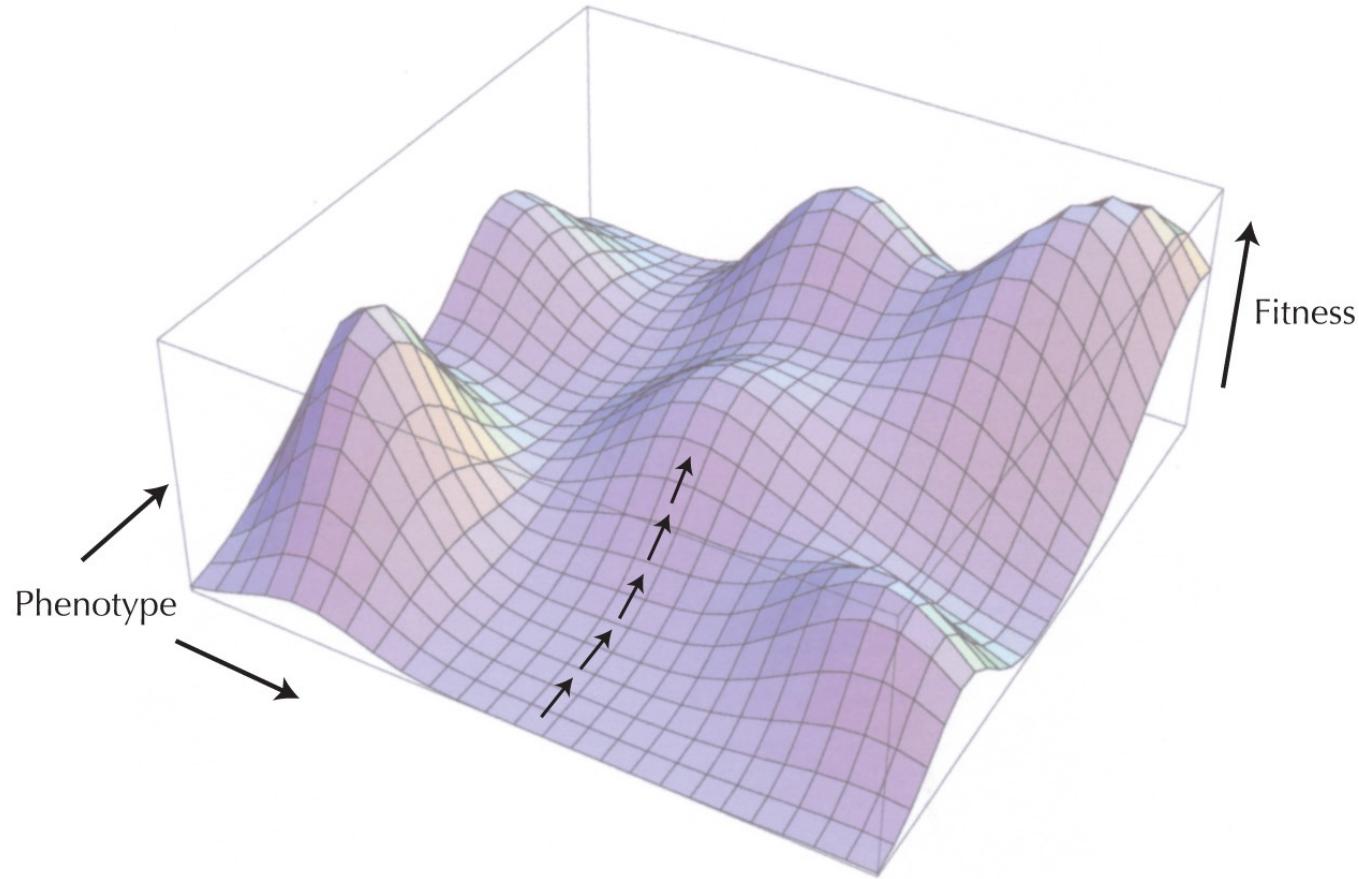
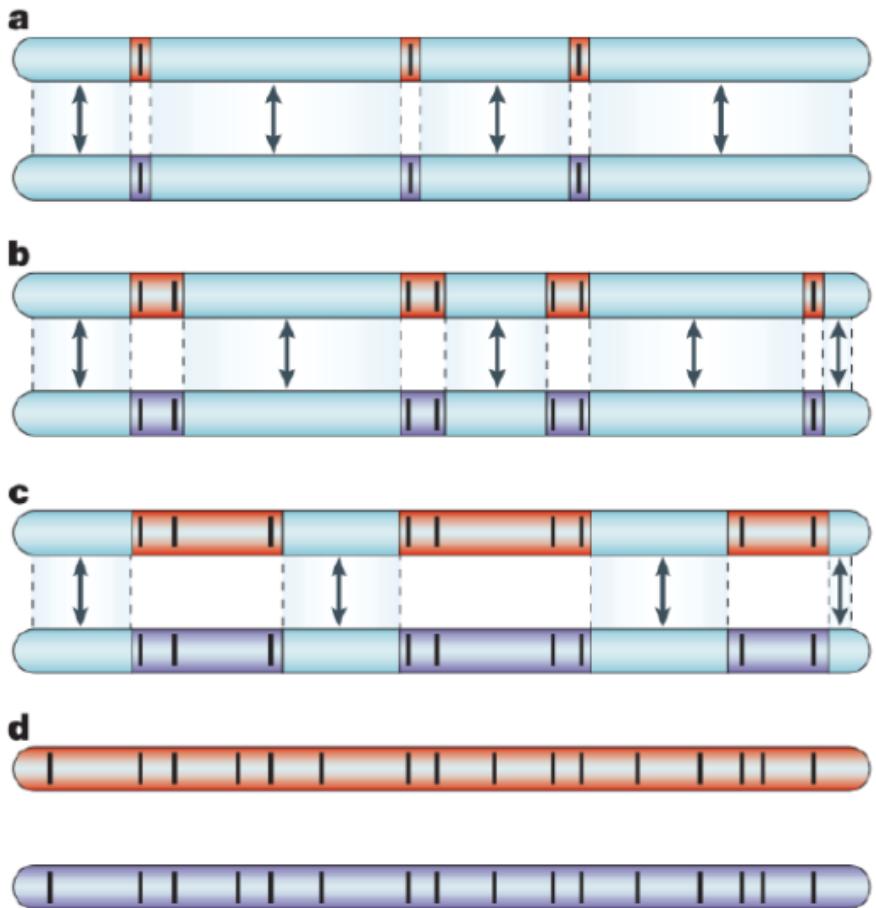


Selection



The genic concept of speciation

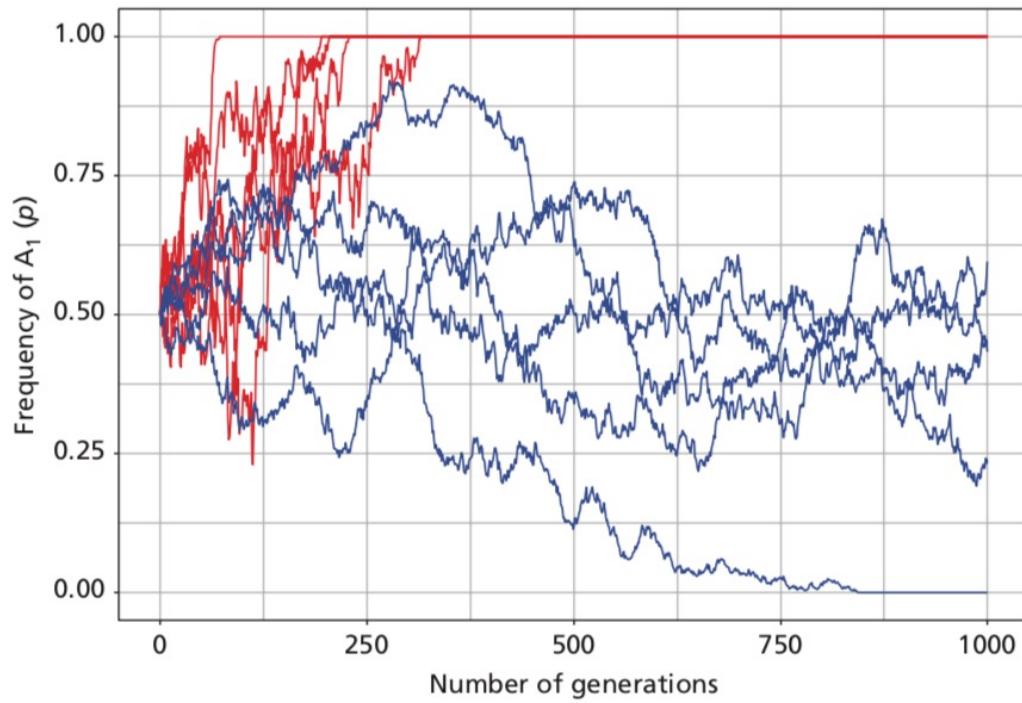


Divergent loci resist gene flow

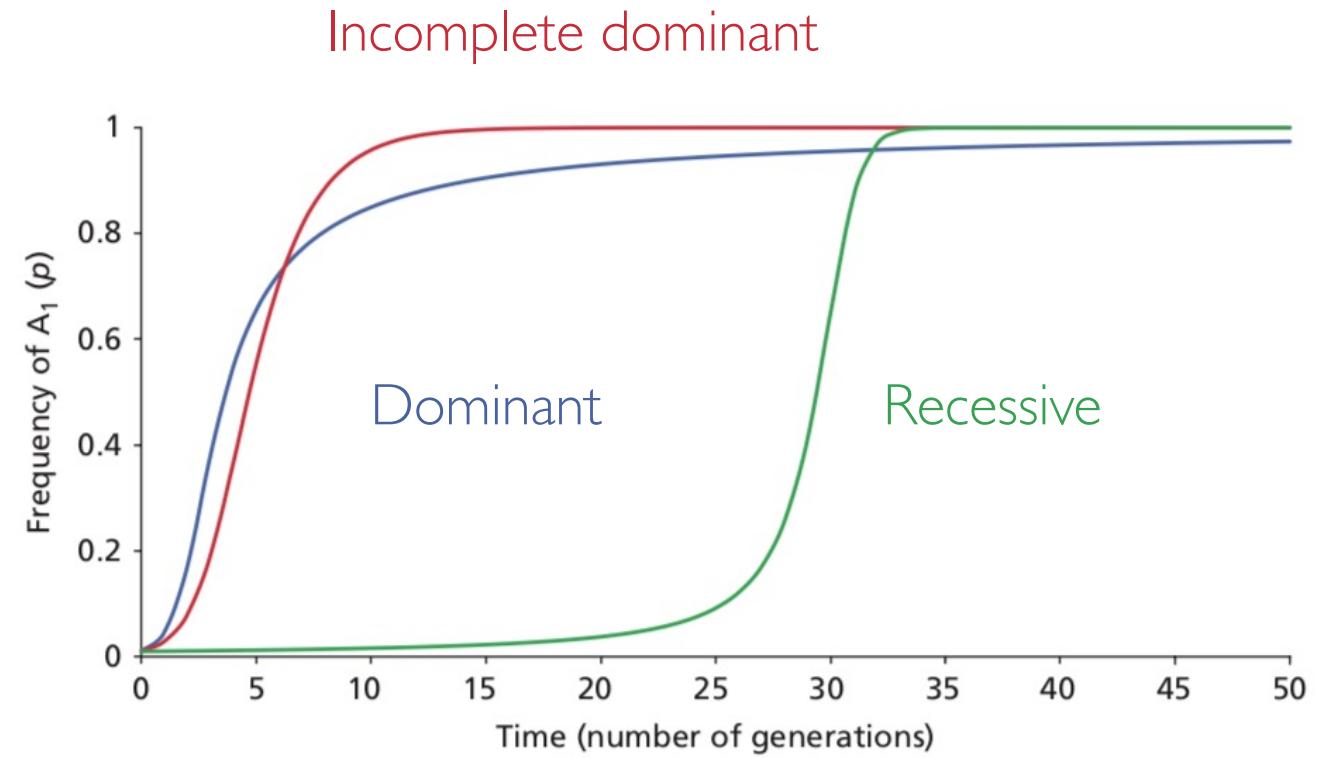
Gene flow continues but
linkage builds and divergent
regions grow

