# Using the the site frequency spectrum to infer demography and natural selection
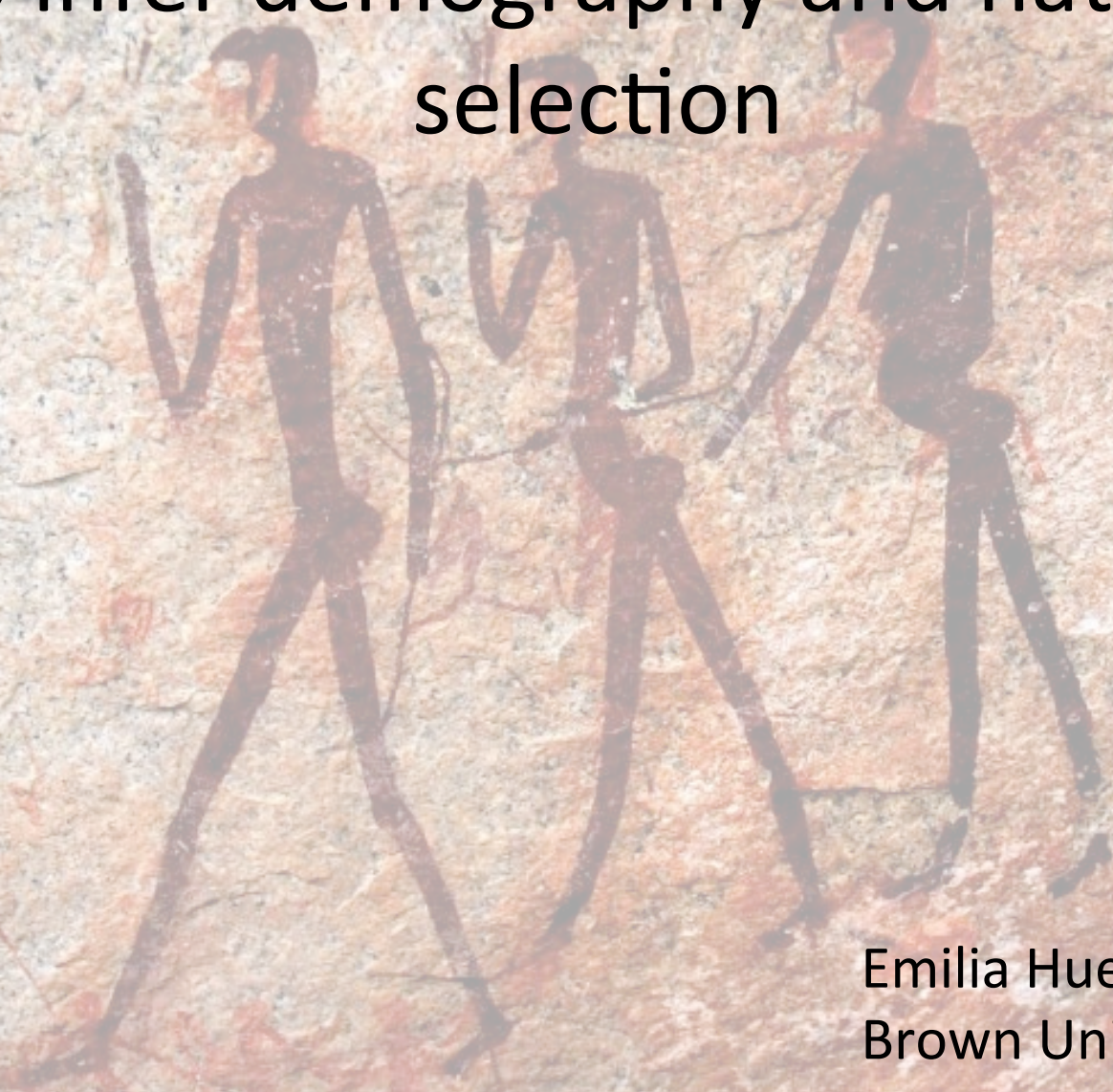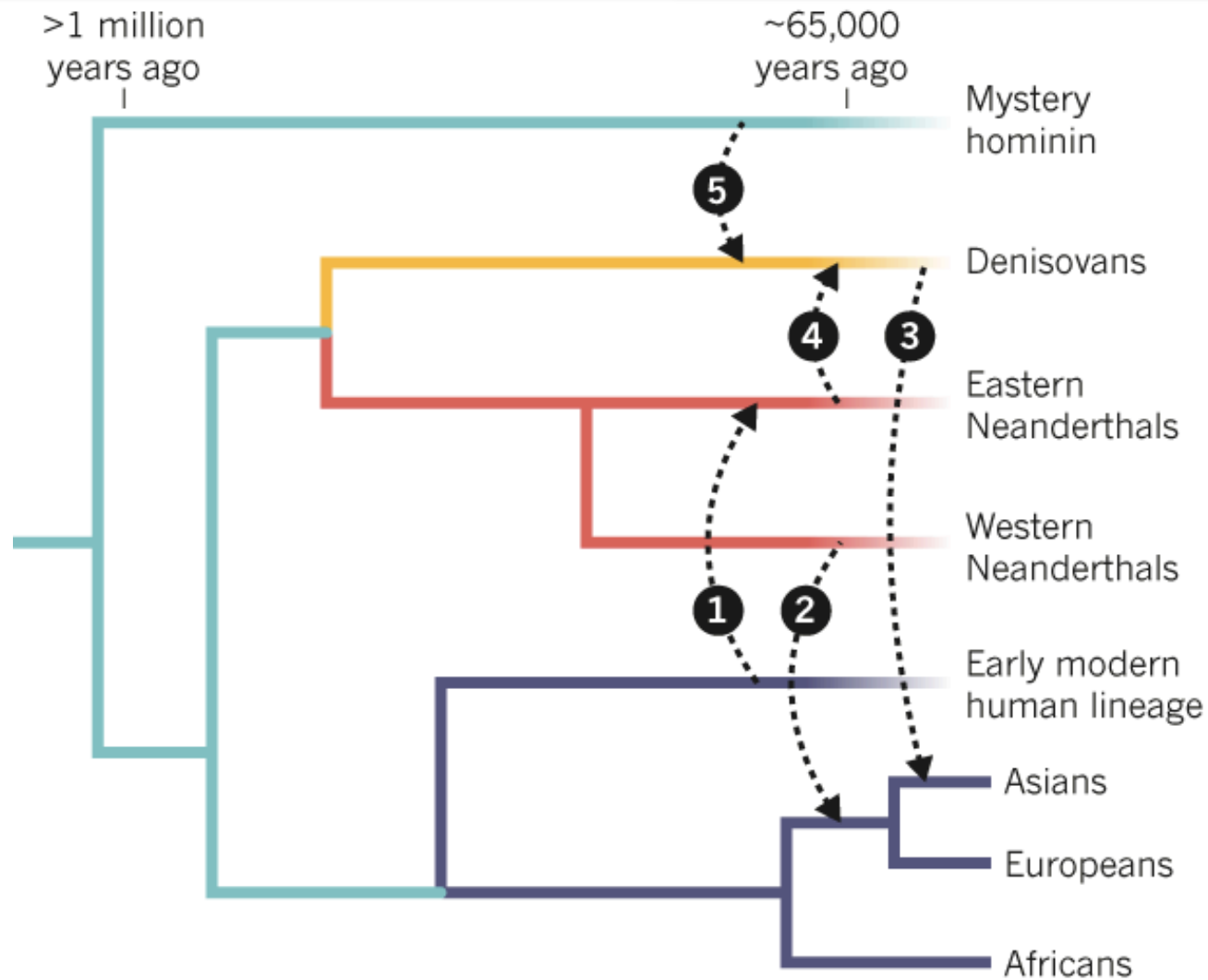
Emilia Huerta-Sanchez
Brown University

# Recent time scale



>1 million years ago

~65,000 years ago

Mystery hominin

Denisovans

Eastern Neanderthals

Western Neanderthals

Early modern human lineage

Asians

Europeans

Africans

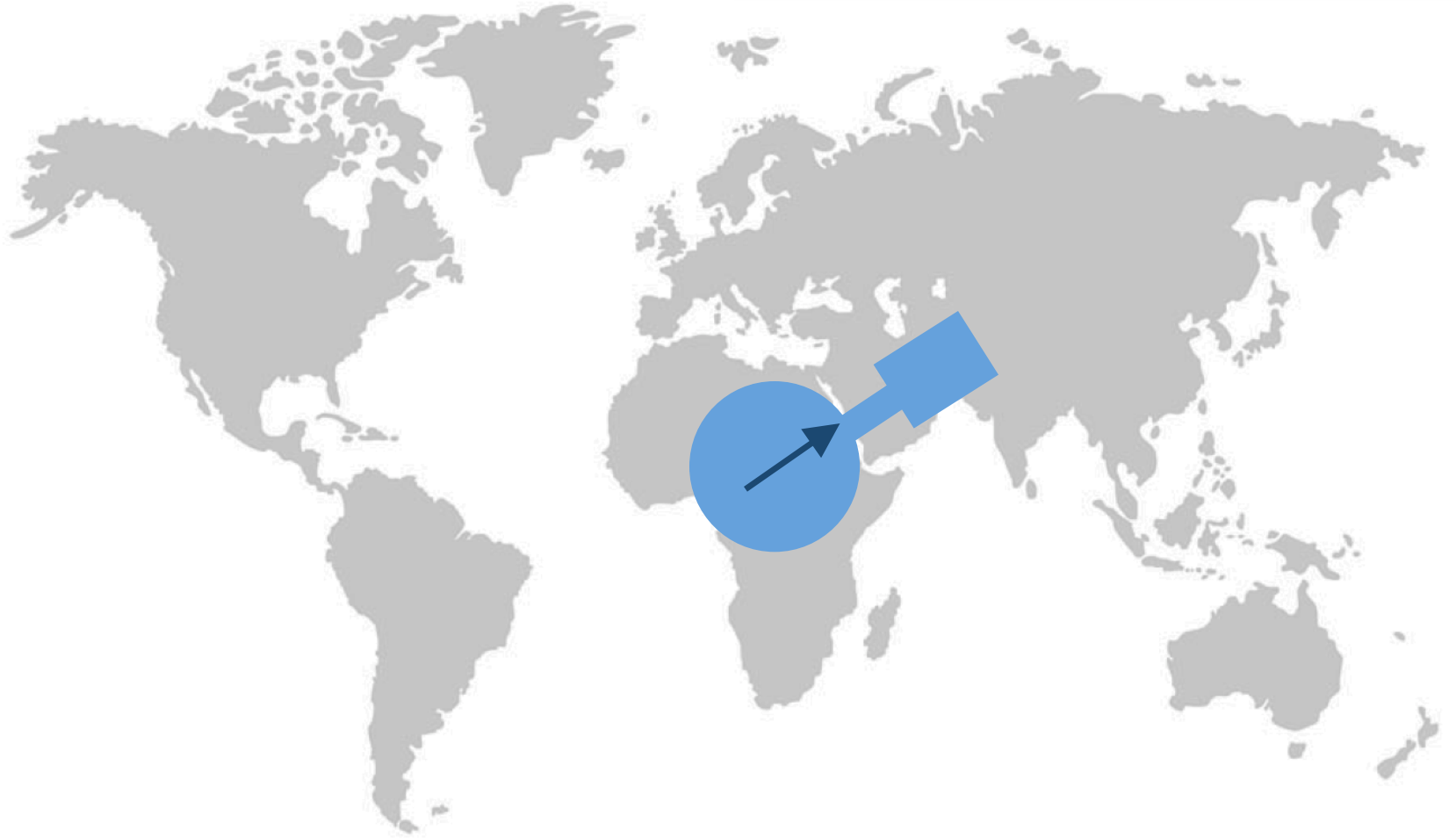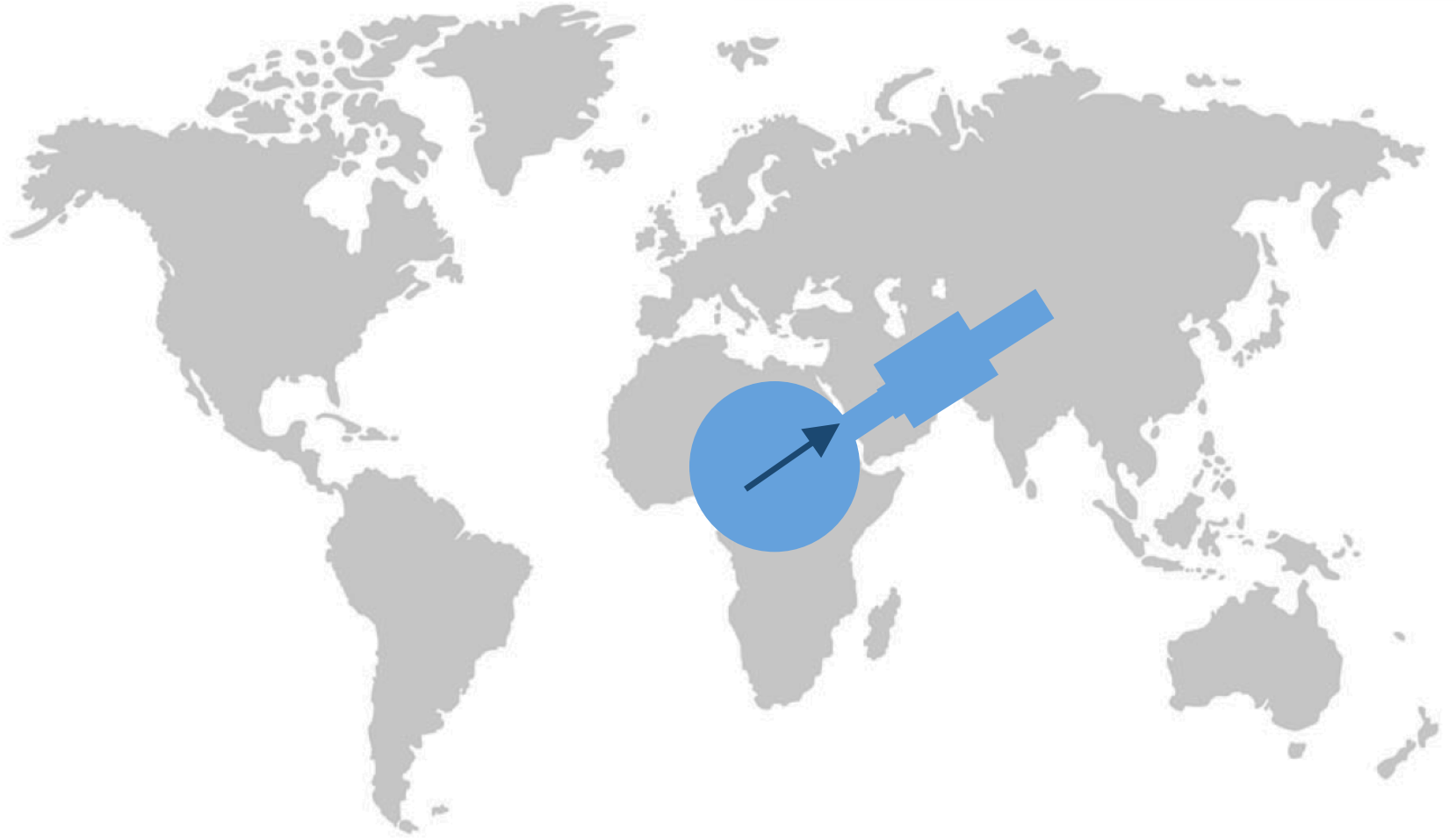···· Interbreeding episode/event

