

# Evolutionary Genomics 2: Population Genomics



Simon Martin

Department of Zoology  
University of Cambridge

# Outline

- Questions in population genomics
- Statistics, inference and examples
  - Genetic diversity and population size
  - Population subdivision and speciation
  - Selection
- Study design
  - Sampling and sequencing

# Questions in population genomics

- Diversity and population structure
  - How much genetic variation exists in a population?
  - How and why is a species subdivided into populations?
- Demographic history and gene flow
  - How has a population size and distribution changed over time?
  - How much migration and gene flow occurs between populations / species?

# Questions in population genomics

- Speciation
  - Which genes/events cause two species to become distinct and remain distinct?
  - What are the relative roles of adaptive and neutral forces in speciation?
  - What is a species?
- Adaptation
  - Where in the genome is natural selection acting?
  - What is the source of beneficial genetic variation?

# Diversity and allele frequencies

- Heterozygosity: Proportion of the genome that is heterozygous

TAGATCGTCCAGATCGAAGTAGCCCCCTTTCGCTGATCTCGTGCCTAAGTAGATCATGATACT  
TAGGTCGTCCAGATCGATCTAGCCCCCTTTCGCTGAGCTCGTGCTTAAGTAGATTATGATAAT

- Inbreeding
- Selection
- Population size
- Humans ~0.001, Fruit Flies ~0.01

# Measuring genetic diversity in a population

- $S$ : number of variable sites
- $\Pi$  (nucleotide diversity): Average number of differences between any pair of sequences

TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGA**T**CTCGTGC**C**TAAGTAGAT**T**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTG**C****T**TAAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAGCTCGTG**C****T**TAAGTAGAT**C**ATGATAAT  
TAG**G**TCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTG**C**TAAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTG**C**TAAGTAGAT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCCCTTTTCGCTGAGCTCGTG**C****T**TAAGTAGAT**T**ATGATAAT

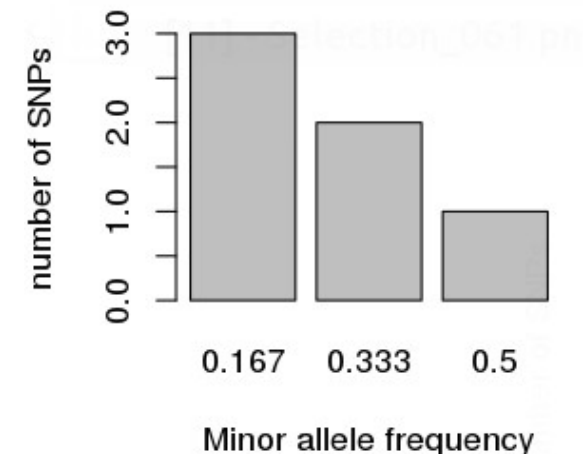
- Selection
- Population size
- If no inbreeding  $\Pi \sim$  heterozygosity

# Measuring genetic diversity in a population

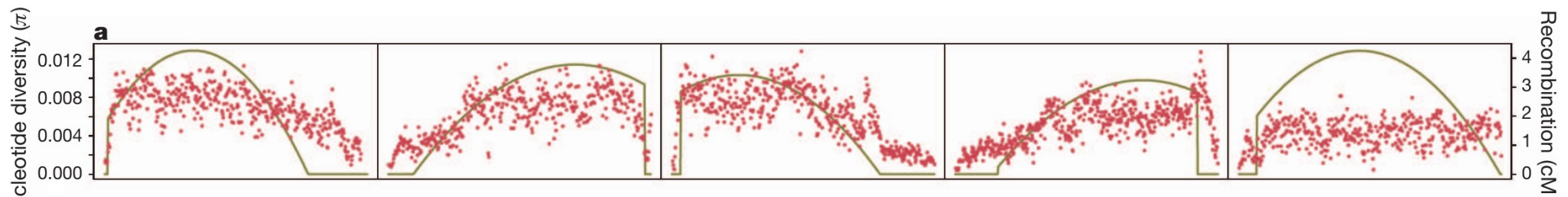
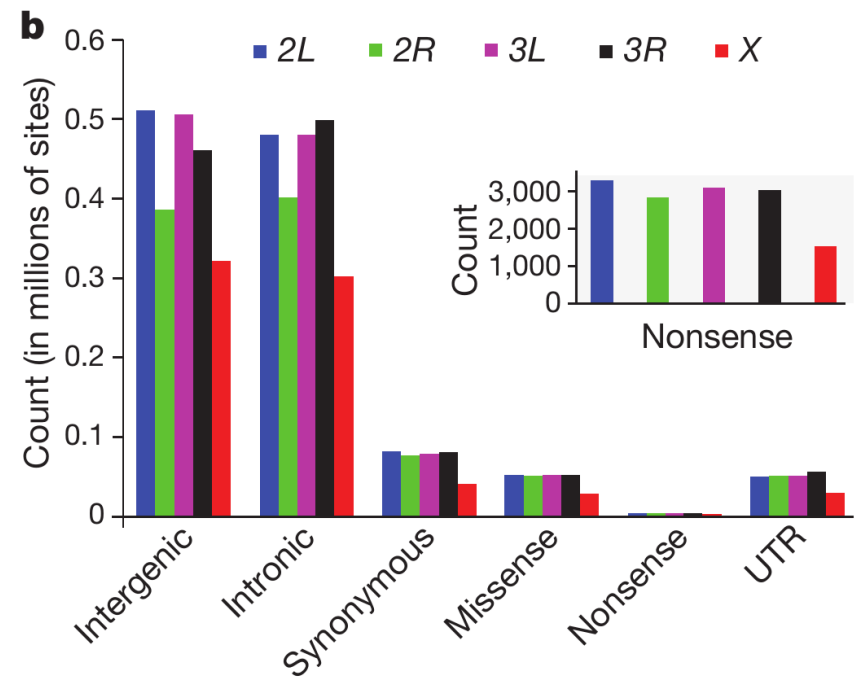
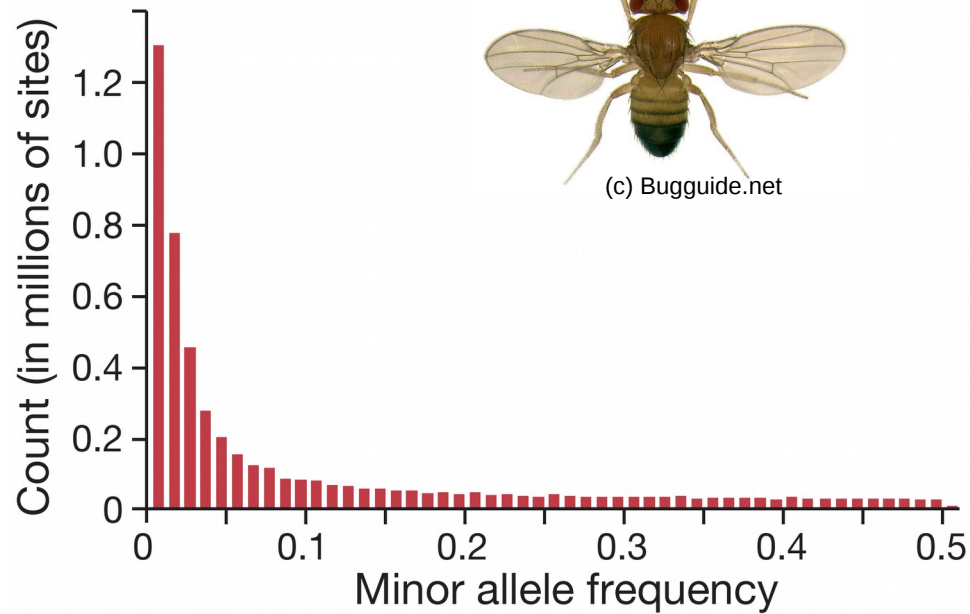
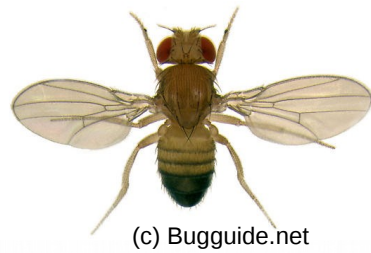
- Site frequency spectrum (SFS) number of occurrences of variants at each frequency

TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAT**T**CTCGTGCCTAAGTAGAT**T**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGC**T**TAAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAG**G**CTCGTGC**T**TAAGTAGAT**C**ATGATAAT  
TAG**G**TCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGC**C**TAAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGC**C**TAAGTAGAT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGC**T**TAAGTAGAT**T**ATGATAAT

