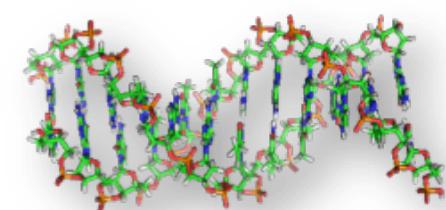
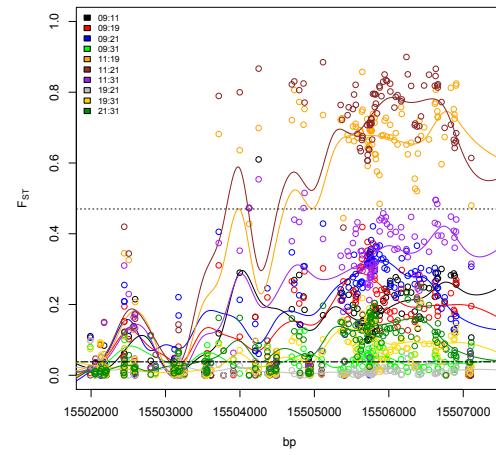


# Detection of the genomic signature of selection

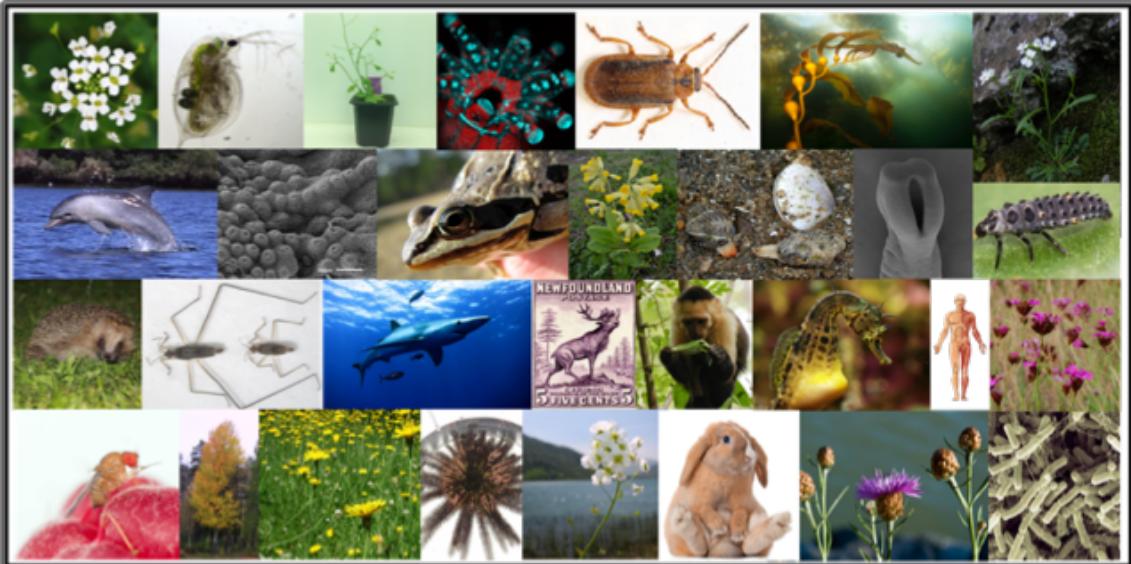
Martin C. Fischer  
Institute of Integrative Biology  
ETH Zürich

[martin.fischer@env.ethz.ch](mailto:martin.fischer@env.ethz.ch)

September 17<sup>th</sup>, 2018



# Diversity and adaptation



Why is there so much variation among and within species?



How do organisms adapt to their environment?



# Genetic diversity

- Evolutionary changes at the molecular level are caused by...

## Neutral processes

- Mutations



- Genetic drift



- Population history

## Mutations

- ❖ Small scale:

- Point mutations (SNPs)
- Deletions
- Insertions

- ❖ Large scale:

- Copy number variation
- Translocations
- Inversions

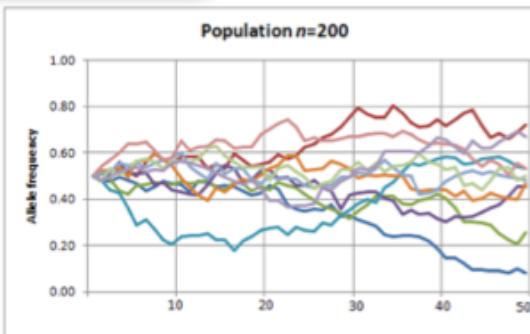
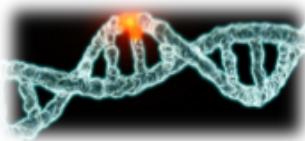


# Genetic diversity

- Evolutionary changes at the molecular level are caused by...

## Neutral processes

- Mutations



- Genetic drift

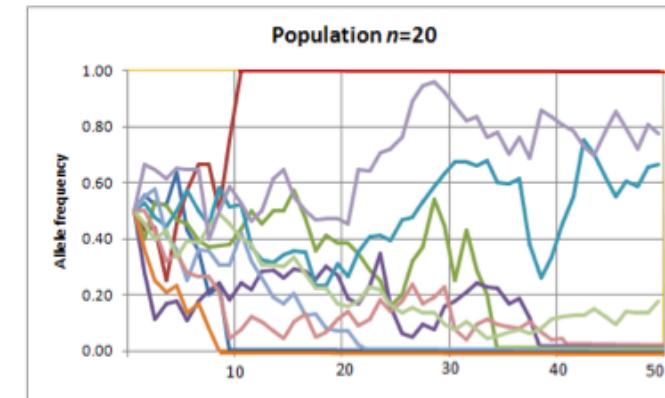


- Population history

## Genetic drift

- ❖ Change of allele frequencies over generations in a population due to random sampling

- ❖ Population size:  
Drift is largest in small populations

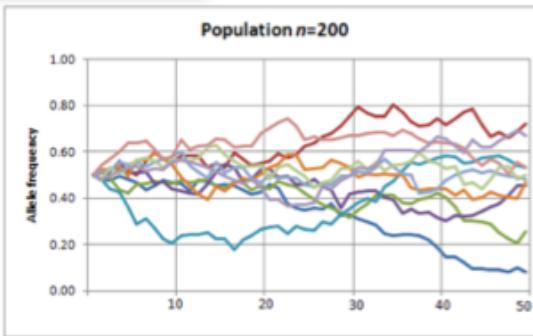
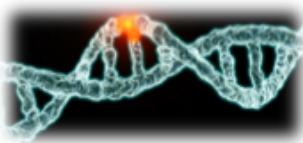


# Genetic diversity

- Evolutionary changes at the molecular level are caused by...

## Neutral processes

- Mutations



- Genetic drift



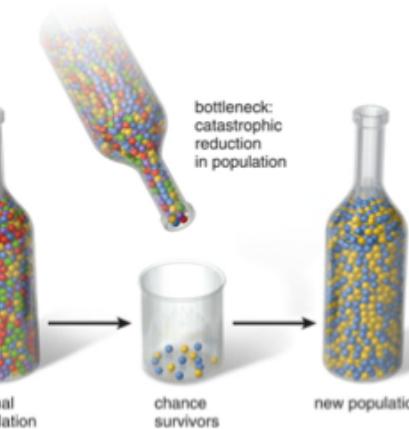
- Population history

## Population history

- ❖ Demography



- ❖ Bottleneck



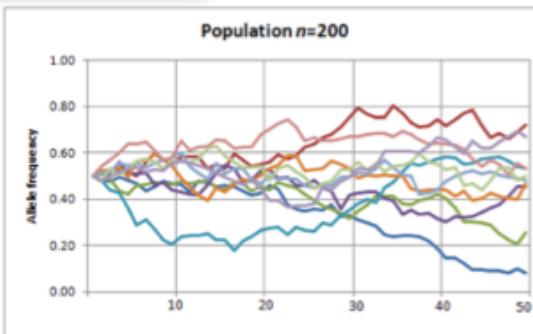
- ❖ Migration / gene-flow

# Genetic diversity

- Evolutionary changes at the molecular level are caused by...

## Neutral processes

- Mutations



- Genetic drift

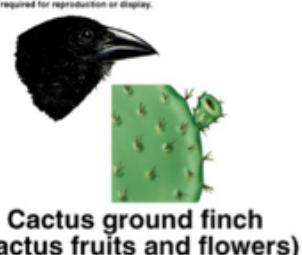


- Population history

## Selection (natural, ...)



Large ground finch (seeds)



Cactus ground finch (cactus fruits and flowers)

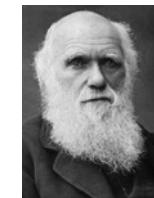


Vegetarian finch (buds)



Woodpecker finch (insects)

Darwin's Finches

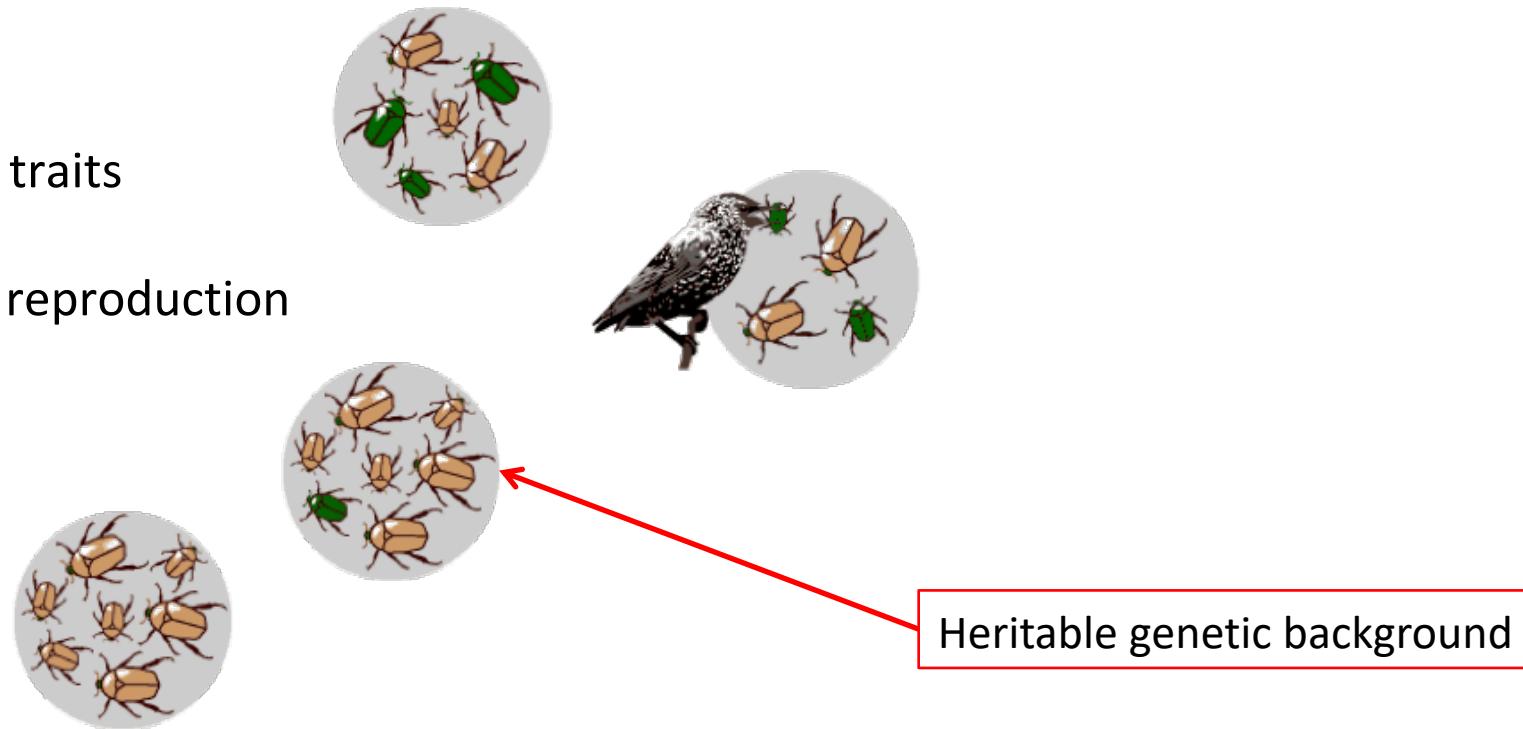


# Natural selection

Natural selection is one of the basic mechanisms of evolution

## Population of beetles:

1. There is variation in traits
2. There is differential reproduction
3. There is heredity
4. End result



# The genomic signature of selection

## Where do we find the genomic signature of selection ?

- ❖ *Arabidopsis thaliana*: 0.15 Gb
- ❖ *Homo sapiens*: 3.2 Gb



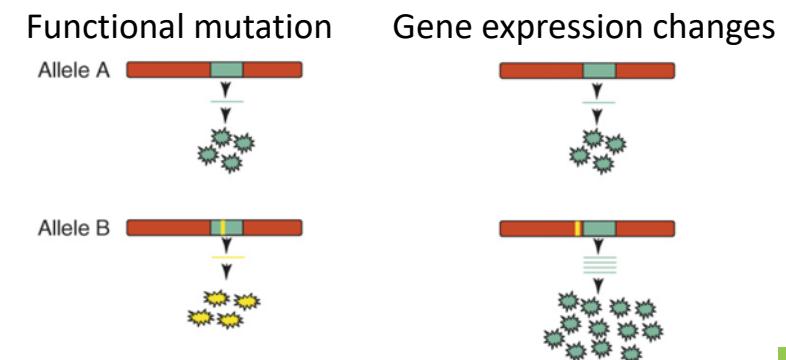
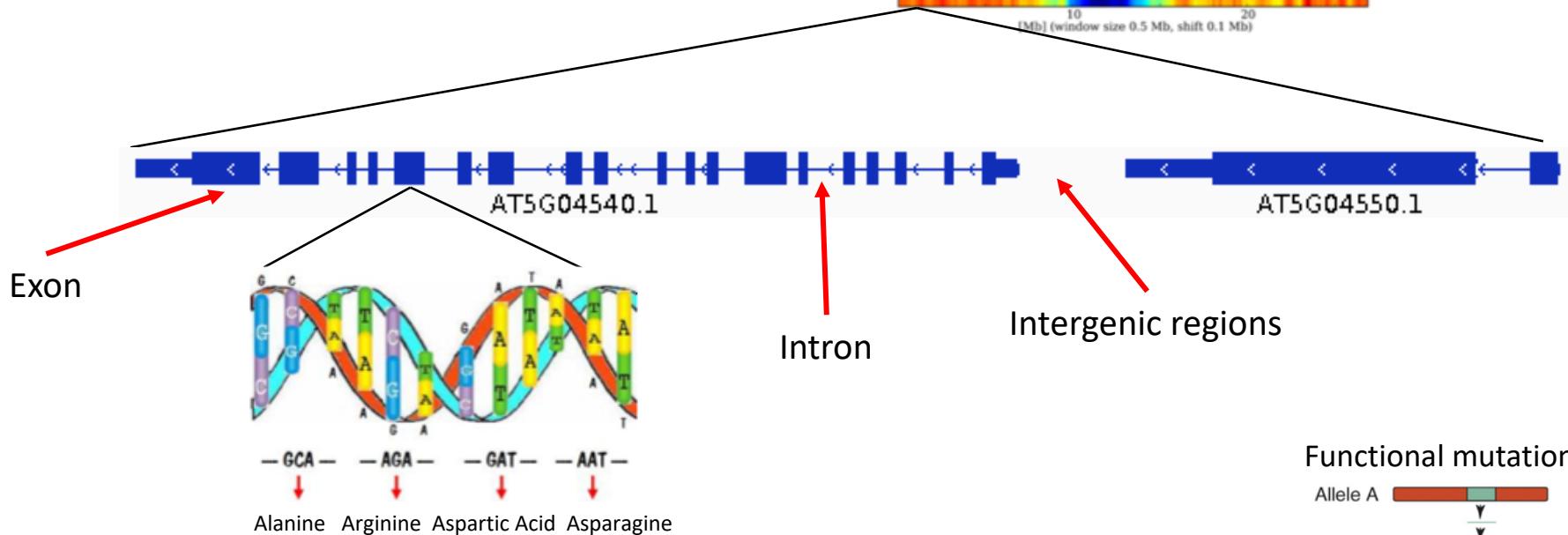
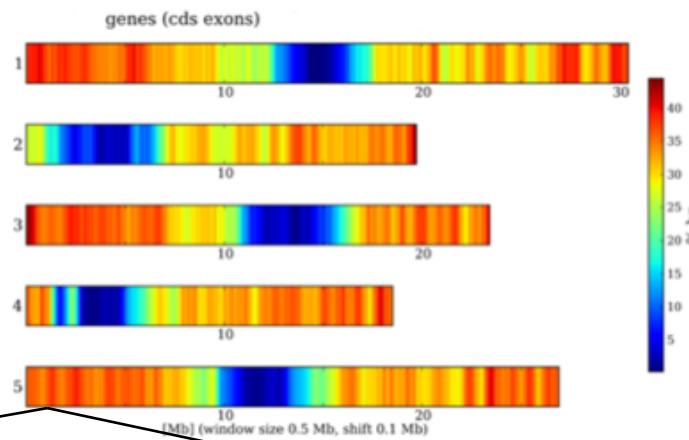
- ❖ A4 => Font 10



