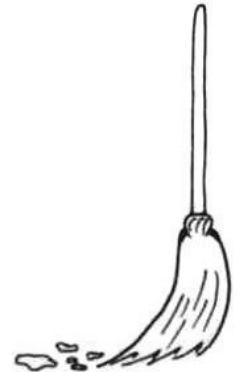
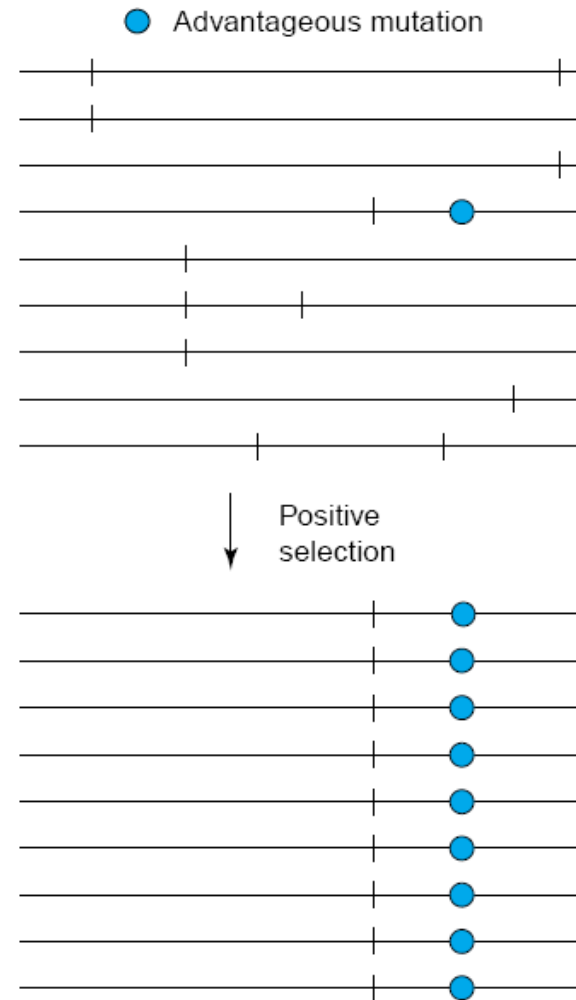


# Detection of selection from molecular data

# Recent positive selection

## Selective sweep

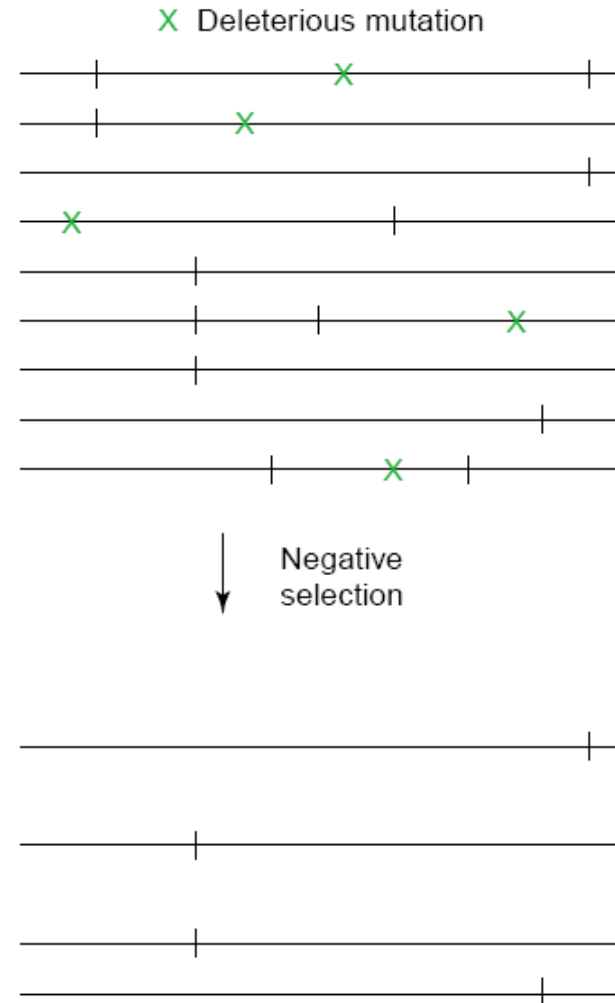
- Reduction of genetic variability around positive mutation.
- Increase of linkage disequilibrium around positive mutation.
- Changes in allele frequency spectra
- Selection not much longer than  $\sim N_e$  generations ago (in humans  $\sim 250\,000$  years).



# Recent negative selection

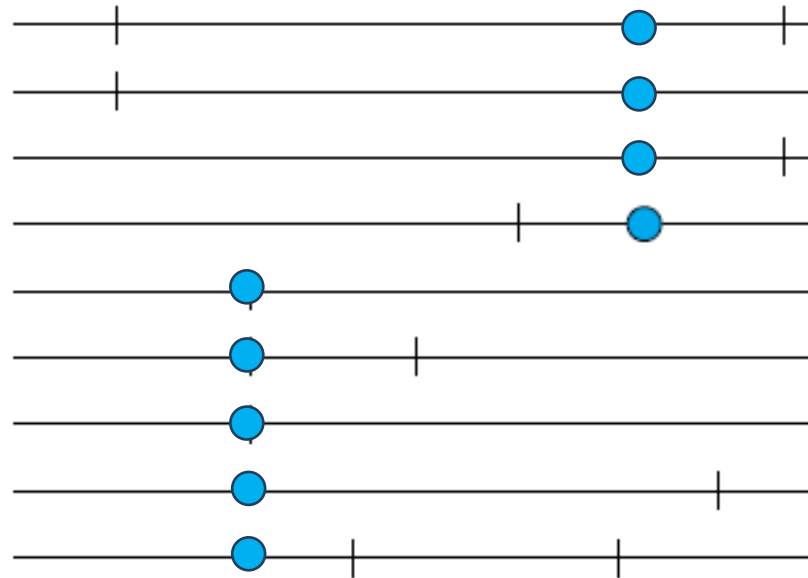
## Background selection

- Reduction of genetic variation and change in allele frequency spectra around the negative mutation, but not so marked as in case of positive selection.



# Balancing selection

- Increases levels of genetic variation
- Maintain relatively high frequencies of alternative alleles.



# Methods of detection of recent positive selection (selective sweep)

## Hudson-Kreitman-Aguadé (HKA) test

$$\theta = 4N_e\mu$$
$$D = 2\mu t$$

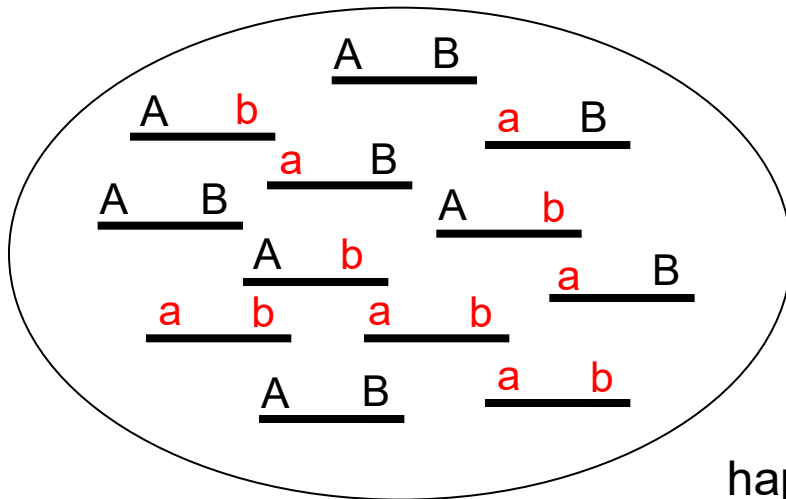
- Compares levels of within species polymorphism ( $\theta$ ) and between species divergence ( $D$ ) in multiple loci.
- For neutral sequences  $\theta/D$  ratio should be constant.
- Recent positive selection reduces  $\theta$ , but does not affect  $D$ .
- Positive HKA test could be caused by selection for linked gene with the studied gene.

	$\theta$	$D$
Lokus 1	$\theta_1$	$D_1$
Lokus 2	$\theta_2$	$D_2$
Lokus 3	$\theta_3$	$D_3$
Lokus 4	$\theta_4$	$D_4$

Software: <http://genfaculty.rutgers.edu/hey/software#HKA>

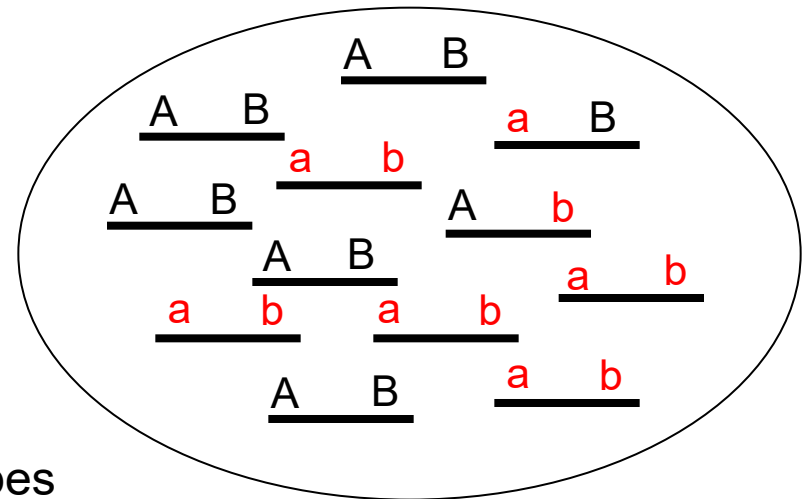
# Linkage disequilibrium

- Certain combinations of alleles in two or more loci occur in higher frequencies than we would expect based on their random combination.



