

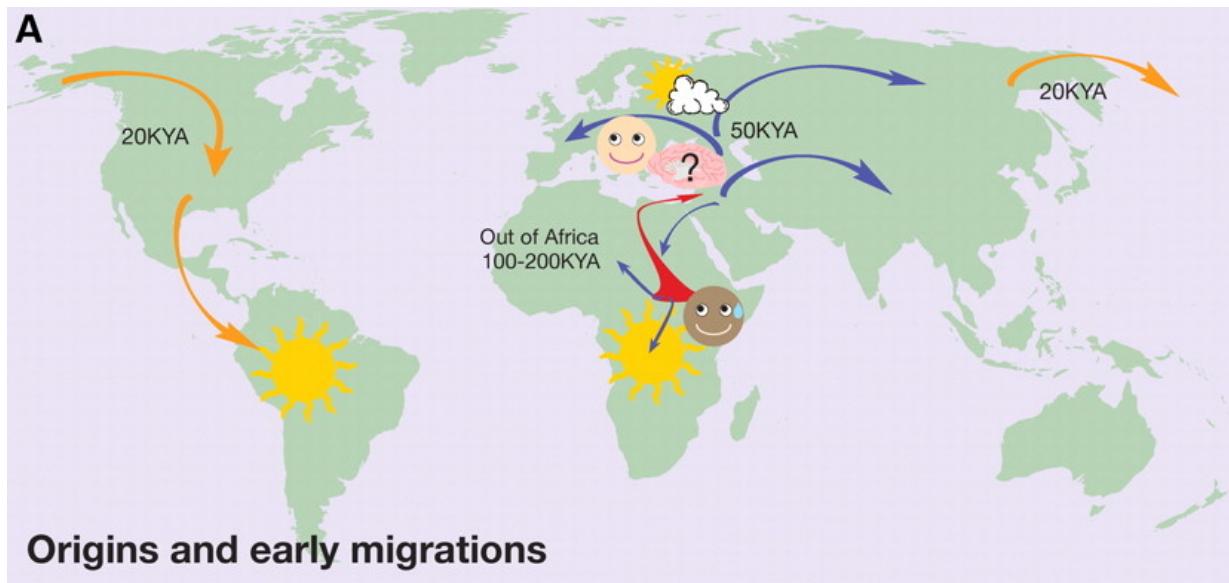
Evolutionary genomics

Data analysis module - Day 4

Detecting signatures of natural selection

April 16th 2015

Evolution and adaptation

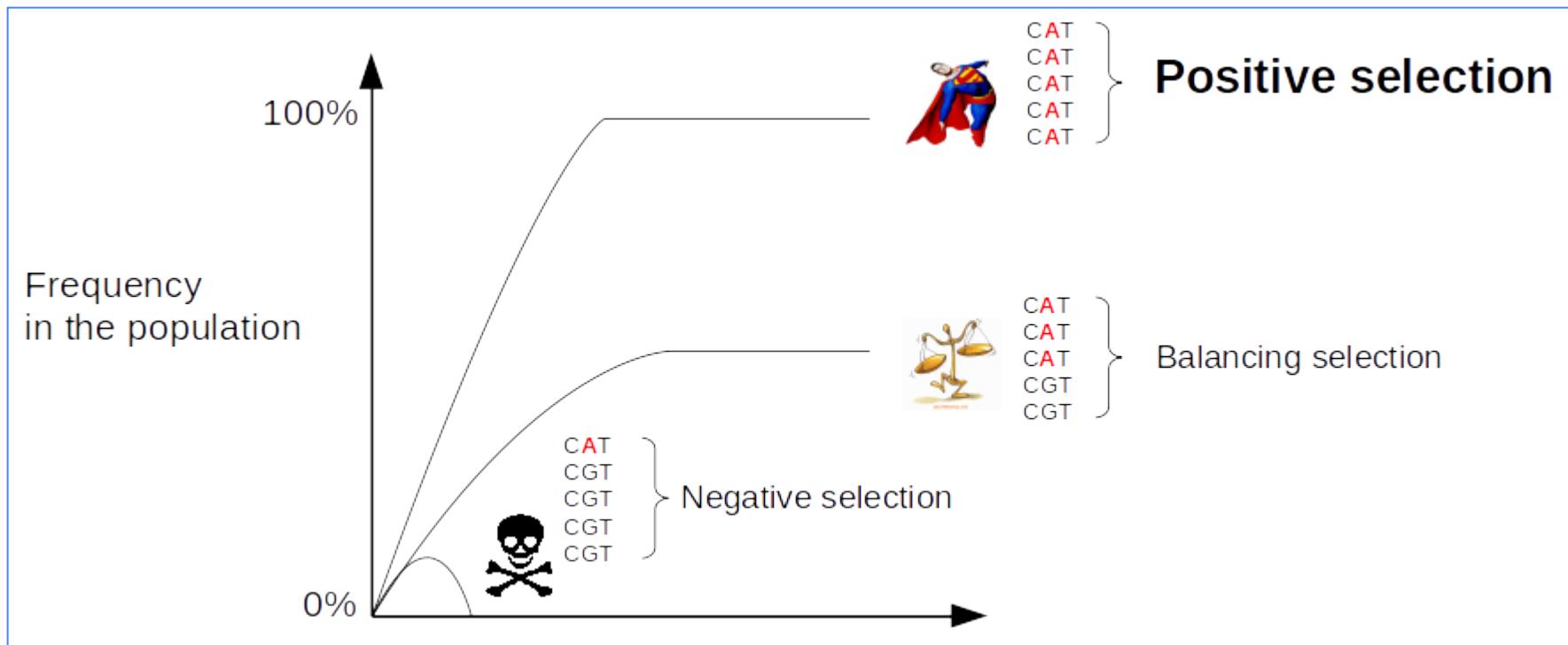


- After the dispersal out-of-Africa, colonisation of a wide range of environments
- Genetic adaptation to:
- Climate (e.g. UV light exposure)
- Diet (e.g. milk consumption)
- Pathogens

Natural selection

Heritable traits that increase the fitness of the become more common

- 1) Mutations arise randomly and evolve according to their effect on the fitness of the carrier



- 2) Sites targeted by natural selection are likely to harbour important functionality

Relevance to human health

- The presence of a **functional** variant is the prerequisite for selection.
- These variants might have played some adaptive role and have a role in human **disease**/condition.
- Recent studies have identified genes related to skin pigmentation, metabolic traits or immune defense targeted by natural selection in the human genome.

OPEN  ACCESS Freely available online

PLOS GENETICS

Differential Susceptibility to Hypertension Is Due to Selection during the Out-of-Africa Expansion

J. Hunter Young¹, Yen-Pei C. Chang¹, James Dae-Ok Kim¹, Jean-Paul Chretien¹, Michael J. Klag¹, Michael A. Levine², Christopher B. Ruff¹, Nae-Yuh Wang¹, Aravinda Chakravarti¹

Hypertension: variants that allowed water and sodium retention and increased vascular reactivity now cause hypertension.

Obesity and insulin resistance: thrifty variants became unfavorable with a shift in diet.

OPEN  ACCESS Freely available online

PLOS GENETICS

Adaptations to Climate in Candidate Genes for Common Metabolic Disorders

Angela M. Hancock¹, David B. Witonsky¹, Adam S. Gordon¹, Gidon Eshel², Jonathan K. Pritchard¹, Graham Coop¹, Anna Di Rienzo^{1*}

Published May 25, 2009

JEM

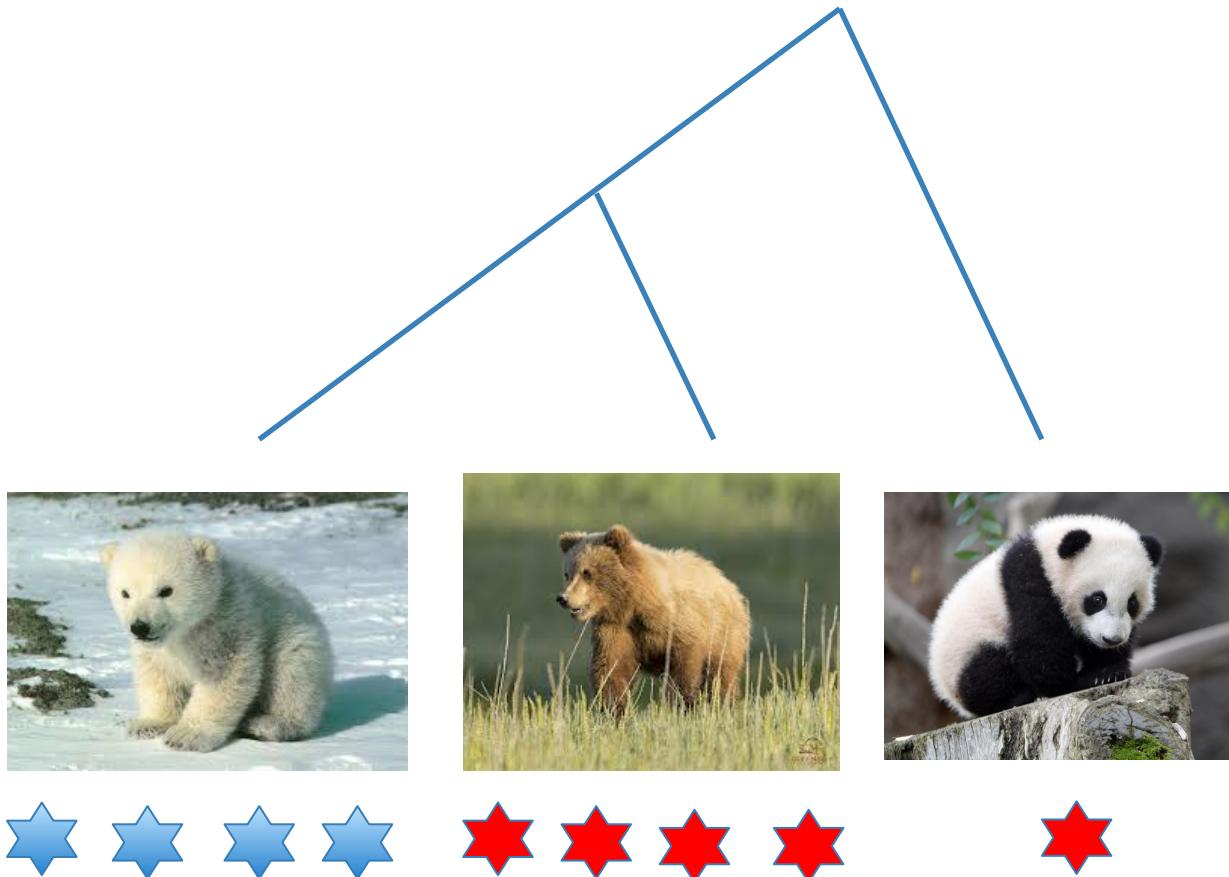
Parasites represent a major selective force for interleukin genes and shape the genetic predisposition to autoimmune conditions

Matteo Fumagalli,^{1,2} Uberto Pozzoli,¹ Rachele Cagliani,¹ Giacomo P. Comi,³ Stefania Riva,¹ Mario Clerici,^{4,5} Nereo Bresolin,^{1,3} and Manuela Sironi¹

Hygiene hypothesis: lack of exposure to pathogens in early life determines immune imbalances and predisposes to atopic and autoimmune diseases.

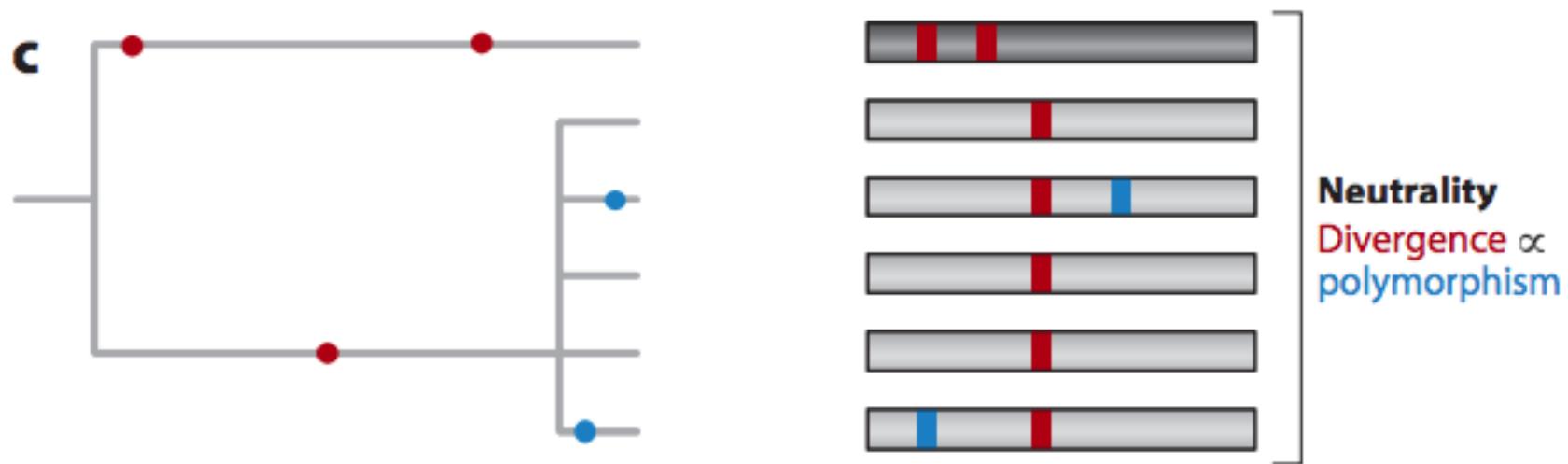
Methods to infer selection (I)

- between-species:
events in the deep past, macro-evolutionary trends,
selection between species

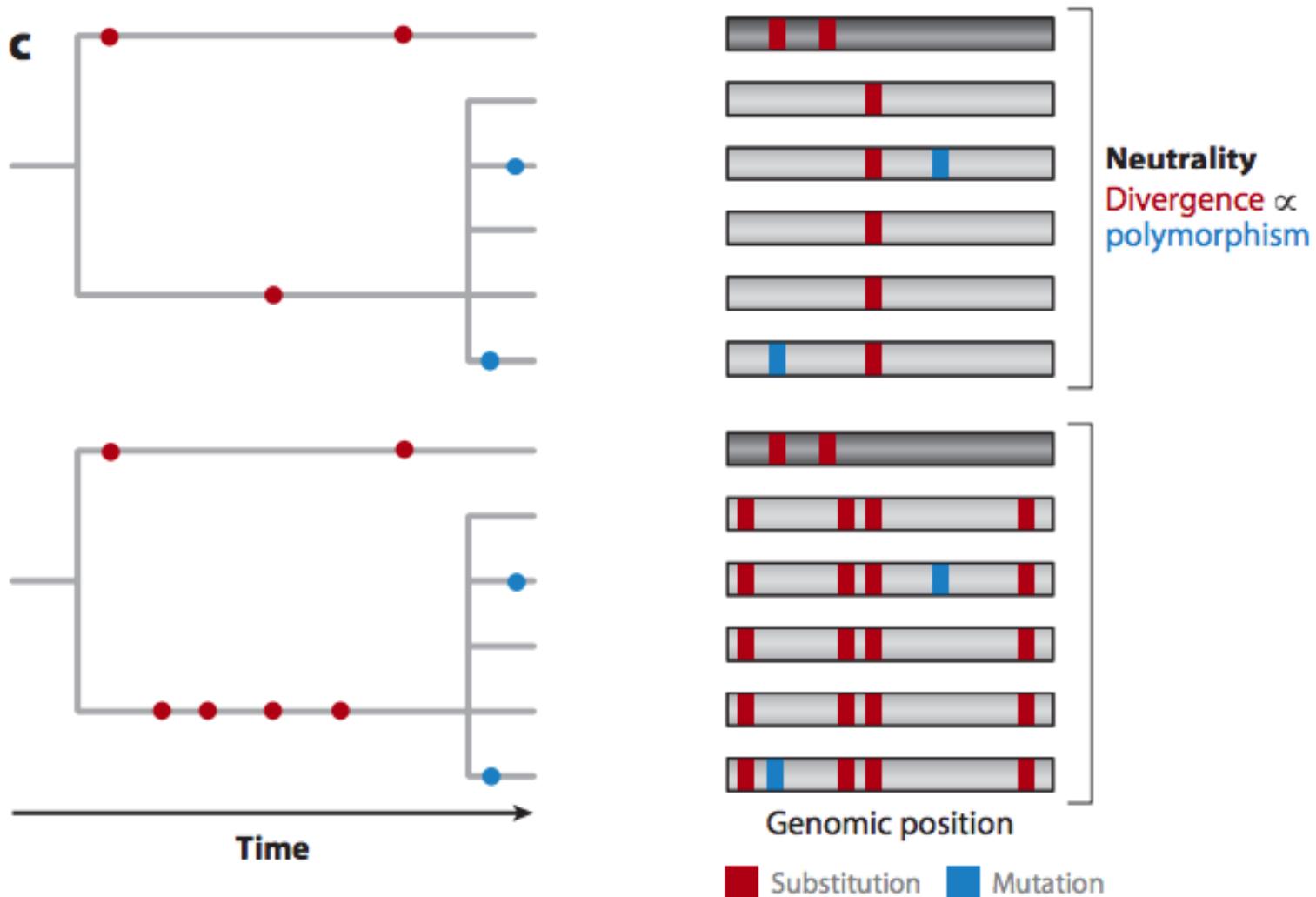


HKA test

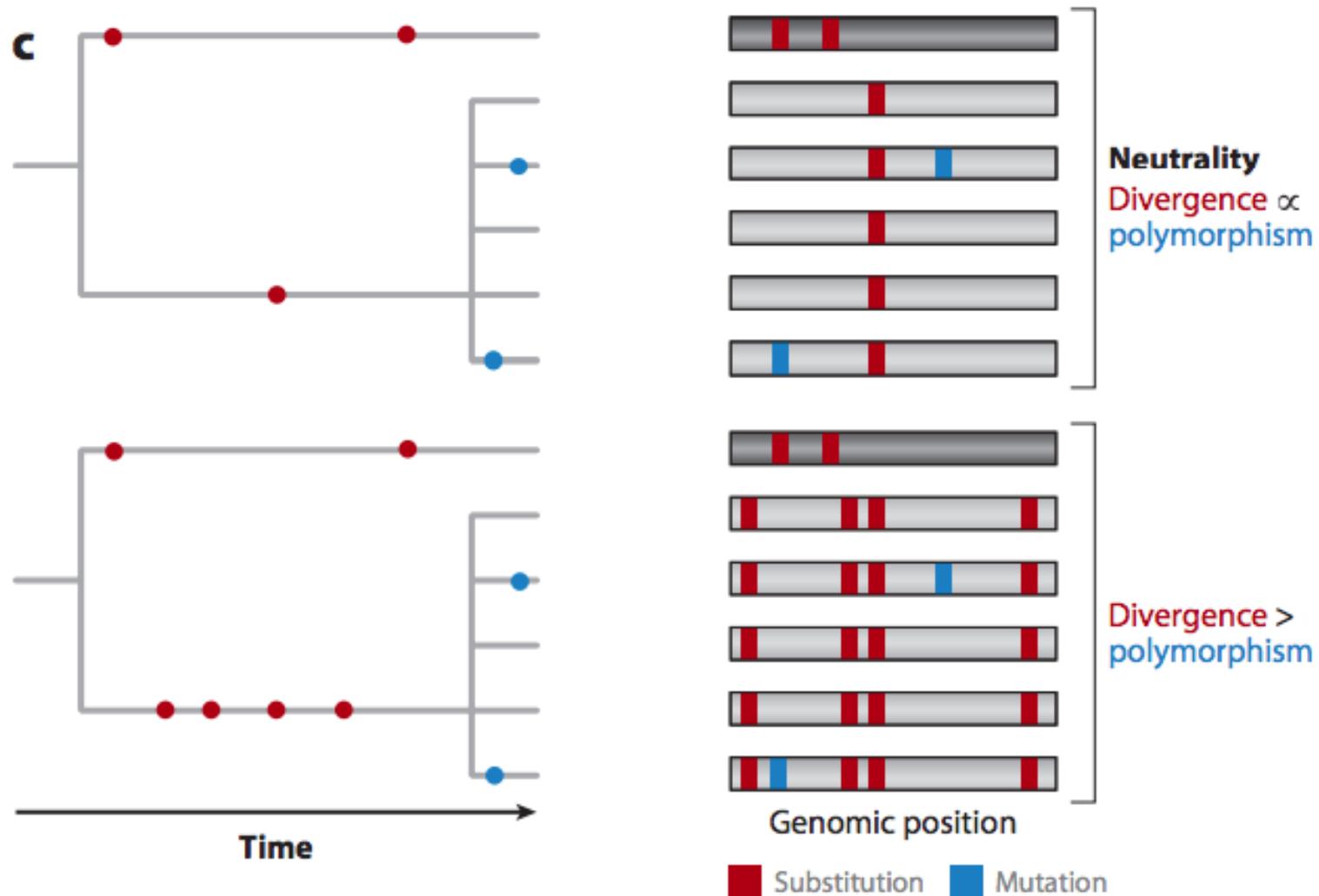
Polymorphism levels depends on local mutation rates: under neutral evolution, the amount **of within- and between-species diversity** is expected to be similar across all loci in the genome.