- Selection
- Changes in population size
- Hybridisation

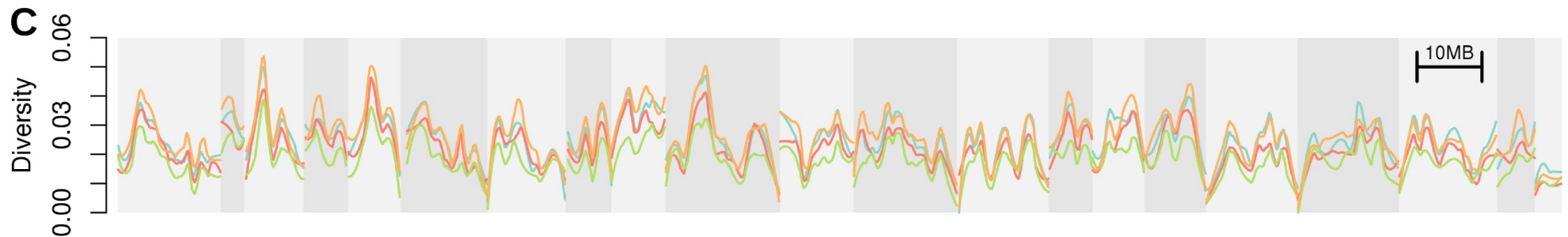
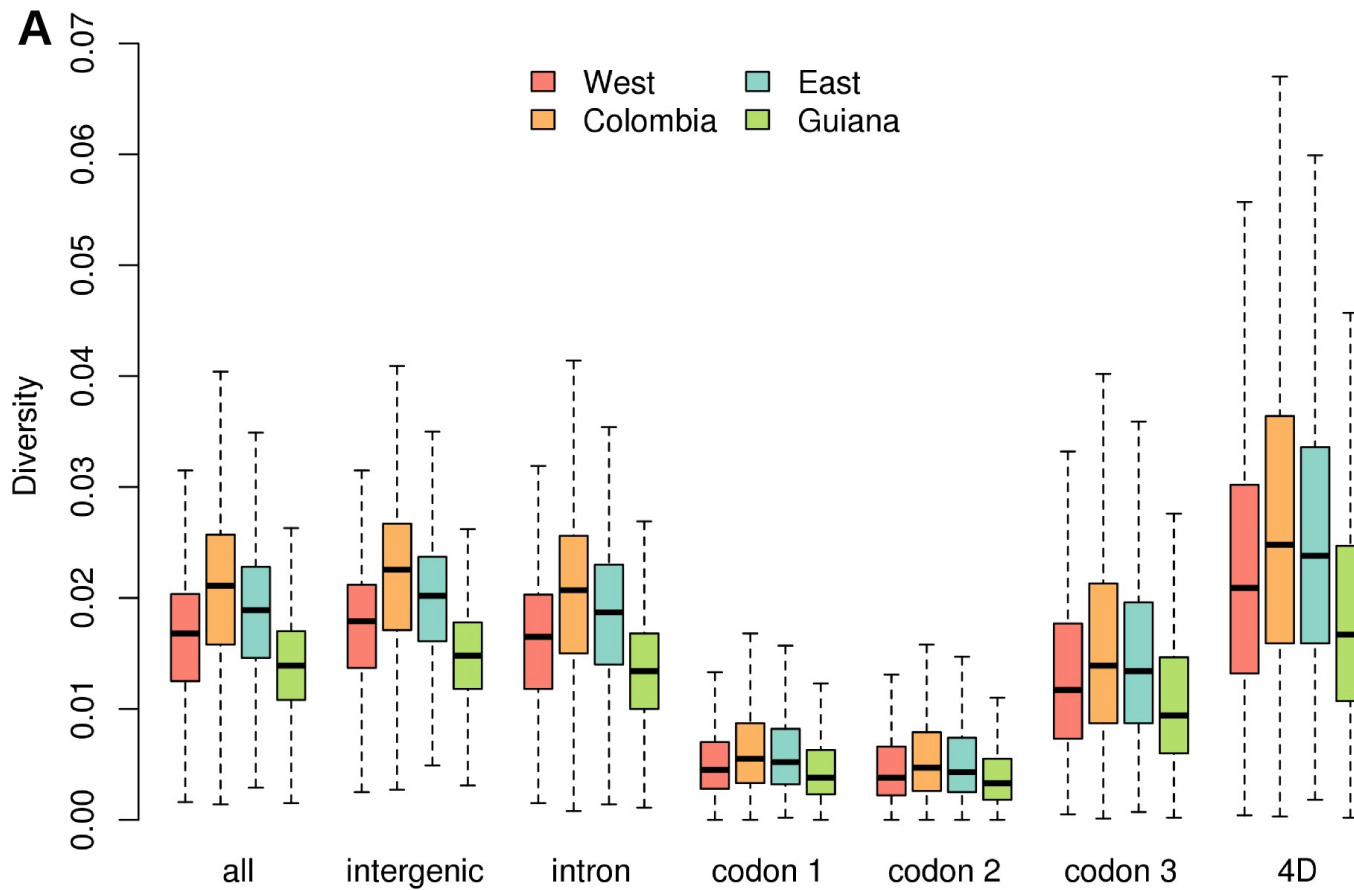


# Measuring genetic diversity in a population



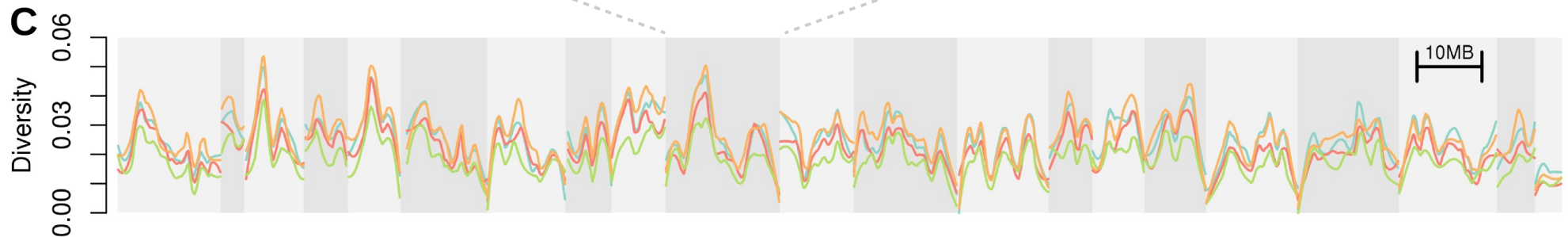
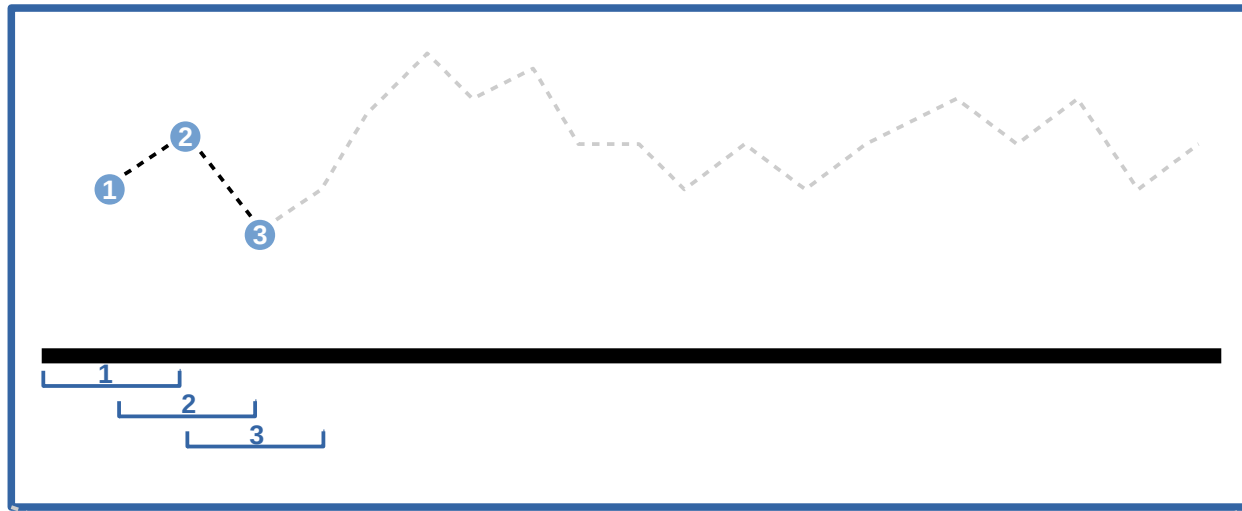