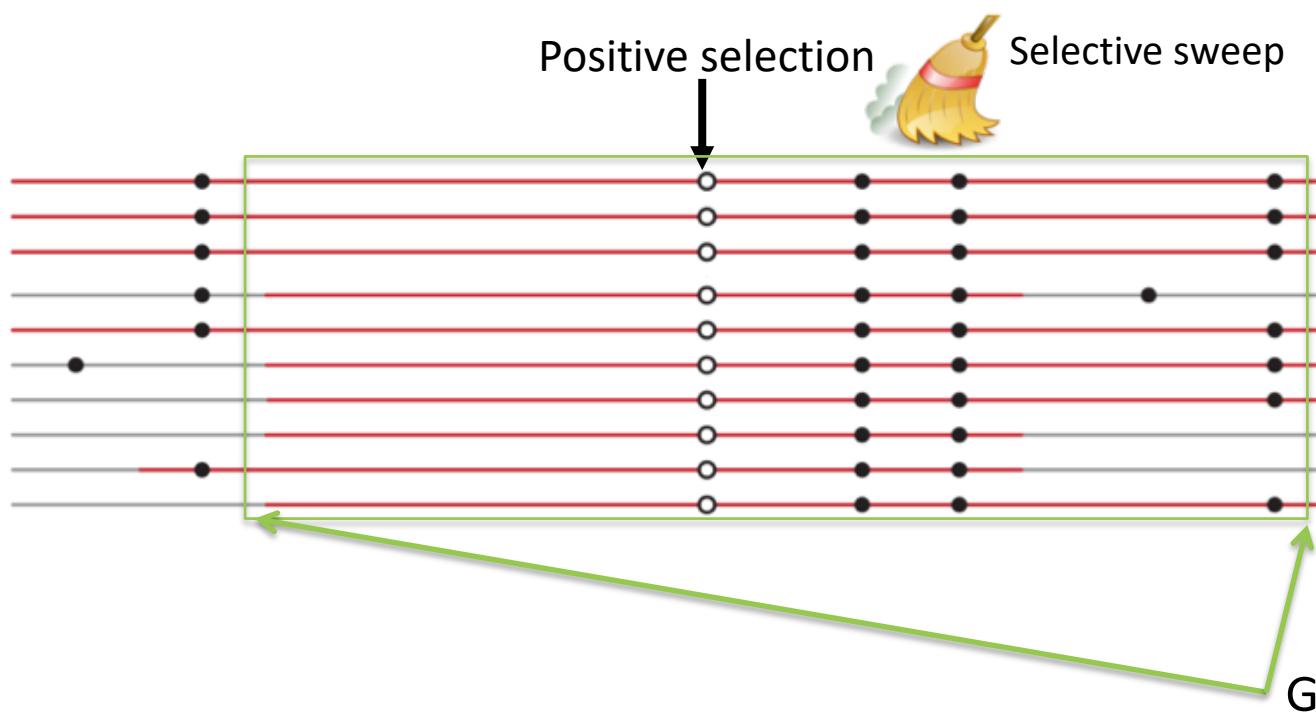
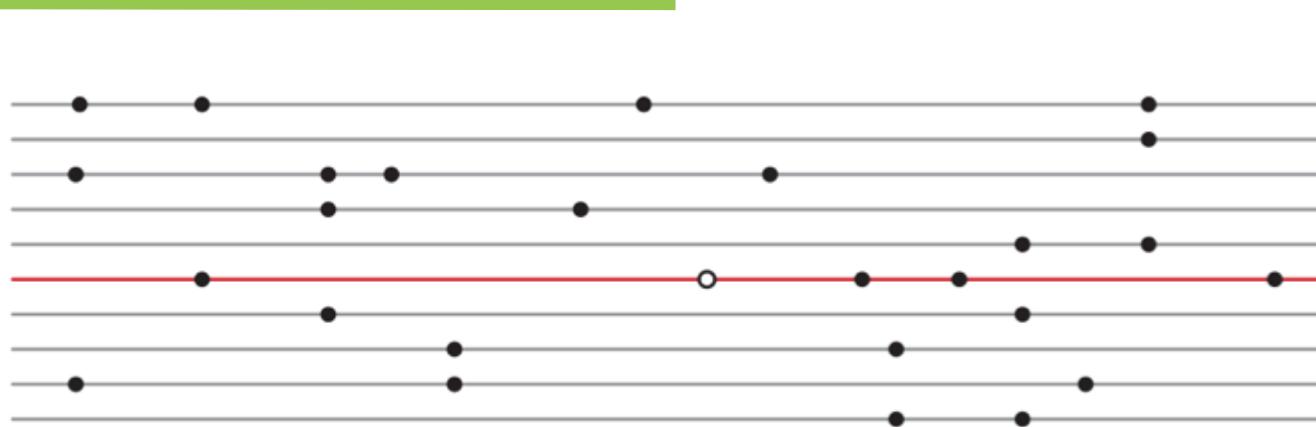
- ❖ *A. thaliana* ~26,000 pages
  - ❖ *Homo sapiens* ~532,000 pages
- ~26 Zurich phone books  
~532 Zurich phone books

# The genome and where selection acts

	Human	<i>A. thaliana</i>
Genome	3.165 Gb	0.157 Gb
Chromosomes	22, XY	5
Chr length	50-250 Mb	20-34 Mb
Genes	~25,000	~28,000



# Selective sweeps



## Selective sweep:

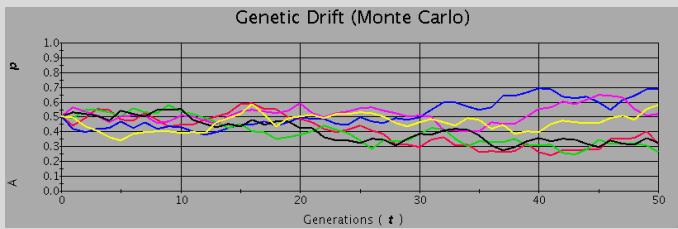
- An allele that increases fitness arises and 'sweeps' to fixation in a population

- Recombination
- Linkage disequilibrium
- Long haplotypes

# Locus under selection behaves differently

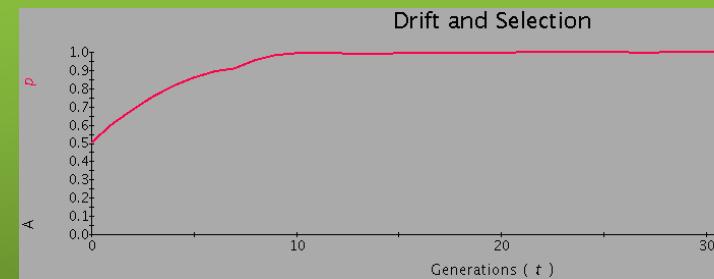
## Neutral processes:

- ❖ Genetic drift
- ❖ Demographic history
- Affect all loci similarly



## Selection:

- Affect only single locus



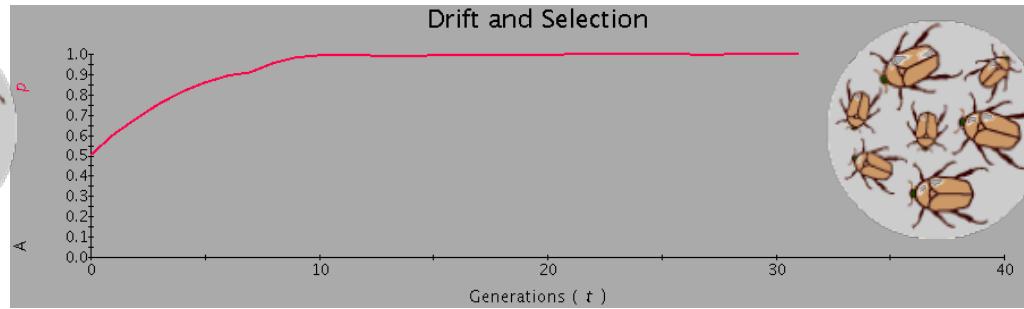
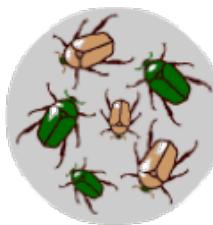
- Outlier detection



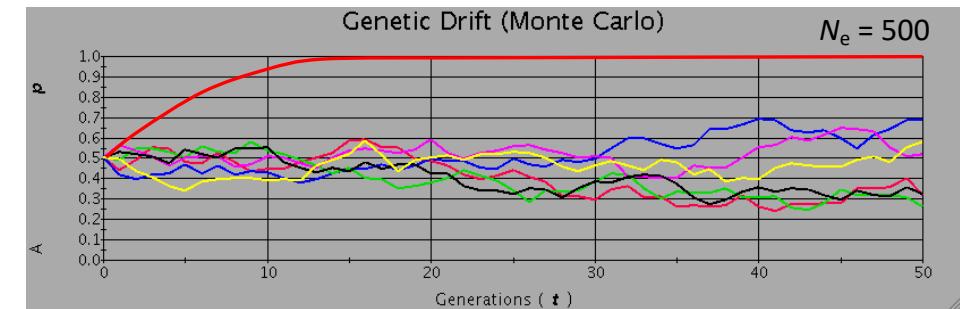
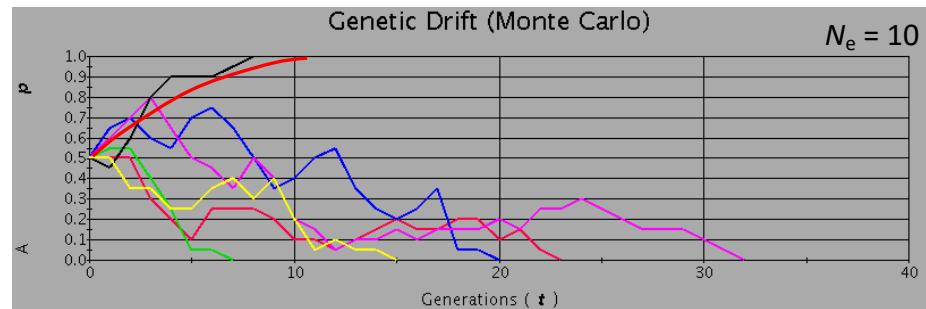
Outlier compared to the rest of the genome

# Population size matters....

## ❖ Selection



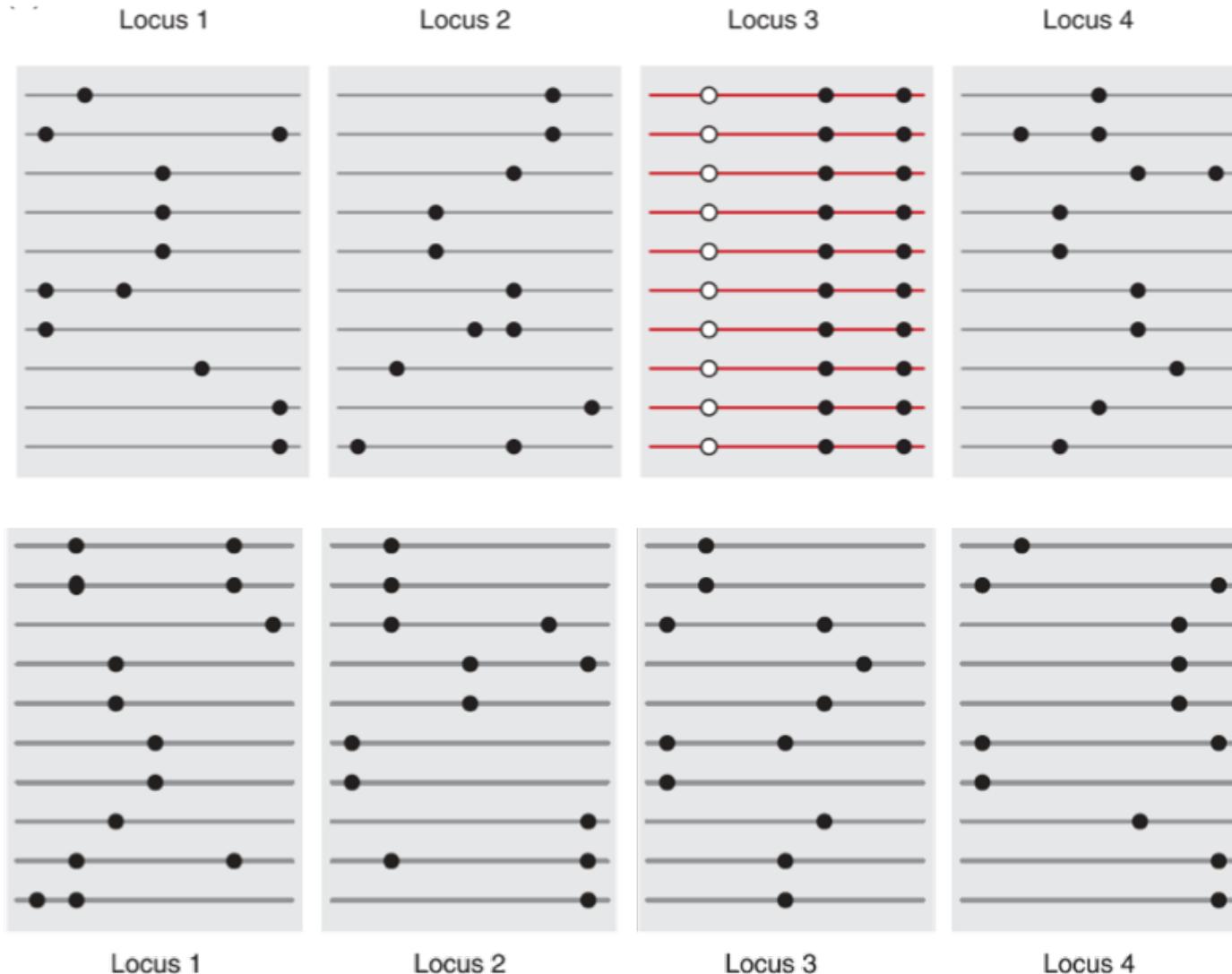
## ❖ Drift



❖ Selection coefficient  $s \gg (1/2N_e)$

- Population under selection needs to have a minimum  $N_e$  to overcome drift, or  $s$  needs to be very strong
- Selection acts most efficient on large populations