Complete reproductive
isolation evolves

Selection and drift at the genetic level

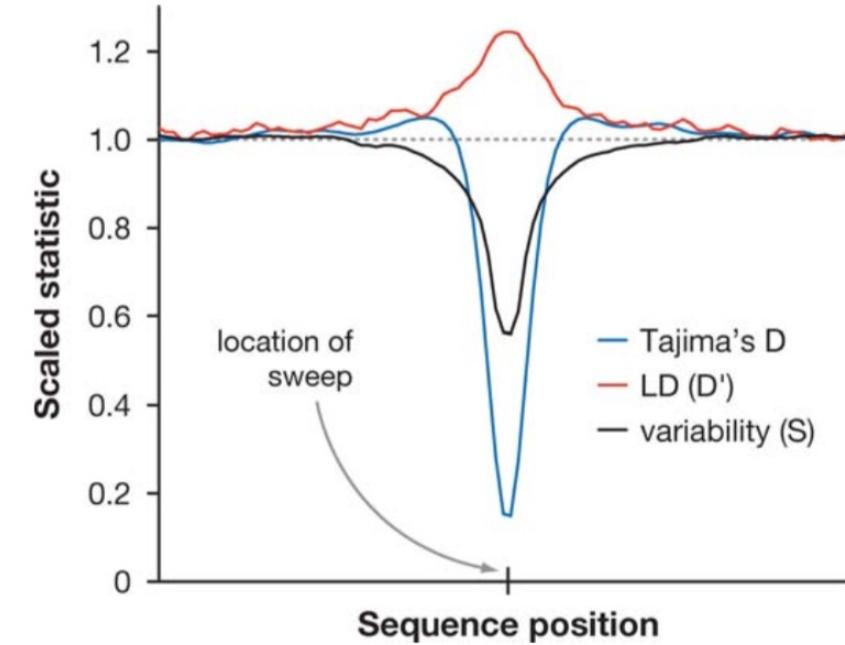
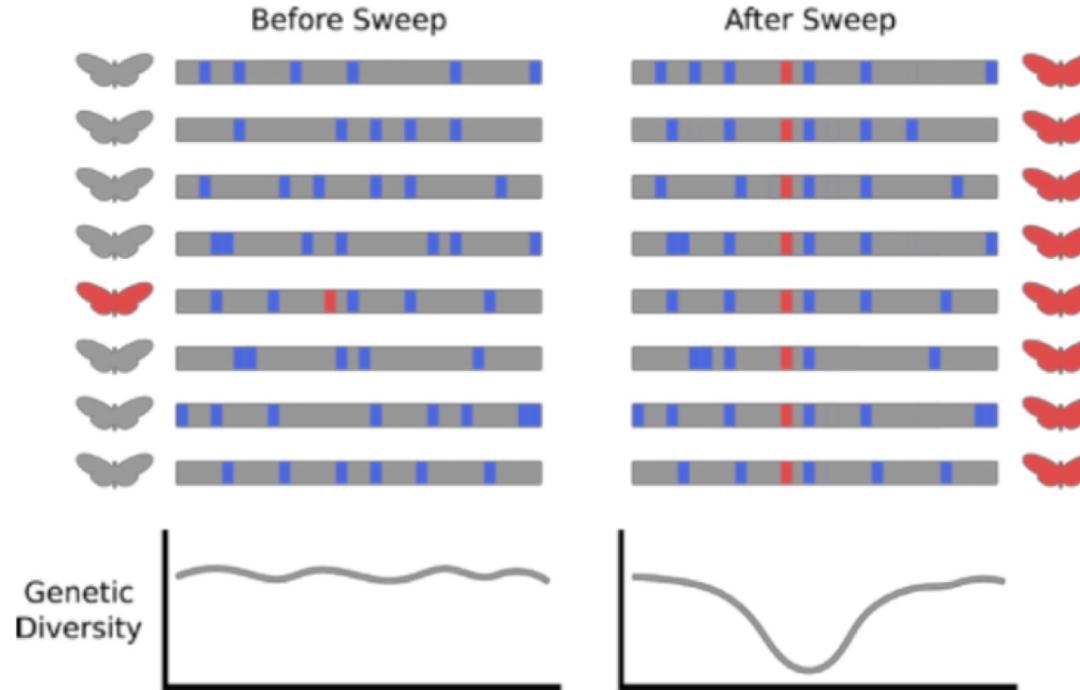


Genetic drift



Directional selection

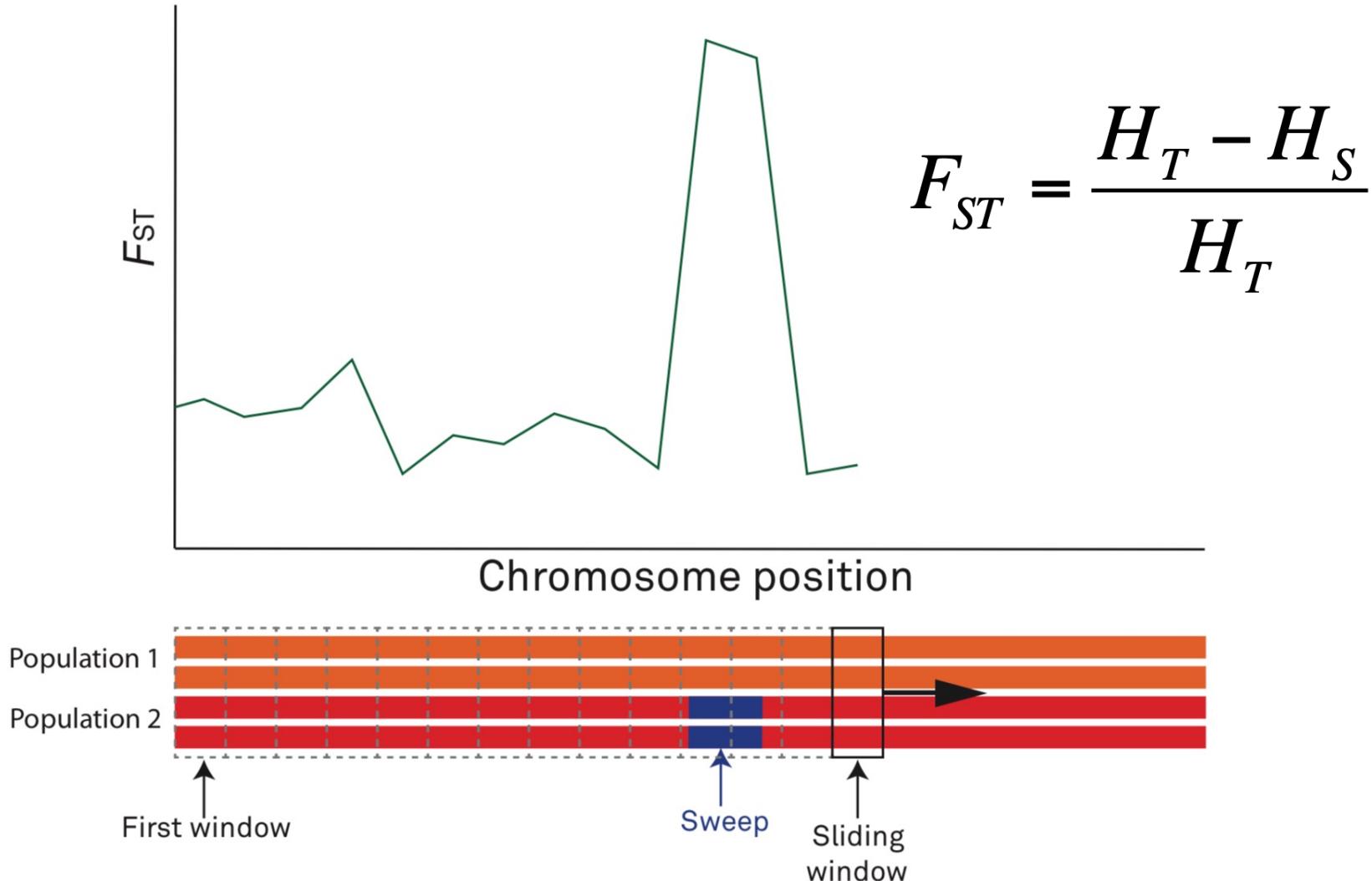
Genetic signatures of selective sweeps



Nielsen (2005) *Nat Rev Gen*

- Reduced genetic diversity/variation
- Increased linkage disequilibrium
- Increased genetic differentiation compared to other populations

Genetic differentiation to detect selection



An idea predating genomic data

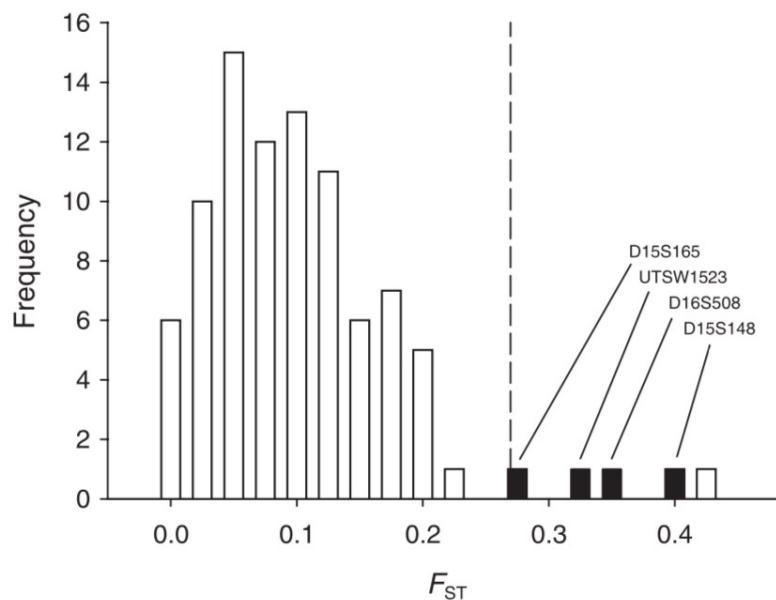
DISTRIBUTION OF GENE FREQUENCY AS A TEST OF THE
THEORY OF THE SELECTIVE NEUTRALITY OF
POLYMORPHISMS^{1,2}

