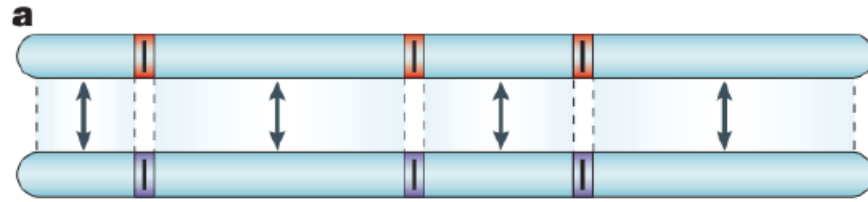
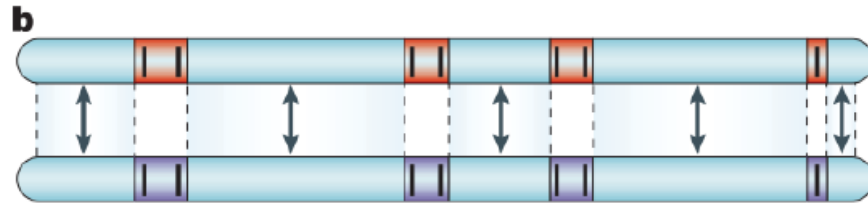


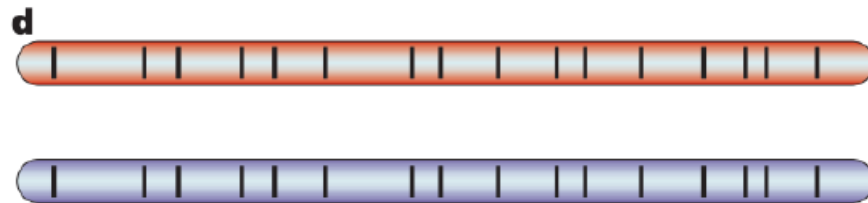
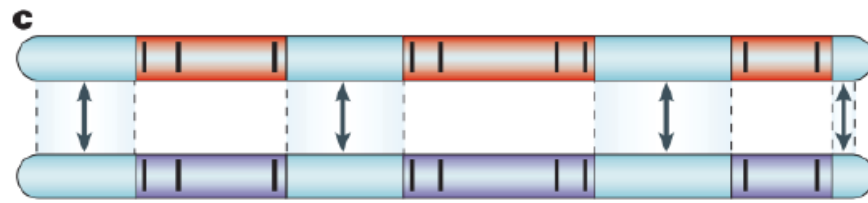
The genic concept of speciation