# Measuring genetic diversity in a population



# Measuring genetic diversity in a population

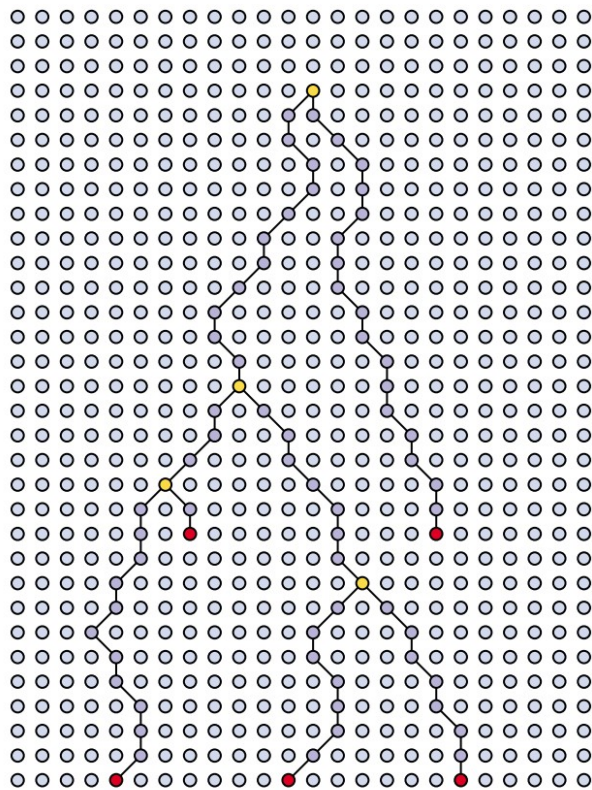
Sliding windows



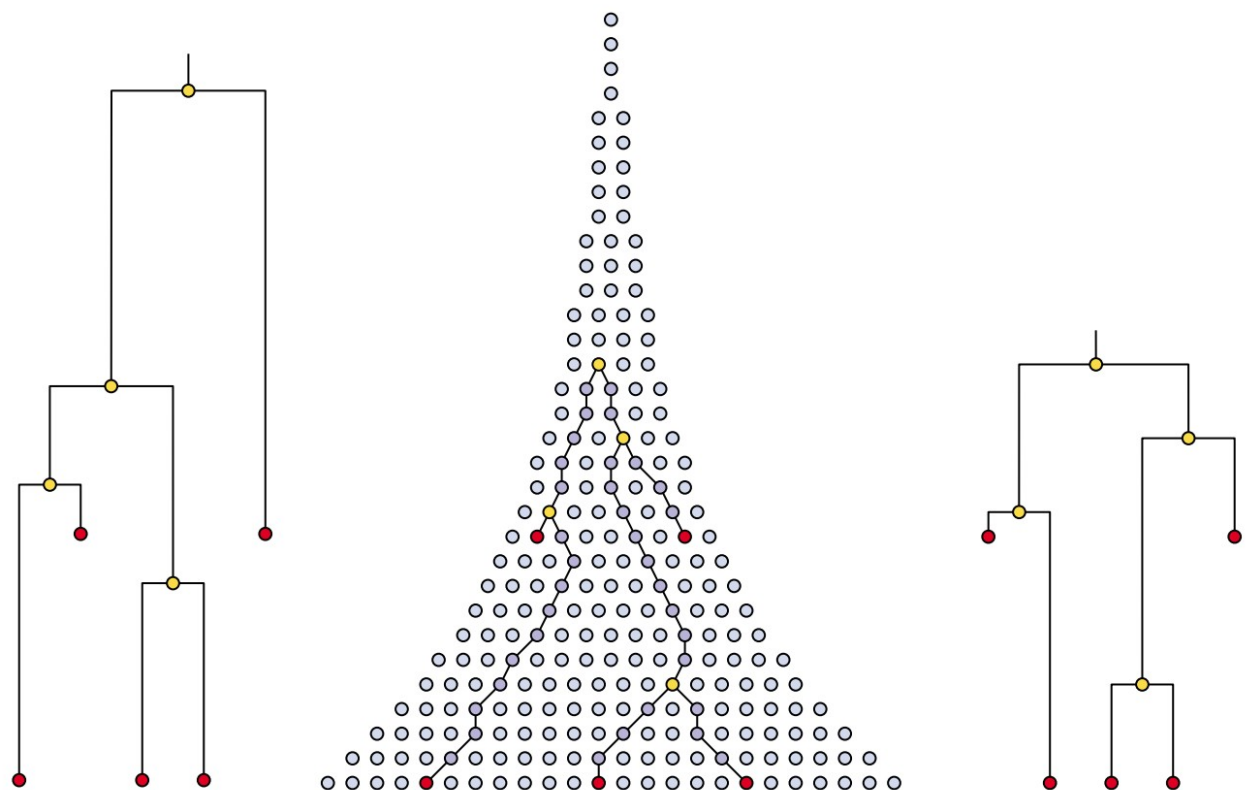
# Inferring population size from diversity data

**Coalescence:** relatedness between individuals relates to population size

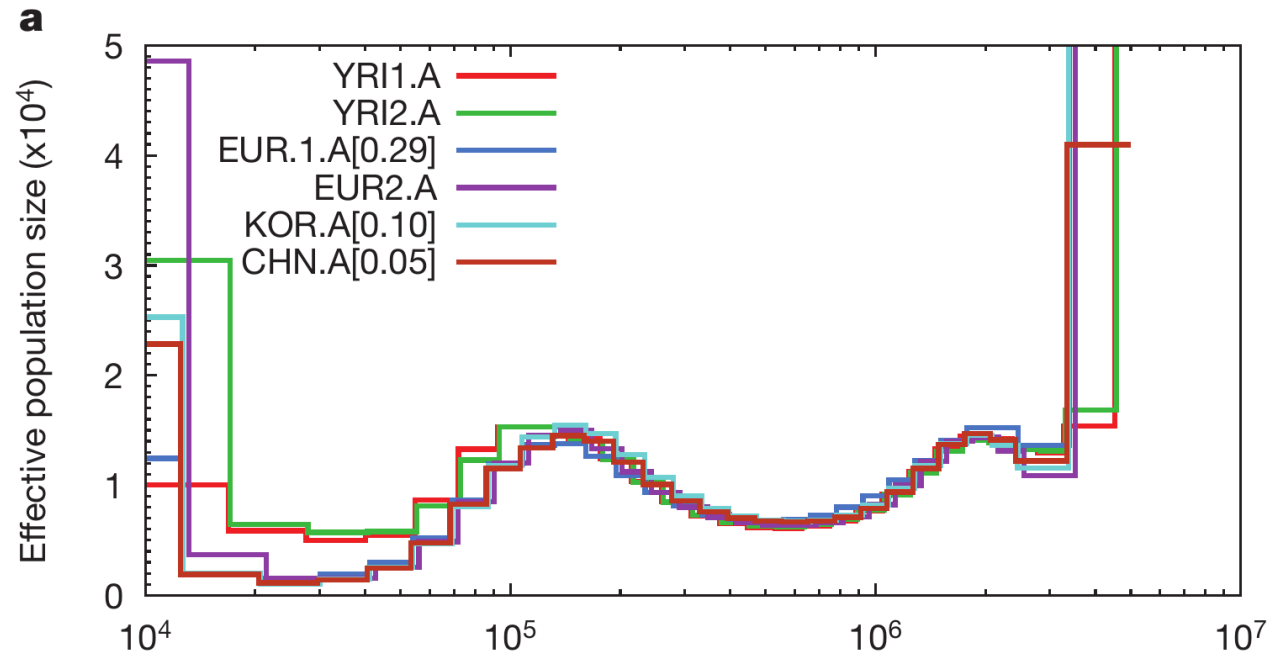
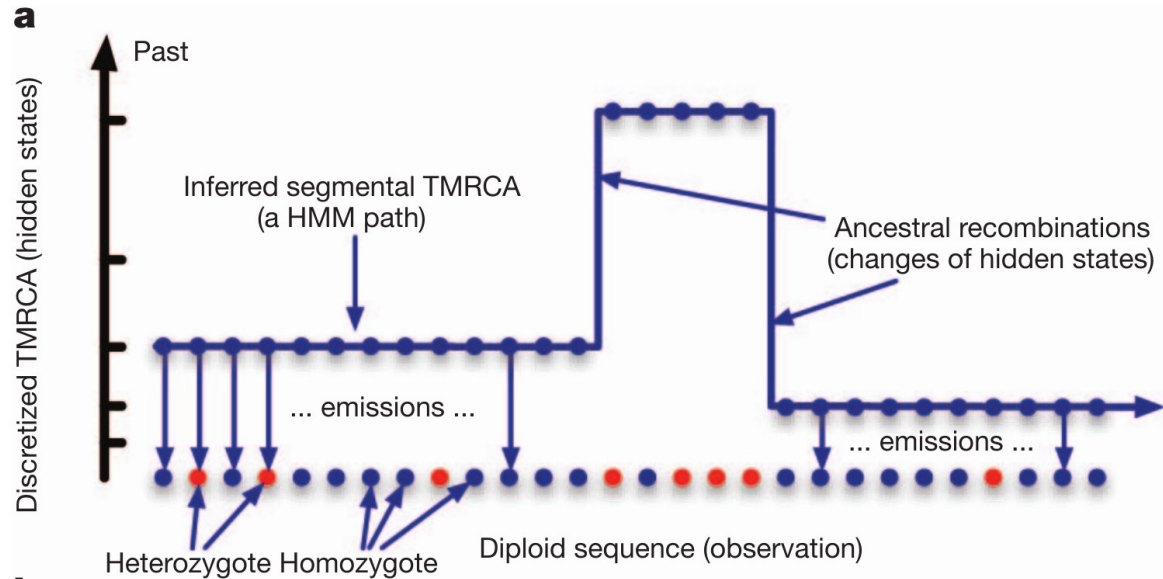
(a)



(b)



# Inferring population size from diversity data



# Population subdivision

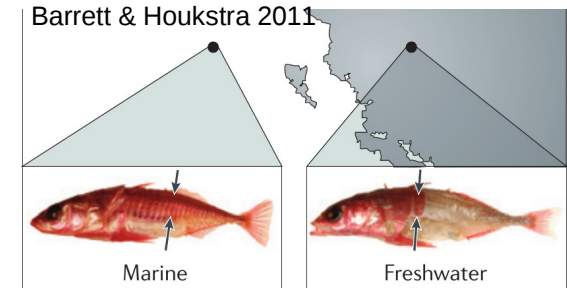
- Differences in allele frequency between sub-populations
- We can use genetic data to identify boundaries between sub-populations

TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAT**T**CTCGTGCC**T**AAGT**C**GATT**T**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAG**G**CTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAG**G**TCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGCC**T**AAGT**C**GATT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCCCTTTTCGCTGAG**G**CTCGTGCT**T**TAAGT**C**GATTATGATAAT

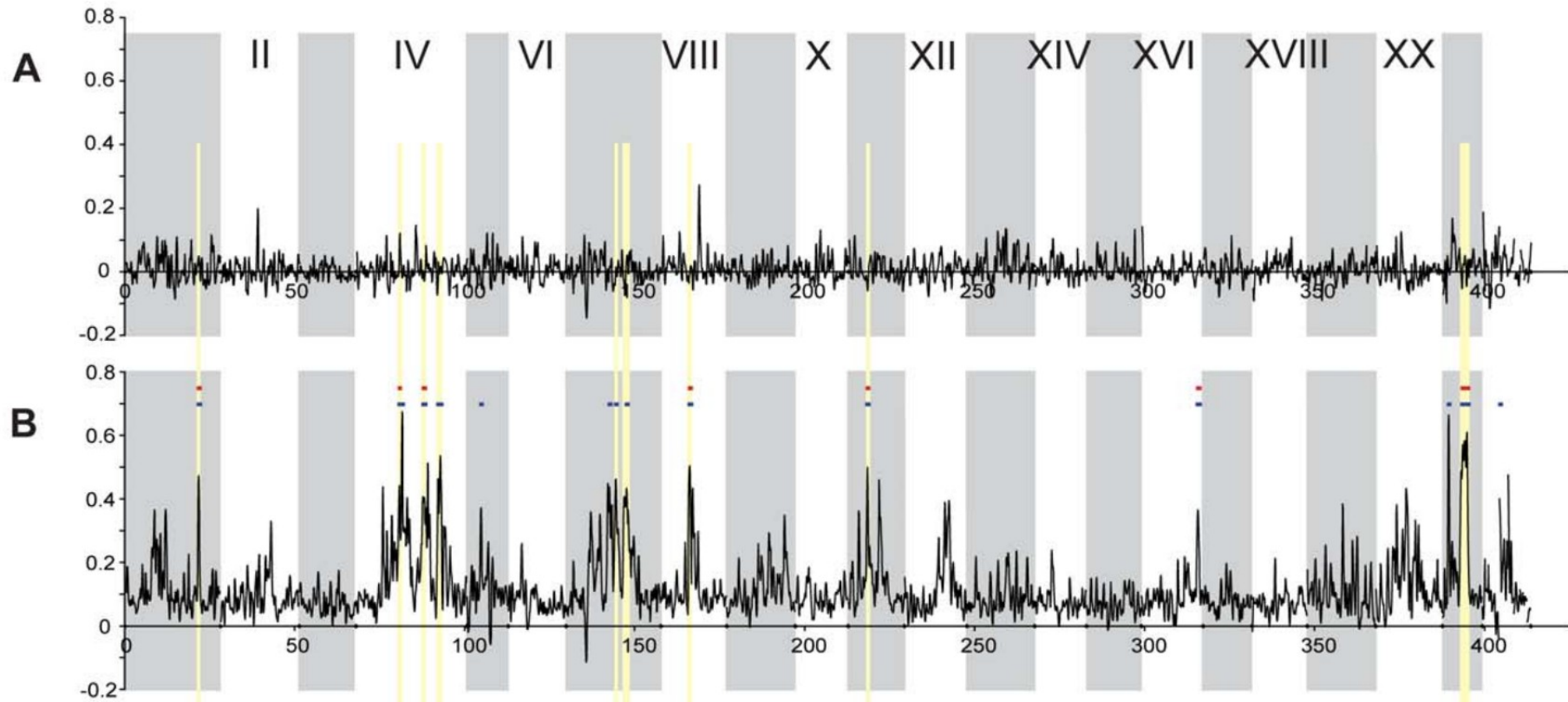
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAG**G**CTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCT**T**TAAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAG**G**CTCGTGCC**T**AAGTAGATT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAG**G**CTCGTGCC**T**AAGTAGATTATGATA**C**T