R. C. LEWONTIN AND JESSE KRAKAUER

$$E[\text{var}(F_{ST})]:\text{var}(F_{FST})$$

Lewontin & Krakauer (1973) Genetics

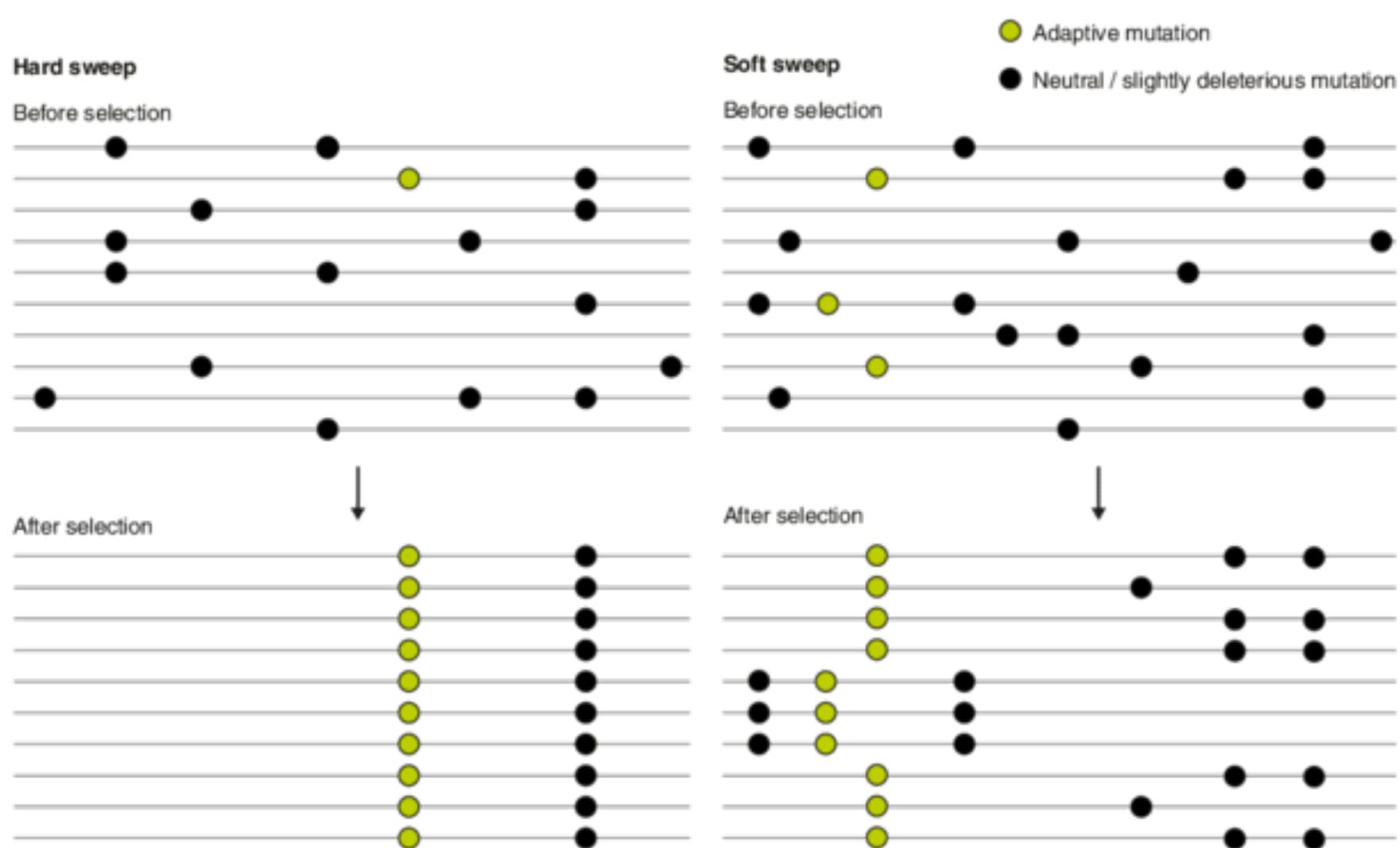
Microsatellite loci putatively under selection in humans



Outlier detection

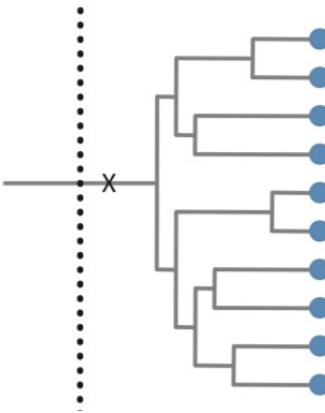
Storz et al (2004) MBE

Hard sweep and soft sweep

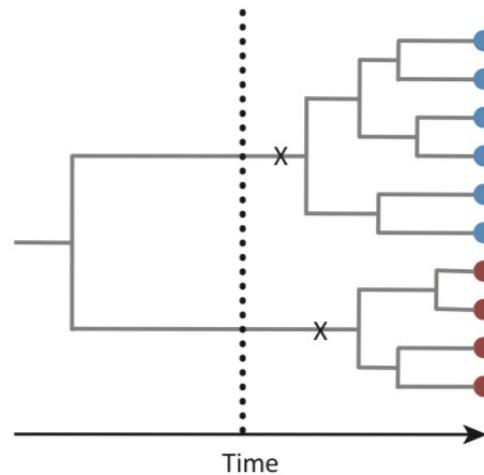


Types of sweep: hard vs soft

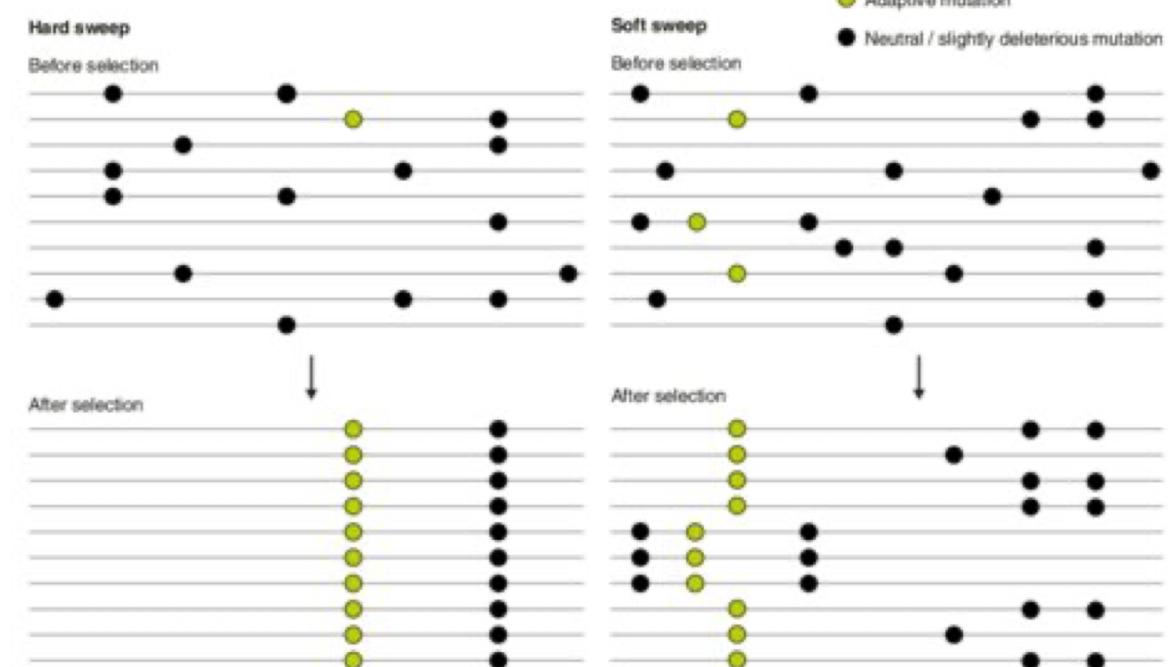
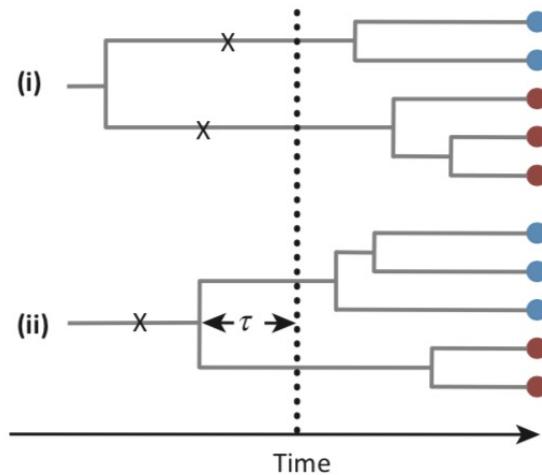
(A) Classic hard sweep



(B) Soft sweep (*de novo* mutations)