AB.... 25%  
 ab.... 25%  
 aB.... 25%  
 Ab.... 25%

LINKAGE EQUILIBRIUM



AB.... 40%  
 ab.... 40%  
 aB.... 10%  
 Ab.... 10%

LINKAGE DISEQUILIBRIUM

# Linkage disequilibrium (D)

$D = \text{observed} - \text{expected frequency of haplotypes}$

Expected frequencies of haplotypes (i.e. random combinations) are given only by allele frequencies.

haplotype	expected frequency
AB	$p_1q_1$ 8%
ab	$p_2q_2$ 48%
Ab	$p_1q_2$ 32%
aB	$p_2q_1$ 12%

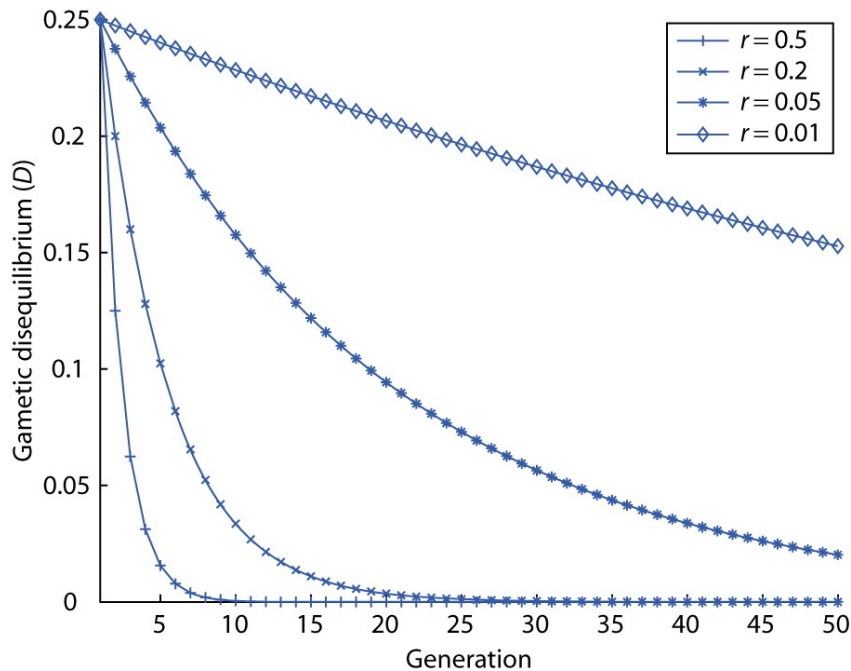
$p_1$  ... frequency of A 40%  
 $p_2$  ... frequency of a 60%  
 $q_1$  ... frequency of B 20%  
 $q_2$  ... frequency of b 80%

$D = 0$  linkage equilibrium  
 $D > 0$  či  $D < 0$  linkage disequilibrium

$D' = D / D_{\max}$  mezi 0 a 1

# Linkage disequilibrium (D)

- The level of linkage disequilibrium depends inversely on recombination rate ( $r$ ) and effective population size ( $N_e$ ).
- Lower levels of linkage disequilibrium in larger populations.
- Linkage disequilibrium between loci tend to decay with time (as recombination events accumulate between loci).



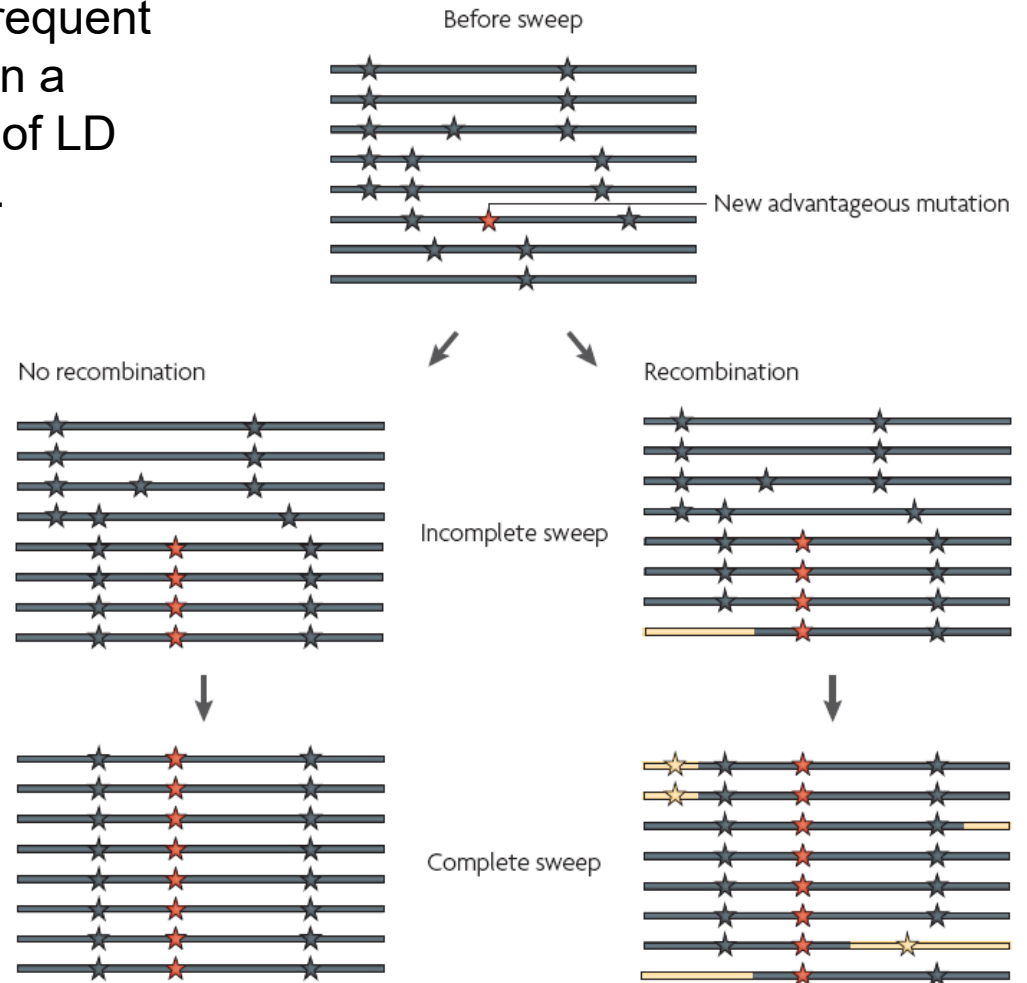
$$D = \frac{1}{4N_e r}$$

Population recombination rate



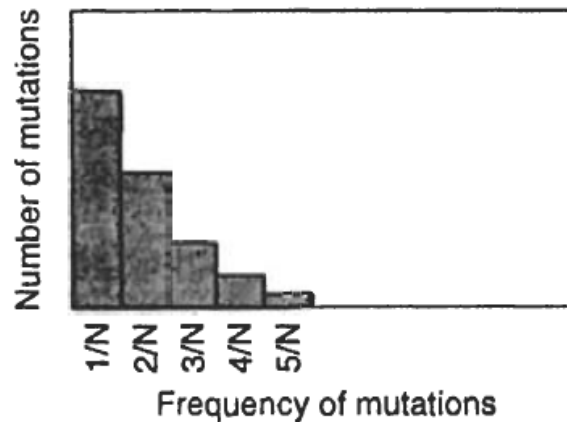
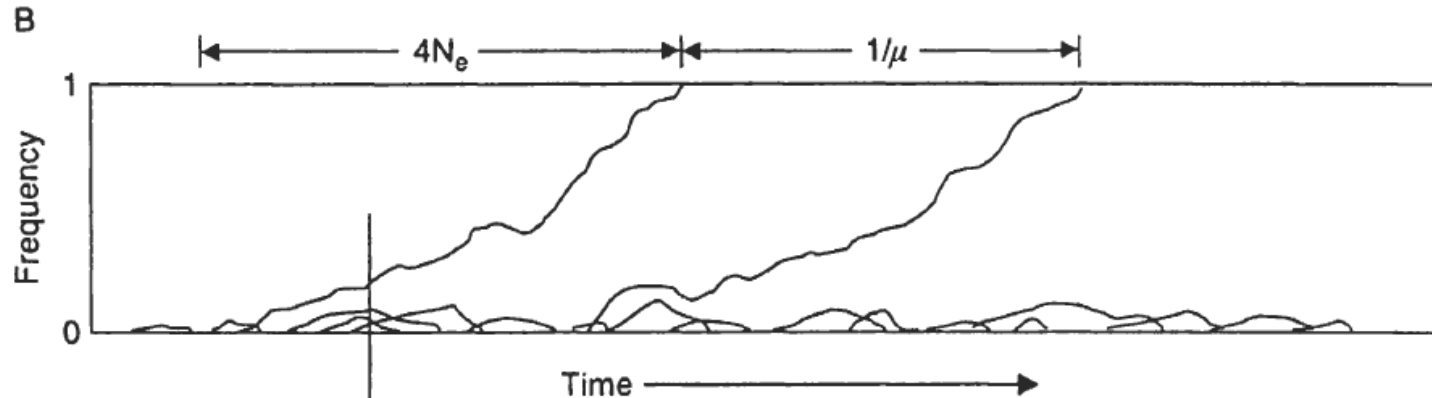
# Detection of recent positive selection based on levels of linkage disequilibrium (LD)

- Ideally detection of relatively frequent haplotype (or haplotype fixed in a subpopulation) with high level of LD compared to other haplotypes.



# Tajima's D test

- Based on distribution of allele frequencies



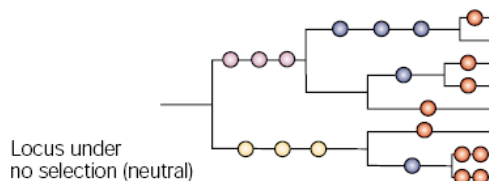
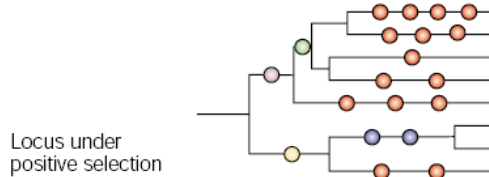
SITE FREQUENCY SPECTRUM

Distribution of allele frequencies in the case of neutrality.