©nature

# Population dispersal

# Population dispersal

# Population dispersal

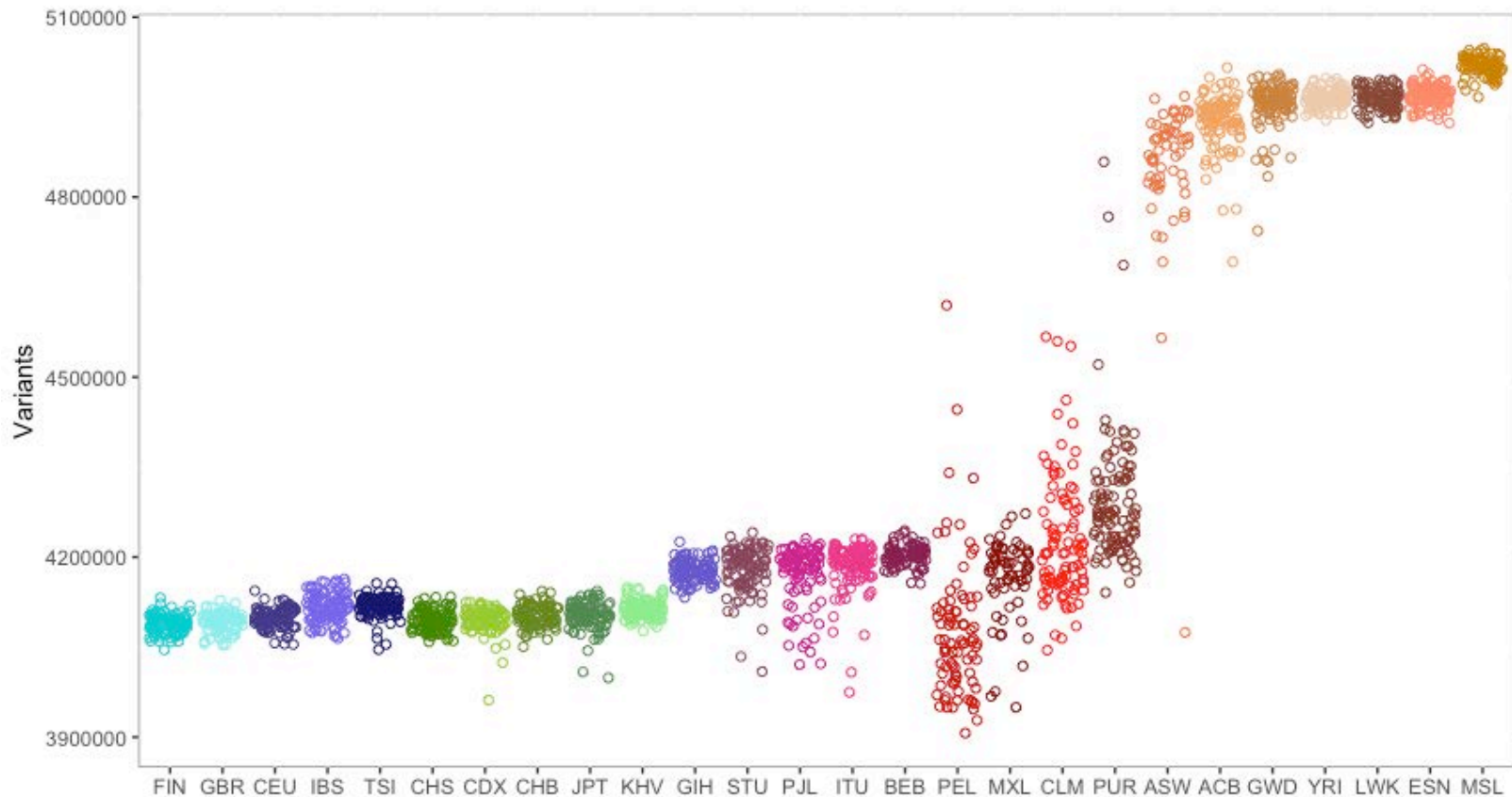# Population dispersal

# Heterozygosity
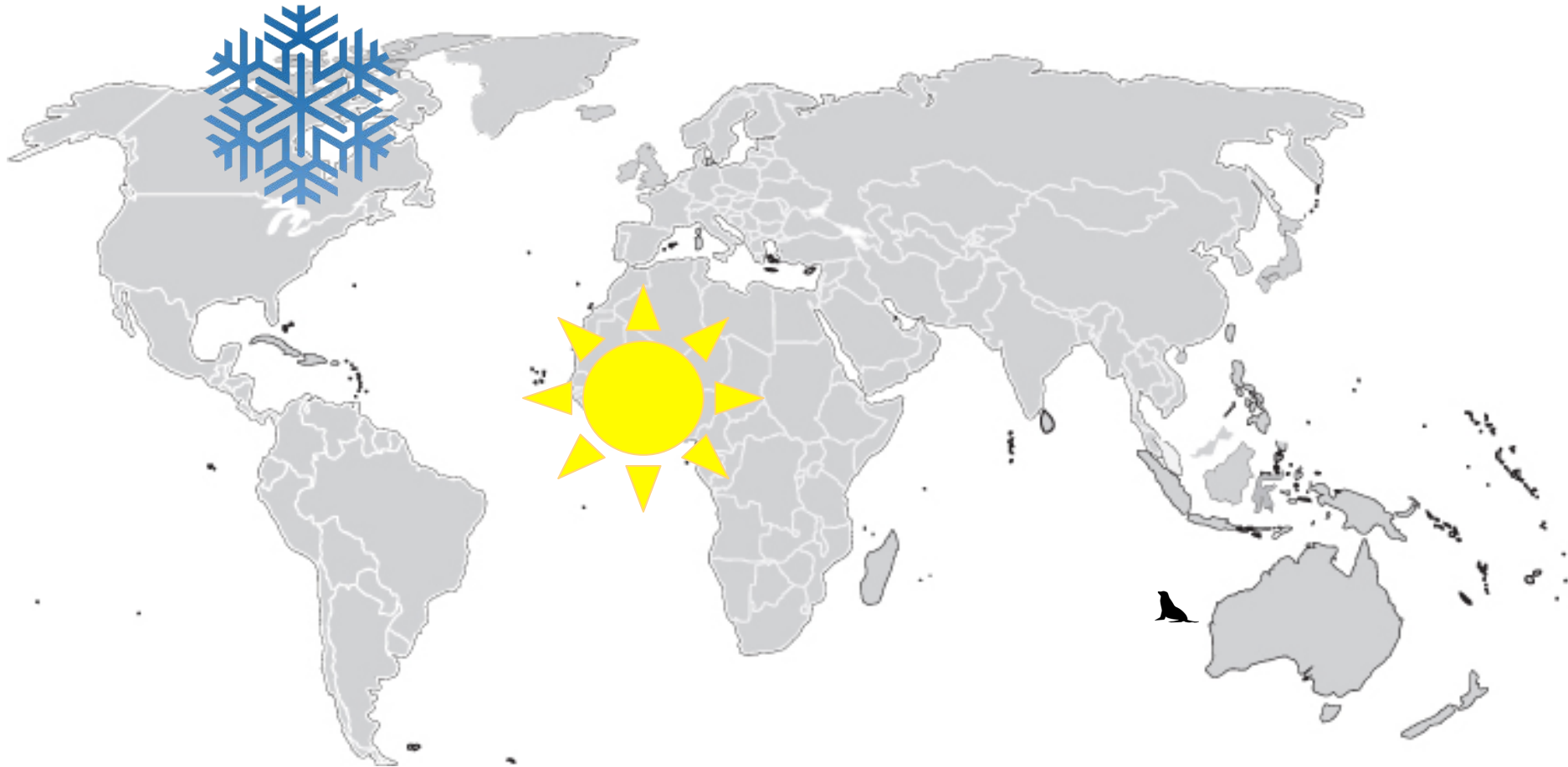
DeGiorgio et al. 2009

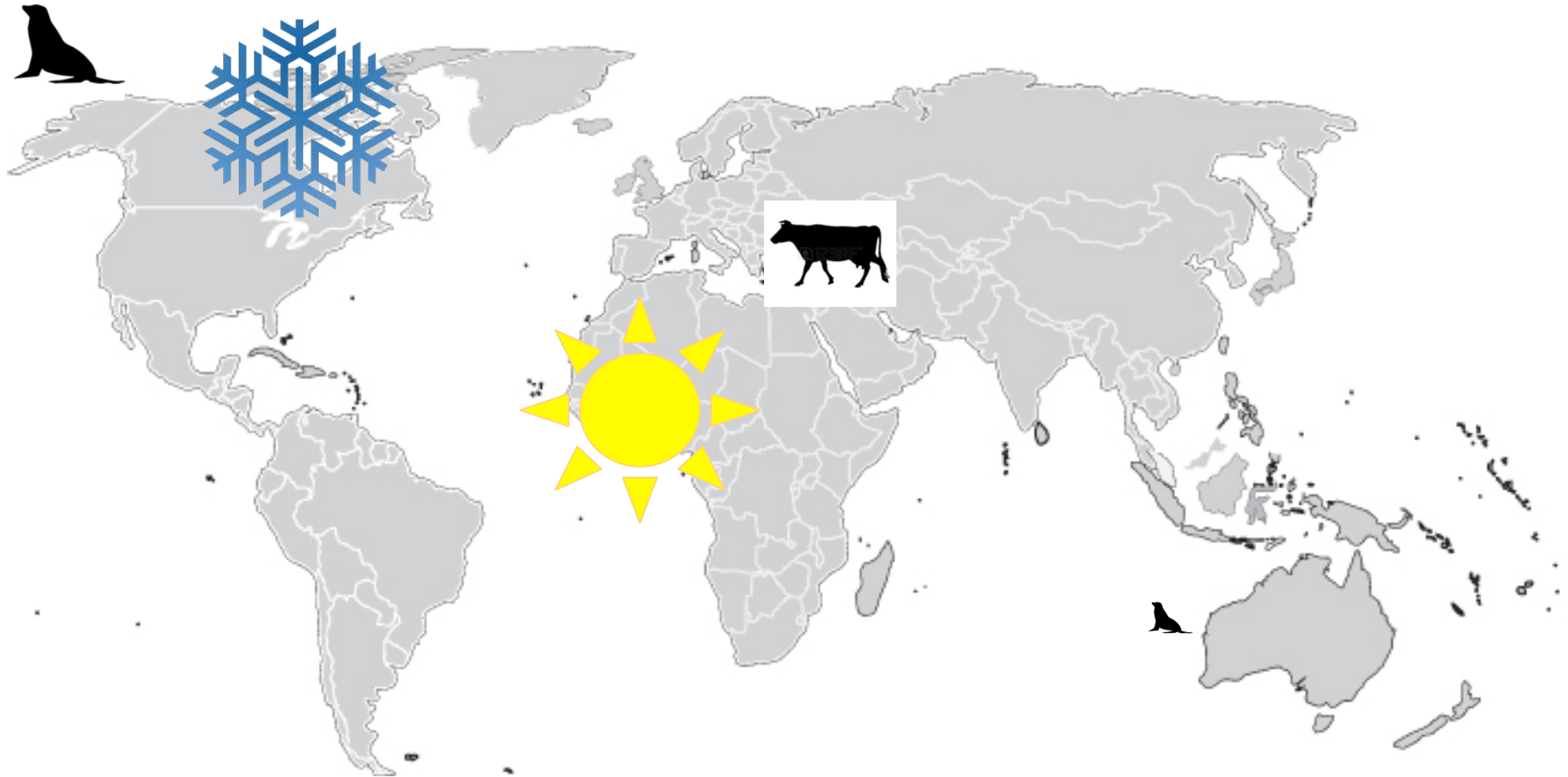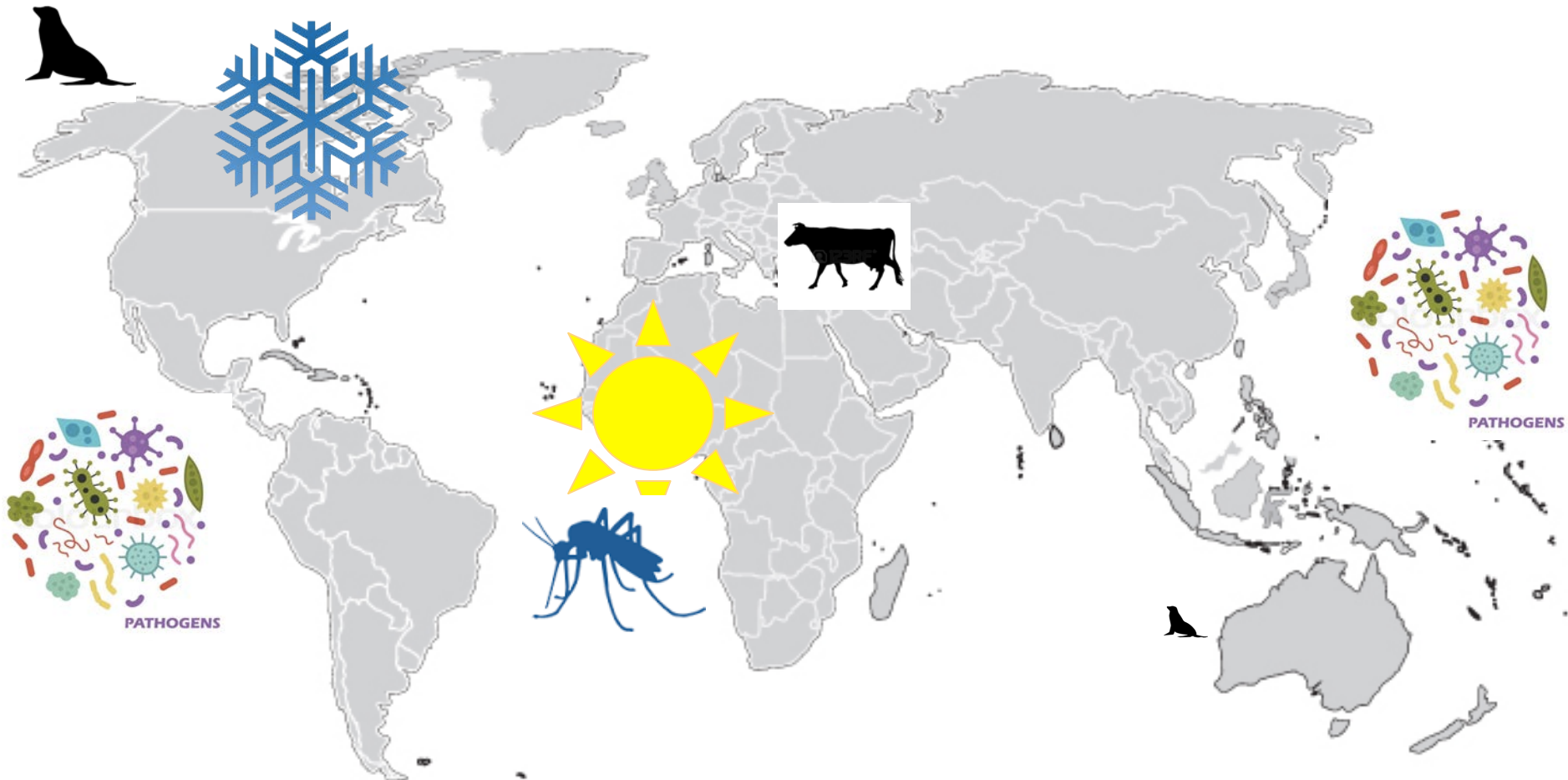# Number of variants per individual

