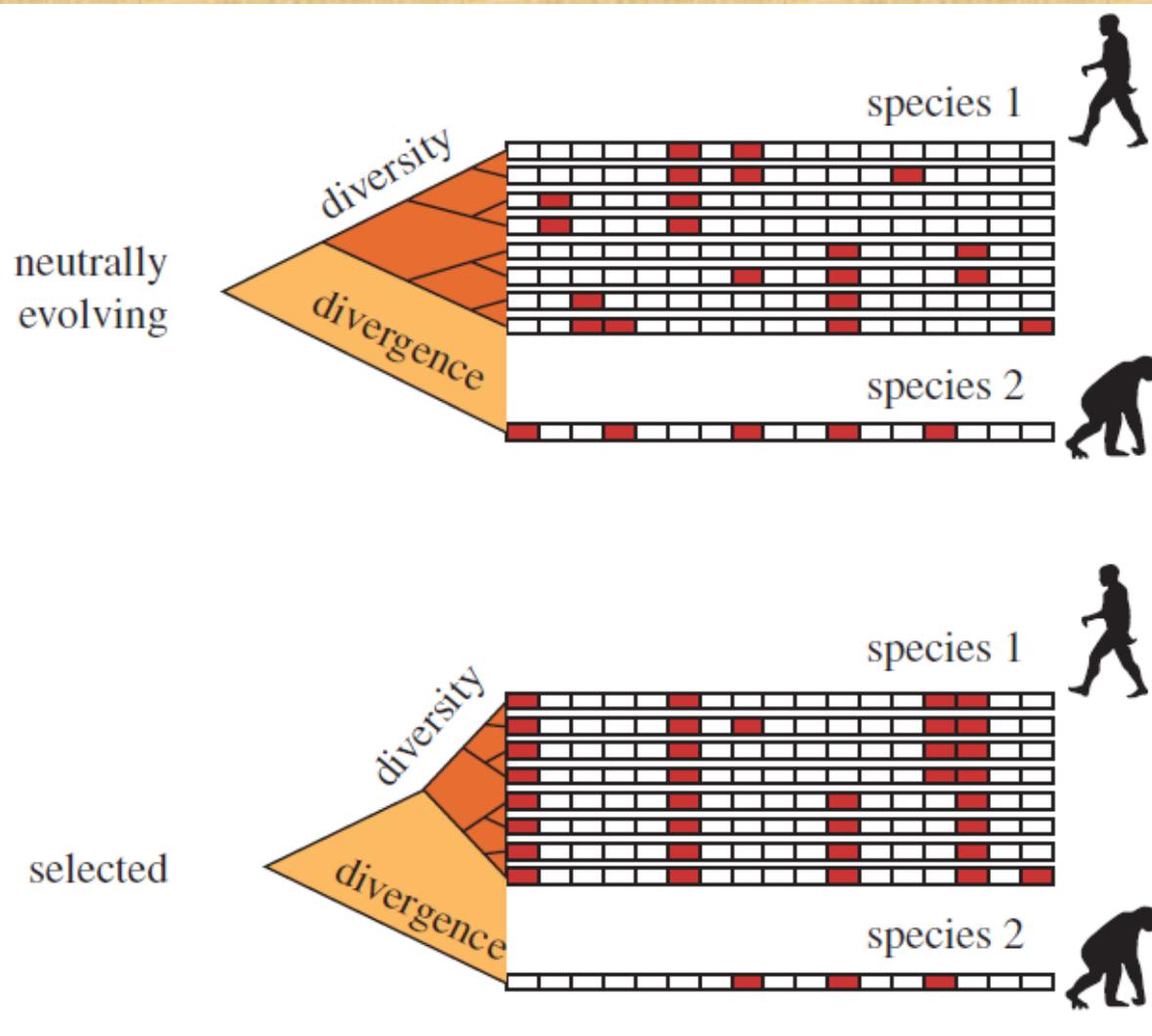
Demography can influence LD, Tajima's D...

Tajima's D as a test for local selection



Locally, less genetic variation than expected given the number of individuals

Testing for the ratio between interlineage divergence vs intralineage variation



For coding genes, the rate of syn vs non syn mutations gives information on the selection process

If neutral evolution,
interlineage divergence in non synonymous site = interlineage divergence in
synonymous site

$$D_n / D_s = 1$$

Above 1: positive selection on new non syn mutations

The **McDonald-Kreitman test** adds a comparison between interlineage divergence and intralinear polymorphism