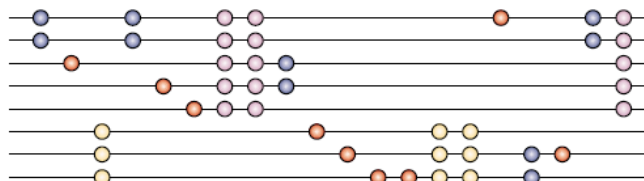
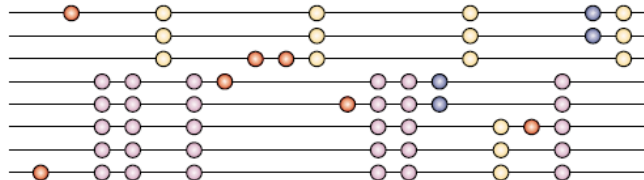
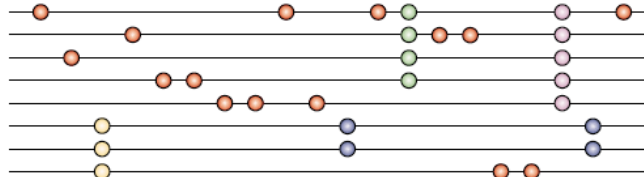
# Tajima's D test

- Positive and negative selection leads to increased proportion of rare alleles.  
Balancing selection leads to the increased proportion of relatively frequent alleles.

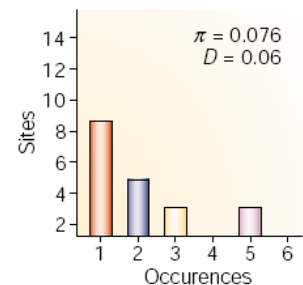
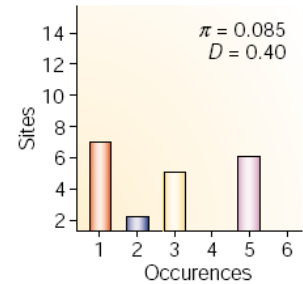
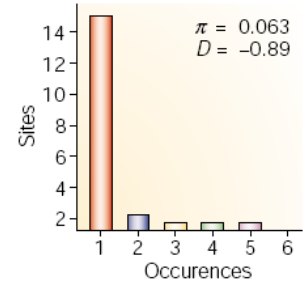
**a** Genealogies



**b** Haplotypes



**c** Site frequency spectra



# Tajima's D test

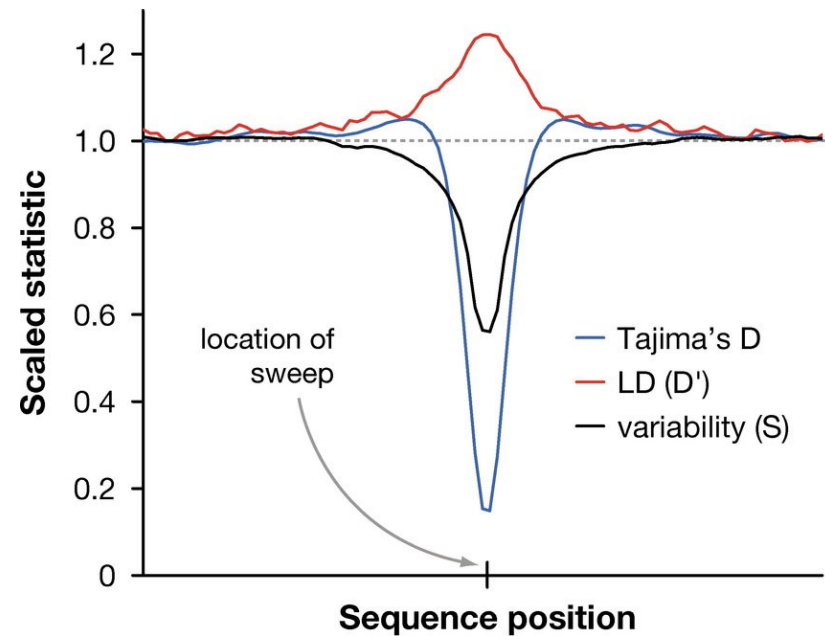
- **Tajima's D =  $(\pi - \theta) / \sqrt{S(\pi - \theta)}$**
- For neutral sequences:  $\theta = \pi$ .  $D = 0$ .
- **Relatively more rare alleles  $\theta > \pi$ .  $D < 0$ .  
(positive or negative selection)**
- **Relatively more frequent alleles  $\pi > \theta$ .  $D > 0$ .  
(balancing selection)**

Population expansion has similar effect on allele frequency spectrum as positive and negative selection.

Bottle-neck has similar effect as balancing selection.

How to differentiate the effect of selection and demographic factors?

# SweepFinder



Nielsen R. 2005.  
Annu. Rev. Genet. 39:197–218



RASMUS NIELSEN

SweepFinder

SweepFinder is a program implementing the method described in Nielsen et al. 2005. Genomic scans for selective sweeps using SNP data. *Genome Research* 1566-1575. It can be used to detect the location of a selective sweep based on SNP data. It will also estimate the frequency spectrum of observed SNP data in the presence of missing data.

## MOLECULAR ECOLOGY

Molecular Ecology (2016) 25, 142–156

doi: 10.1111/mec.13351

DETECTING SELECTION IN NATURAL POPULATIONS: MAKING SENSE OF  
GENOME SCANS AND TOWARDS ALTERNATIVE SOLUTIONS

### Detecting recent selective sweeps while controlling for mutation rate and background selection

CHRISTIAN D. HUBER,\*†‡ MICHAEL DEGIORGIO,§¶ INES HELLMANN\*\* and  
RASMUS NIELSEN††

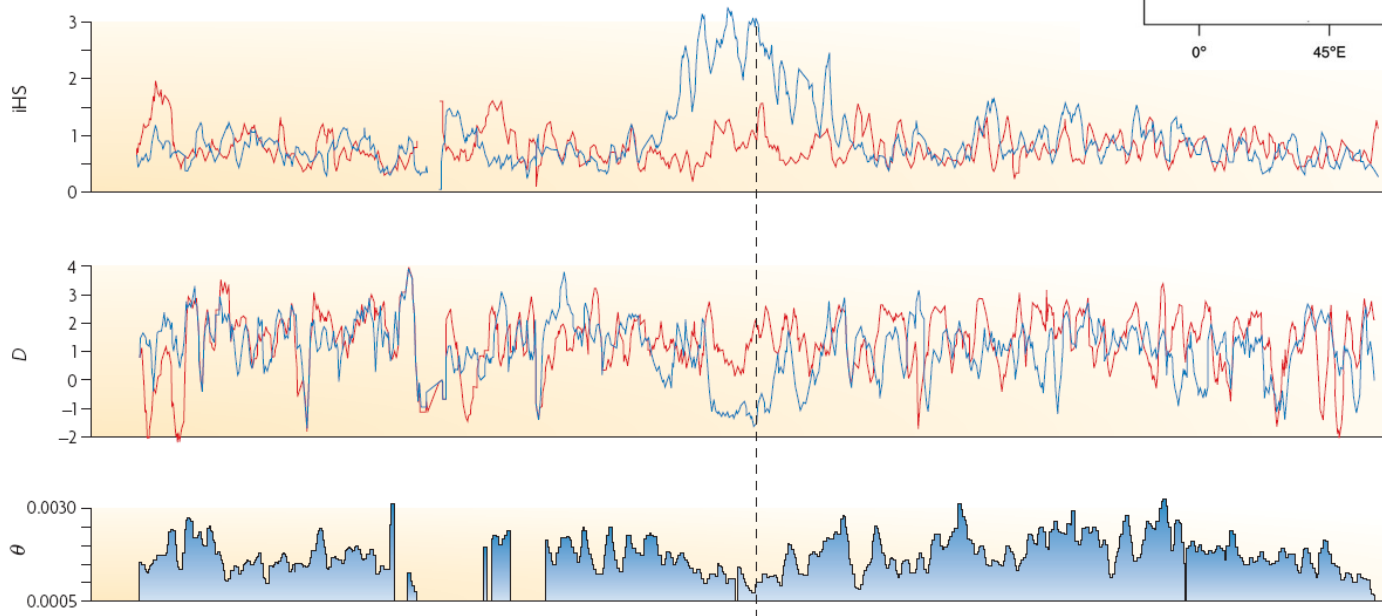
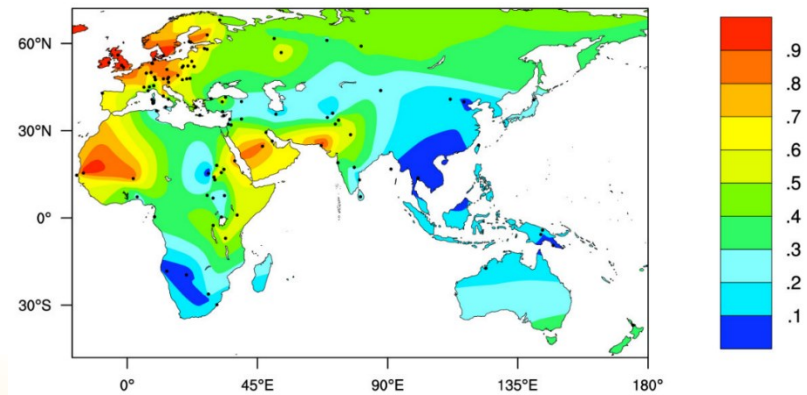
# Selective sweeps in human genome

## Lactase persistence

- Several independent mutations in LCT gene, allows digestion of milk in adulthood.
- Connected with the spread of pastoral farming.
- $s \sim 1,4-15\%$



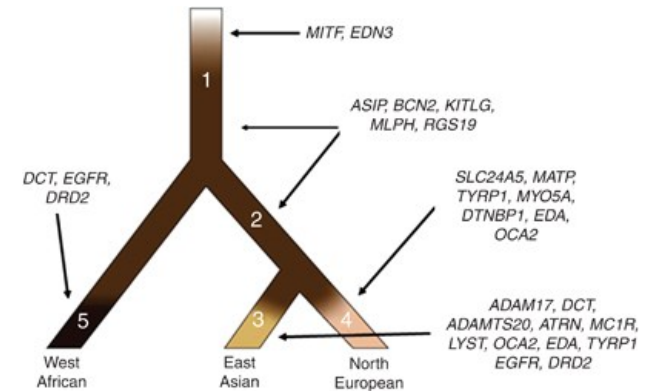
Occurrence of lactase persistence



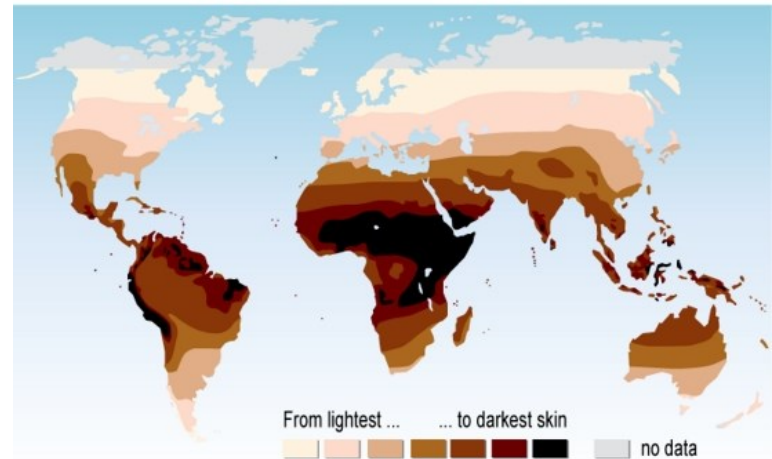
— Europe  
— Asia

# Skin color

- Genes affecting skin color in humans: *HERC2*, *SLC45A2*, *TYR*
- Selection for lighter skin in Europe.
- $s \sim 2-10\%$



Skin colour map for indigenous people  
 Predicted from multiple environmental factors



Source: Chaplin G.®, *Geographic Distribution of Environmental Factors Influencing Human Skin Coloration*, *American Journal of Physical Anthropology* 125:292–302, 2004; map updated in 2007.