Divergent loci resist gene flow



Gene flow continues but
linkage builds and divergent
regions grow

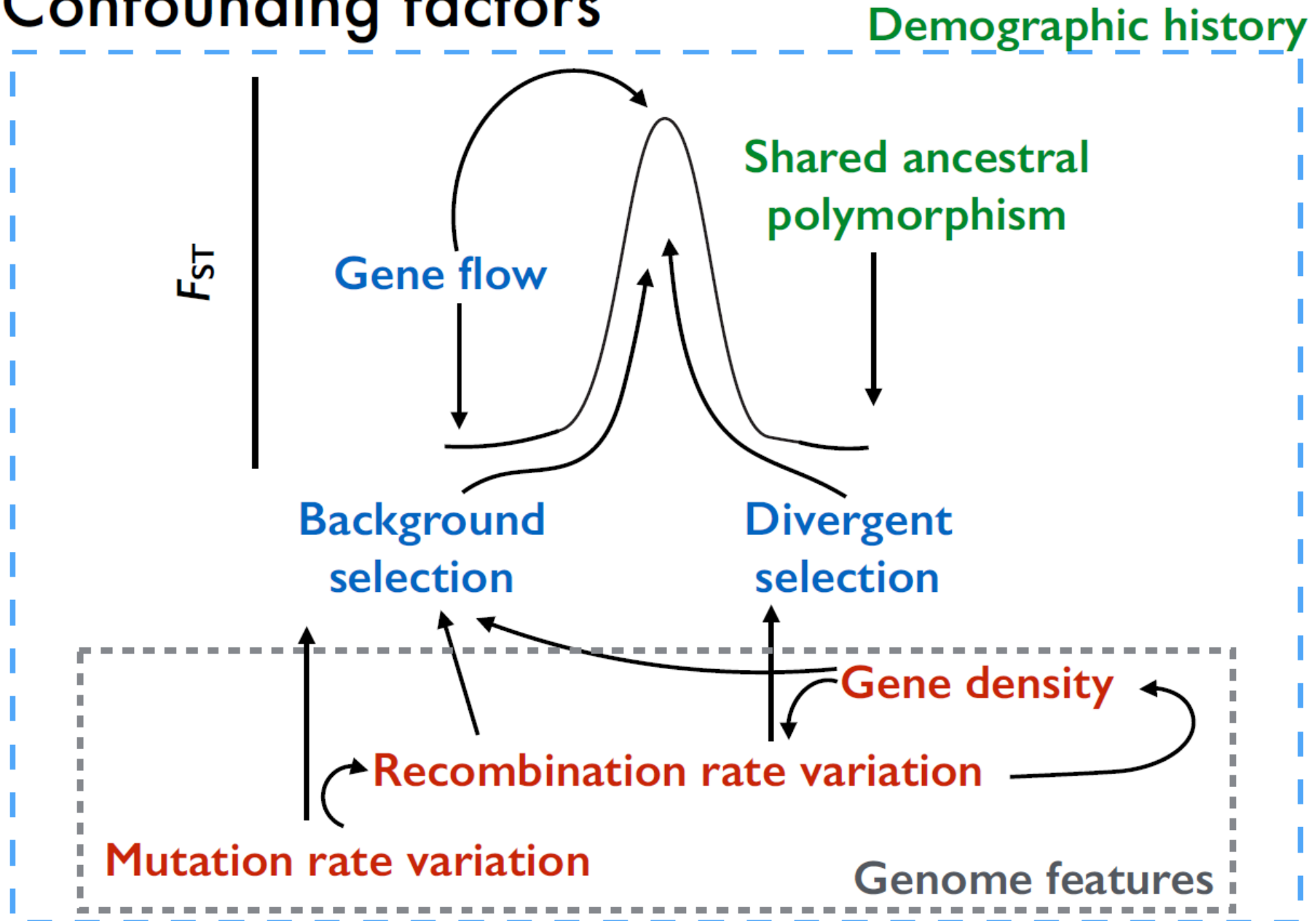


Complete reproductive
isolation evolves

Signatures and statistics to detect candidate barrier loci

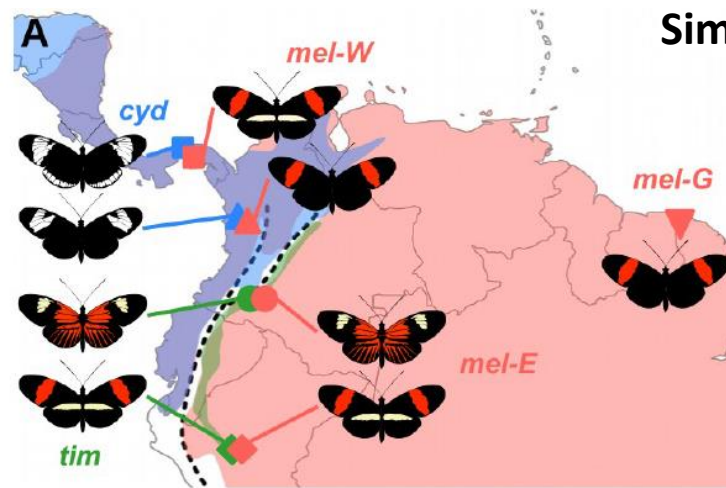
- Locally restricted gene flow
 - Reduced f_d
 - G_{min}
- Increased differentiation and potentially divergence
 - Increased F_{st}
 - Increased d_{xy}
- Selective sweep signals in one or both populations
 - Increased Haplotype Length, e.g. iHS and $XP-EHH$
 - Reduced π
 - Negative Tajima's D

