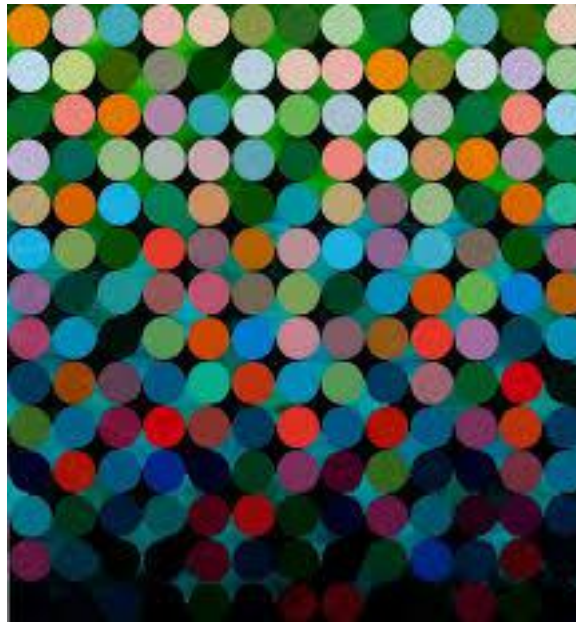


## How to detect selection in DNA sequences?



C. Lemaire  
Norwich, April 5-6th 2017

# MUTATION, RECOMBINATION AND SELECTION : POPULATION GENOMICS INFERENCES

 Polymorphic sites identification

 Polymorphism rate estimation

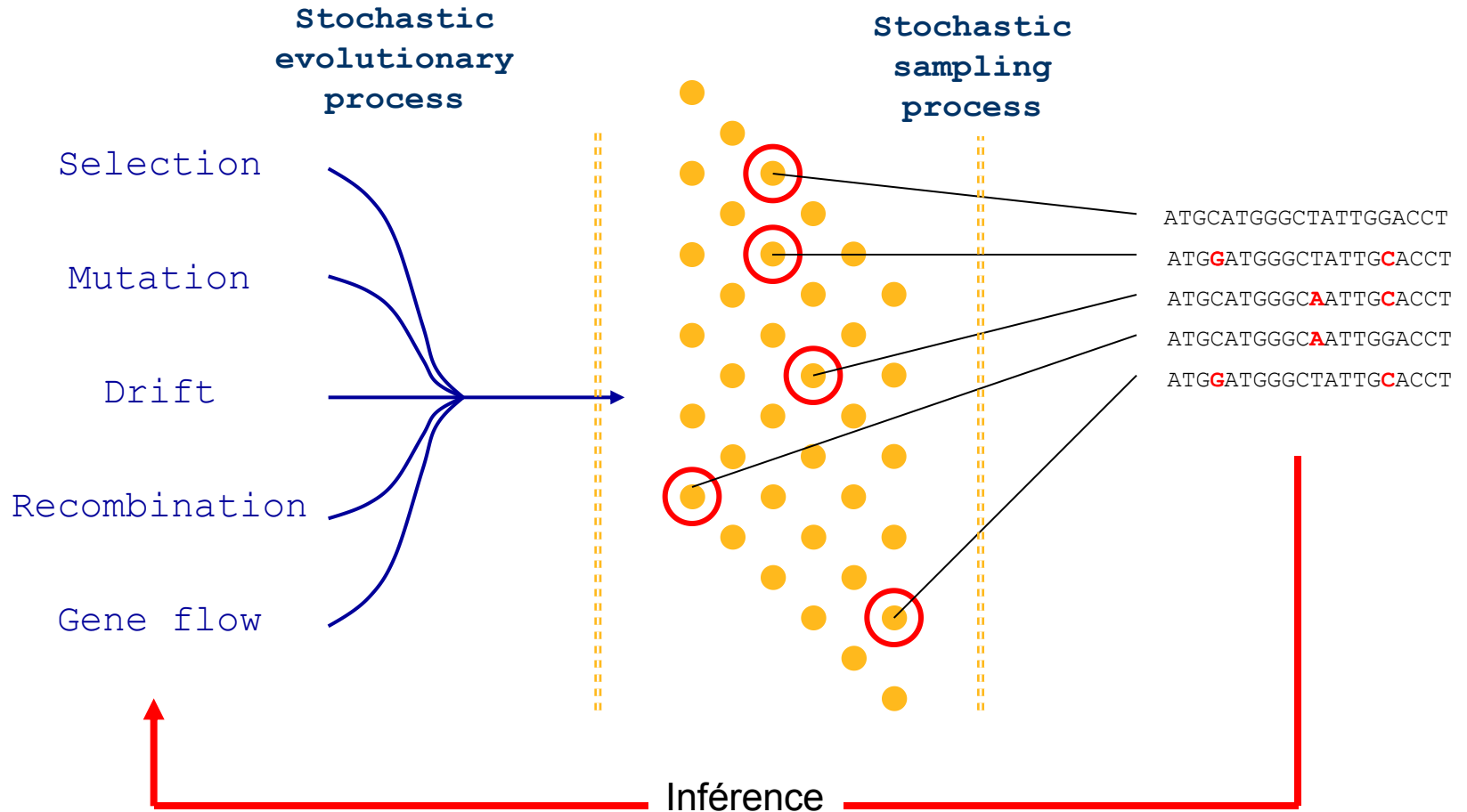
 What forces shape polymorphism?

# INFERENCES IN POPULATION GENOMICS

Evolutionary parameters

Whole population

Sample



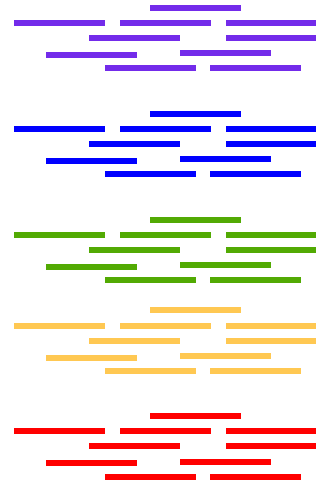
# Polymorphism identification



**DNA or RNA  
Extraction**



**Sequencing**  
454  
Illumina, ...



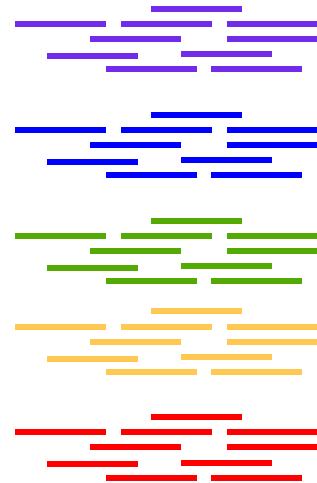
# Polymorphism identification



**DNA or RNA  
Extraction**



**Sequencing**  
454  
Illumina, ...



**SNP calling  
Indel calling**



*de novo* Assembling  
ou  
Alignment on reference



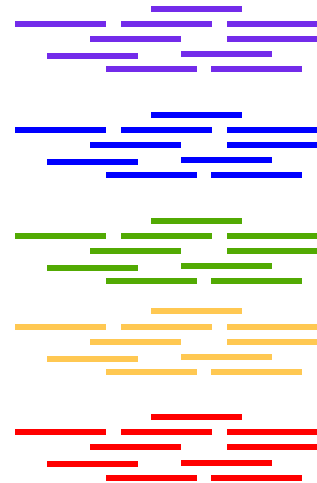
# Polymorphism identification



**DNA or RNA  
Extraction**



**Sequencing**  
454  
Illumina, ...



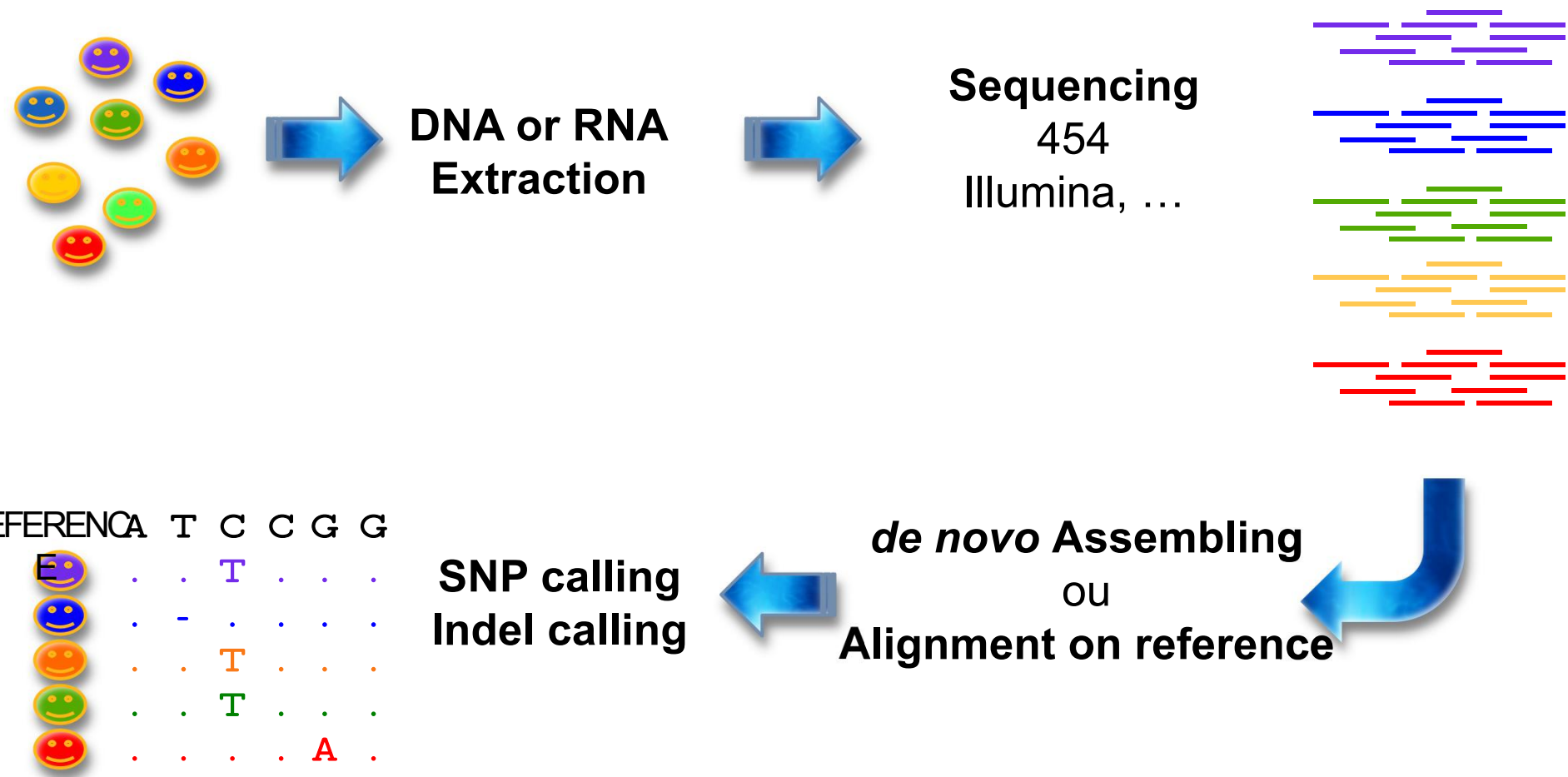
*de novo* Assembling  
ou  
Alignment on reference