# Local adaptations

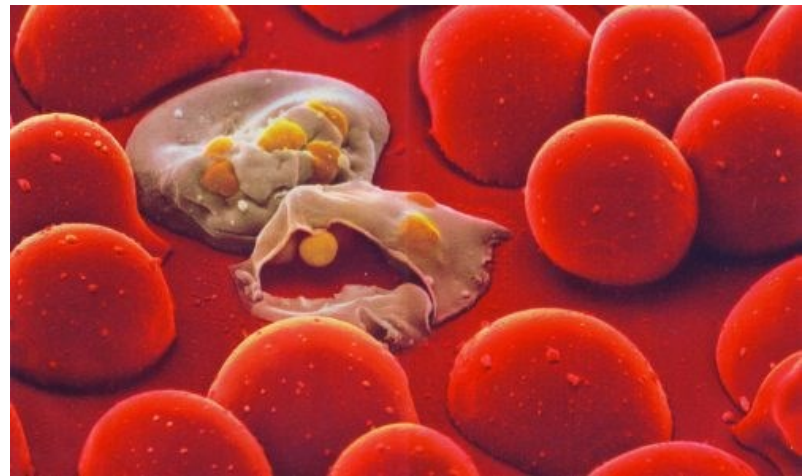
## Inuits in Greenland

- Selective sweeps around genes *FADS1*, *FADS2*. Desaturases, metabolism of fatty acids.



## Resistance against malaria

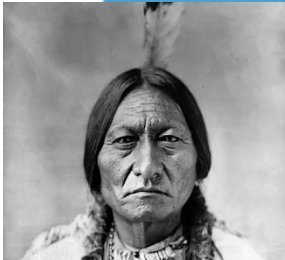
- Mutations in genes *GYP A* and *GYP B* coding for receptors on the red blood cells, *Plasmodium falciparum* uses to get into the cells.





# Detection of genes for local adaptations based on levels of genetic differentiation

- Genes for local adaptations have increased levels of genetic differentiation ( $F_{ST}$ ) between populations.

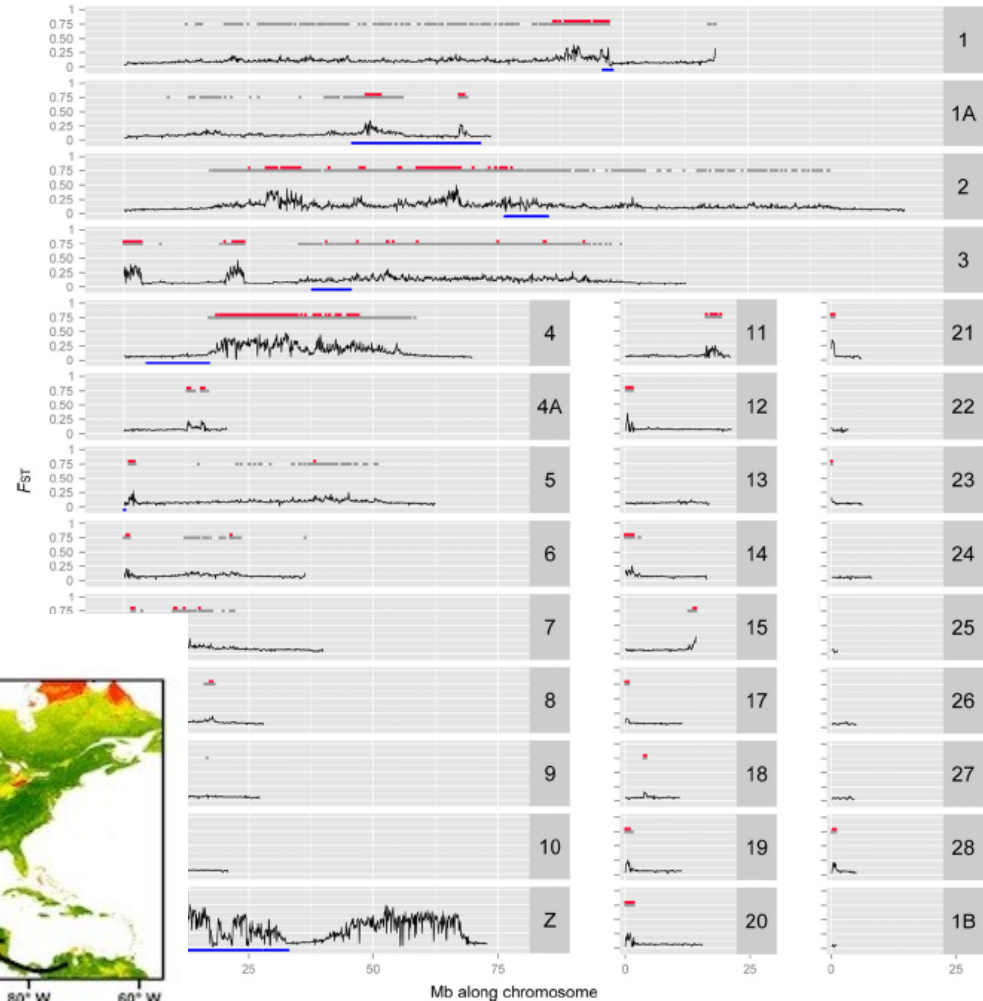
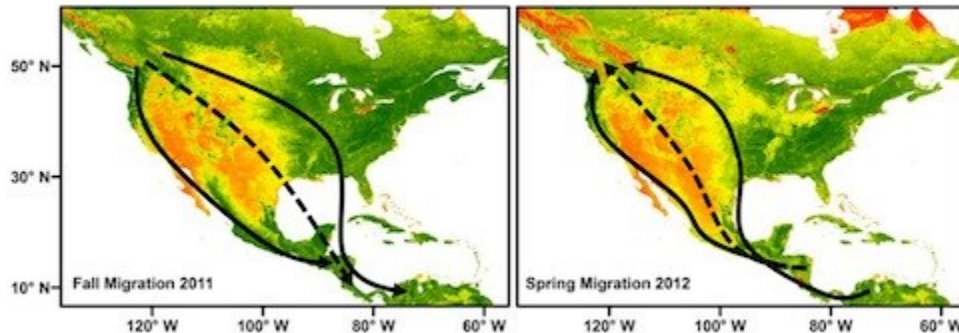


# Genomic islands of increased $F_{ST}$

- Positive selection leads to increased differences in allele frequencies between populations.

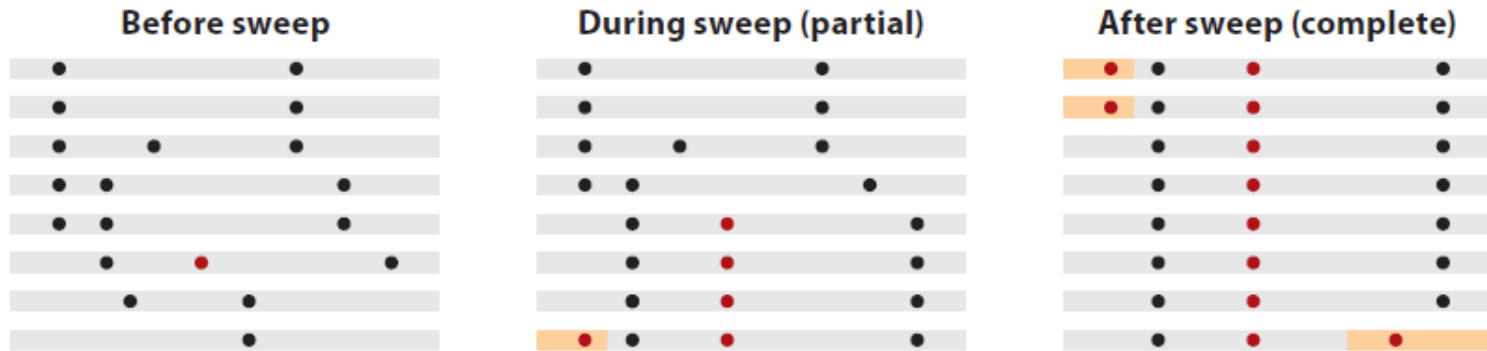
## Islands of differentiation in Swainson's thrush