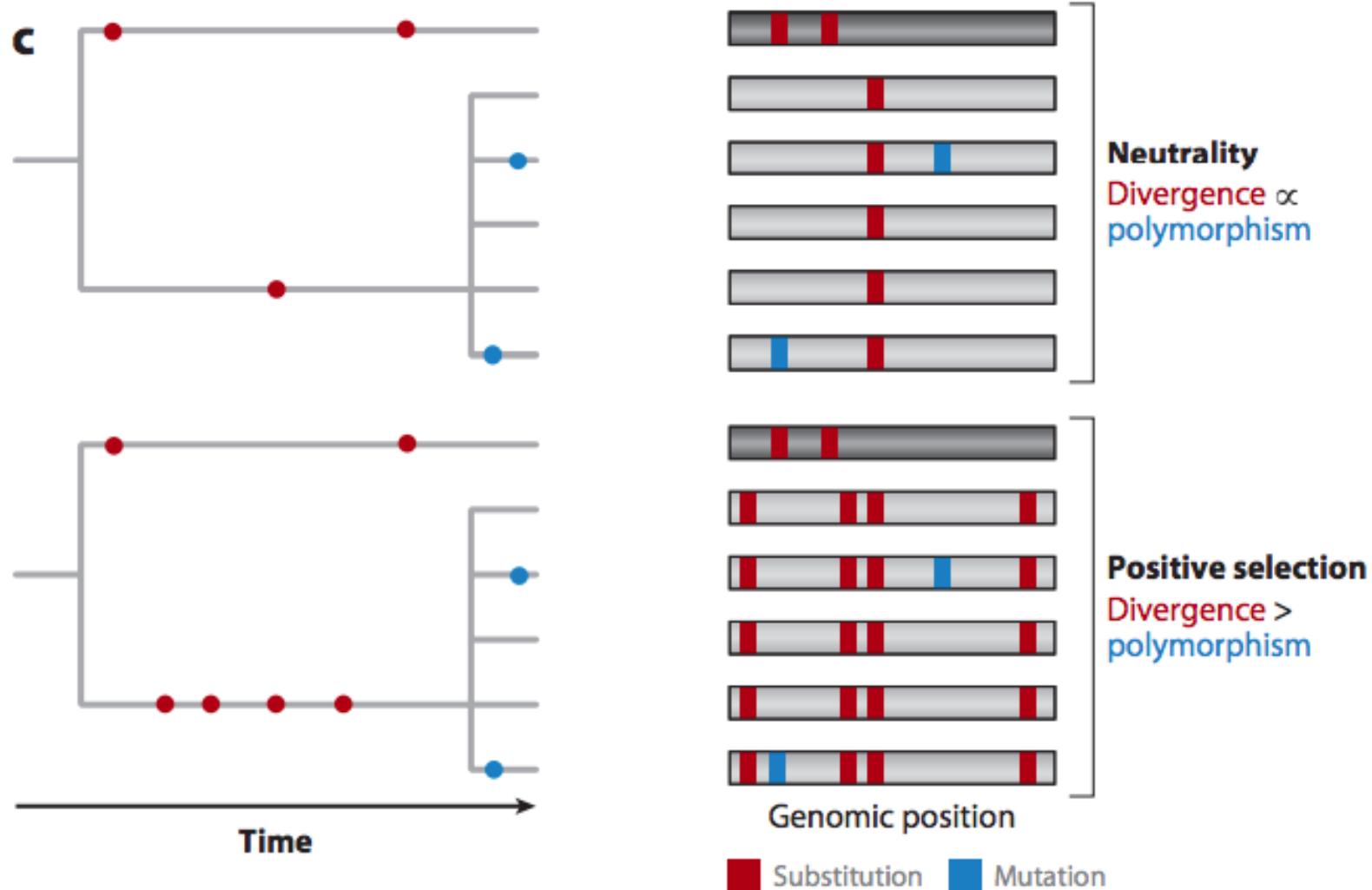
HKA test



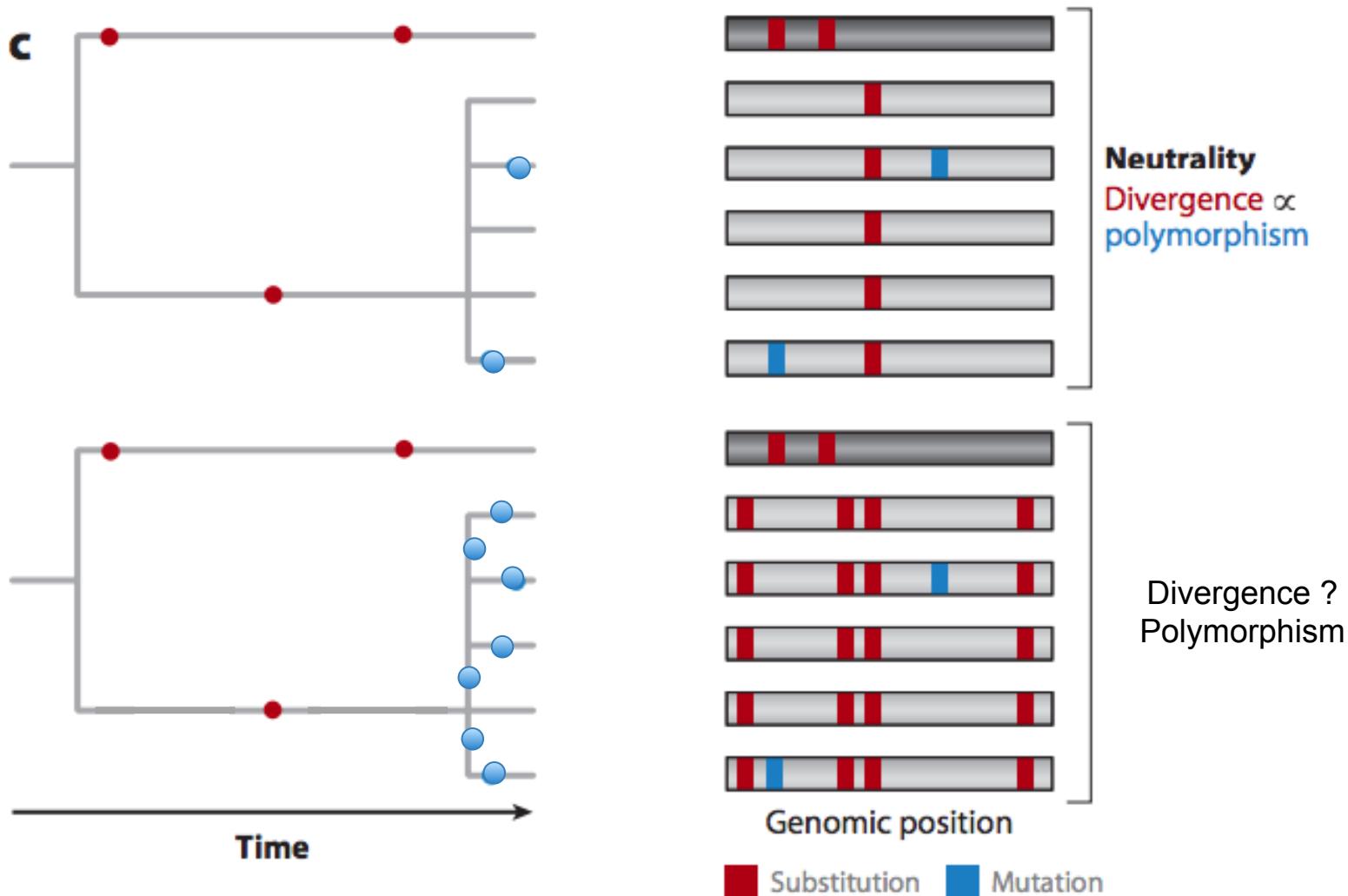
HKA test



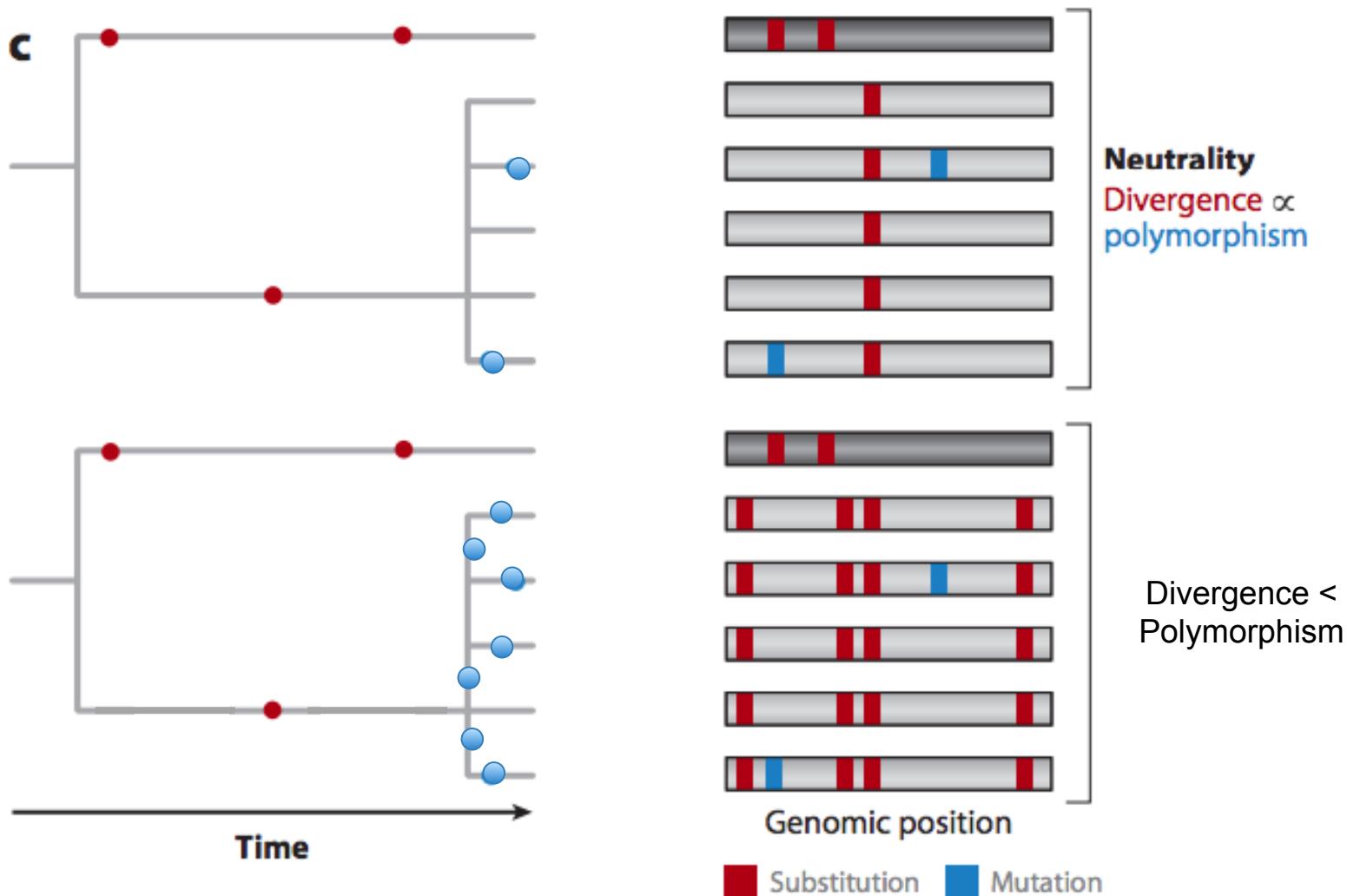
HKA test – positive selection



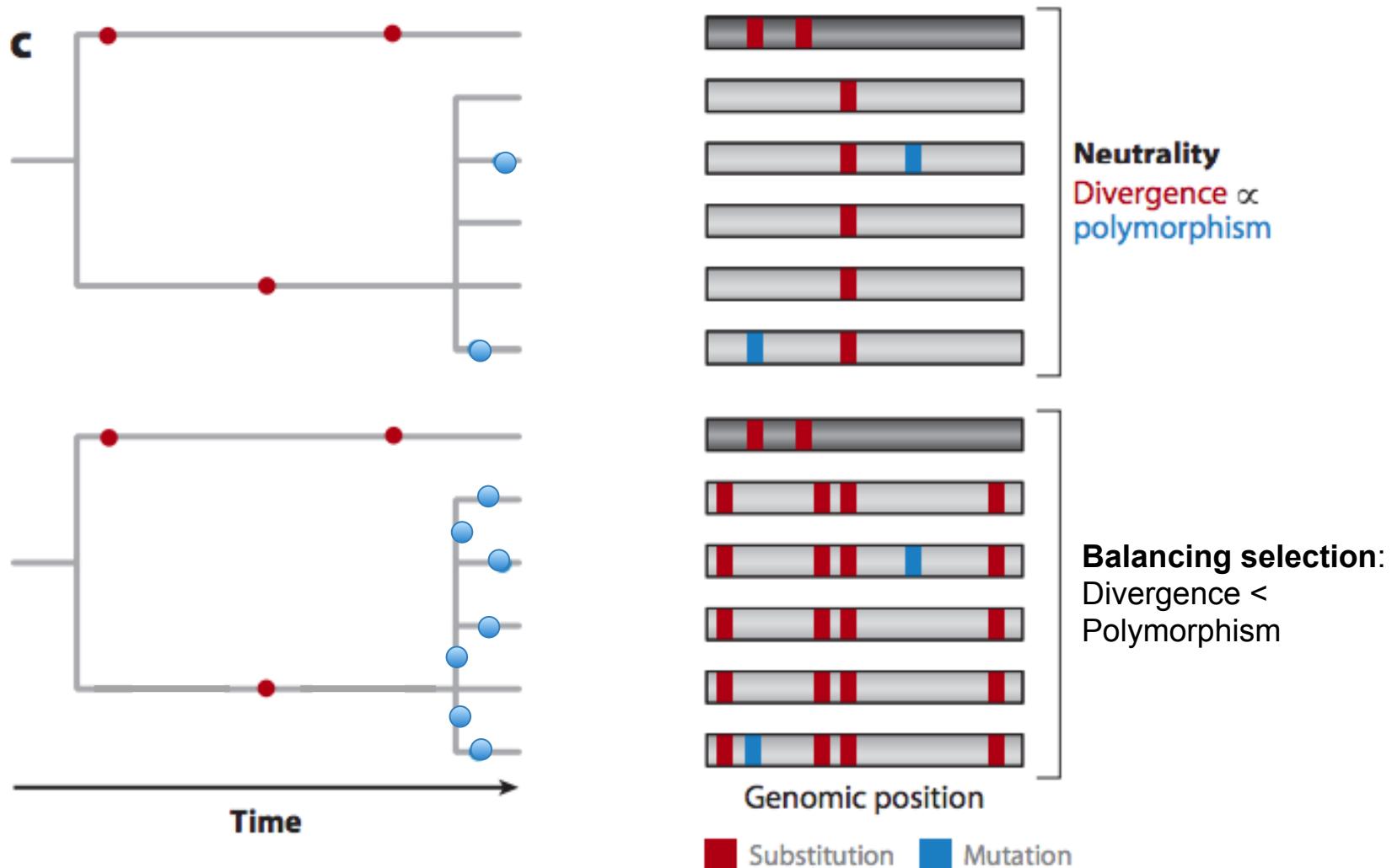
HKA test



HKA test



HKA test – balancing selection



HKA test

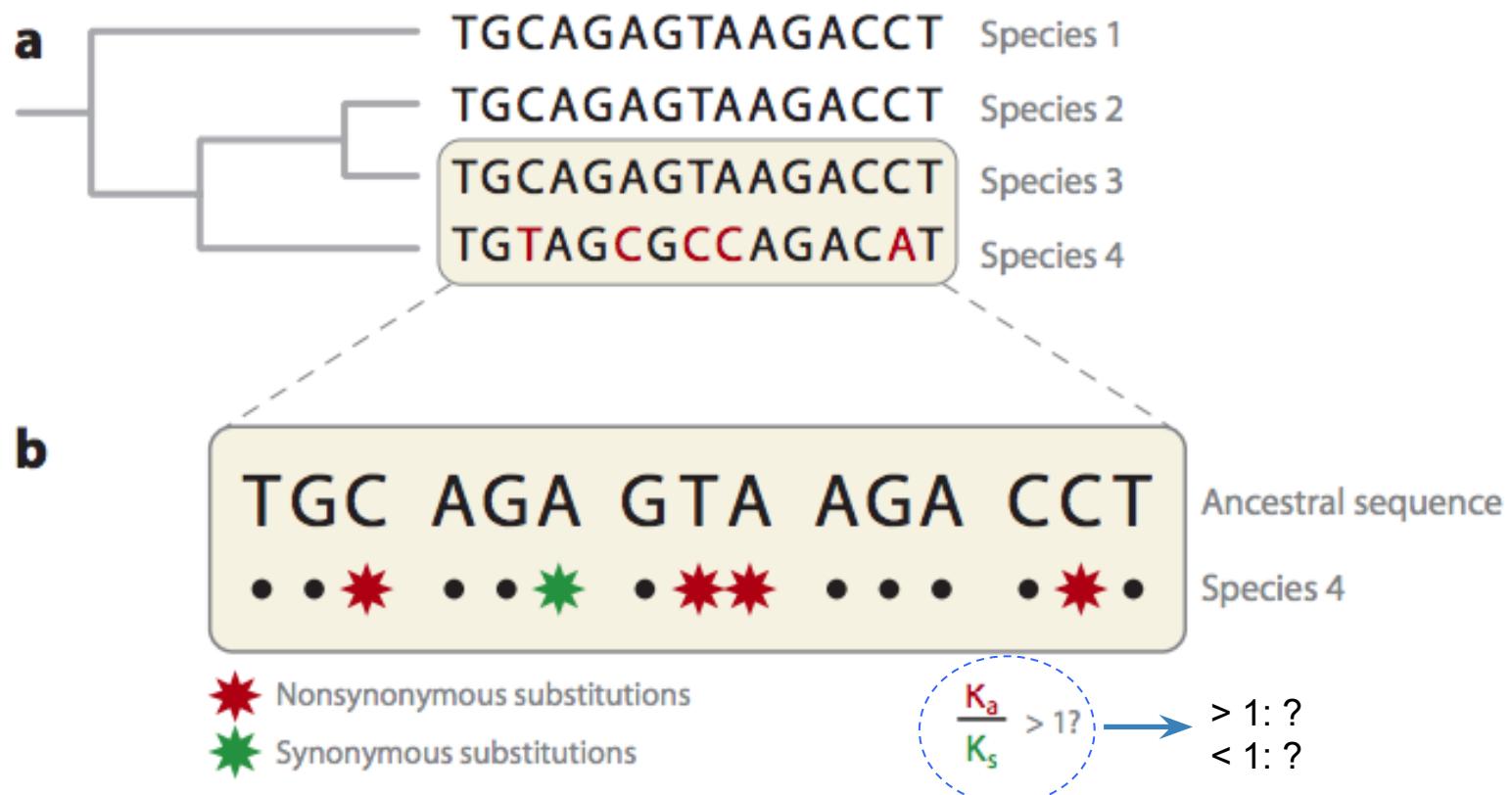
Hudson-Kreitman-Aguadè (HKA, Hudson et al. 1987) test

		Under investigation	Control / Neutrally evolving
		Gene1	Gene2
Polymorphic		G1-poly	G2-poly
	Fixed	G1-fixed	G2-fixed

Contingency table 2x2: chi-square test

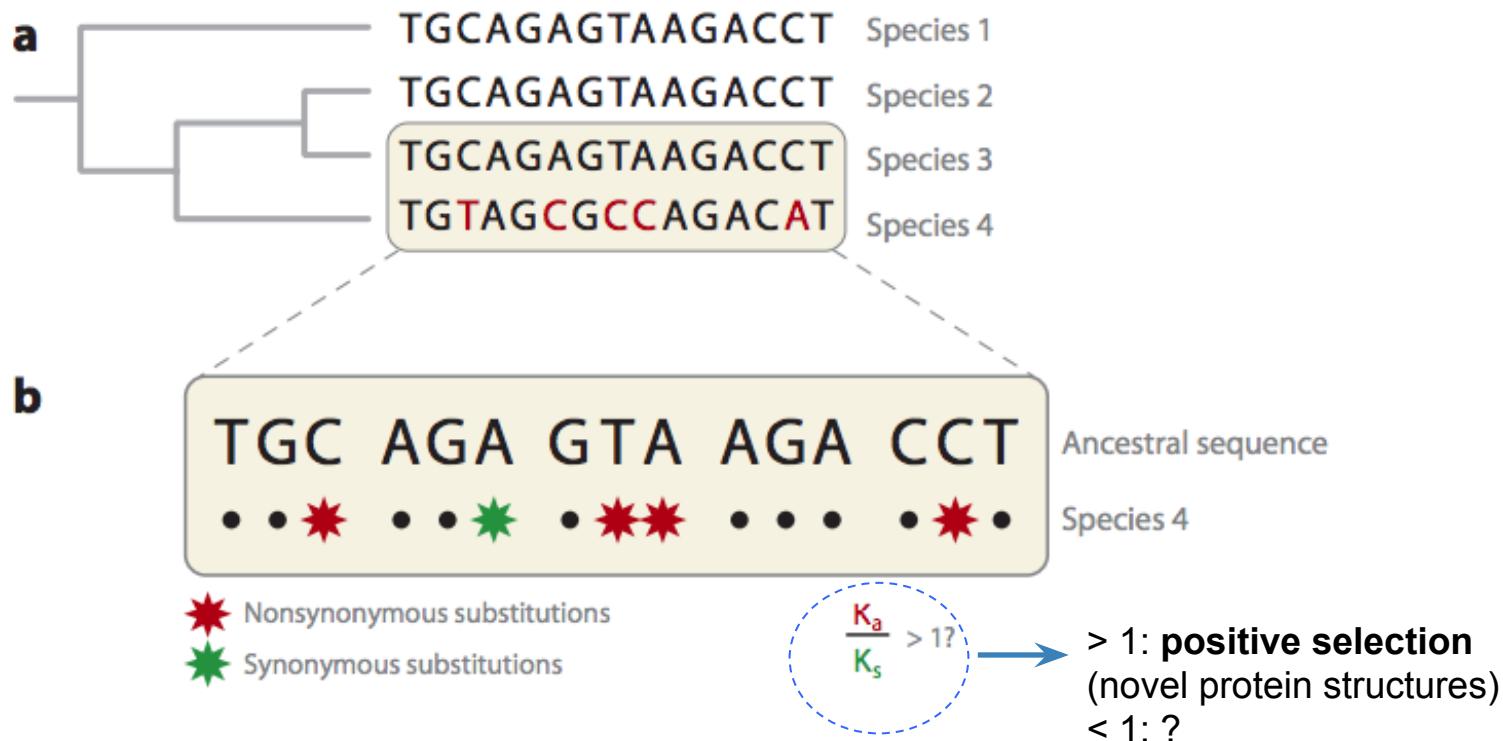
d_N/d_S

Comparison between rate of **nonsynonymous** substitutions and rate of **synonymous** substitutions (also called K_a/K_s or ω)



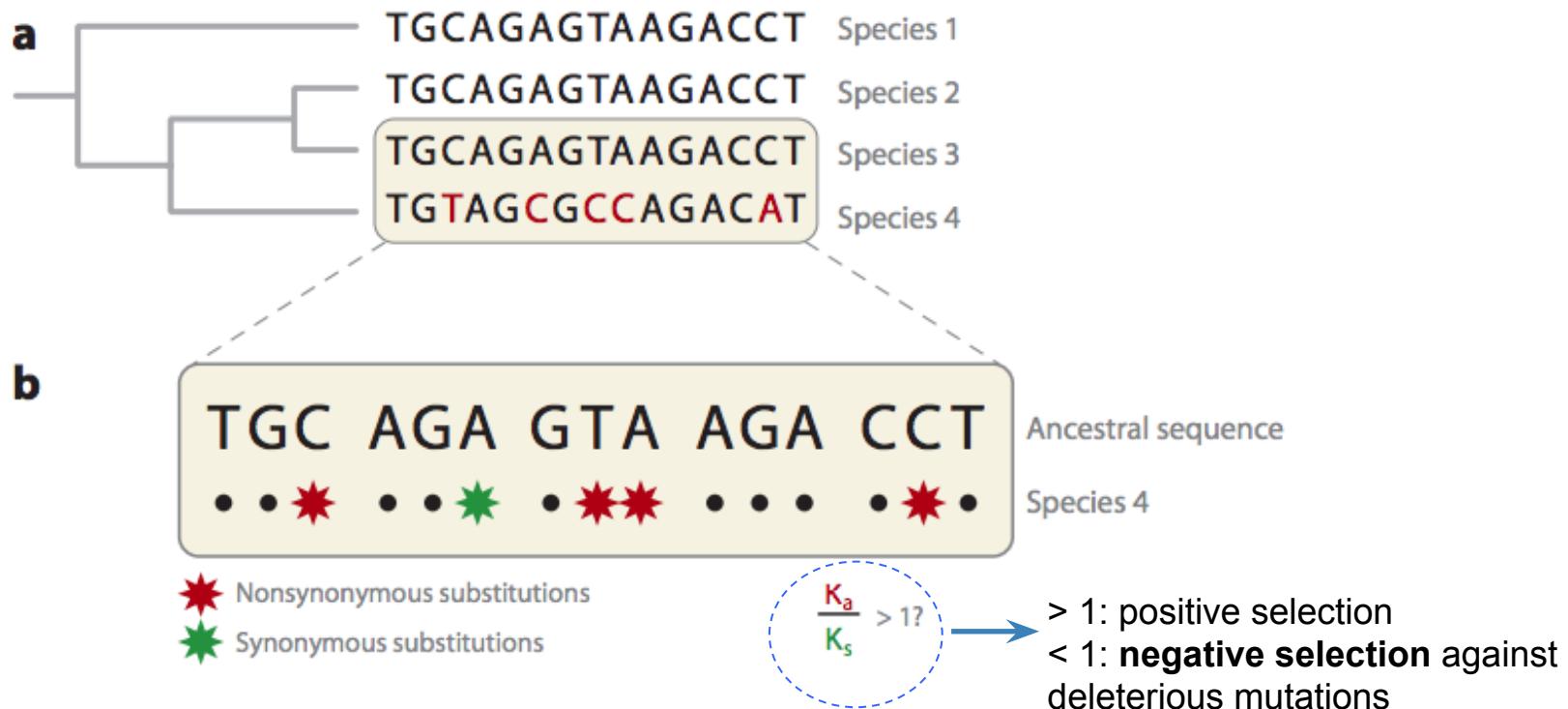
d_N/d_S

Comparison between rate of **nonsynonymous** substitutions and rate of **synonymous** substitutions (also called K_a/K_s or ω)



d_N/d_S

Comparison between rate of **nonsynonymous** substitutions and rate of **synonymous** substitutions (also called K_a/K_s or ω)



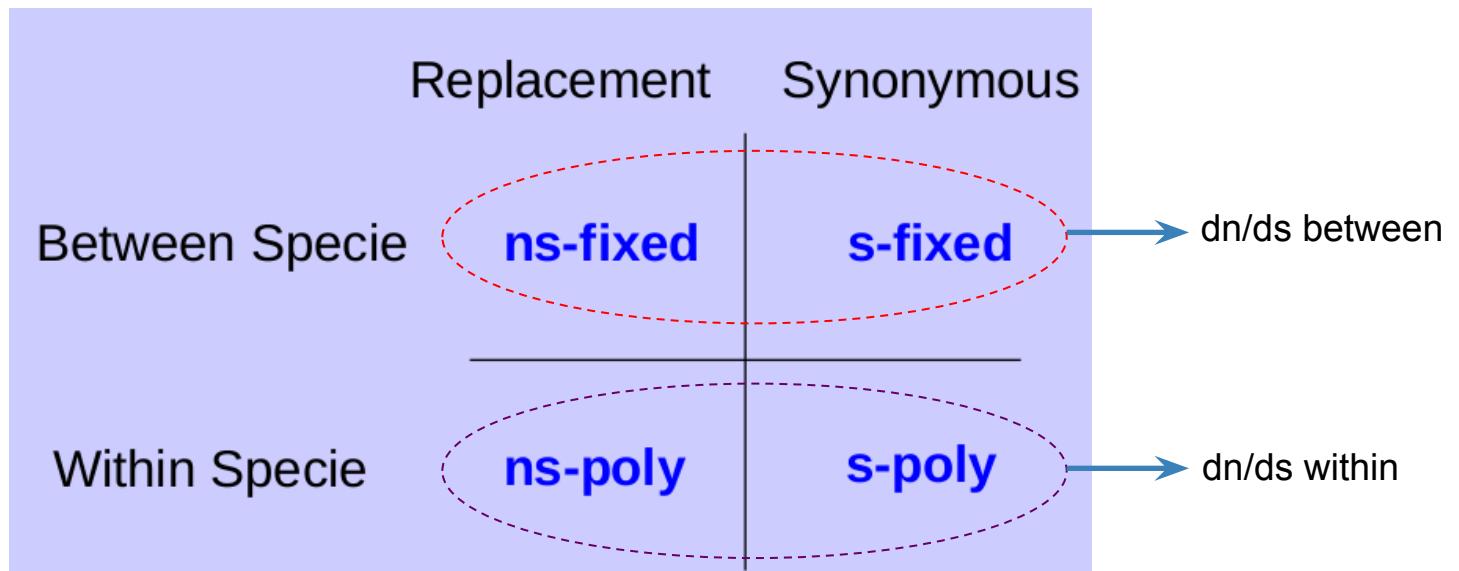
MK test

The MKT compares the amount of variation within a species to the divergence between species at **two types of sites**, one of which is putatively neutral and used as the reference to detect selection at the other types of sites.