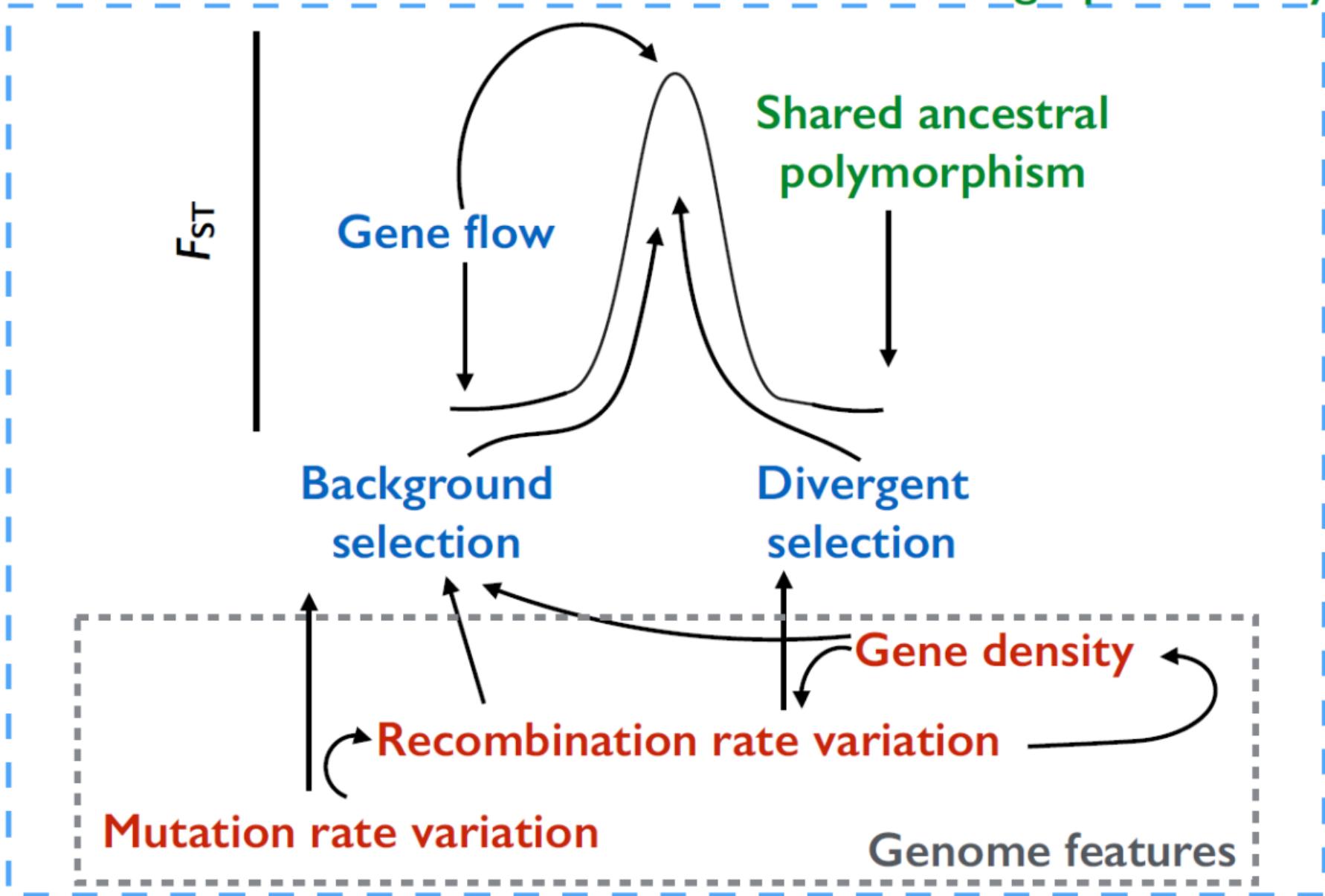


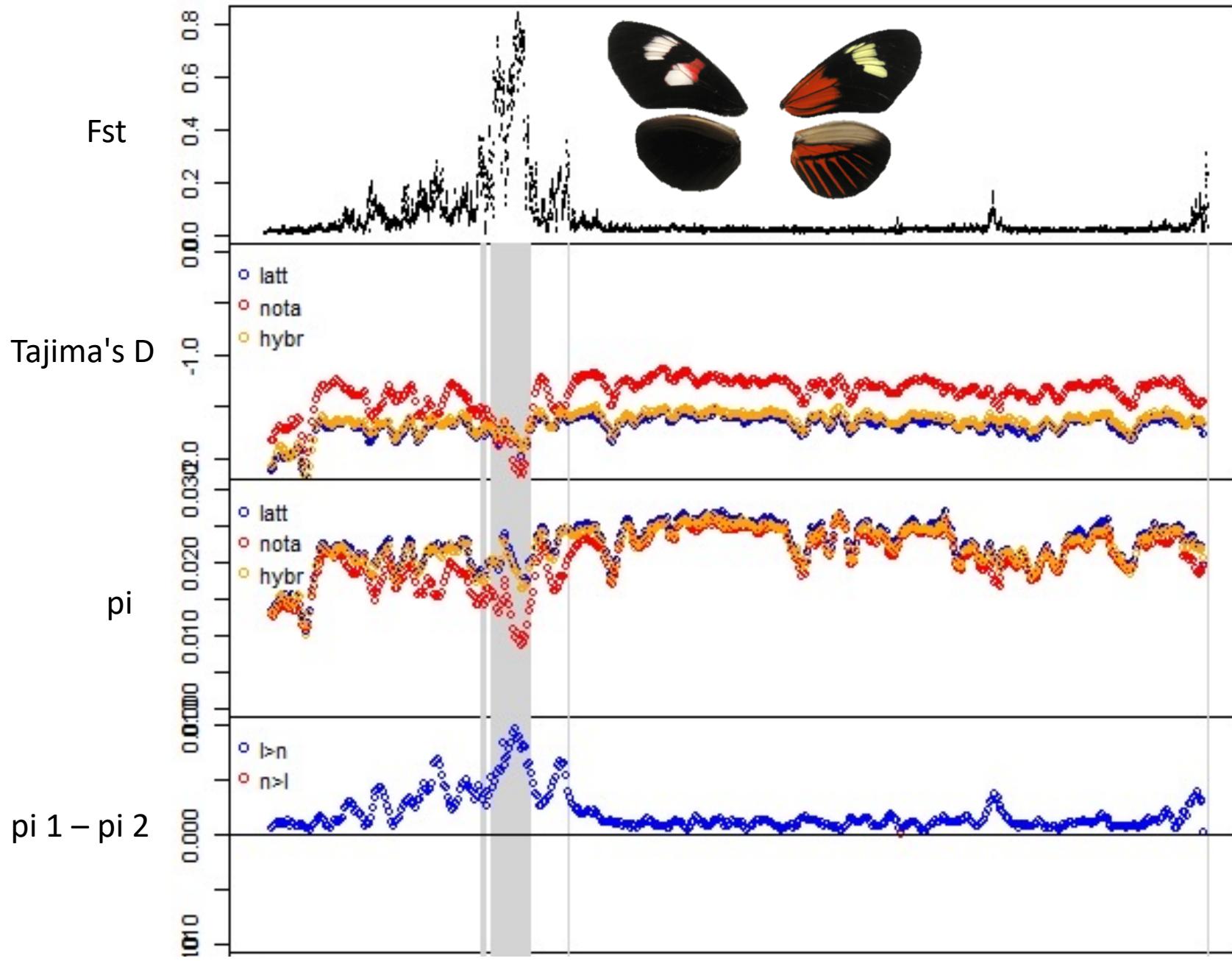
(C) Soft sweep (standing variation)