# Population subdivision – a genome-wide view

- $F_{ST}$  is a measure of differences in allele frequencies between populations



$F_{ST}$  between two marine populations (A) or between marine and freshwater populations (B) (Hohenlohe et al. 2010)

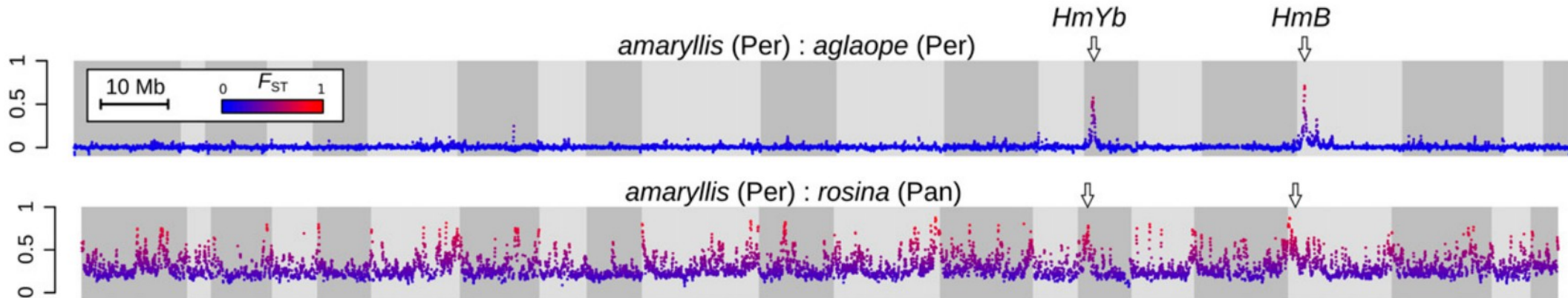
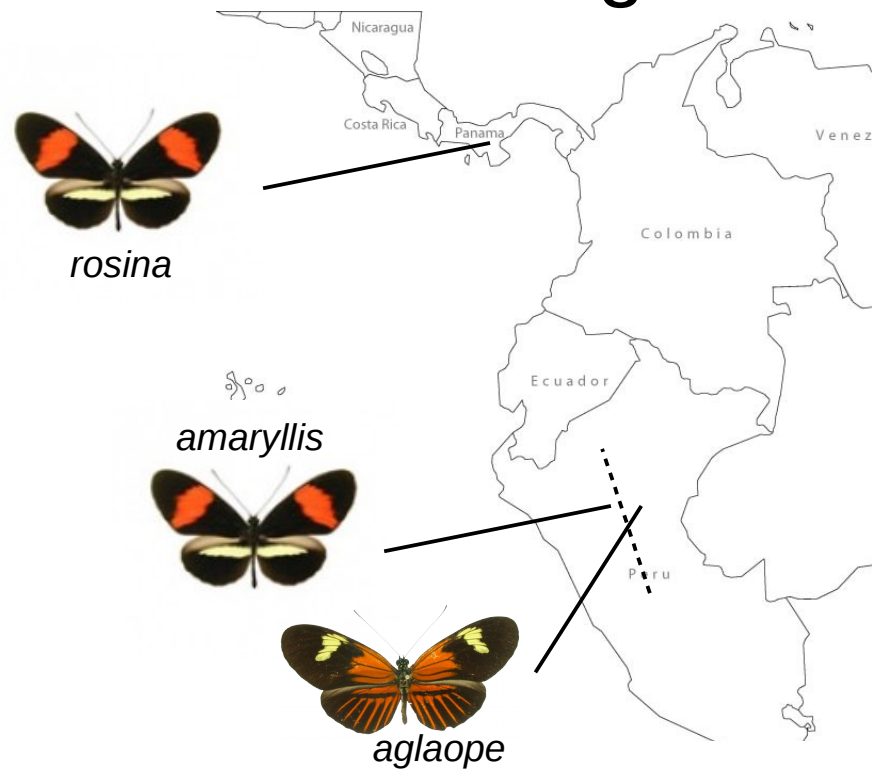




# Population subdivision – a genome-wide view

Parapatric wing-pattern races of *Heliconius melpomene* are only strongly differentiated at two wing patterning loci.

(Martin et al. 2013  
Genome Research)



# Population subdivision – identifying populations using genomic data

- Given a genetic dataset, can we infer whether there are distinct populations?

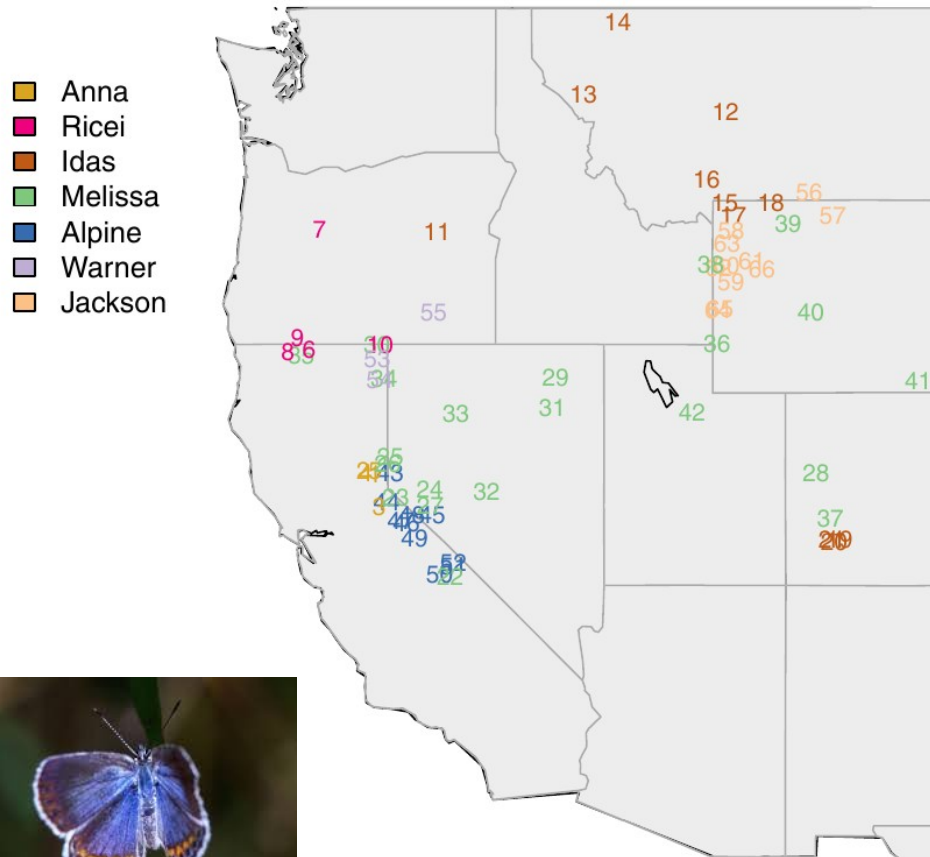
TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGA**T**CTCGTGC**C**TAAGT**C**GATT**T**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAG**G**TCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTG**C**TAAGT**C**GATT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GATTATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGA**T**CTCGTGC**C**TAAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTG**C**TAAGTAGAT**C**ATGATAAT  
TAGATCGTCCAGATCGATCTAGCCC**G**TTTTCGCTGA**T**CTCGTGCT**T**TAAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGA**T**CTCGTG**C**TAAGTAGAT**C**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTG**C**TAAGTAGAT**T**ATGATAAT  
TAGATCGTC**G**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTG**C**TAAGTAGATTATGATA**C**T