# Selective sweeps



Single locus is affected

➤ Recombination

Skewed allele frequency

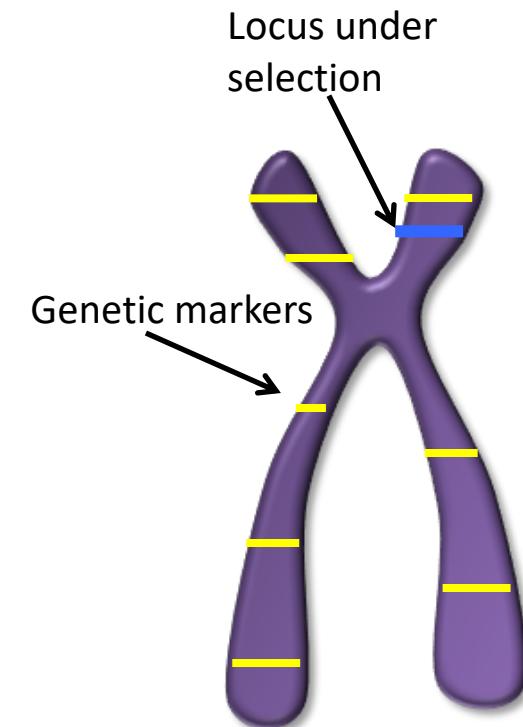
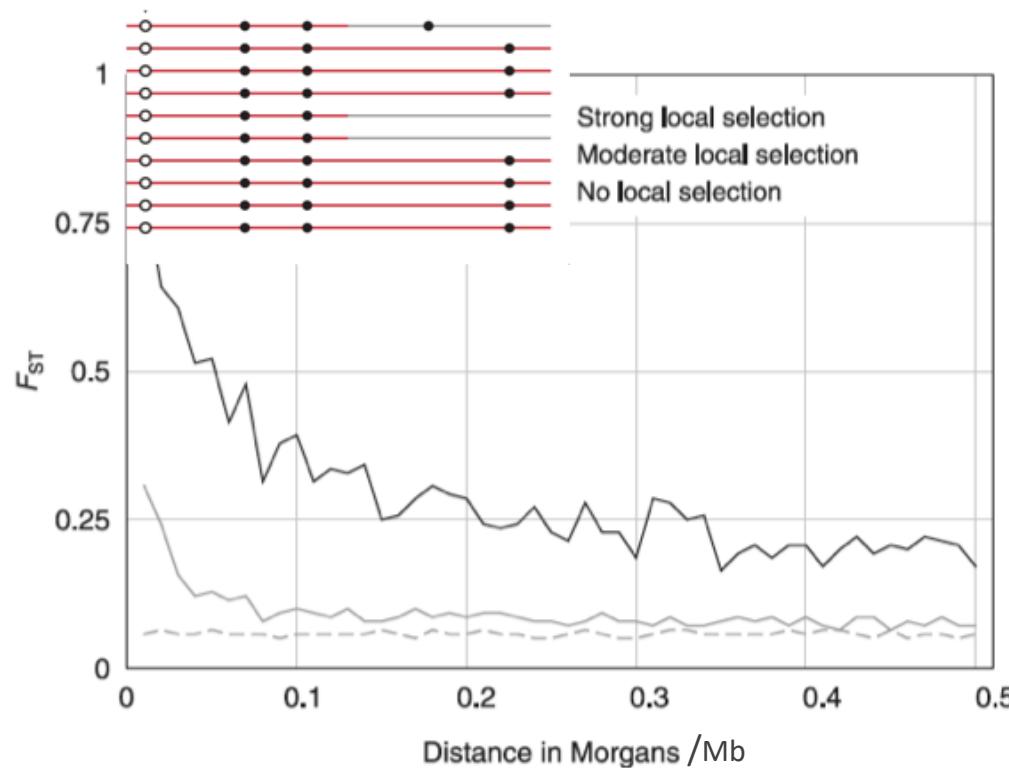
Reduced diversity ( $\pi, S$ )

Strong locus specific differentiation

➤  $F_{ST}$ -outlier

➤  $F_{ST} = \pi_{\text{total}} - \pi_{\text{subpopulations}} / \pi_{\text{total}}$

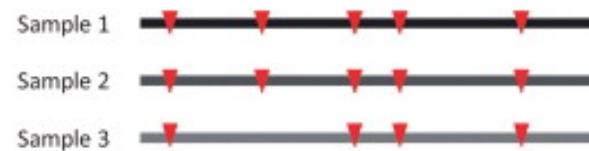
# Genetic hitchhiking



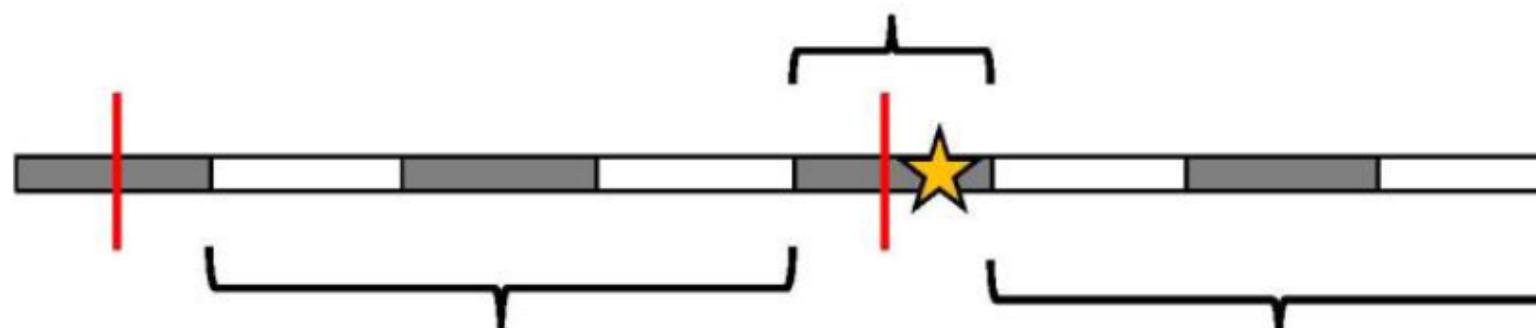
- ❖ The signature of selection decays with increasing distance from the locus under selection  
=> **genetic hitchhiking**
- ❖ Recombination

# Marker density

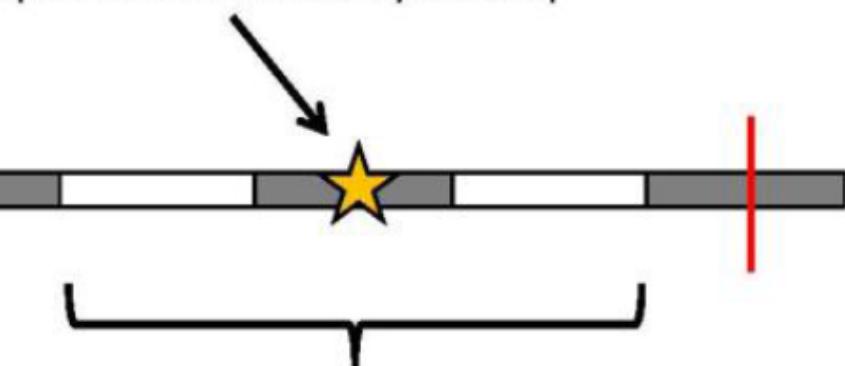
❖ E.g. reduced representation sequencing (e.g. ddRAD-seq)



RAD-tag in LD with adaptation SNP



Adaptation SNP missed by RADseq

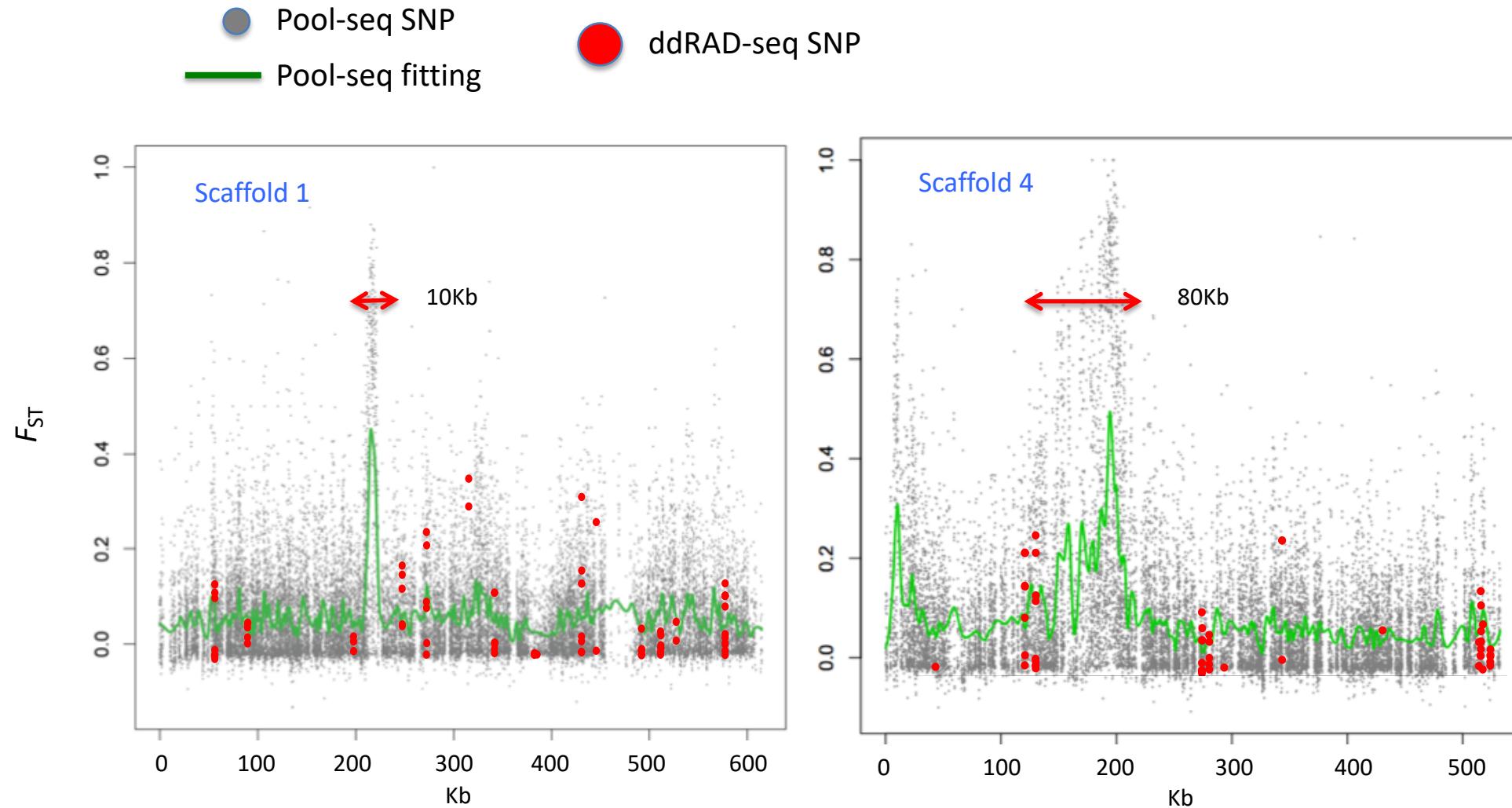


Genomic region missed by RADseq

Genomic region missed by RADseq

Genomic region missed by RADseq

# Marker density is important to detect outliers!



RAD-seq data can't detect most of the 'islands of selection'



# Four flavors of selection

---

## Positive (directional)

- 'New' (non-synonymous) mutations selected for
- Evolution of novel protein function

## Diversifying

- Geographically restricted selection (e.g. due to spatial variation in climate)

## Balancing

- Maintenance of multiple alleles within-population
- E.g. heterozygote advantage (sickle cell anemia), frequency dependent selection

## Negative (purifying)

- New (non-synonymous) mutations selected against
- Retention of existing protein function

# Case studies: Detect the genomic signature of selection

## Different methods to detect selection:

- ✧ Reduced level of genetic variation (e.g.  $\pi$ )
- ✧ Linkage disequilibrium ( $LD$ )
- ✧ Skew of allele frequency spectra (Tajima's  $D$ )
- ✧ Locus specific population differentiation
  - $F_{ST}$ -outlier approach
  - Population Branch Statistics (PBS)
- ✧ Environmental association analysis (EAA)
- ✧ Low coverage sequencing