	Replacement	Synonymous
Between Species	ns-fixed	s-fixed
Within Species	ns-poly	s-poly

MK test

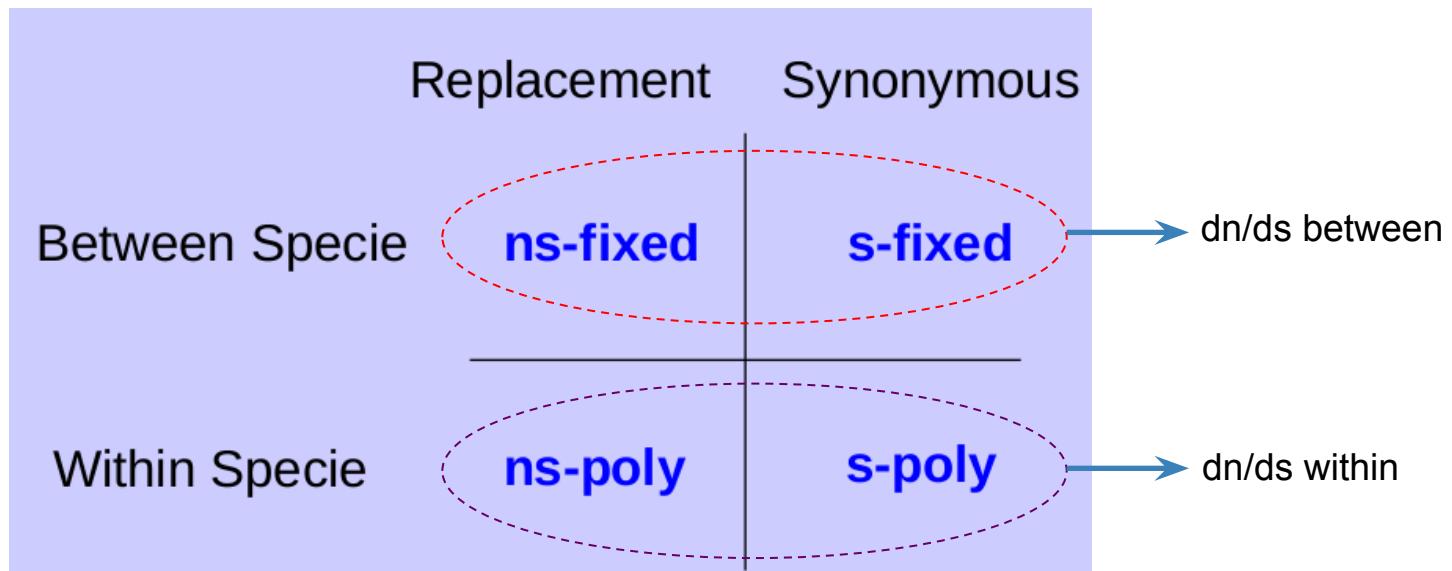
The MKT compares the amount of variation within a species to the divergence between species at **two types of sites**, one of which is putatively neutral and used as the reference to detect selection at the other types of sites.



MK test

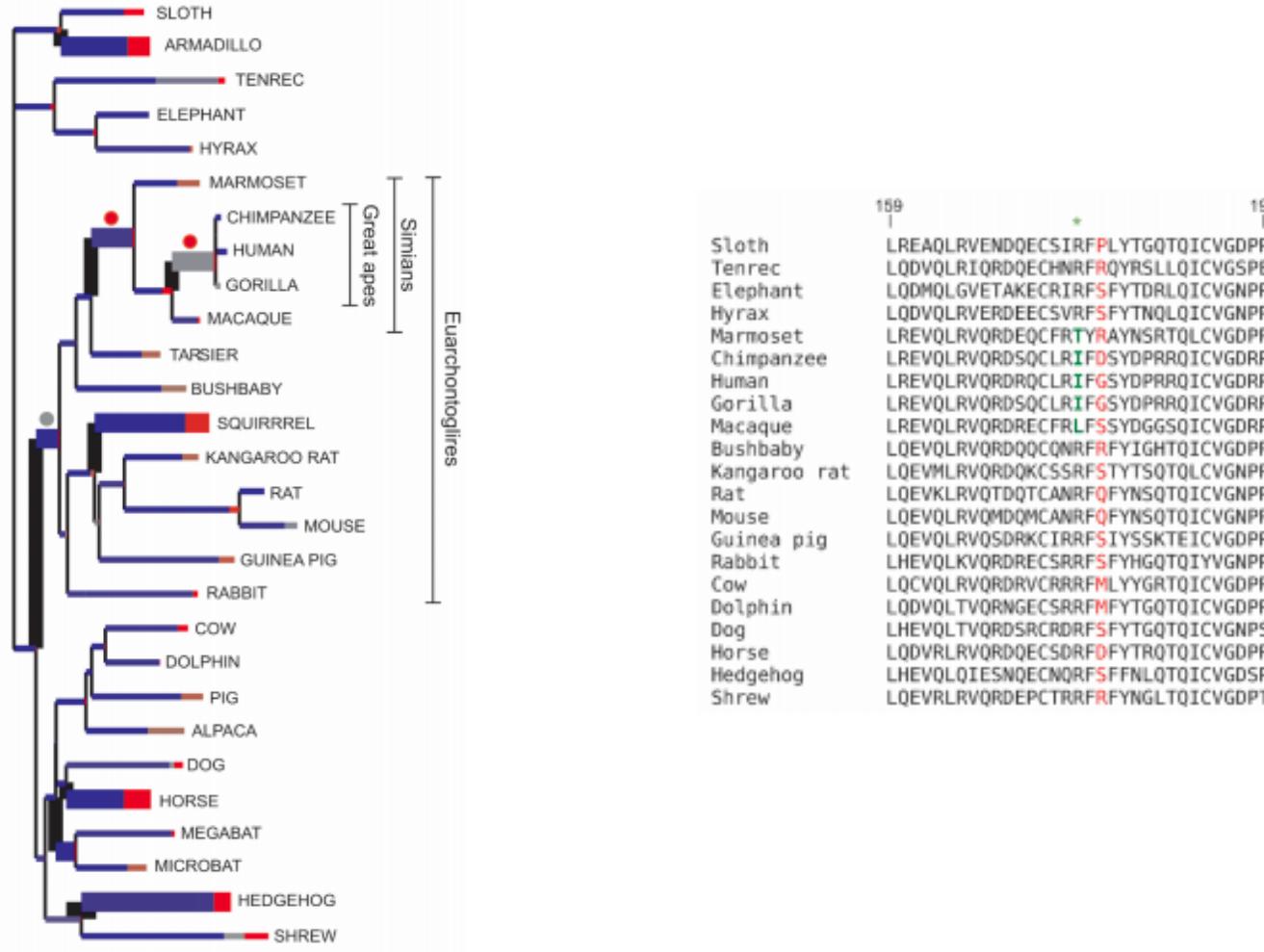
Under neutrality, these rates should be equal.

If the between-species ratio exceeds the within-species ratio indicates **positive selection** between species (a larger within-species value suggests balancing selection or excess of weakly deleterious variants).



Comparative genomics

Identify multiple-species conserved genomic elements which show an accelerated rate of substitution in a particular lineage.



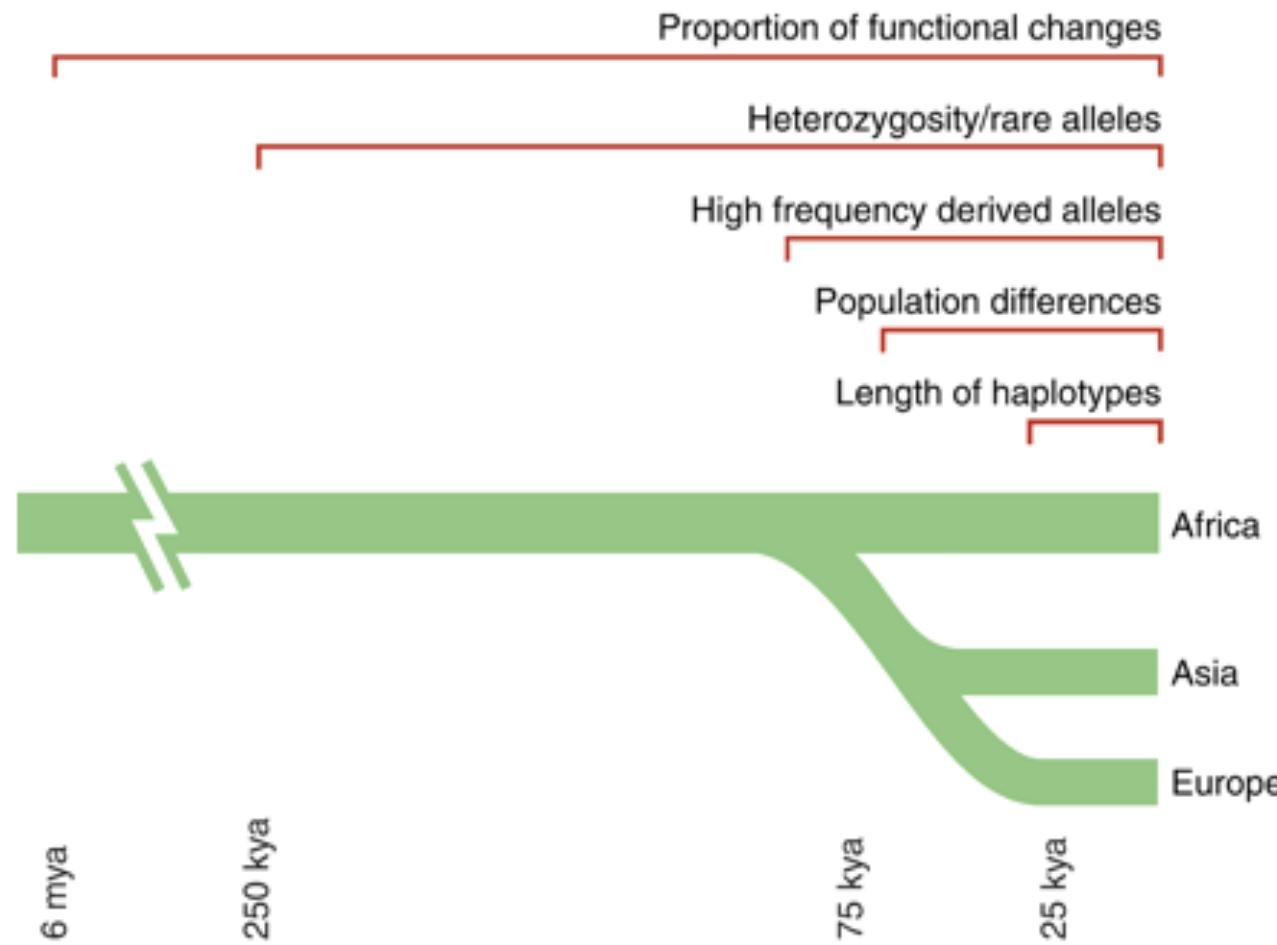
Methods to infer selection (II)

- within-species:

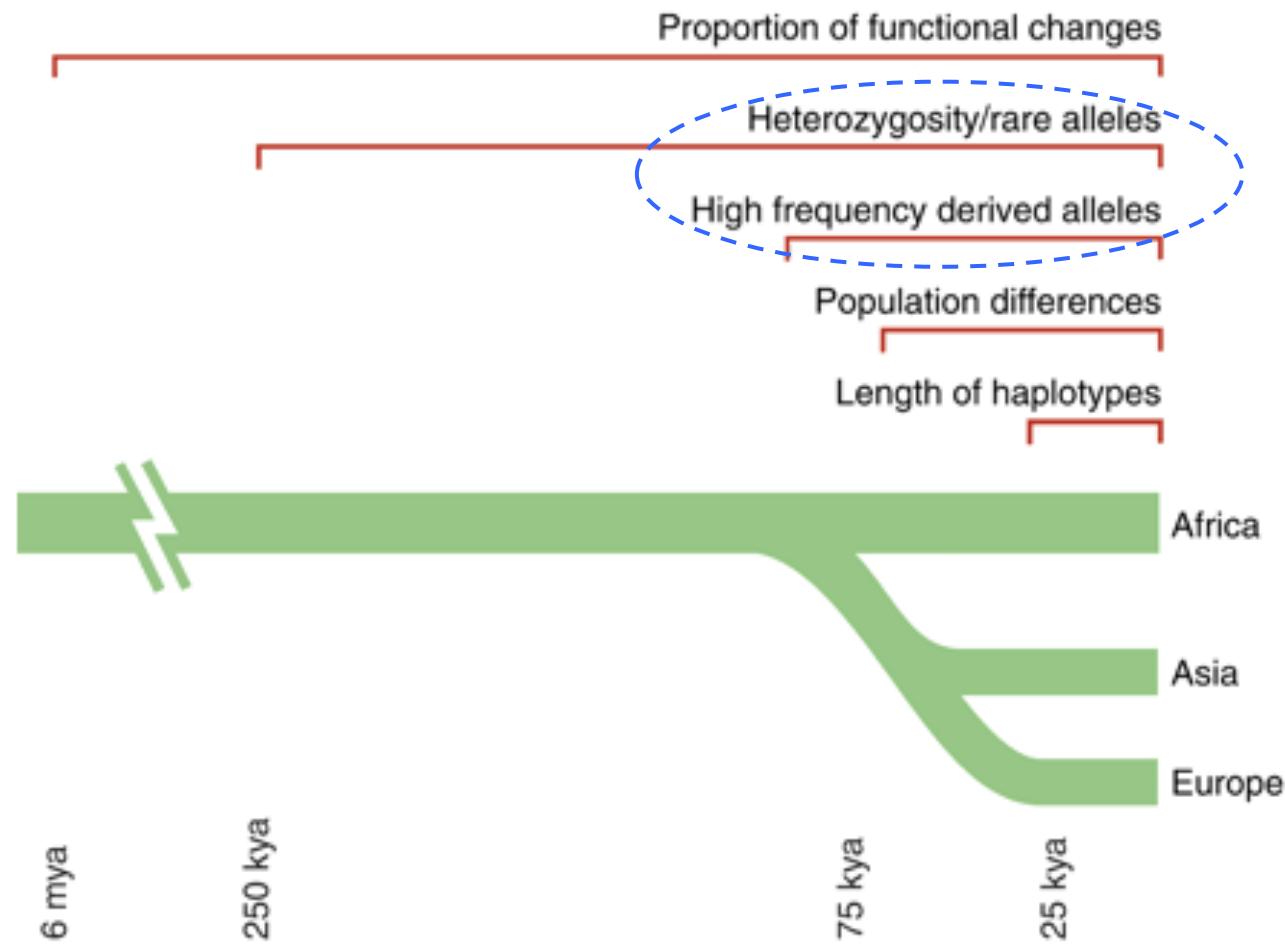
Micro-evolutionary events between populations, local adaptation