Standing variation and independent *de novo* mutations

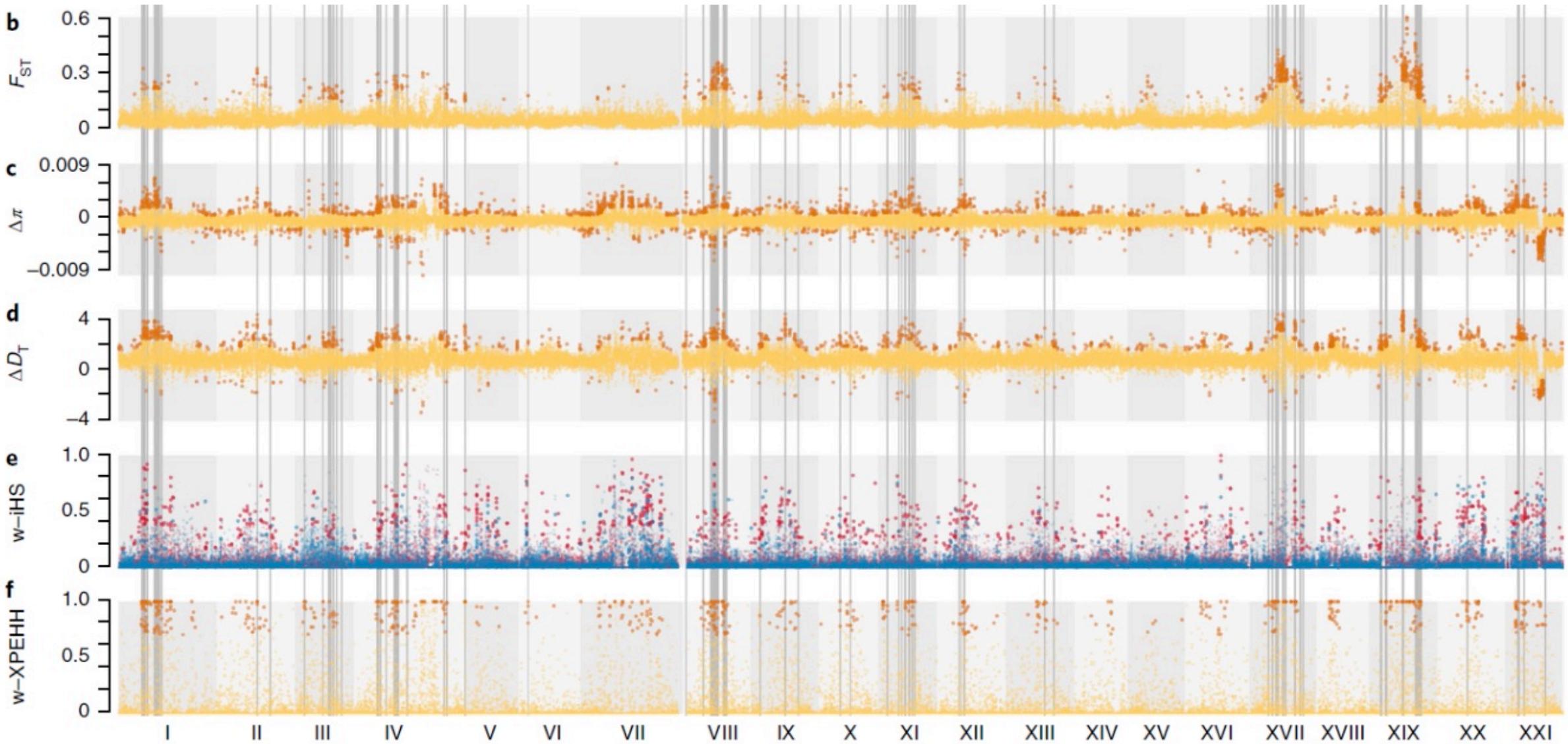
Confounding factors

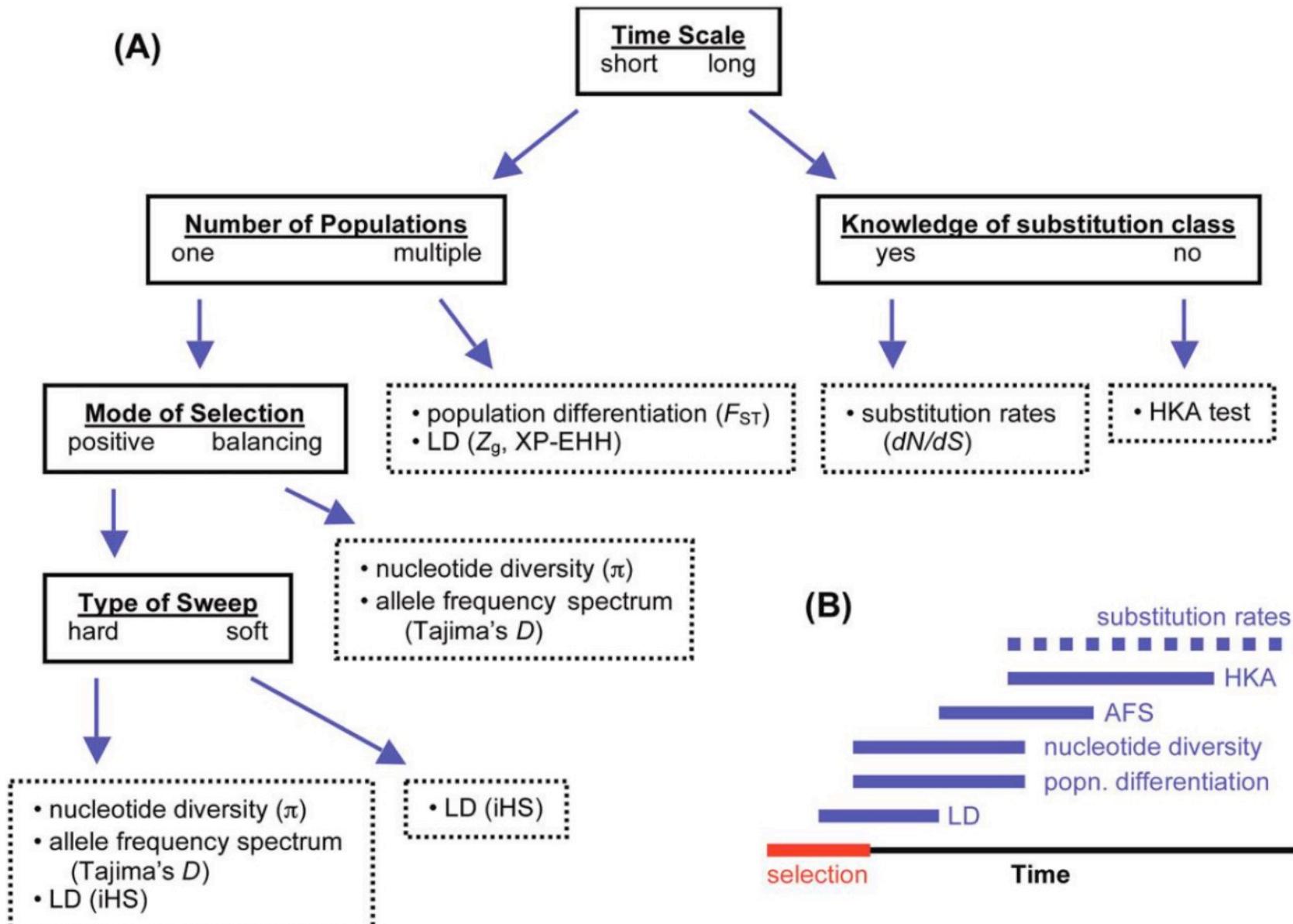


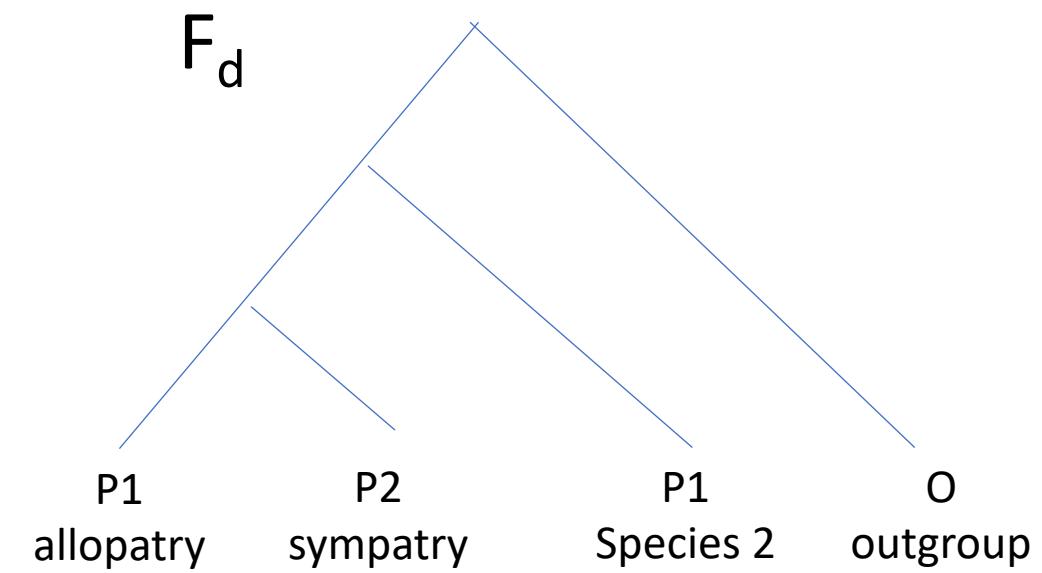
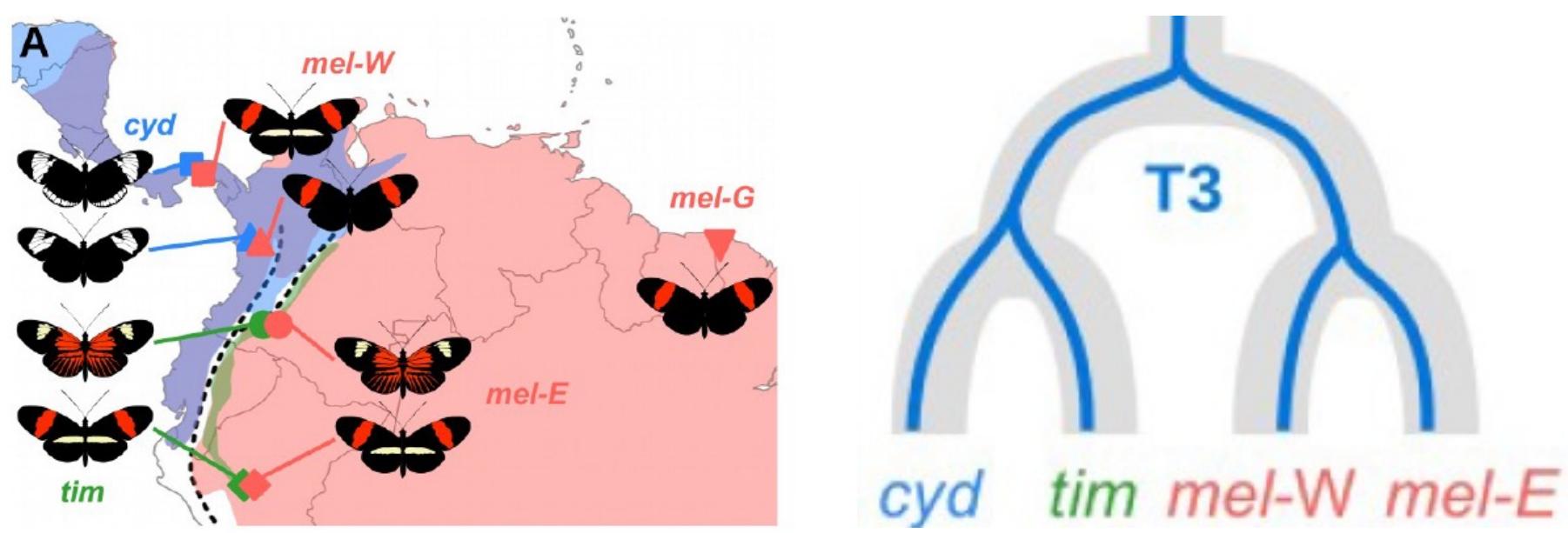


Use different statistics

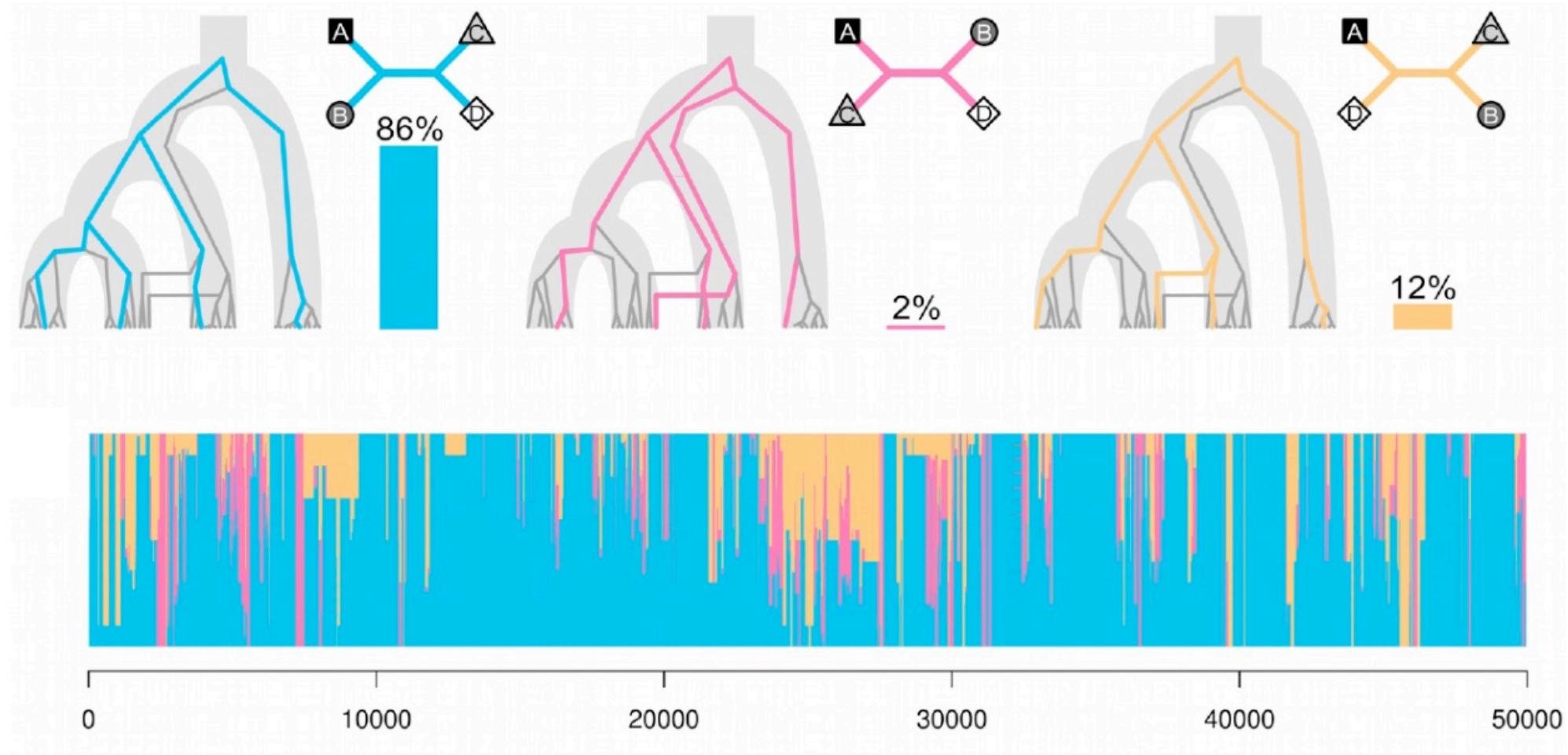
Marques et al., 2018, Nature Ecol Evol







TWISST: Visualizing gene trees across the genome

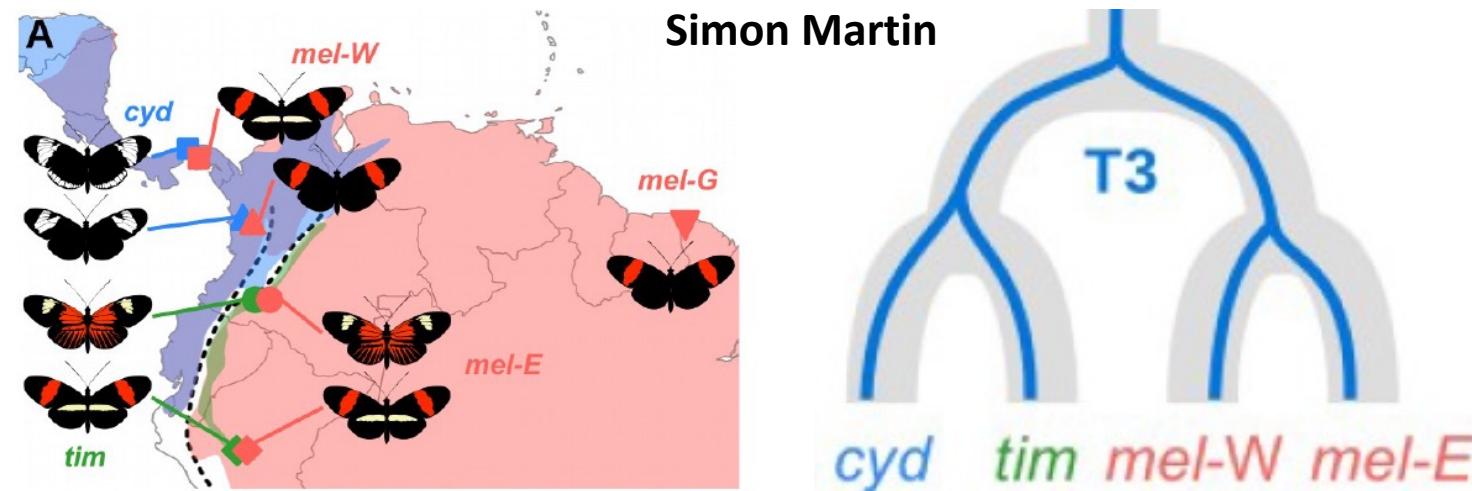


Signatures and statistics to detect candidate barrier loci

- Locally restricted gene flow
 - Reduced f_d
 - Gmin
- Increased differentiation and potentially divergence
 - Increased Fst -> PBS
 - Increased dxy
 - increased $\Delta\pi$
- Selective sweep signals in one or both populations
 - Increased Haplotype Length, e.g. iHS and XP-EHH
 - Reduced π
 - Negative Tajima's D

Detecting regions under divergent selection and barriers to gene flow

- If high gene flow -> F_{ST} is a good measure for detecting regions under divergent selection or barrier loci
- If the taxa are divergent enough -> d_{xy} may work best, particularly if levels of gene flow are not very high and in cases of secondary contact Ideally correct for differences in pi with an outgroup.
- If the taxa are very young and gene flow is not very high, f_d or TWISST might help if allopatric and parapatric populations were sequenced