Confounding factors

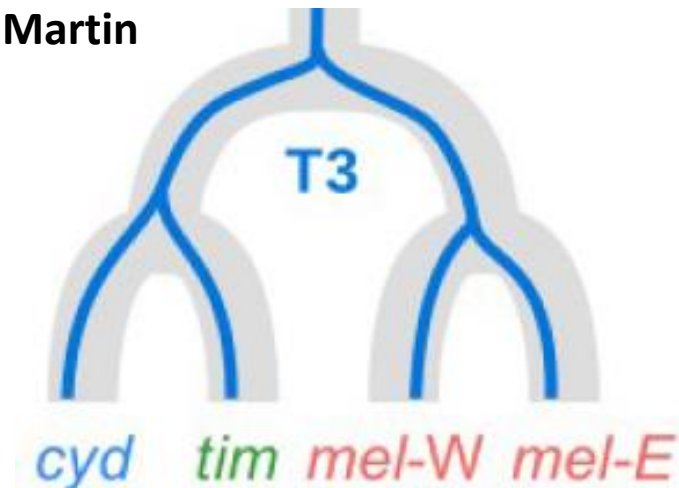


Detecting regions under divergent selection and barriers to gene flow

- 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



Simon Martin



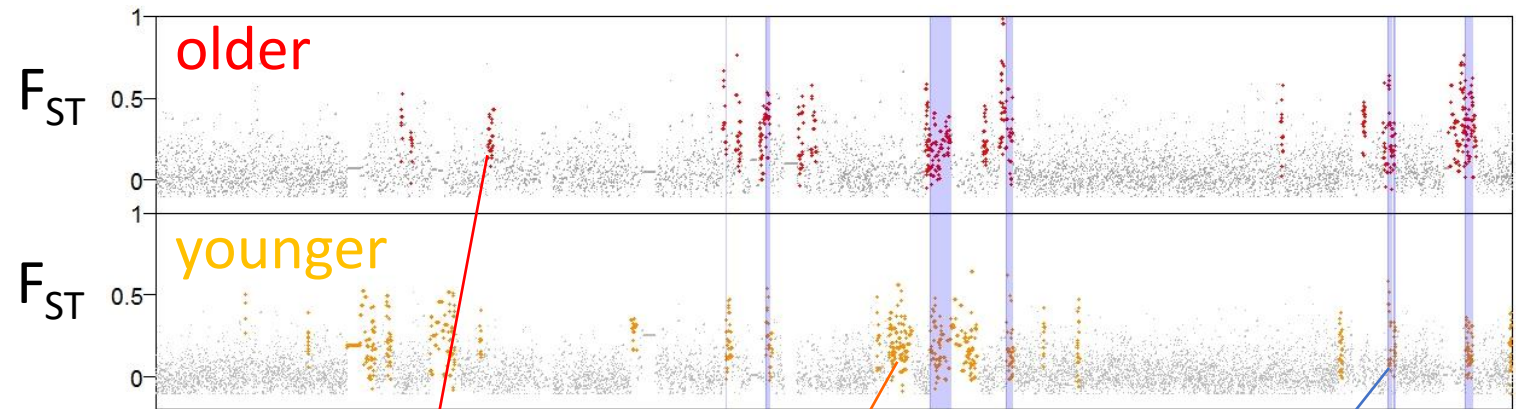
Detecting regions under divergent selection and barriers to gene flow

- 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.

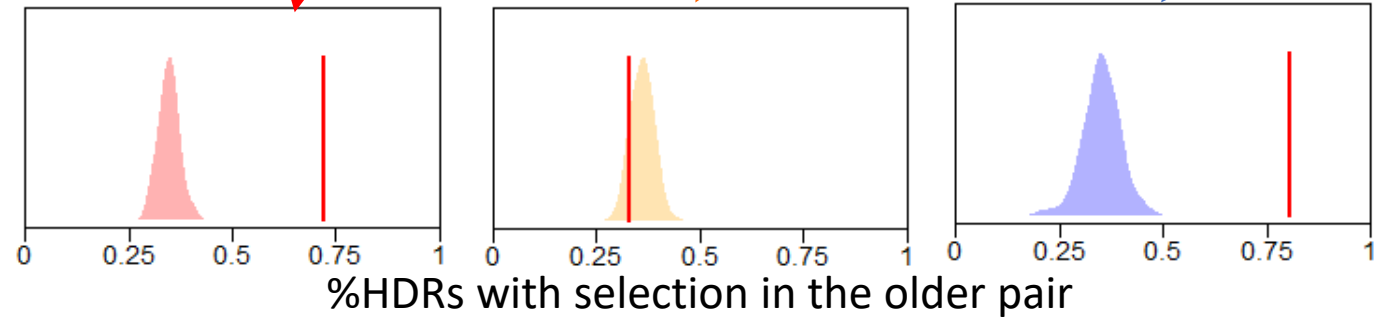
Enrichment of selection statistics support the action of selection

Selection statistics:

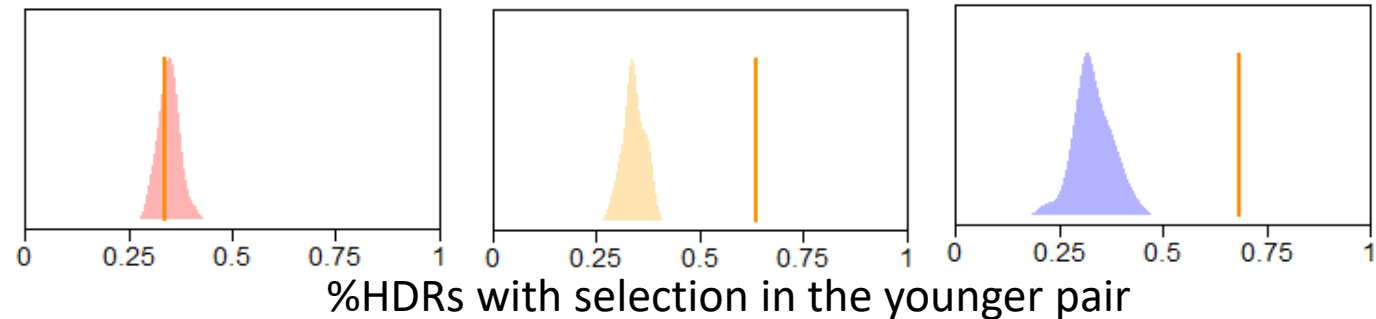
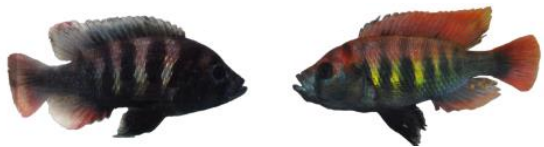
d_{xy}
Tajima's D
 $\Delta\pi$
XP-EHH
iHS