## Comparison:

- ✧ Within a population
- ✧ Among populations
- ✧ Among species

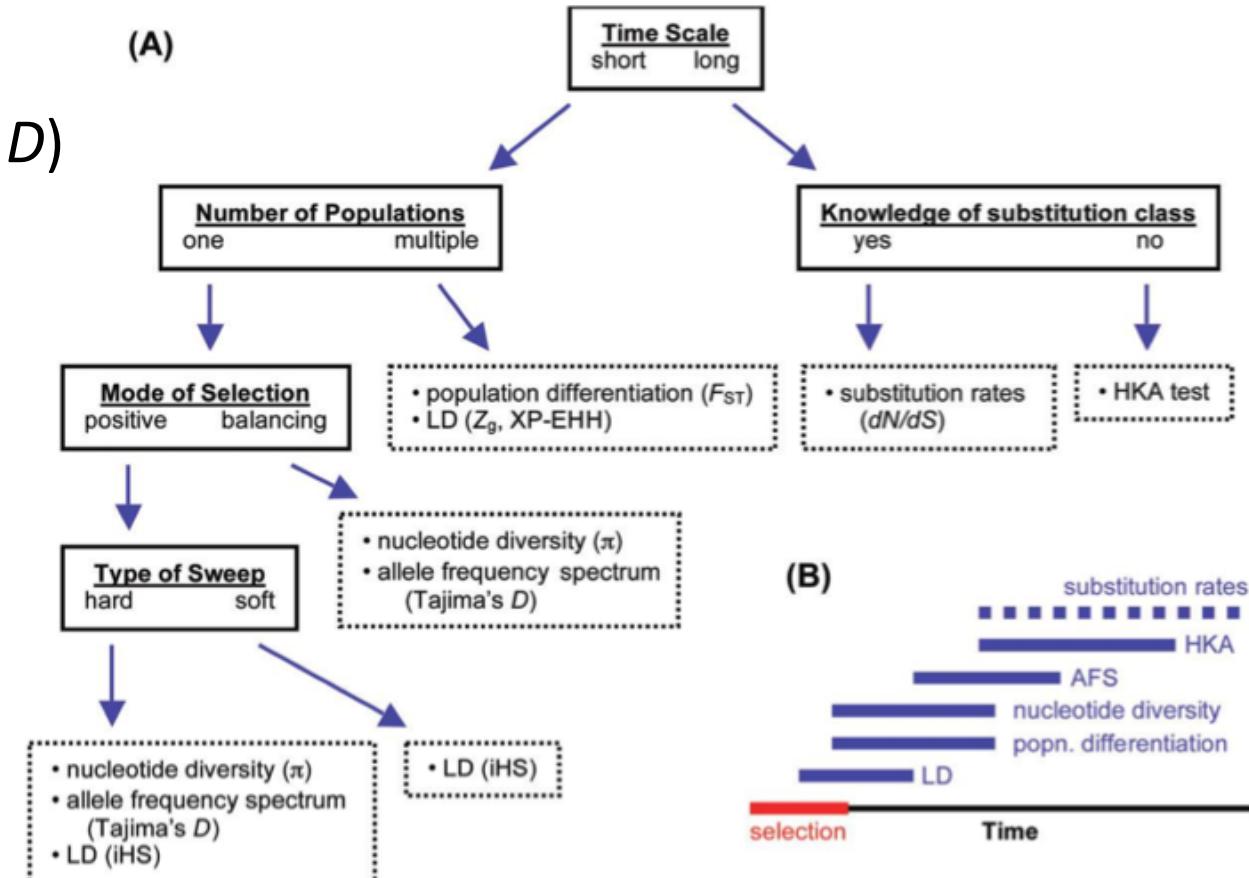
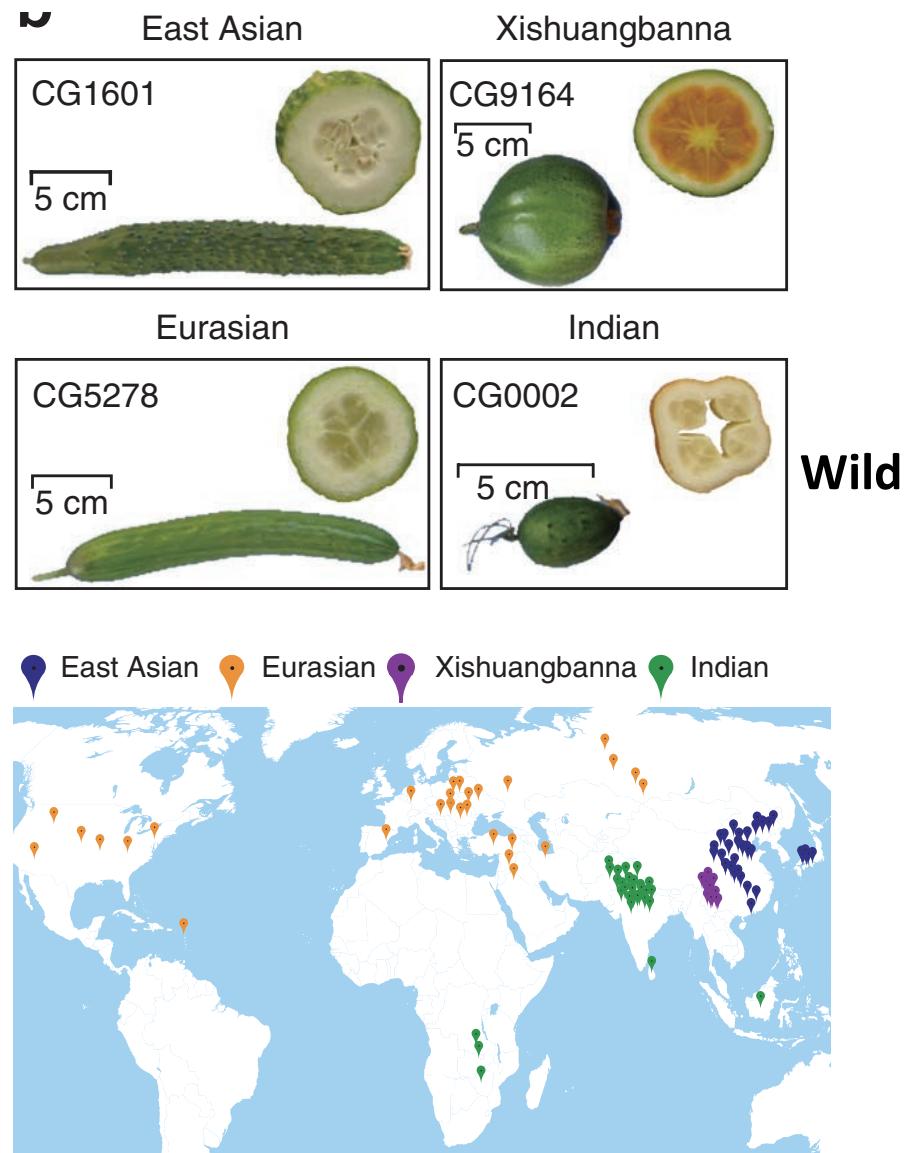


Figure 1 of Hohenlohe *et al.* (2010)

# Reduced level of genetic variation ( $\pi$ )

## Cucumber

- ❖ 3 cultivated (C) and 1 wild (W) cucumber groups
- ❖ Morphologically different
- ❖ Genome re-sequencing (n=115)
- ❖  $\pi$ : nucleotide diversity
  - mean number of nucleotide substitutions per site between any two randomly selected DNA sequences in a population



# Reduced level of genetic variation ( $\pi$ )

Domestication sweep

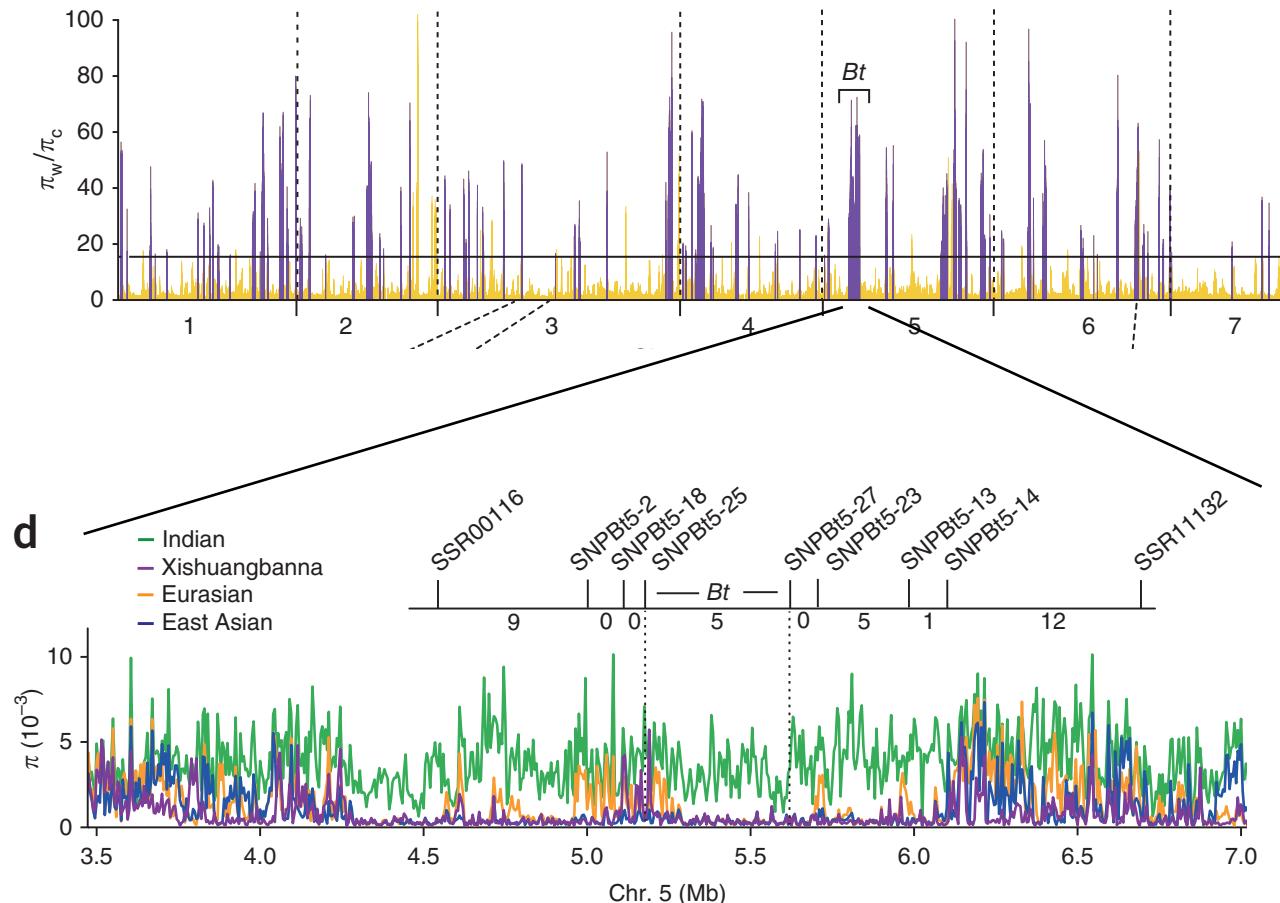
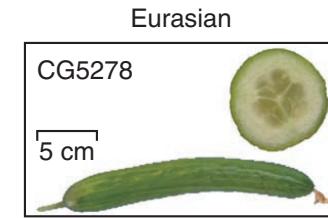
- ❖ Comparing  $\pi_w/\pi_c$
- ❖ 112 regions detected

❖ **Bt** locus

- Fruit bitterness

❖ Reduced  $\pi$  in cultivated cucumbers

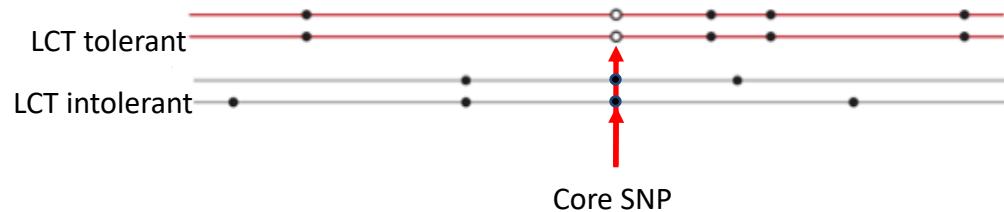
- ❖ 2 Mb in length (Bt: 442 kb)



# Linkage disequilibrium (LD)

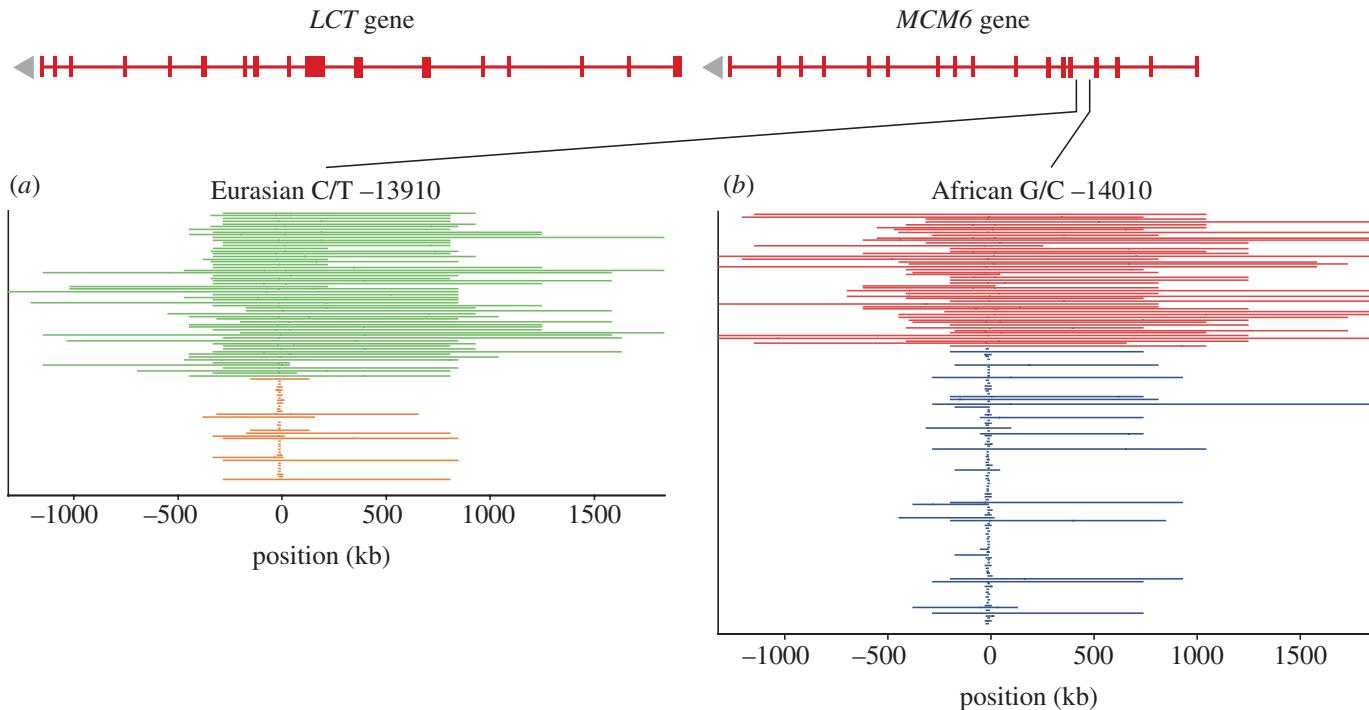
EHH: extended haplotype homozygosity

- Measures the decay of homozygosity from a 'core' SNP



LCT: lactose persistent gene in humans

- 2 Mb haplotypes
- Independent evolution in Africa and Europe



# Skew of allele frequency spectra (Tajima's $D$ )

## Tajima's $D$

- ❖ Normalized difference between  $\pi$  and segregating sites ( $S$ ,  $\Theta_W$ )