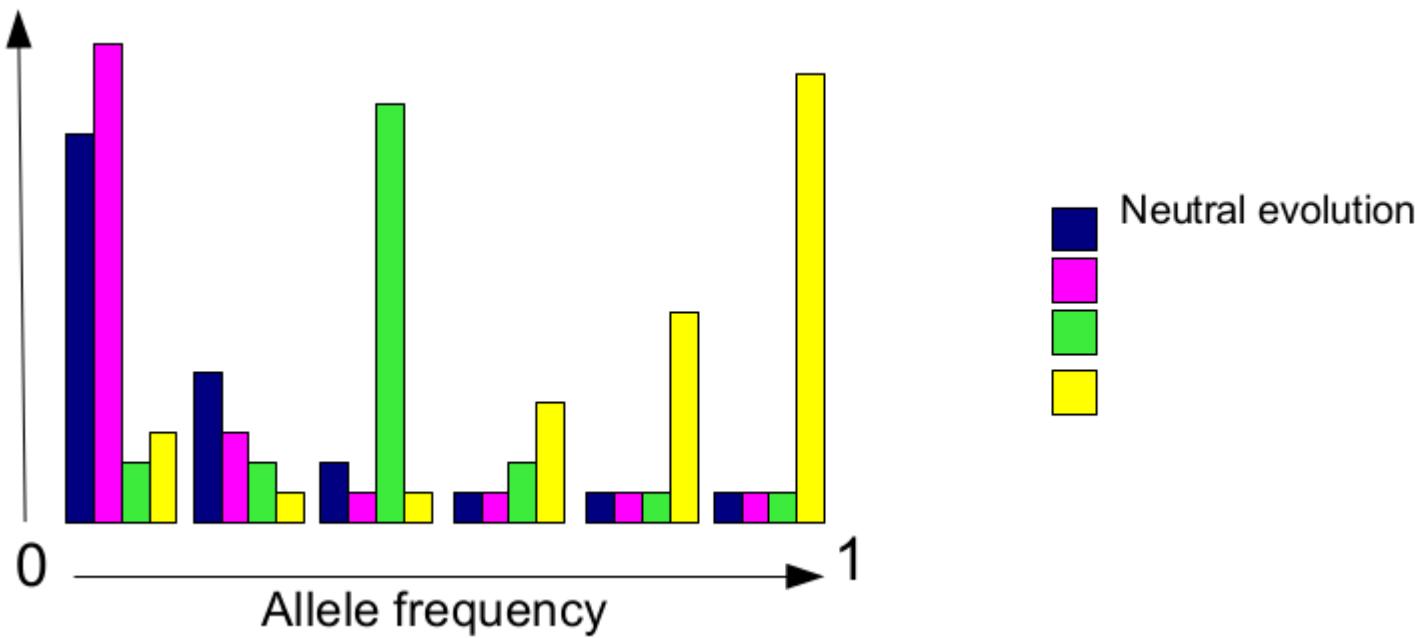
Methods to infer recent selection



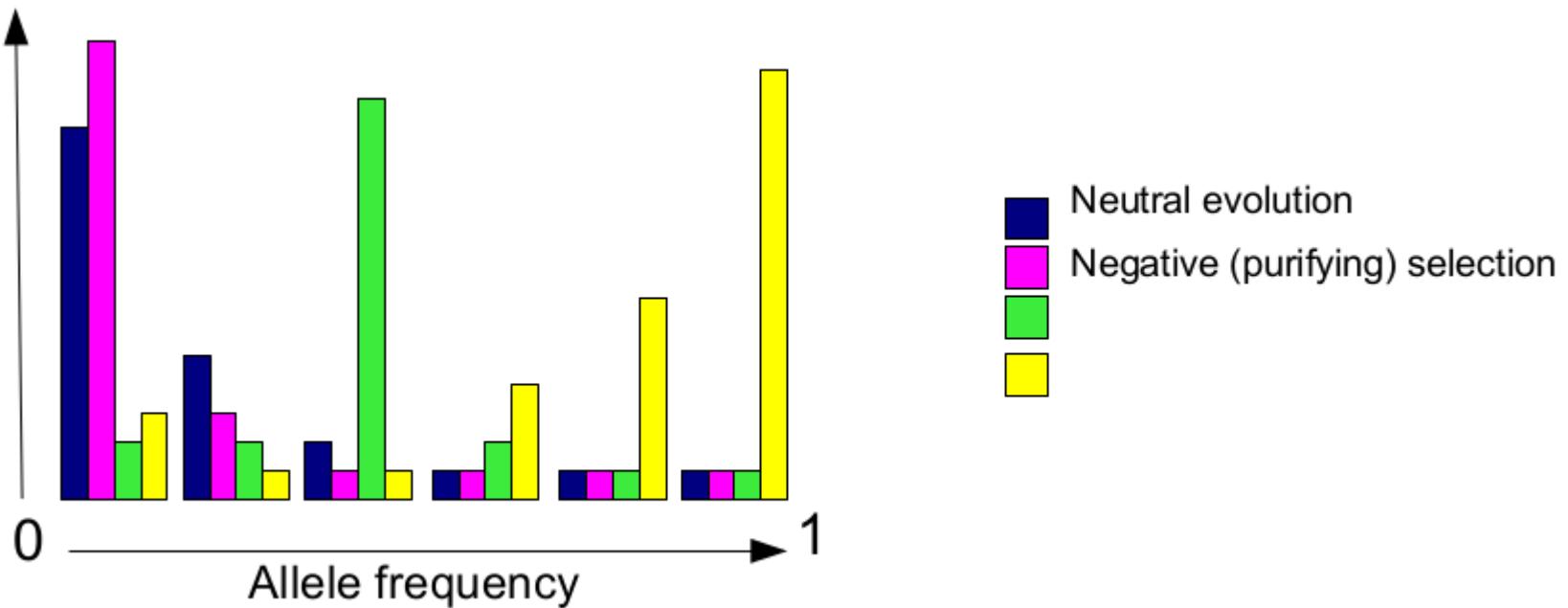
Methods to infer recent selection



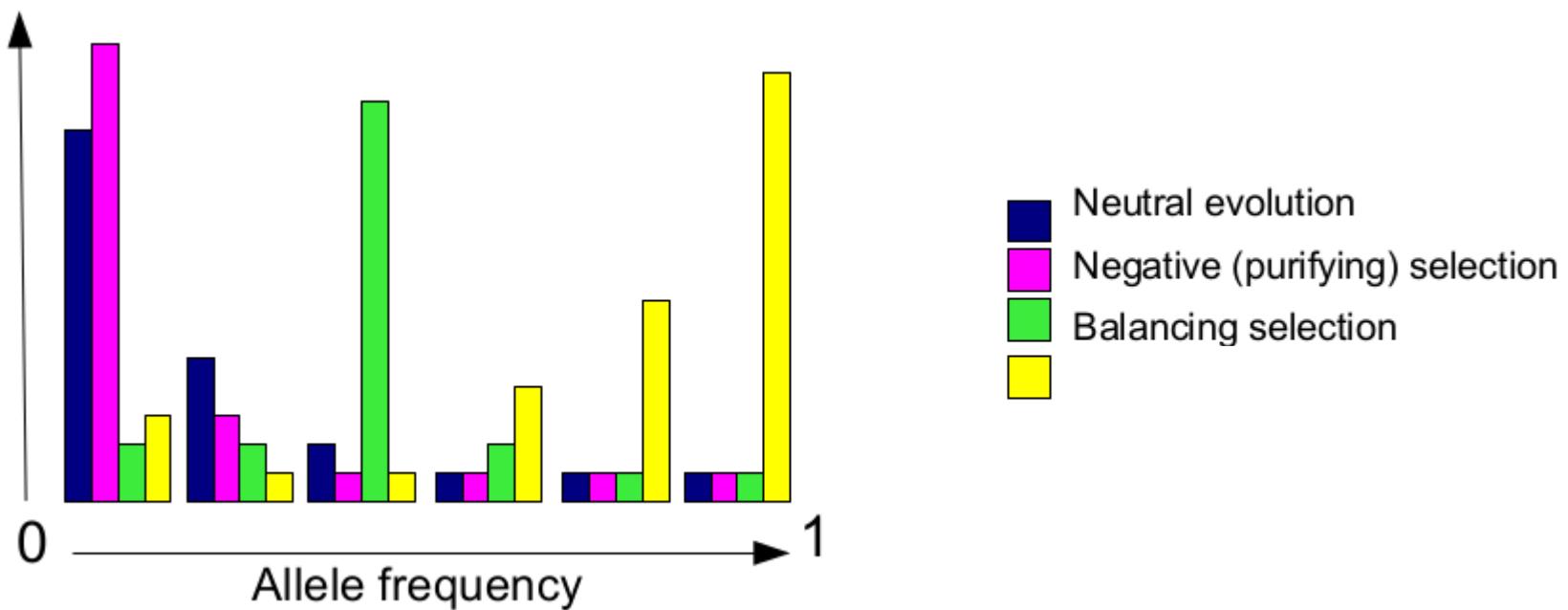
SFS-based statistics



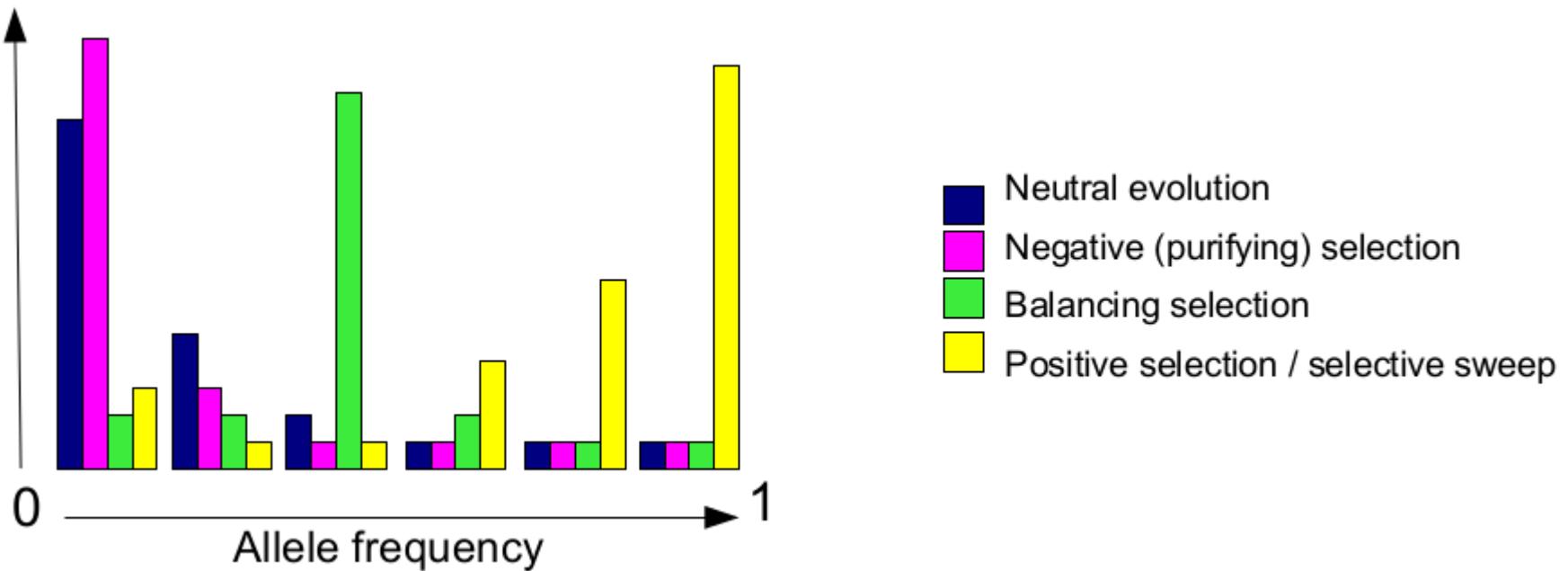
SFS-based statistics



SFS-based statistics



SFS-based statistics

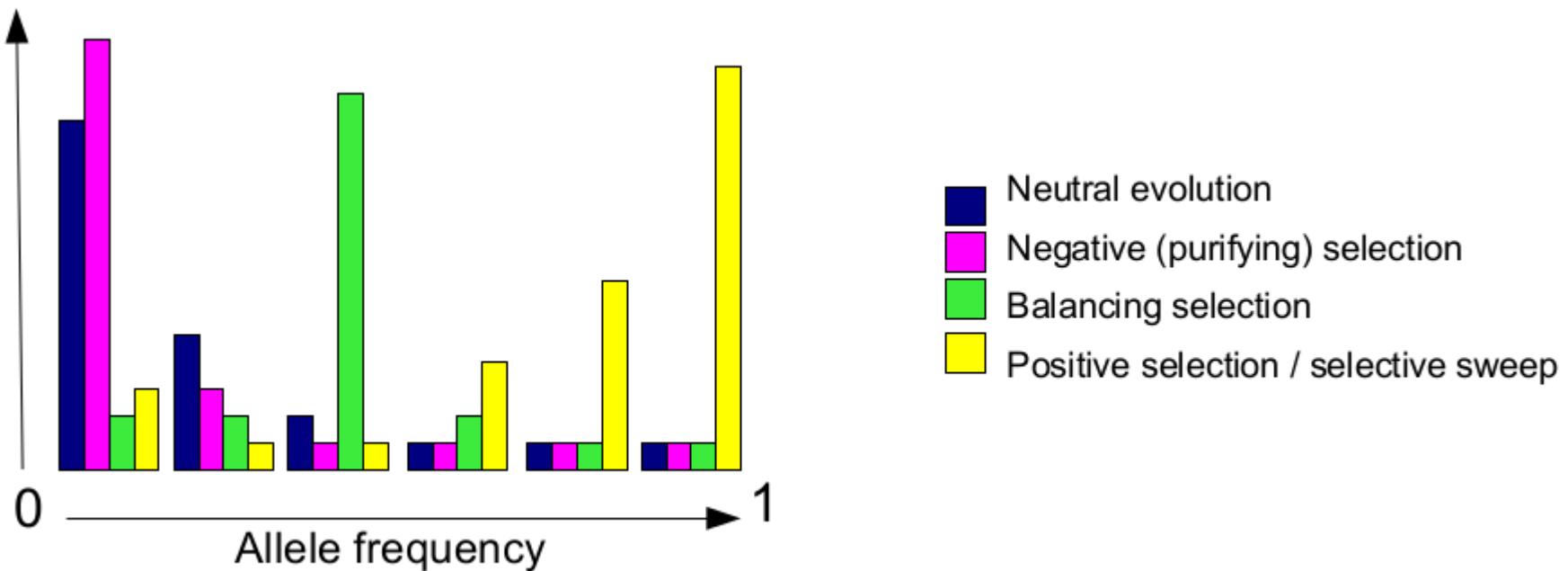


SFS-based statistics

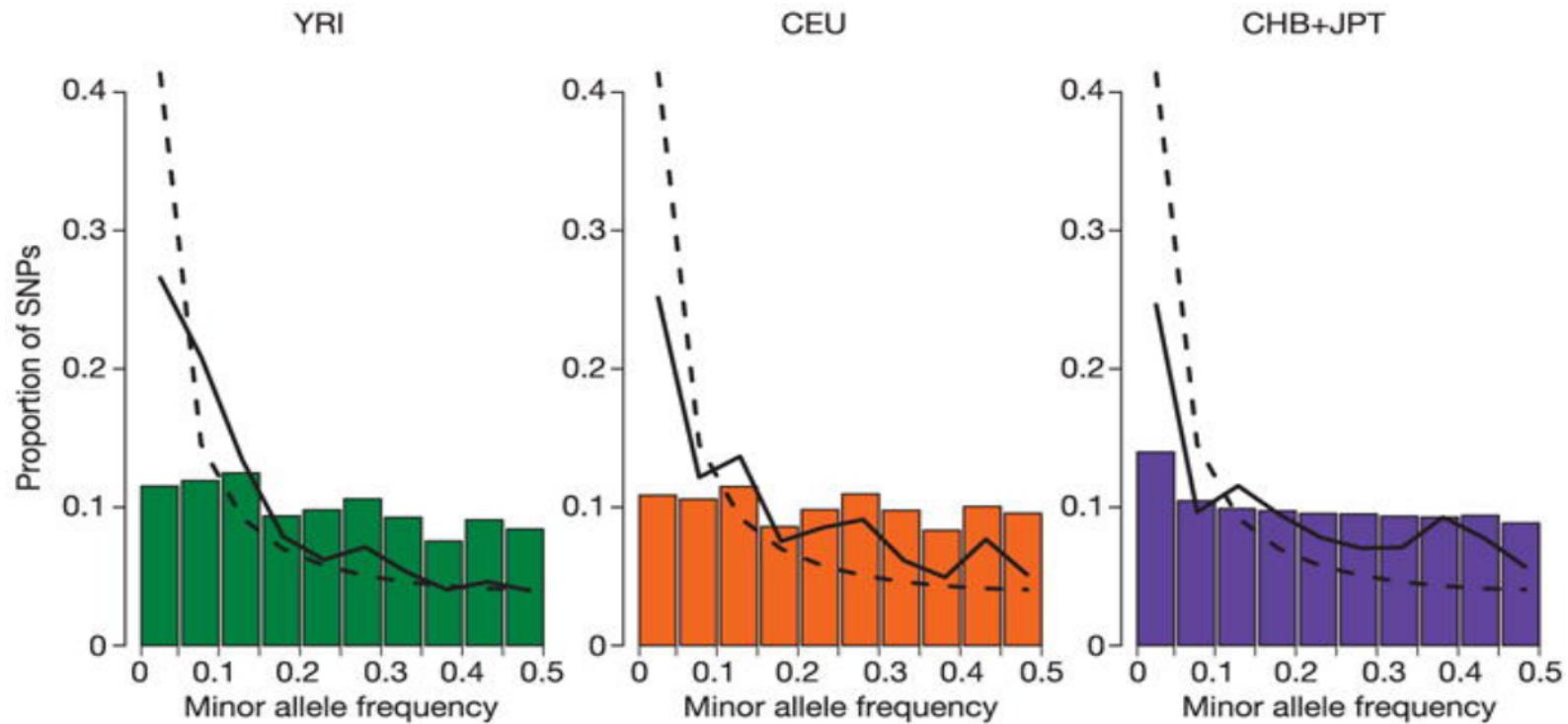
Tajima's D (1989) is the standardized difference of two nucleotide diversity estimates π and θ_w .

- Positive values of D indicate an excess of intermediate frequency alleles;
- Negative values of D indicate an excess of low frequency alleles.

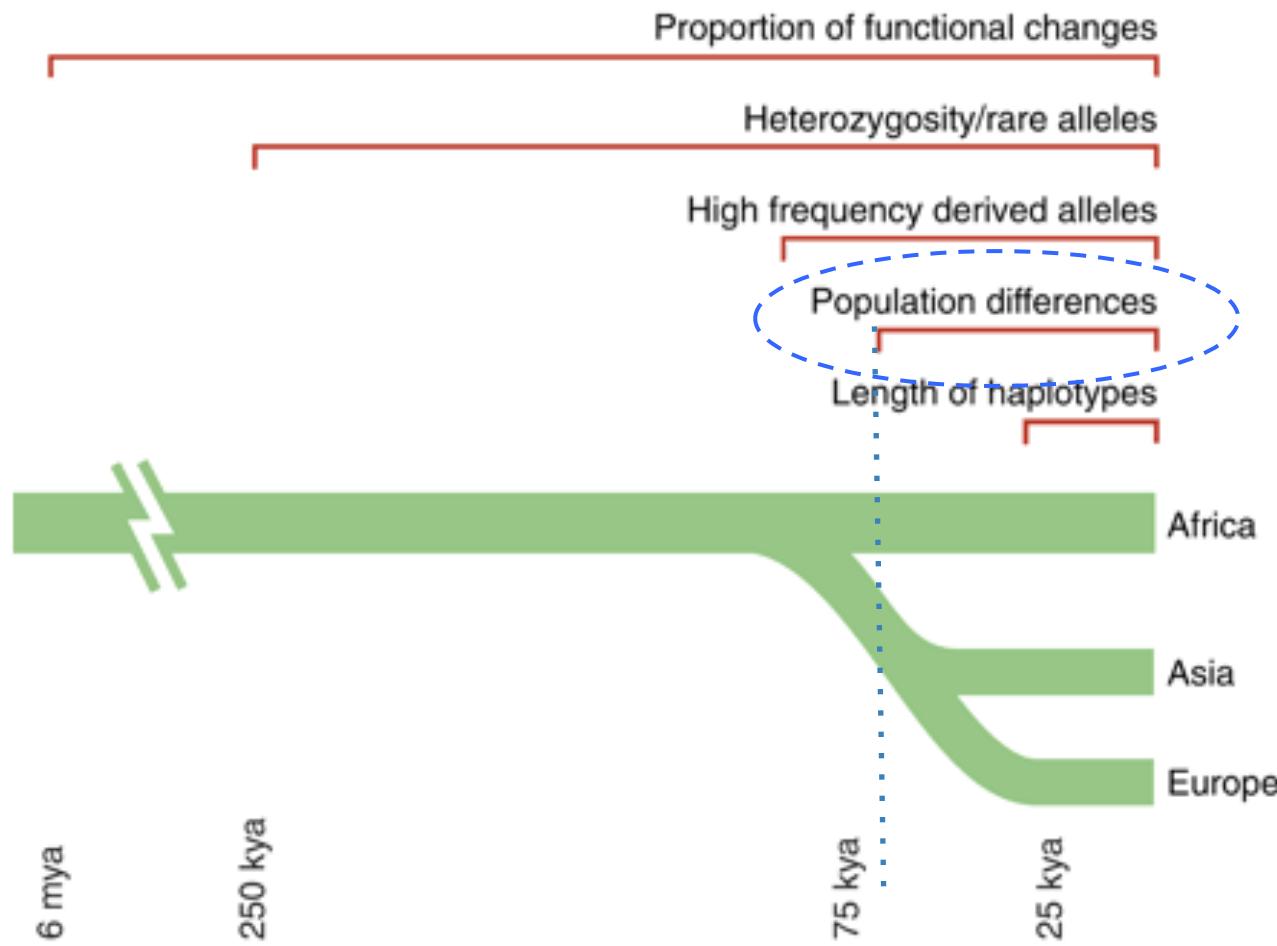
Fu and Li (1993) extended Tajima's test by including information regarding the genealogy of the sample.



Ascertainment bias

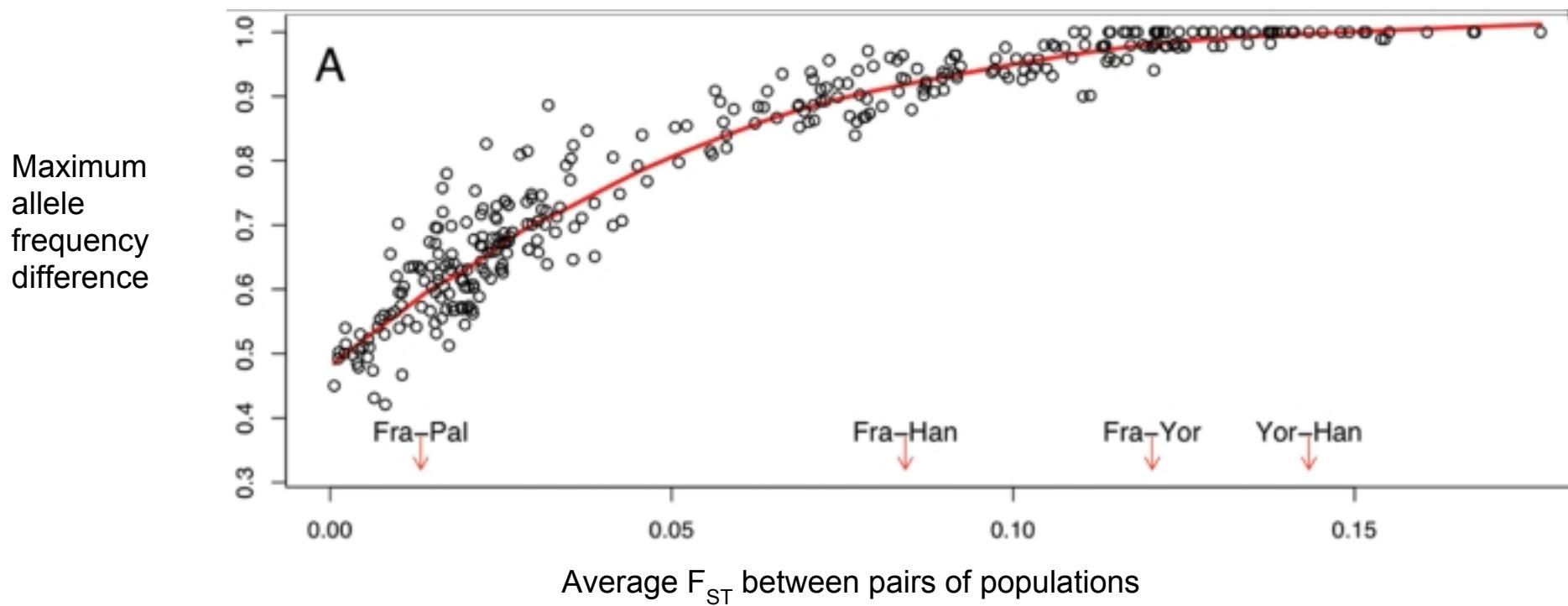


Methods to infer recent selection



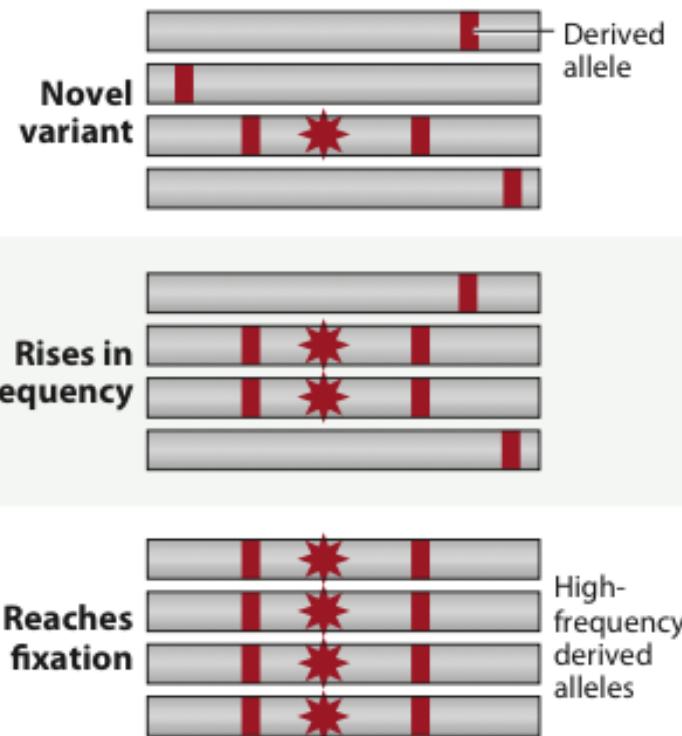
F_{ST}

Method-of-moments or maximum-likelihood estimators of F_{ST} .