# Positive natural selection
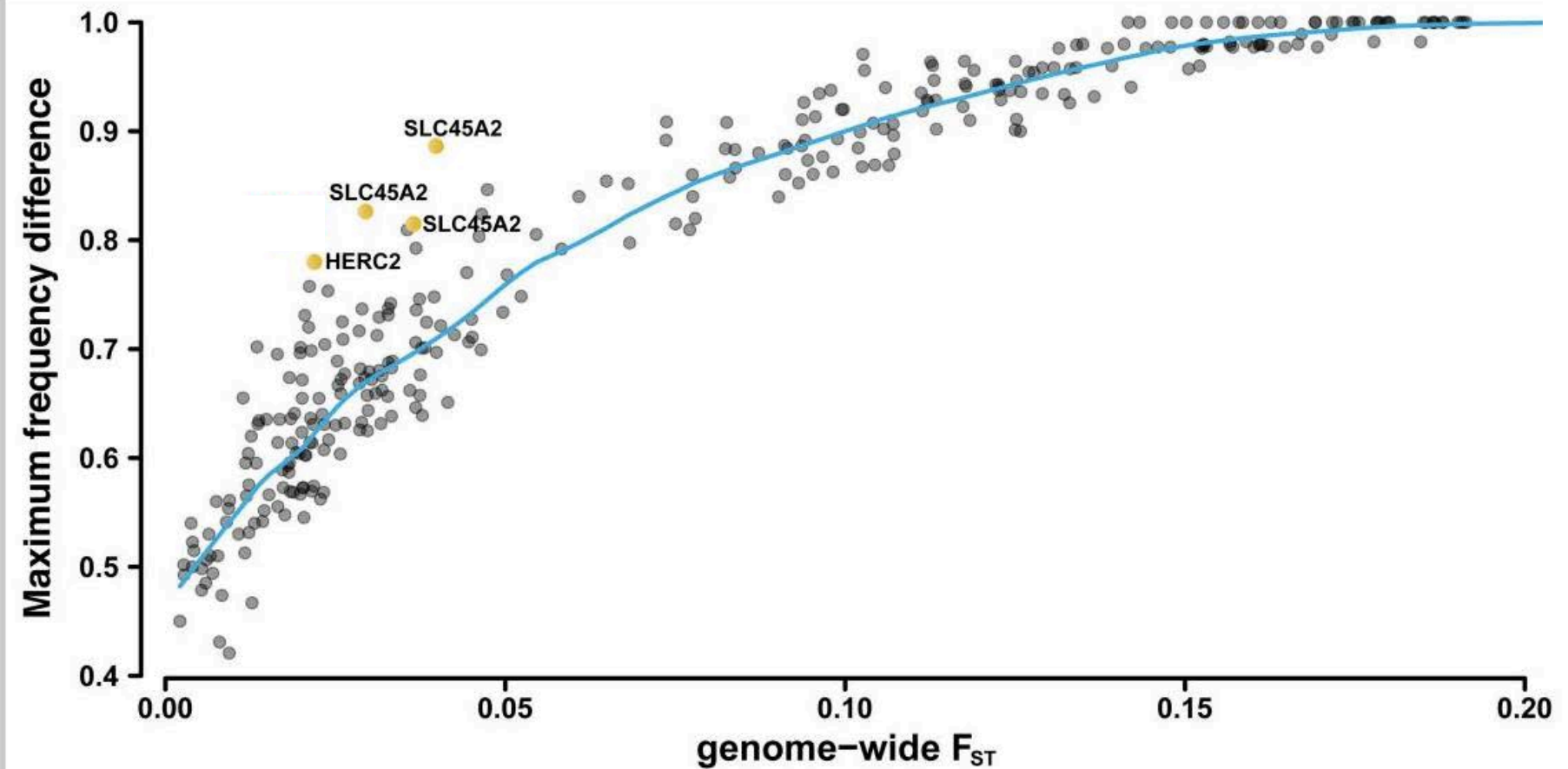
# Positive natural selection

# Positive natural selection

Positive natural selection

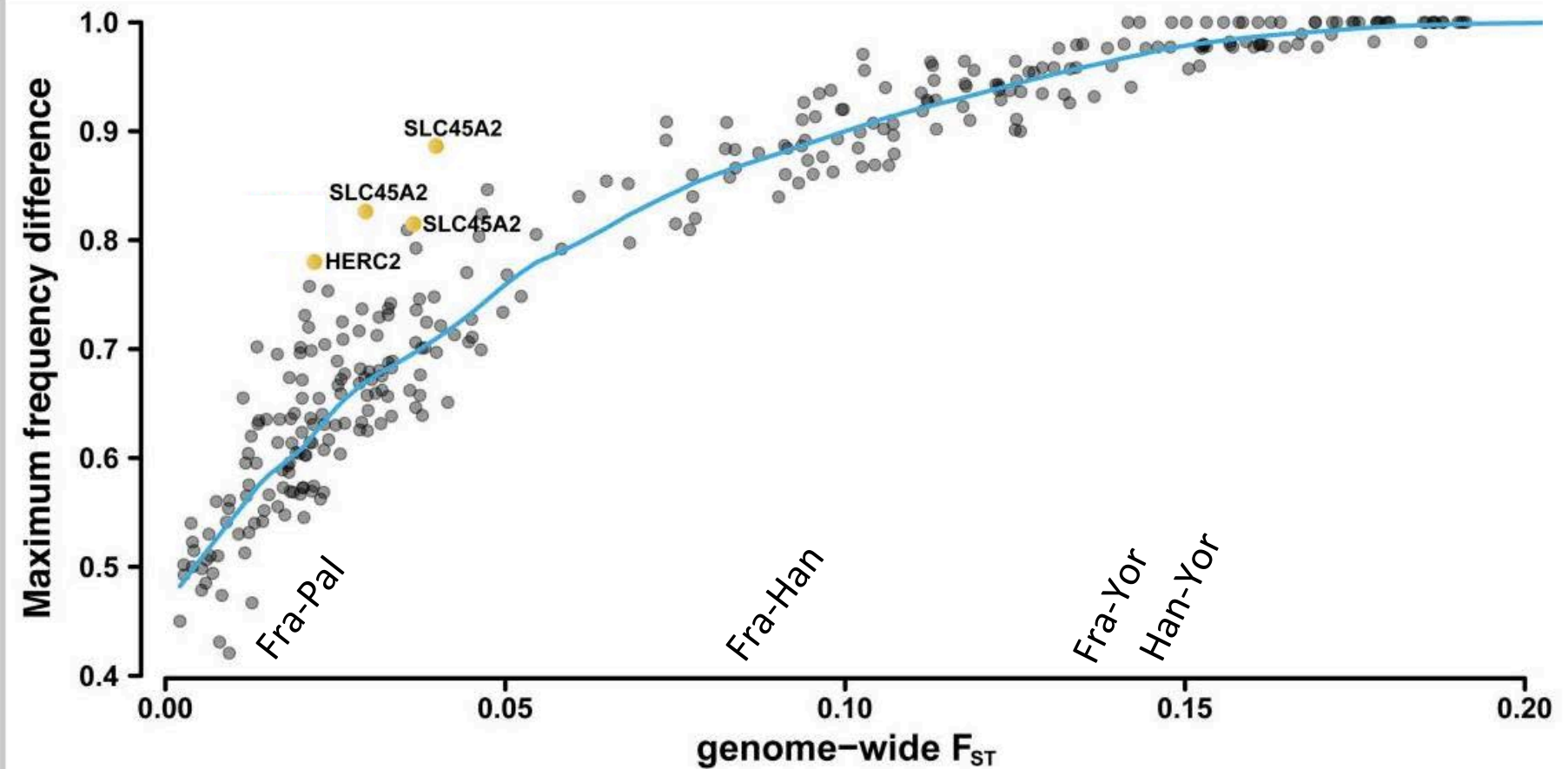# Genetic Differentiation

Redrawn from Coop et al. 2009

# Genetic Differentiation

Redrawn from Coop et al. 2009

# Tons of modern human DNA sequence data

What are the relative contributions of genetic drift and natural selection?
What are the genes under positive selection?

# Tons of modern human DNA sequence data

What are the relative contributions of genetic drift and natural selection?
What are the genes under positive selection?