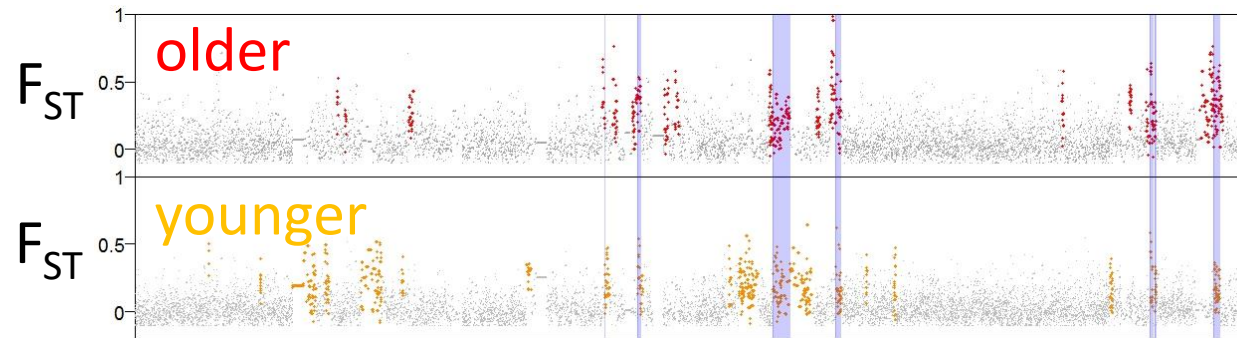
Selection in the older pair



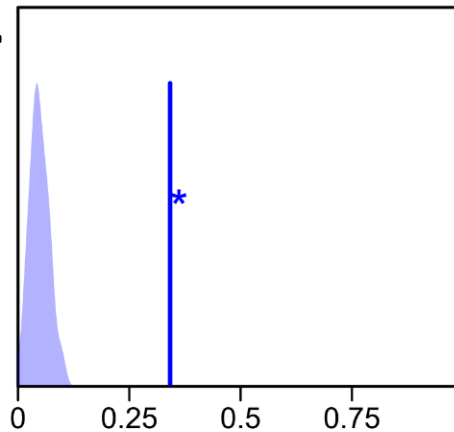
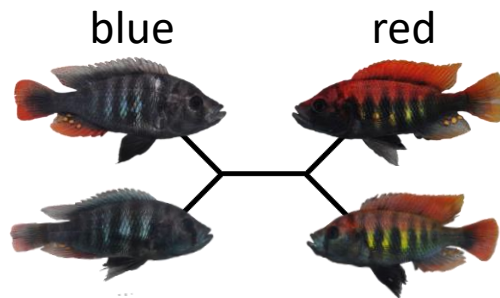
Selection in the younger pair



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



Species group by color



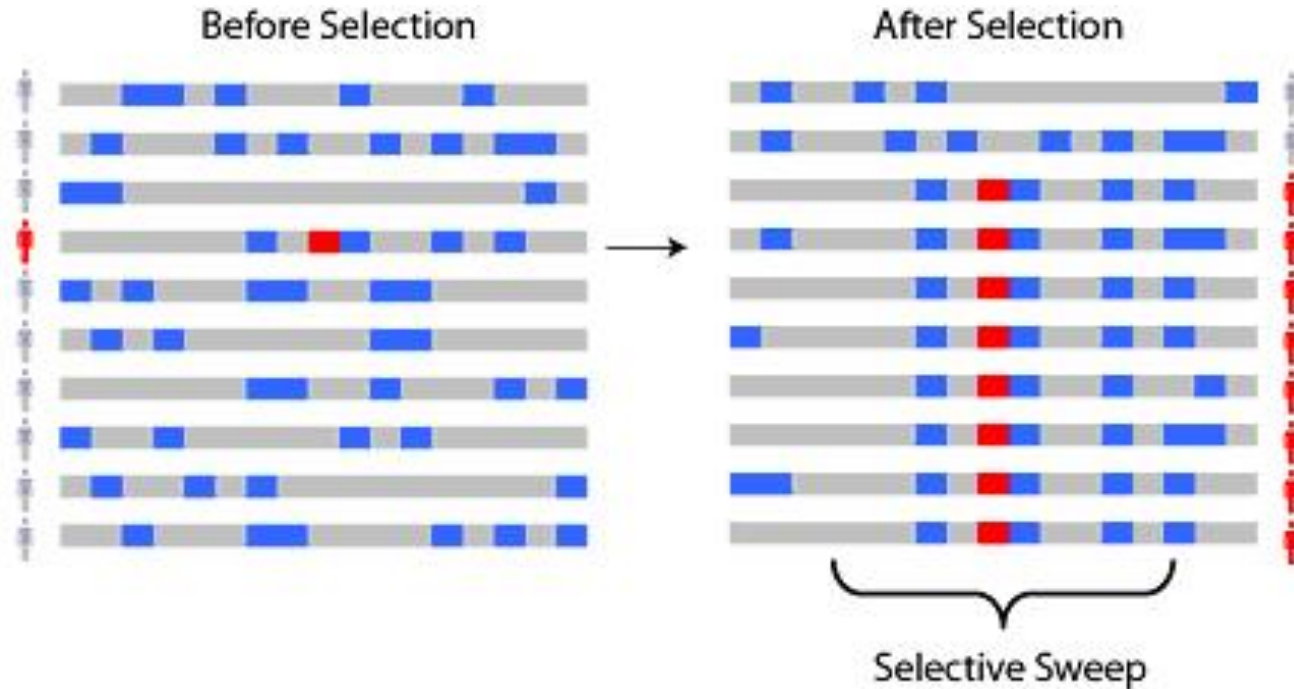
%shared HDRs with color topology

TWISST

(Martin & Van Belleghem, 2017)