**SNP calling  
Indel calling**



# Polymorphism identification



# Polymorphism identification

*de novo* Assembling  
or  
Alignment on reference




SNP calling  
Indel calling

 Read quality

 Mate pair alignment

 Coverage thresholds min & max

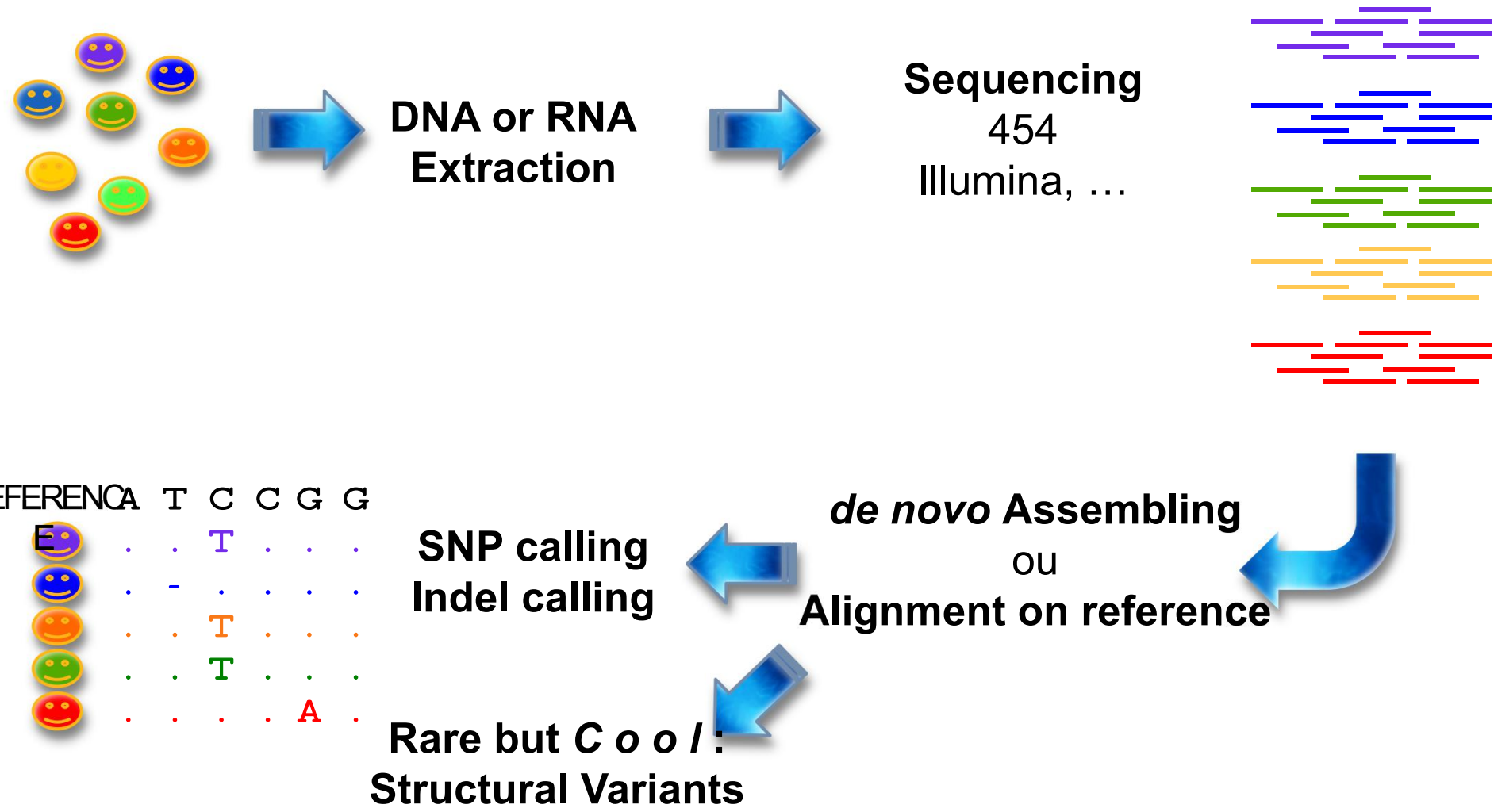
 Max divergence rate (e.g. BWA 8% + 2 indels)



Keep information on overall number of sites (i.e. **non filtered**) for polymorphism rates estimates



# Polymorphism identification






# Polymorphism identification

## Divergence islands detection

 *de novo* assembling > alignment on référence (e.g. limit <8% divergence for BWA)









## CONCLUSIONS

-  Polymorphic sites are mainly « Single Nucleotide Polymorphism » or SNPs .
-  WARNING: non covered regions and filters for SNP calling for polymorphism rates estimates.
-  Coverage analysis: « Structural Variants » : duplications, deletions, insertion, divergence islands, that could have strong phenotypic effects.

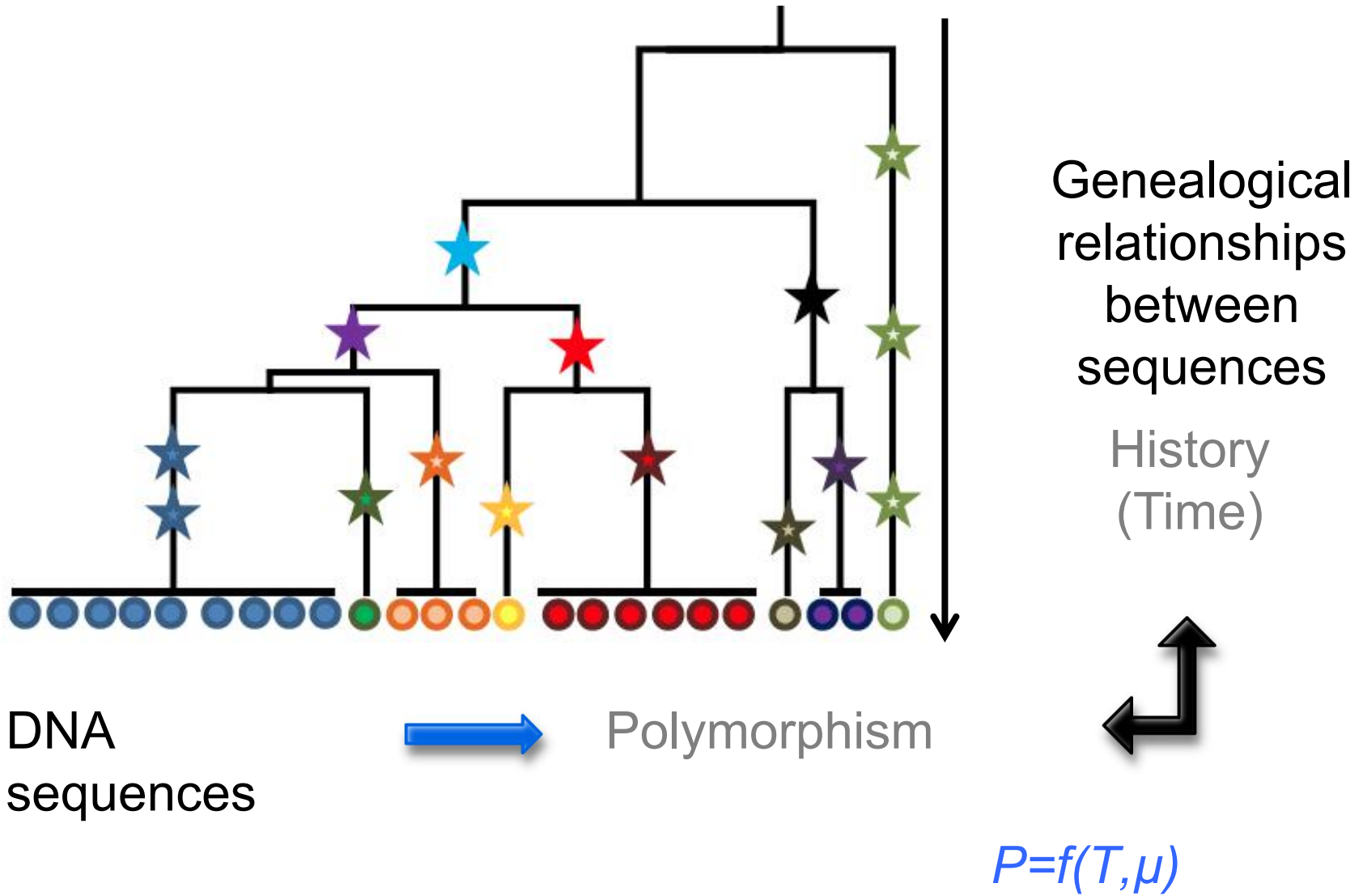
# Polymorphism rate estimation

## POLYMORPHISM RATE

	CC GCA GAG TTA CTA ATC GA	$N = 6$	individuals
	C <b>G</b> GCA GAG TTA CTA ATC GA	$L = 21$	sites on alignment
	CC GCA <b>A</b> AG TTA C <b>C</b> A AT <b>T</b> GA	$S = 5$	polymorphic sites
	CC GCA GAG TTA C <b>C</b> A ATC GA		
	CC GCA <b>A</b> AG TTA CTA ATC GAG <b>G</b>		
	CC GCA <b>A</b> AG TTA CTA ATC GA		

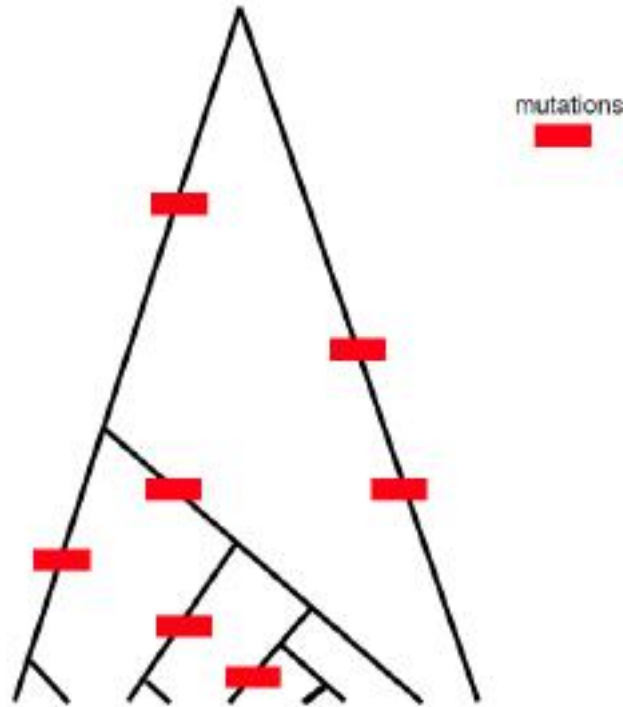
How to quantify polymorphism?

# Polymorphism is a combination of mutation and time



# Genealogical relationships between sequences

MRCA: Most Recent Common Ancestor



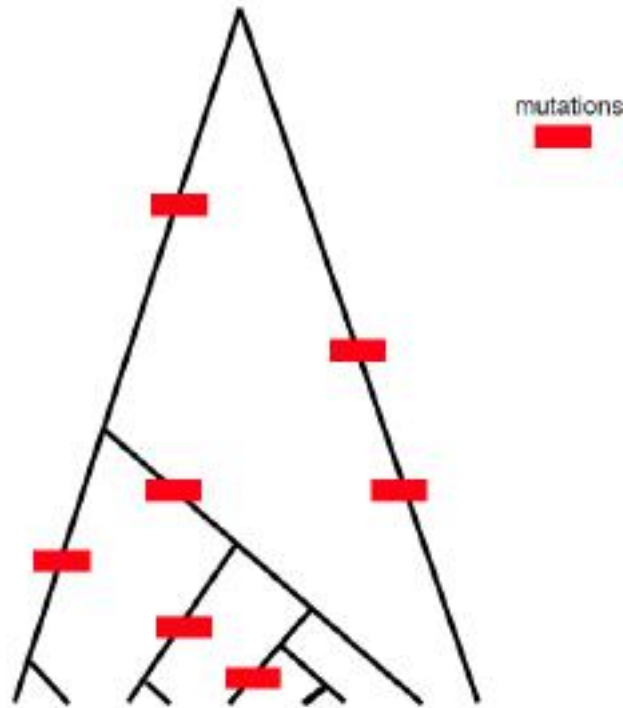
The coalescent, *i.e.* the genealogy of alleles at a locus, depends on evolutionary forces that shape polymorphisms.  
Null hypothesis in evolution.