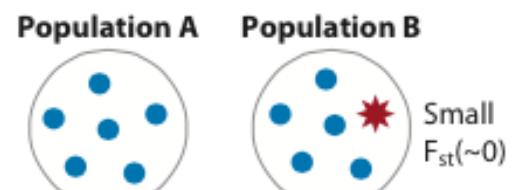


F_{ST}

b Frequency spectrum

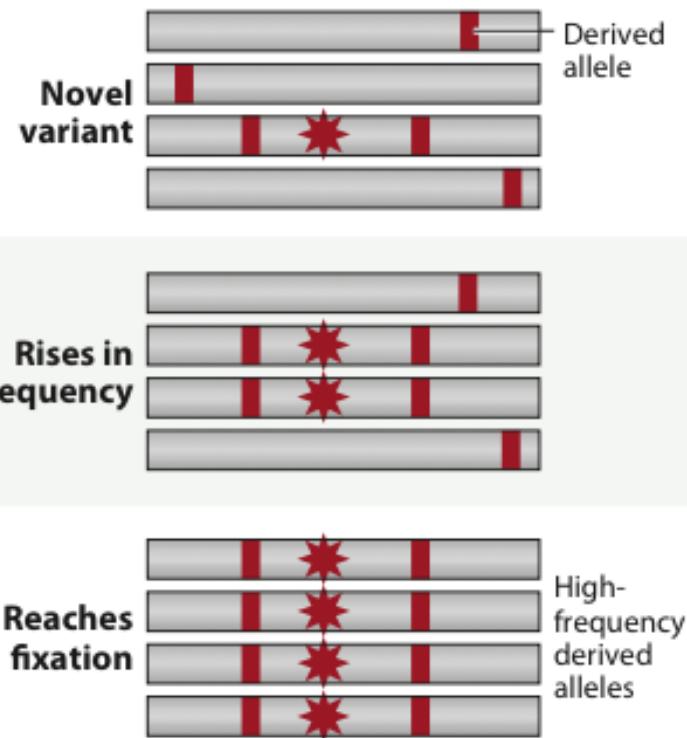


d Population differentiation

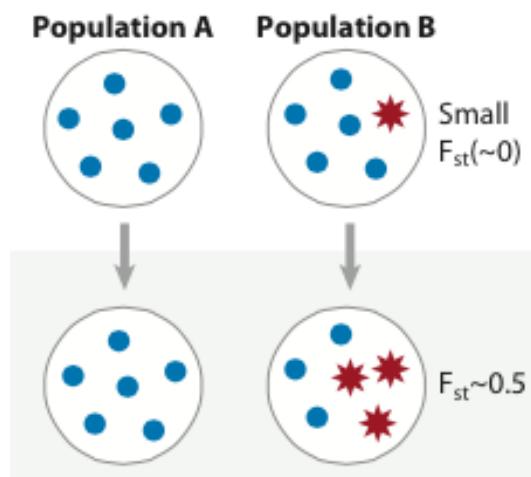


F_{ST}

b Frequency spectrum

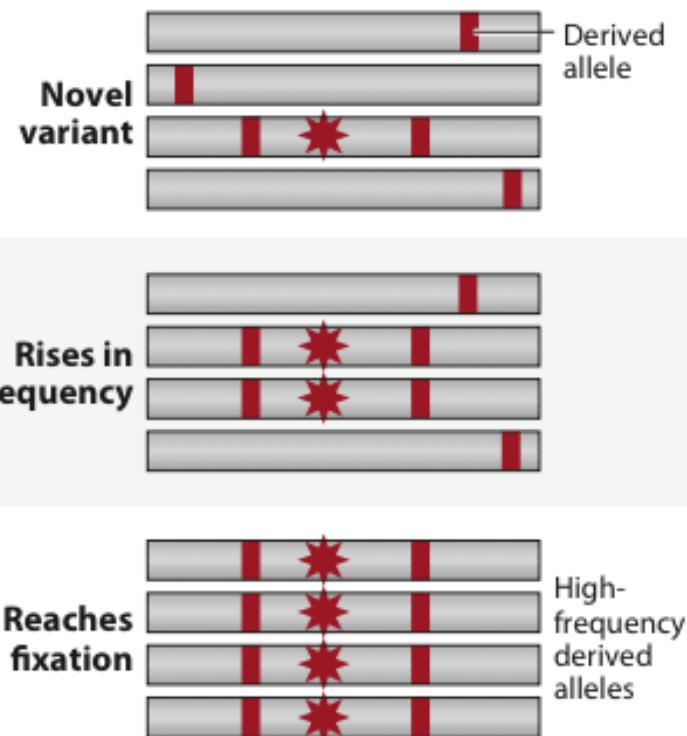


d Population differentiation

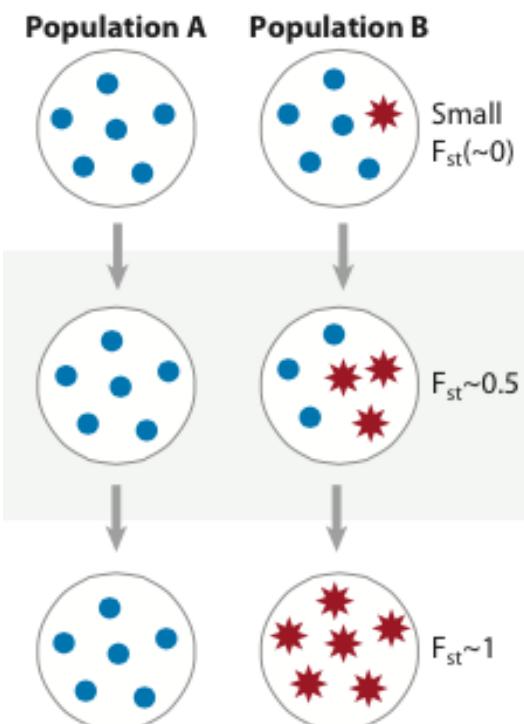


F_{ST}

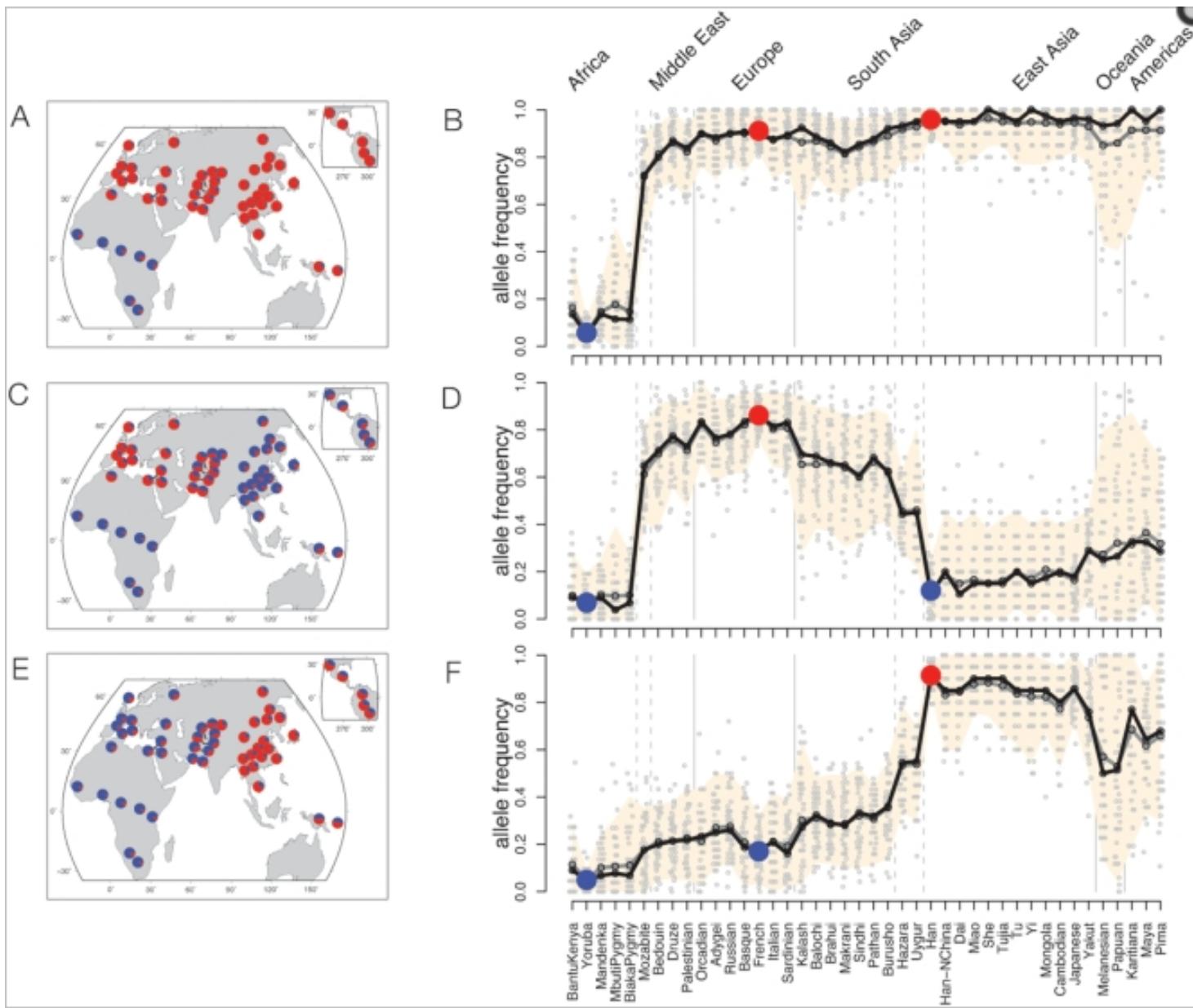
b Frequency spectrum



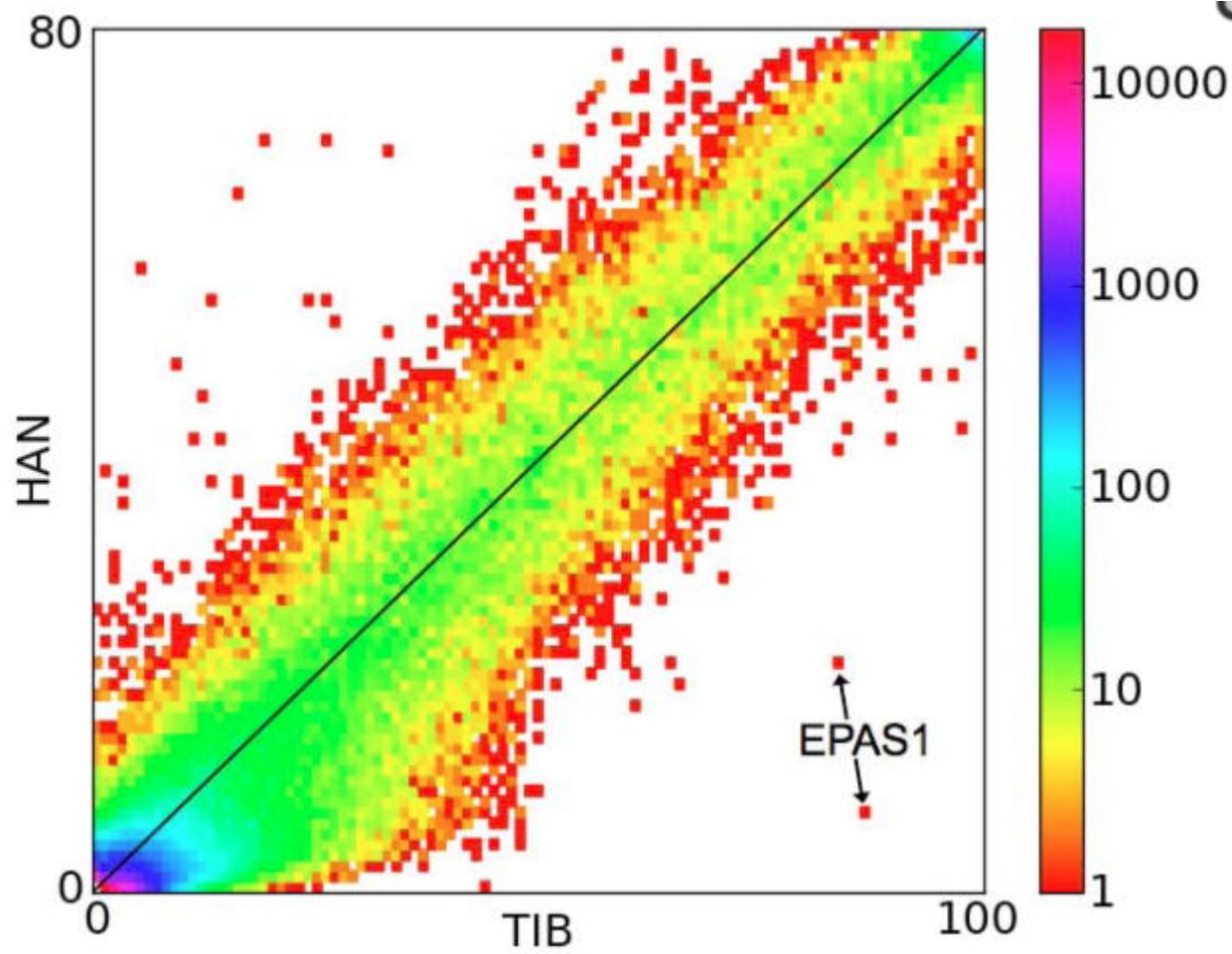
d Population differentiation



Genetic differentiation

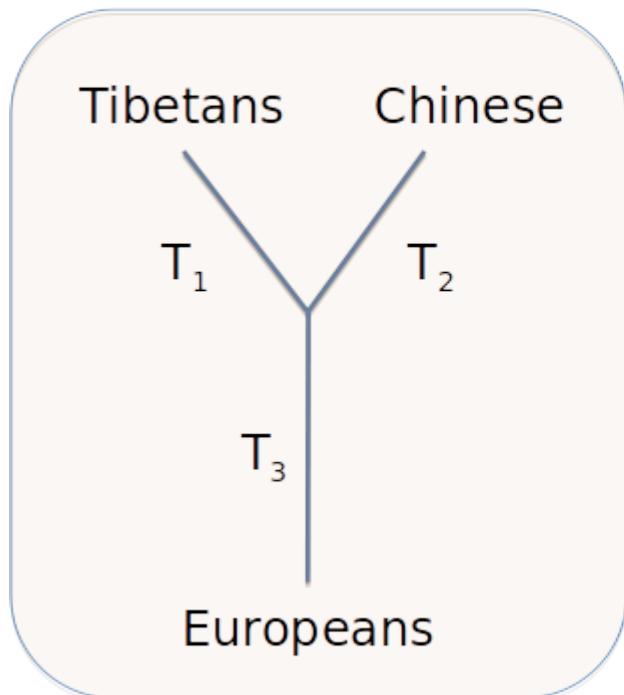


2D-SFS

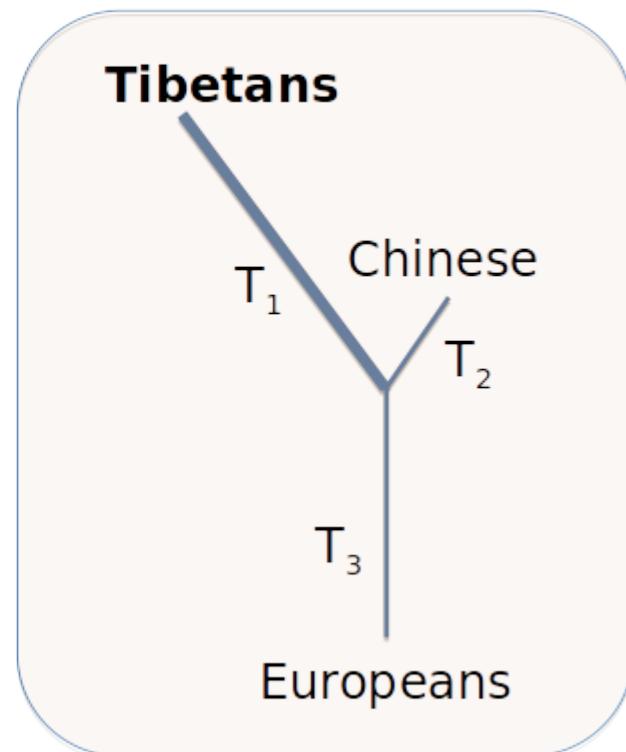


Population Branch Statistic

Neutral evolution



Positive selection

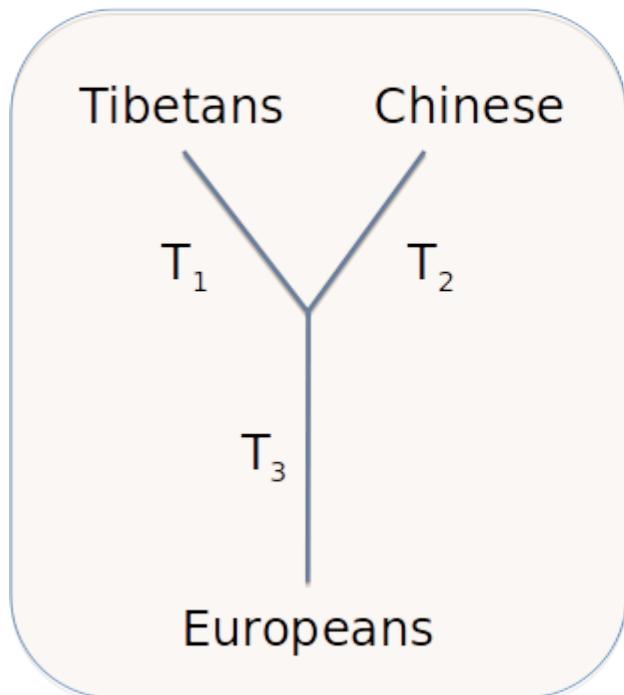


Population Branch Statistic (PBS) ≈ ?

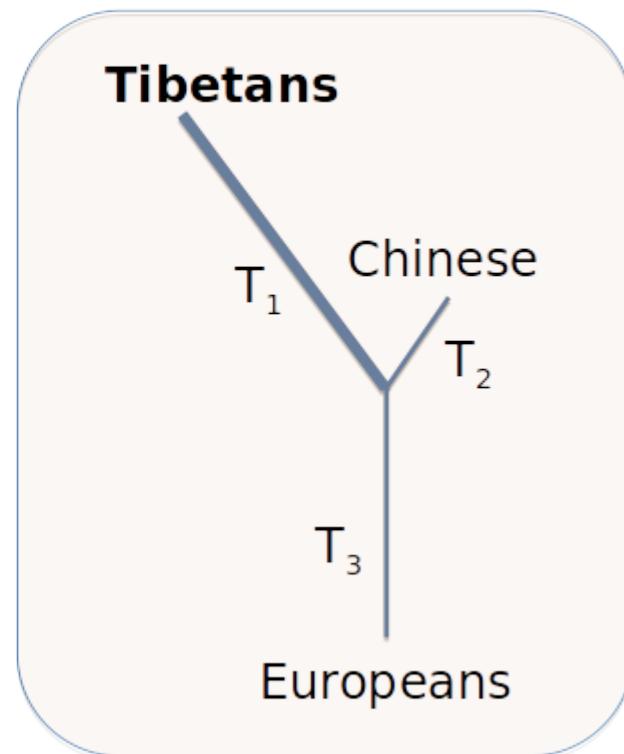
Yi et al. Science 2010

Population Branch Statistic

Neutral evolution

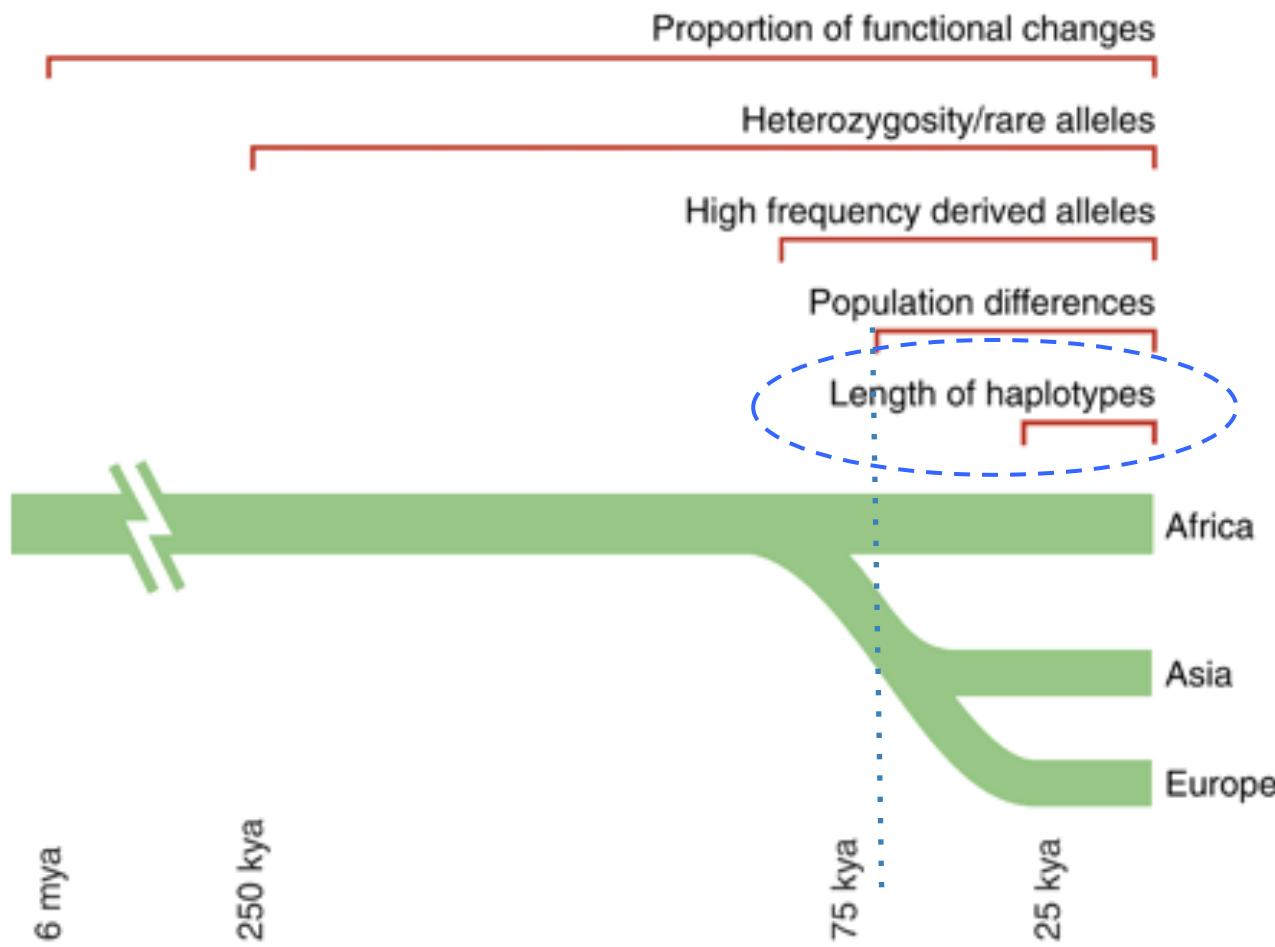


Positive selection



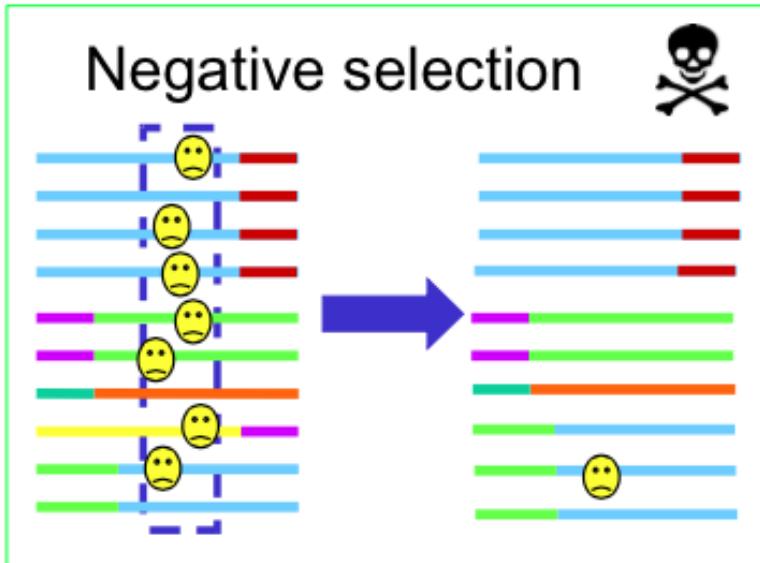
$$\text{Population Branch Statistic (PBS)} \approx T_{12} + T_{13} - T_{23}$$

Methods to infer recent selection

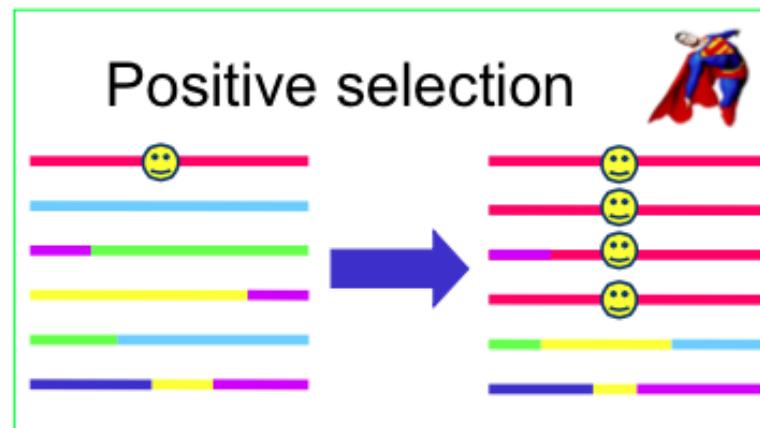


Selective sweep

Effect of genetic hitchhiking

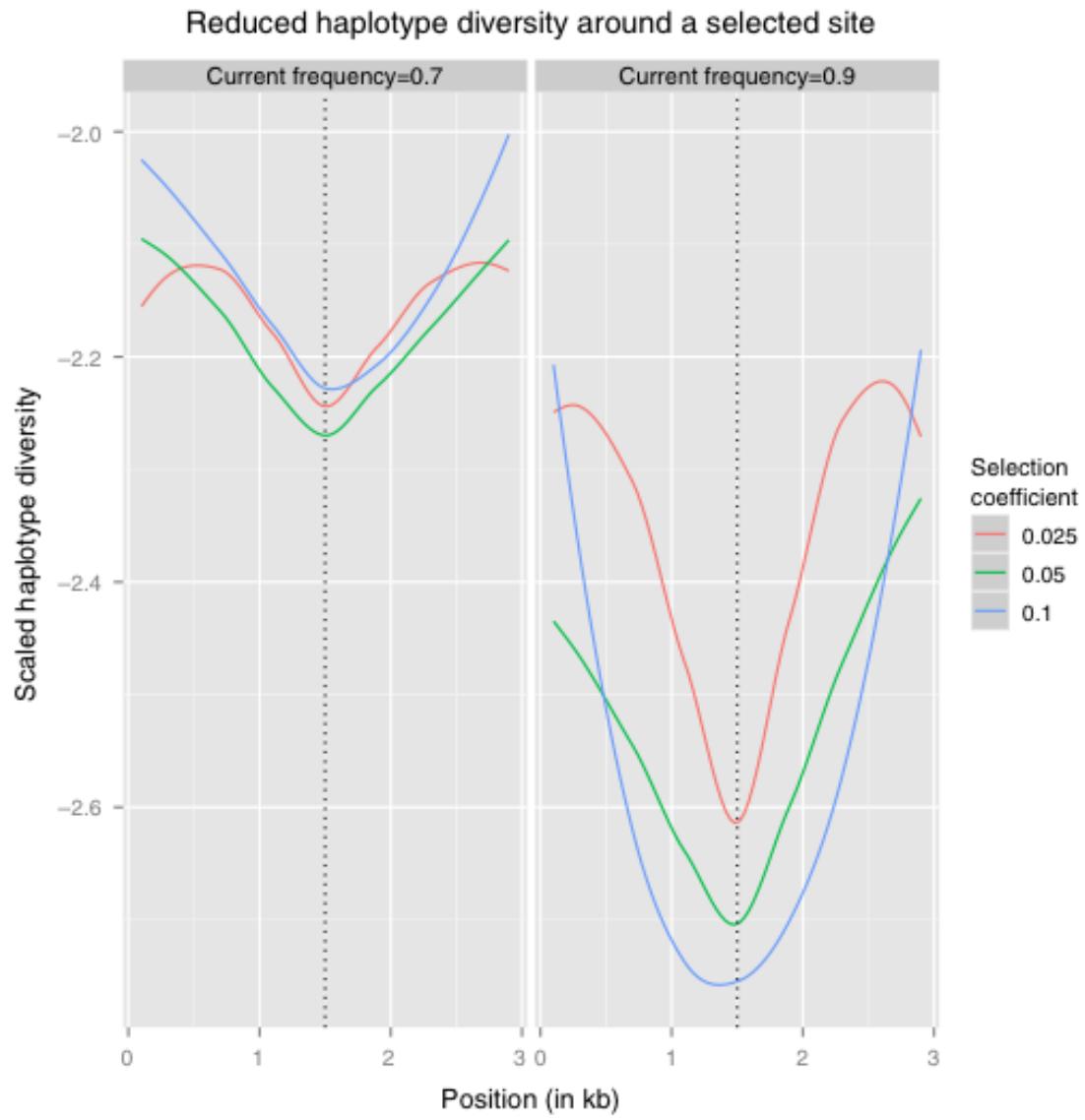


- Selective sweep: the beneficial allele is driven to fixation
- Partial sweep: when the beneficial allele hasn't reached fixation yet.