- Population size has not been constant in time

- The effect of demography can mimic the effect of natural selection

# DNA sequence data

← Sites sequenced →

↑ individuals →
↓ individuals ←

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

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

**Chimp reference**

# DNA sequence data

← Sites sequenced →

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

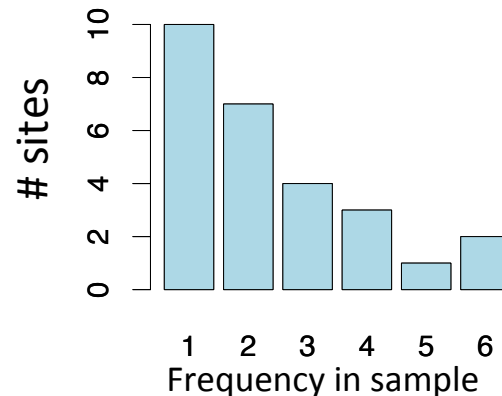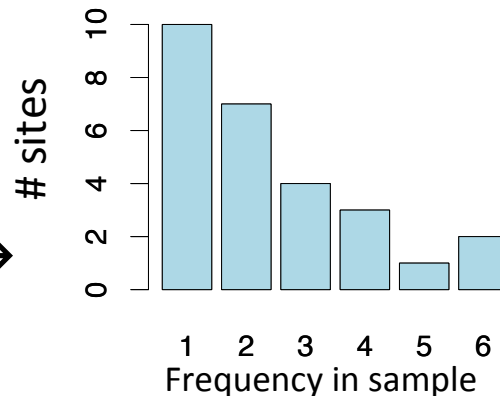0 2 1 1 1 3 3 6 1 1 1 2 1 2 1 4 2 1 2 4 1 3 3 0 2 4 0 0 6 2 0 5

# DNA sequence data

← Sites sequenced →

# DNA sequence data

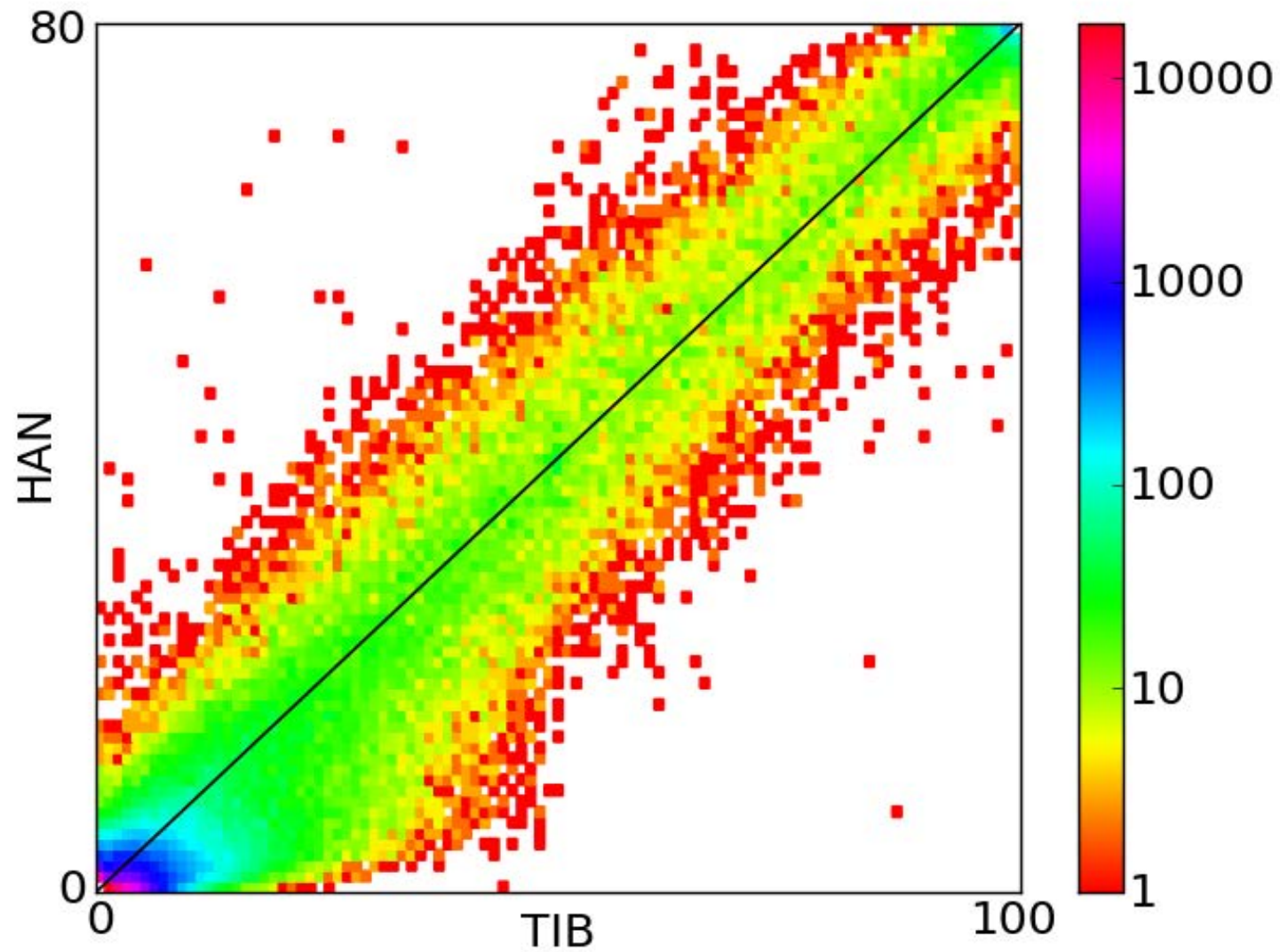← Sites sequenced →

individuals →

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

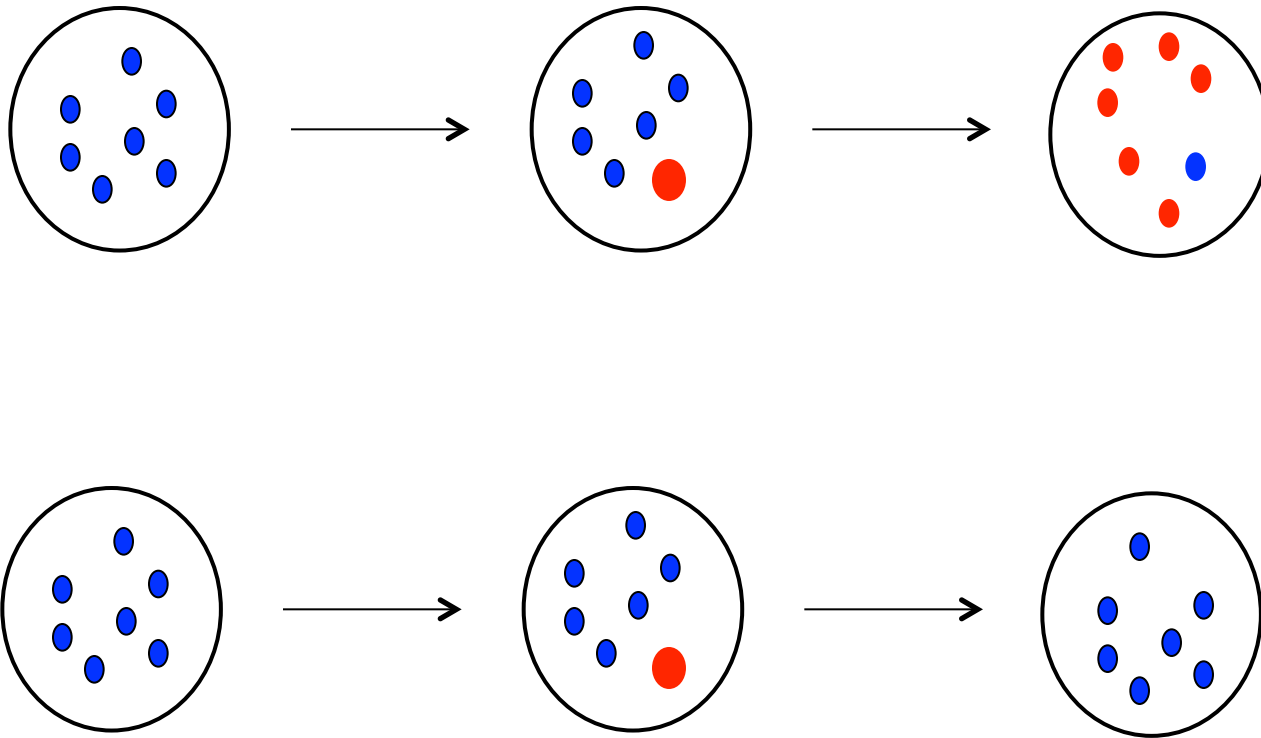0 2 1 1 1 3 3 6 1 1 1 2 1 2 1 4 2 1 2 4 1 3 3 0 2 4 0 0 6 2 0 5
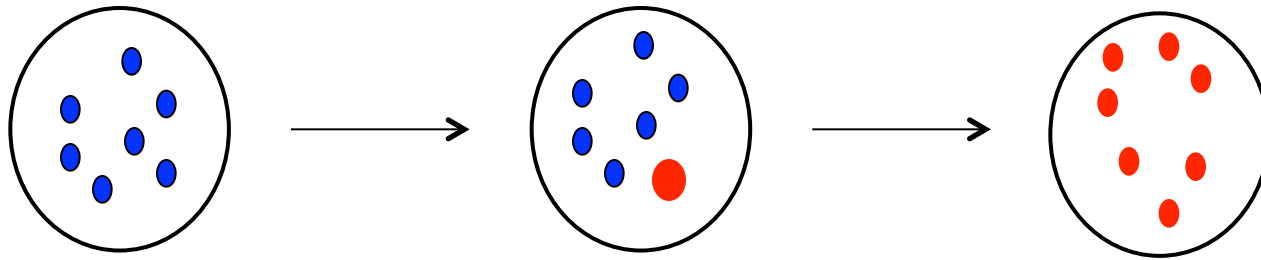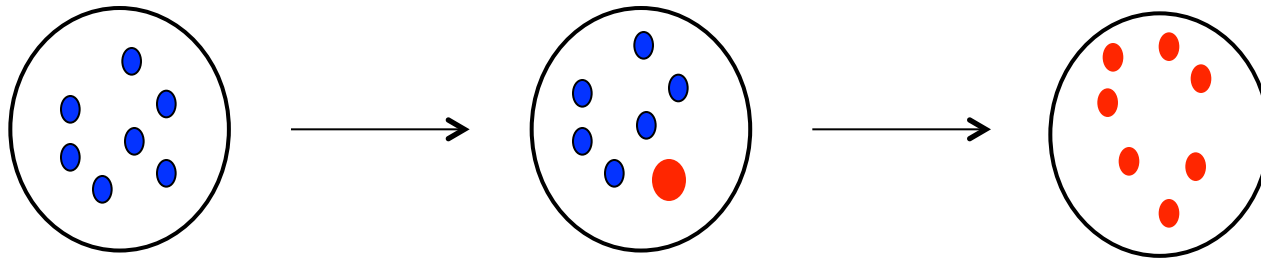
The site frequency spectrum →
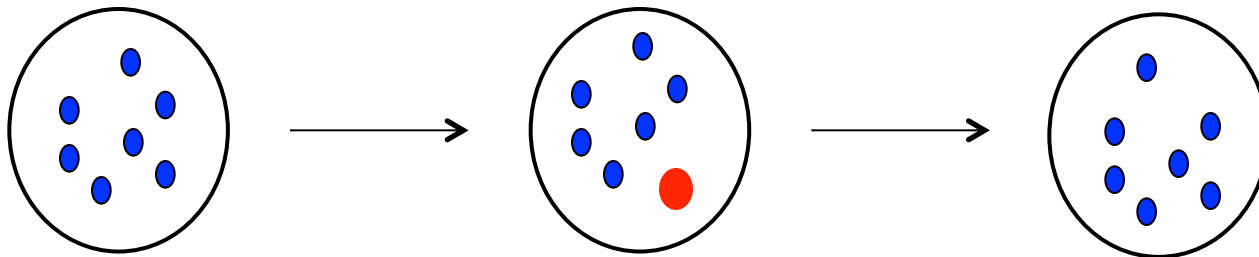
# By chance …