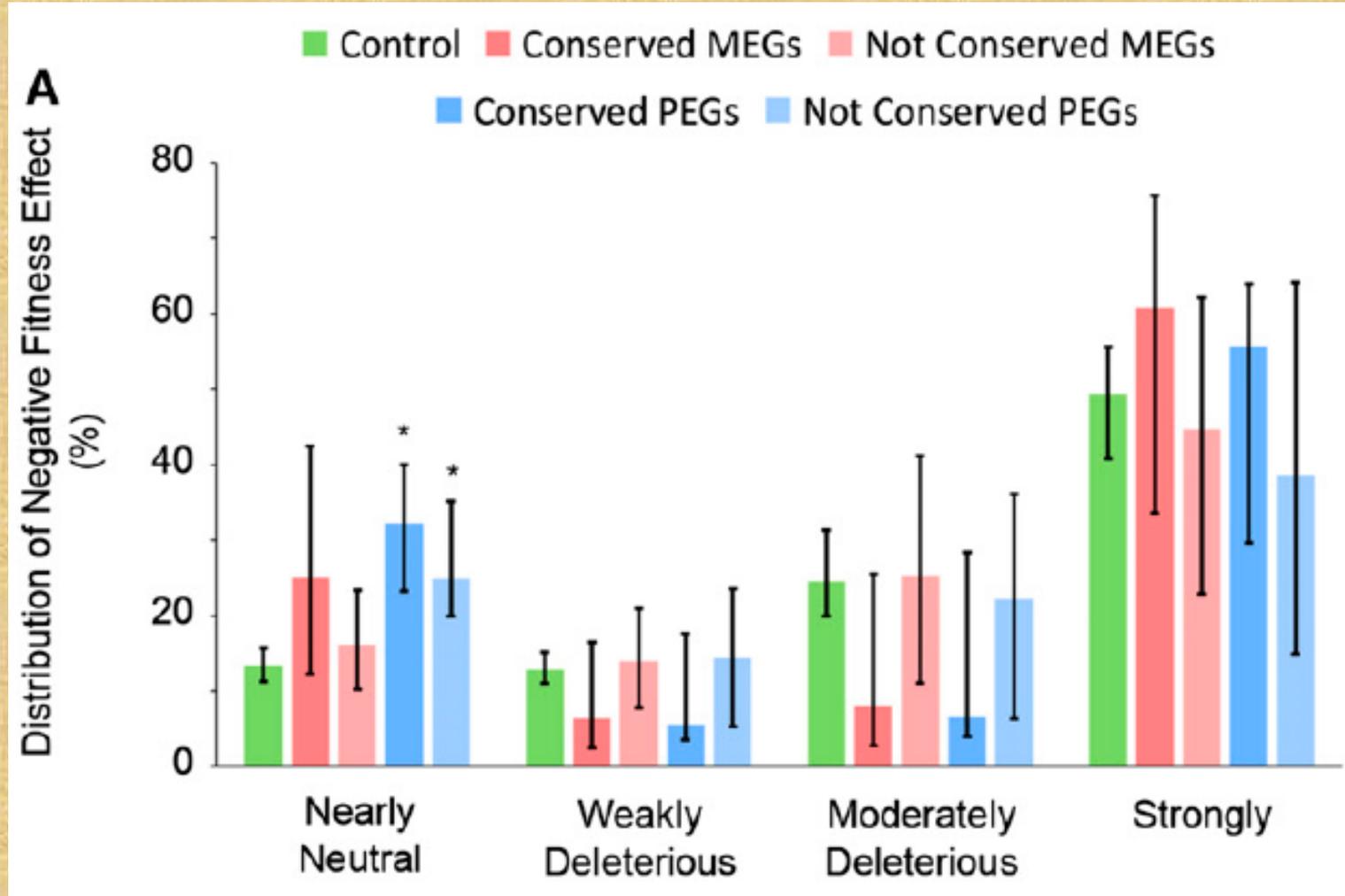
Makes a ratio of P_n / P_s by D_n / D_s

$$\alpha = 1 - \frac{D_s P_n}{D_n P_s}$$

If $\alpha = 0$, divergence = polymorphism = neutral

If $\alpha = 1$, divergence \gg polymorphism = positive selection

The distribution of fitness effect for candidate genes



The strength of **purifying selection**: $N_e S$

F_{st} scans

Intralineage variation



Interlineage divergence

Intralineage variation

> or =

Interlineage divergence

→ LOW F_{st}

Intralineage variation

< or =

Interlineage divergence

→ HIGH F_{st}

LOW F_{st}

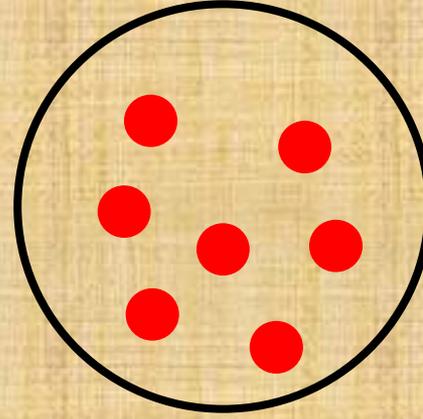
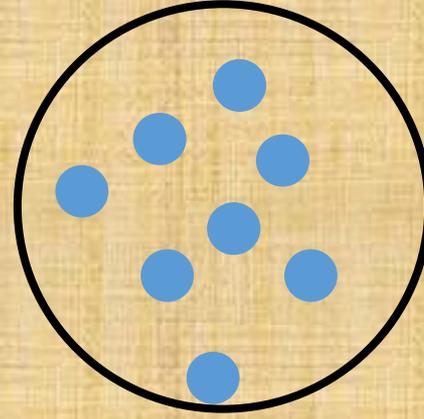
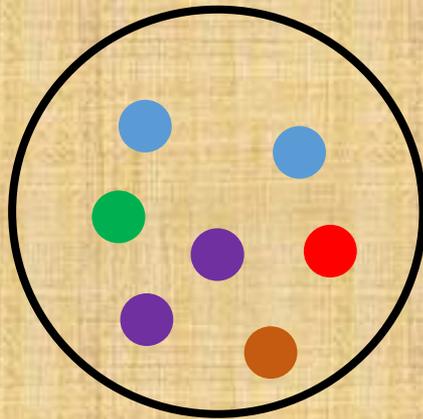
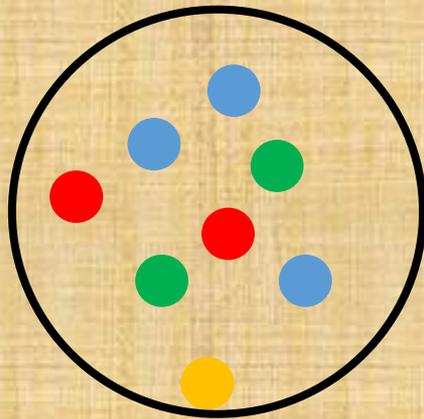
HIGH F_{st}