Detecting regions under divergent selection and barriers to gene flow

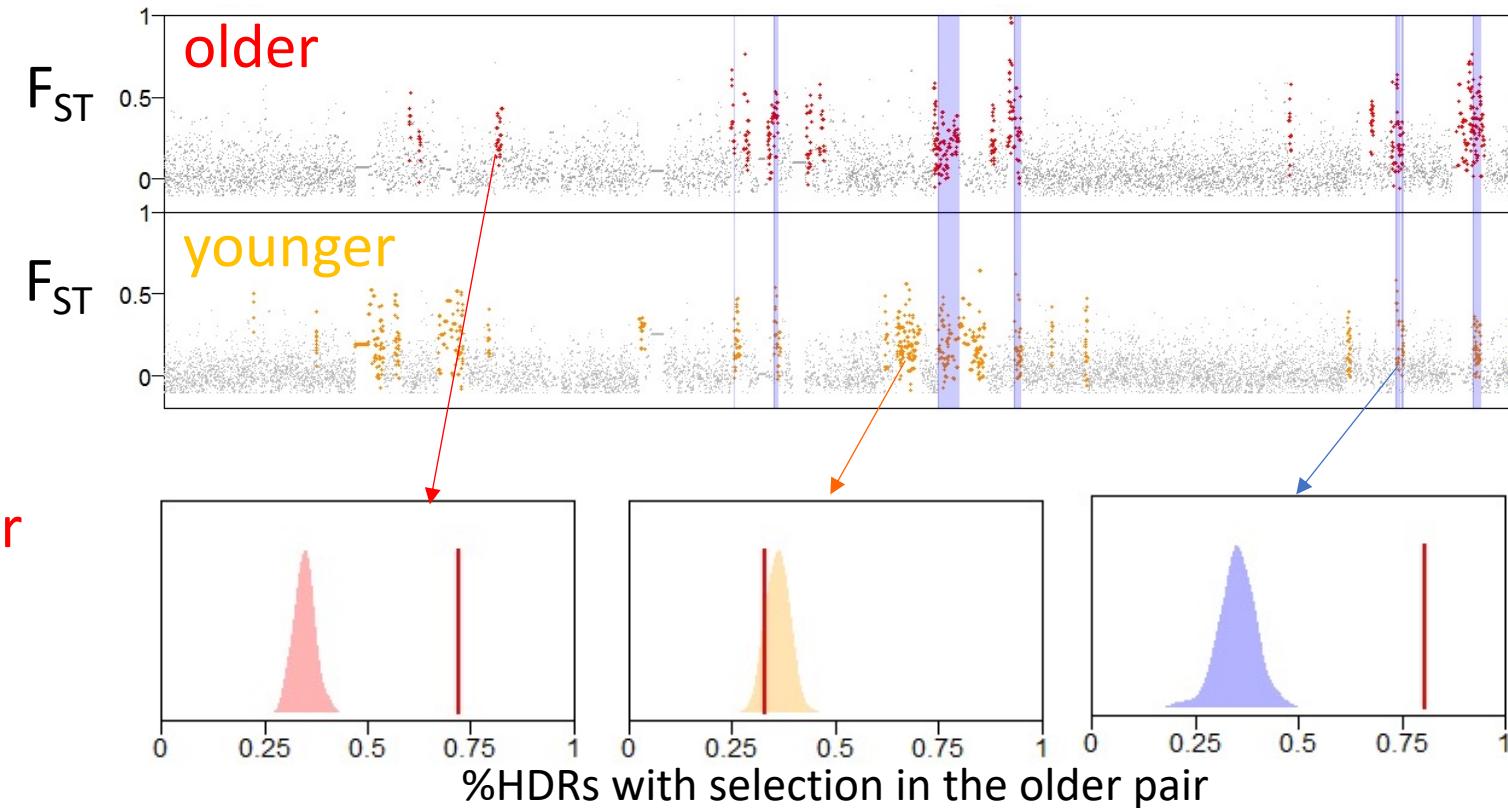
- If rates of gene flow between the two taxa compared is high, F_{ST} is a good measure for detecting regions under divergent selection or barrier loci
- If the taxa are divergent enough, d_{xy} may work best, particularly if levels of gene flow are not very high and in cases of secondary contact. Ideally correct for differences in π with an outgroup.
- If the taxa are very young and gene flow is not very high, f_d or TWISST might help if allopatric and parapatric populations were sequenced
- If there is no gene flow, it is better to search for signatures of selective sweeps (e.g. iHS, XP-EHH, Tajima's D). However, inferring if these regions are involved in speciation is difficult.

Enrichment of selection statistics support the action of selection

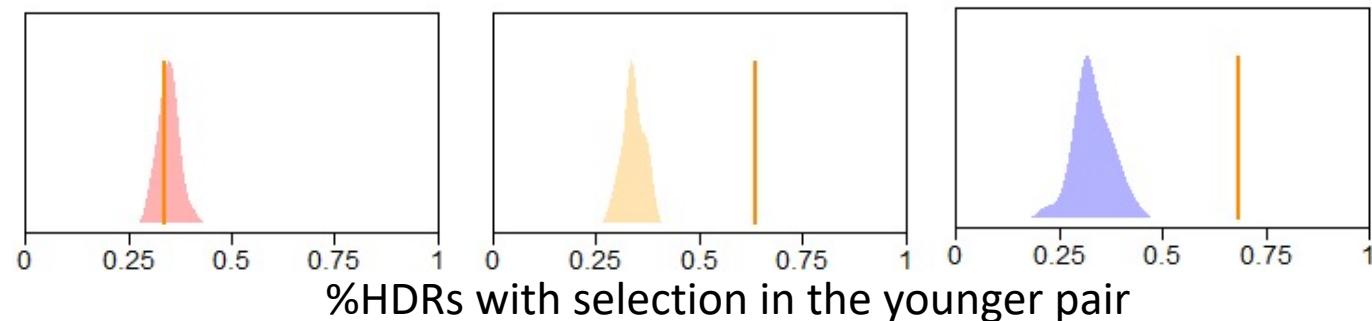
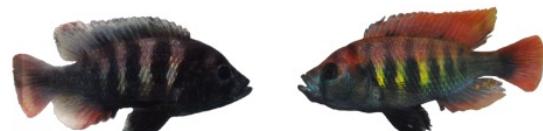
Selection statistics:

d_{xy}
Tajima's D
 $\Delta\pi$
XP-EHH
iHS

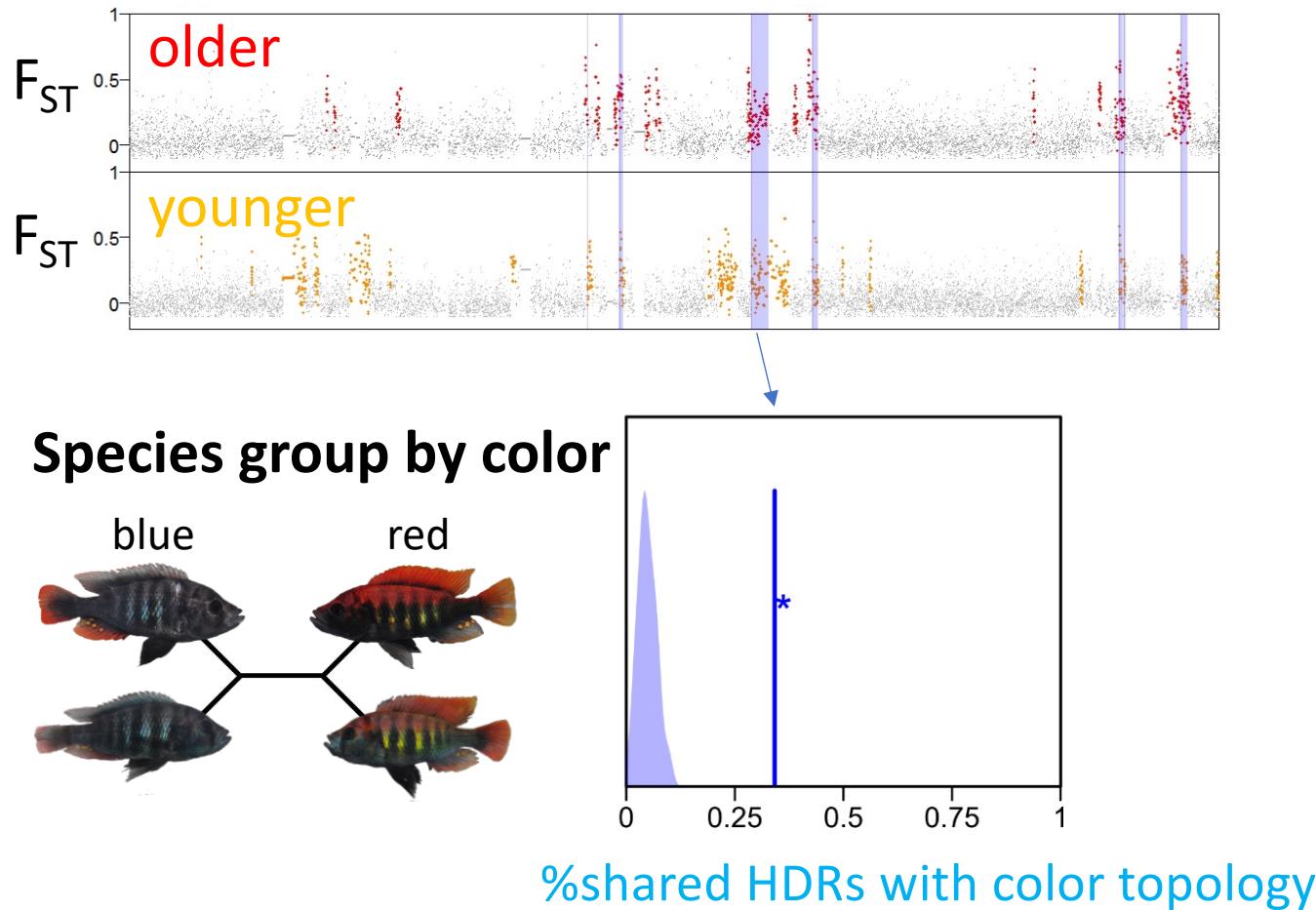
Selection in the older pair



Selection in the younger pair



Example: Highly differentiated regions shared by both species pairs show parallel allele frequency differences

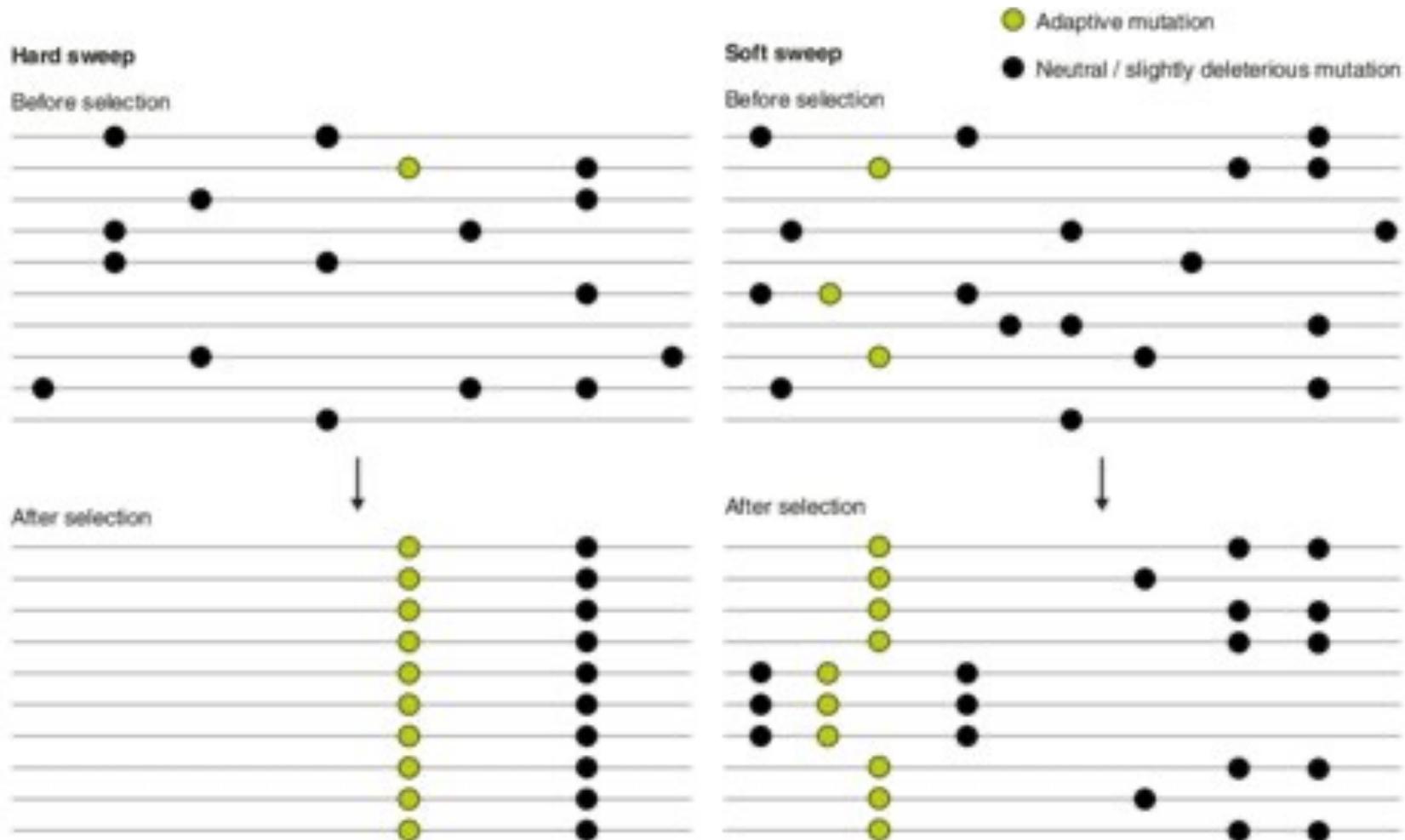


TWISST

(Martin & Van Belleghem, 2017)