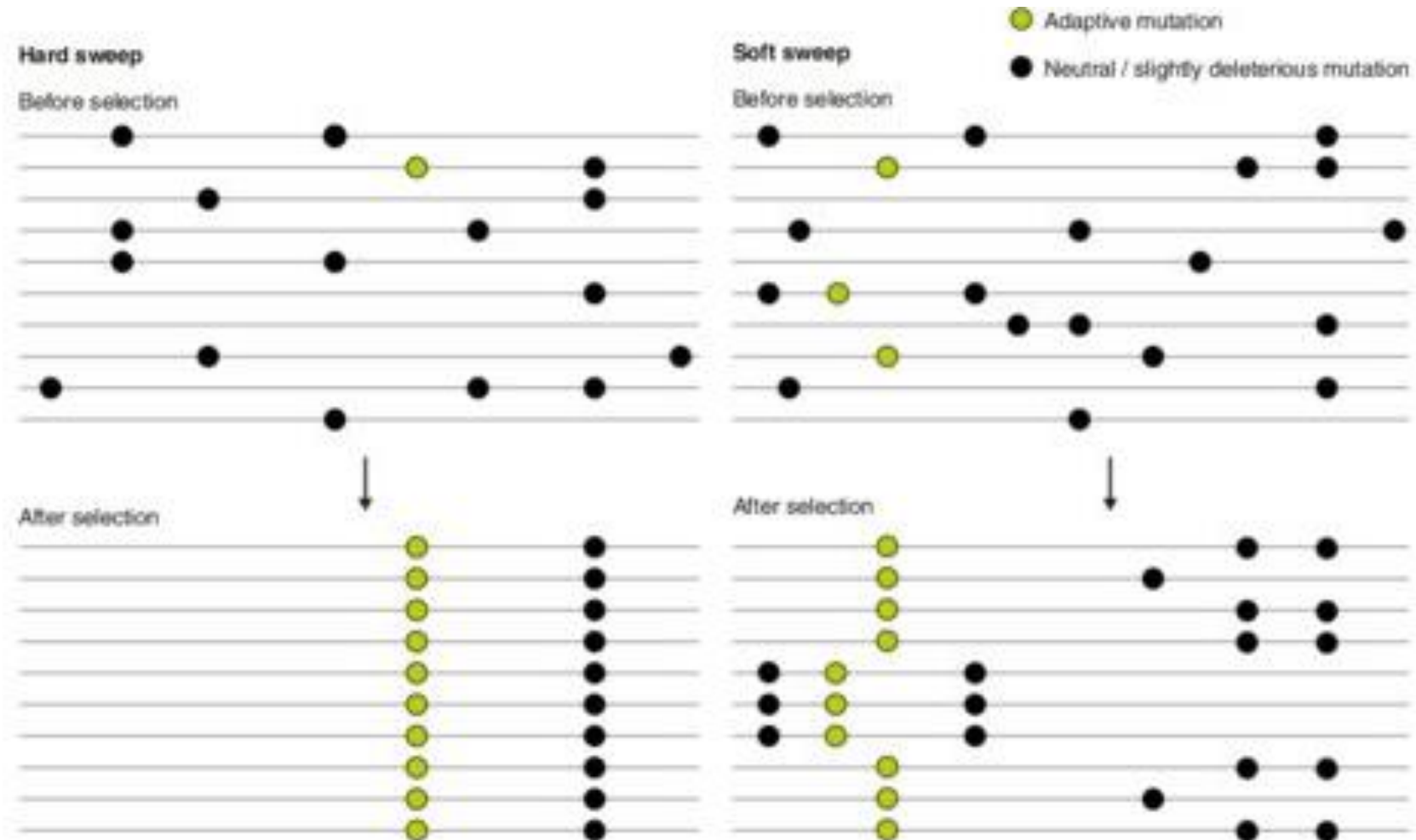
Meier *et al.*, 2018, MBE

Selective sweep signatures

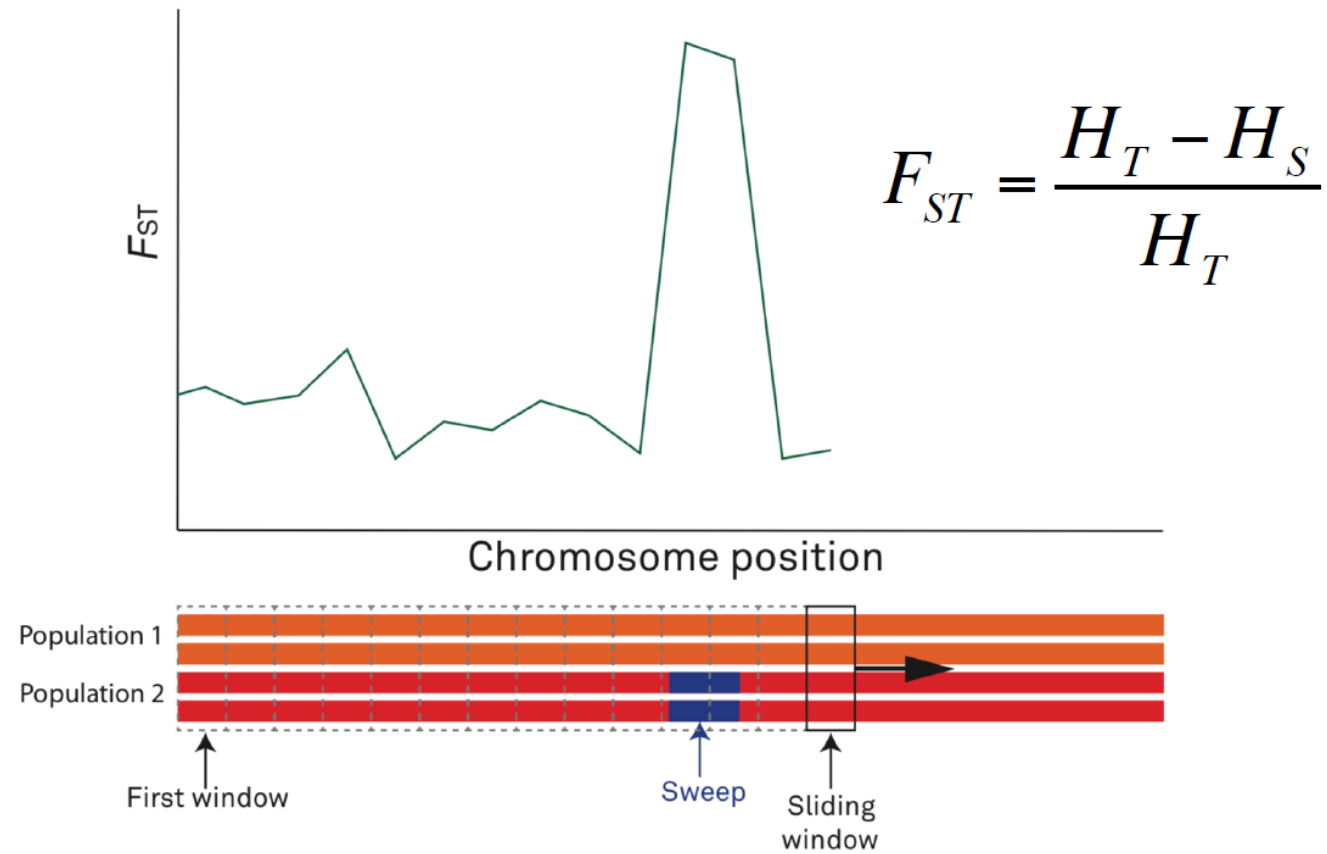


- Increased haplotype length
- Reduced genetic variation
- Increased genetic differentiation to another population