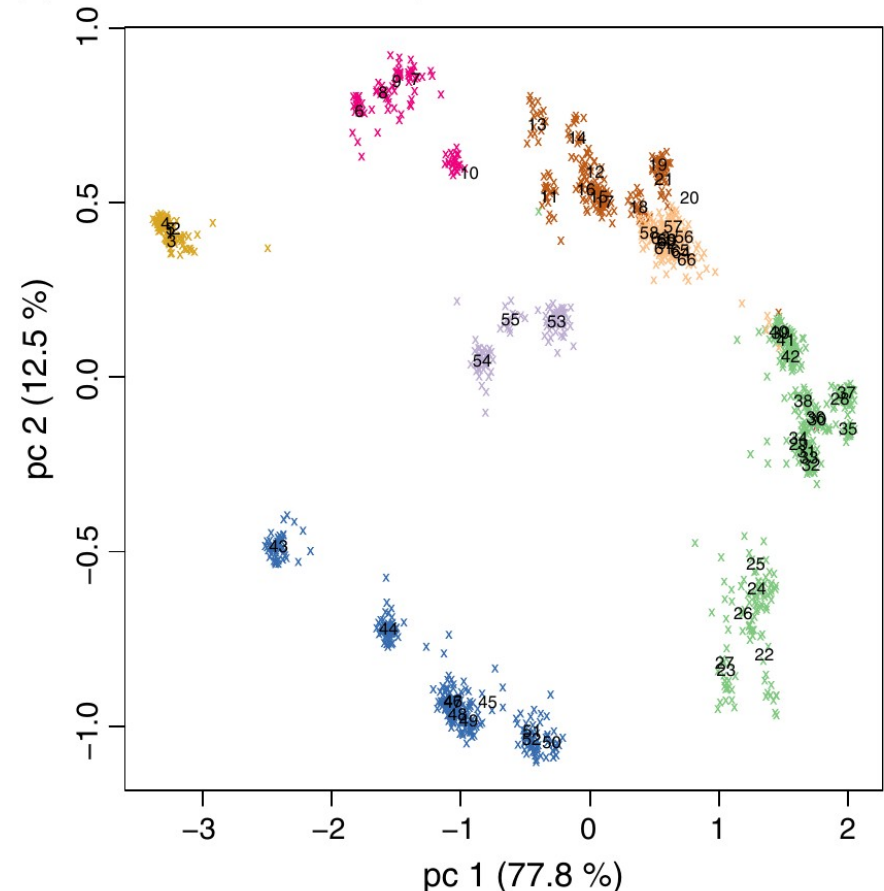
# Population subdivision – identifying populations using genomic data

- PCA (Principle Components Analysis) simplifies huge genetic datasets into two (or more) dimensions

(a) Sample locations

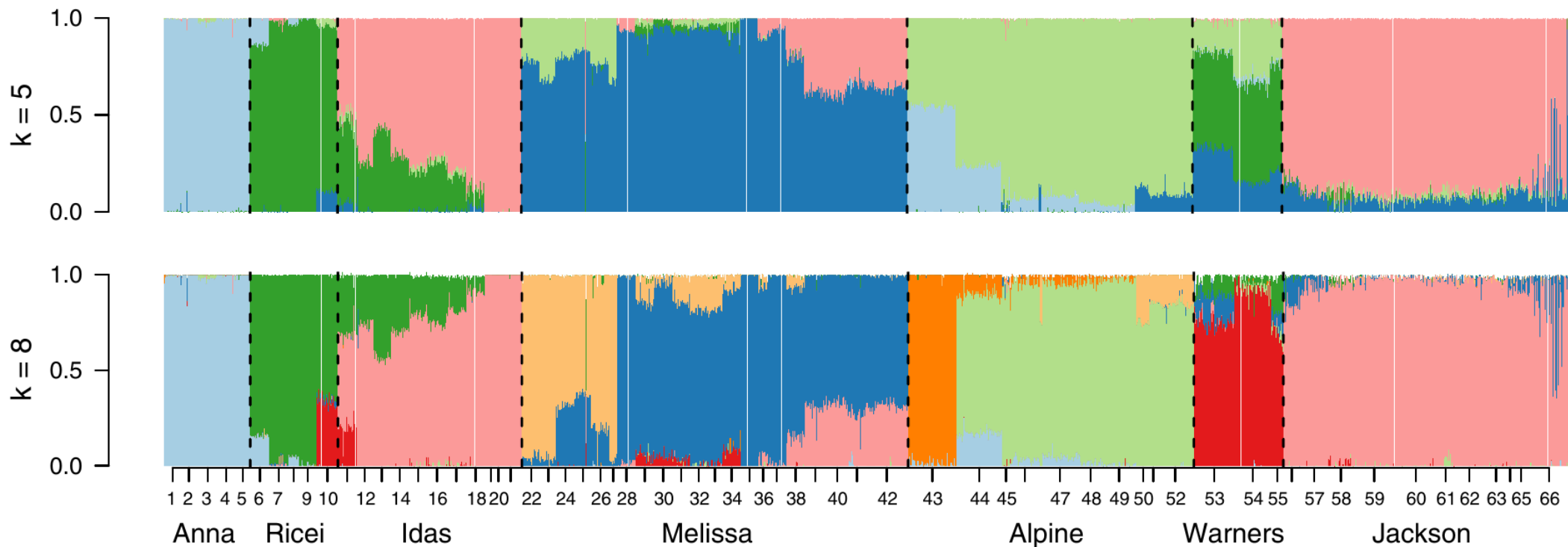


(b) Common variant PCA plot

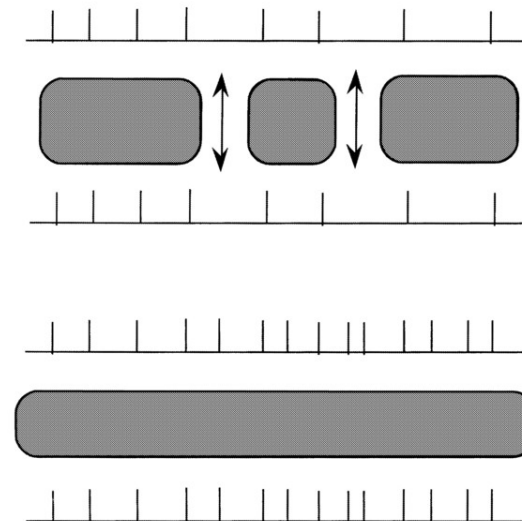
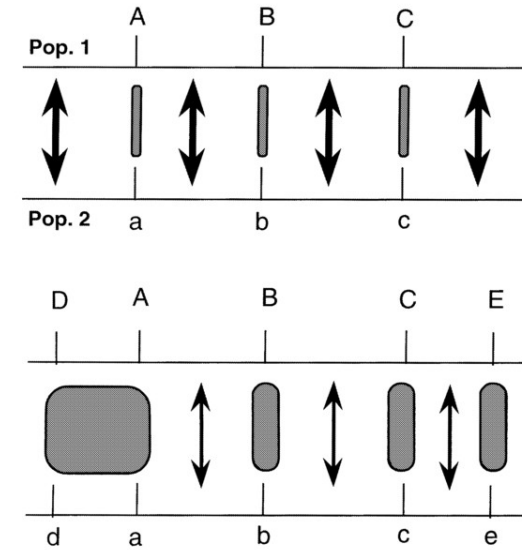
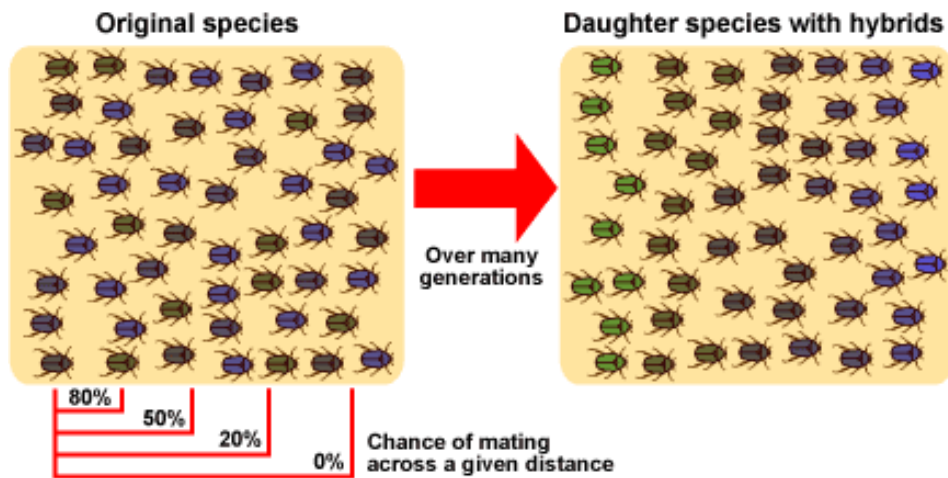
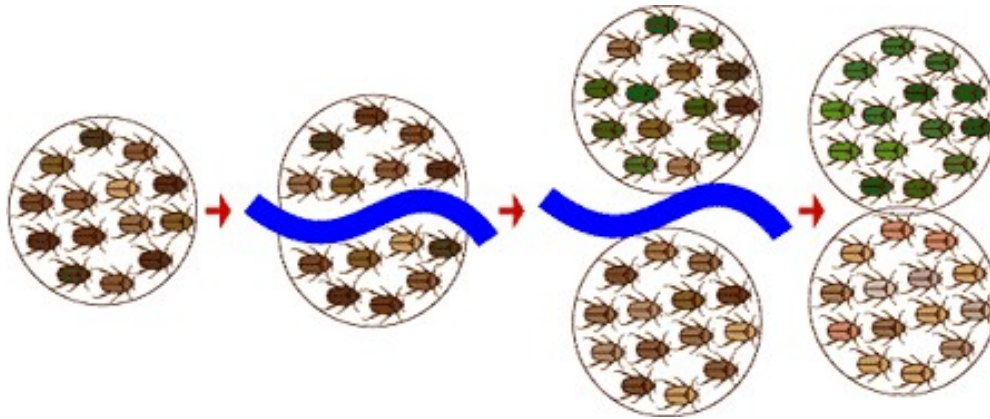