- ❖  $d = \pi - \Theta_W$

$$D = \frac{d}{\sqrt{V(d)}}$$

- ❖ **Balancing selection**

excess intermediate-freq. SNPs:

$$\pi > \Theta_W, +D$$

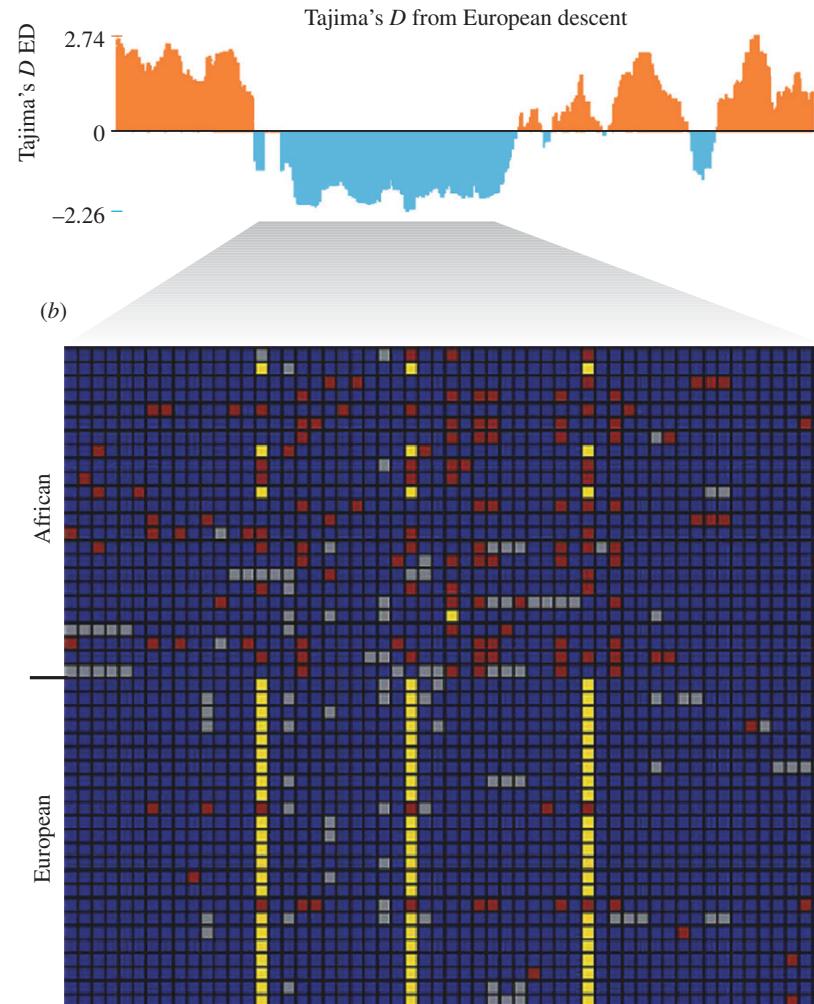
- ❖ **Positive selection**

excess of low-frequency SNPs:

$$\pi < \Theta_W, -D$$

- ❖ **Correct for demographic effects**

- Compare values against genome-wide estimates



# Skew of allele frequency spectra

Human *CLSPN* gene

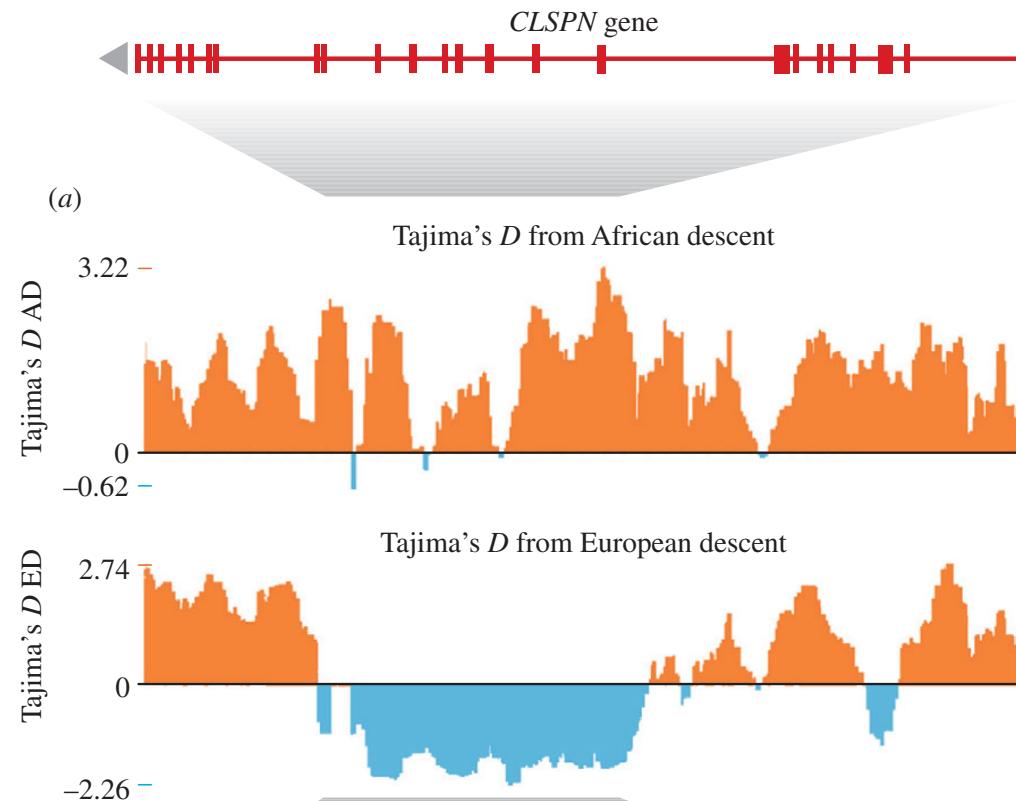
❖ Inferred from dense SNP data

❖ Tajima's *D* plot

❖ Positive selection in European ( $-D$ )

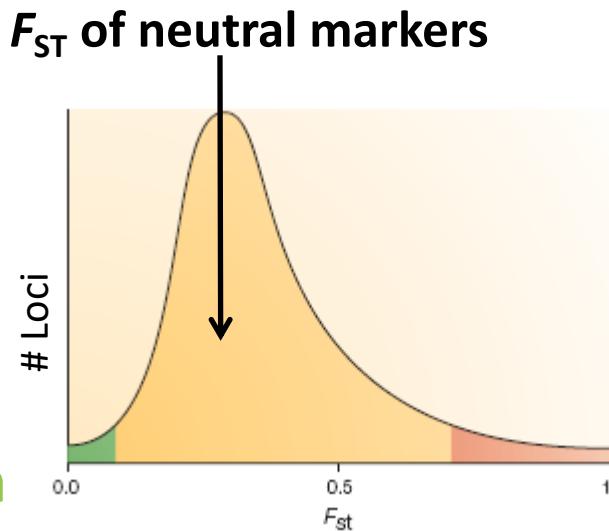
❖ 1.5 Mb

❖ Unknown function!



# Model free $F_{ST}$ -outlier approach

- >3 populations
- Screen many loci (>10,000; up to whole genome)
- **Outliers:** e.g. 95% quantile of  $F_{ST}$  distribution



## Balancing selection

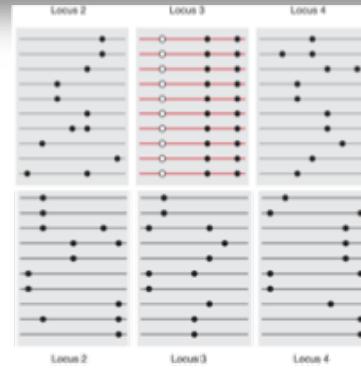
- relatively uniform frequencies across populations  
⇒ low  $F_{ST}$  values

Fischer *et al.* 2014 PLoS One

## Positive selection

- increased level of differentiation among populations  
⇒ high  $F_{ST}$  values

Fischer *et al.* 2011/2013 MolEcol



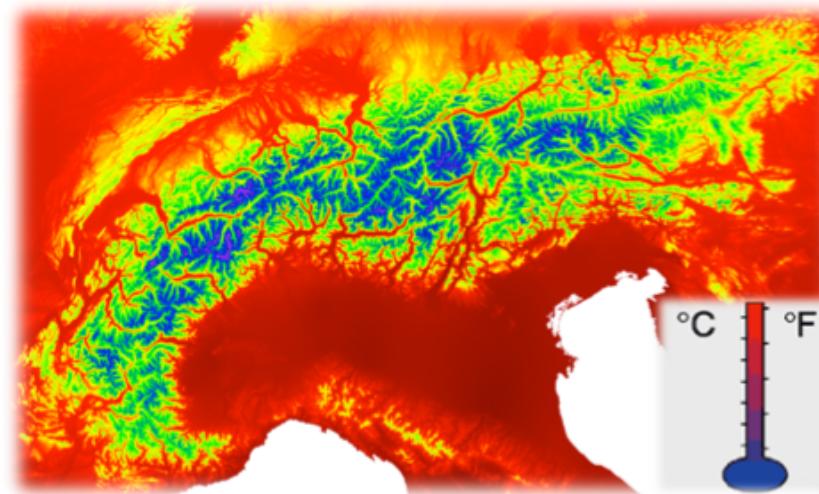
Luikart *et al.* 2003



- ❖ Highly heterogeneous
- ❖ Strong environmental gradients
- Genetic basis of adaptation

Population genomic footprints of selection and associations with climate in natural populations of *Arabidopsis halleri* from the Alps

MARTIN C. FISCHER,\* CHRISTIAN RELLSTAB,† ANDREW TEDDER,‡ STEFAN ZOLLER,§  
FELIX GUGERLI,† KENTARO K. SHIMIZU,‡ ROLF HOLDEREGGER\*† and ALEX WIDMER\*



# Study organism

---

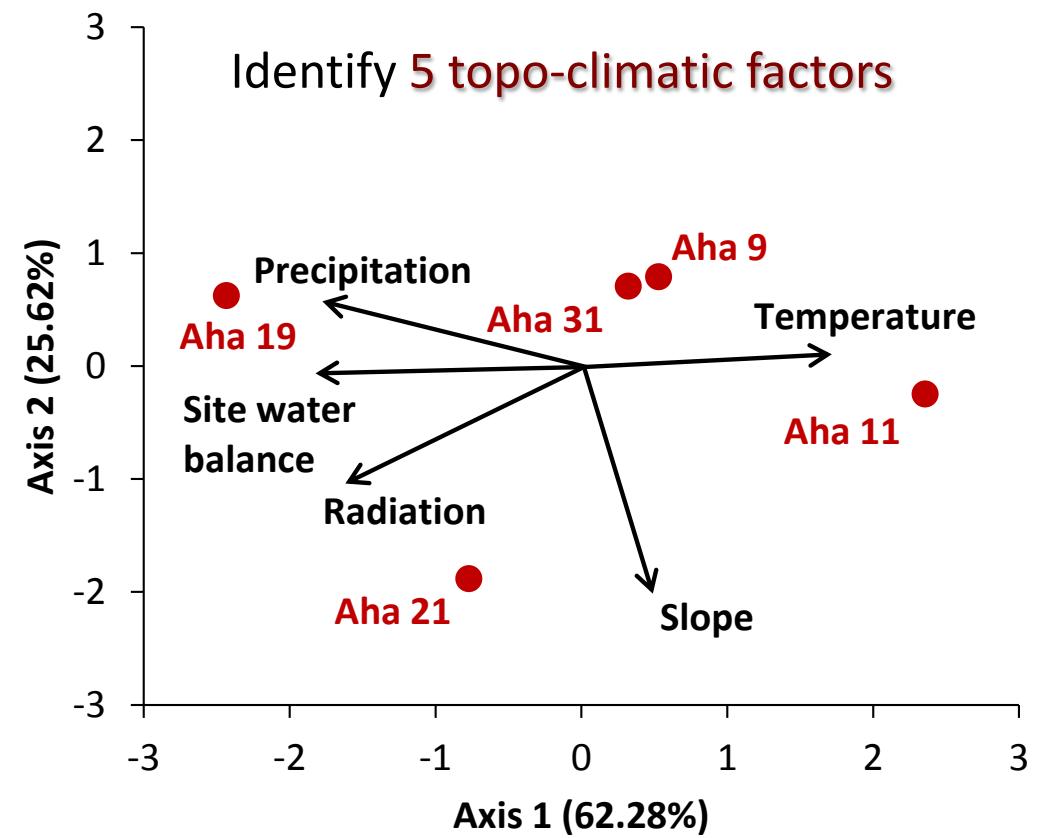
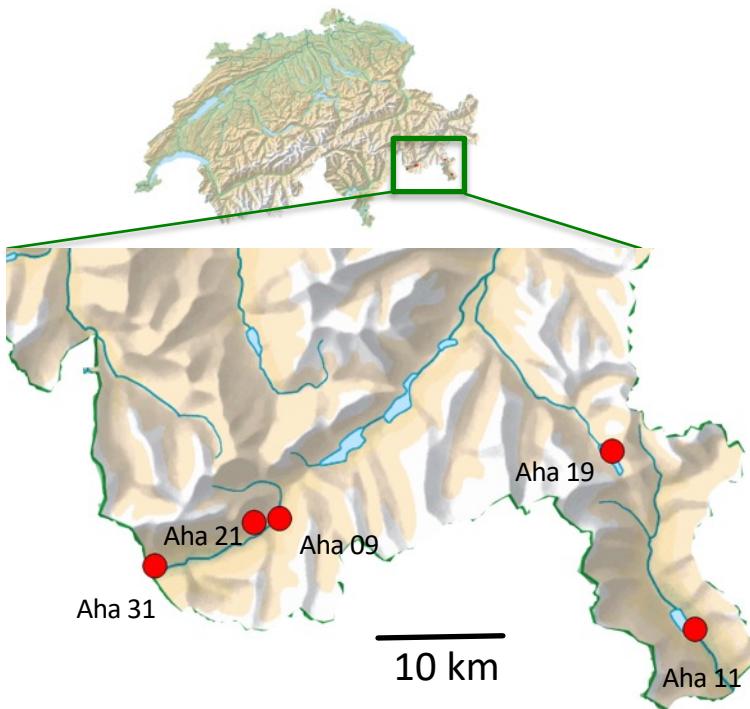
## *Arabidopsis halleri*

- ❖ Close relative of the model organism *A. thaliana*
- ❖ Genome size 255 Mbp
- ❖ Strictly outcrossing
- ❖ 300 – 2400 m a.s.l.



# Heterogeneous Alpine environments

- ❖ 5 populations in close vicinity (2 – 45 km)
- ❖ 20 individuals each
- ❖ Cover wide range of abiotic environmental condition
  - E.g. 790 – 2308 m a.s.l.