Hard vs soft sweep



Sliding window estimates to detect selection



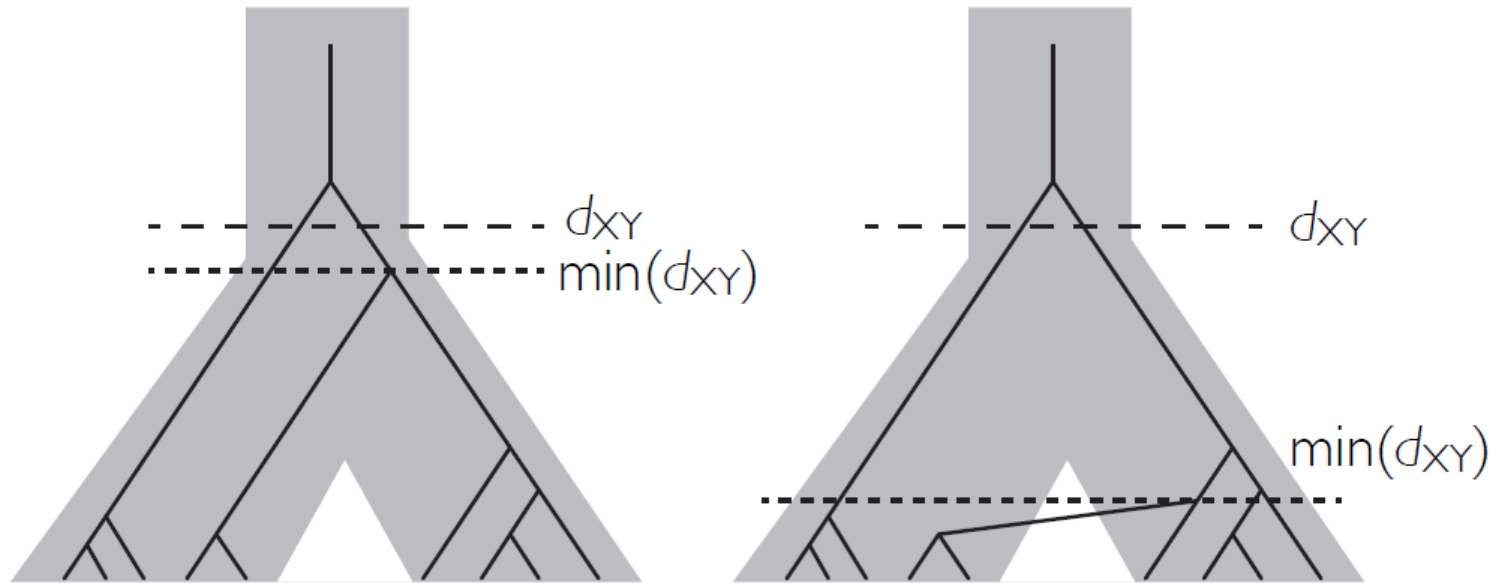
Absolute measures of divergence

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

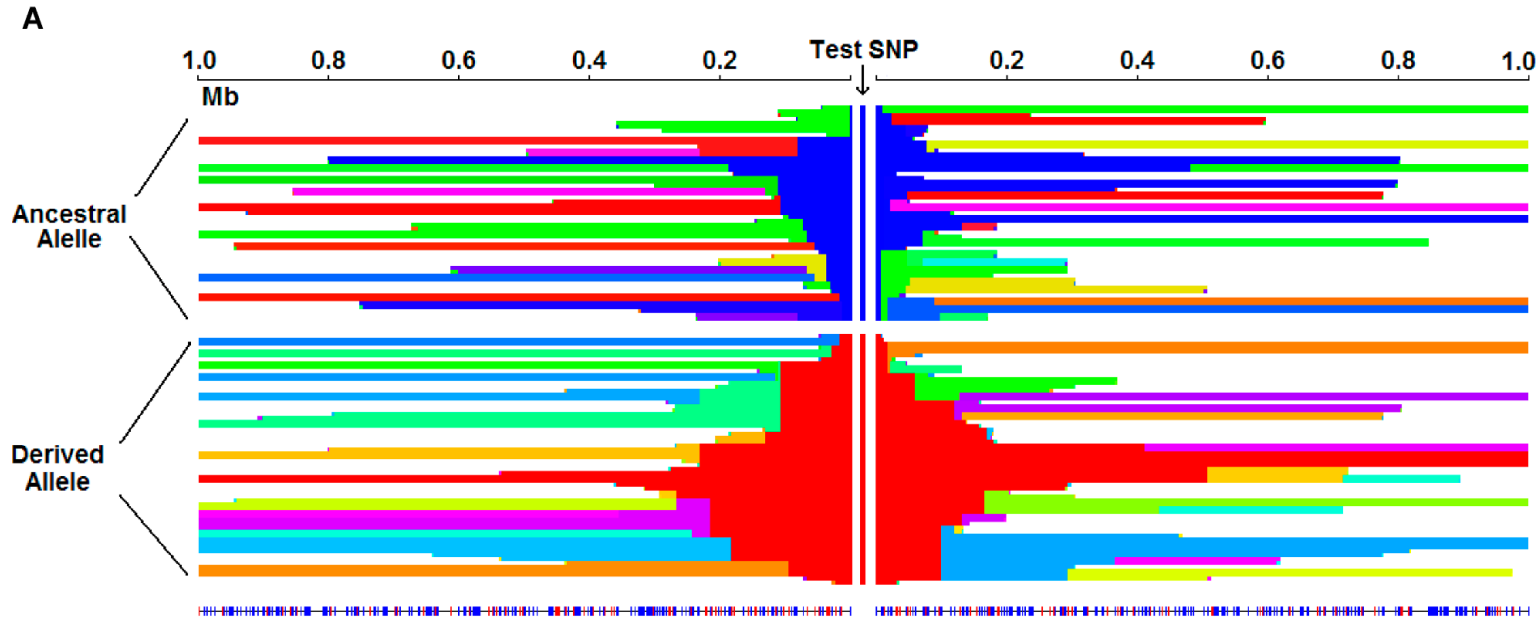
Average number of pairwise differences between two populations