(Achaz, Introduction à la coalescence, 2005)

More details on the coalescence theory on:

[http://www.sfu.ca/biology/courses/bisc869/869\\_lectures/MHP\\_Coalescent.pdf](http://www.sfu.ca/biology/courses/bisc869/869_lectures/MHP_Coalescent.pdf)

## MRCA: Most Recent Common Ancestor



(Achaz, Introduction à la coalescence, 2005)

**Number of mutations**  
Poisson approximation of  
Binomial

$$P(k/t) = e^{-\mu t} \frac{(\mu t)^k}{k!}$$

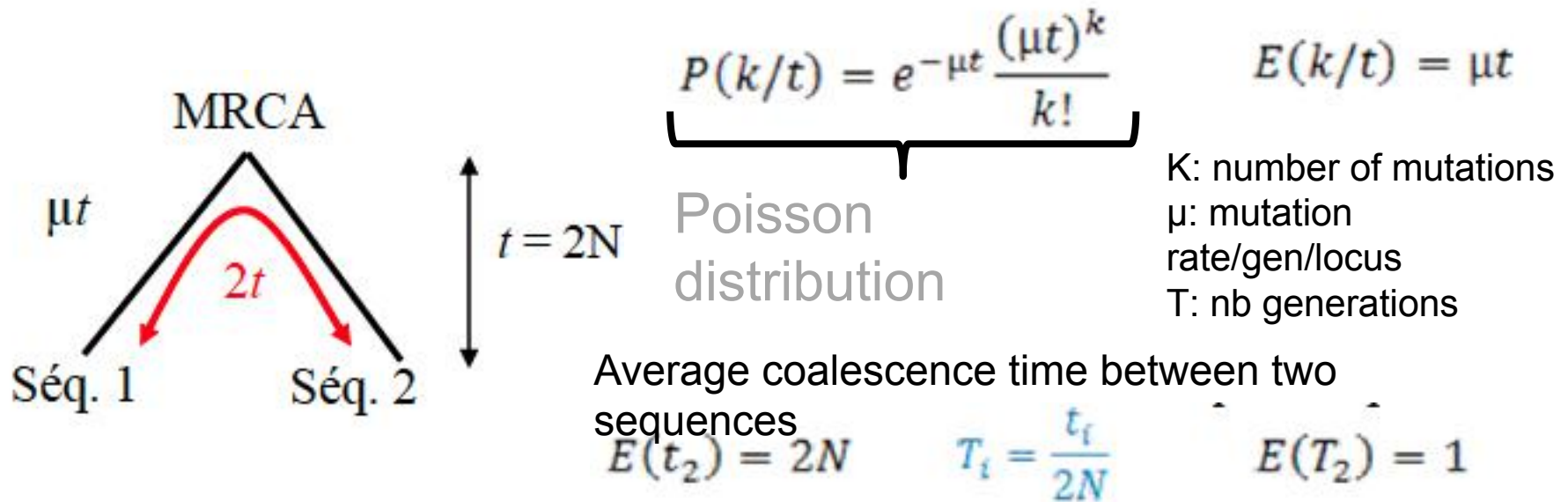
Poisson  
distribution

$$E(k/t) = \mu t$$

K: number of mutations  
 $\mu$ : mutation  
rate/gen/locus

# Genealogical relationships between sequences

What is the expected number of differences between two sequences?



$$E(k_2) = 2 \times \mu \times E(t_2)$$

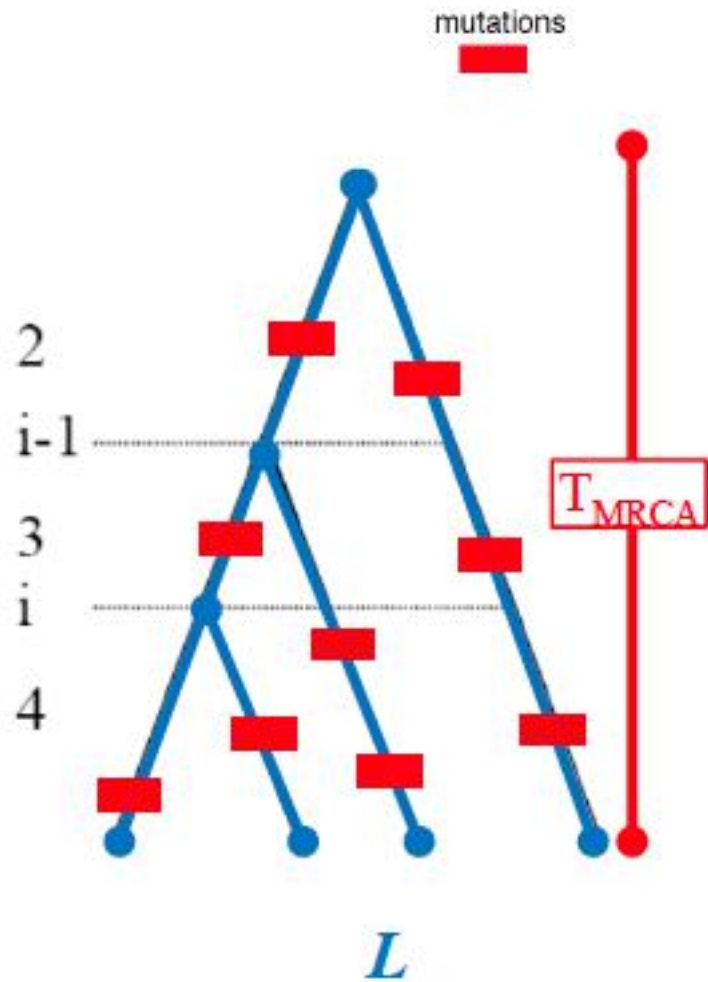
$$E(k_2) = 2 \times \mu \times 2N = 4N\mu = \theta$$

$q$  Is the populational mutation rate.

Easier to estimate than  $N$  and  $\mu$



# The Watterson's $q_s$



Length of a genealogy

$$L = a_n \times 4N \quad \text{with} \quad a_n = \sum_{i=1}^{n-1} \frac{1}{i}$$

Expected number of polymorphic sites

$$E(S) = \mu \times L$$

$$E(S) = \mu \times a_n \times 4N$$

$$E(S) = a_n \theta$$

$$\theta_s = \frac{S}{a_n}$$

## The Tajima's $P$

$P$  is the expected average number of differences between any pairs of sequences in the population

$$\Pi = \frac{1}{C_n^2} \sum_{i=1}^{n-1} \sum_{j=i+1}^n d_{ij}$$

$$E(d_{ij}) = E(k_2) = 2 \times \mu \times 2N = 4N\mu = \theta$$

$$\theta_{\Pi} = E(\Pi)$$

# Polymorphism rate estimation

## POLYMORPHISM RATE







😊	CC	GCA	GAG	TTA	CTA	ATC	GAA	$N = 6$	individuals		
😐	CC	G	GCA	GAG	TTA	CTA	ATC	GAA	$L = 21$	sites on alignment	
😊	CC	GCA	A	AG	TTA	C	A	ATT	GAA	$S = 5$	polymorphic sites
😊	CC	GCA	GAG	TTA	C	A	ATC	GAA			
😐	CC	GCA	A	AG	TTA	CTA	ATC	GAG			
😊	CC	GCA	A	AG	TTA	CTA	ATC	GAA			

$$P = \frac{\sum_{i=1}^N \sum_{j=1}^N d_{ij}}{L \cdot N \cdot (N-1) / 2}$$

$$Q = \frac{S}{L \cdot \sum_{i=1}^N 1/i}$$

# Polymorphism rate estimation

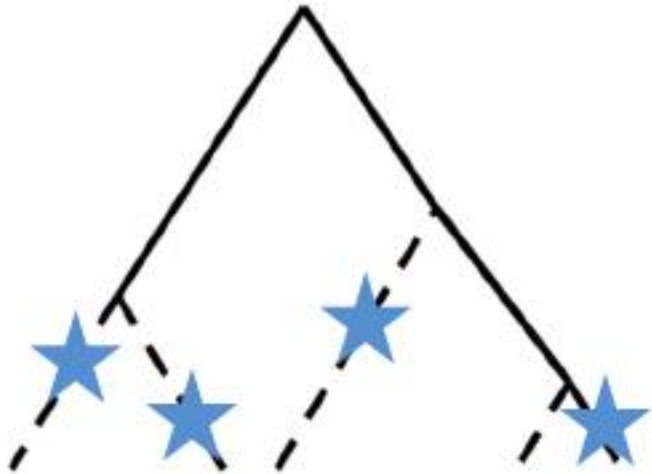
## POLYMORPHISM RATE

	CC GCA GAG TTA CTA ATC GAA	$N = 6$	individuals
	C <b>G</b> GCA GAG TTA CTA ATC GAA	$L = 21$	sites on alignment
	CC GCA <b>A</b> AG TTA C <b>C</b> A ATT GAA	$S = 5$	polymorphic sites
	CC GCA GAG TTA C <b>C</b> A ATC GAA		
	CC GCA <b>A</b> AG TTA CTA ATC GAG <b>G</b>		
	CC GCA <b>A</b> AG TTA CTA ATC GAA		

$$P = \frac{\sum_{i=1}^N \sum_{j=1}^N \frac{p_{ij}}{L \cdot N \cdot (N-1)/2}}{P = 0.102}$$

$$Q = \frac{S}{\sum_{i=1}^L \frac{1}{i}} \quad Q = 0.104$$

## The Fu & Li's $q_{he}$



Internal (solid) and external (dashed) in a genealogy.

Estimator based on the mutations found on external branches.

Present in one copy: **singleton**

$$\eta_e = L_n \mu = 4N\mu$$

$$E(\eta_e) = \theta$$

# Polymorphism estimation

## Theoretical expectations of the polymorphism rate

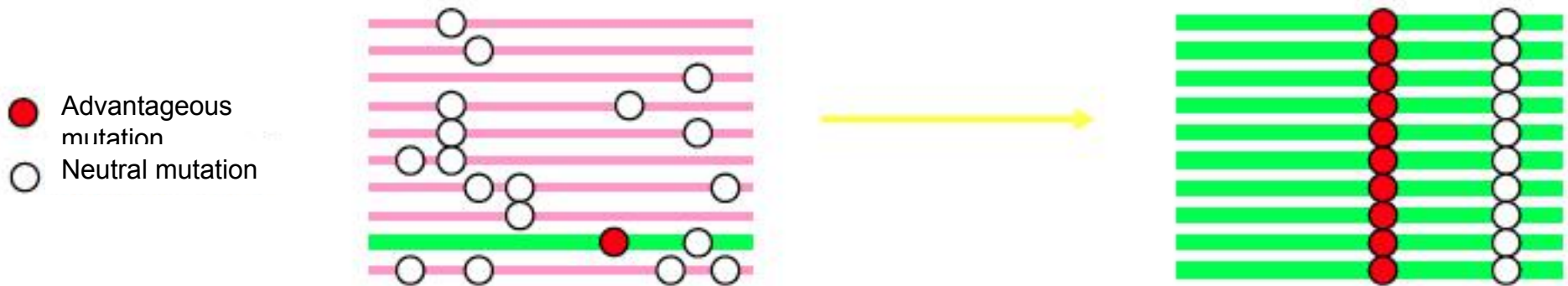
 Neutral mutations

 Constant size population

$$P = Q = 4 N_e \mu \text{ diploid population}$$

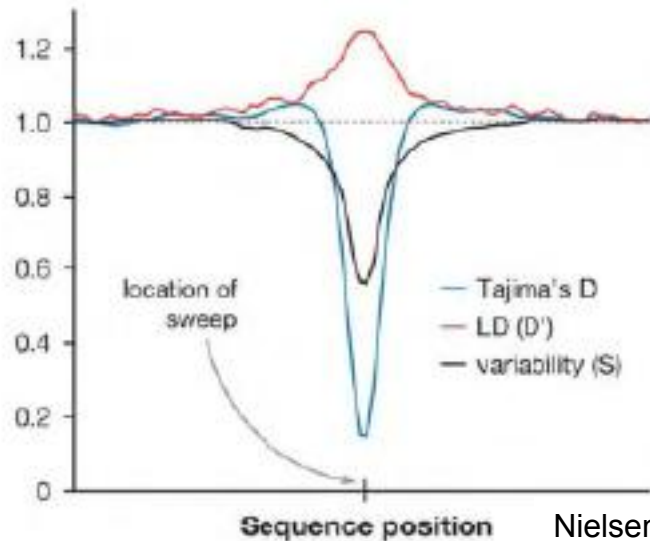
$$P = Q = 2 N_e \mu \text{ haploid population}$$

# Positive selection



An advantageous mutation will increase in frequency and eliminate all polymorphism around

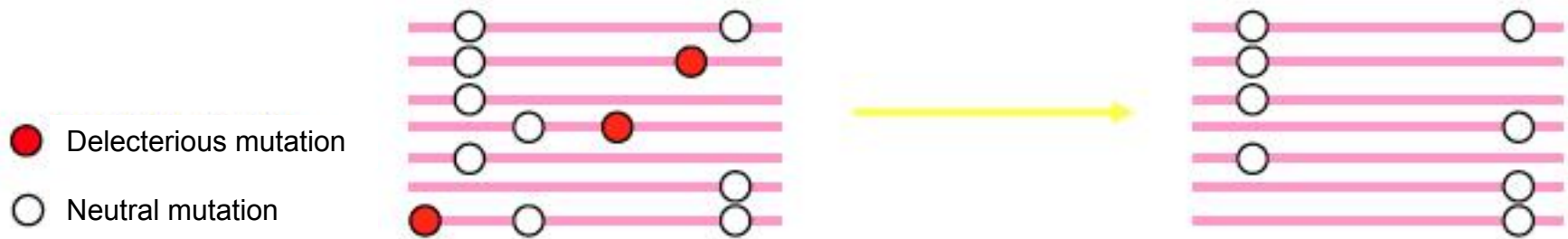
Selective sweeps



Nielsen (2005)