Lineage A

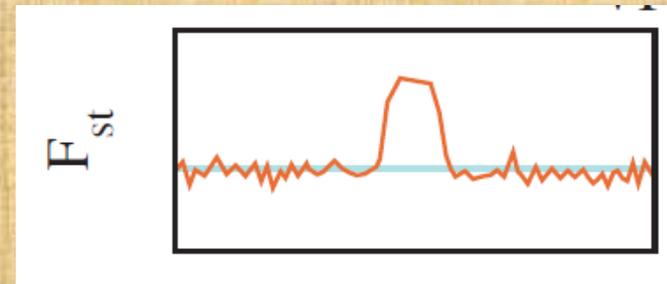
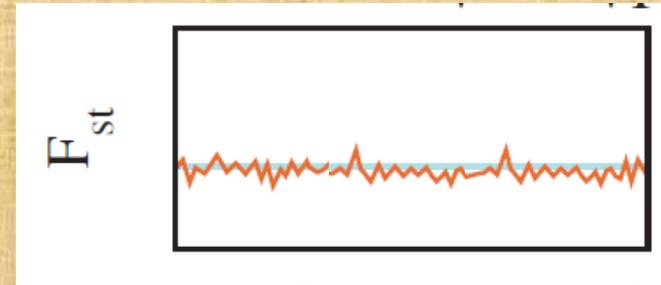
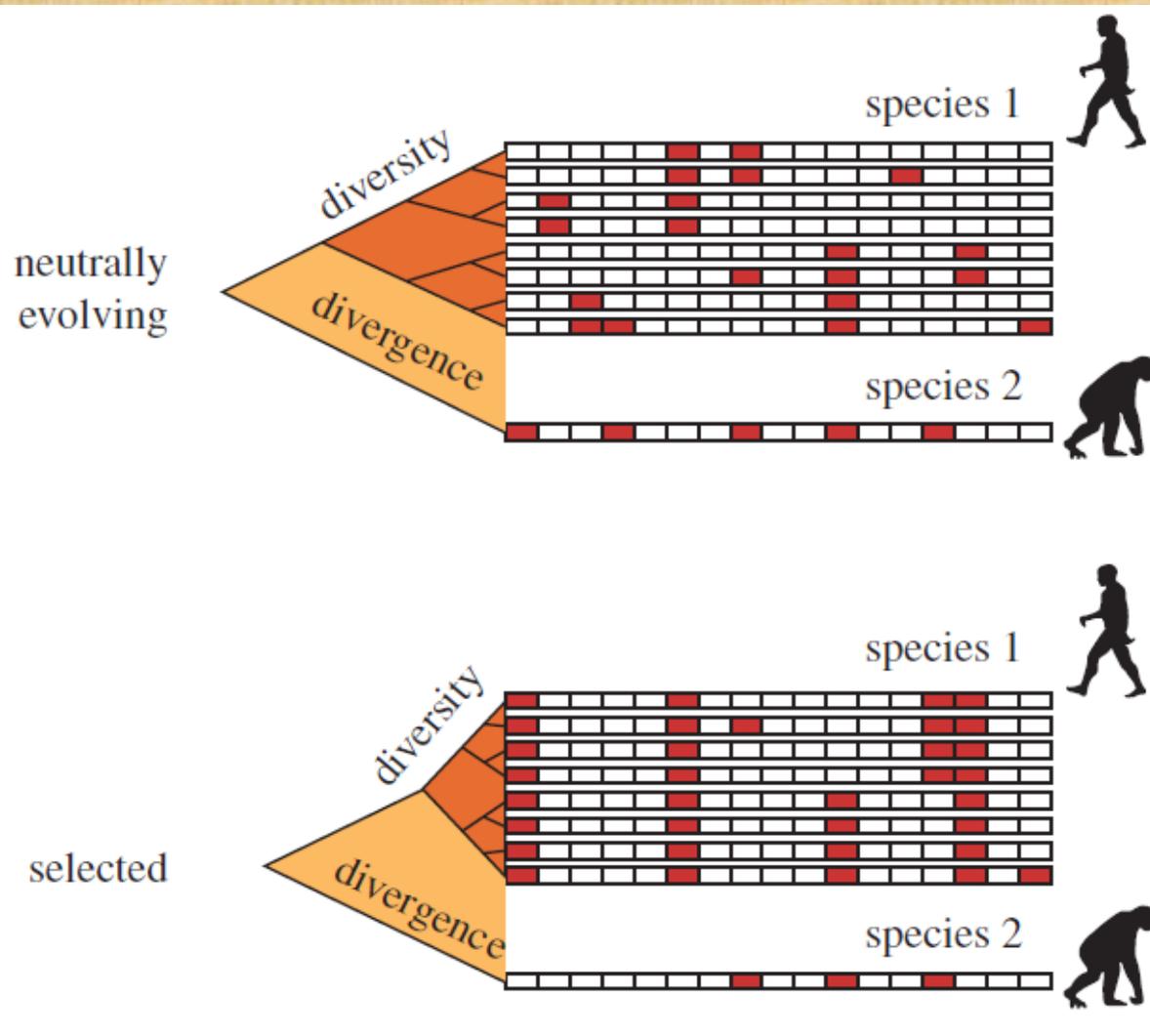
Lineage B