# Population subdivision – identifying populations using genomic data

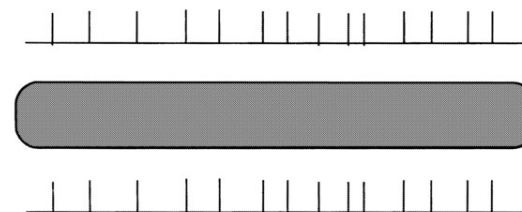
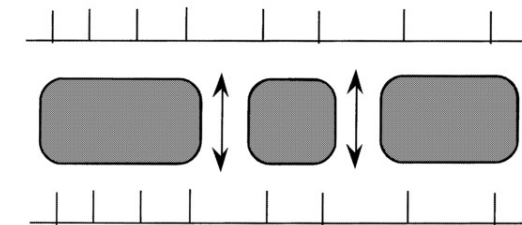
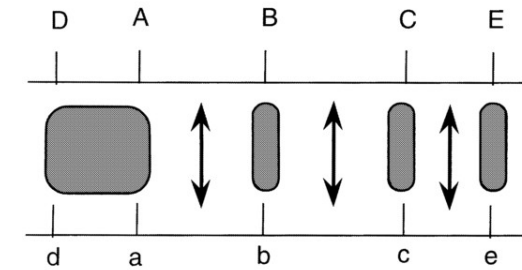
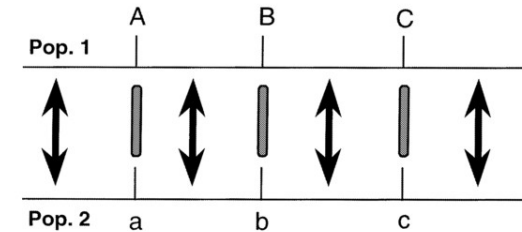
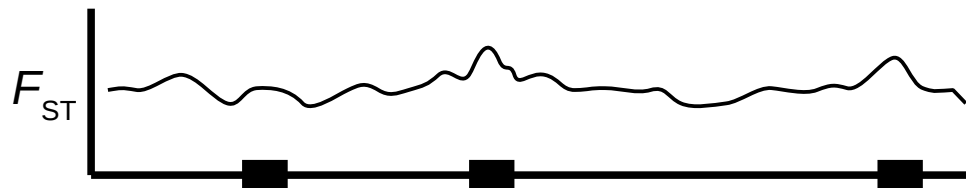
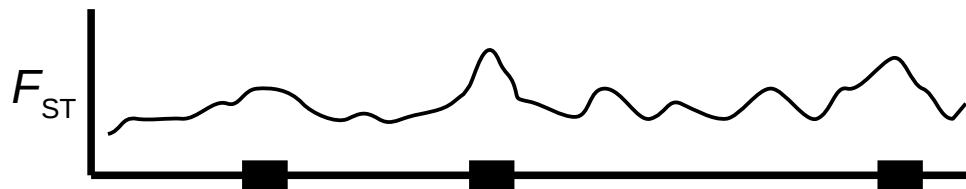
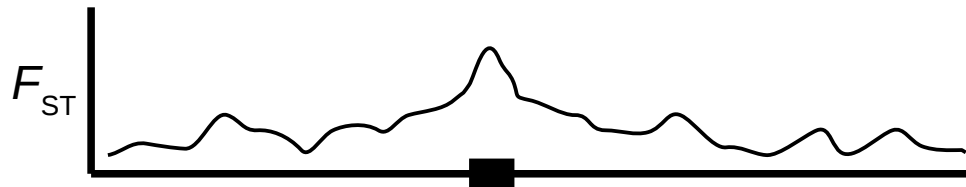
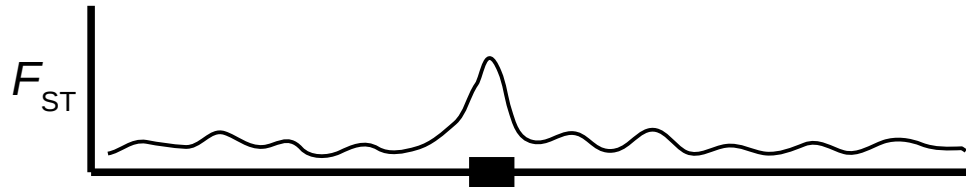
- STRUCTURE (Pritchard et al. 2000) estimates the likelihood that each sample falls into each of  $k$  clusters



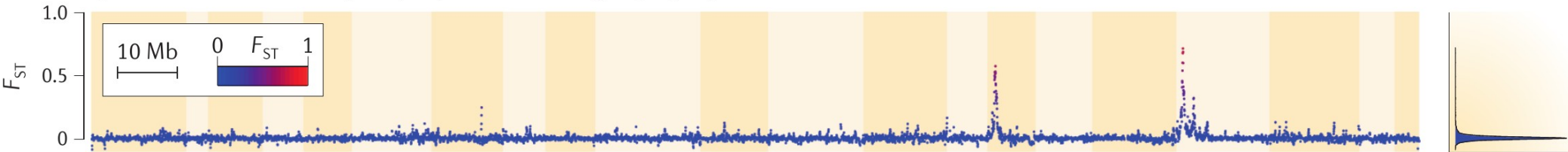
# Speciation and gene flow



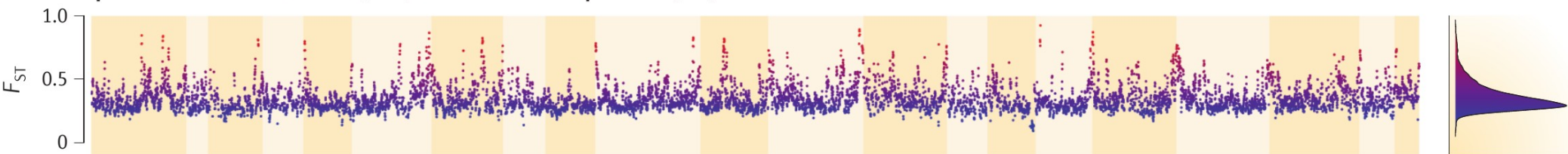
# Speciation and gene flow



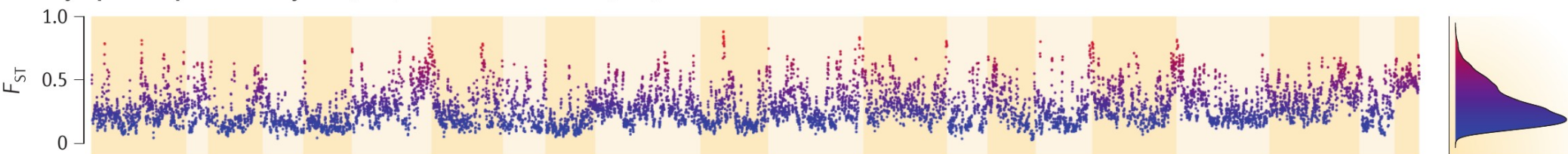
**Aa** Parapatric races: *H. m. amaryllis* (Per) versus *H. m. aglaope* (Per)



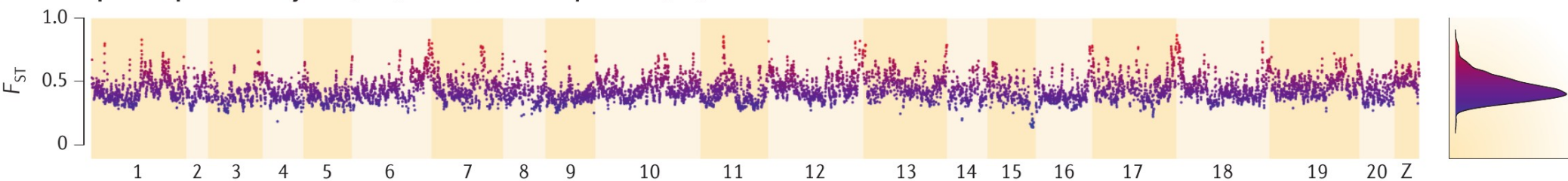
**Ab** Allopatric races: *H. m. rosina* (Pan) versus *H. m. melpomene* (FG)



**Ac** Sympatric species: *H. cydno* (Pan) versus *H. m. rosina* (Pan)

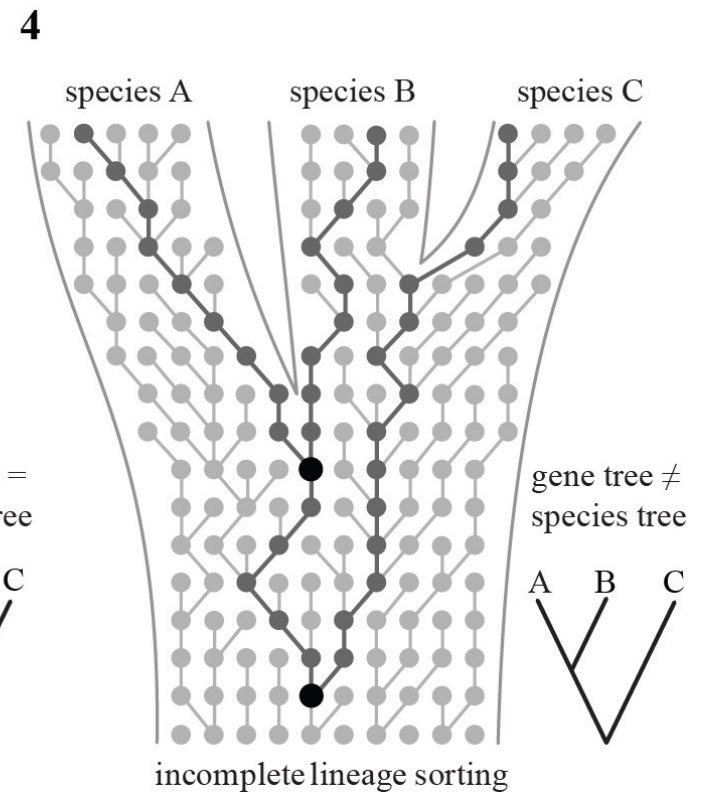
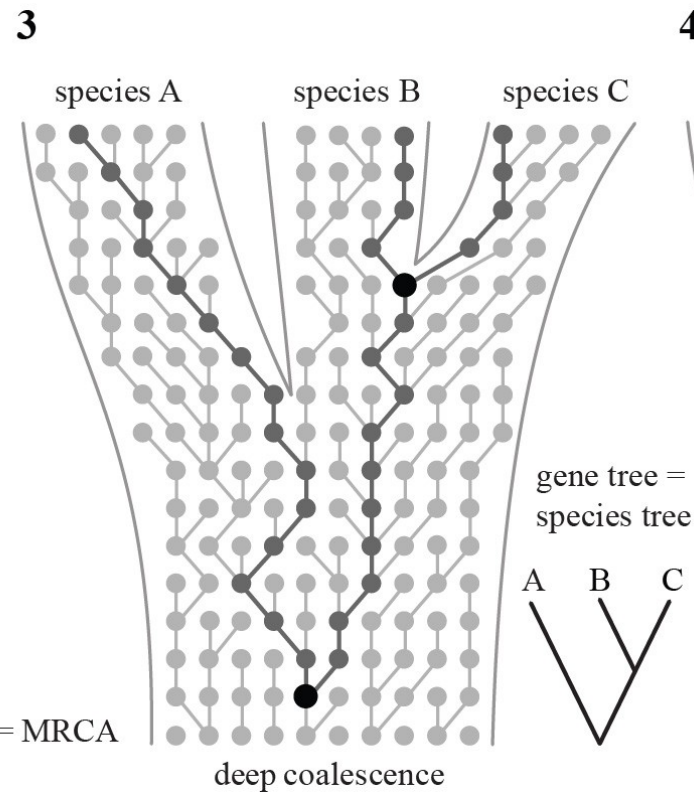
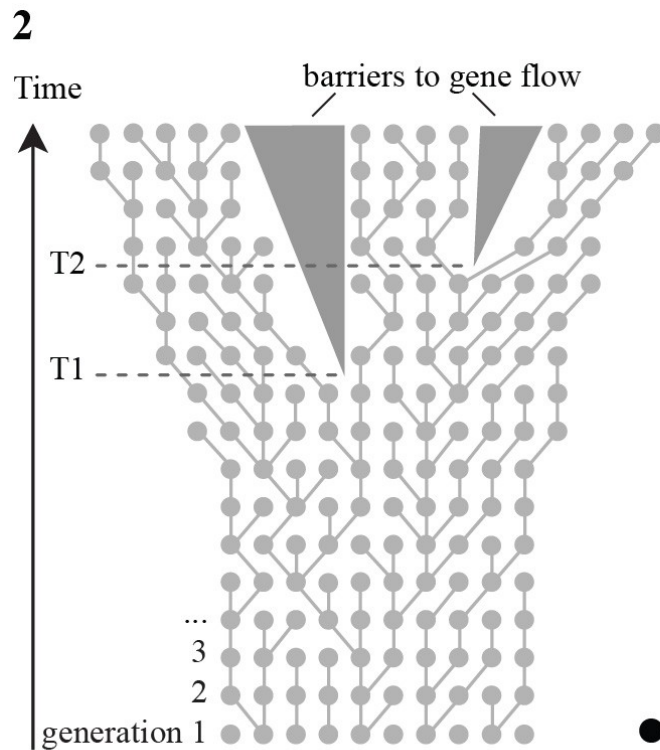