# Purifying selection

## Selection against deleterious mutations



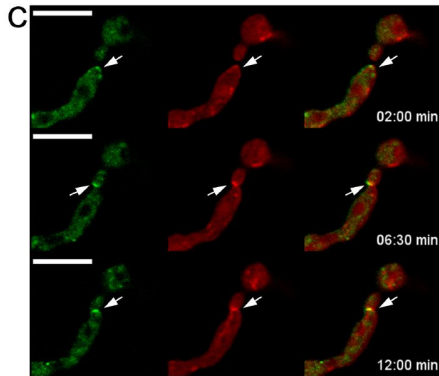
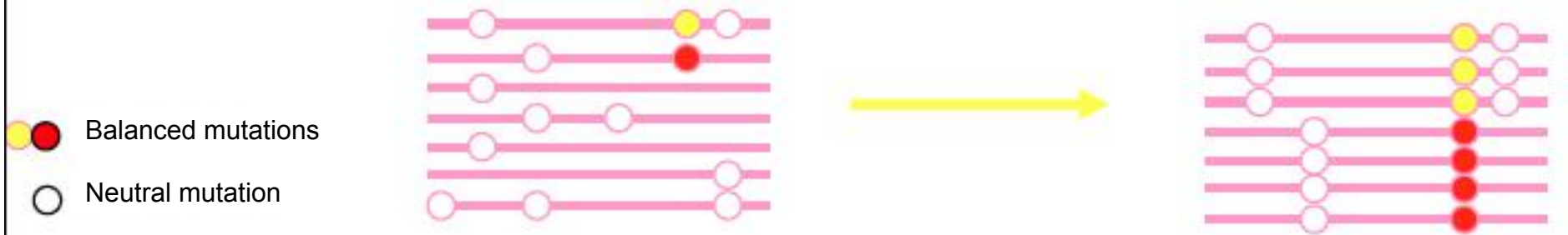
Related to gene functional constraints: Housekeeping genes

Selection against change!!

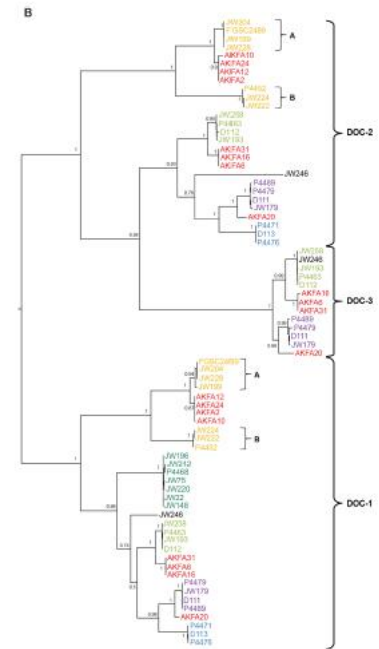


# Balancing selection

Increases the variability into the population



Recognition genes  
doc-1 doc-2 and  
doc-3 in *Neurospora*  
*sp.*



Neutrality tests are based on comparisons between estimators of  $q$

Watterson's $q_S$	$\theta_S = \frac{S}{a_n}$	$a_n = \sum_{i=1}^{n-1} \frac{1}{i}$
Tajima's $q_P$	$\theta_\pi = E(\Pi)$	
Fu & Li's $q_{he}$	$\theta_{\eta_e} = E(\eta_e)$	

$$T = \frac{\widehat{\theta}_1 - \widehat{\theta}_2}{\sqrt{\text{var}(\widehat{\theta}_1 - \widehat{\theta}_2)}}$$

## Tajima's D (1983)

Tajima's  $q_P$

$$\theta_{\Pi} = E(\Pi)$$

Watterson's

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

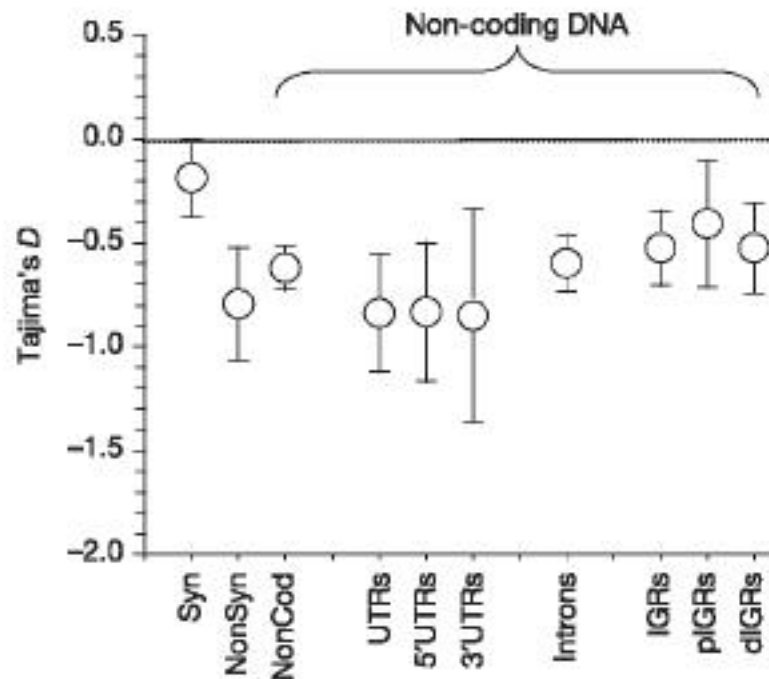
$$D = \frac{\Pi - \frac{S}{a_n}}{\sqrt{\text{var}\left(\Pi - \frac{S}{a_n}\right)}}$$

$D < 0$  Excess of rare variants -> Population expansion / selective sweep / purifying selection

$D > 0$  Deficit of rare variants -> Population decline / Population structure/ balancing selection

# Adaptive evolution of non-coding DNA in *Drosophila*

Peter Andolfatto<sup>1</sup>



**Figure 1 | Mean Tajima's  $D$  values for coding and non-coding DNA.** Means across loci are given with bars indicating two standard errors. The expectation of  $D$  under the neutral model is shown as a dotted line. Syn, synonymous sites; NonSyn, non-synonymous sites; NonCod, pooled non-coding DNA.

## Fu & Li's F (1993)

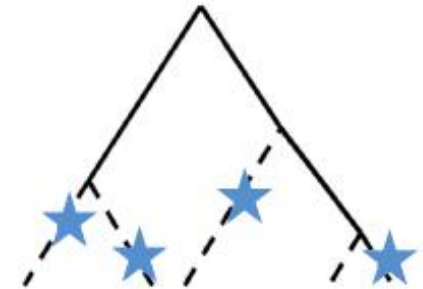
Tajima's  $d_P$

$$\theta_{\Pi} = E(\Pi)$$

Fu & Li's  $d_{he}$

$$\theta_{\eta_e} = E(\eta_e)$$

$$F = \frac{\Pi - \eta_e}{\sqrt{\text{var}(\Pi - \eta_e)}}$$



$D$  Should be less affected by selection than  $d_{he}$

Detection of recent selective sweeps

## Fay & Wu's H (2000)

Tajima's  $d_P$

$$\theta_\Pi = E(\Pi)$$

$\xi_i$

Number of mutations  
present  $i$  times

$$\theta_\xi = iE(\xi_i)$$

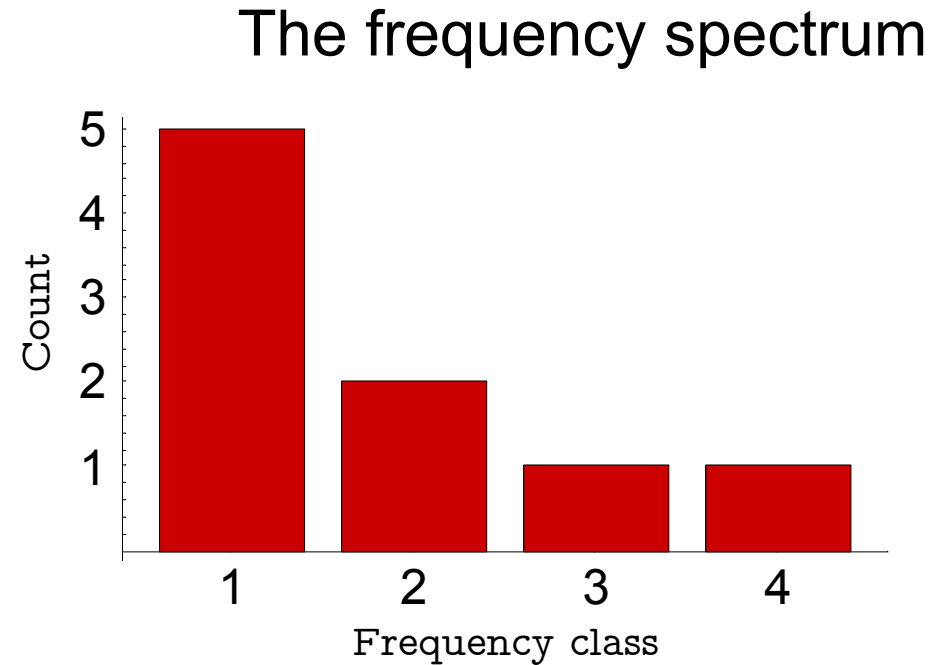
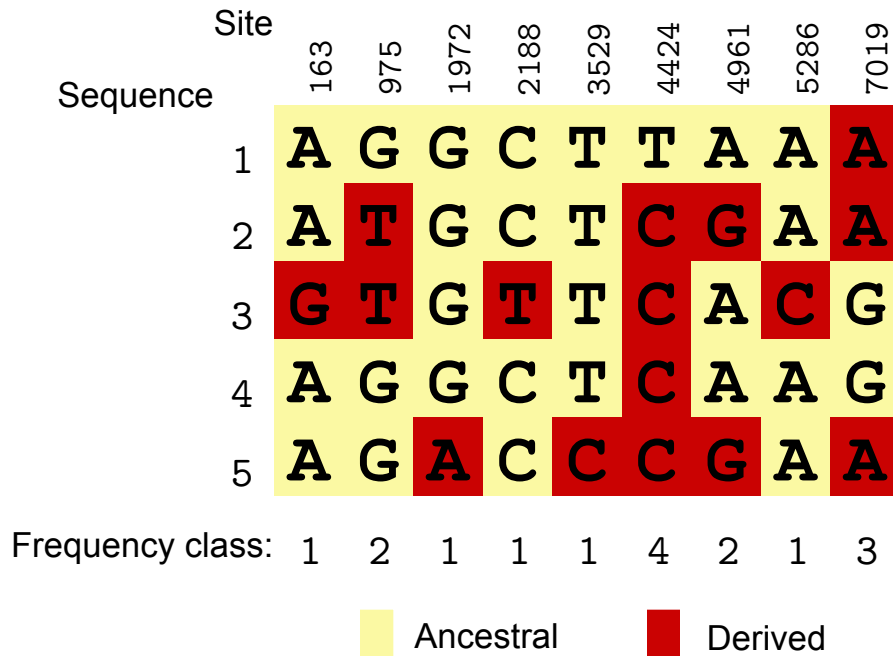
For  $i = 1$   $\theta_{\eta_e} = E(\eta_e)$

$$\theta_H = \frac{2}{n(n-1)} \sum_{i=1}^{n-1} i^2 \xi_i$$







$$H = \frac{\Pi - \frac{2}{n(n-1)} \sum_{i=1}^{n-1} i^2 \xi_i}{\sqrt{\text{var} \left( \Pi - \frac{2}{n(n-1)} \sum_{i=1}^{n-1} i^2 \xi_i \right)}}$$