**Positive Selection**

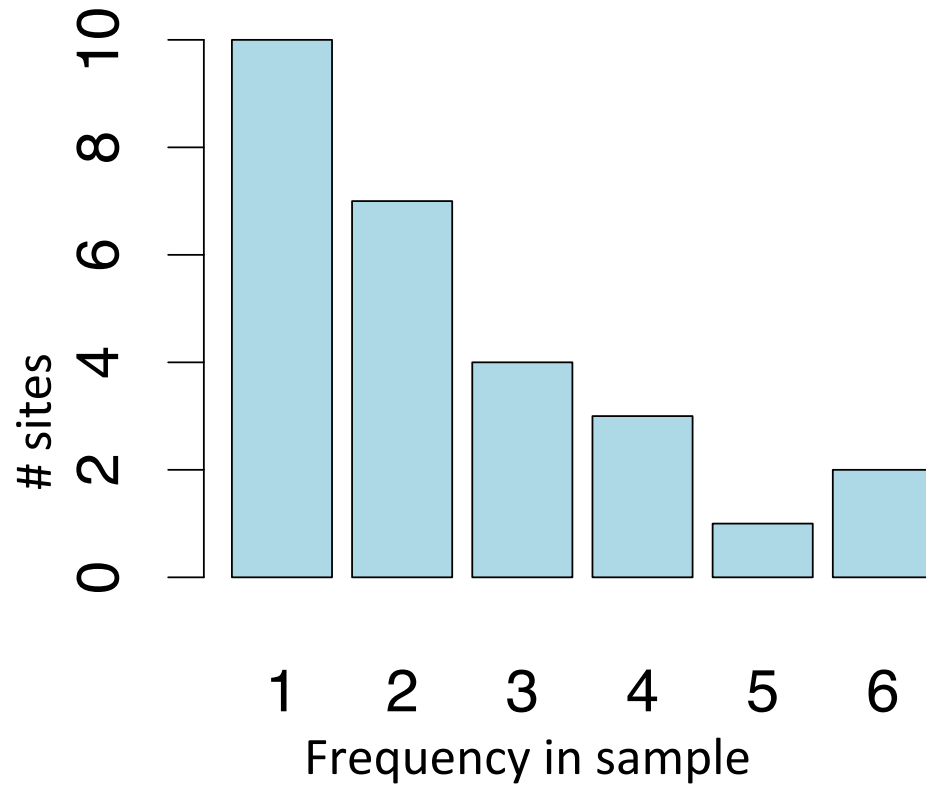**Positive Selection**



**Negative Selection**

# How is the SFS shaped by difference processes?

# The Wright-Fisher model



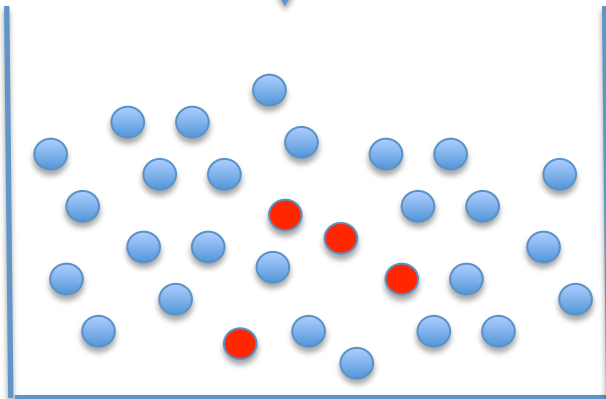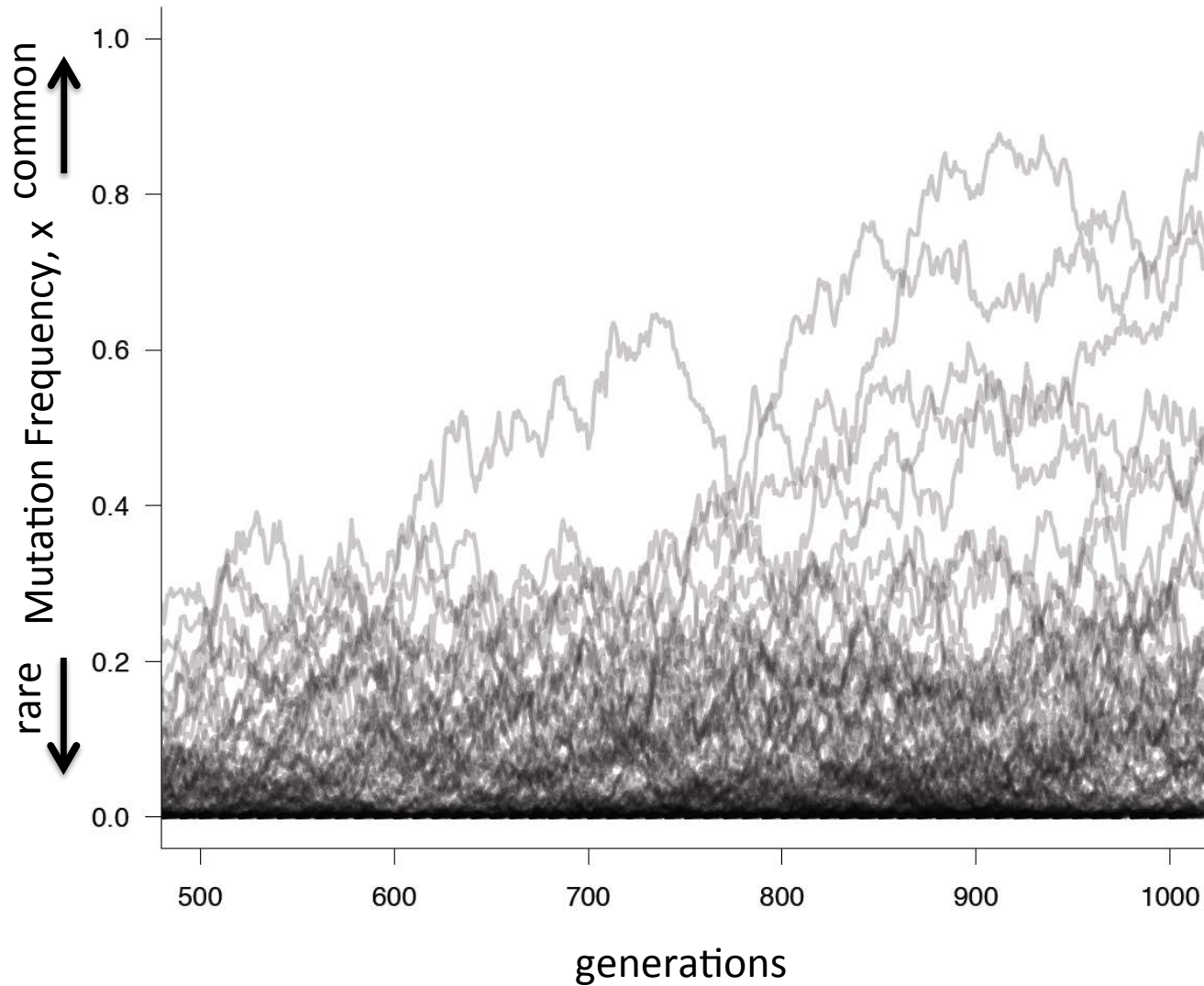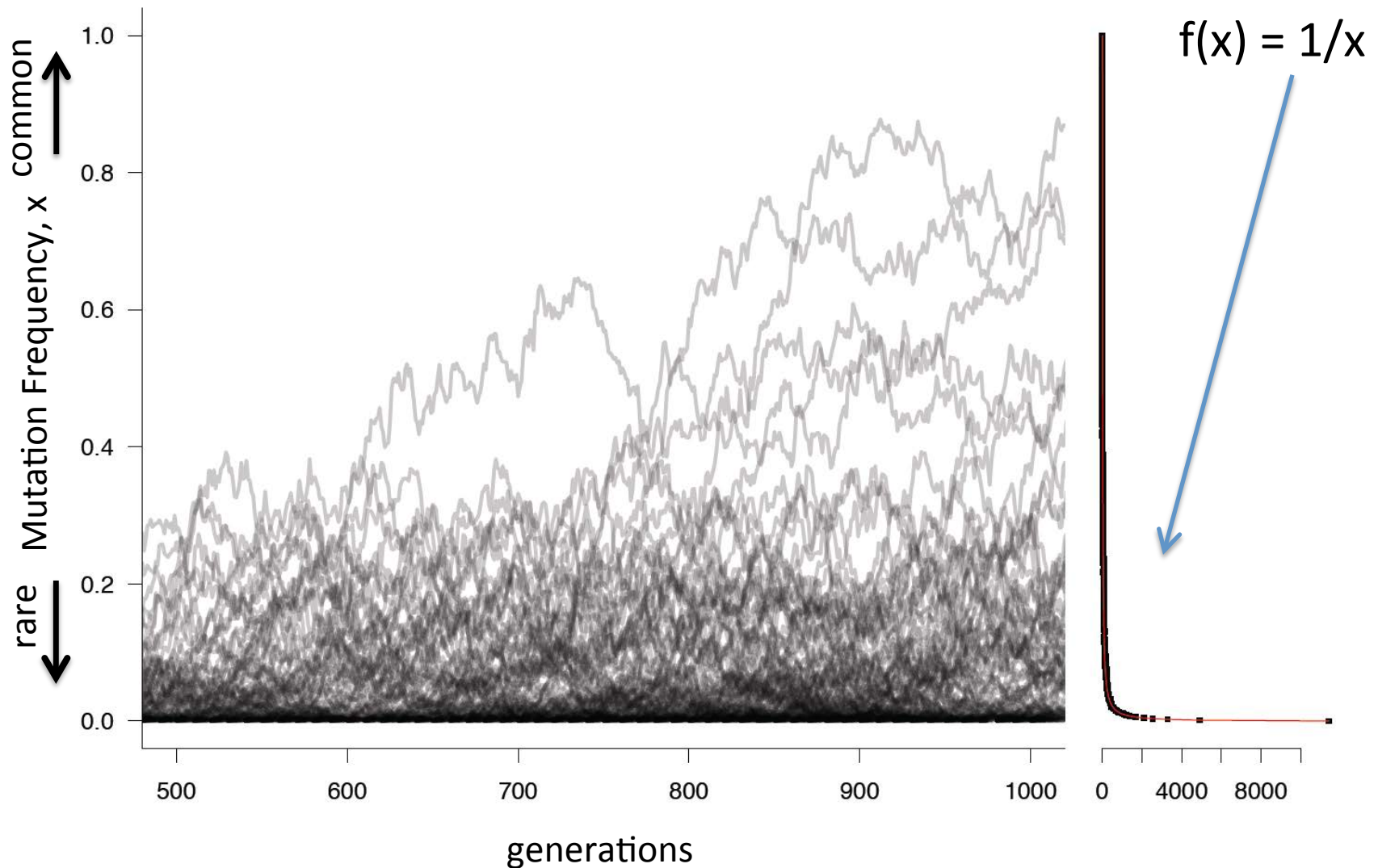**X(t)** = number of individuals carrying mutation at time t

**X(t)** is a **Markov chain** with transition probabilities:

$$P_{ij} = Bin(2N, i/2N)$$

Time t

Time t+1

# If mutations are arriving at Poisson times
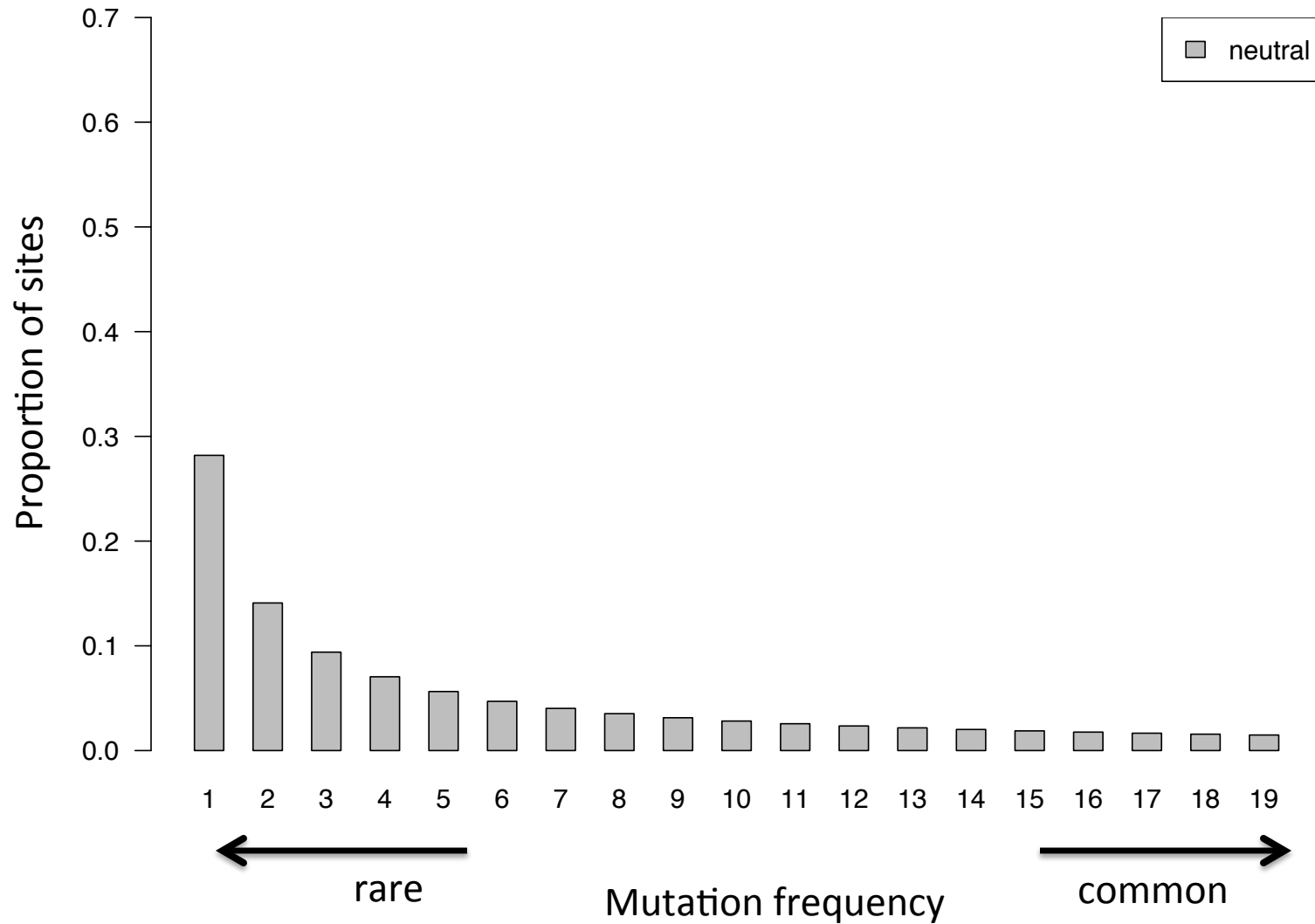
# If mutations are arriving at Poisson times