Meier *et al.*, 2018, MBE

Hard vs soft sweep



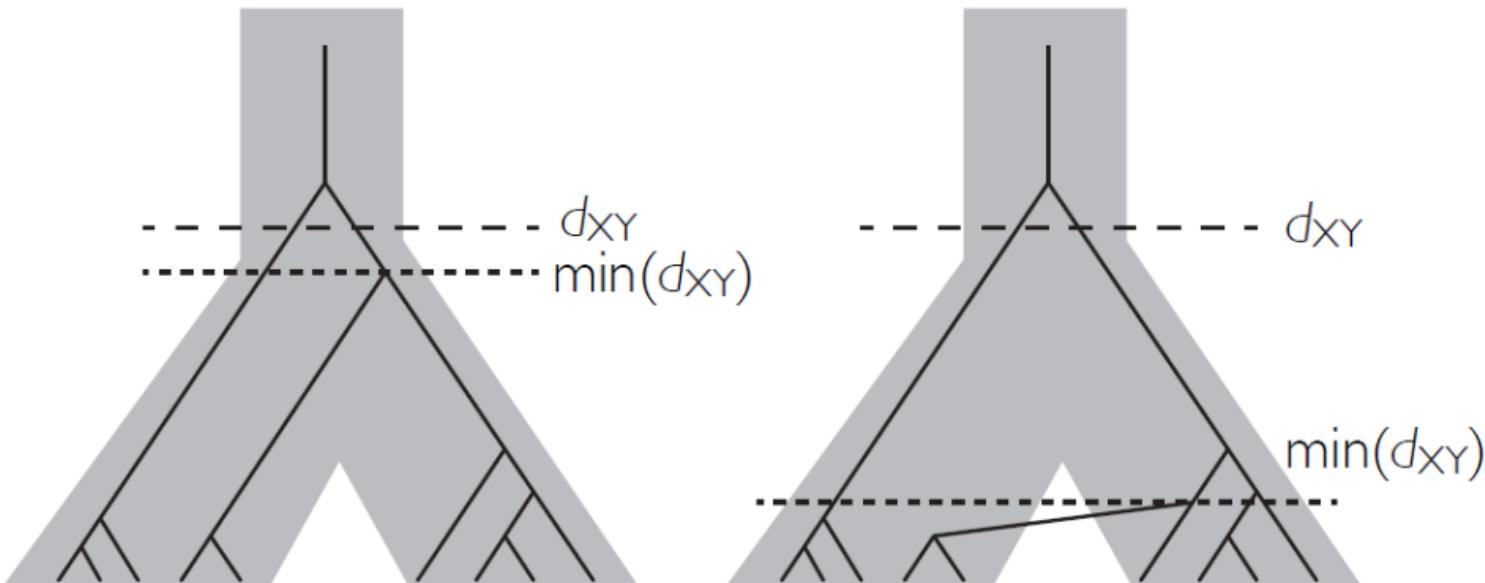
Absolute measures of divergence

$$d_{XY} = \sum_{ij} x_i y_j d_{ij}$$

Average number of pairwise differences between two populations

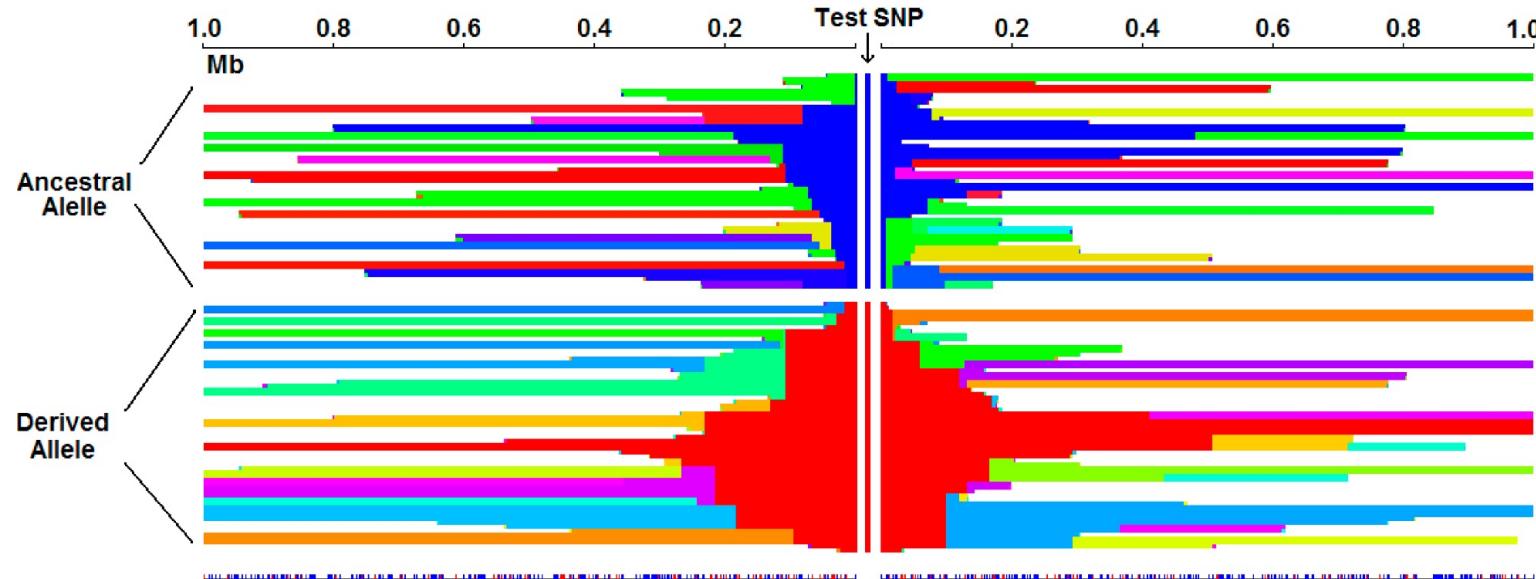
Pop A	Pop B
ACTGTC	ATTAGC
ATTGTC	ACTGGC
ACTGTC	ACTAGC
ATTGTC	ATTAGC

Here d_{XY} is
0.375



iHS and XP-EHH

A



iHS: within a population

iHS (integrated haplotype score) compares haplotype lengths **within a population**

-> an allele under selection will lead to increased haplotype length relative to other haplotypes in the same region

-> useful to detect **ongoing/incomplete sweeps**

XP-EHH: between populations

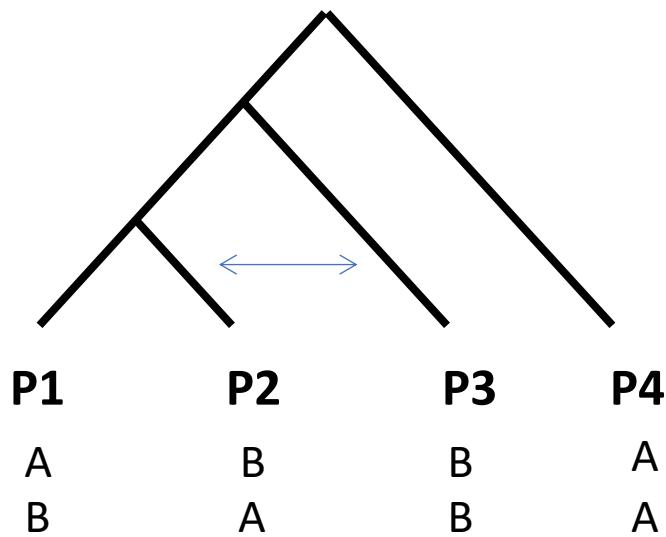
XP-EHH (cross population extended haplotype homozygosity) compares haplotype lengths **between populations**

-> a population that had a sweep has increased haplotype lengths relative to the haplotypes in the other population in the same region

-> most powerful with **complete sweeps** restricted to one population

Sliding window introgression: f_d

f_d can be applied to smaller number of ABBA and BABA sites than D and is thus ideal for sliding windows. ABBA and BABA patterns are computed from allele frequencies and the f test of the four populations is standardized by the maximum value it could get which would be the scenario of complete mixing between P2 and P3. P2 and P3 are thus both set to PD which is the taxon with higher derived allele frequency of P2 and P3.



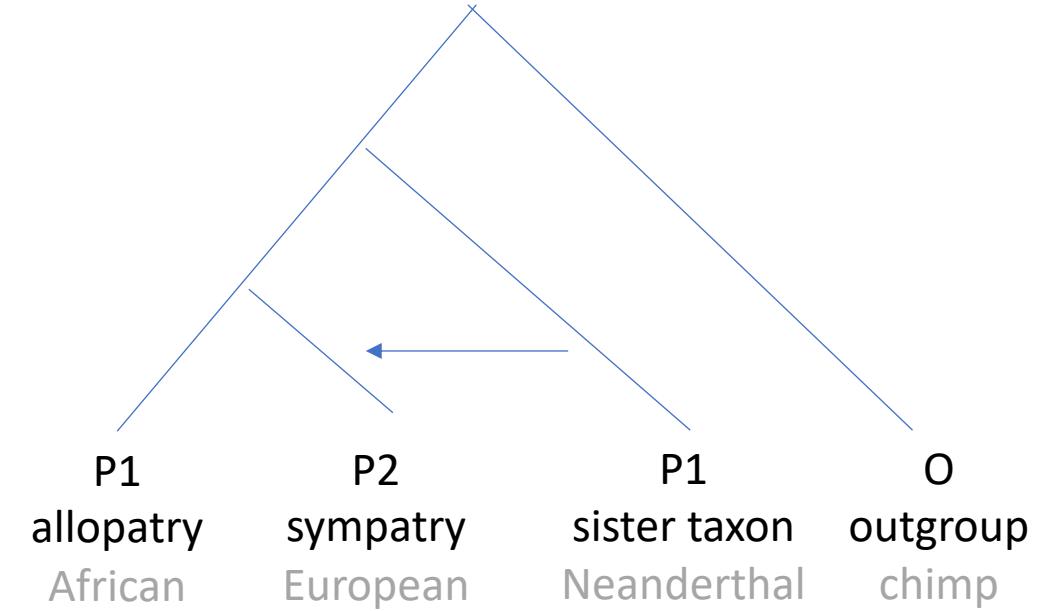
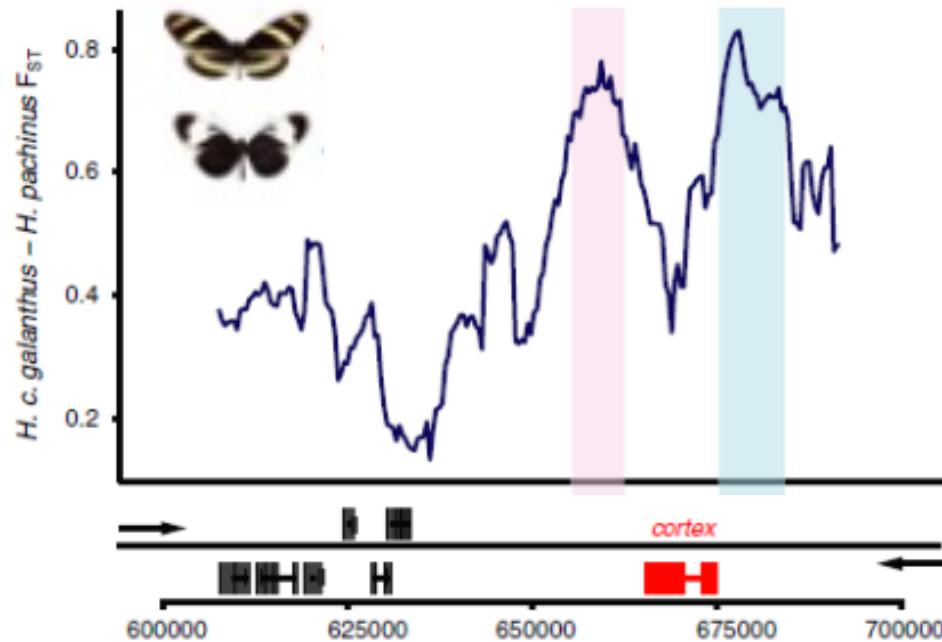
$$C_{\text{ABBA}}(i) = (1 - \hat{p}_{i1})\hat{p}_{i2}\hat{p}_{i3}(1 - \hat{p}_{i4})$$

$$C_{\text{BABA}}(i) = \hat{p}_{i1}(1 - \hat{p}_{i2})\hat{p}_{i3}(1 - \hat{p}_{i4})$$