Detection of recent selective sweeps

# The frequency spectrum: an example



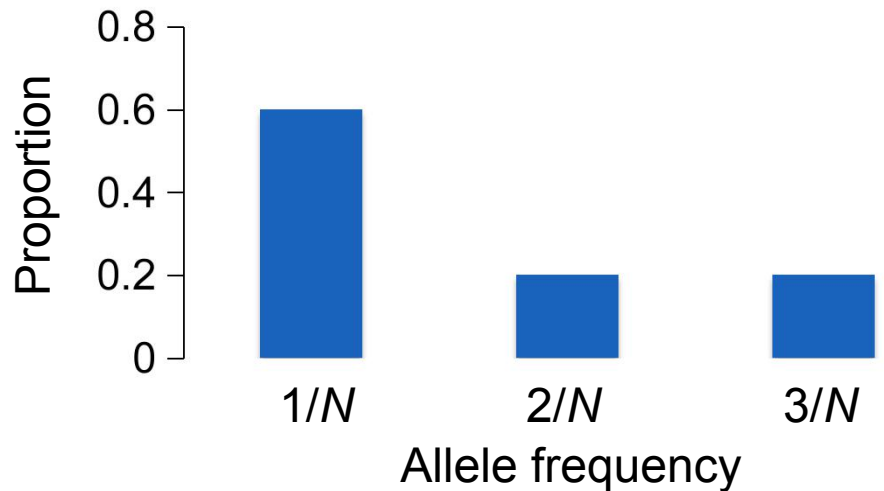
**SITE FREQUENCY SPECTRUM : SFS**




CC GCA GAG TTA CTAATC GAACG GCA GAG TTA CTAATC GAACC GCA AAG TTA CCAATT GAACC GCA GAG TTA CCAATC GAACC GCA AAG TTA CTAATC GAGCC GCA AAG TTA CTAATC GAA

$N = 6$   
 $L = 21$   
 $S = 5$

individuals  
sites on alignment  
polymorphic sites



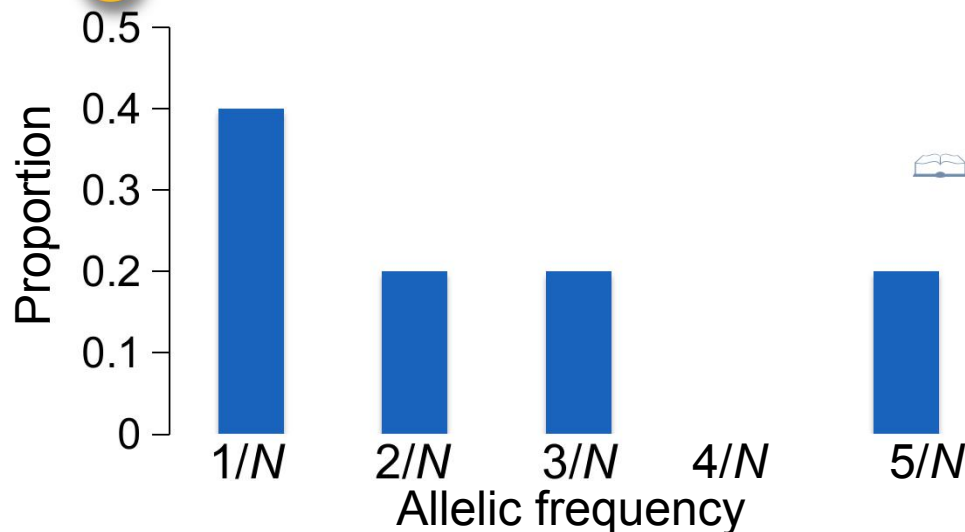
 SFS can be « *folded* » when we only take into account frequencies from 1/N to N/2



# Polymorphism estimation

## ***SITE FREQUENCY SPECTRUM : SFS***

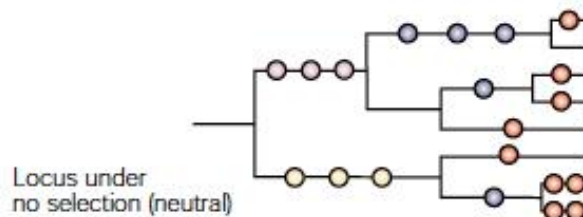
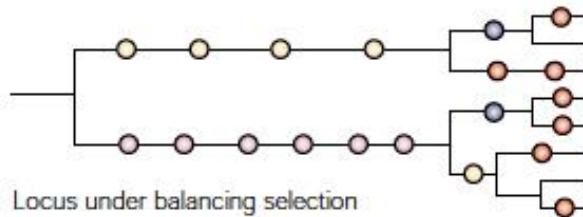
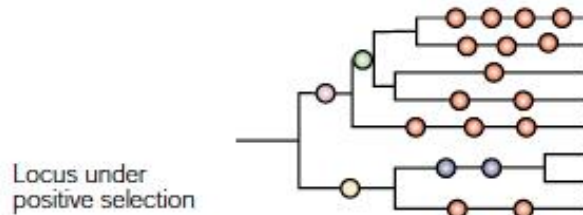
☹ CC GCA GAG TTA CTA ATC GA A  $N = 6$  individuals  
☹ CC **G** GCA GAG TTA CTA ATC GA A  $L = 21$  sites on alignment  
😊 CC GCA **A** AG TTA **C** A ATT **T** GA A  $S = 5$  polymorphic sites  
😊 CC GCA GAG TTA **C** A ATC GA A  
😊 CC GCA **A** AG TTA CTA ATC GA **G**  
😊 CC GCA **A** AG TTA CTA ATC GA A



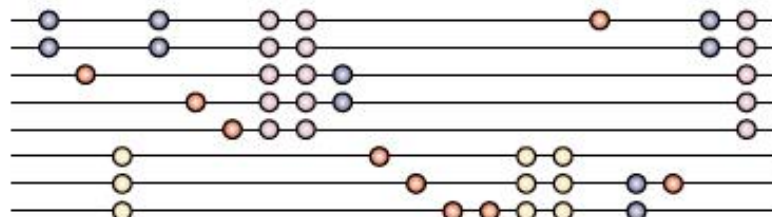
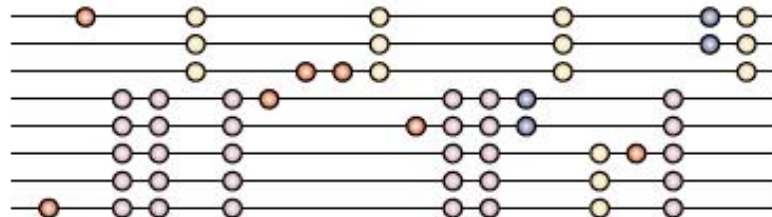
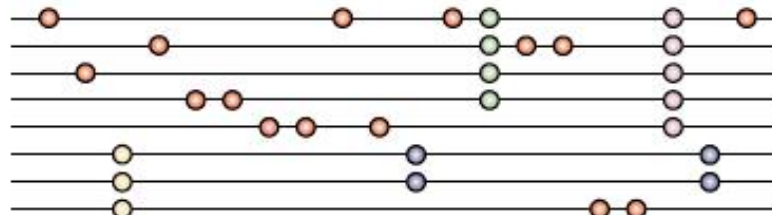
📖 le SFS can be « *folded* » or « *unfolded* » if we know the ancestral state

# Evolutionary forces like selection shape the SFS

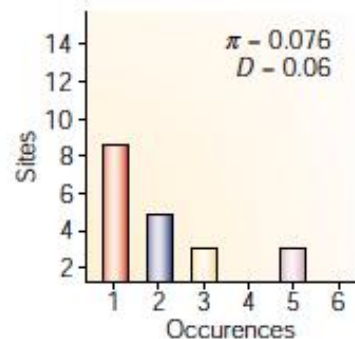
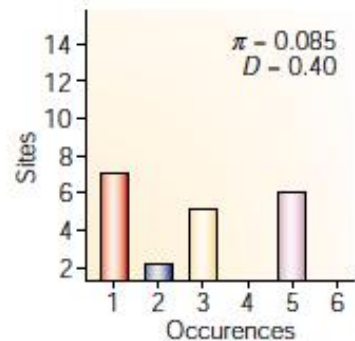
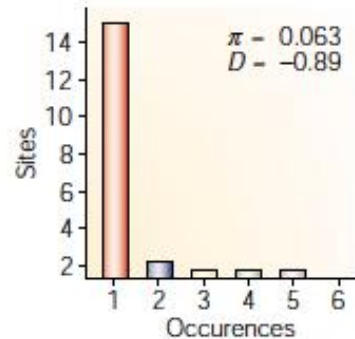
**a Genealogies**



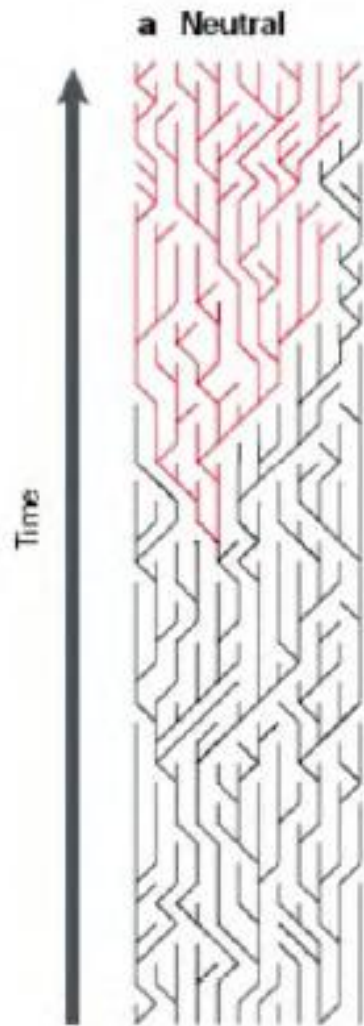
**b Haplotypes**



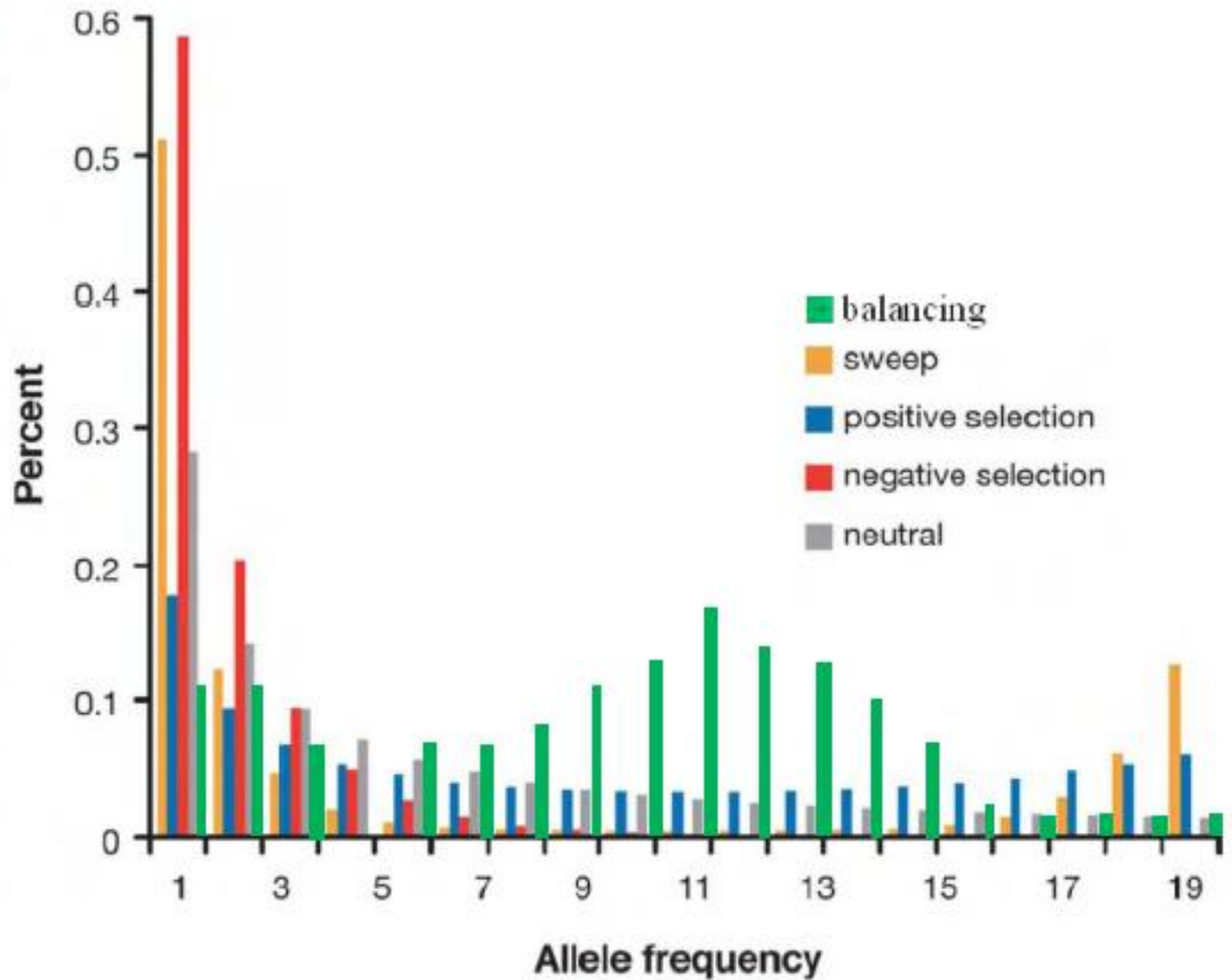
**c Site frequency spectra**



# Neutrality

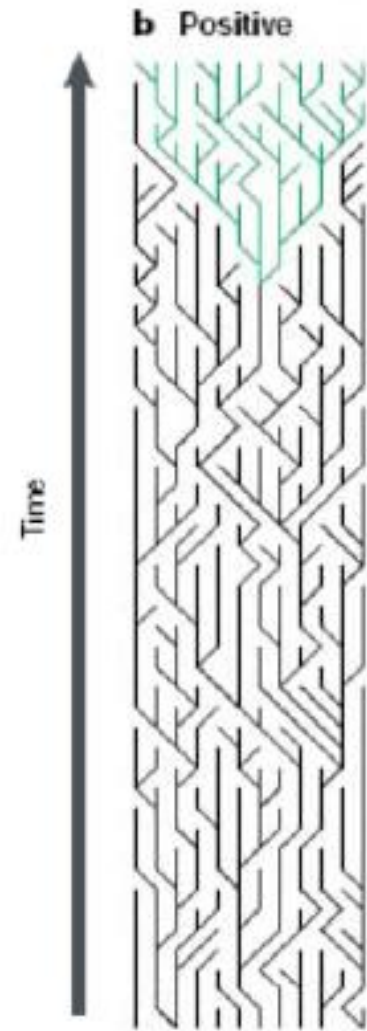


(Bamshad & Wooding)