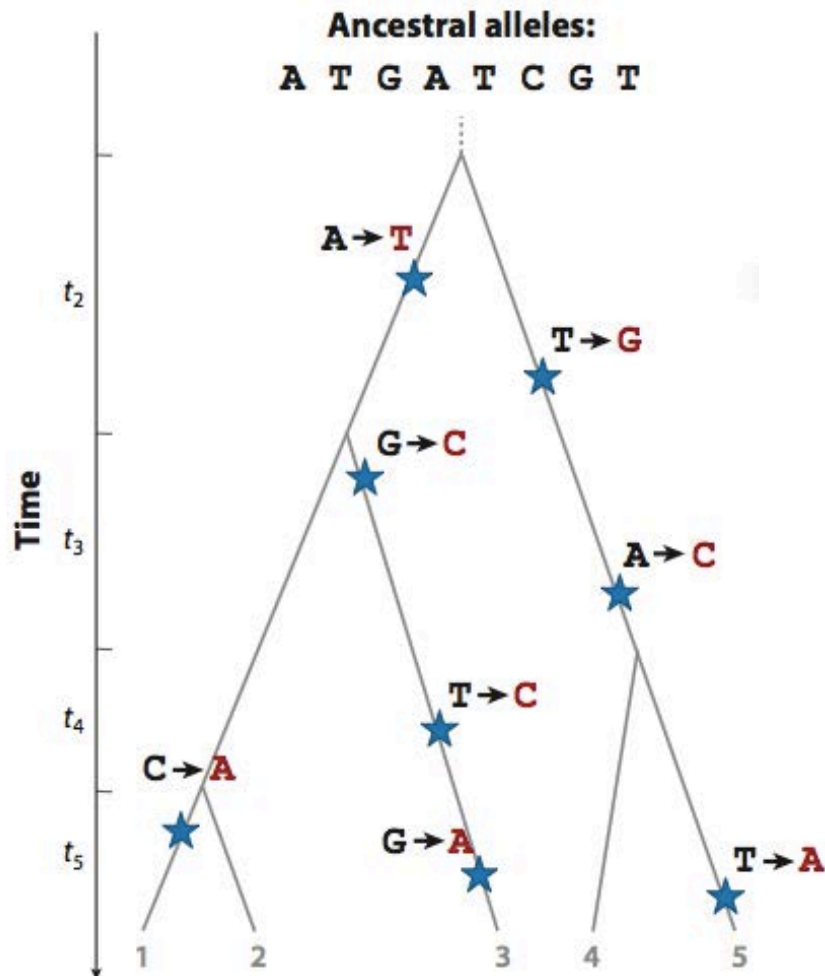rare ← Mutation Frequency, $x$ → common

generations

$f(x) = 1/x$

The expected number of variant sites at frequency *i*

$$\theta F(i, \gamma) = \theta \int_0^1 f(q) \Pr(i \mid q) dq$$

$$= \theta \int_0^1 \frac{1}{q} \binom{n}{i} q^i (1-q)^{n-i} dq$$
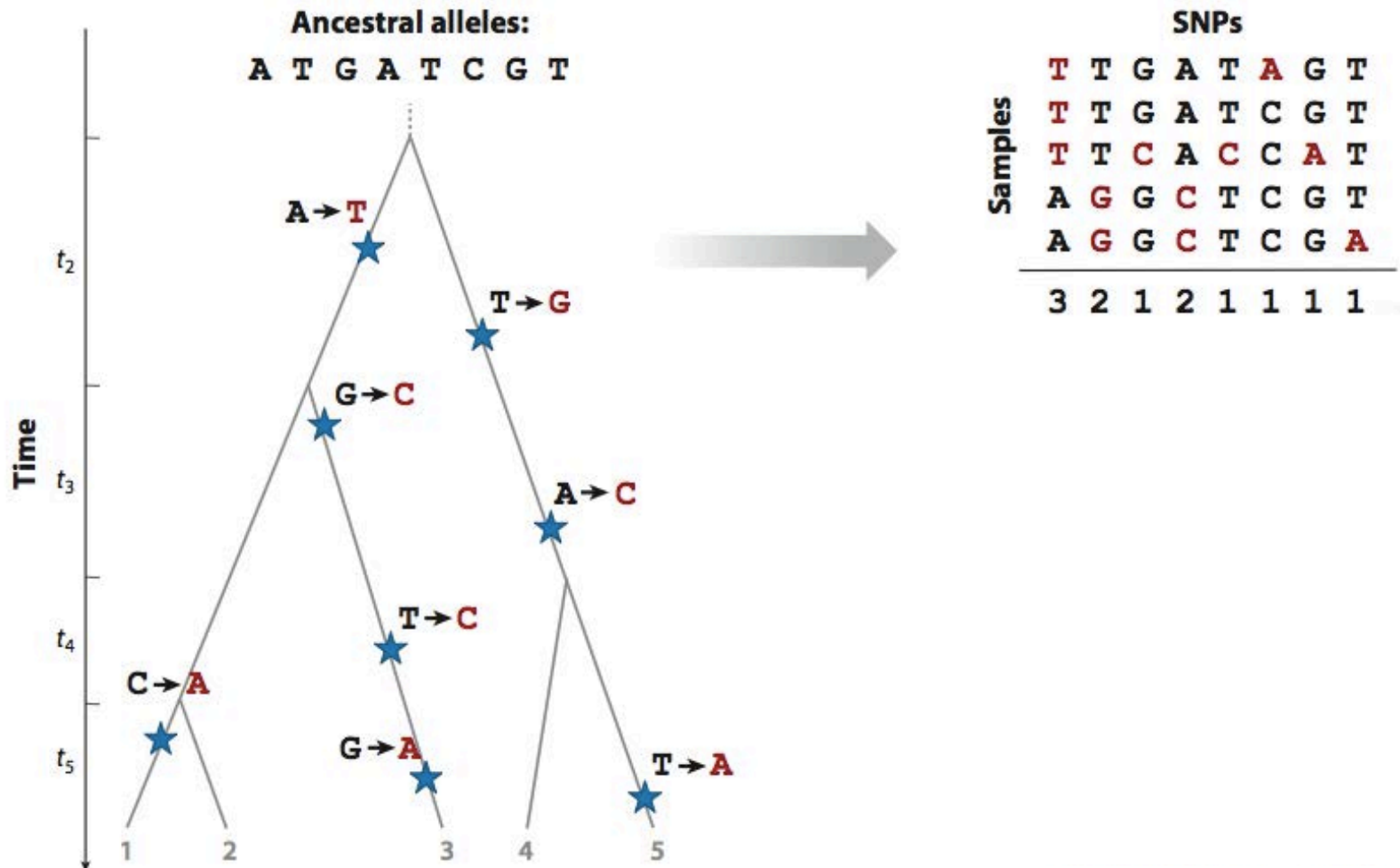
$$= \frac{\theta}{i}$$

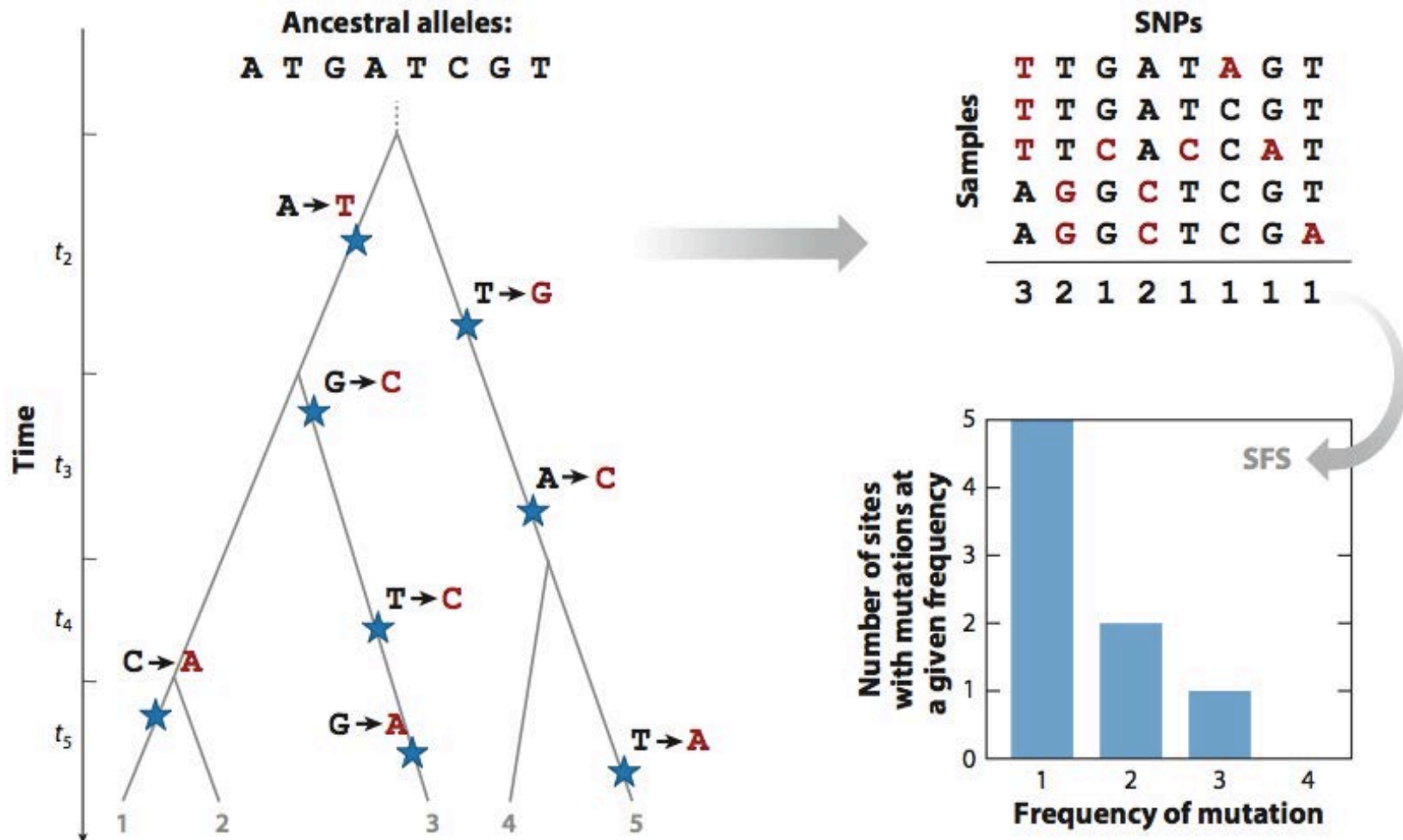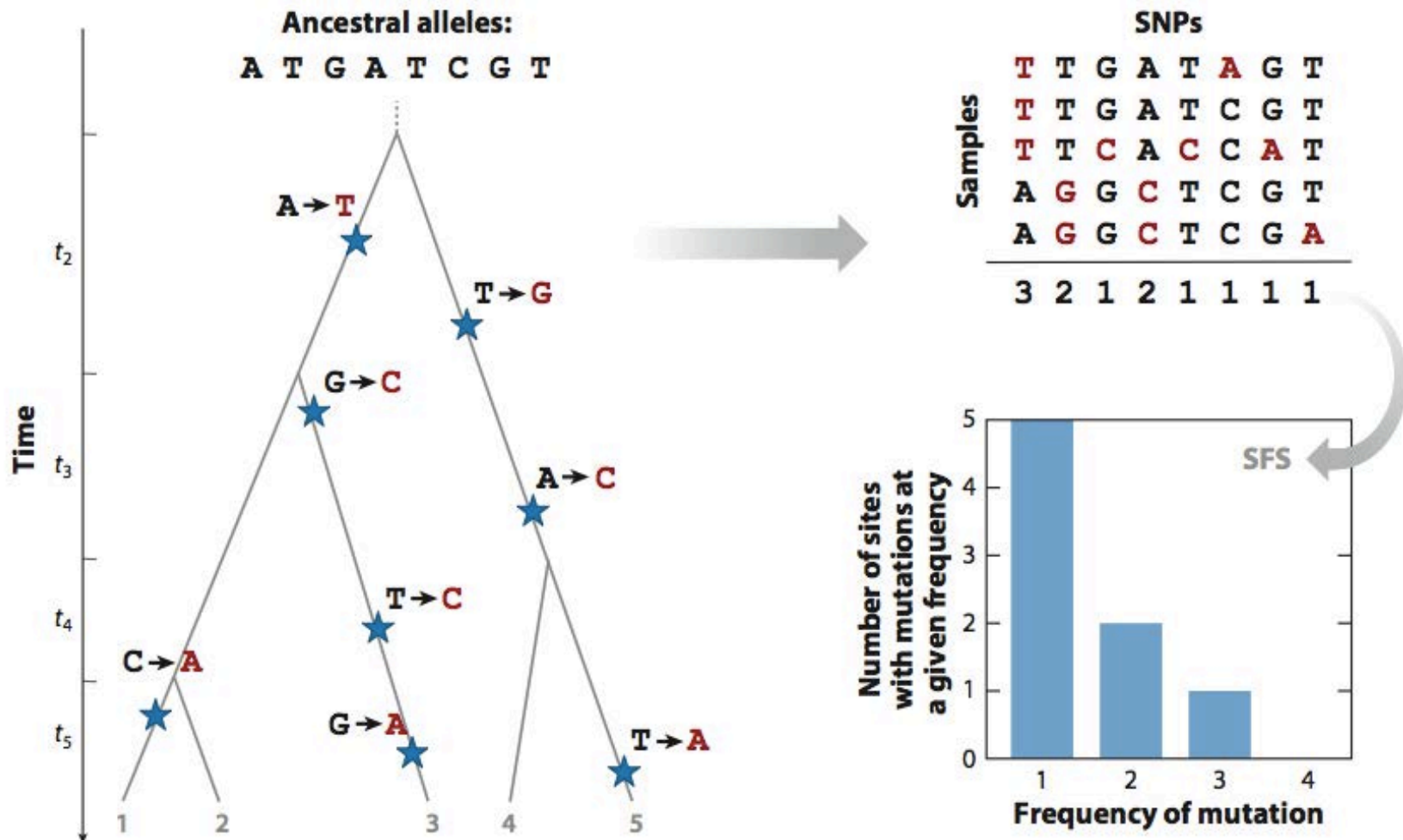# The site frequency spectrum

# The site frequency spectrum from the coalescent

# The site frequency spectrum from the coalescent

# The site frequency spectrum from the coalescent

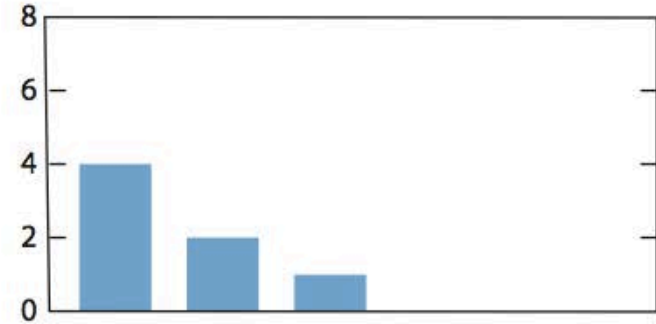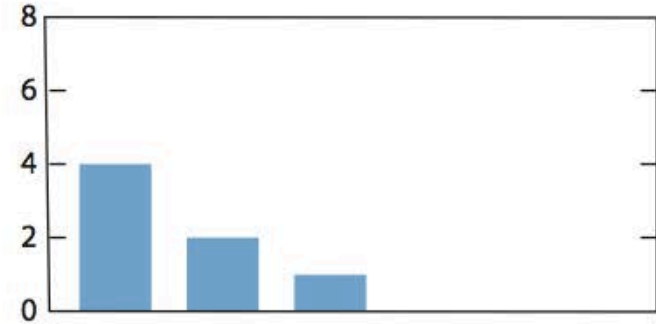Beichman A, Huerta-Sanchez E, Lohmueller K. Annual Review of Ecology, Evolution, and Systematics (2019)