Lineage A

Lineage B



F_{st} scans



Plan

Demography and other factors influencing genetic variation

Bottlenecks, migration, mating system etc

N_e and efficacy of selection

LD/Recombination rate

Testing for demographic confounding factors

Nucleotide diversity, heterozygosity

Testing for demographic events: Tajima's D

Estimating the "genetic load" and efficacy of selection: SFS and DFE

Detecting natural selection vs neutrality (and demography) in the genome

types of selection

selective sweeps

dN/dS and PN/PS

DFE and SFS

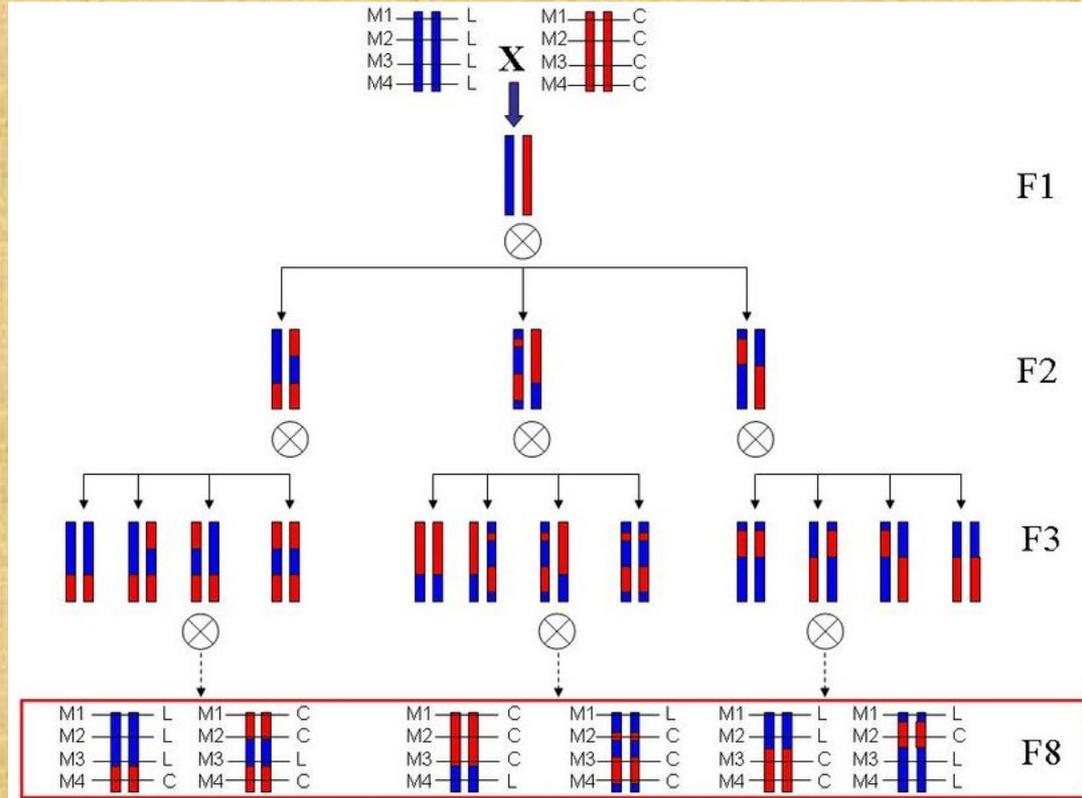
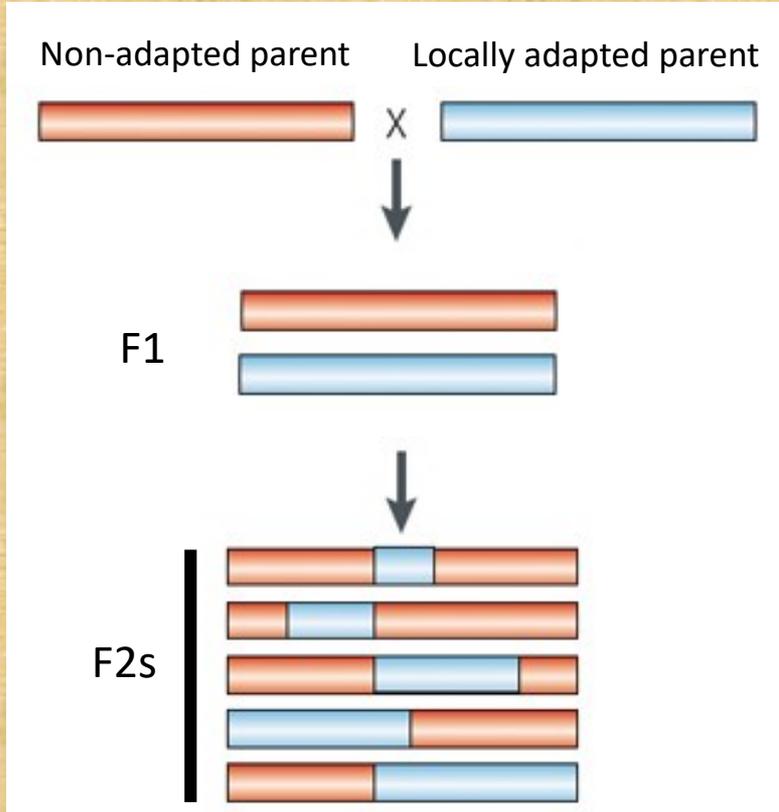
Fst scans