Pop A	Pop B
A C T G T C	A T T A G C
A T T G T C	A C T G G C
A C T G T C	A C T A G C
A T T G T C	A T T A G C

Here d_{XY} is
0.375



iHS and XP-EHH



iHS: within a population

XP-EHH: between populations

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**

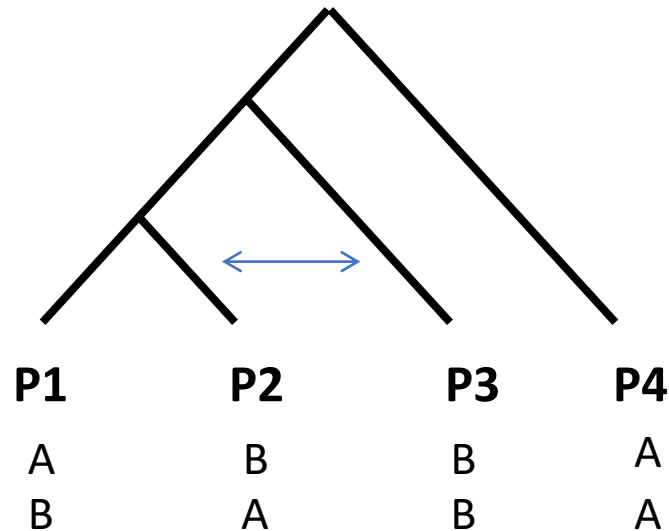
XP-EHH (cross population extended haplotype homozygosity) compares haplotype lengths **between populations**

-> 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.



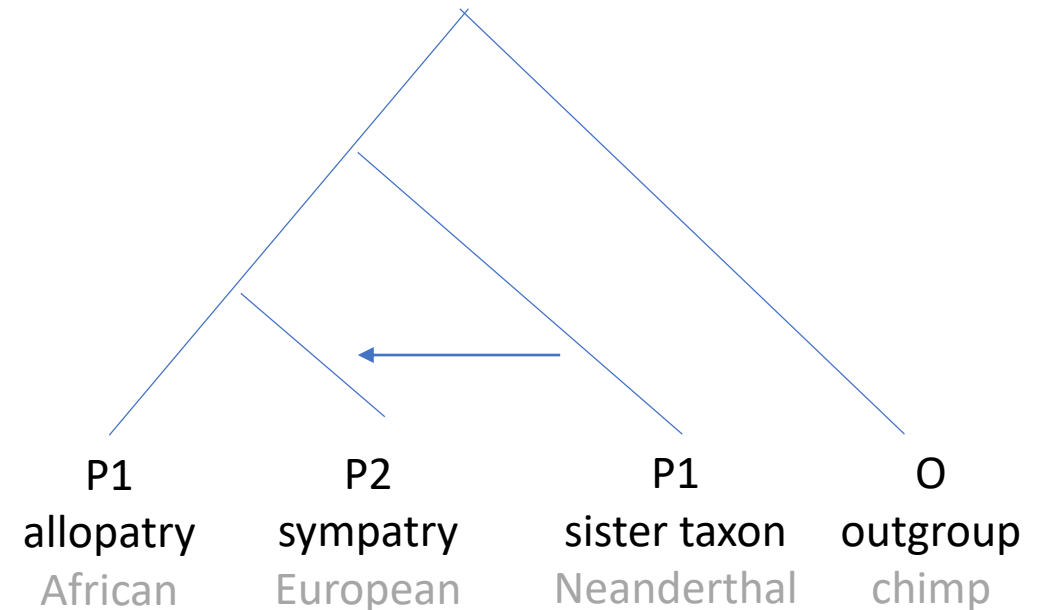
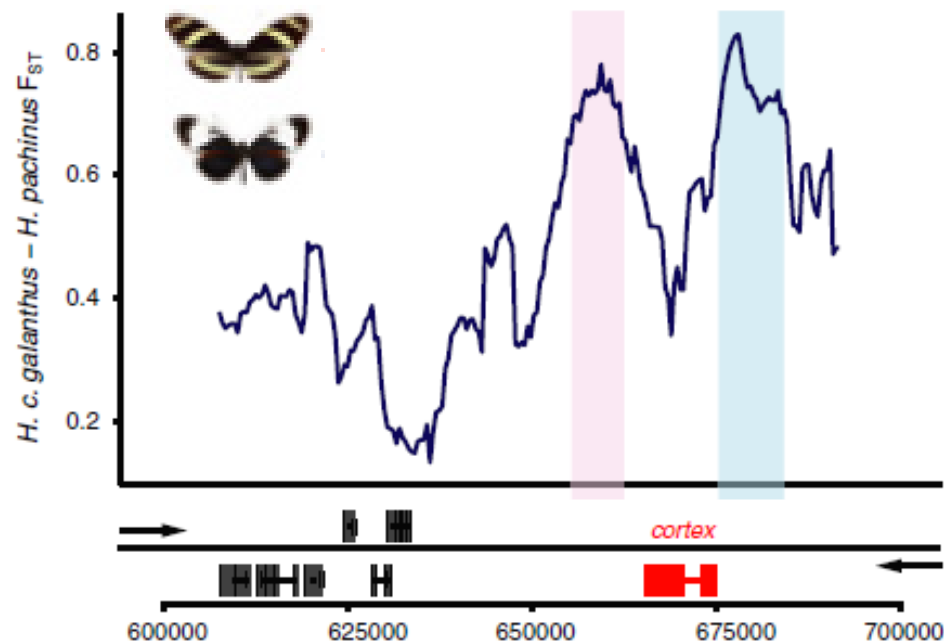
$$C_{ABBA}(i) = (1 - \hat{p}_{i1})\hat{p}_{i2}\hat{p}_{i3}(1 - \hat{p}_{i4})$$

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

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

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