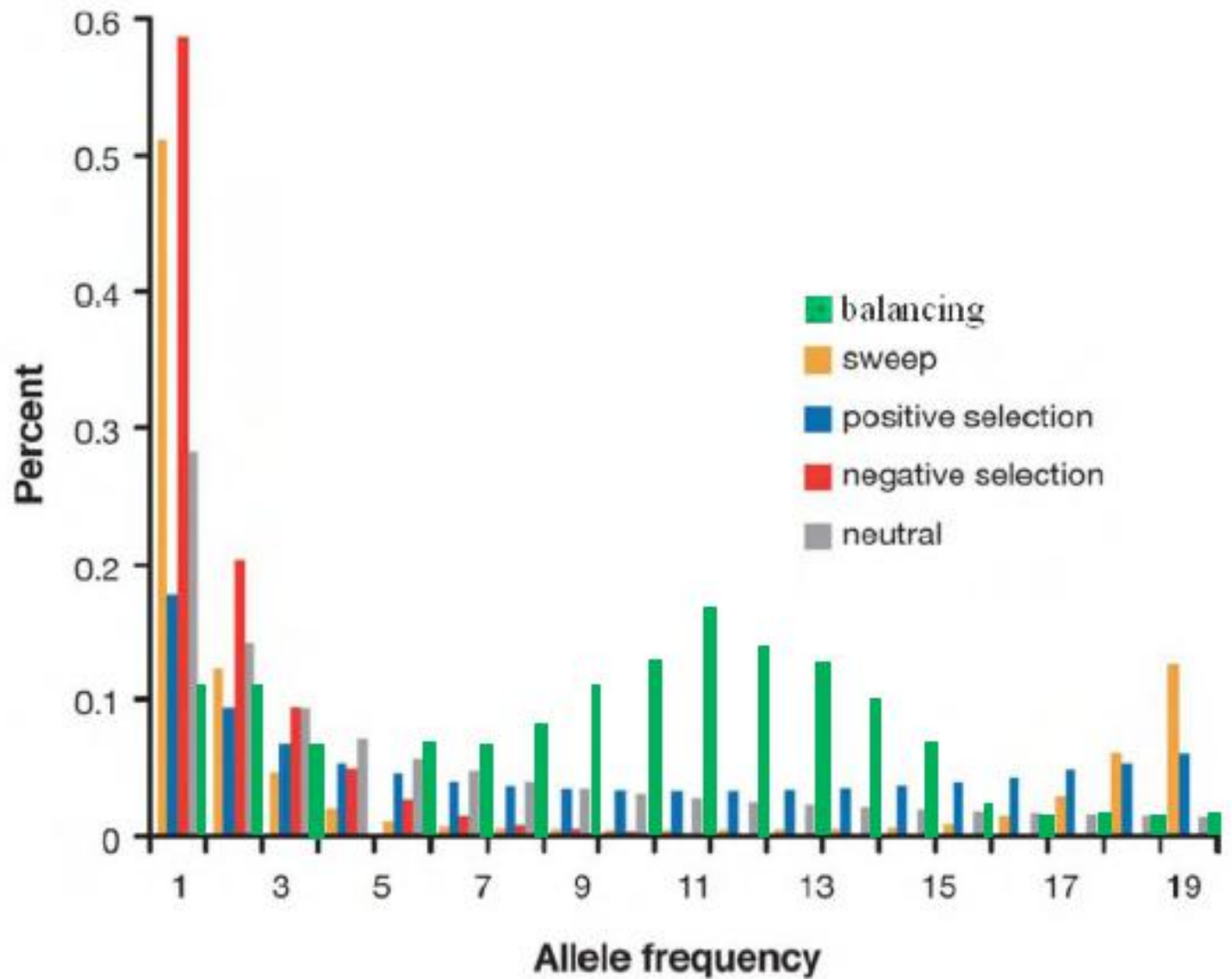


(d'après Nielsen, 2005)

# Positive selection

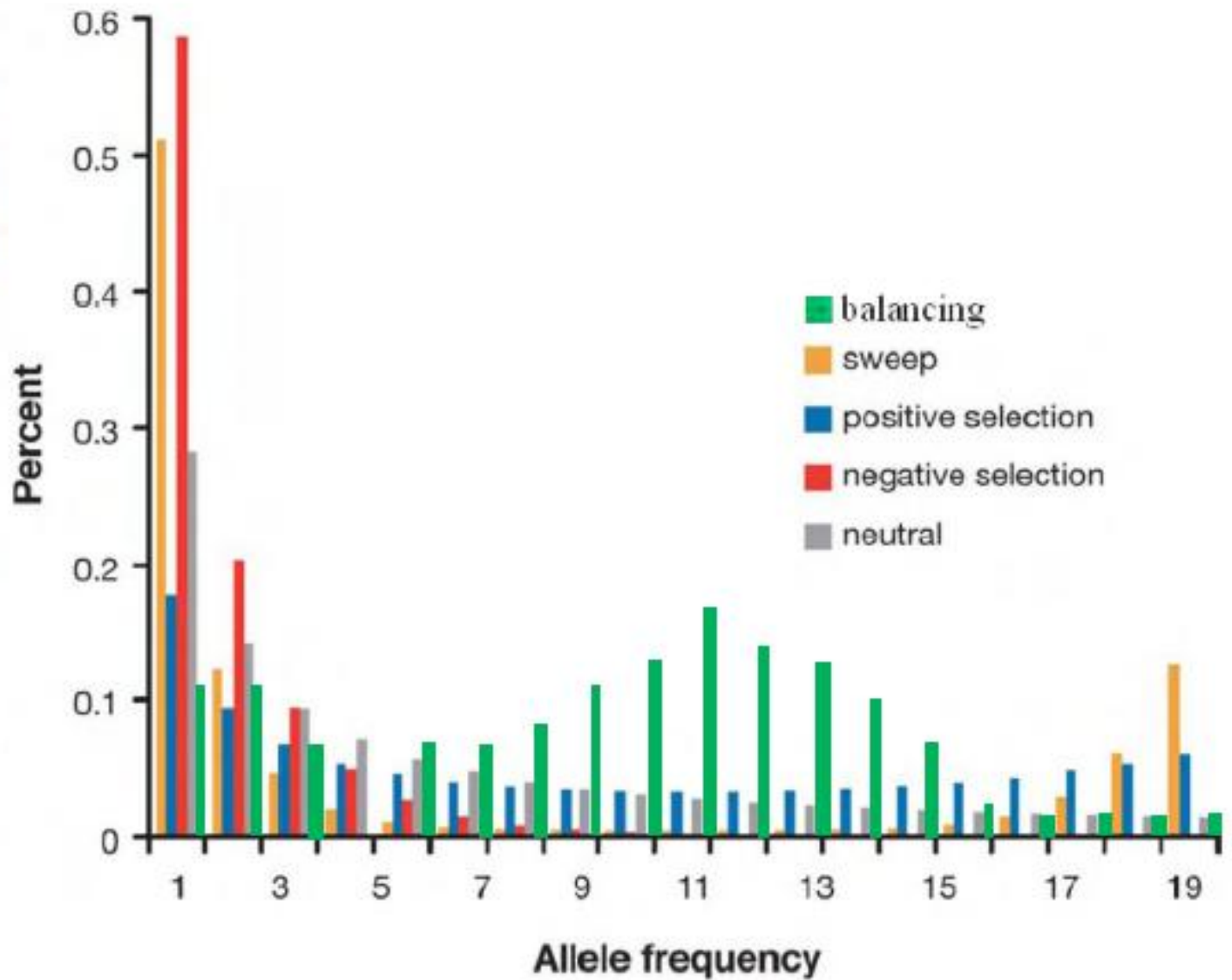


(Bamshad & Wooding)



(d'après Nielsen, 2005)

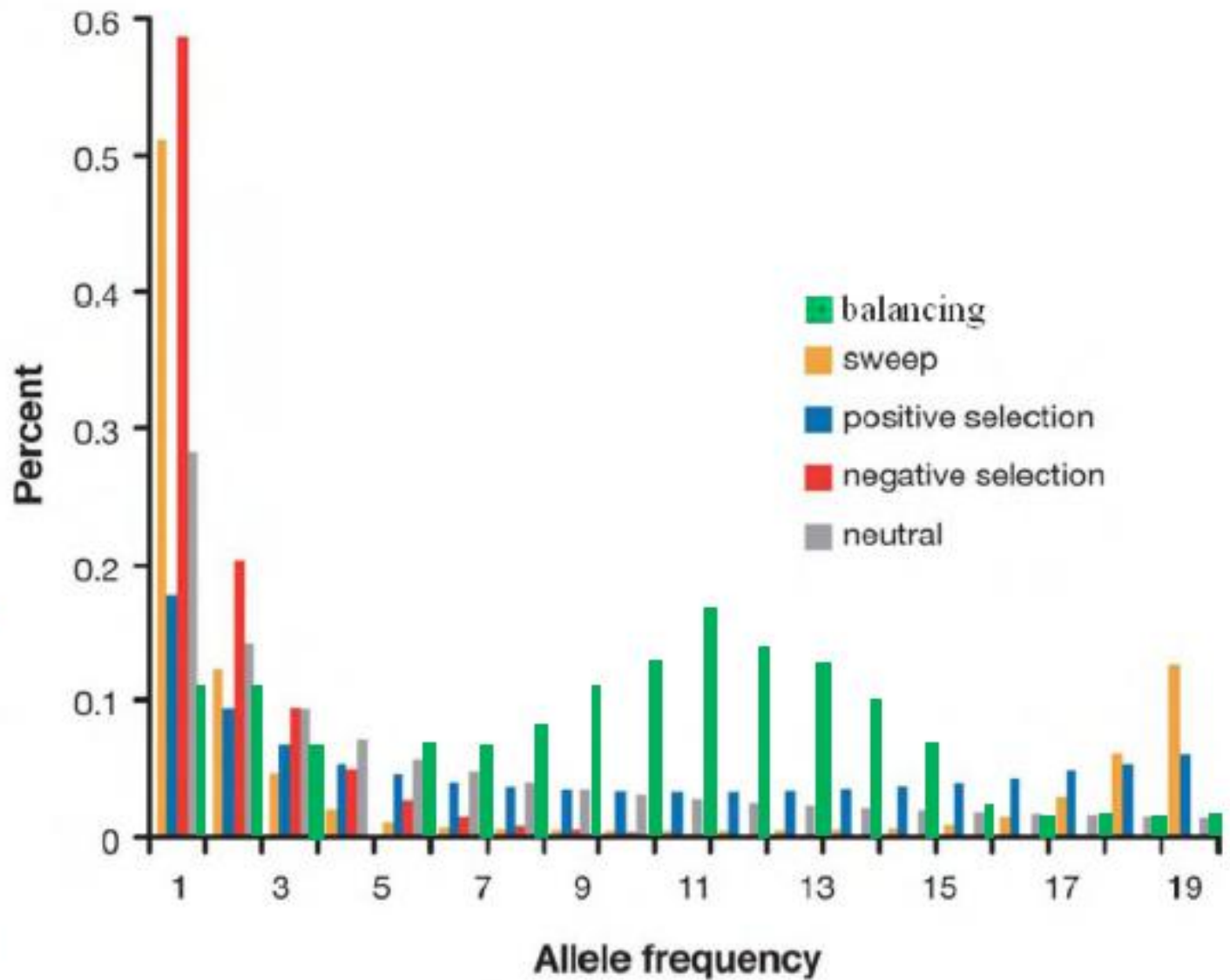
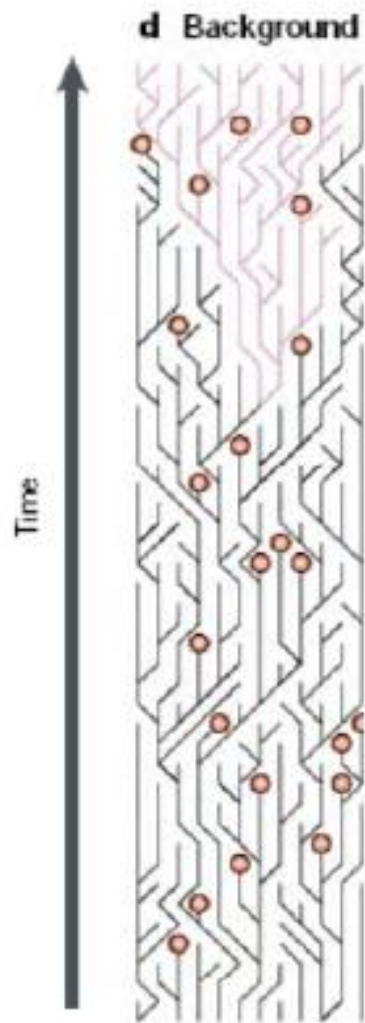
# Balancing selection