Adaptomics methods involving additional aspects than just allele frequencies

Genetic vs phenotypic association: QTLs and GWAS

Environmental associations

Genetic vs phenotypic association: QTLs

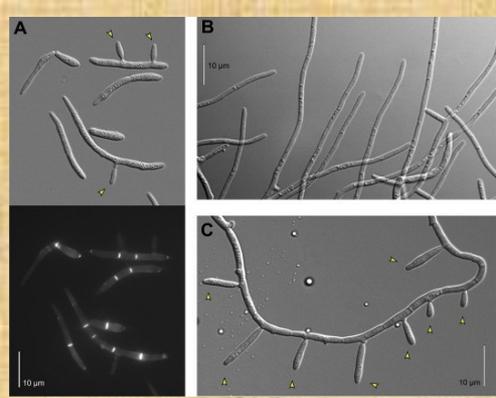


- + : fast to obtain
- : heterozygosity
- : few recombination events (one generation) → large QTLs

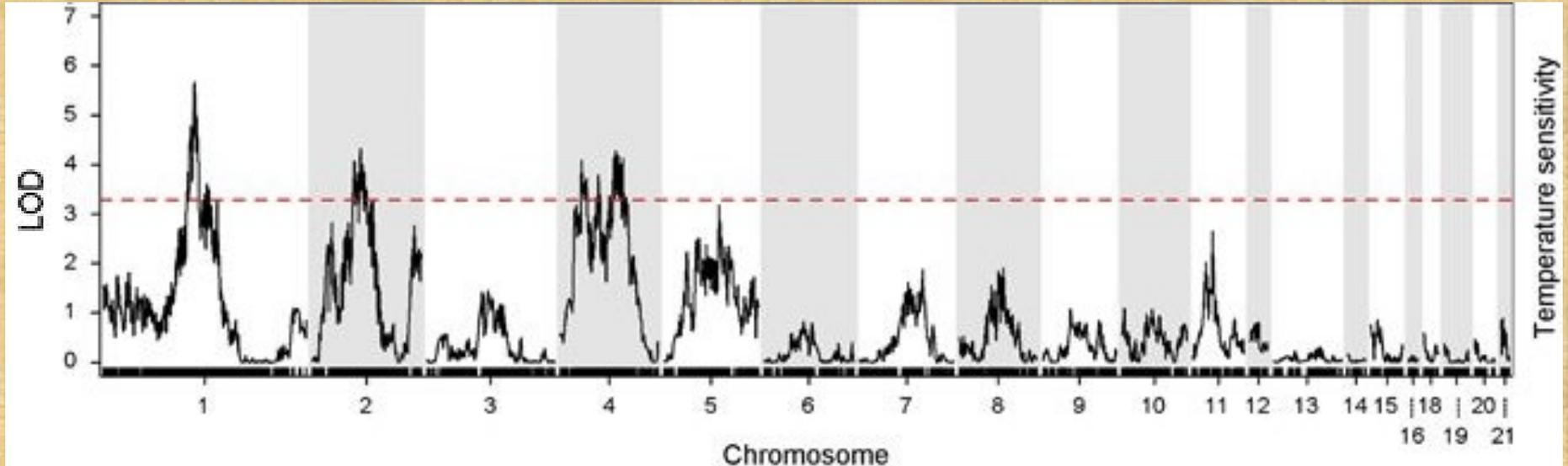
- + : more recombination events (>5 generations) = more combinations possible = narrow QTLs
- + : more homozygous = clearer phenotypes (dominance)
- : more time and effort consuming

Phenotyping of the F2s/RILs + Sequencing

Genetic vs phenotypic association: QTLs

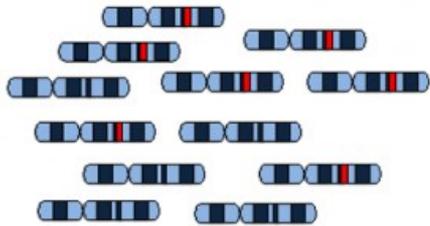


Zymoseptoria tritici

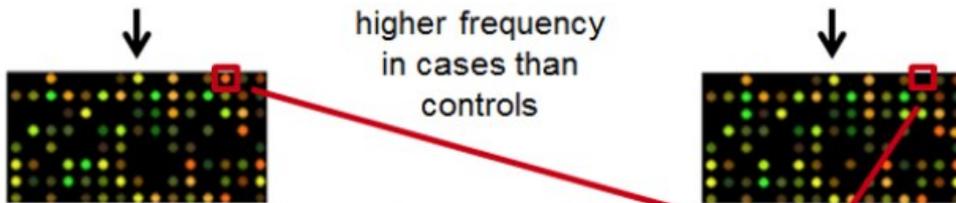
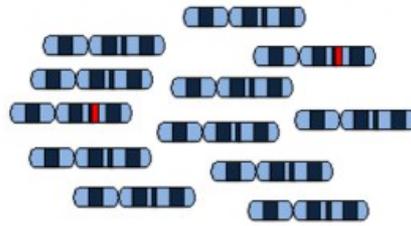