- Different migration routes

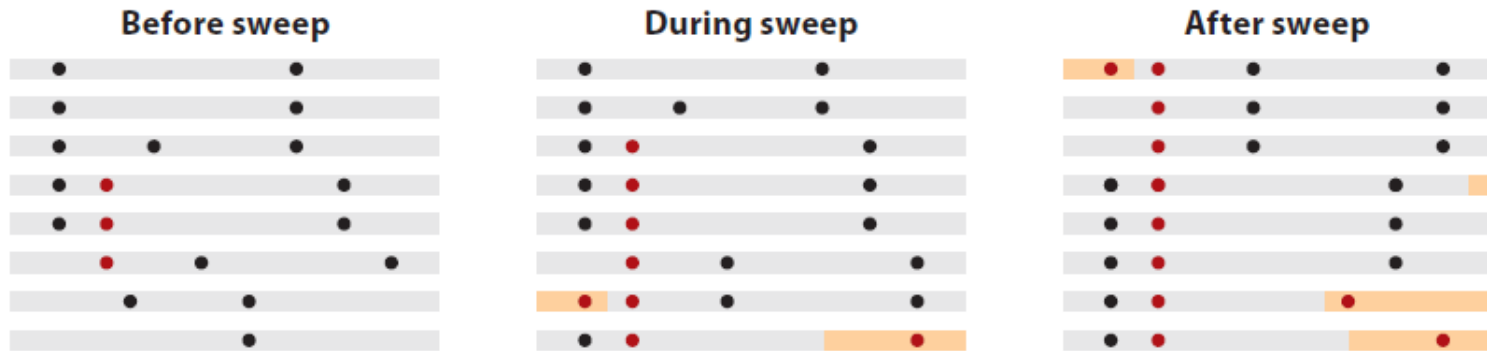


# Adaptation from standing genetic variation (soft selective sweep)

## a Hard sweep



## b Positive selection on standing variation



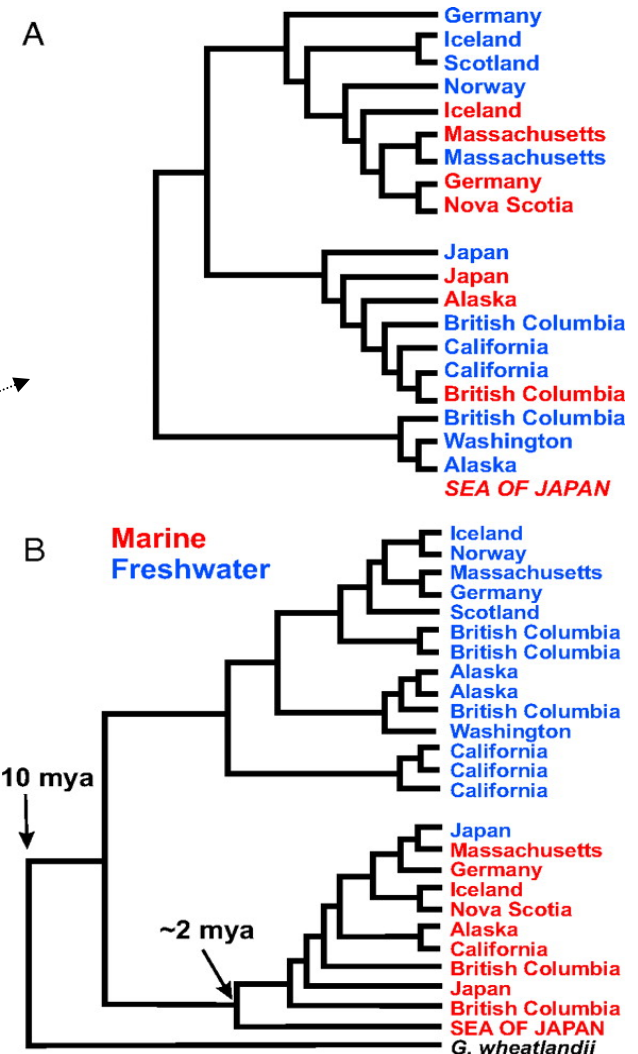
# Adaptation from standing genetic variation and parallel evolution

## Three spine stickleback

marine form (up) a freshwater form (down)

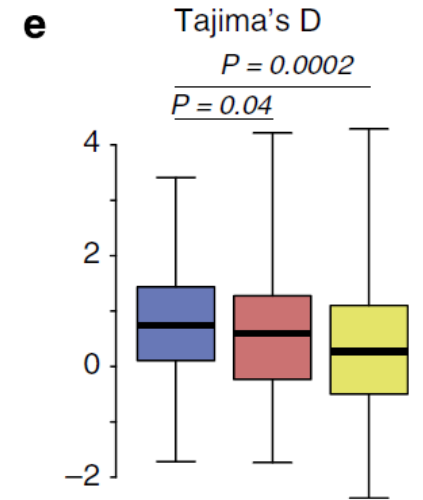
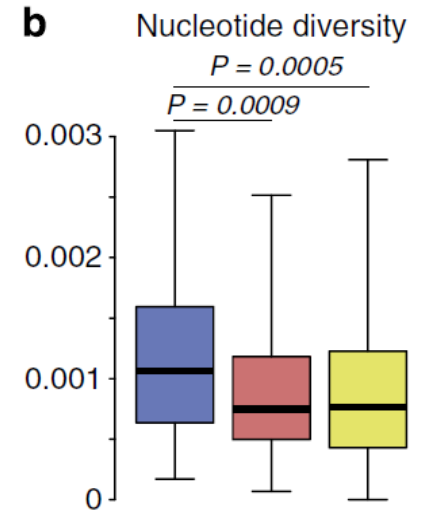
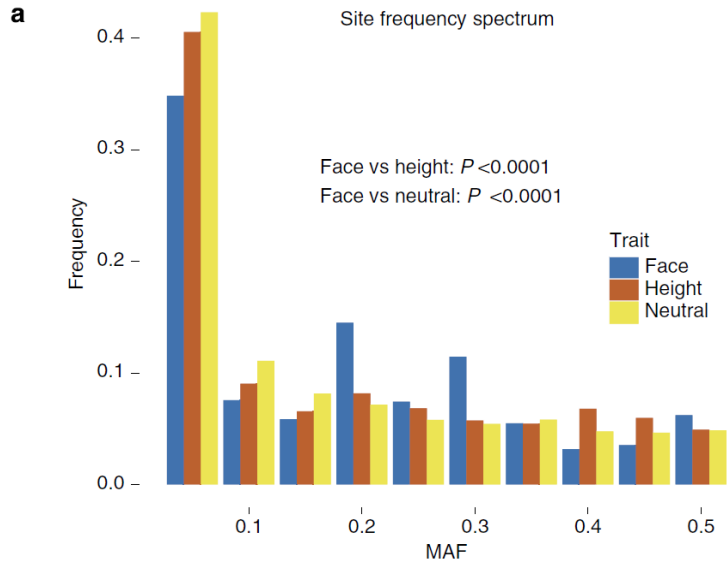


- Freshwater form arose multiple times independently by colonization of new rivers.
- Phylogeny based on *Ectodisplasin (Eda)* gene underlying some differences between forms.
- The same pattern holds for other genes underlying phenotypic differences between forms. These genes often in inversions.



# Balancing selection

## Maintains variability in human faces





# Detection of recurrent positive selection

Repeated fixation of advantageous mutations in the same locus increases the ratio between non-synonymous substitutions ( $K_A$ ) and synonymous substitutions ( $K_S$ ).

$$K_A/K_S = 1$$

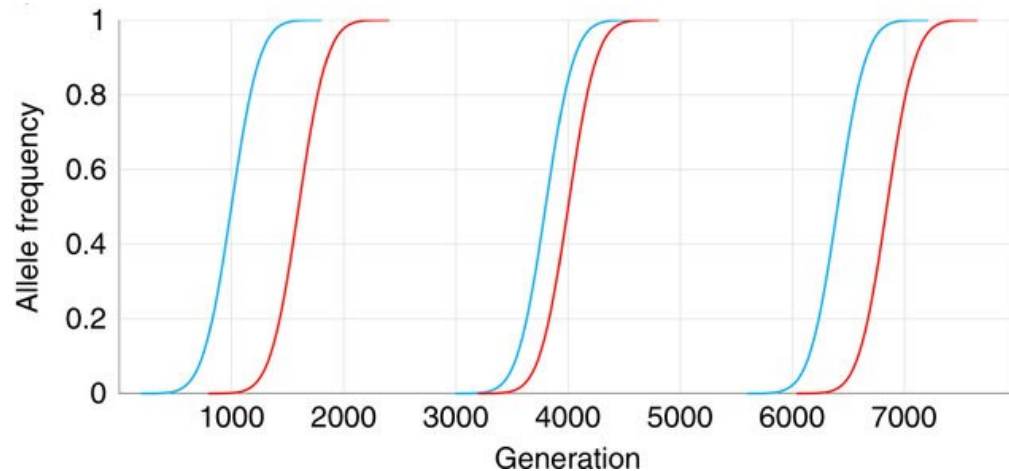
neutral evolution

$$K_A/K_S > 1$$

positive selection



"The Red Queen has to run faster and faster in order to keep still where she is. That is exactly what you all are doing!"



# Detection of long-term negative selection

Longterm negative selection decreases the ratio between non-synonymous substitutions ( $K_A$ ) and synonymous substitutions ( $K_S$ ).

$$K_A/K_S = 1$$

neutral evolution

$$K_A/K_S < 1$$

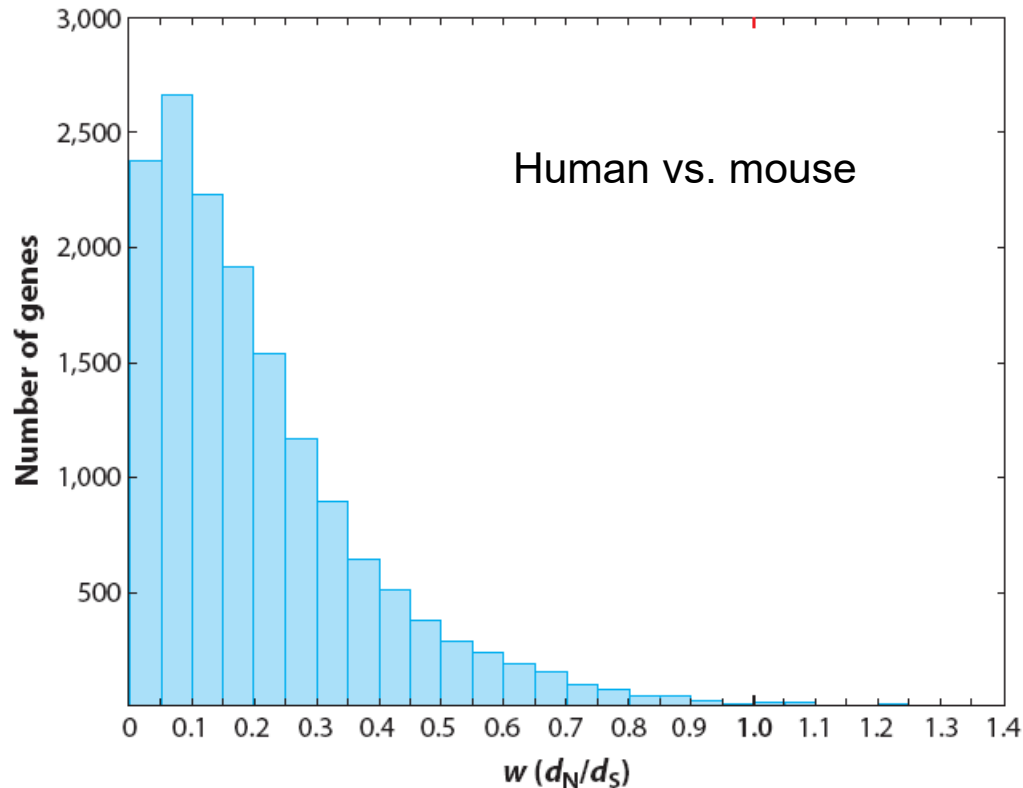
negative selection

**Ultra conserved genetic elements:**  
extremely low substitution rate for  
non-synonymous substitutions.



## $K_A/K_S$ test ( $d_N/d_S$ test)

- Most genes have  $K_A/K_S \ll 1$ .  
Under negative selection.
- Genes with high  $K_A/K_S$  often associated with **reproduction, immunity**, in mammals in **olfactory sense** (e.g. OBP genes).