(Bamshad & Wooding)

(d'après Nielsen, 2005)

# Purifying/background selection

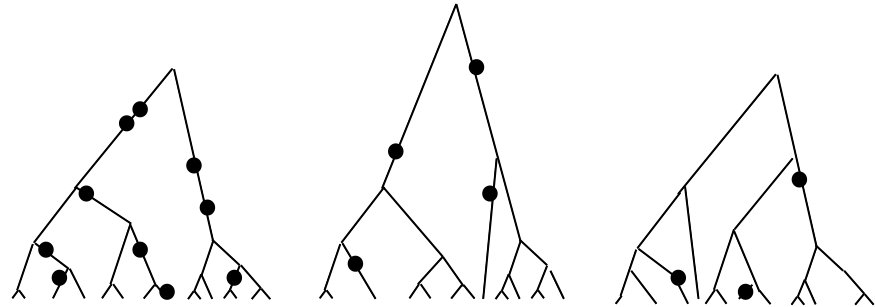


(Bamshad & Wooding)

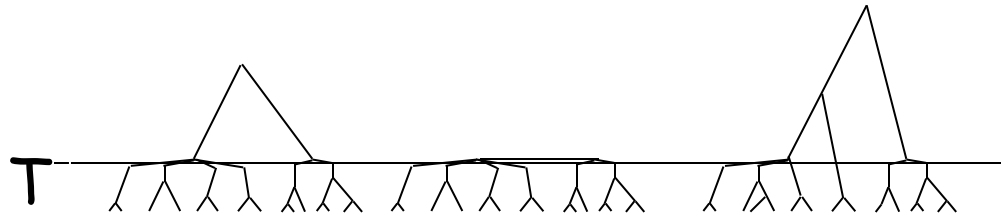
(d'après Nielsen, 2005)

# Demography and selection can be confounded

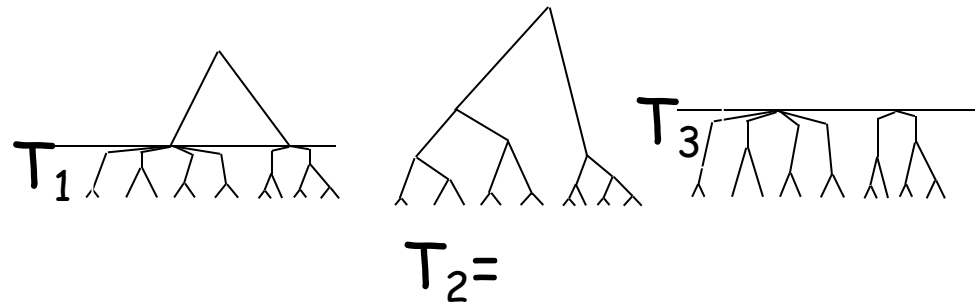
**M<sub>1</sub>: neutral, constant size**  
 $p$  parameters ( $q_1, \dots, q_p$ )



**M<sub>2</sub>: bottleneck**  
 $p+2$  parameters ( $T, S, q_1, \dots, q_p$ )



**M<sub>3</sub>: selective sweep**  
 $3p$  parameters  
( $T_1, S_1, q_1, \dots, T_p, S_p, q_p$ )

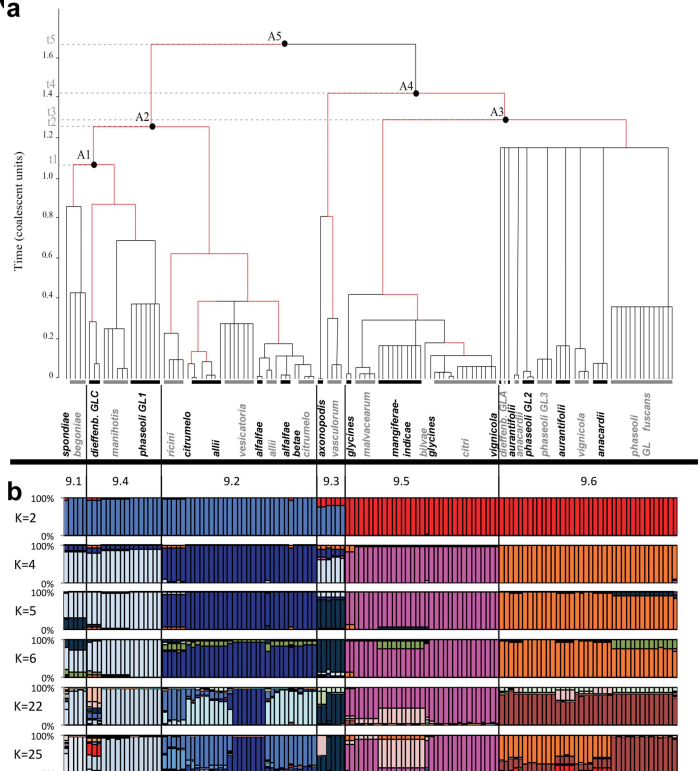




assume ONE population (perturbation of the SFS)

## Important to figure out how many populations are in my dataset

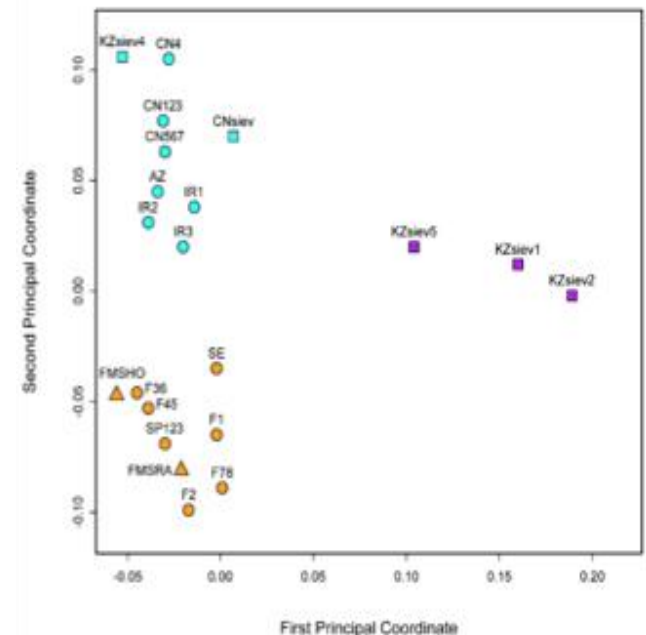
# Use of STRICTIF



*Xanthomonas axonopodis*

or

# PCA



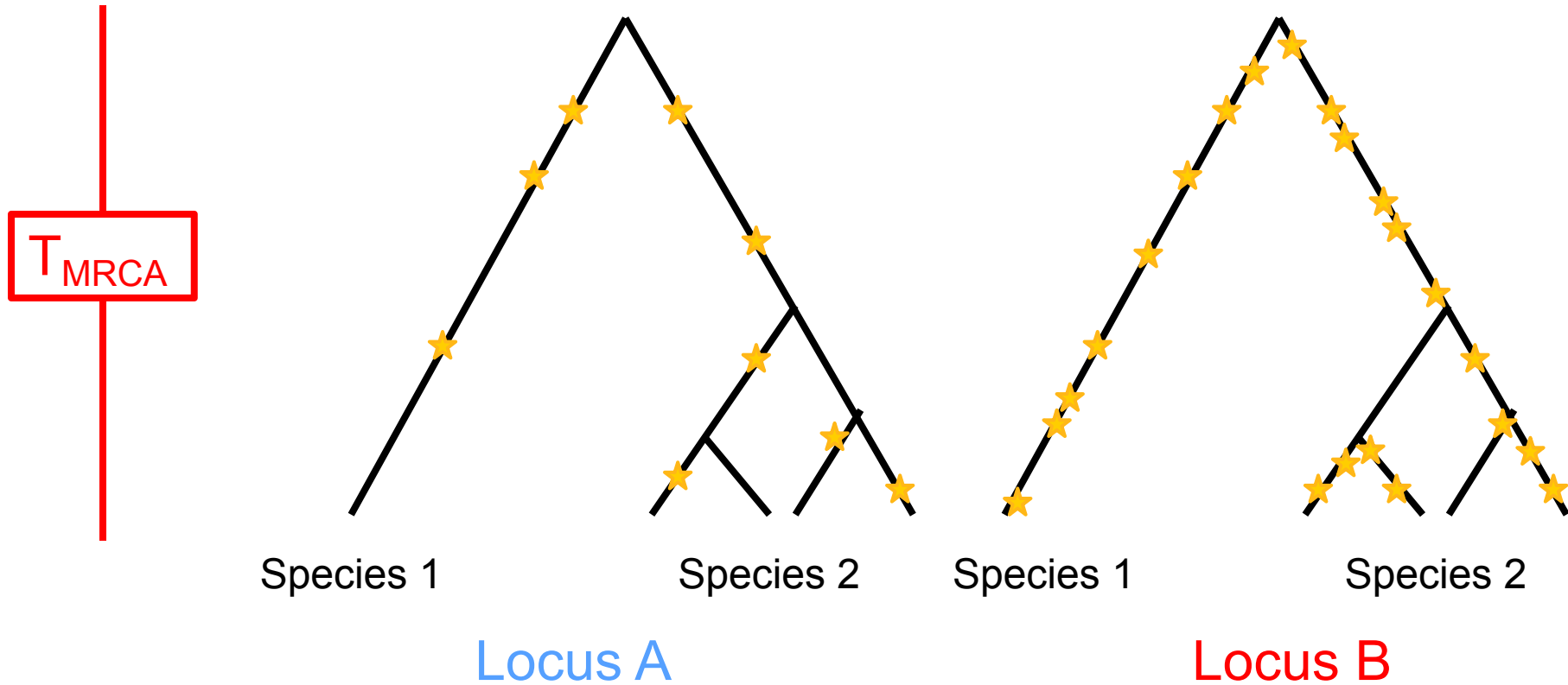
*Venturia inaequalis*



There are other tests based on polymorphism and divergence

To be continued....