LOD score: the probably of association between phenotype and genotype

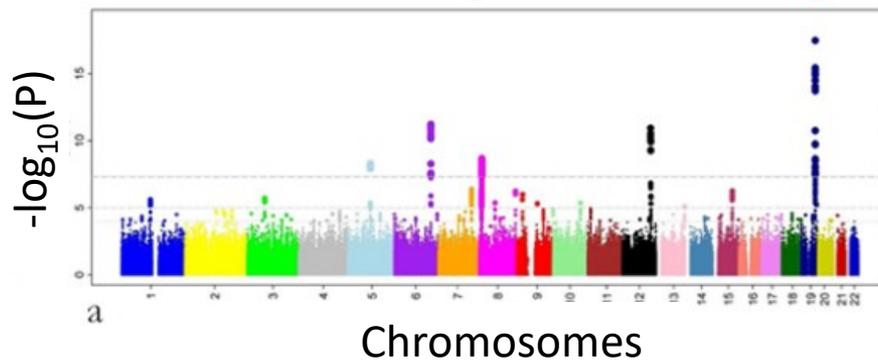
Non-adapted pop



Locally adapted pop



higher frequency
in cases than
controls

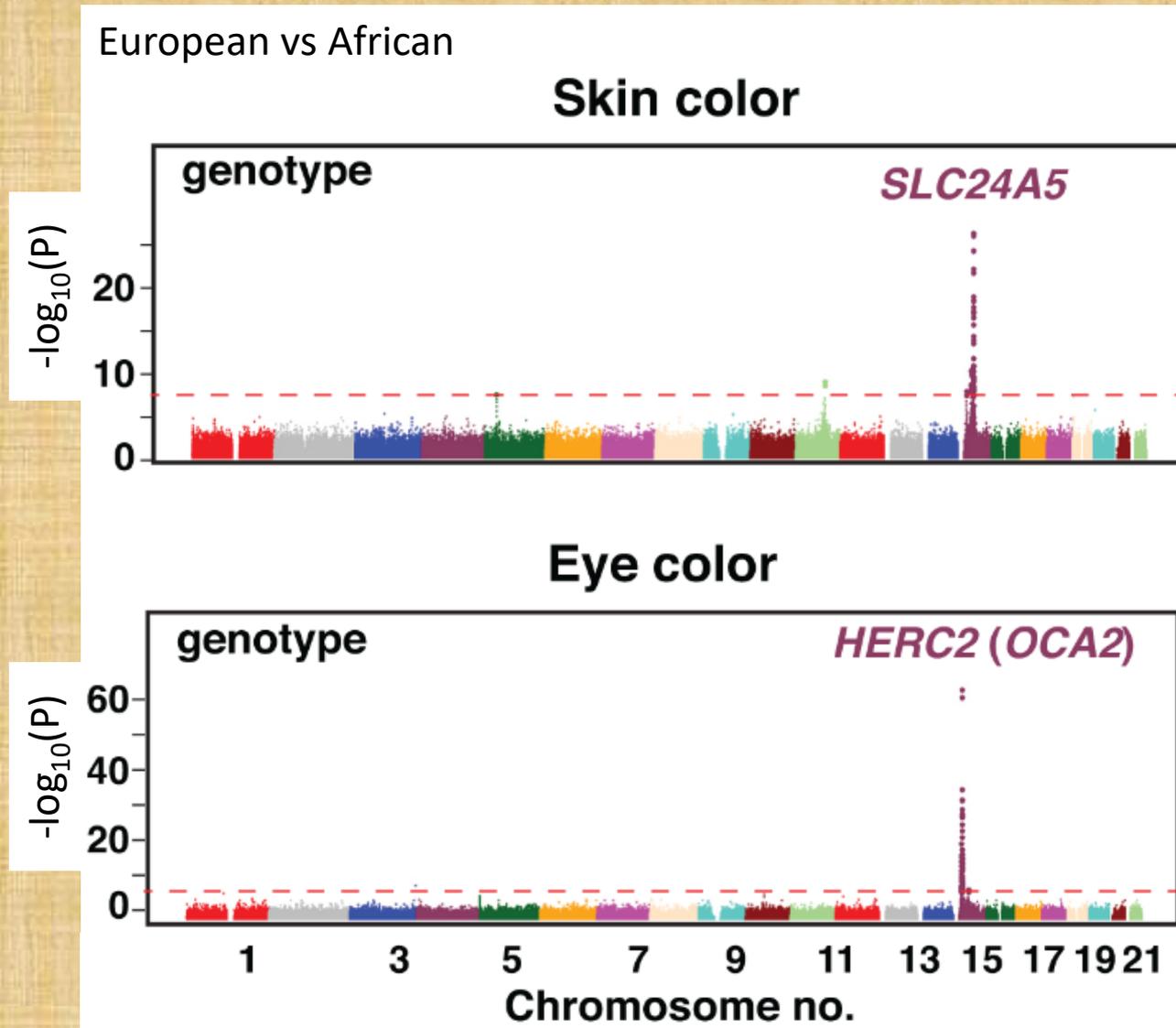