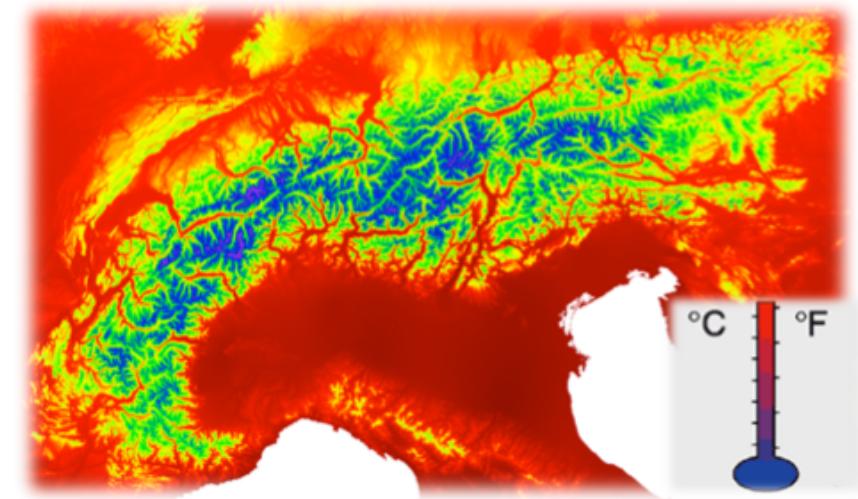


# Adaptive genomics

---

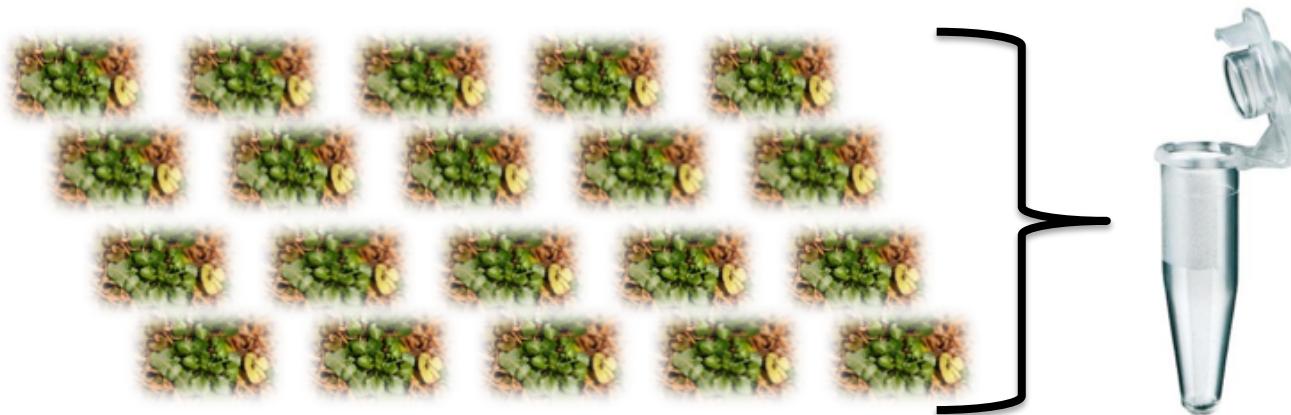
Many abiotic factors are changing on short distance

- ❖ Expected **signature of selection?**
  - No signature of selection
  - Some genes with major effects
  - Many different genes under selection
  
- ❖ Which **genes** are involved in **adaptation?**
  - Gene functions?
  - Abiotic factors?
  
- Population genomic approach



# Whole genome re-sequencing

- No bias from insufficient marker density or distribution
- ❖ **Pool-Seq**; pooled population approach
  - Cost effective => 5 libraries
  - Reduces amount of DNA required
- ❖ 20 diploid genomes pooled per population



# Accuracy of Pool-Seq approach

OPEN  ACCESS Freely available online

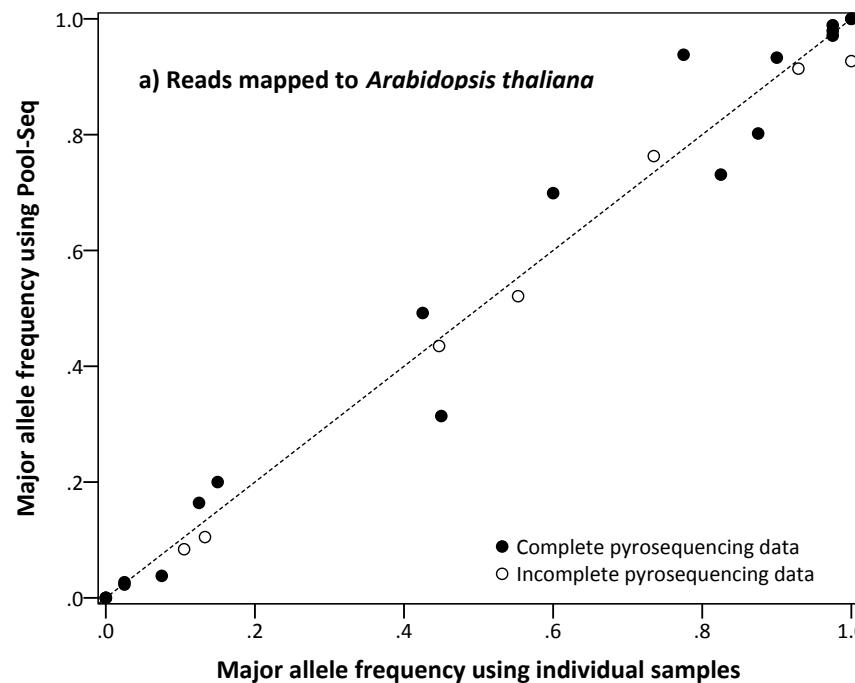


## Validation of SNP Allele Frequencies Determined by Pooled Next-Generation Sequencing in Natural Populations of a Non-Model Plant Species

Christian Rellstab<sup>1</sup>, Stefan Zoller<sup>2</sup>, Andrew Tedder<sup>3</sup>, Felix Gugerli<sup>1</sup>, Martin C. Fischer<sup>4\*</sup>

**1** Biodiversity and Conservation Biology, Swiss Federal Research Institute WSL, Birmensdorf, Switzerland, **2** Genetic Diversity Centre, ETH Zürich, Zürich, Switzerland, **3** Institute of Evolutionary Biology and Environmental Studies and Institute of Plant Biology, University of Zürich, Zürich, Switzerland, **4** Institute of Integrative Biology, ETH Zürich, Zürich, Switzerland

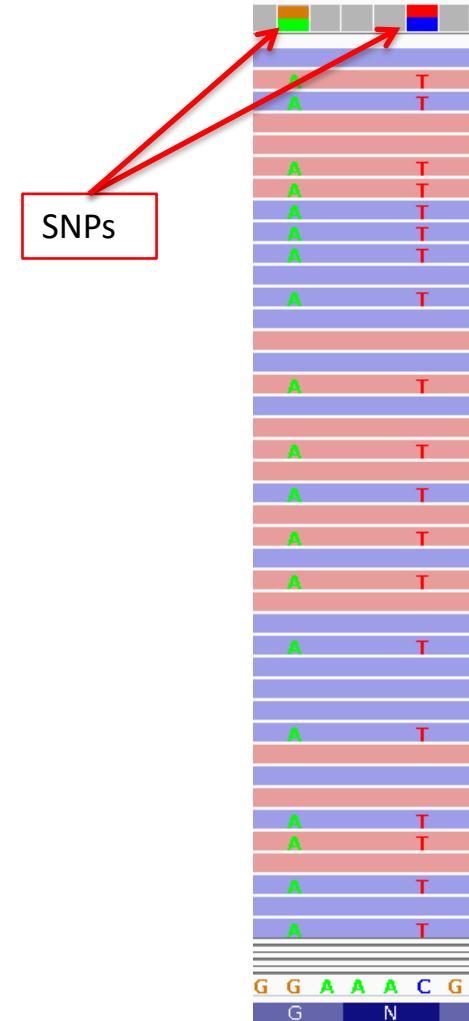
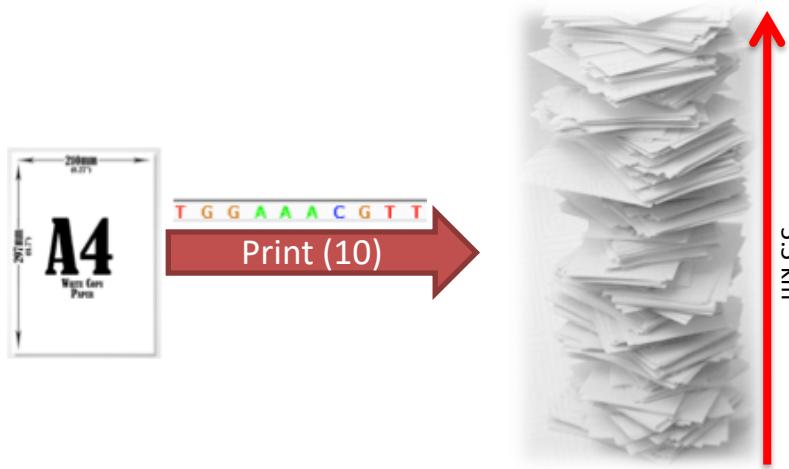
- ❖ Individual SNP genotyping
- ❖ 3 populations
- ❖ 9 SNPs validated
- ❖ PyroMark
- ❖  $R^2 = 0.98$



# Population genomics

## ❖ Mapping of reads onto *A. thaliana* ref. genome

- BWA (Li & Durbin 2009)
- ~60x coverage
- >120 billion bases



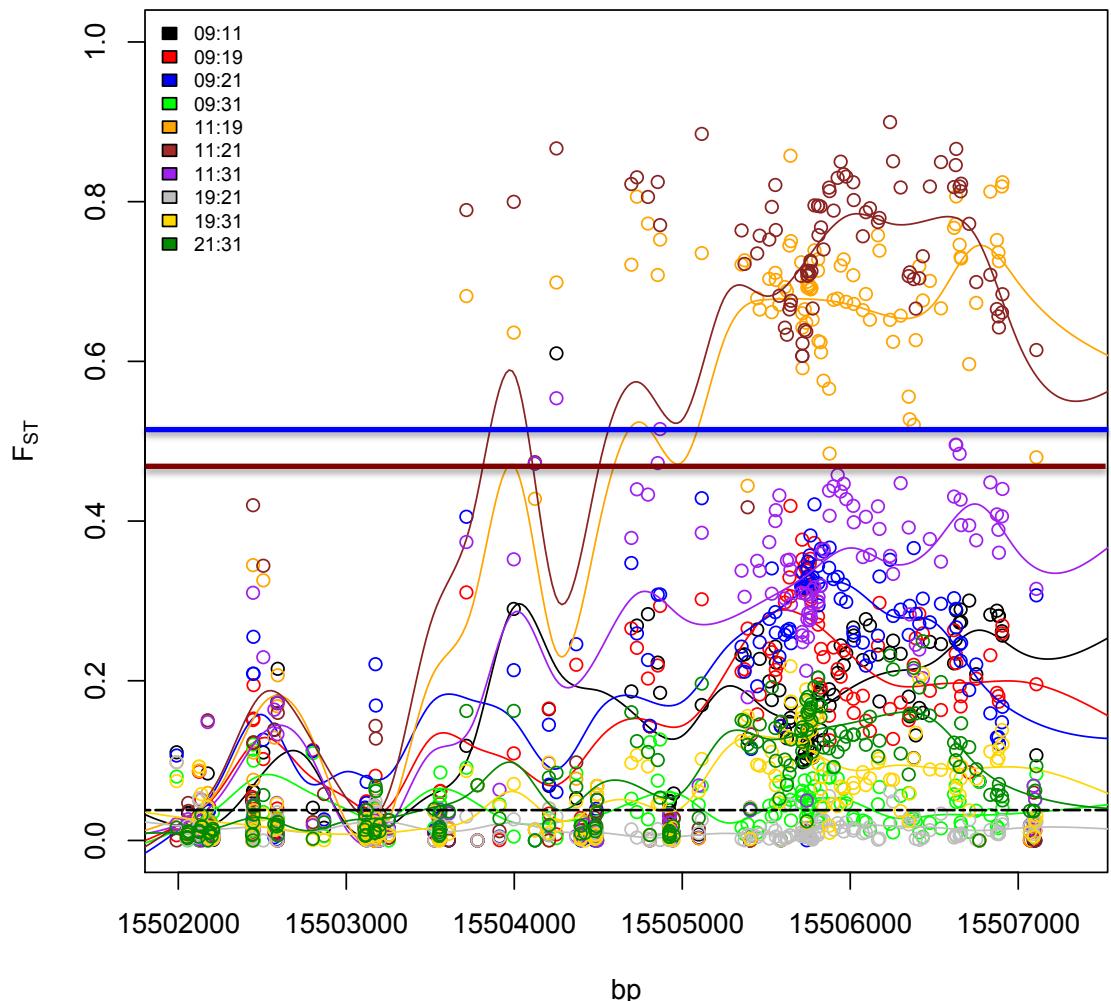
## ❖ PoPopulation2 (Kofler *et al.* 2011)

- Identified **2,091,957 SNPs**
- 25,764 genes covered
- Genome-wide pairwise  $F_{ST}$ 
  - Mean: 0.038

# $F_{ST}$ -outlier detection

## 3 step model free approach:

- ❖ 0.1% of highest  $F_{ST}$  regions:  $F_{ST} > 0.47$ 
  - $F_{ST}$  sliding windows approach of 500 bp
- ❖ 0.1% of highest  $F_{ST}$  SNPs:  $F_{ST} > 0.54$ 
  - Corrected for population structure
- ❖ Corrects for coverage
  - Fisher's exact test
  - $p < 2.39 \times 10^{-9}$
  - Strong allele frequency differences