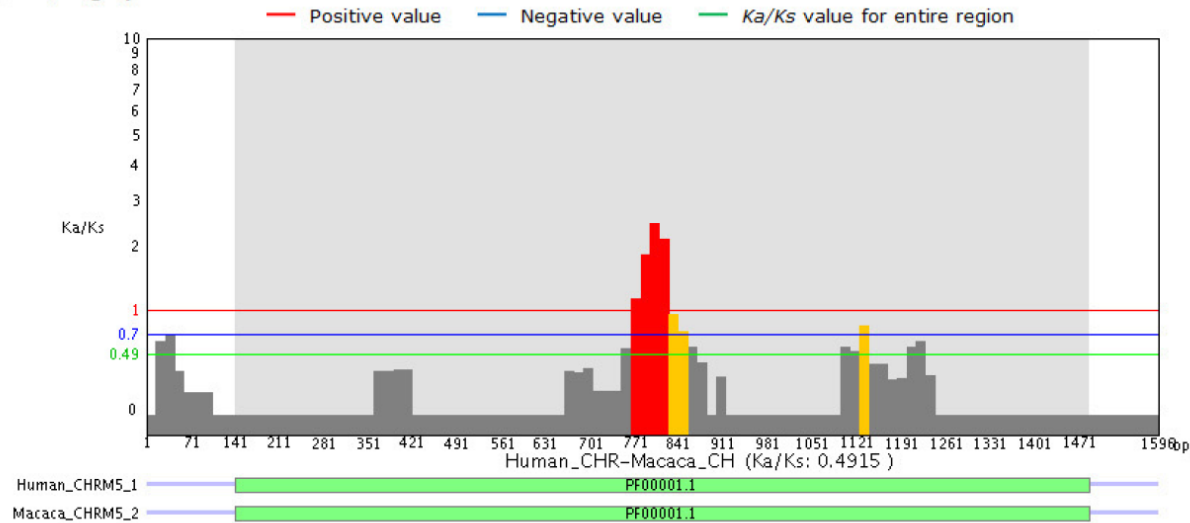




# $K_A/K_S$ test – site specific analysis

- Sliding window analysis – allows to calculate  $K_A/K_S$  for different parts of the genes.
- $K_A/K_S$  can be calculated for individual codons if we compare multiple species.

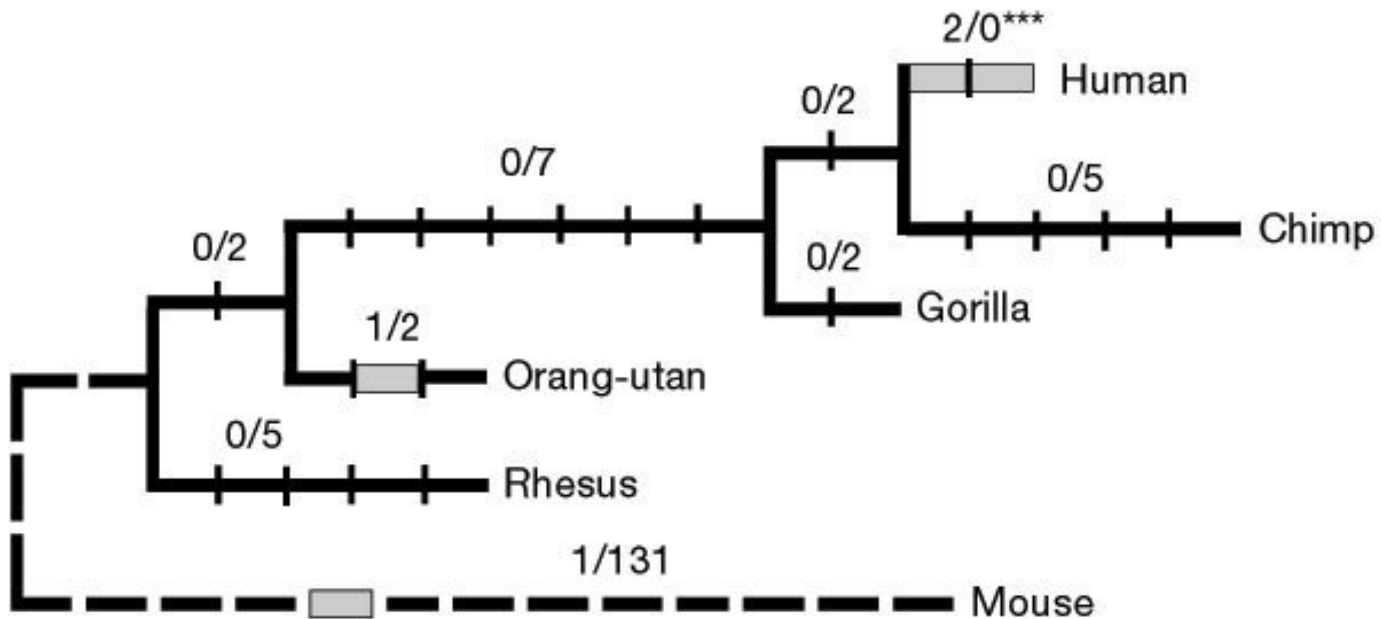
## ● Ka/Ks graph



## $K_A/K_S$ test – lineage specific analysis

FOXP2 gene for speech, positive selection in human lineage

No. of nonsynonymous/ synonymous substitutions.



## McDonald-Kreitman (MK) test

- Compares numbers of nonsynonymous and synonymous substitutions within species (P) and between species (D) for individual genes.
- Neutral genes  $Ds/Ps = Dn/Pn$ .
- Positively selected genes  $Dn/Pn \gg Ds/Ps$ .

	D	P
Synonymous	Ds	Ps
Nonsynonymous	Dn ↑	Pn

Positive selection

MK test allows to estimate proportion of aminoacid substitutions driven by positive selection ( $\alpha$ ).

**Organisms with higher  $N_e$  have higher  $\alpha$ .  
Selection in larger populations is more efficient.**

- *Drosophila* ( $N_e \sim 10^6$ )  $\alpha > 50\%$
- *Mus musculus castaneus* ( $N_e \sim 5 \times 10^5$ )  $\alpha \sim 40-60\%$
- *Populus tremula* ( $N_e \sim 10^5$ )  $\alpha \sim 30-40\%$
- *Mus musculus domesticus* ( $N_e \sim 10^5$ )  $\alpha \sim 13\%$
- Human ( $N_e \sim 10^4$ )  $\alpha \sim 10 - 20\%$

# Genetic basis of adaptations

## Importance of coding vs. regulatory changes

Sean B. Carroll



Mutations in regulatory regions of genes are more important.

Most so far detected mutations underlying adaptations are in coding regions of the gene.



Hopi E. Hoekstra



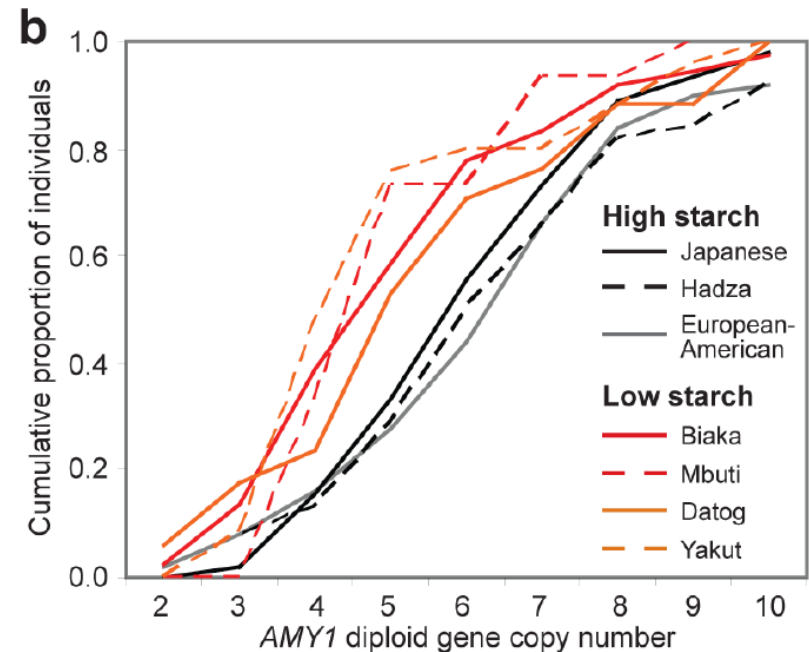
Jerry A. Coyne

Hoekstra HE and Coyne JA (2007). The locus of evolution: evo devo and the genetics of adaptation. *Evolution*.