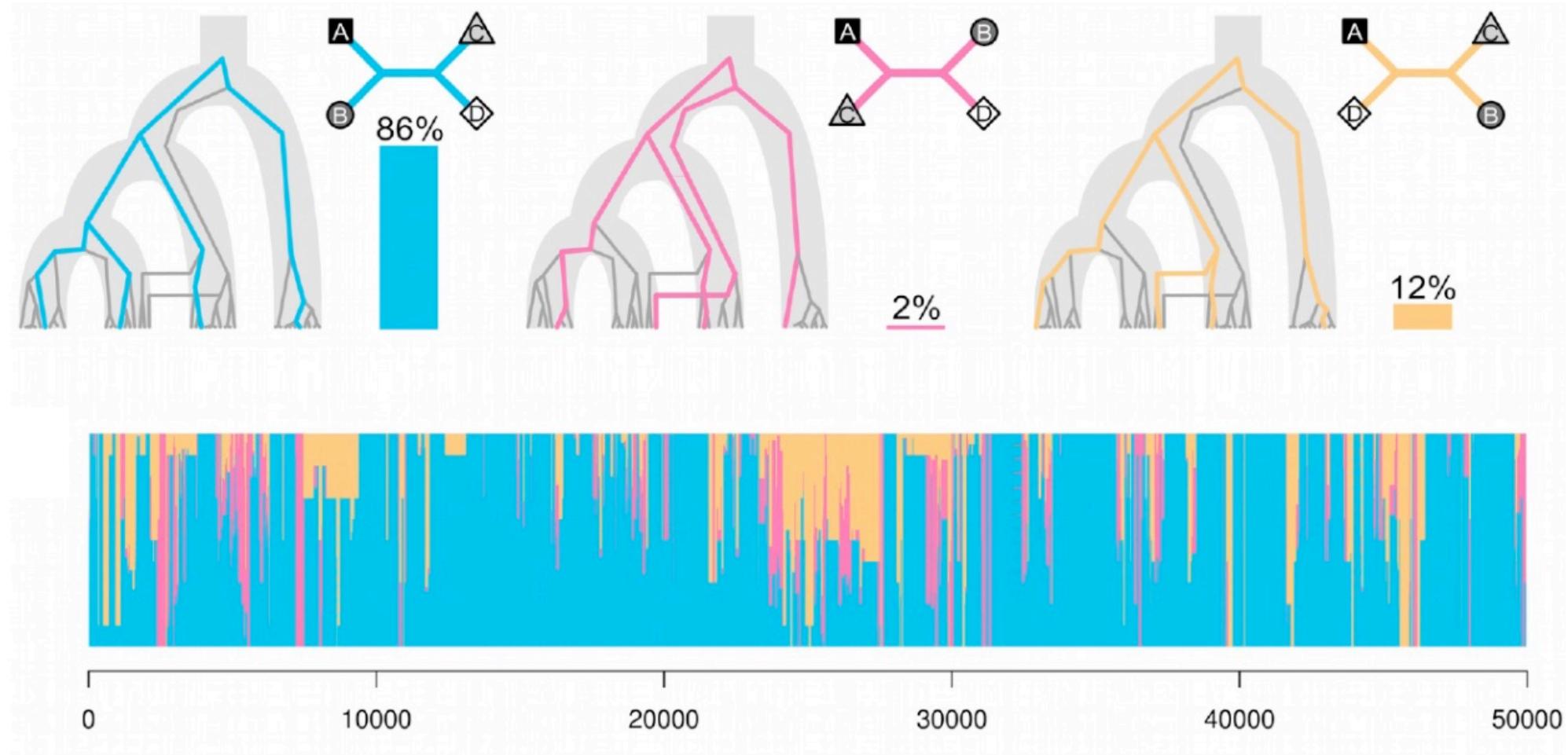
$$\hat{f}_d = \frac{S(P_1, P_2, P_3, O)}{S(P_1, P_D, P_D, O)}$$

PD=P2 or P3
(taxon with higher
derived allele frequency)

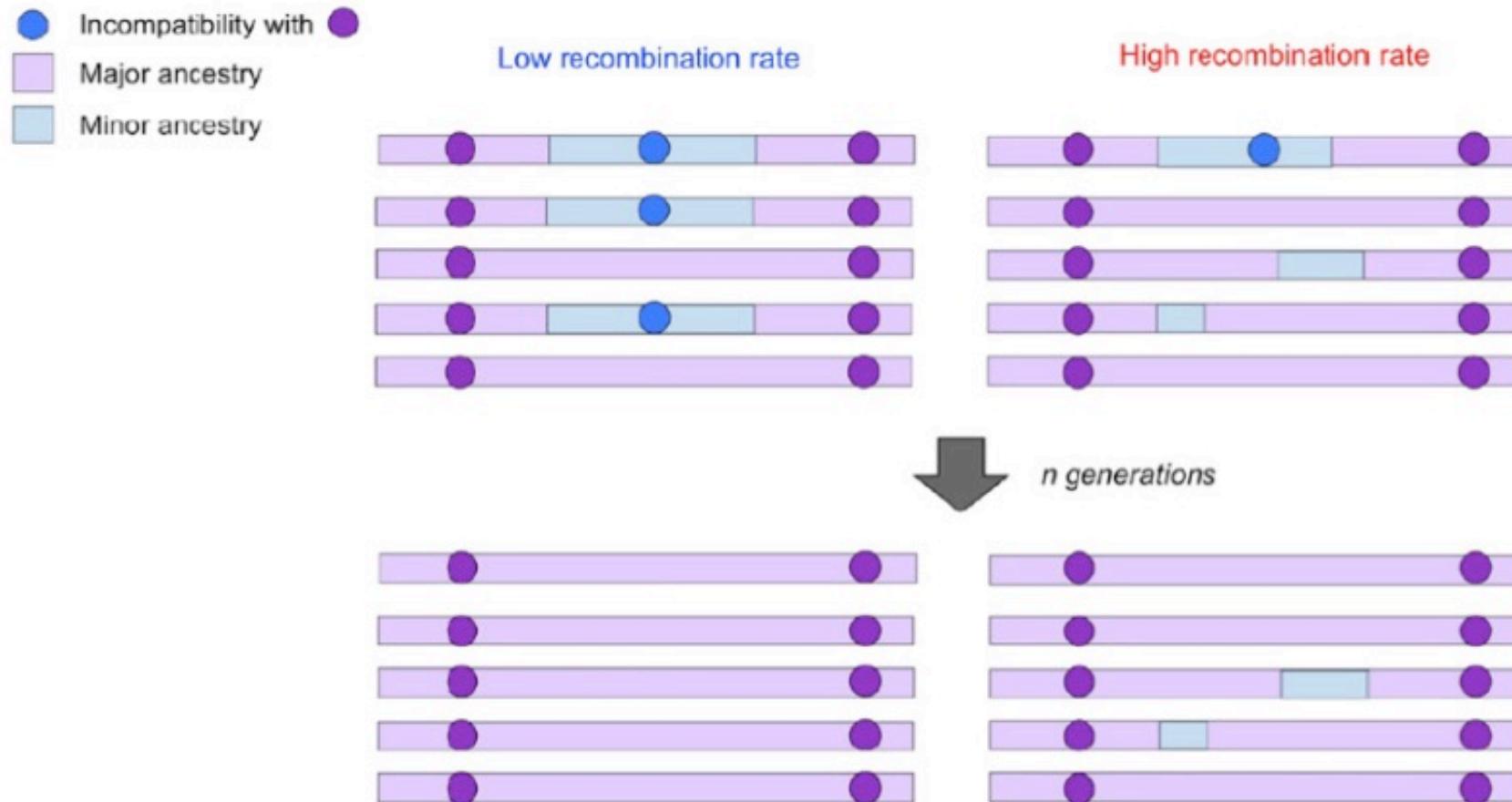
f_d can be used to find regions of reduced gene flow
if allopatric and sympatric populations exist
or alternatively, of adaptive introgression



TWISST: Visualizing gene trees across the genome

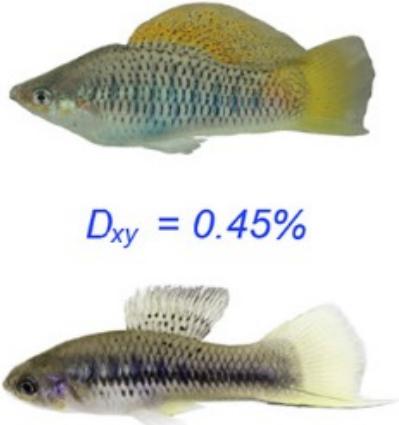


Less introgression in regions of low recombination

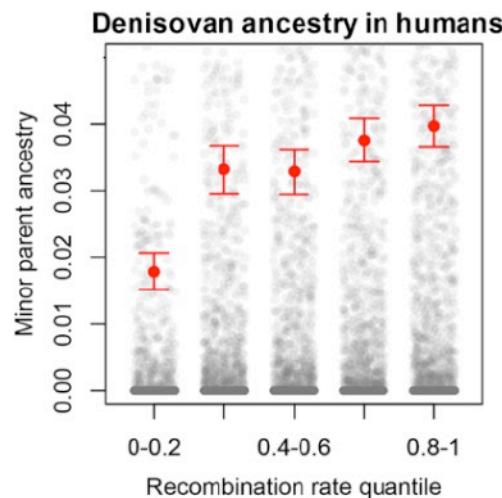


Empirical evidence

Minority parent
of hybrids between X.
malinche and X.
birchmanni



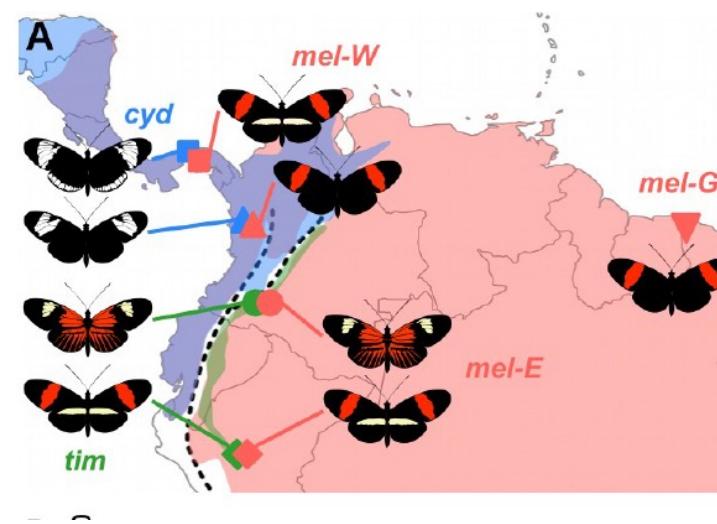
$$D_{xy} = 0.45\%$$



Schumer et al., 2018
Juric et al., 2016

Denisovan or
Neanderthal ancestry
in humans

Introgession from
sympatric species in
Heliconius



Martin et al., accepted Plos Biol