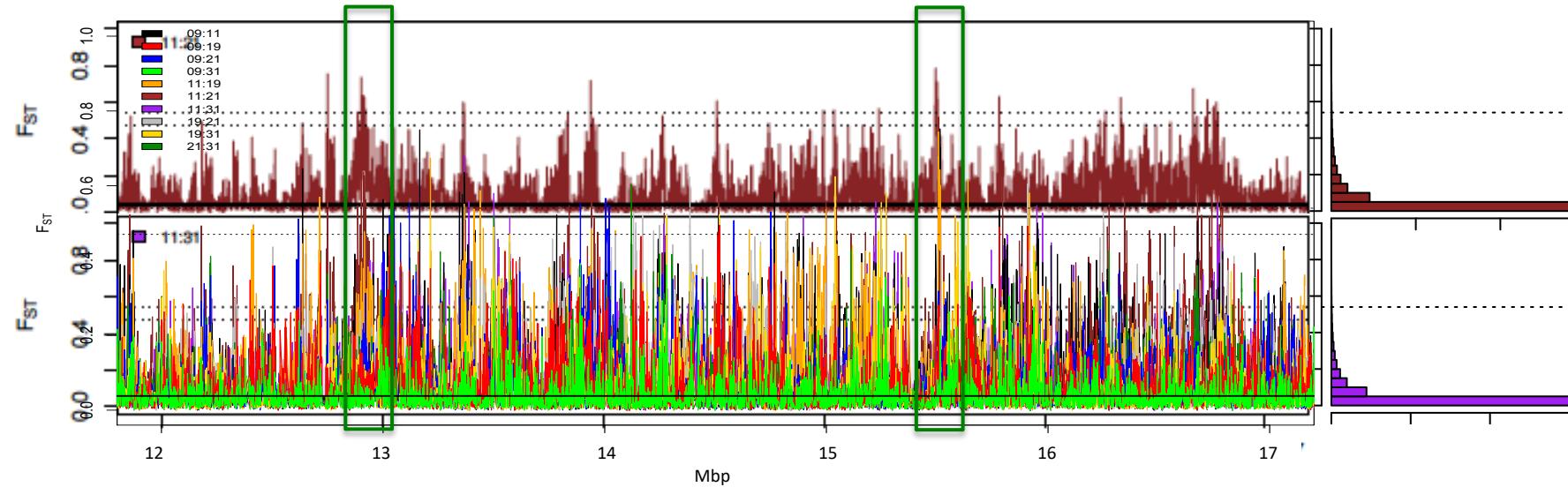


# Genomic Signature of Selection

❖ 4282 strong outlier SNPs

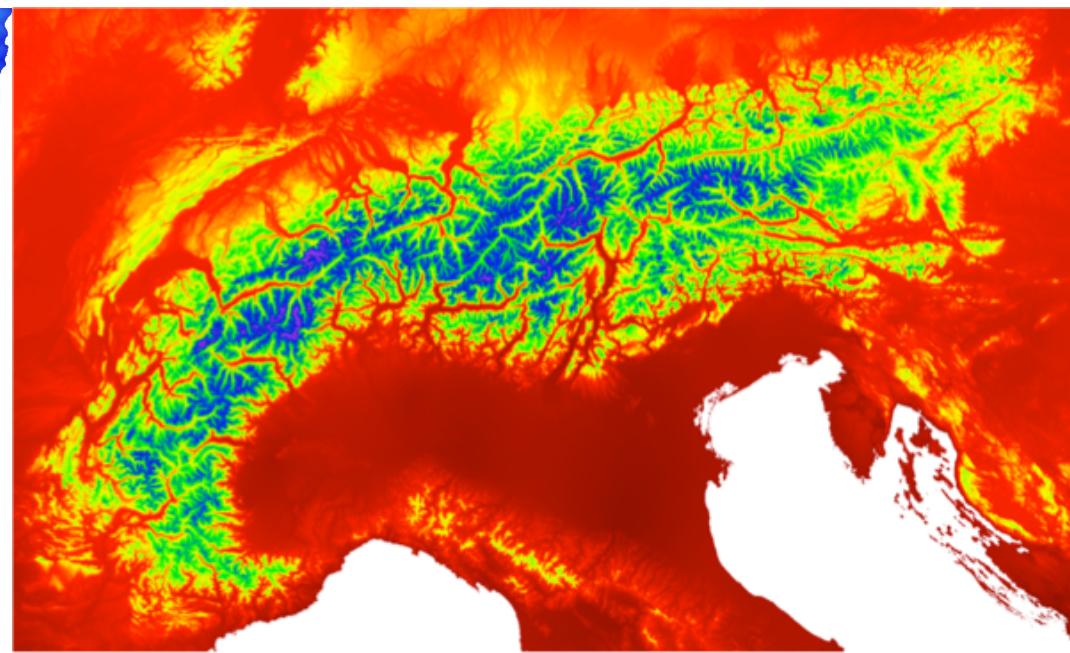
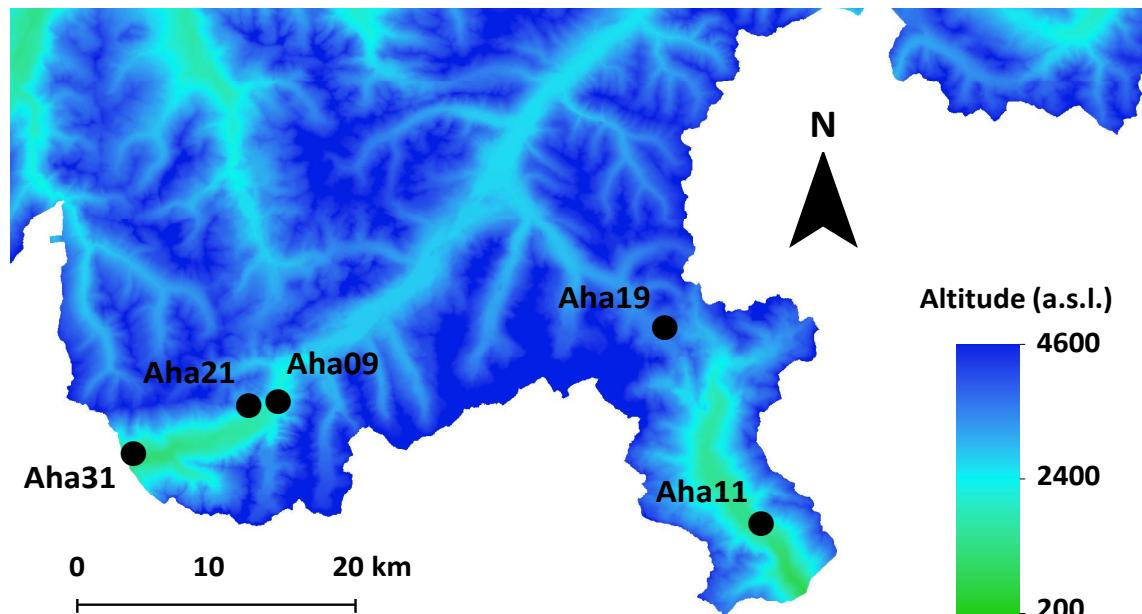
➤ 0.2% of all SNPs

❖ 571 outlier genes



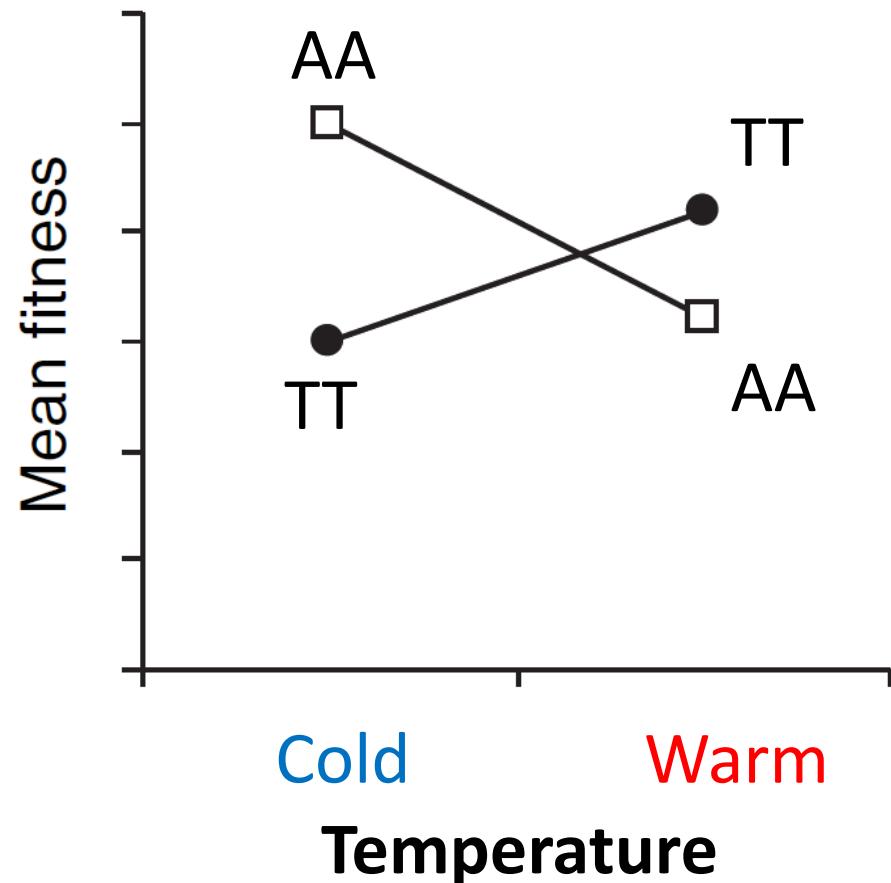
# Gene functions

- ❖ 571 outlier genes
  - 139 genes (24%) unknown function
- ❖ Identify genes involved in environmental adaptation



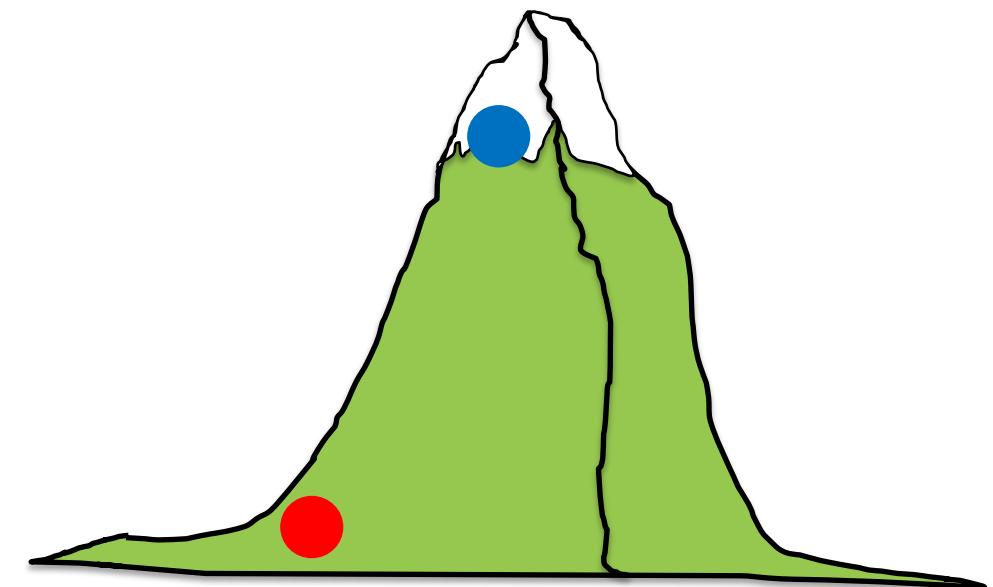
# Environmental association analyses

- Population from cold habitat
- Population from warm habitat



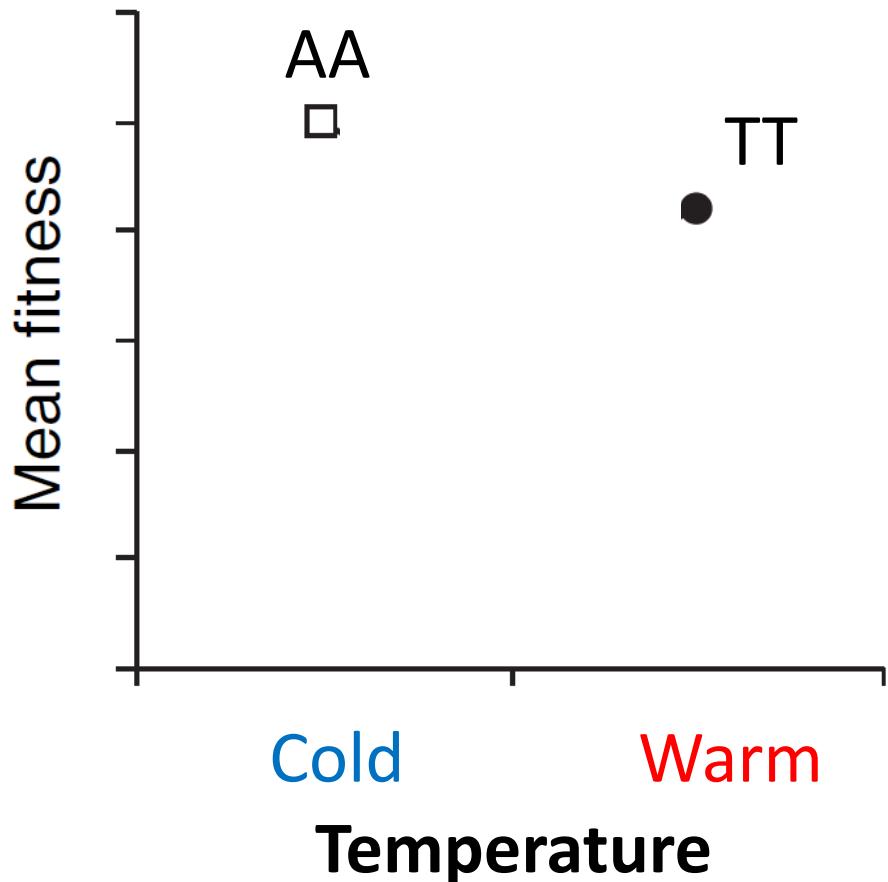
ACGTGTAGCAGTAGC**AT**GATGCTGATCGATT  
ACGTTAGCAGTAGC**TT**GATGCTGATCGTTT

Fitness difference is genetically controlled



# Environmental association analyses

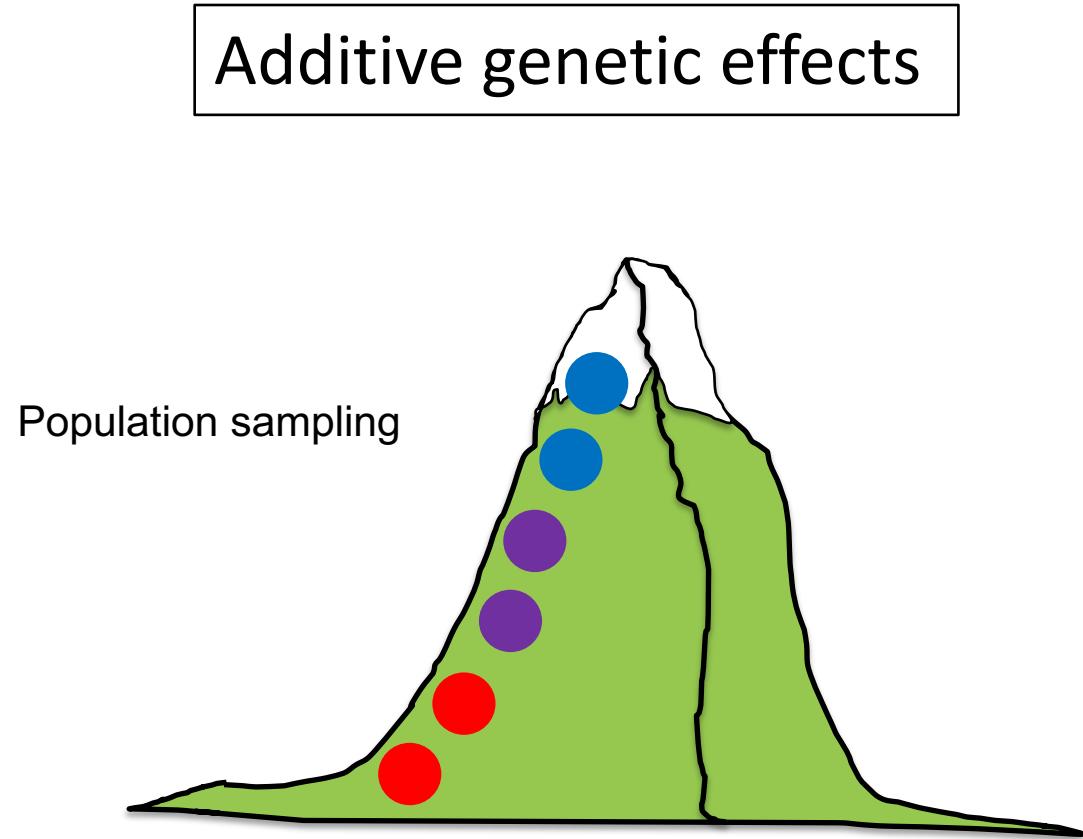
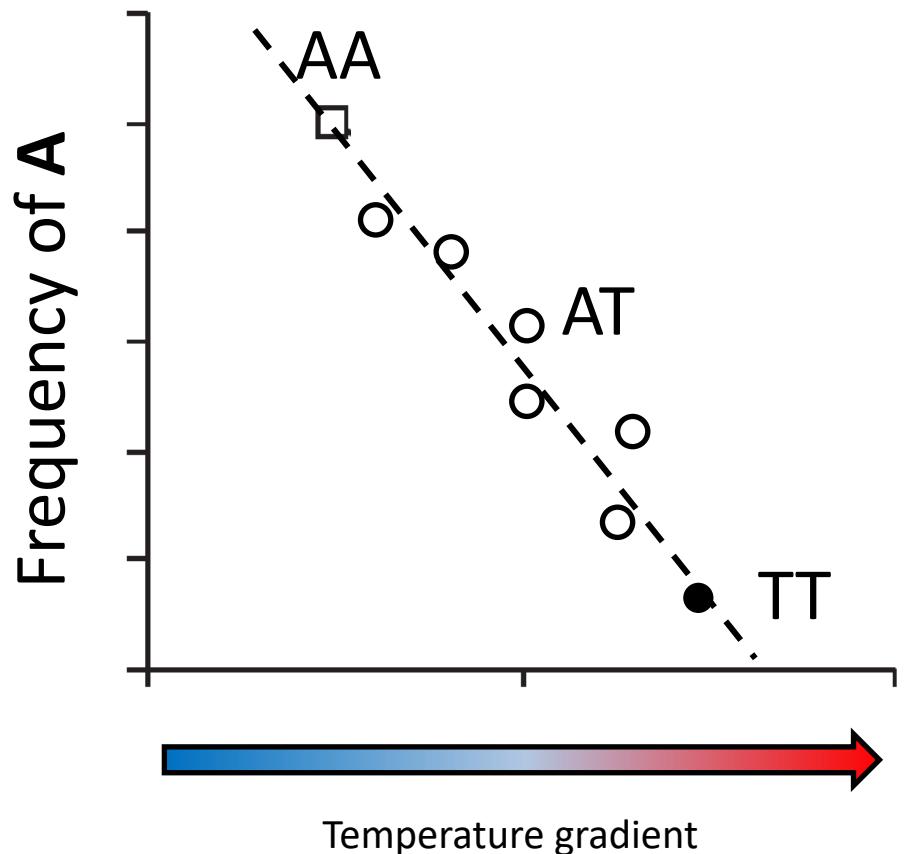
- Population from cold habitat
- Population from warm habitat