Genotype-Phenotype association

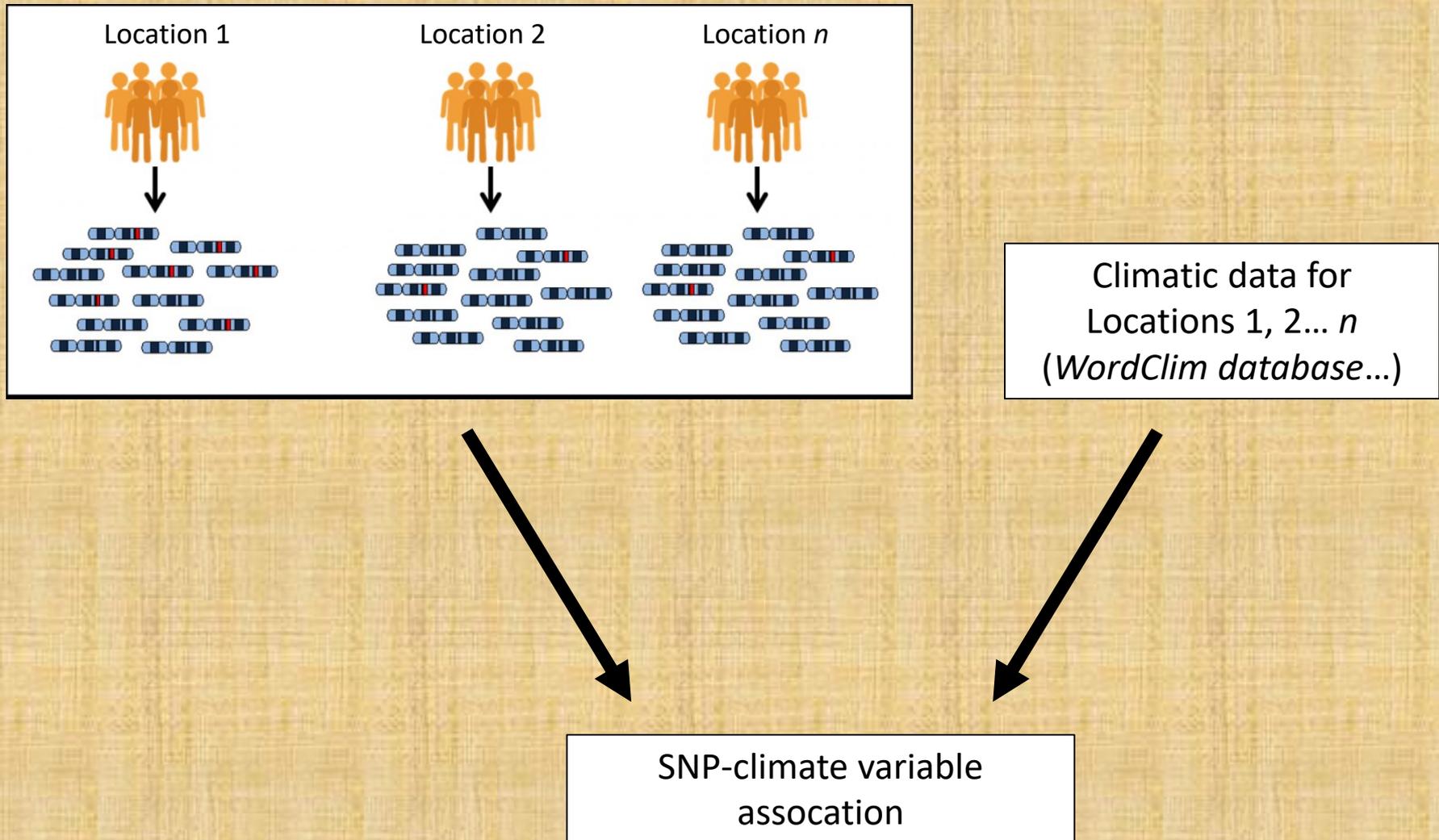
Better resolution than QTLs: one can infer the effect of a single SNP on the phenotype

Important factors to take into account: demography, especially population structure