**Ad** Allopatric species: *H. cydno* (Pan) versus *H. m. melpomene* (FG)





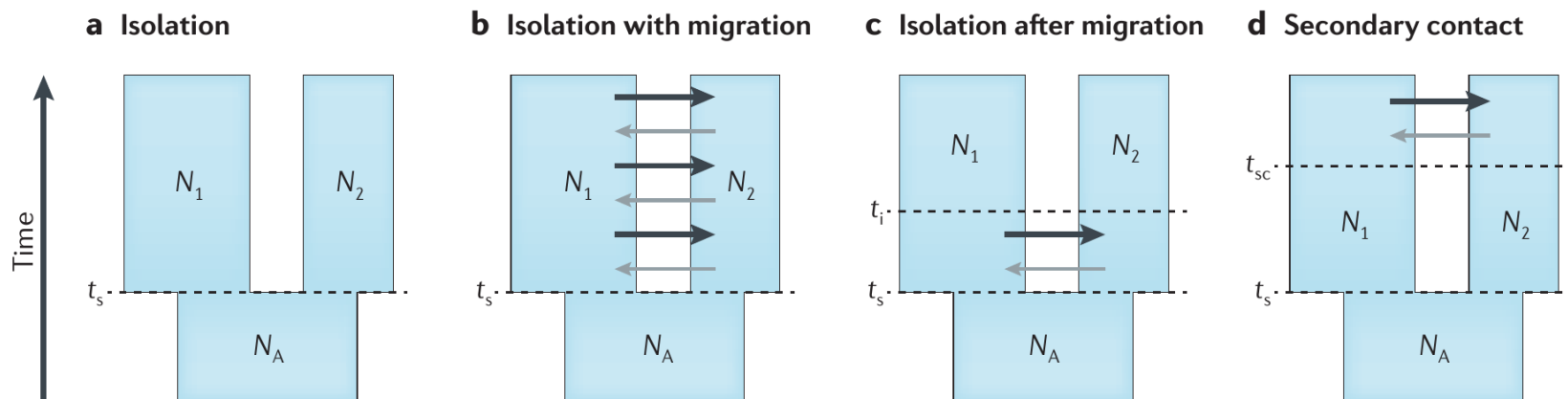
# Detecting gene flow



# Detecting gene flow

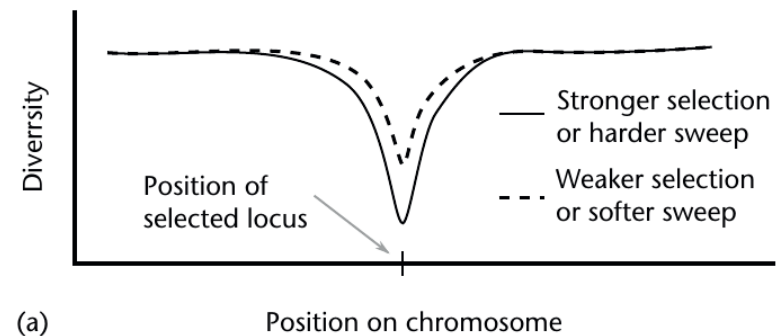
Huge genomic data sets allow sensitive fitting of evolutionary models.

TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAT**T**CTCGTGCC**T**AAGT**C**GAT**T**ATGATAAT  
 TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
 TAGATCGTCCAGATCGA**A**CTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
 TAG**G**TCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**C**ATGATAAT  
 TAGATCGTCCAGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCC**T**AAGT**C**GAT**T**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCCCTTTTCGCTGAGCTCGTGCT**T**TAAGT**C**GAT**T**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
 TAGATCGTCCAGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCT**T**TAAGTAGAT**C**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCC**G**TTTTCGCTGAT**T**CTCGTGCC**T**AAGTAGAT**C**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTGCC**T**AAGTAGAT**T**ATGATAAT  
 TAGATCGT**C**AGATCGATCTAGCCC**G**TTTTCGCTGAGCTCGTGCC**T**AAGTAGATTATGATA**C**T

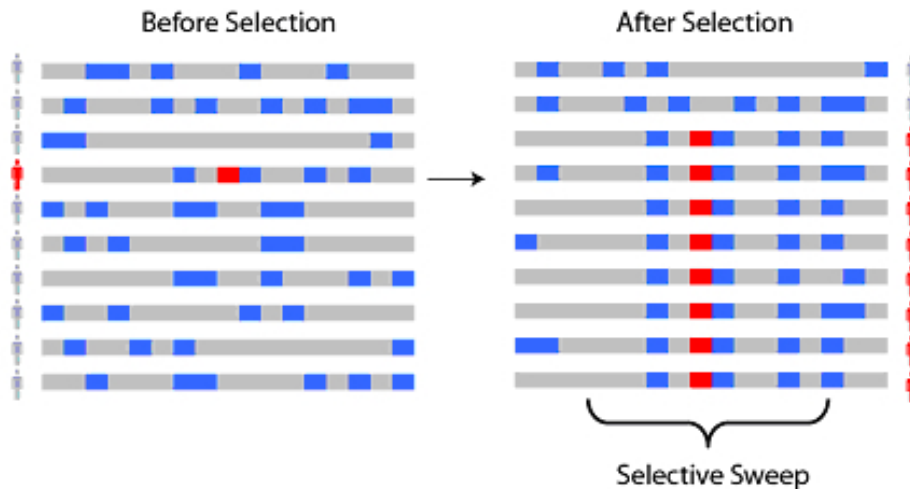


# Selective sweeps

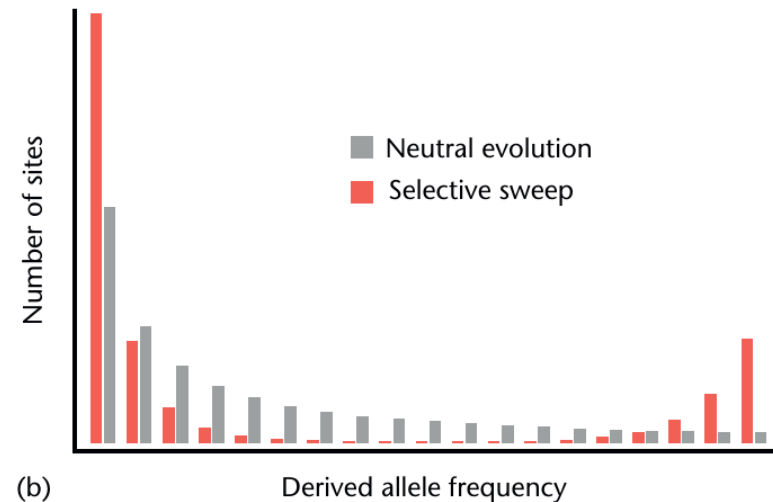
- Strong selection rapidly fixes a beneficial allele in the population
- Causes reduced diversity at the selected site and a change in the site frequency spectrum