# Adaptive evolution by gene duplications

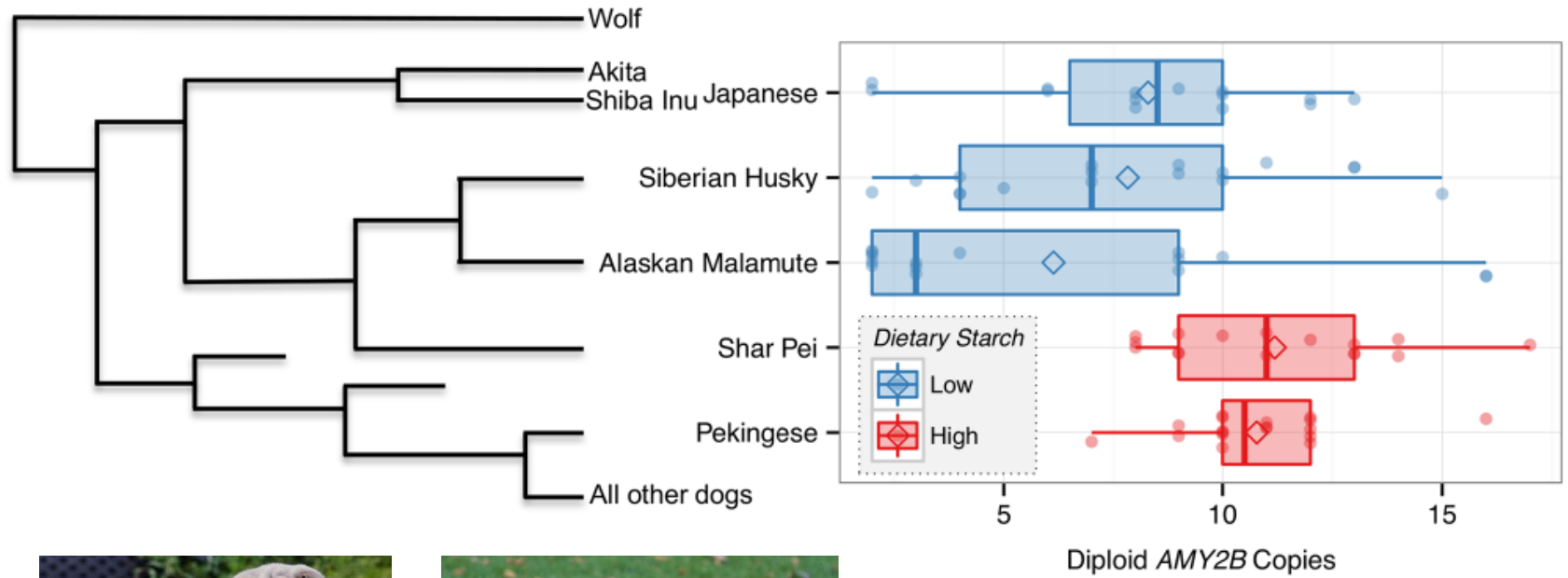
## Salivary amylase gene (*AMY1*)

- Digest starch.
- One of the most variable gene in terms of gene copies in the genome.
- More copies of *AMY1* gene more amylase enzyme in saliva.
- Populations with starch rich diet have higher number of *AMY1* gen copies.



Perry et al. 2007, *Nature Genetics*

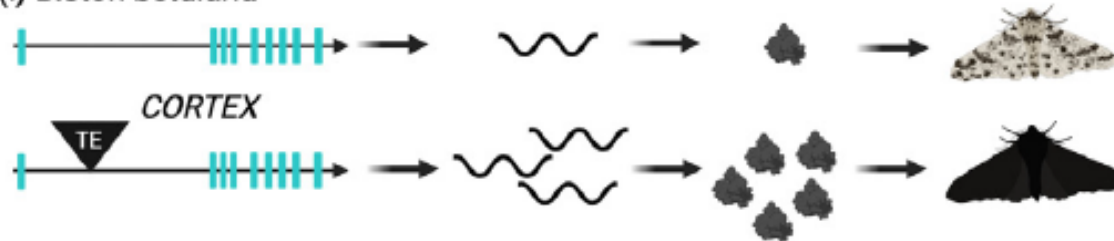
# Number of AMY2B amylase gene copies is related to the diet of domestic dogs



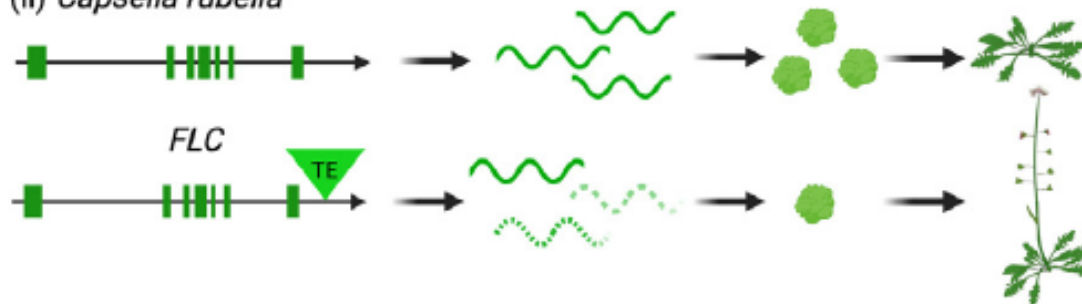
# The role of structural variations in adaptations

## (A) Transposon insertion (TE)

(i) *Biston betularia*

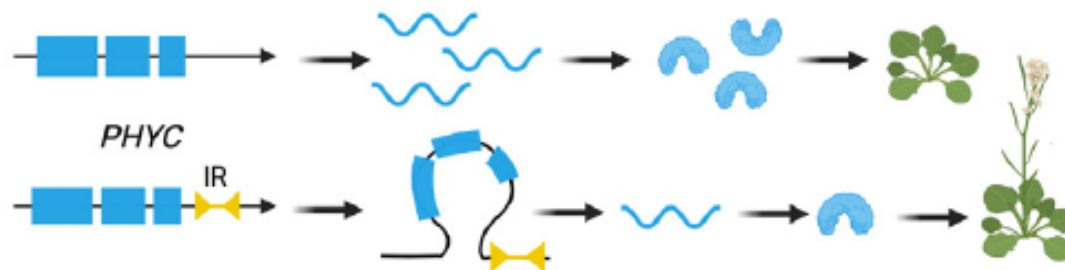


(ii) *Capsella rubella*



## (B) Insertion-deletion mediated chromatin topology change

*Arabidopsis thaliana*

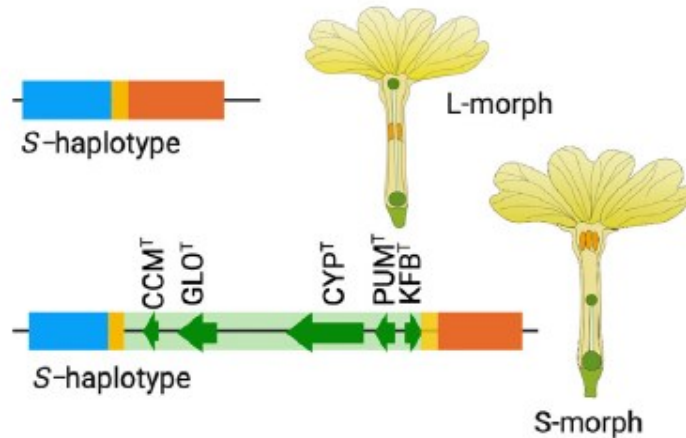




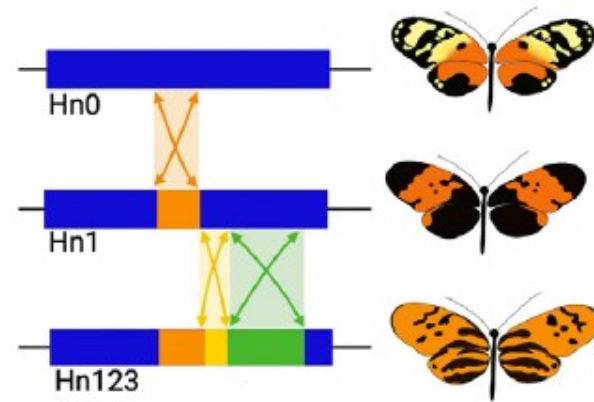
# The role of structural variations in adaptations