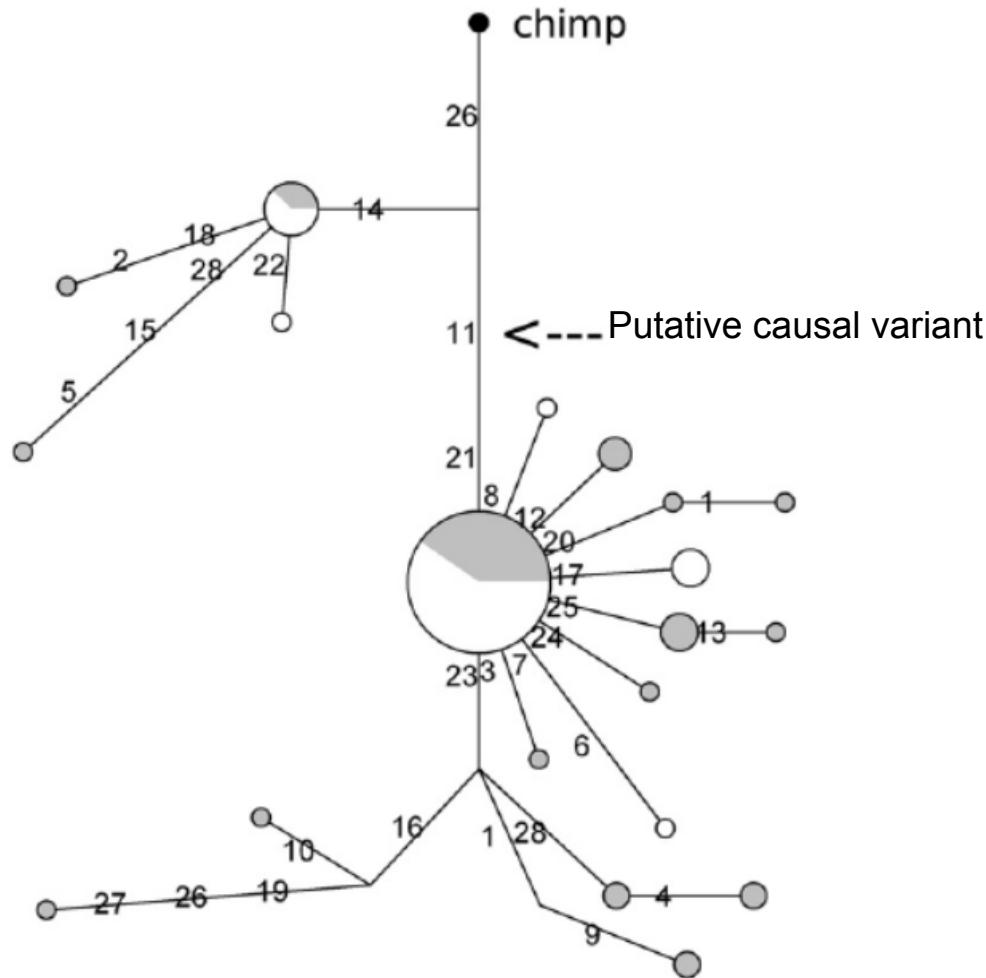
Haplotype diversity under selection

- Reduction of polymorphisms
- Decay of haplotype diversity is proportional to selective strength
- Recombination breaks the signal
- The effect is stronger in proximity of the selected site

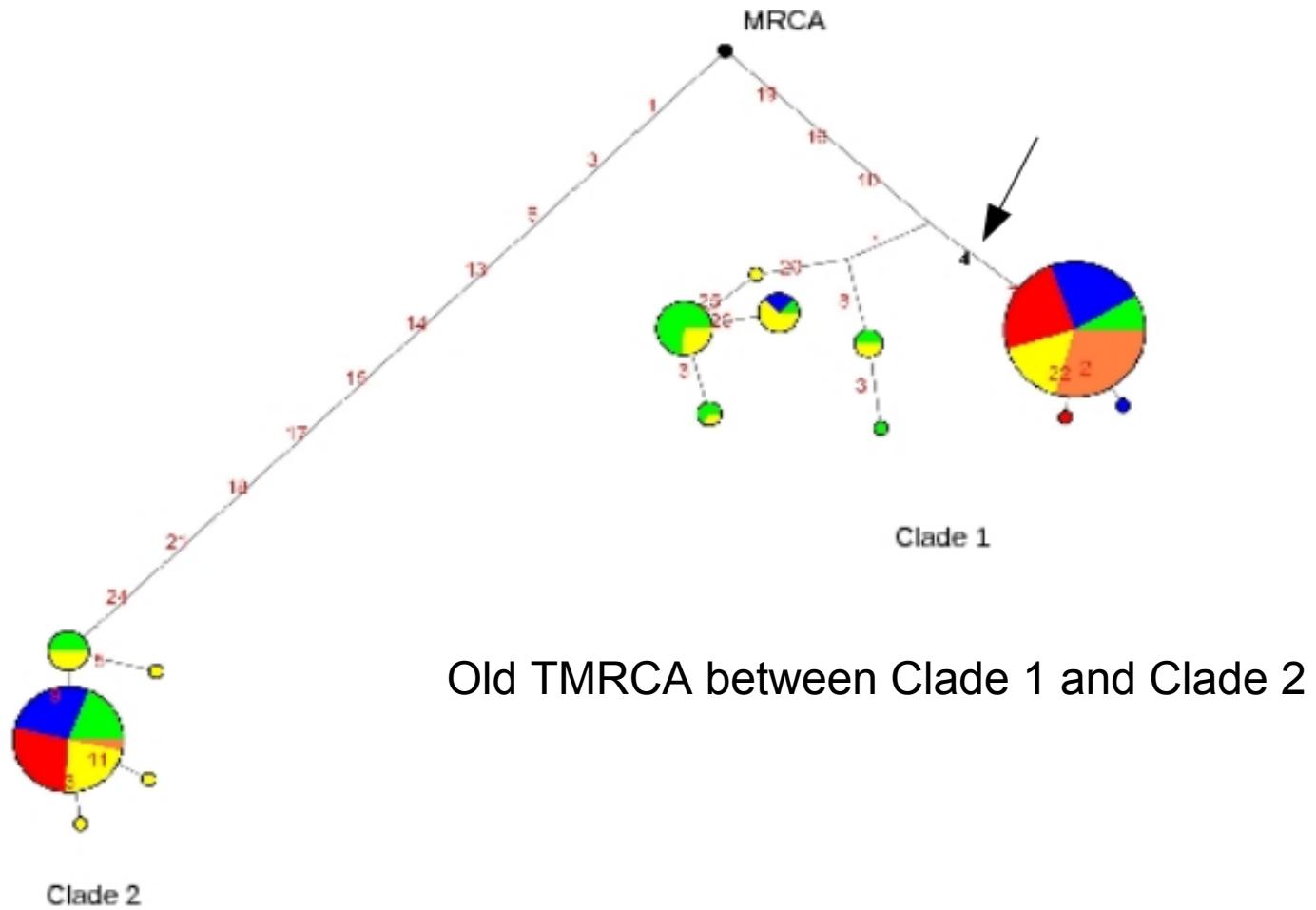


Haplotype genealogy

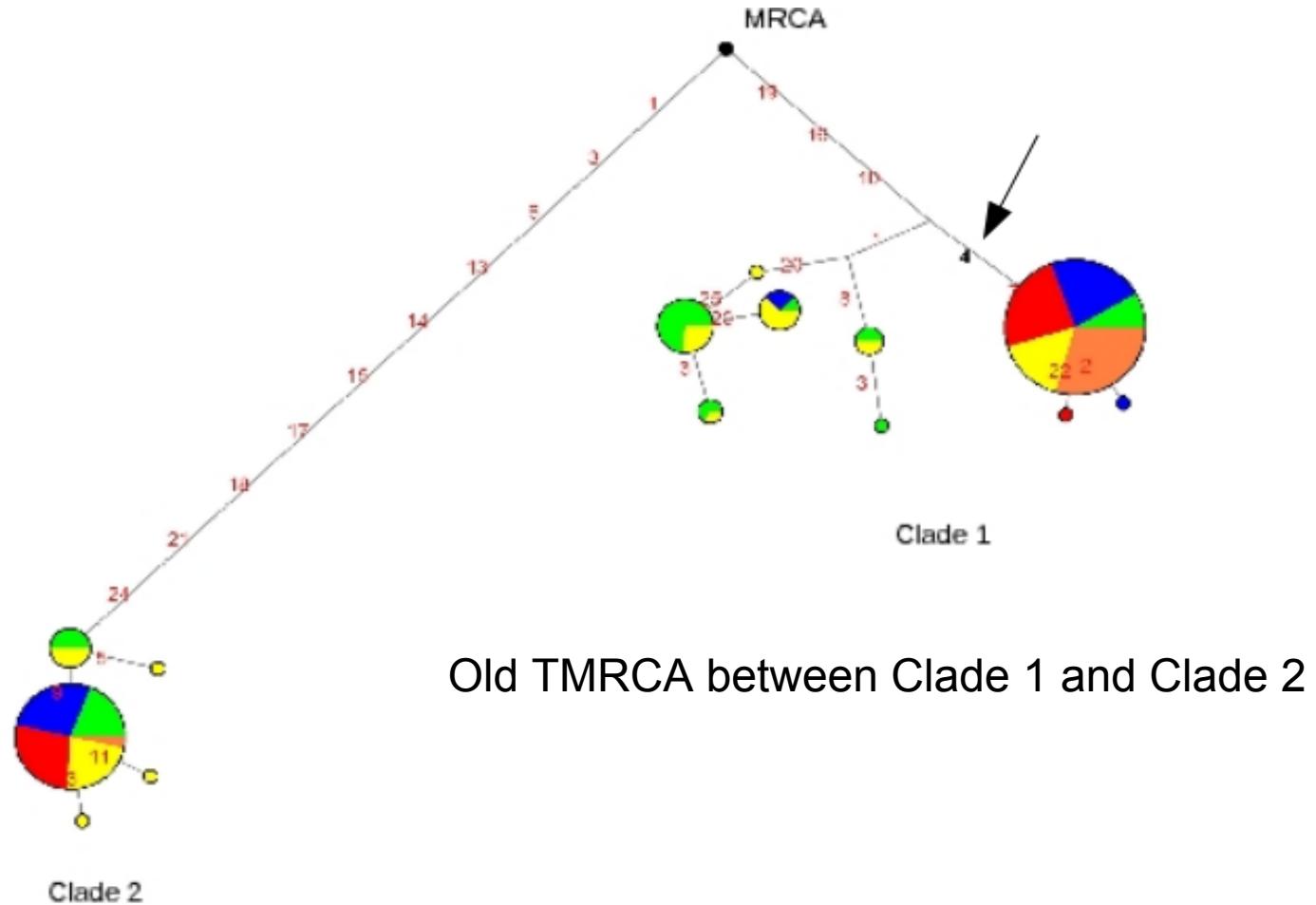
Median-joining network



Haplotype genealogy (under ? selection)

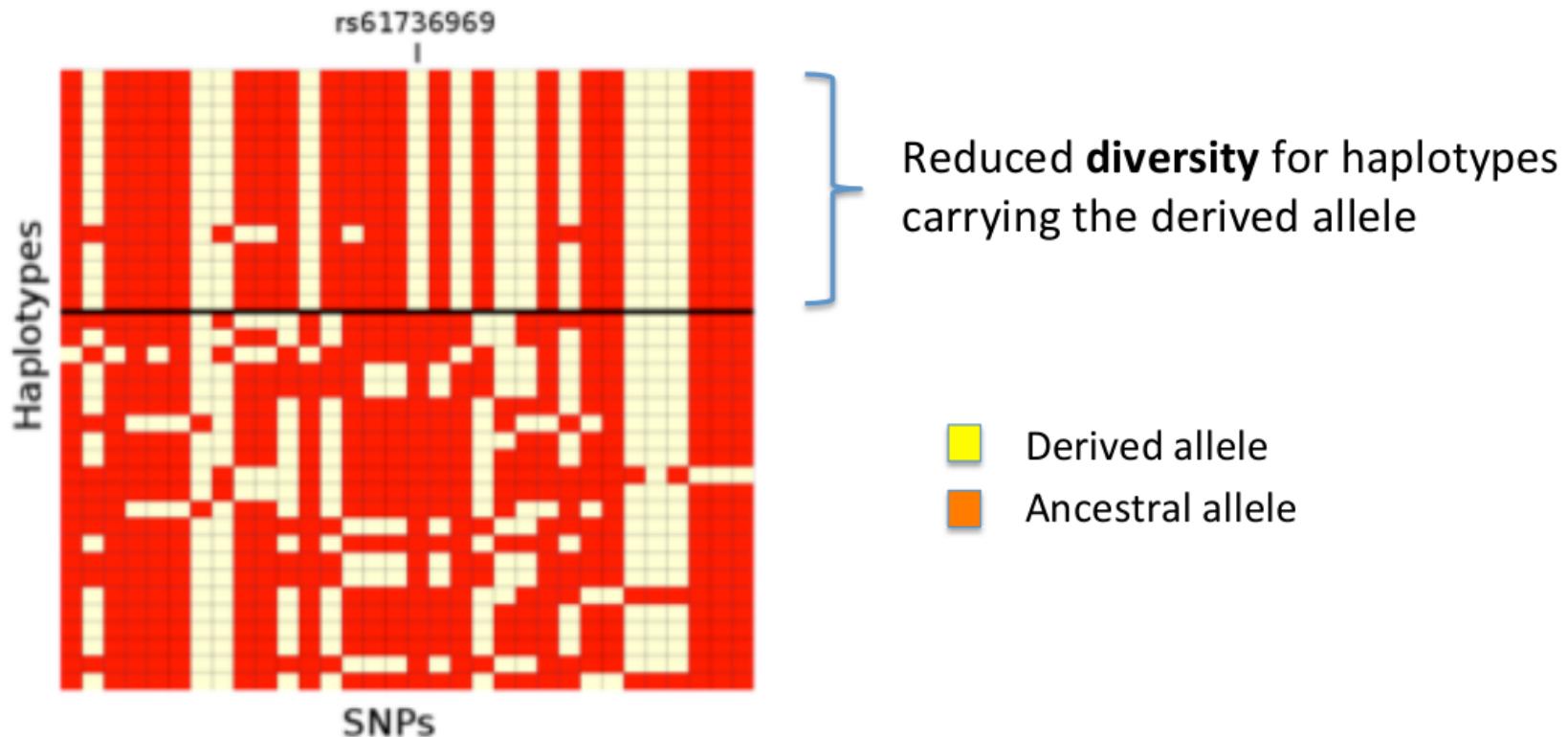


Haplotype genealogy (under balancing selection)

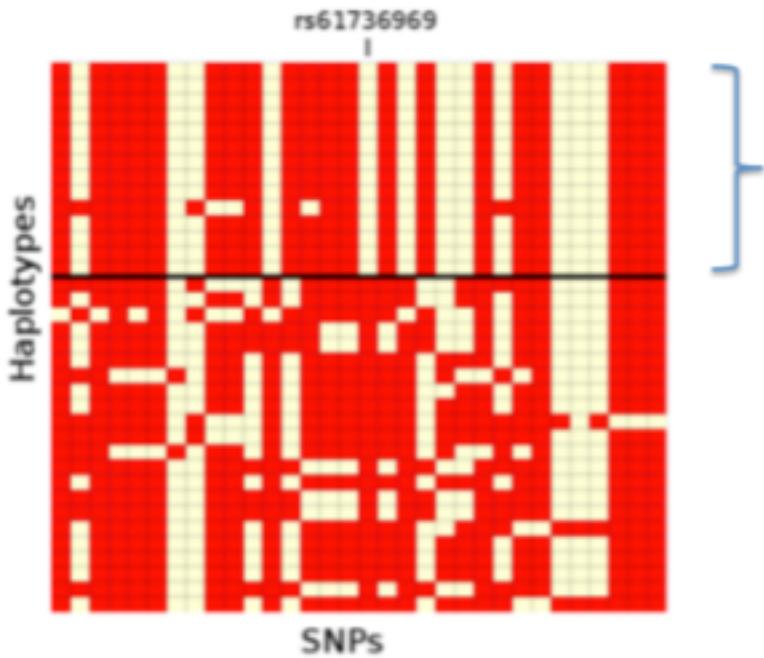


Haplotype diversity

Quantification of the haplotype diversity



Haplotype diversity

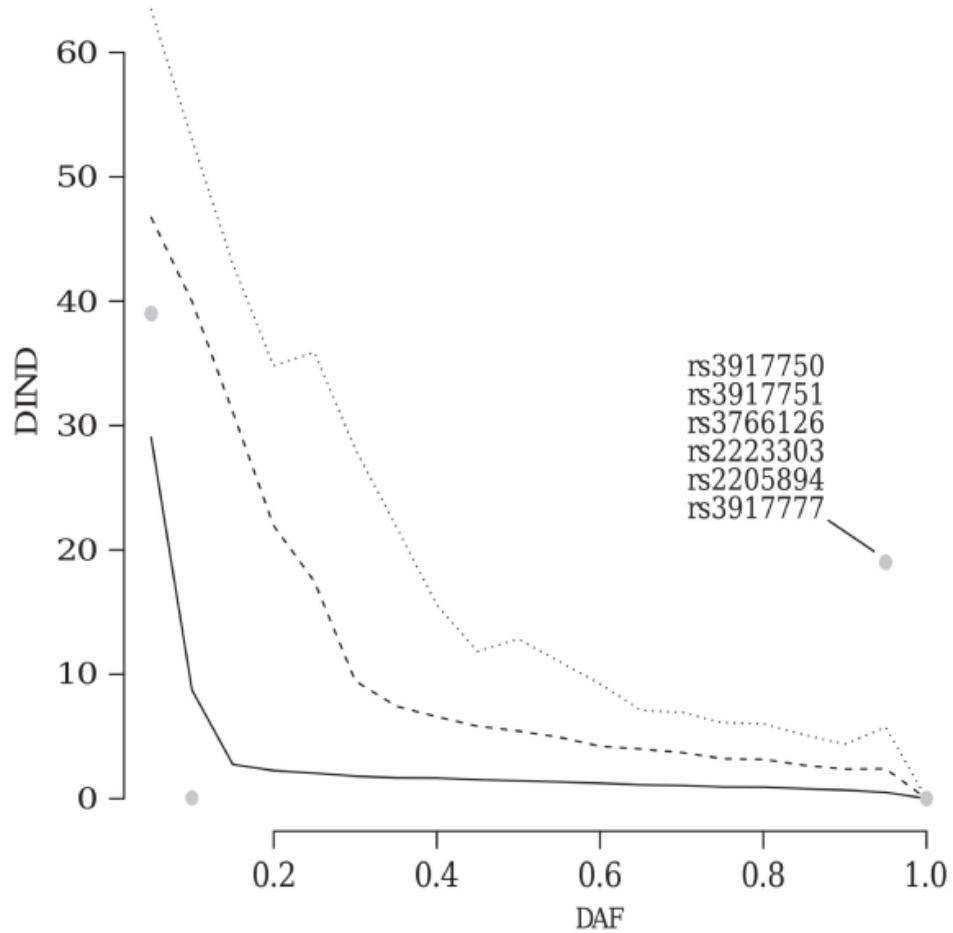


DIND (Derived Intra-allelic Nucleotide Diversity):

$$i\pi_{D,x} = \frac{\sum_{k=1}^{n_D-1} \sum_{l=k+1}^{n_D} d_{kl}}{n_D C_2}$$

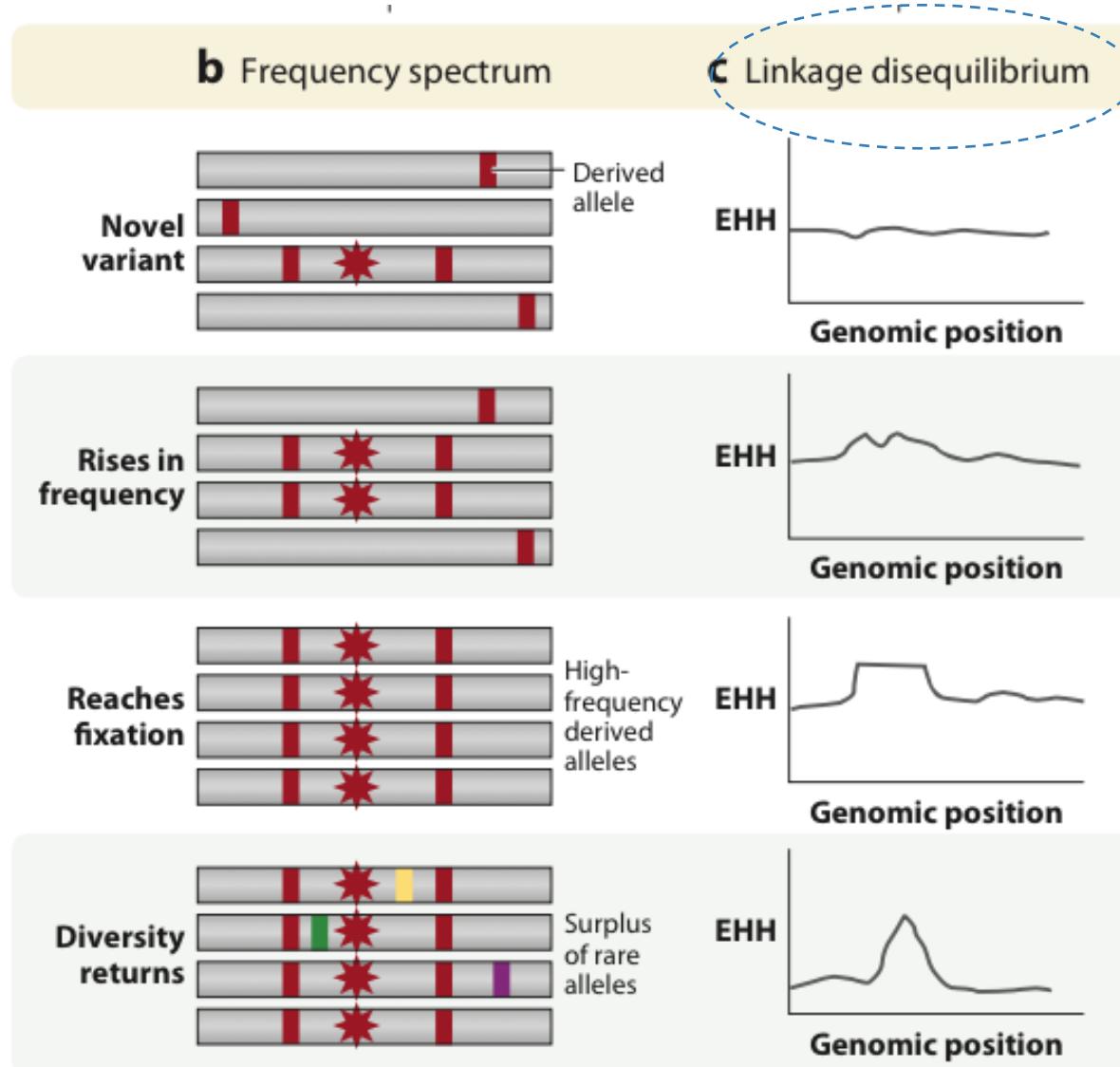
DIND ~ Diversity of haplotypes carrying the derived allele / Diversity of haplotypes carrying the ancestral allele

DIND

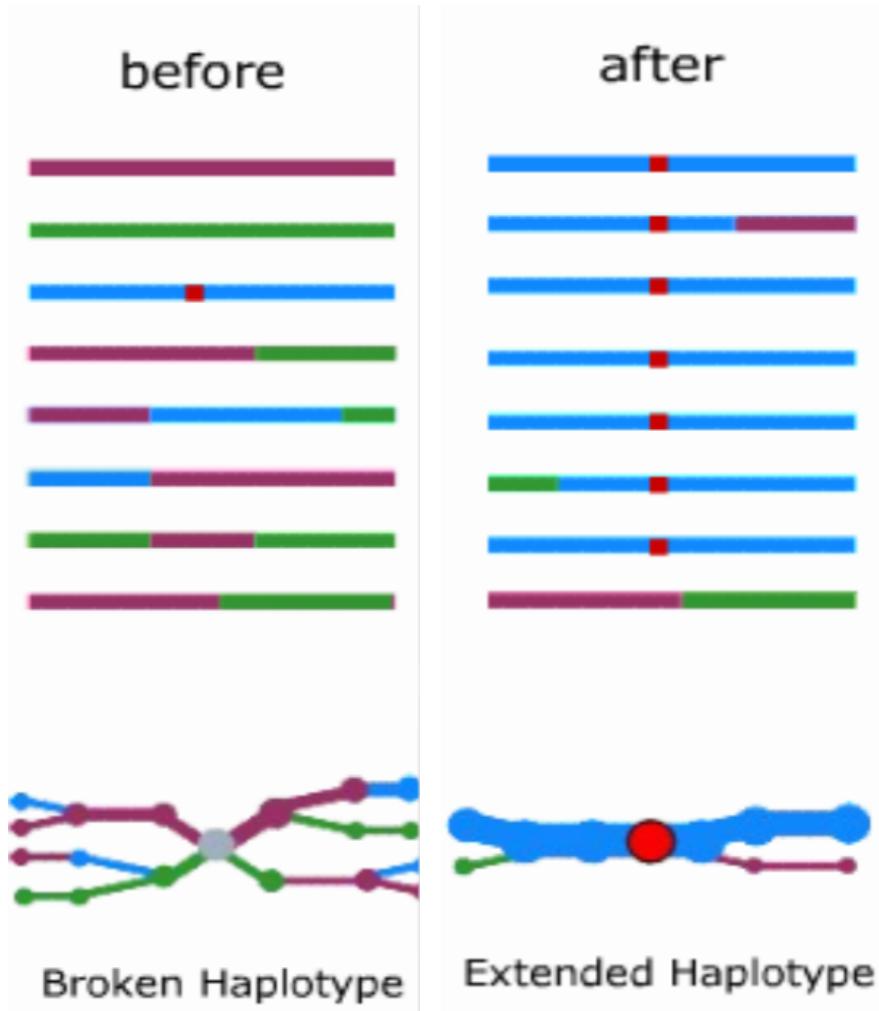


DIND is useful to identify the causal variant by looking at outliers in the DIND/DAF distribution

Extended Linkage Disequilibrium



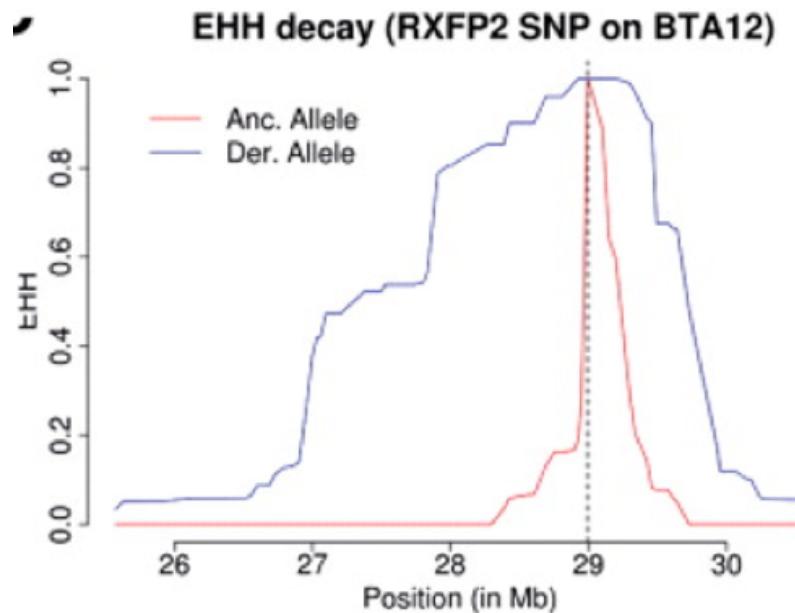
Extended Haplotype Homozygosity (EHH)



LD-based approaches:

- identifying variants under partial or incomplete selective sweep
- EHH from a core allele to a specified distance
- EHH decreases with distance by the action of recombination

EHH-derived statistics



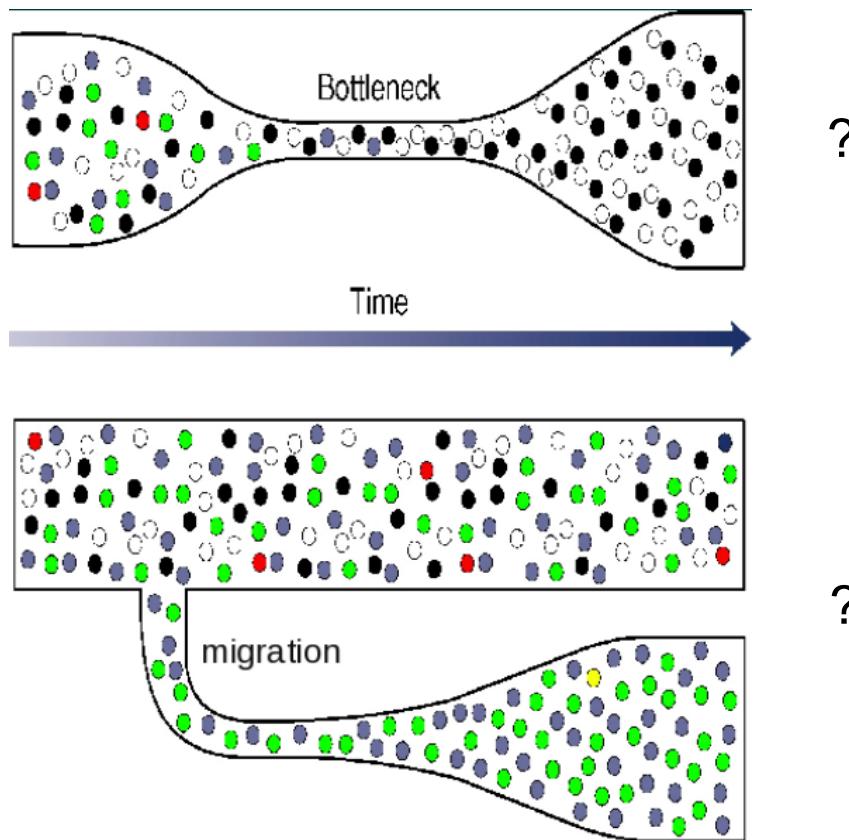
- ❑ **iHS** compares the area under the curve defined by EHH for the derived and ancestral variants as one travels further in genetic distance from the core region
- ❑ cross-population extended haplotype homozygosity (**XP-EHH**) compares haplotype lengths between populations.