(a)



nature.com

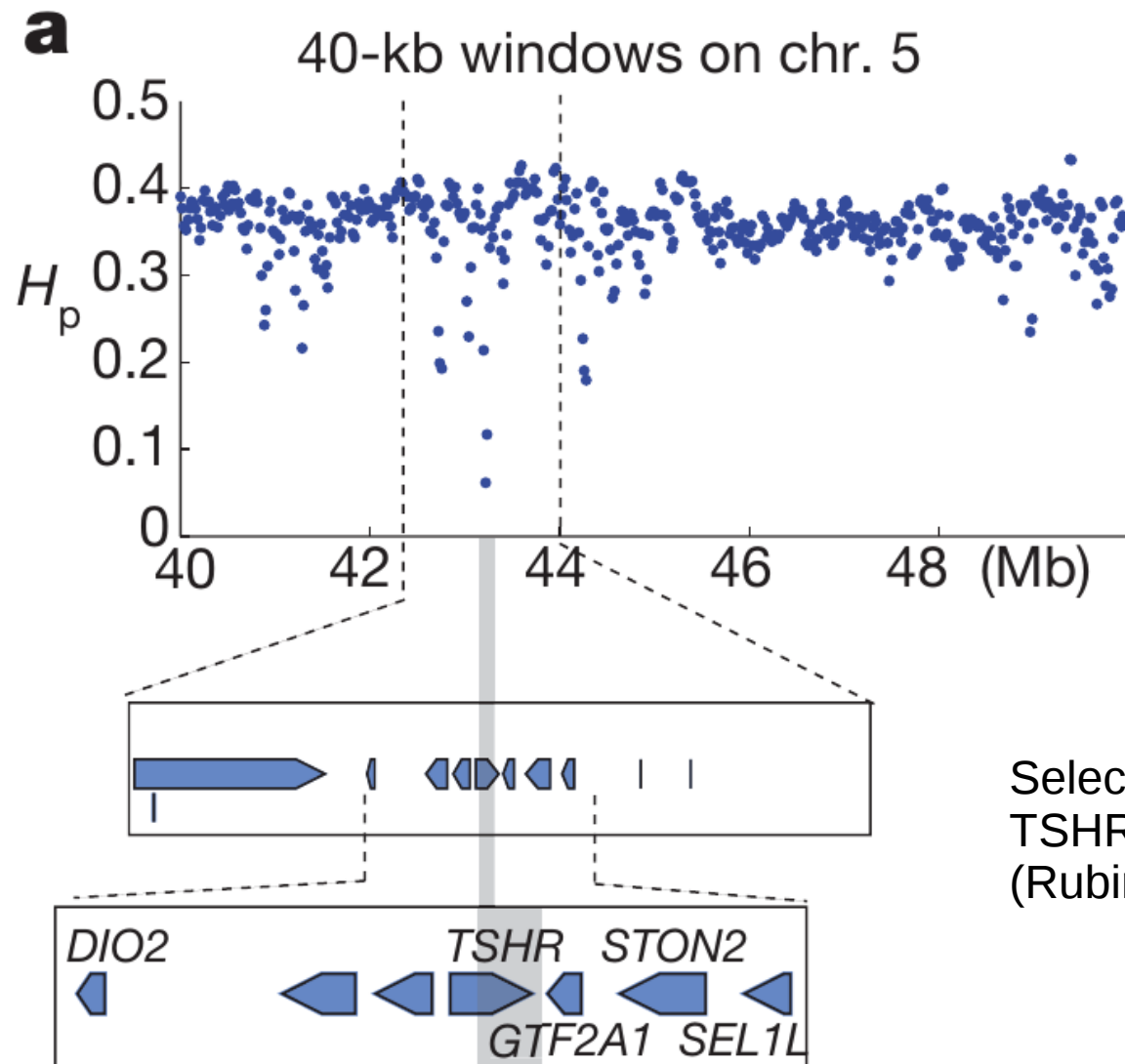


(b)

Martin et al. 2013 eLS

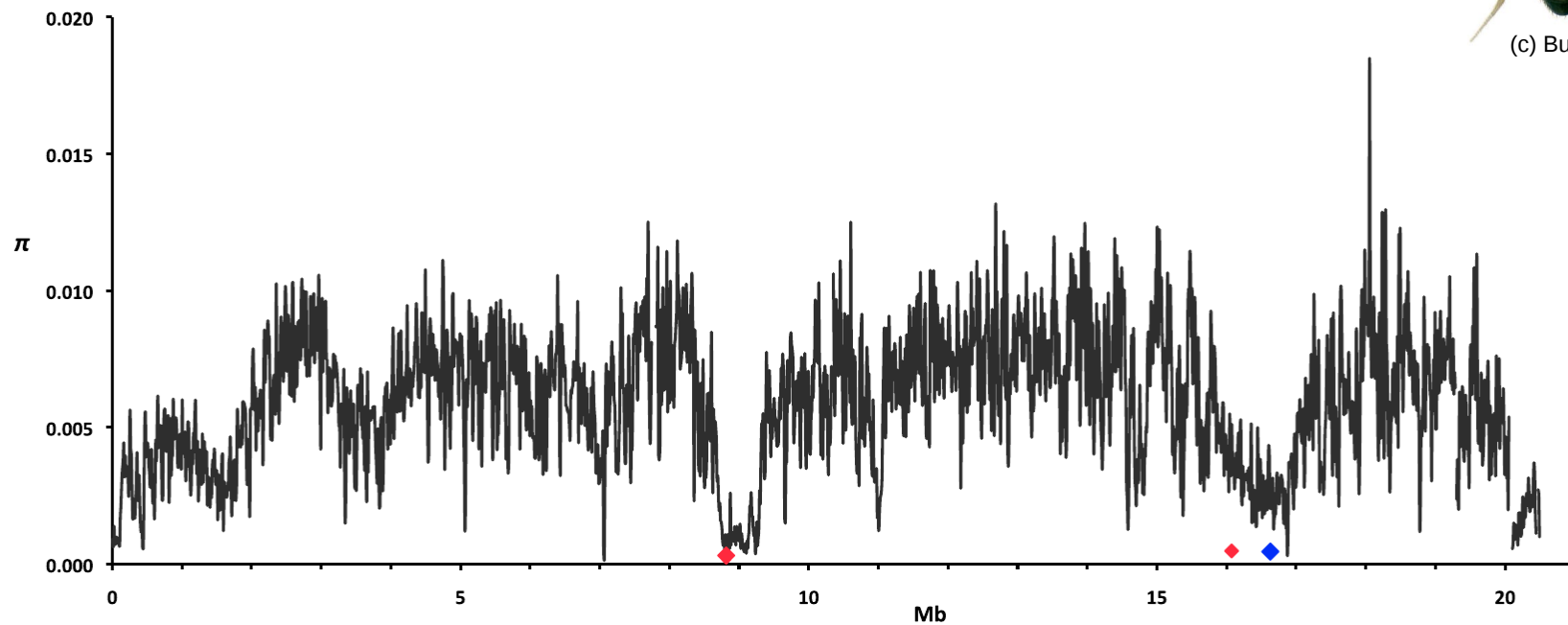


# Selective sweeps



Selective sweep at the  
TSHR gene in chickens.  
(Rubin et al. 2012 Nature)

# Selective sweeps



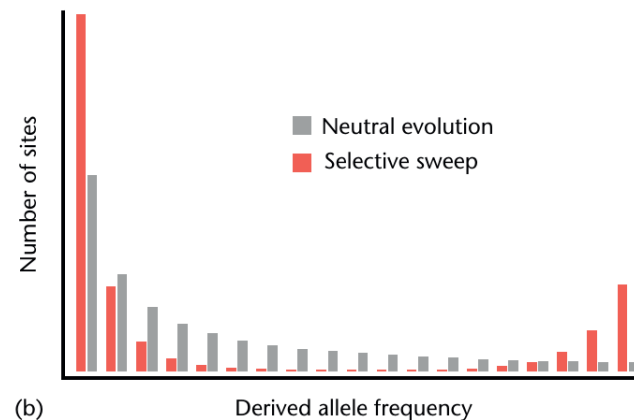
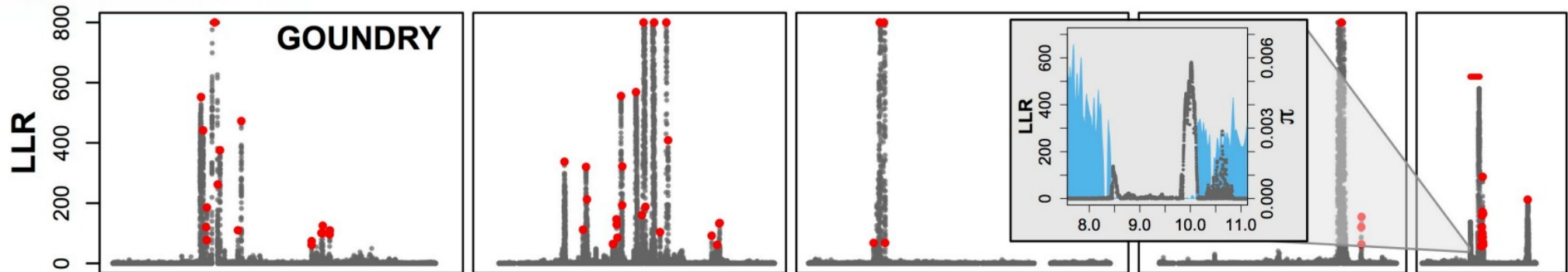
**Figure 5.** Nucleotide diversity ( $\pi$ ) along the *D. mauritiana* X chromosome. The location of genes potentially causing the two selective sweeps are indicated: (large red diamond) *MDox/Dox*; (large blue diamond) *Odsh*; (small red diamond) *E(Dox)*. Nucleotide diversity ( $\pi$ ) is plotted in nonoverlapping 10-kb windows.

# Selective sweeps

- **Sweepfinder** (Nielsen et al. 2005) identifies regions showing a strong skew in the site frequency spectrum.

















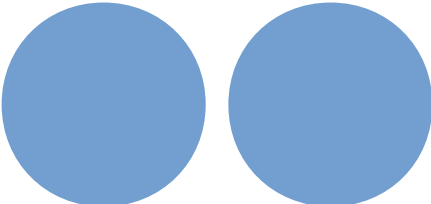



**Figure 1:**

Selective sweeps in *Anopheles gambiae* – Crawford et al. 2014













# Designing a population genomic study

- Number of samples

			Demography	Selection	Association
10s - 100s					
5-10					
1-5					
pools					

# Designing a population genomic study

- Type of sequencing

			Demography	Selection	Association
Whole Genome					
Targeted Capture					
Sub-genomic (e.g. RADseq)	