## Supergenes

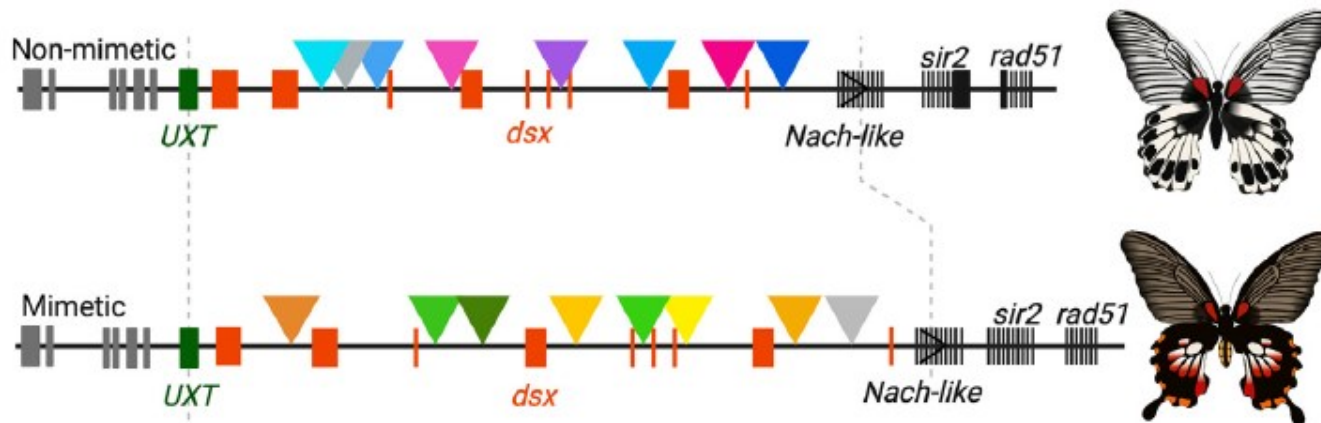
(A) *Primula vulgaris*



(B) *Heliconius numata*



(C) *Papilio memnon*



# Adaptations underlied by epigenetic changes

nature  
ecology & evolution

ARTICLES

<https://doi.org/10.1038/s41559-018-0569-4>

## An epigenetic mechanism for cavefish eye degeneration

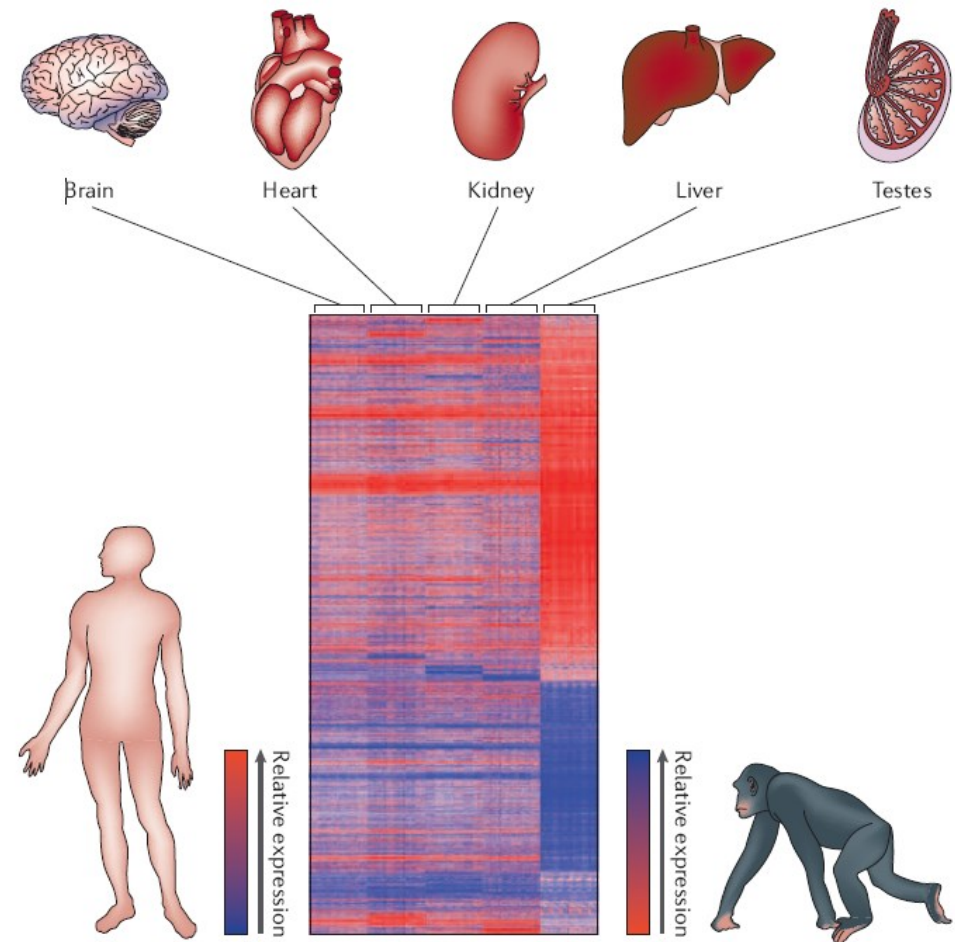
Aniket V. Gore<sup>1\*</sup>, Kelly A. Tomins<sup>1</sup>, James Iben<sup>2</sup>, Li Ma<sup>3</sup>, Daniel Castranova<sup>1</sup>, Andrew E. Davis<sup>1</sup>, Amy Parkhurst<sup>1</sup>, William R. Jeffery<sup>3</sup> and Brant M. Weinstein<sup>1\*</sup>

Coding and non-coding mutations in DNA contribute significantly to phenotypic variability during evolution. However, less is known about the role of epigenetics in this process. Although previous studies have identified eye development genes associated with the loss-of-eyes phenotype in the Pachón blind cave morph of the Mexican tetra *Astyanax mexicanus*, no inactivating mutations have been found in any of these genes. Here, we show that excess DNA methylation-based epigenetic silencing promotes eye degeneration in blind cave *A. mexicanus*. By performing parallel analyses in *A. mexicanus* cave and surface morphs, and in the zebrafish *Danio rerio*, we have discovered that DNA methylation mediates eye-specific gene repression and globally regulates early eye development. The most significantly hypermethylated and downregulated genes in the cave morph are also linked to human eye disorders, suggesting that the function of these genes is conserved across vertebrates. Our results show that changes in DNA methylation-based gene repression can serve as an important molecular mechanism generating phenotypic diversity during development and evolution.

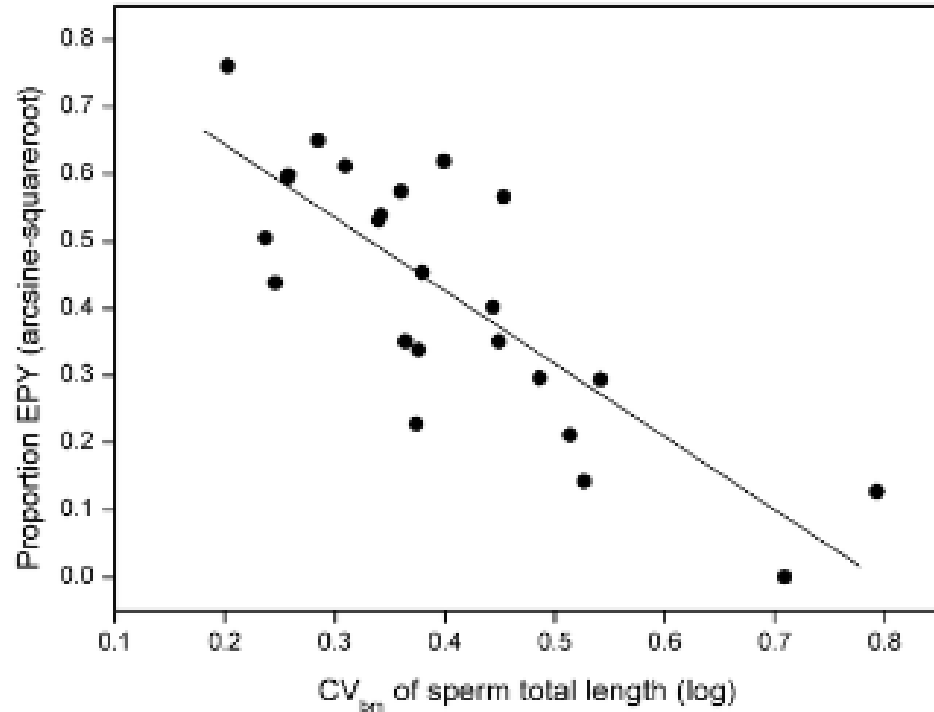
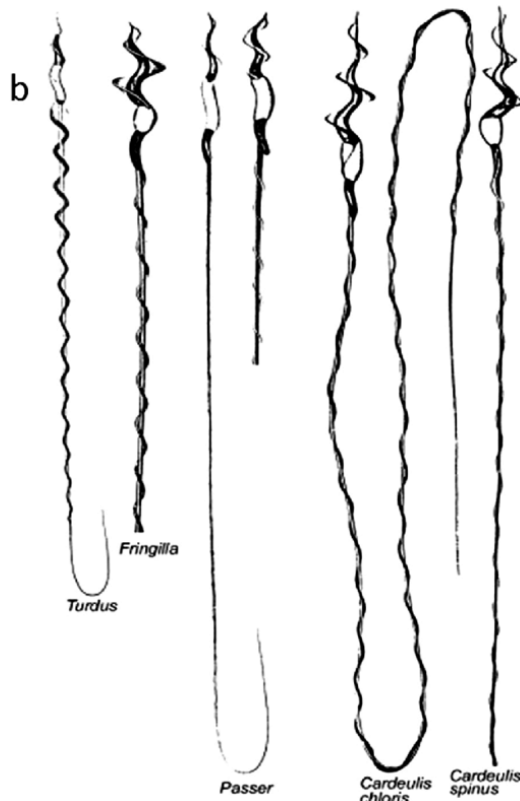


# How to detect positive selection in gene expression data?

- Relatively higher divergence in expression between species compared to variability in expression within species.
- Expression in testes show cca 3-times higher divergence compared to variability within species than somatic tissues.



# Sperm Length Variation as a Predictor of Extrapair Paternity in Passerine Birds

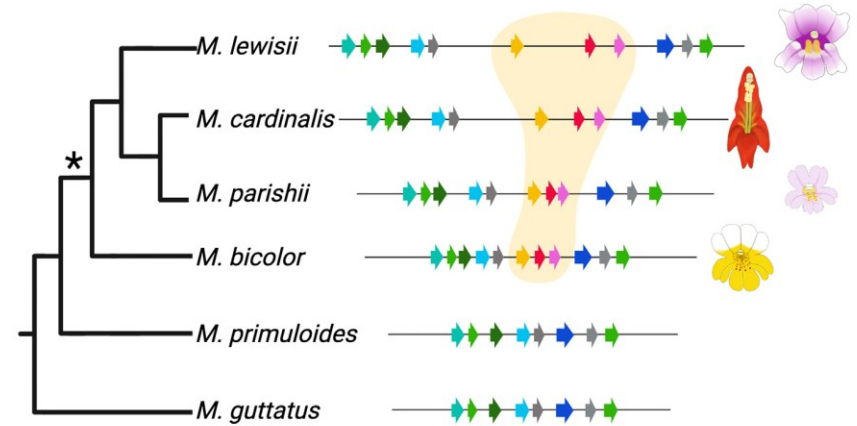




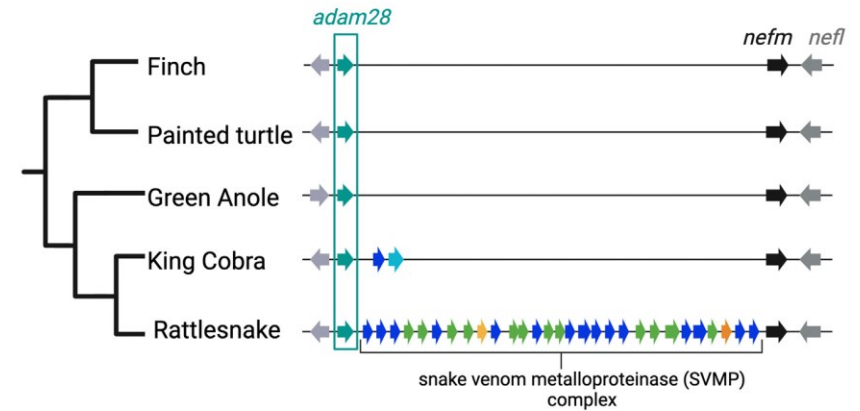
# Taxon-specific genes underlying adaptive traits



(A) YUP-SOLAR-PELAN in *Mimulus*



(B) Snake venom metalloproteinase (SVMP) genes in *Crotalus*



(C) Antifreeze glycoprotein (AFGP) genes in Antarctic notothenioid fish