Software available

- DnaSP (<http://www.ub.edu/dnasp/>)
- Arlequin (<http://cmpg.unibe.ch/software/arlequin35/>)
- BayeScan (<http://cmpg.unibe.ch/software/BayeScan/>)
- libsequence (<http://www.molpopgen.org/software.html>)
- sweep (<http://www.broadinstitute.org/mpg/sweep/>)
- iHS (<http://coruscant.itmat.upenn.edu/software.html>)
- nSL (<http://cteg.berkeley.edu/~nielsen/resources/software/>)
- Pre-computed values (USCS genome browser tables)
- Homemade scripts
- ...

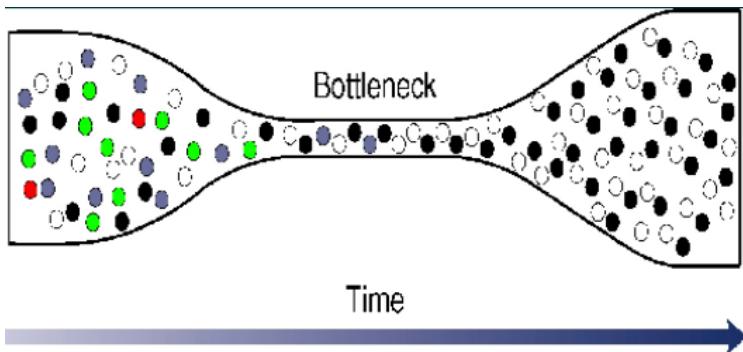
Neutral confounding factors

Various demographic scenarios may result in the rejection of the null hypothesis of neutrality.

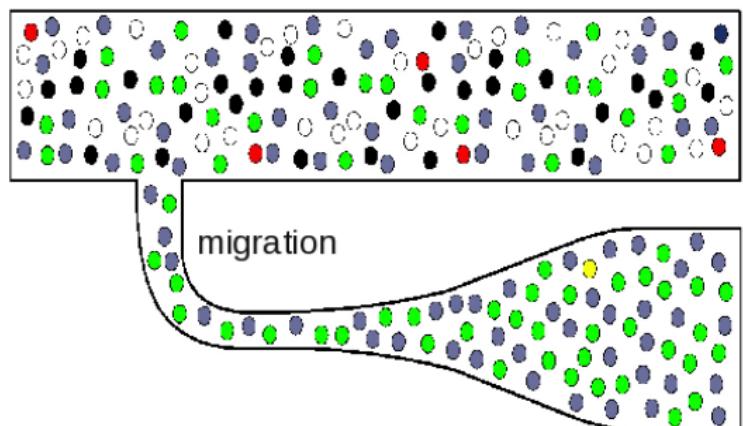


Neutral confounding factors

Various demographic scenarios may result in the rejection of the null hypothesis of neutrality.



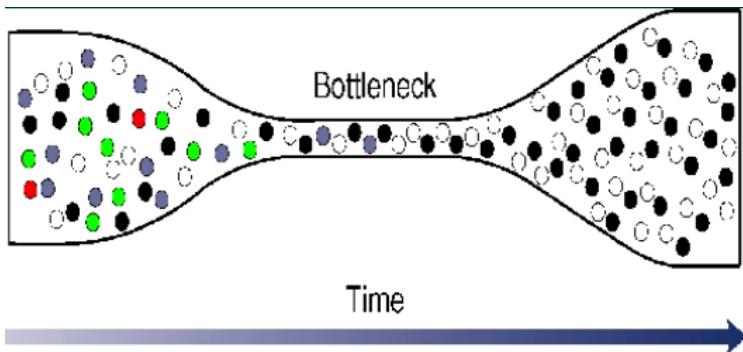
- Bottleneck: excess of intermediate frequency alleles



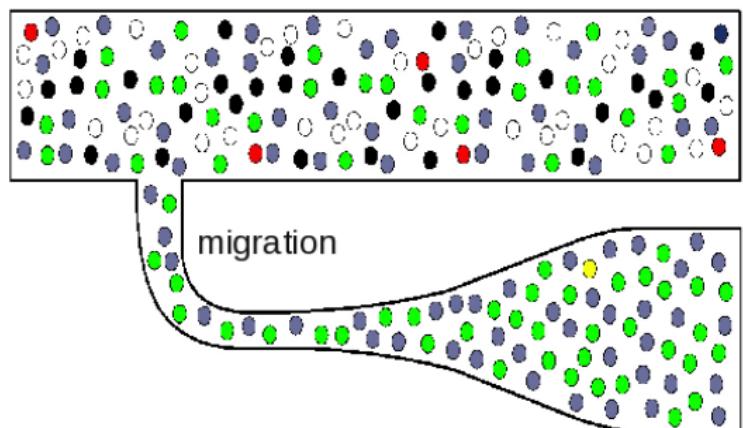
- Founder effect: decreased diversity
- Expansion: ?

Neutral confounding factors

Various demographic scenarios may result in the rejection of the null hypothesis of neutrality.



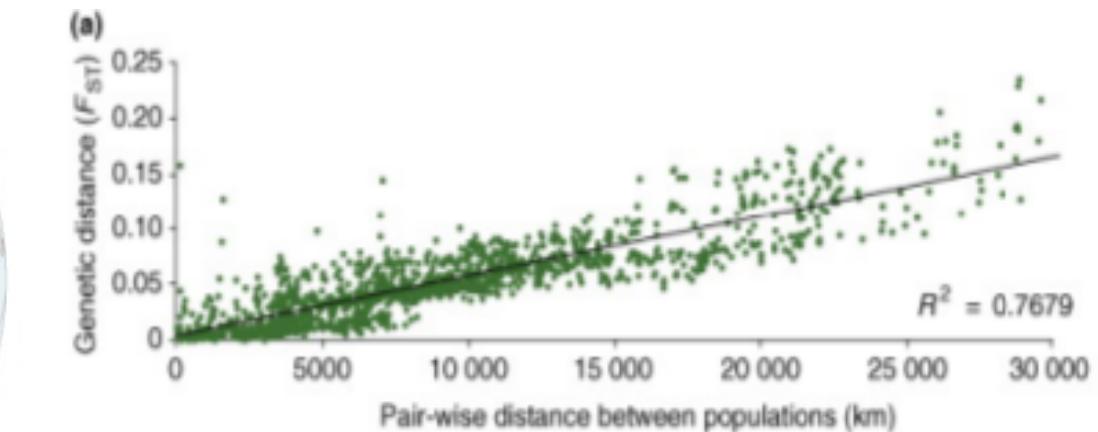
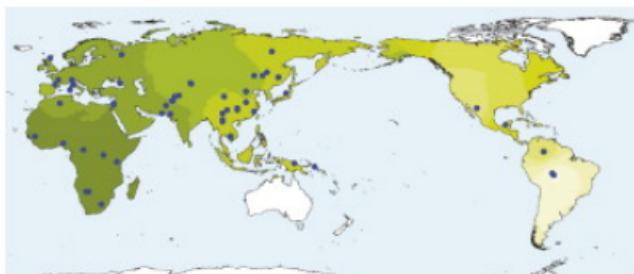
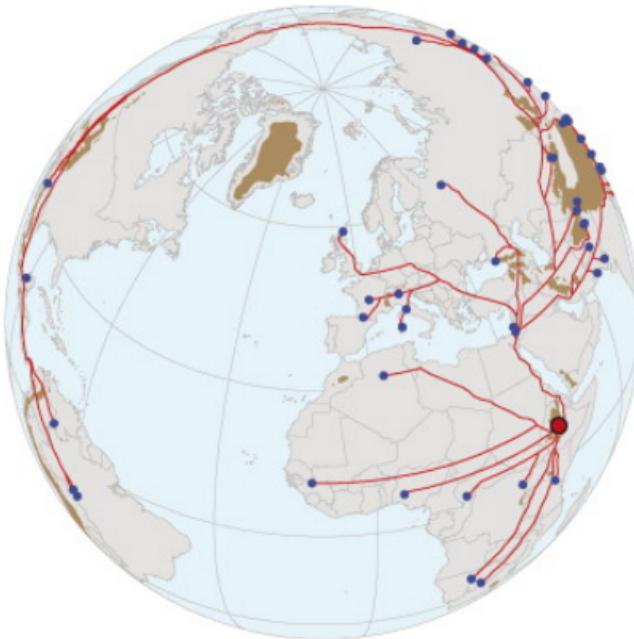
- Bottleneck: excess of intermediate frequency alleles



- Founder effect: decreased diversity
- Expansion: excess of low frequency alleles

Isolation by distance

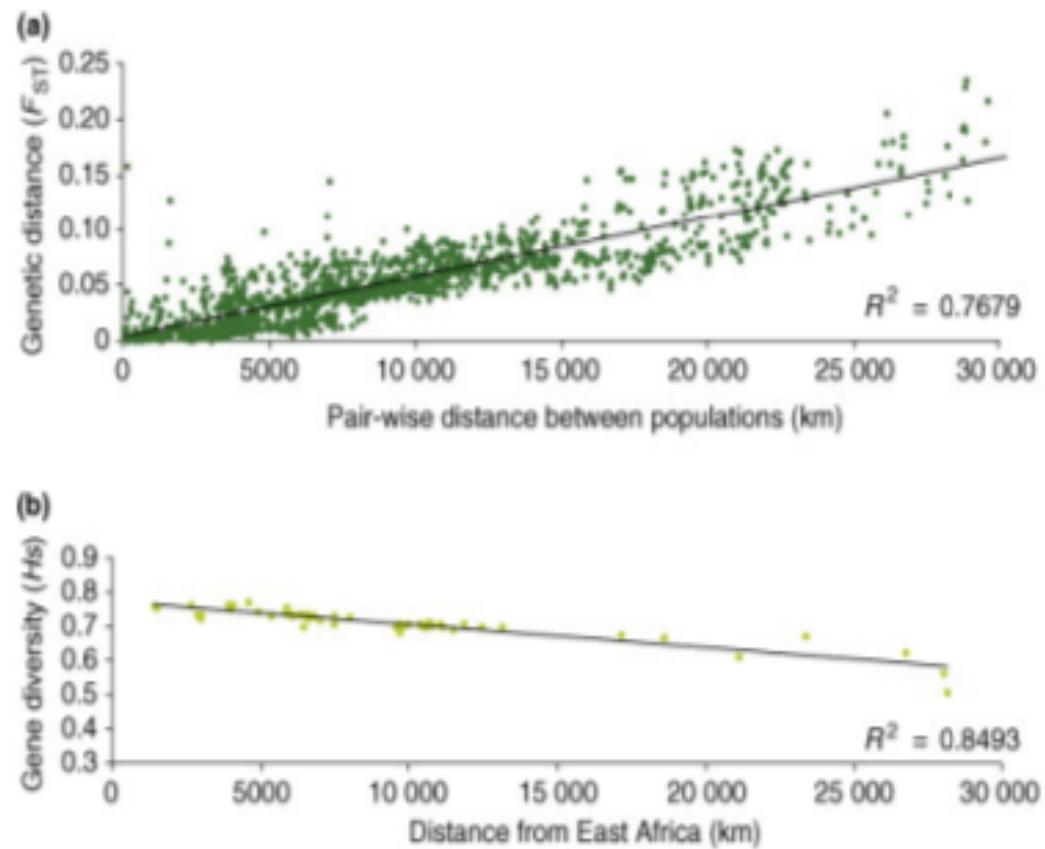
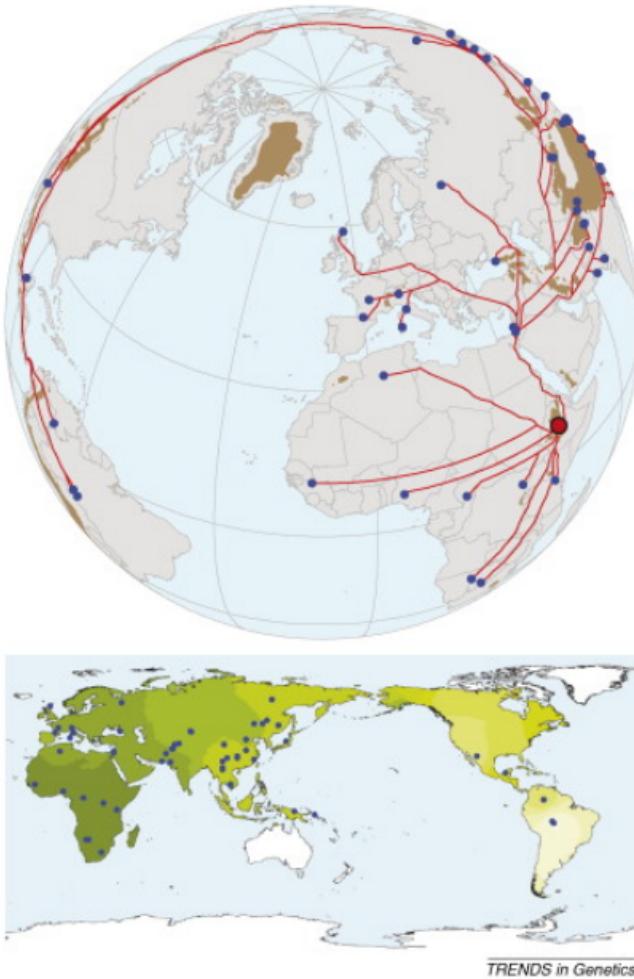
- Human populations are not independent.
- Global human genetic variation is greatly affected by geography.



TRENDS in Genetics

Isolation by distance

- Human populations are not independent.
- Global human genetic variation is greatly affected by geography.

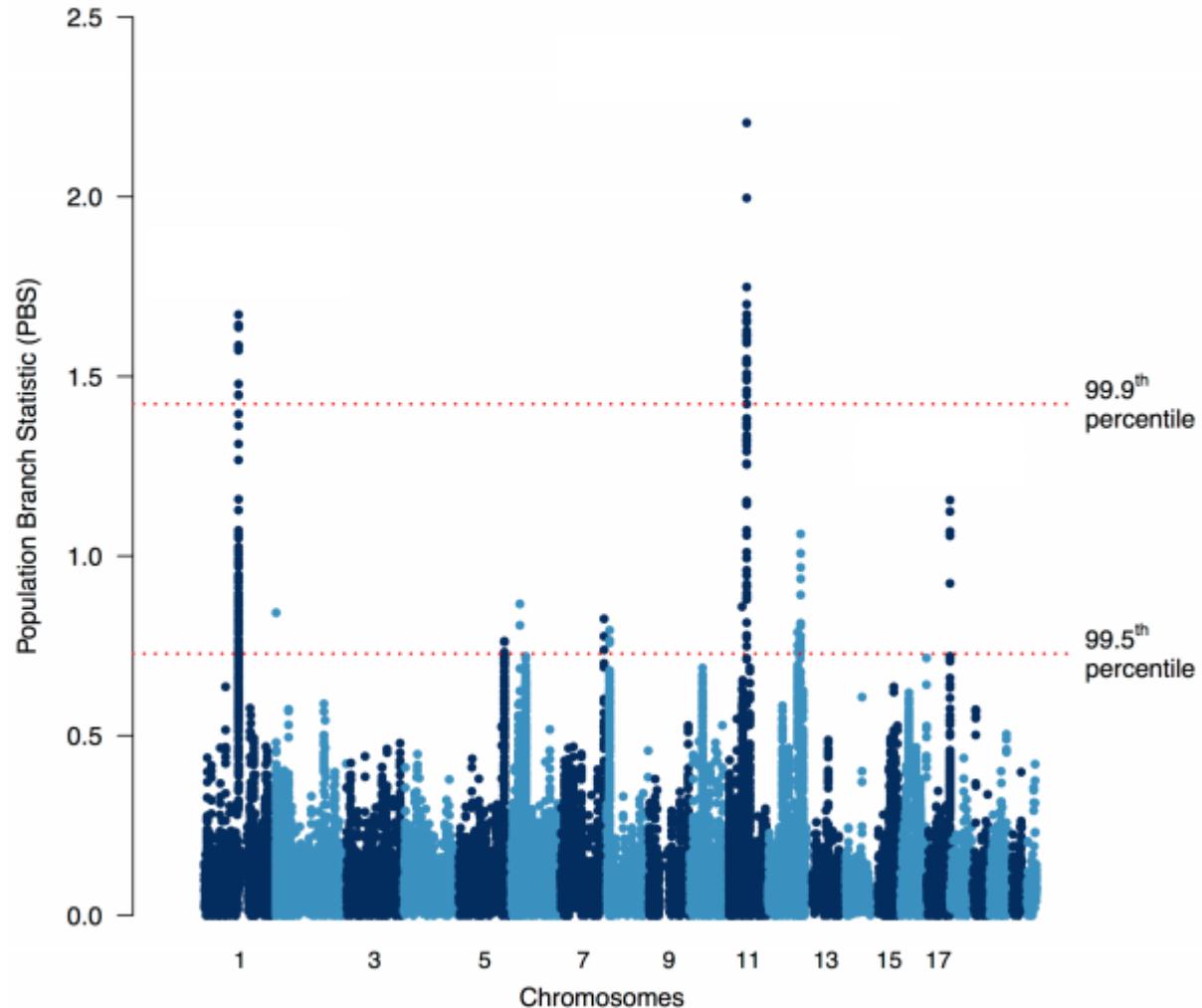


Assessing statistical significance

- 1) **Comparison** with other loci in the genome. While selection acts on a locus, demography equally affects the whole genome. We can calculate the percentile rank of the statistic for the gene/region of interest.
- 2)

Empirical distribution

Outlier approach

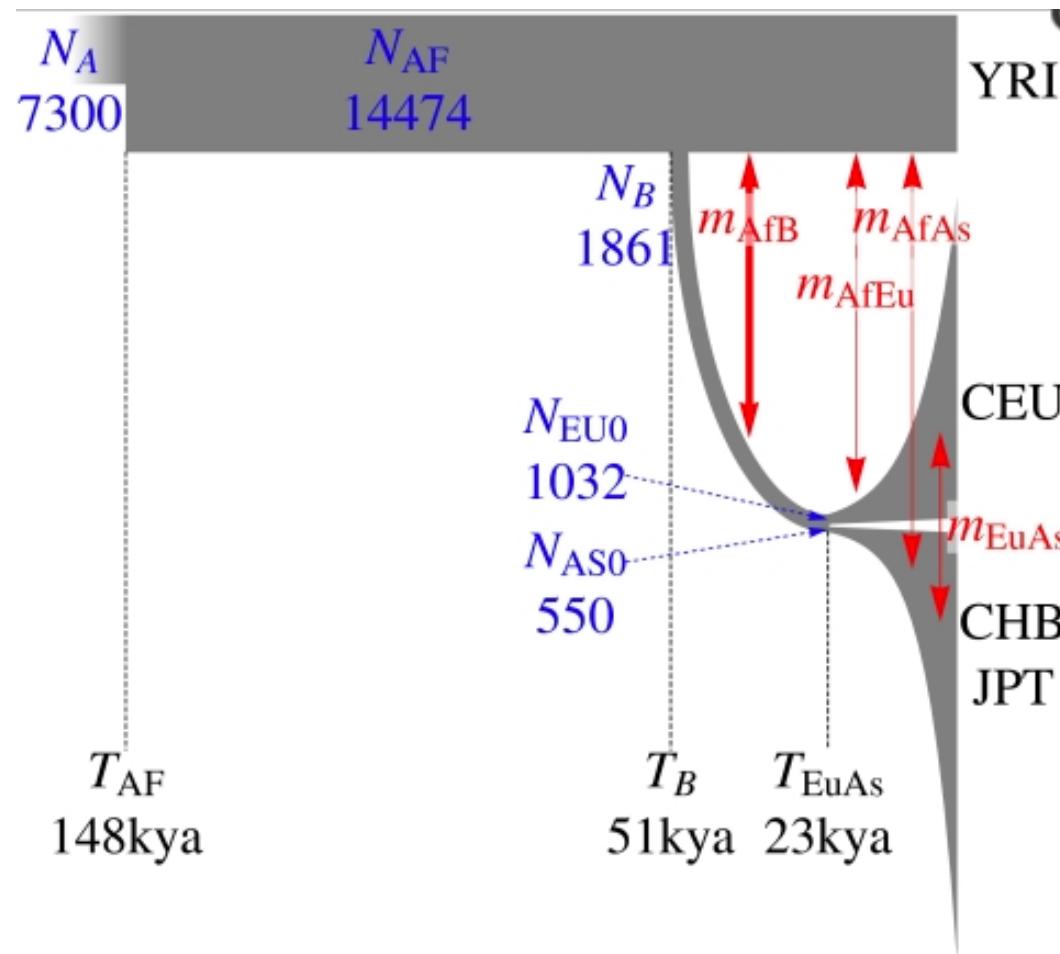


Assessing statistical significance

- 1) **Comparison** with other loci in the genome. While selection acts on a locus, demography equally affects the whole genome. We can calculate the percentile rank of the statistic for the gene/region of interest.
- 2) **Simulations.** Coalescence is the process in which, looking backwards in time, the genealogies of alleles merge at a common ancestor. We can incorporate a demographic model for species' evolution.

Demographic model

Simulations under a neutral demographic model



Adaptation to environments

Detect putative loci targeted by natural selection showing extreme or unusual values of statistics based on site frequencies spectrum, haplotypes structure, population genetic differentiation.

But:

- 1) **power** to detect a selective event only if the process is recent or strong enough or the beneficial allele is nearly fixed to leave a clear footprint;
- 2) unable to underline the **environmental** factor (pathogens, diet or climate) that acted as selective **pressure**.