# The site frequency spectrum from the coalescent

# The effect of demography

# The effect of demography

# The effect of demography

Synonymous sites

# The effect of a bottleneck and population structure

**bottleneck**

**Population structure**

Beichman A, Huerta-Sanchez E, Lohmueller K. Annual Review of Ecology, Evolution, and Systematics (2019)

# The effect of demography

# If mutations are arriving at Poisson times



f(x) = 1/x

generations

The expected number of variant sites at frequency *i*

$$\theta F(i, \gamma) \quad = \int_0^1 \frac{f(q,\gamma)}{2} \cdot \Pr(i \mid q)\, dq$$

$$= \int_0^1 \frac{1 - e^{-2\gamma(1-q)}}{1 - e^{-2\gamma}} \frac{1}{q(1-q)} \binom{n}{i} q^i (1-q)^{n-i} dq$$

The expected number of variant sites at frequency $i$

$$\theta F\left(i,\gamma\right) \quad = \int_0^1 \frac{f(q,\gamma)}{2} \cdot \Pr\left(i \mid q\right) dq$$

$$= \int_0^1 \frac{1-e^{-2\gamma(1-q)}}{1-e^{-2\gamma}} \frac{1}{q(1-q)} \binom{n}{i} q^i(1-q)^{n-i} dq$$

The probability of seeing $x_i$ sites at frequency $i$

$$p\left(X_i = x_i \mid \theta, \gamma\right) = e^{-\theta F(i,\gamma)} \frac{\left(\theta F\left(i,\gamma\right)\right)^{x_i}}{x_i!}.$$

The expected number of variant sites at frequency $i$

$$\theta F(i, \gamma) = \int_0^1 \frac{f(q,\gamma)}{2} \cdot \Pr(i \mid q) \, dq$$

$$= \int_0^1 \frac{1 - e^{-2\gamma(1-q)}}{1 - e^{-2\gamma}} \frac{1}{q(1-q)} \binom{n}{i} q^i (1-q)^{n-i} dq$$

The probability of seeing $x_i$ sites at frequency $i$

$$p(X_i = x_i \mid \theta, \gamma) = e^{-\theta F(i,\gamma)} \frac{(\theta F(i, \gamma))^{x_i}}{x_i!}.$$

The likelihood function

$$L_u(\theta, \gamma \mid x) = \prod_{i=1}^{n-1} e^{-\theta F(i,y)} \frac{(\theta F(i, \gamma))^{x_i}}{x_i!}$$

# If mutations are arriving at Poisson times



f(x) = 1/x

# The site frequency spectrum

# The site frequency spectrum

# The site frequency spectrum

# The site frequency spectrum



**Legend:**
- neutral (grey)
- positive selection (red)
- negative selection (blue)

Y-axis: Proportion of sites (0.0 to 0.7)

X-axis: Mutation frequency (1 to 19)

rare ← → common

# The site frequency spectrum



Legend:
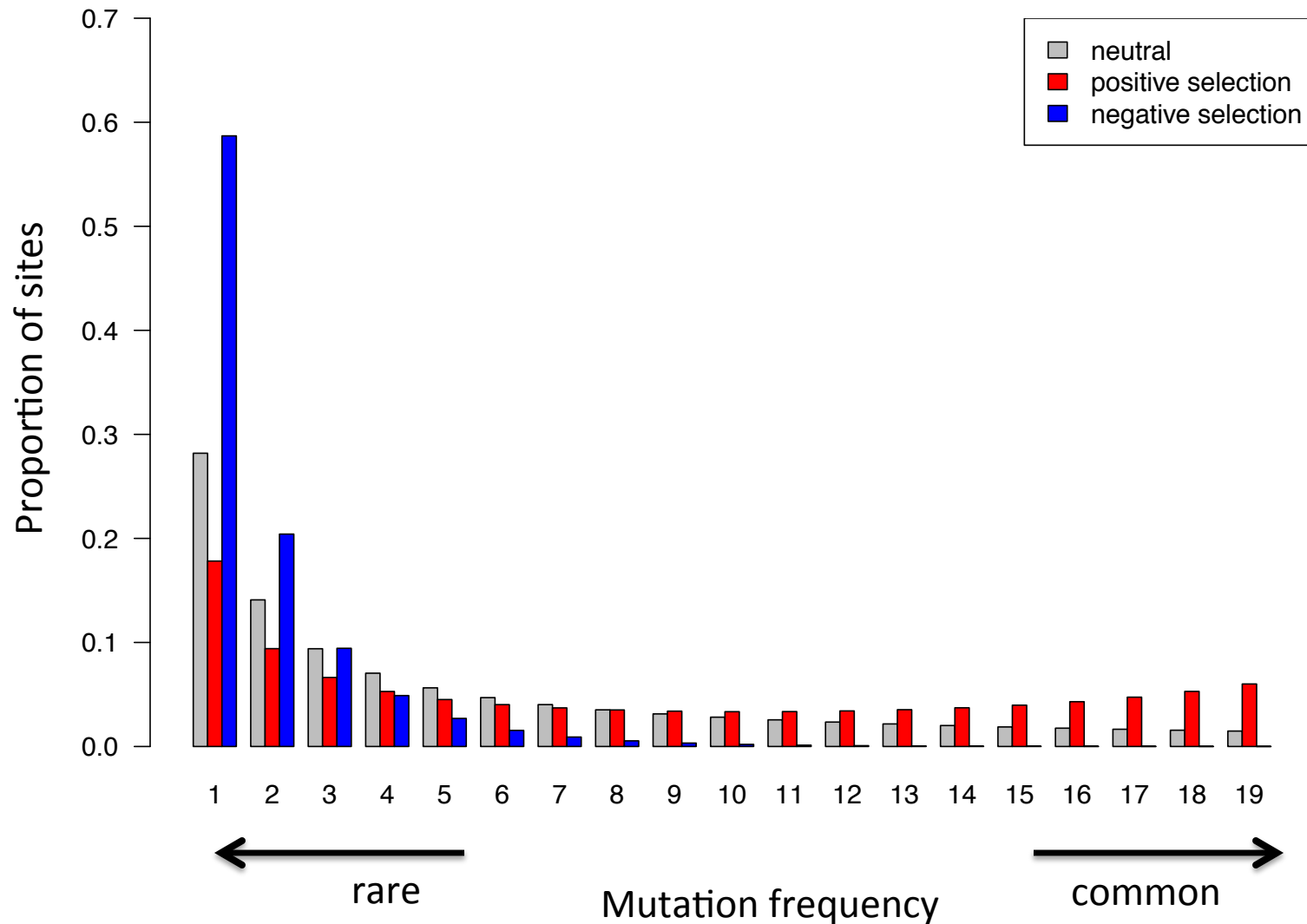- neutral
- positive selection
- negative selection
- fluctuating selection

$$f(x,\beta) = \frac{2}{K(\beta)x(1-x)} \log\left(\frac{1-r_1(\beta)}{x-r_1(\beta)} \bullet \frac{r_2(\beta)-x}{r_2(\beta)-1}\right)$$

Proportion of sites

rare      mutation frequency      common

Huerta-Sanchez et al. (2008)

200 individuals of Danish nationality

# An excess of rare mutations



Li* Y. et al. (2010) Nature Genetics

# An excess of rare mutations

# An excess of rare mutations



**Negative Selection**

Li* Y. et al. (2010) Nature Genetics

# An excess of rare deleterious mutations

# X chromosome



Li* Y. et al. (2010) Nature Genetics

# Distribution of selective effects of non-synonymous mutations

$$f(x) = k\frac{1}{x} + (1-k)\int f_{sel}(x,-\gamma)\,Gamma(\gamma;\alpha,\beta)\,d\gamma$$

A proportion *k* of mutations are neutral

A proportion *(1-k)* are deleterious

Li* Y. et al. (2010) Nature Genetics

# A larger proportion of weakly deleterious mutations



**Estimated gamma densities**

Boyko et al. (2008)
Plos Genetics

Li* Y. et al. (2010) Nature Genetics

$$\tau = t/2N_A$$

$$N_1 = \nu_1 N_A$$

$$\theta = 4N_A\mu$$

$$M = 2N_A m$$

$$N_2 = \nu_2 N_A$$

# Demographic inference



$$\tau = t/2N_A$$

$$N_1 = \nu_1 N_A$$

$$\theta = 4N_A\mu$$

$$M = 2N_A m$$

$$N_2 = \nu_2 N_A$$

$N_A$  $N_{AF}$

$m_{AF\text{-}B}$  $m_{AF\text{-}EU}$

$T_{AF}$  $T_B$  $N_B$  $N_{EU0}$  $T_{EU}$

$m_{EU\text{-}AS}$

$N_{AS0}$  $r_{AS}$

$T_{EU\text{-}AS}$

(also $m_{AF\text{-}AS}$)