Genetic vs phenotypic association: GWAS



Environmental association: principle

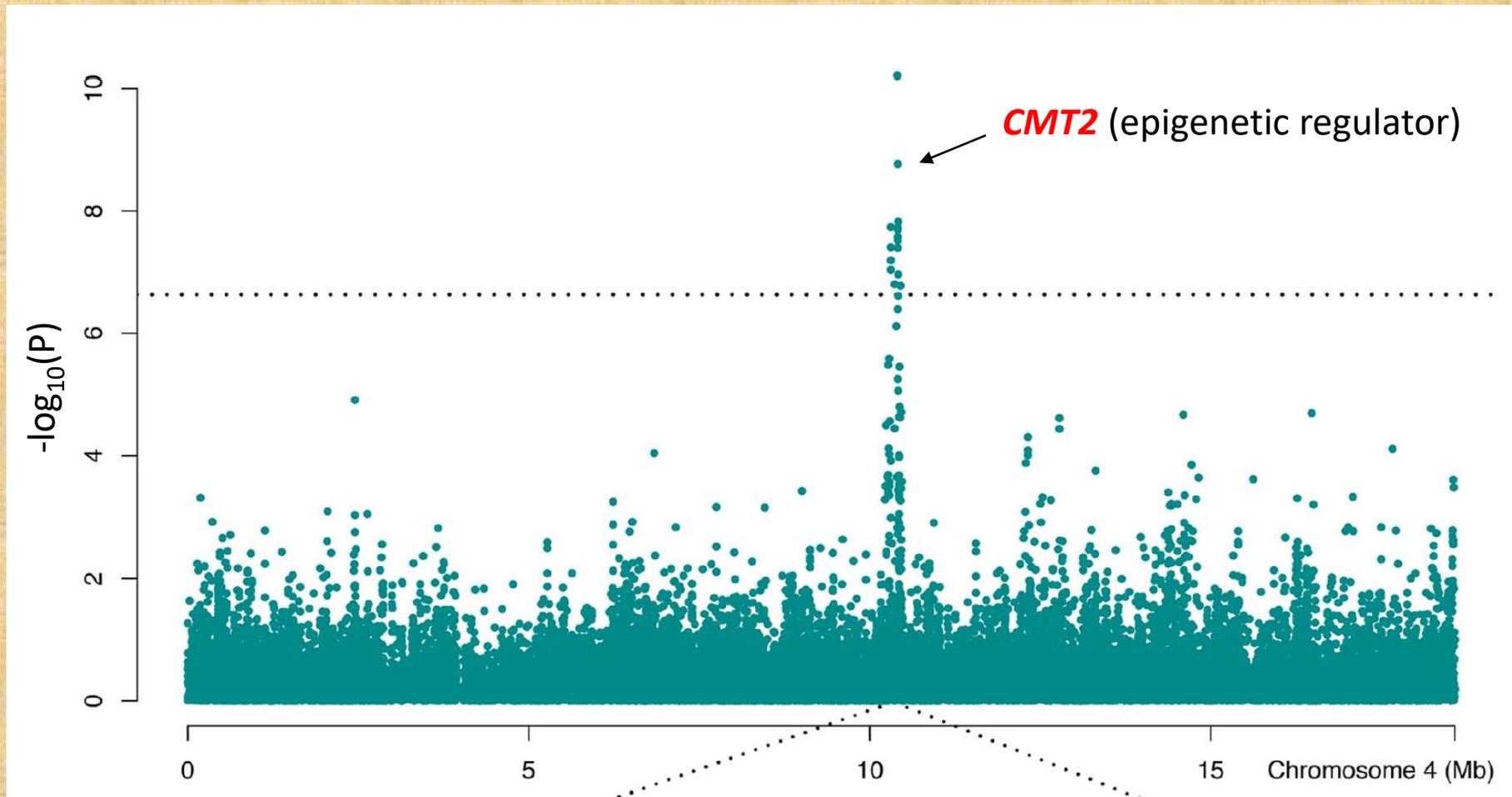


Environmental association: example



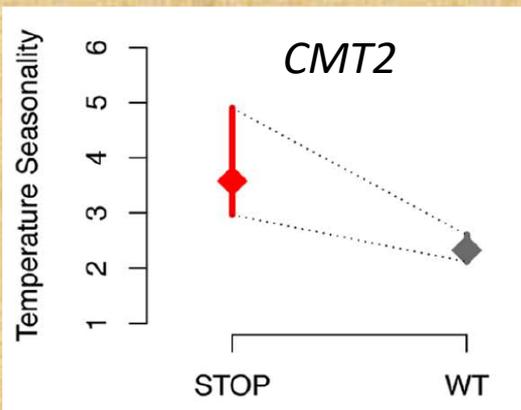
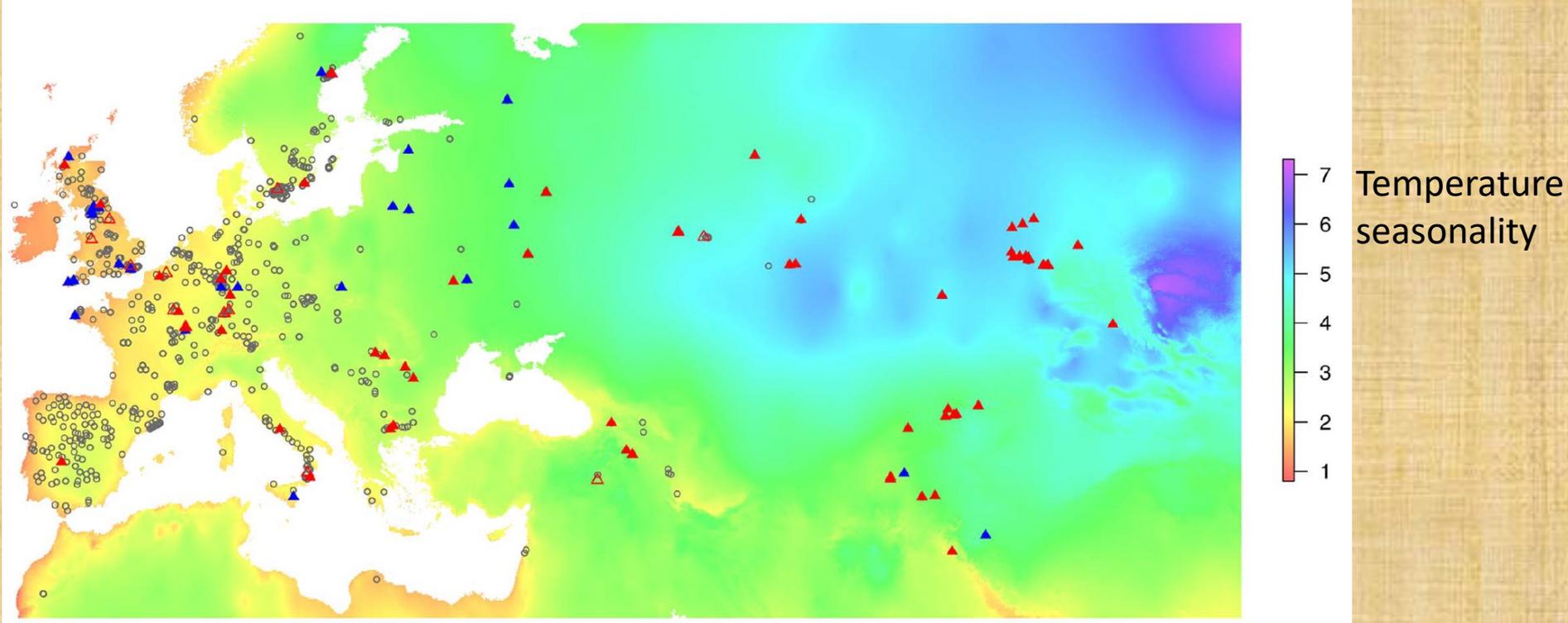
A. thaliana

Association with **temperature seasonality** (difference between summer and winter)



Environmental association: example

○ Functional allele of *CMT2* ▲/▲ Stop-codon allele



Conclusions

Neutral molecular evolution = null hypothesis

Demography can have large effects on molecular evolution

It is necessary to test demographic effects before searching for selection:
when searching for loci under selection, always determine the genomic pattern

Different methods to search for selection, need to be adapted to the model study