Adaptation to environments

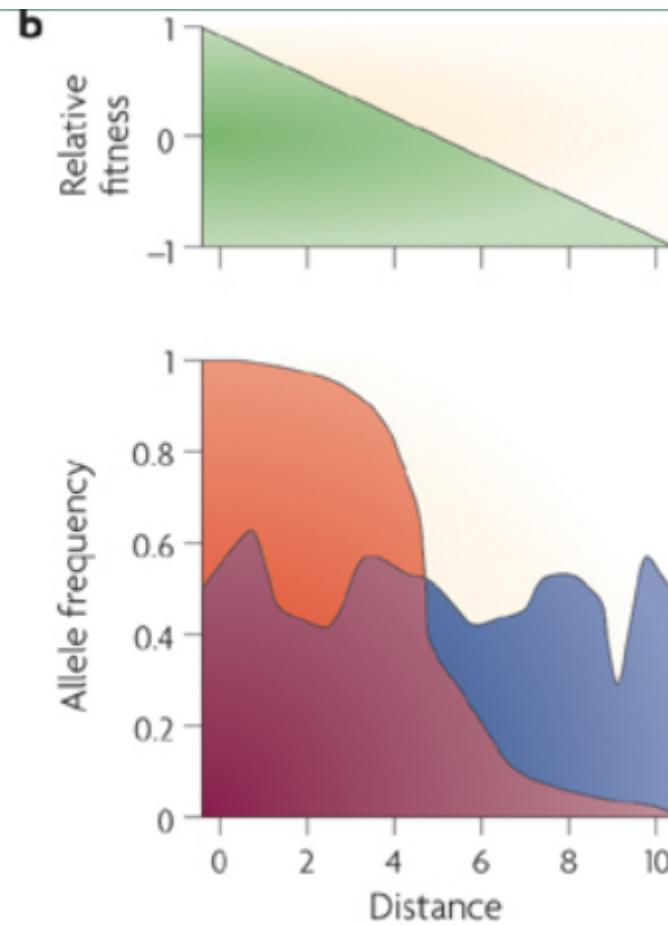
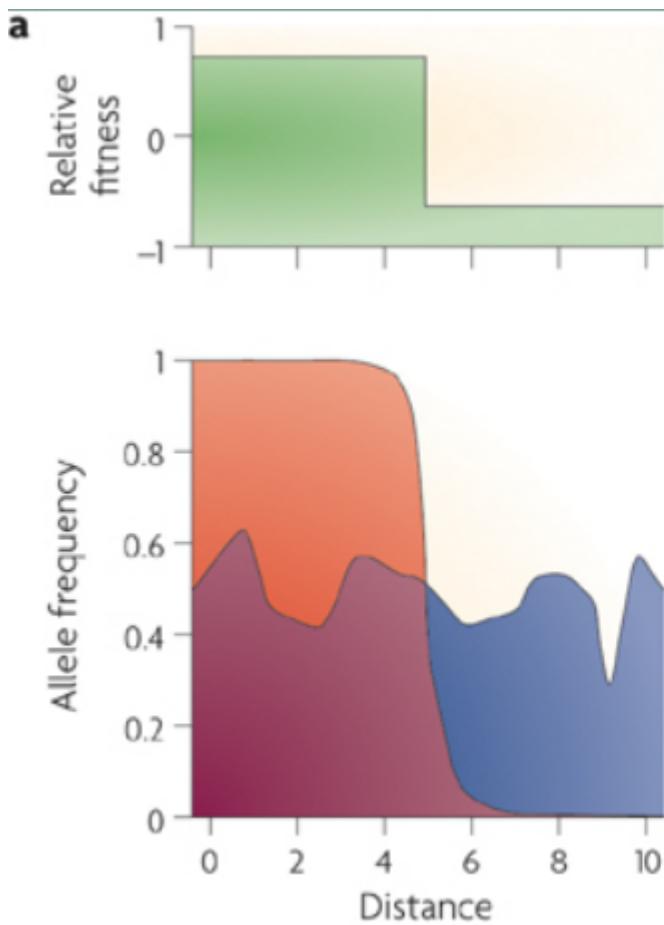
But:

- 1) **power** to detect a selective event only if the process is recent or strong enough or the beneficial allele is nearly fixed to leave a clear footprint;
- 2) unable to underline the **environmental** factor (pathogens, diet or climate) that acted as selective **pressure**.

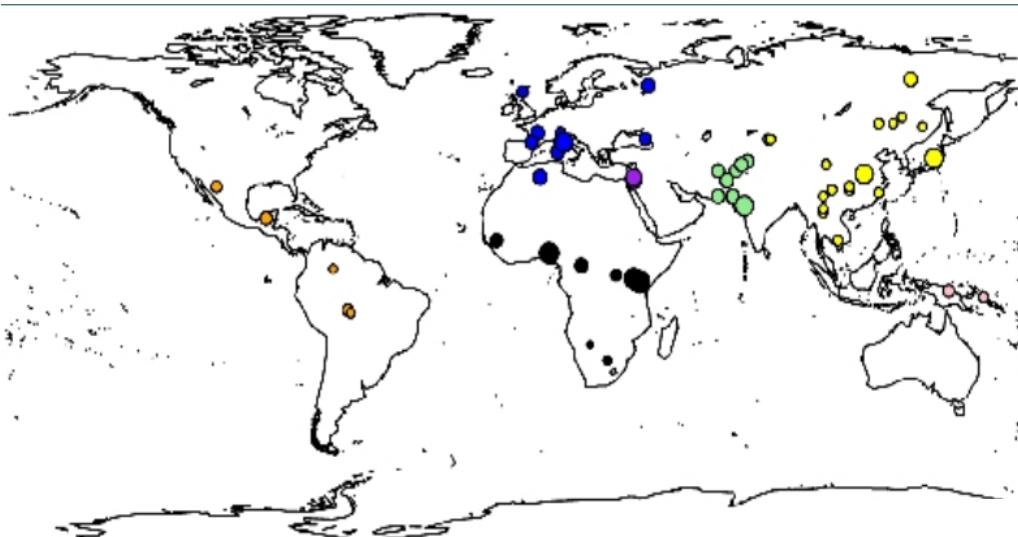
Solutions:

- 1) Alternative models might better explain the complex signatures of adaptation (selective sweeps from standing variation or polygenic adaptation)
- 2) Methods detecting polymorphisms whose frequencies are greatly correlated with environmental variables.

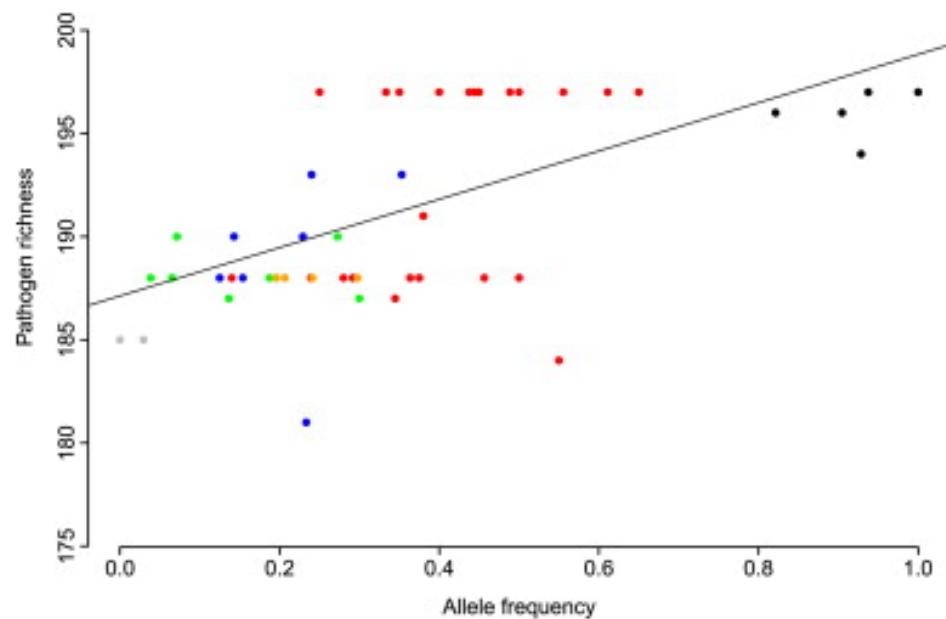
Spatial variation



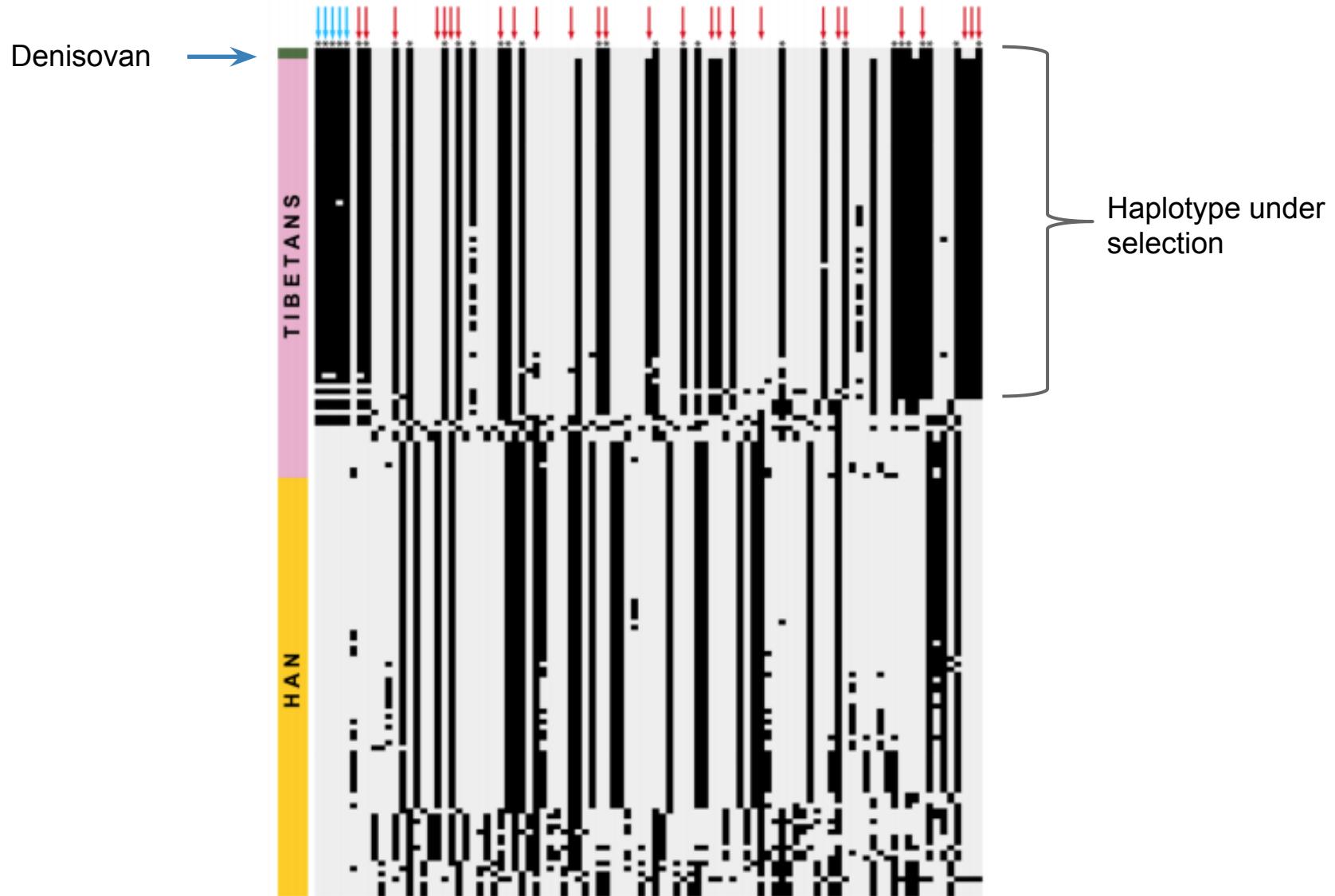
Correlations



- Mantel correlations
- Bayesian methods (BayEnv)
- Mixed models (LFMM)



Adaptive introgression



Adaptive admixture

Figure 2: Tibetans as a mixture of the HA-proxy and Han Chinese-related ancestral populations in the scaffold tree.

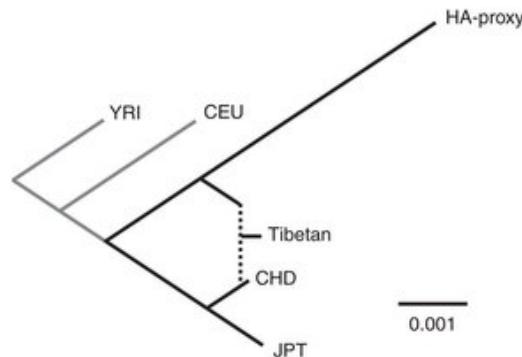
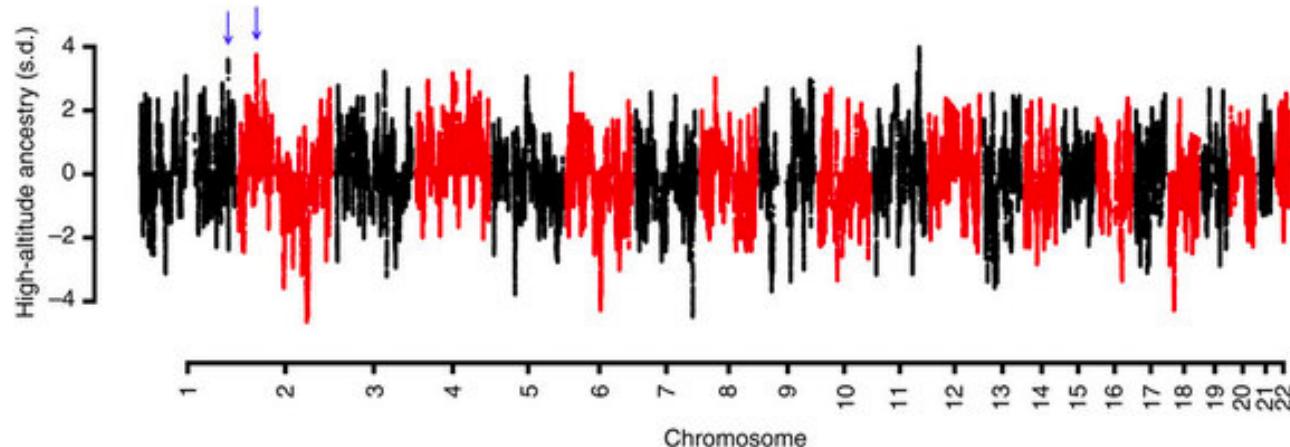


Figure 4: The distribution of high-altitude ancestry proportions across the Tibetan genome.

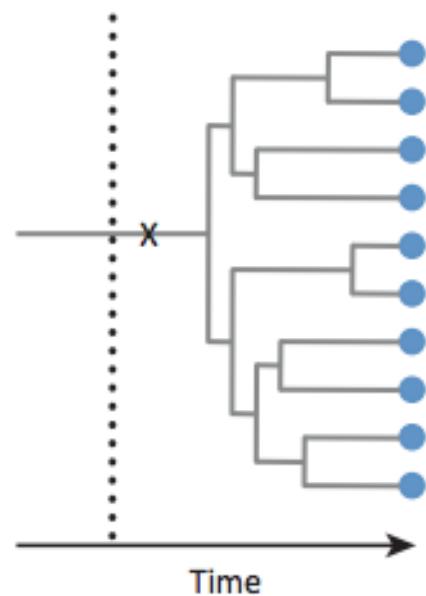


Blue arrows mark the positions of *EGLN1* (in chromosome 1) and *EPAS1* (in chromosome 2).

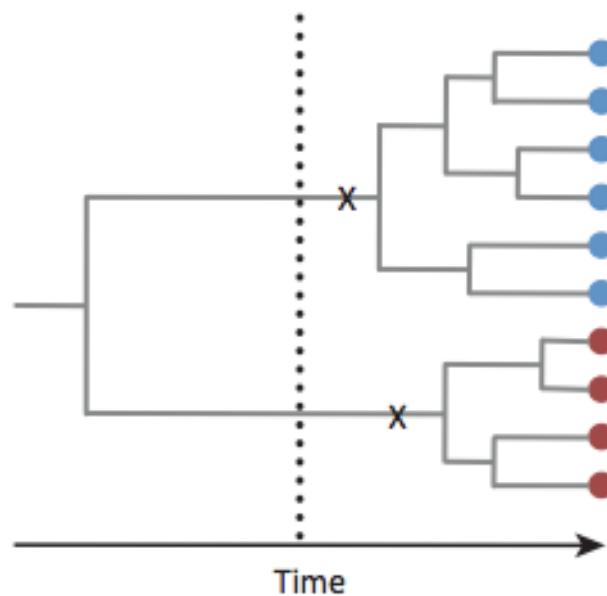
Jeong et al. Nature Comm. 2014

Soft sweep

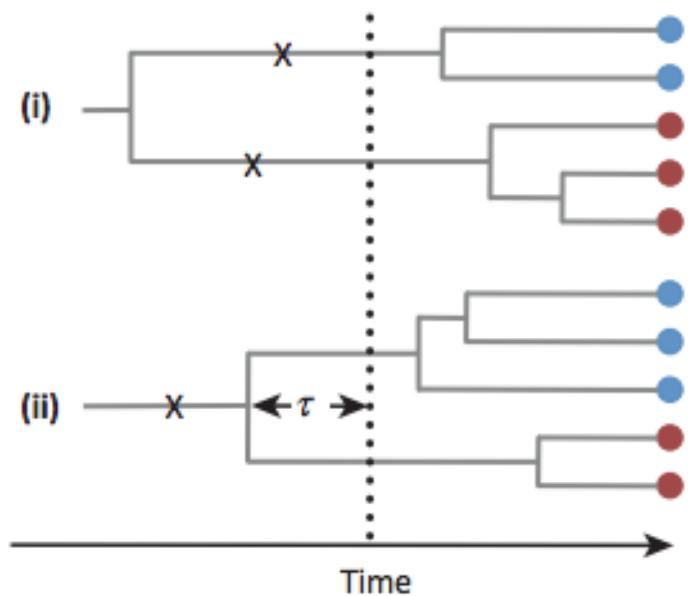
(A) Classic hard sweep



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



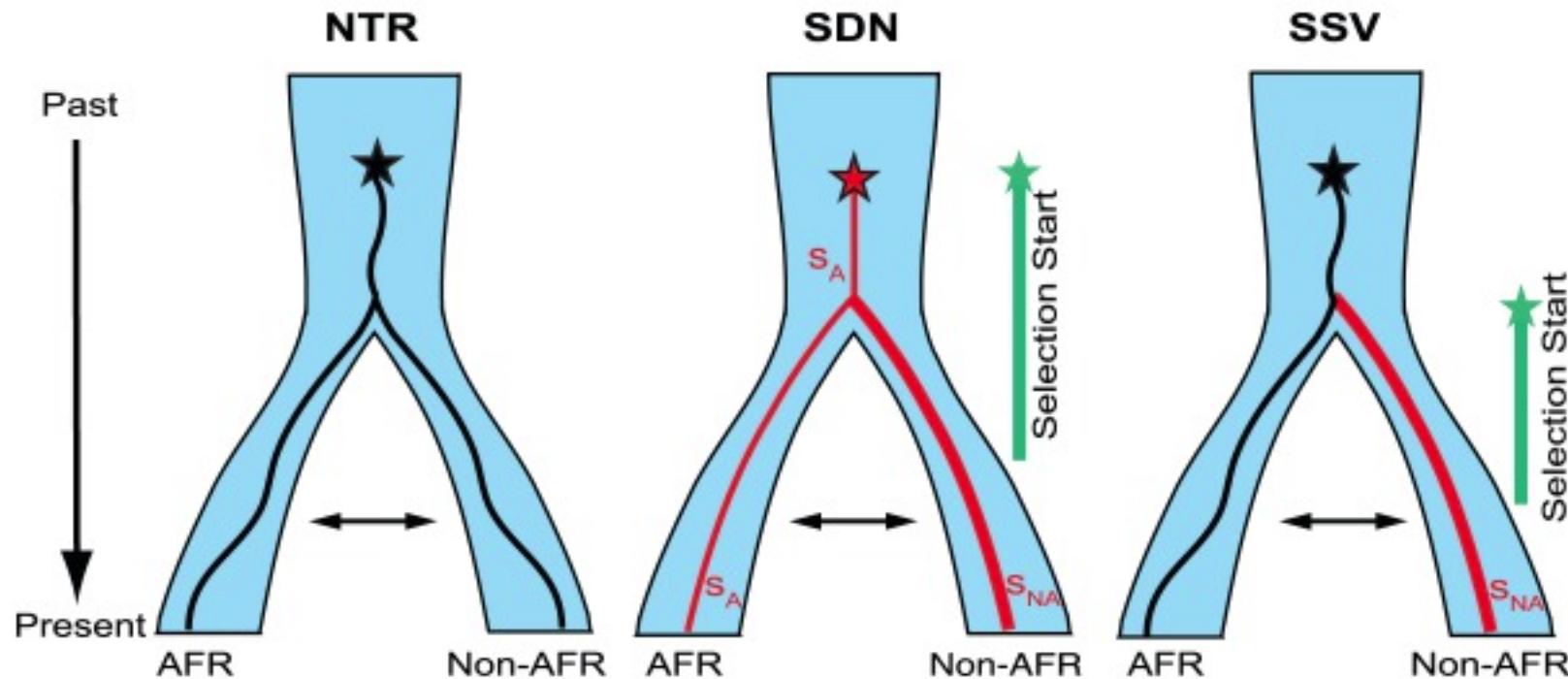
(C) Soft sweep (standing variation)



TRENDS in Ecology & Evolution

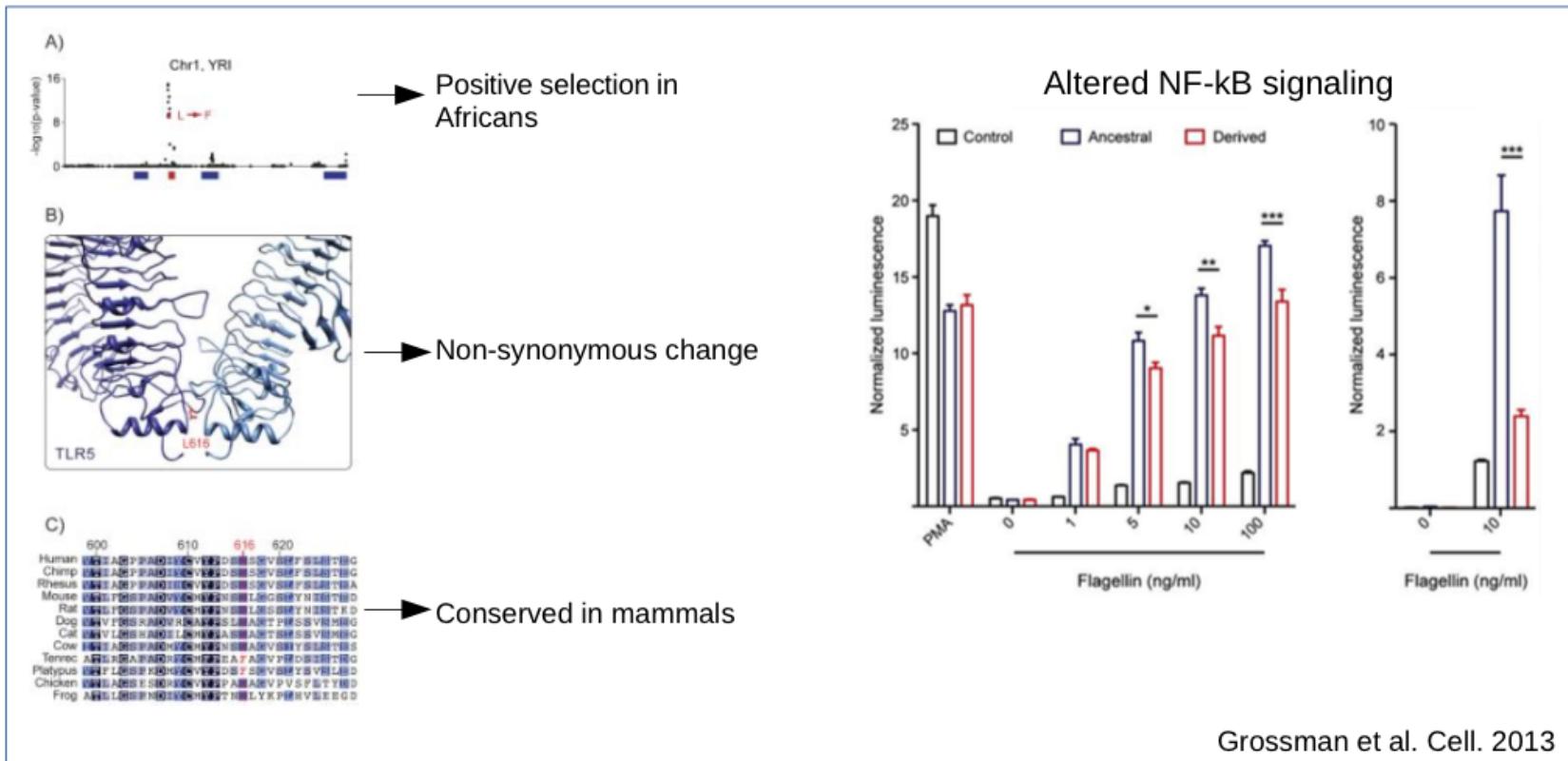
Simulation-based approaches

Approximate Bayesian Computation (ABC)



Future perspectives

Integrative approach with functional validation of candidate adaptive variants



Paper discussion



Cell

Population Genomics Reveal Recent Speciation and Rapid Evolutionary Adaptation in Polar Bears

Shiping Liu,^{1,2,20} Eline D. Lorenzen,^{3,4,20} Matteo Fumagalli,^{3,20} Bo Li,^{1,20} Kelley Harris,⁵ Zijun Xiong,¹ Long Zhou,¹ Thorfinn Sand Korneliussen,⁴ Mehmet Somel,^{3,21} Courtney Babbitt,^{6,7,22} Greg Wray,^{6,7} Jianwen Li,¹ Weiming He,^{1,2} Zhuo Wang,¹ Wenjing Fu,¹ Xueyan Xiang,^{1,8} Claire C. Morgan,⁹ Aoife Doherty,¹⁰ Mary J. O'Connell,⁹ James O. McInerney,¹⁰ Erik W. Born,¹¹ Love Dalén,¹² Rune Dietz,¹³ Ludovic Orlando,⁴ Christian Sonne,¹³ Guojie Zhang,^{1,14} Rasmus Nielsen,^{1,3,15,16,*} Eske Willerslev,^{4,*} and Jun Wang^{1,16,17,18,19,*}