# The Hudson Kreitman & Agadé (HKA) test (Hudson, Kreitman & Agadé, 1987)



Under the neutral hypothesis

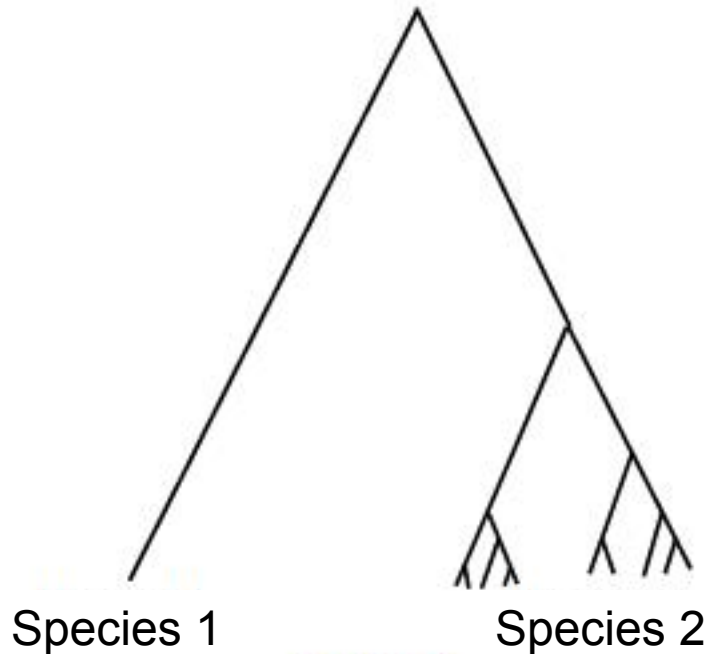
$$D_A = 2T\mu_A$$

$$S_A = 4N\mu_A a_n \\ = \theta_A a_n$$

$$D_B = 2T\mu_B$$

$$S_B = 4N\mu_B a_n \\ = \theta_B a_n$$

$$\frac{\theta_A}{D_A} = \frac{\theta_B}{D_B} = \frac{4N\mu_i a_n}{2T\mu_i} = \frac{4Na_n}{2T} \quad \text{constant } \forall \text{ locus}$$



Locus A

$$D_A = 2T\mu_A$$

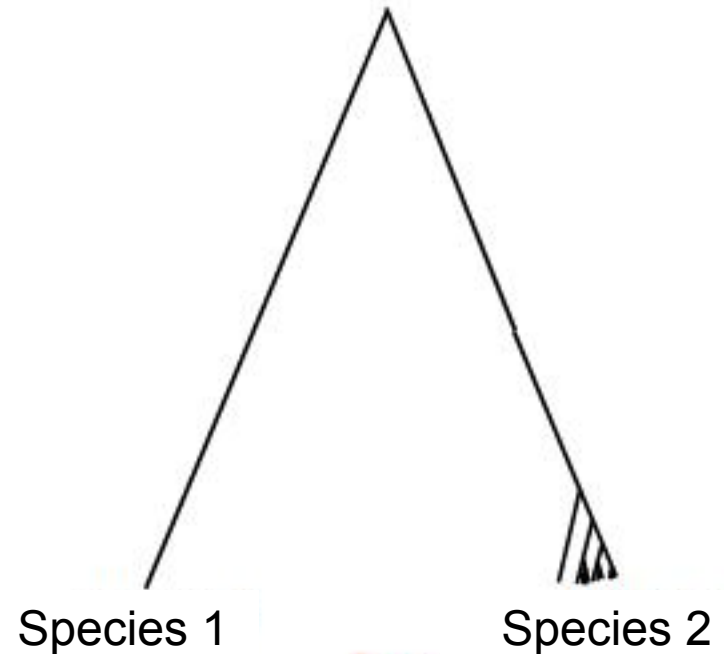
$$S_A = 4N\mu_A a_n$$

$$= \theta_A a_n$$

$$D_B = 2T\mu_B \sim D_A$$

$$S_B = 4N\mu_B a_n$$

$$= \theta_B a_n < S_A$$



Locus B

$$\frac{\theta_A}{D_A} > \frac{\theta_B}{D_B}$$

Reduced polymorphism at locus B probably by selection because it can not be due to:

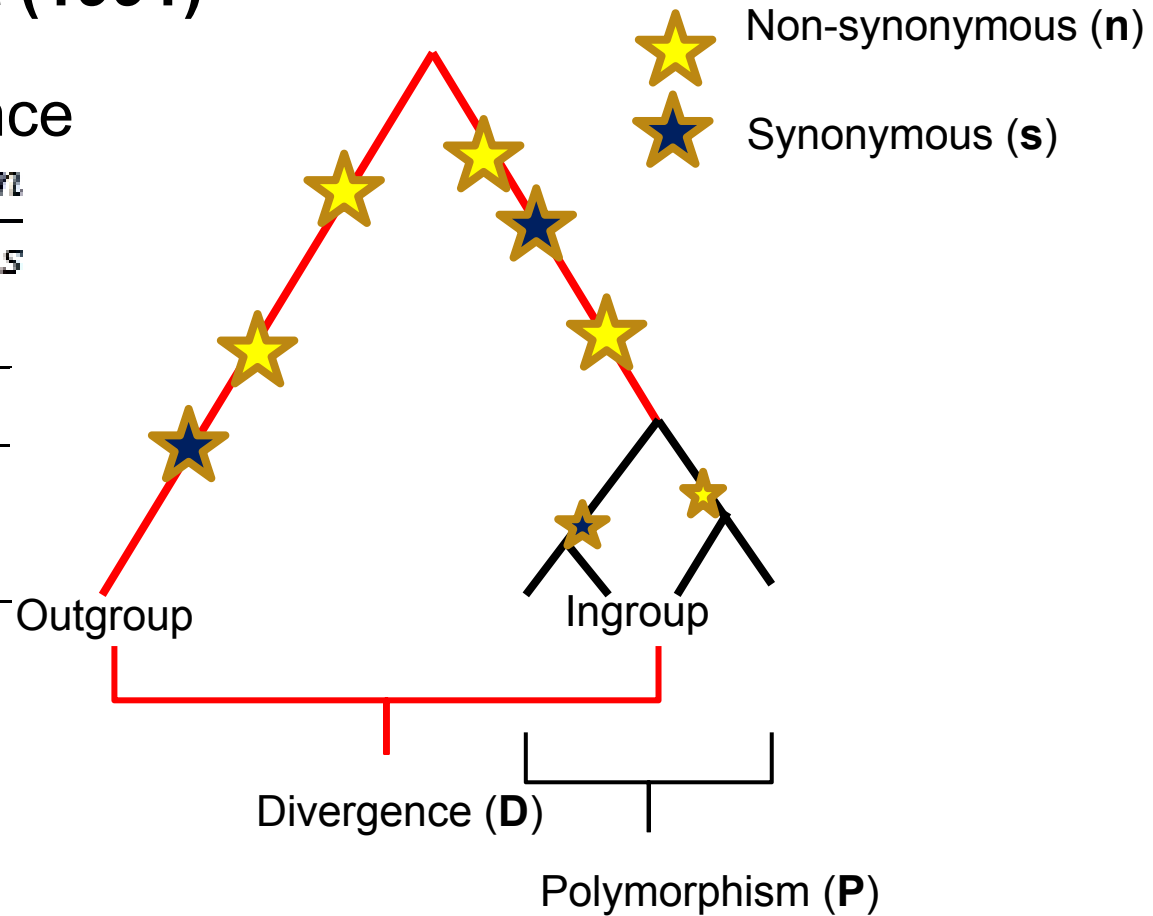
- **reduced population size (too much polymorphism at locus A) or**

# •McDonald Kreitman test (1991)

•Test for adaptive divergence

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

	Divergence	Polymorphism
NS	4	1
S	2	1



Mathematics	Biology
$\alpha = 0$ $P_n/P_s$	Neutral
$\alpha < 0$ $D_n/D_s < P_n/P_s$	<b>Polymorphism (Pn) excess</b>
$\alpha > 0$ $D_n/D_s > P_n/P_s$	<b>Divergence (Dn) excess</b>

Mathematics	Biology
$\alpha = 0$ $P_n/P_s$	Neutral
$\alpha < 0$ $D_n/D_s < P_n/P_s$	Polymorphism ( $P_n$ ) excess
$\alpha > 0$ $D_n/D_s > P_n/P_s$	Divergence ( $D_n$ ) excess

• BUT there is estimation biases

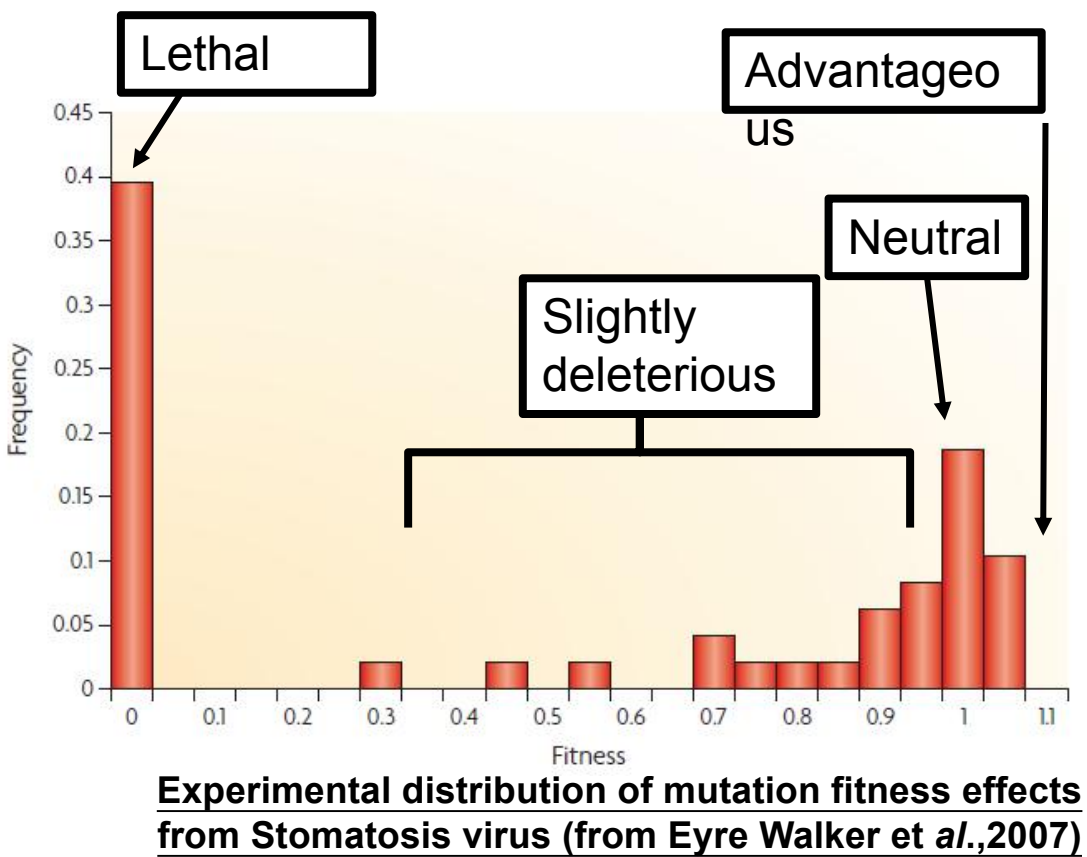
- Distribution of fitness effects

slightly deleterious mutations which segregate within ingroup (increase  $P_n$ )

★ 3 kind of mutation effects :

- **Neutral:**
  - Absence of effects on the fitness
- **Advantageous:**
  - Increase the fitness
- **Deleterious (lethal + Slightly deleterious):**
  - Complex effects (Fitness variation between 0 and 1)

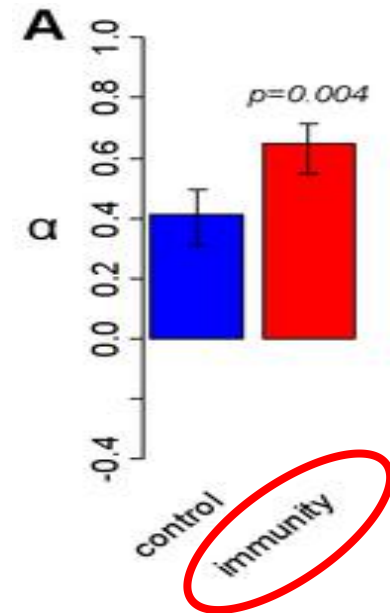
★ Underestimation of  $\alpha$



# Molecular evolution gives us a lot of information

## Example :

- Genes involved in interactions with other organism show higher adaptive rates than others
- Ex: Adaptive evolution for immunity genes in *Drosophila* (Obbard et al., 2009)
  - Adaptation for efficient immunity system against pathogen
- What about pathogen evolution ?



Adaptive rate comparison (From Obbard et al., 2009)