YRI

CEU

CHB

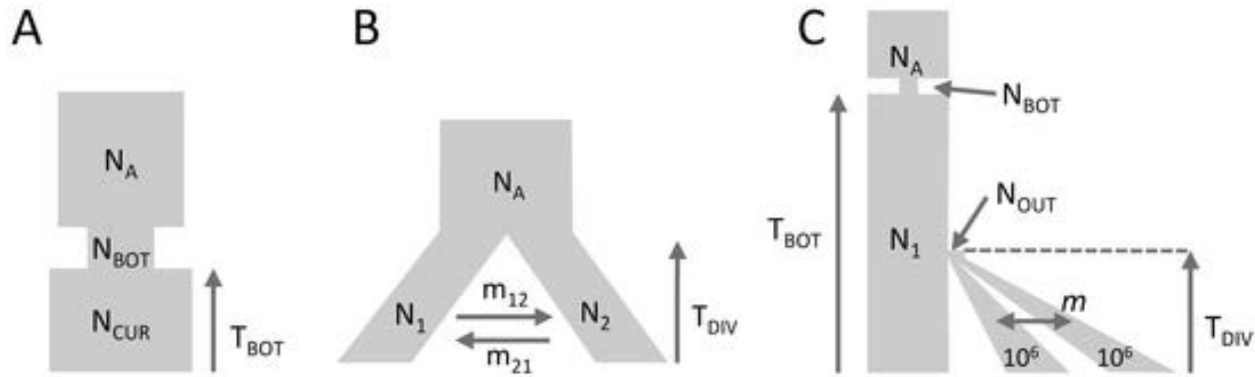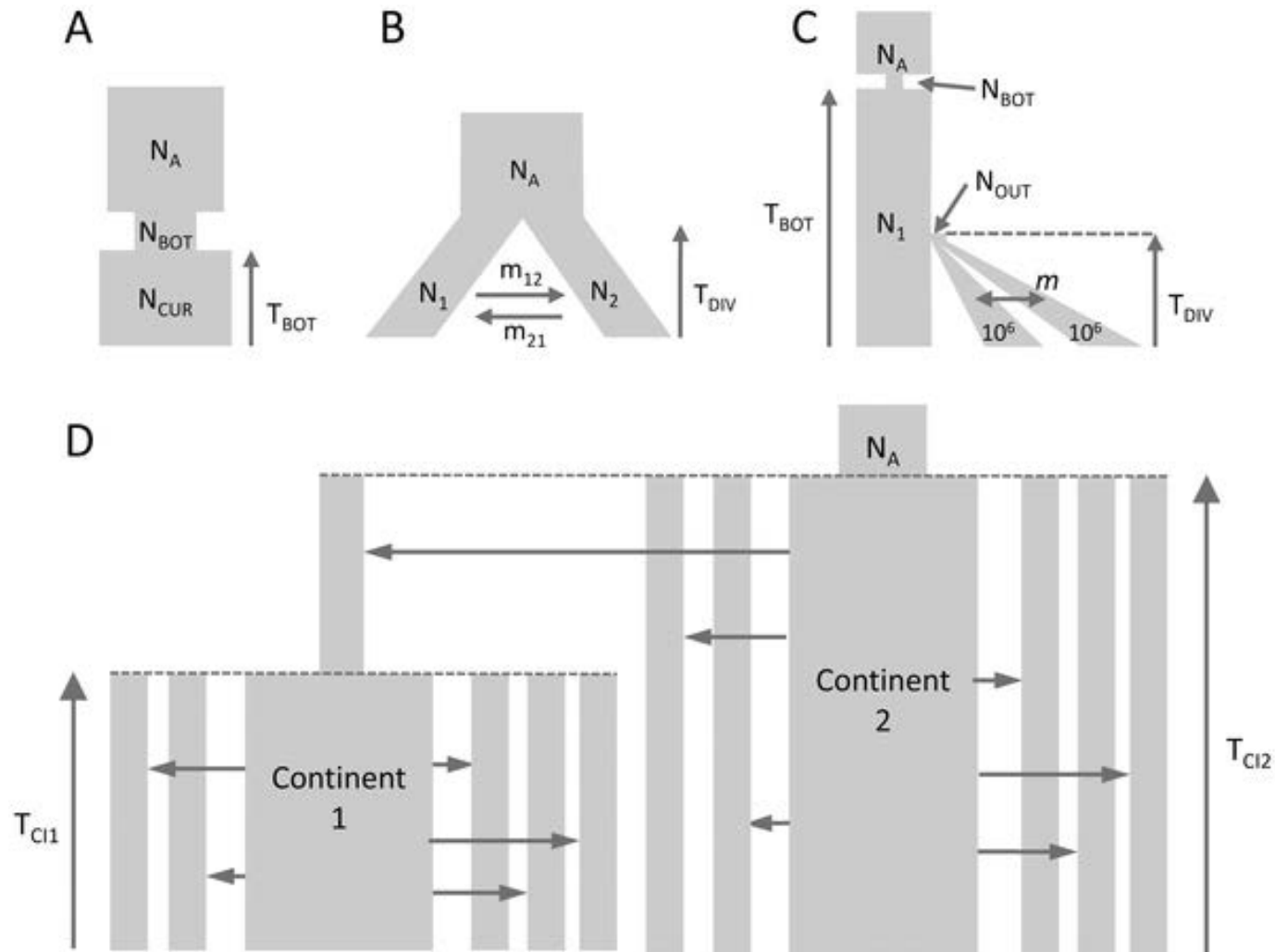$$\frac{d}{dt}f(q,t;\Theta) = \frac{1}{2}\frac{d^2}{dq^2}\{V(q;\Theta)f(q,t;\Theta)\} - \frac{d}{dq}\{M(q;\Theta)f(q,t;\Theta)\}$$
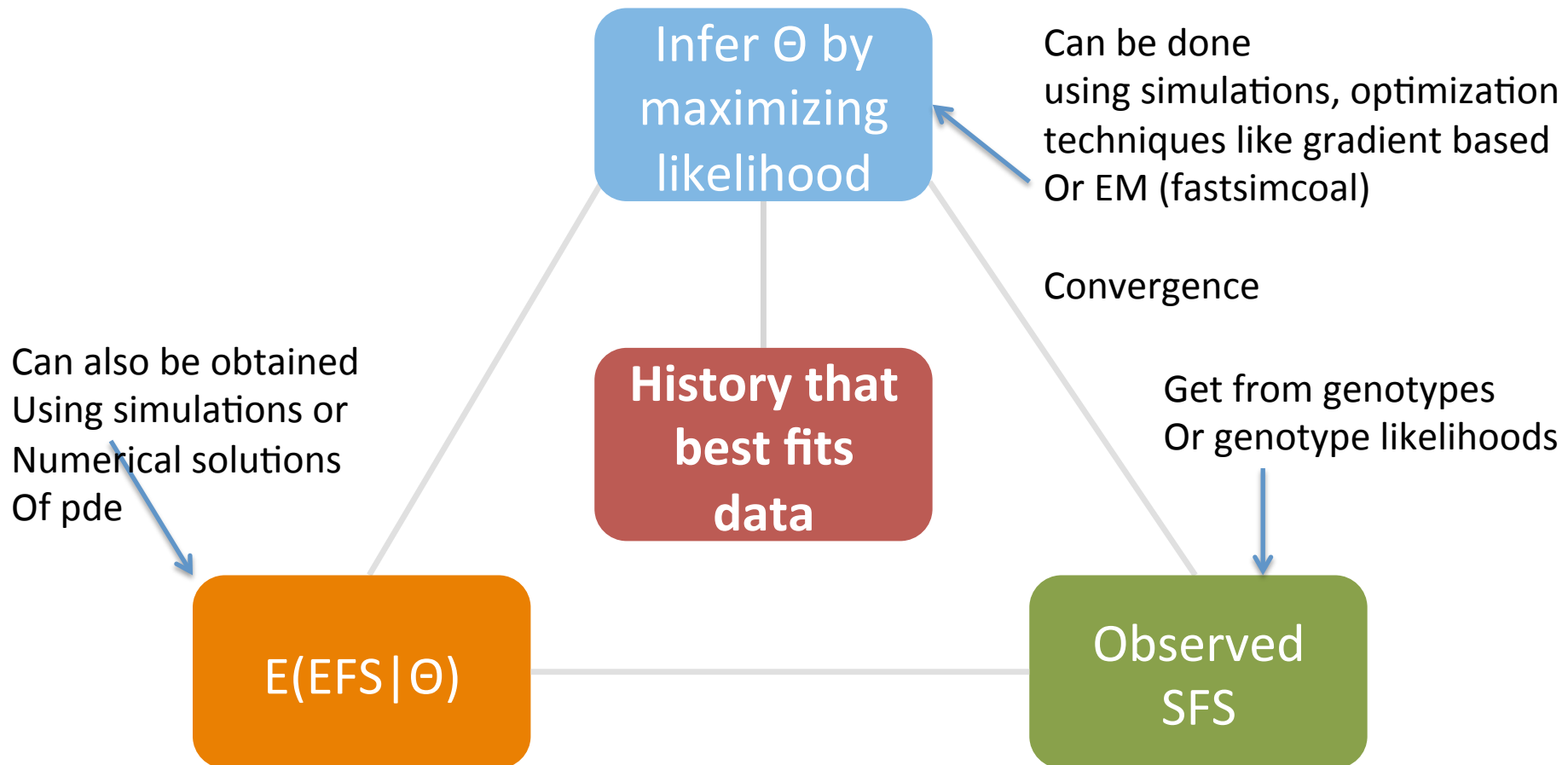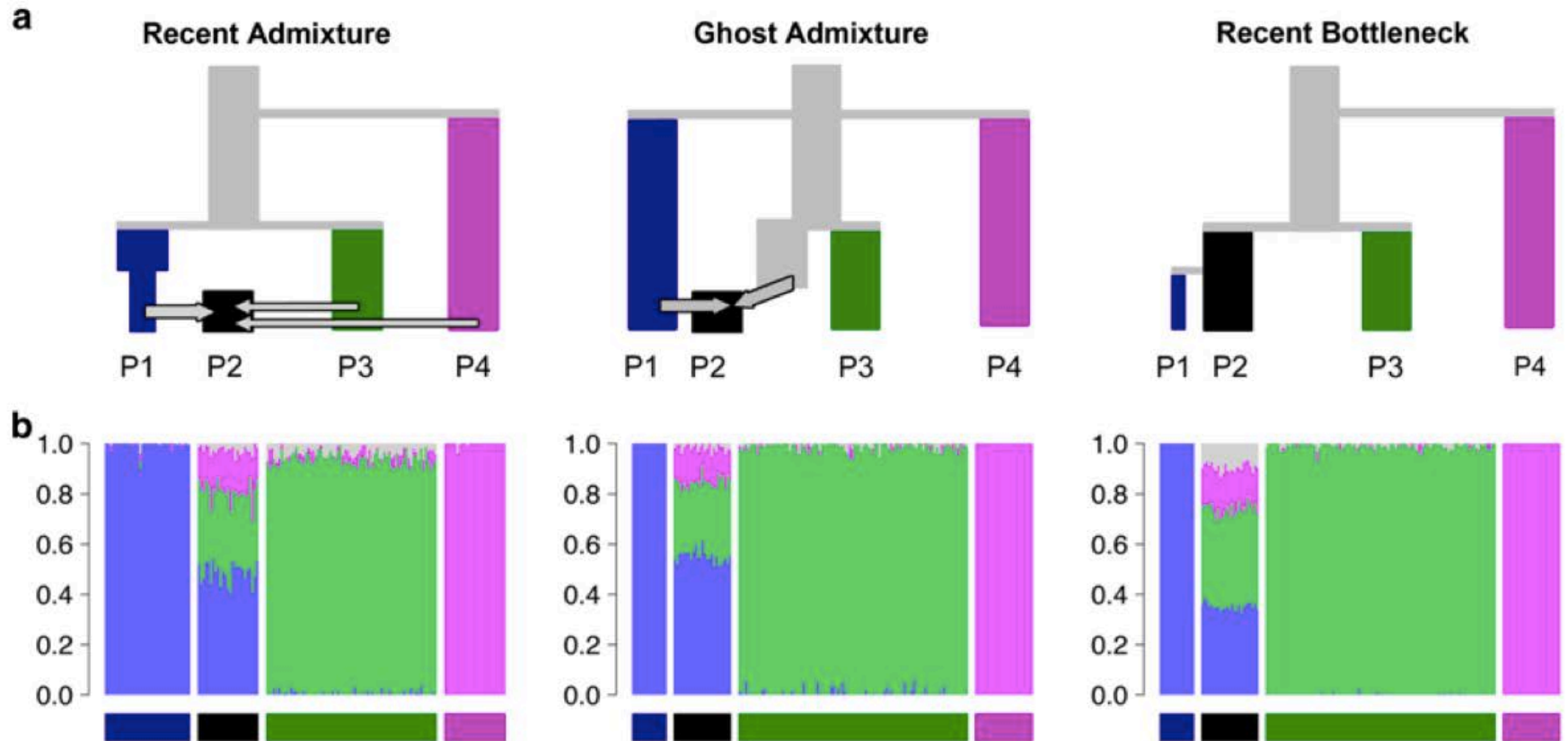
Fastsimcoalsim
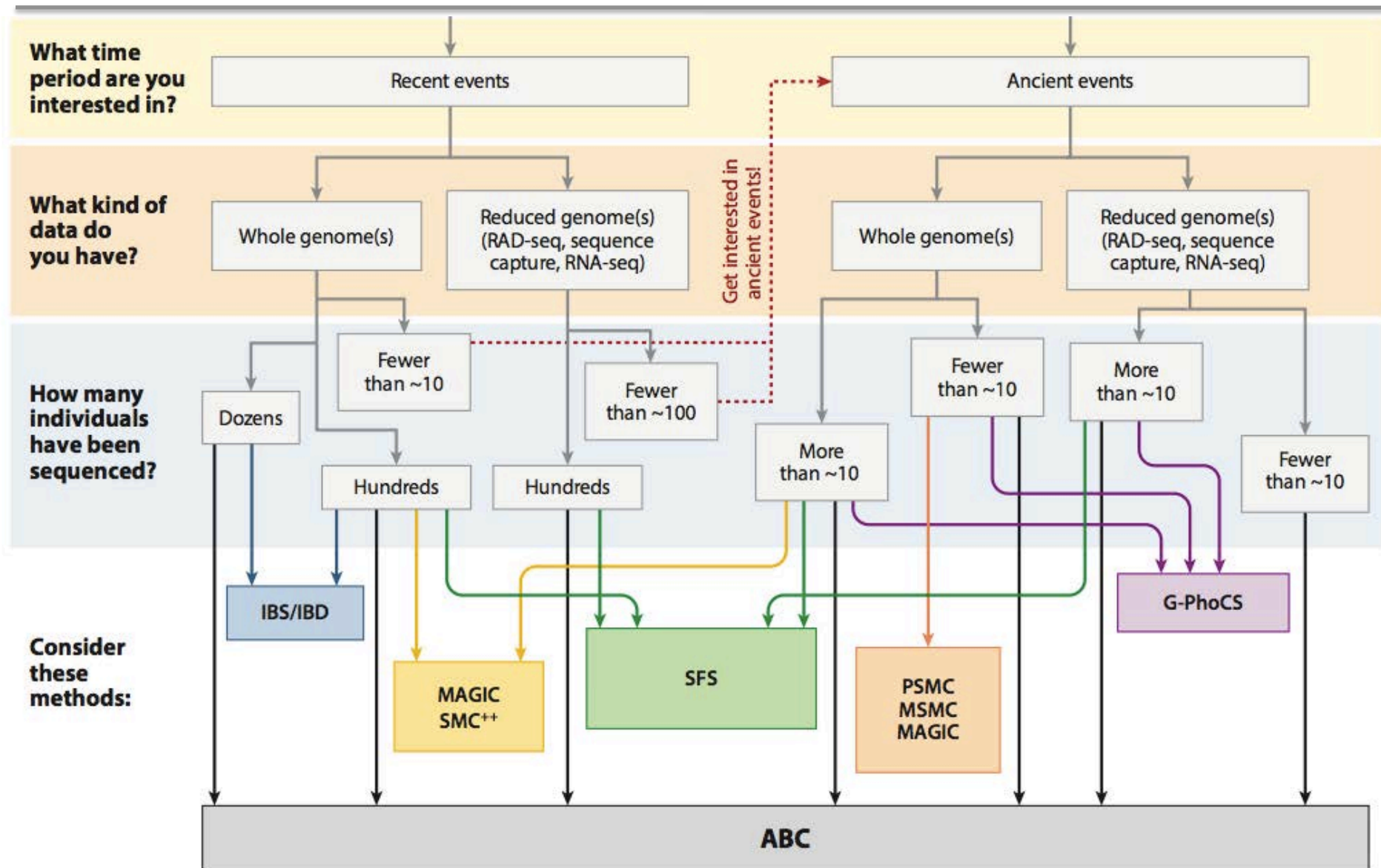Excoffier et al. 2013

# Demographic inference with dadi of fastsimcoal

# Benefits of model-based approach over Structure



A tutorial on how not to over-interpret STRUCTURE and ADMIXTURE bar plots

Daniel J. Lawson, Lucy van Dorp & Daniel Falush ✉

Decision Tree from Beichman et al. (2018) Annual Review of Ecology, Evolution, and Systematics

- When the population size grows, genetic drift slows down and selection becomes more efficient

- A population bottleneck speeds up genetic drift and makes natural selection get less efficient

- Population size changes shift the site frequency spectrum away from the shape SFS($n$)~1/$n$ that is expected under neutrality

- Natural selection also shifts the SFS away from the equilibrium 1/$n$ shape

- The SFS shape can be used to test for selection, but only after correcting for nonequilibrium demography