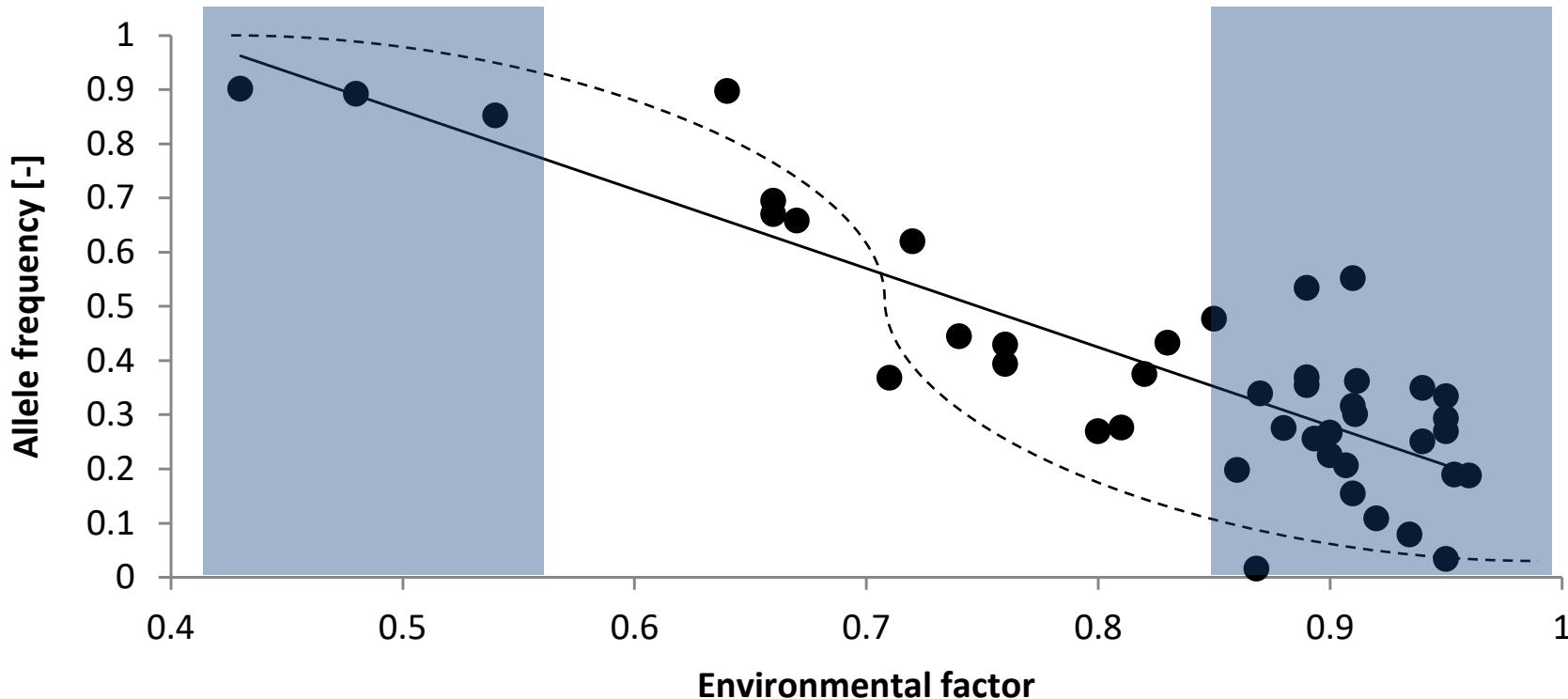
All populations are adapted to their local environmental conditions

# Environmental association analyses

- Population from cold habitat
- Population from warm habitat



# Environmental association analyses



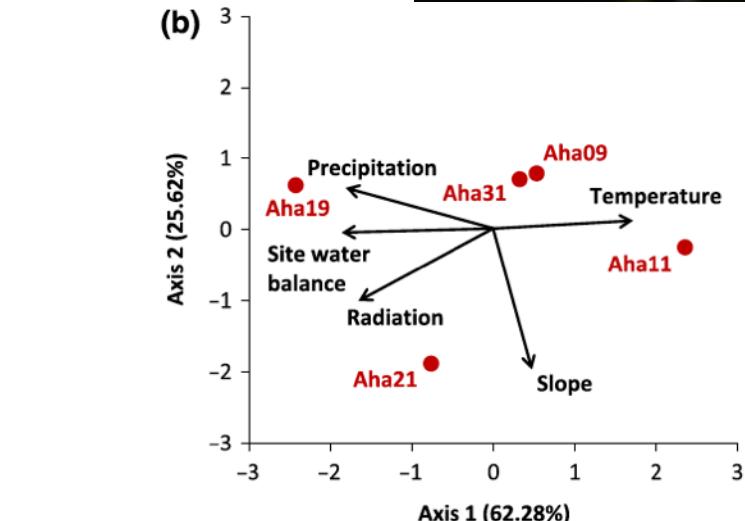
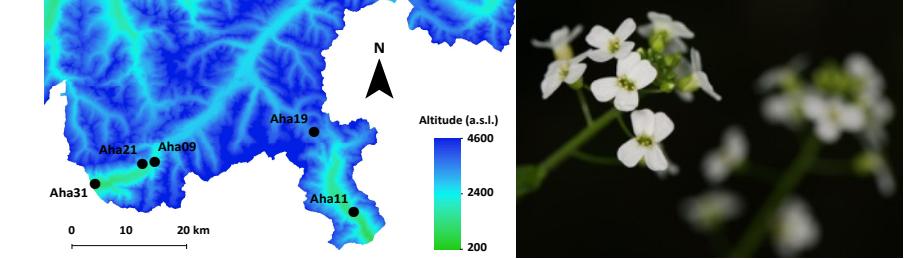
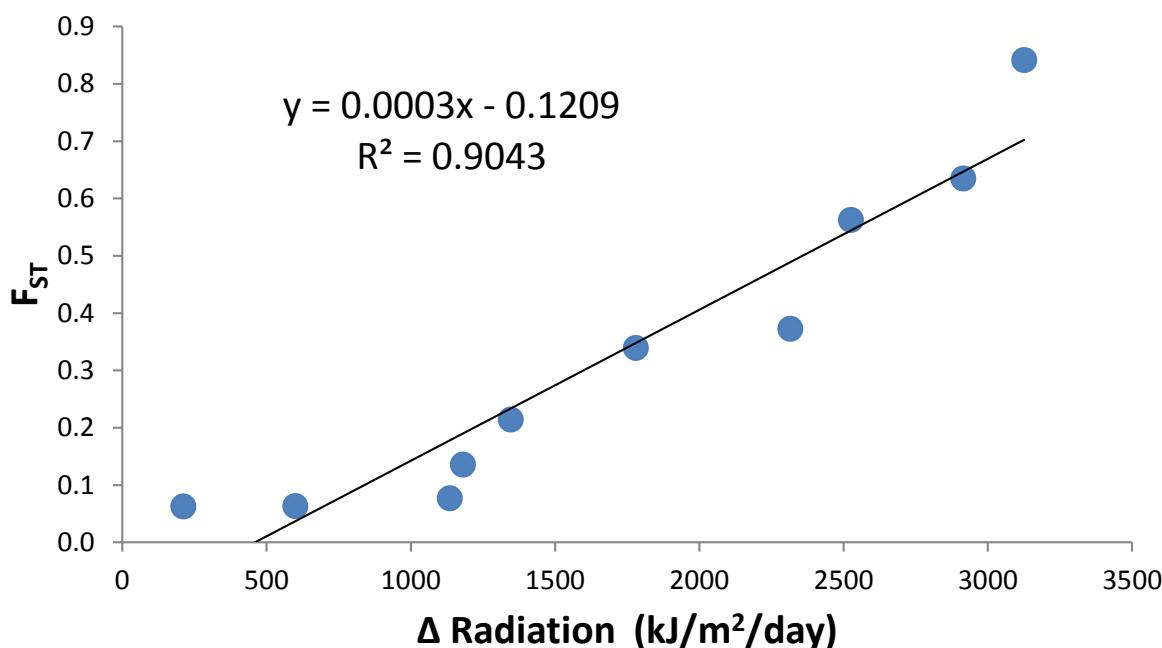
- ❖ Correlation/association of genotype with environmental factors
- ❖ Identify alleles, SNPs, or genes associated with particular environmental factors (temperature, rain, soil, UV, ...)
- ❖ Identify environmental factors that drive local adaptation

# Environmental Associations

## Partial Mantel tests

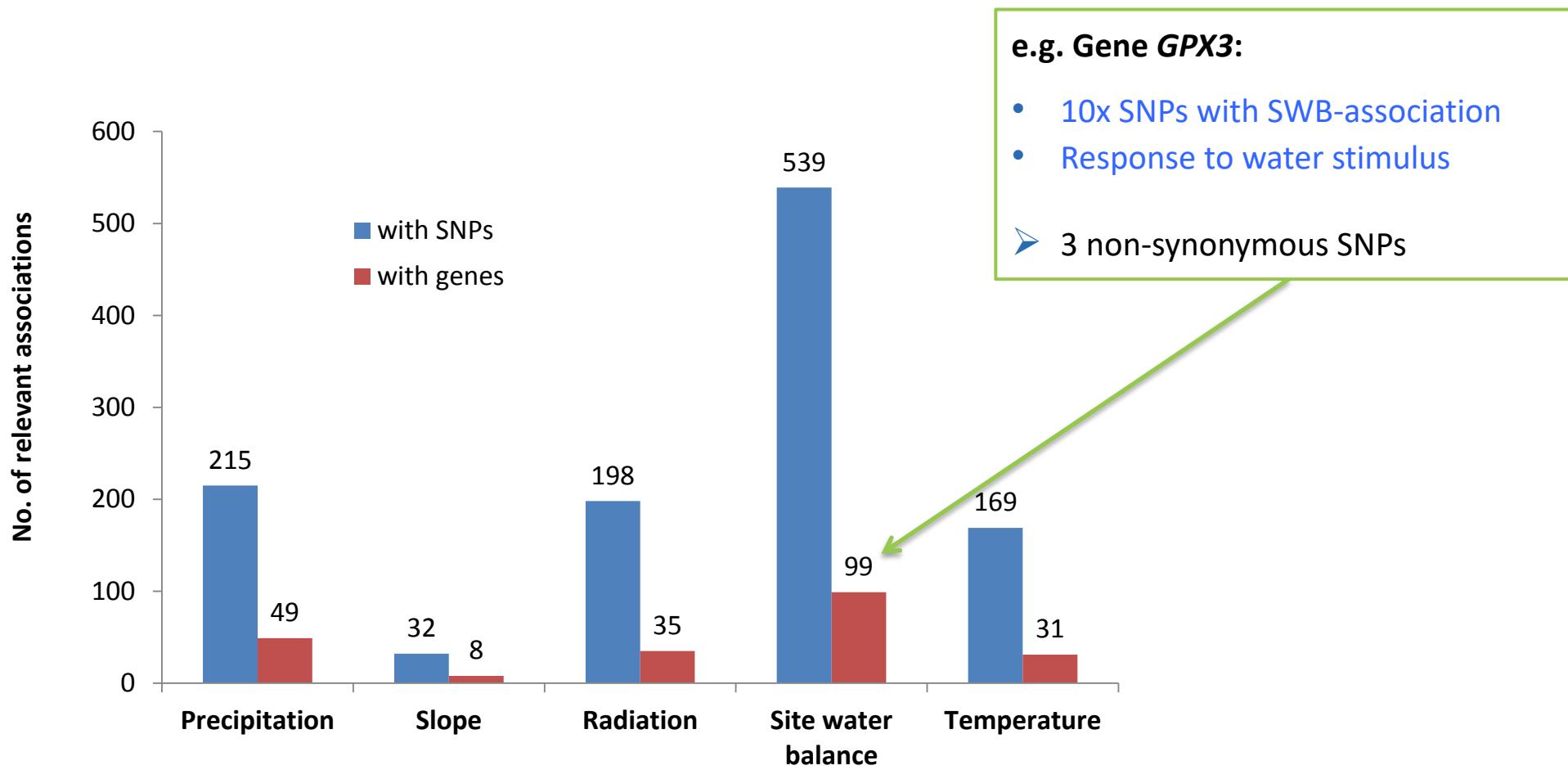
- ❖ Correlates two distance matrices:
  - Pairwise **genetic** distance of **outlier SNPs**
  - Pairwise **climatic** distance of **environmental factors**

- ❖ Controlling for population structure



# Abiotic environmental associations

❖ 175 genes with environmental associations (**30.6%**; out of 571 genes)



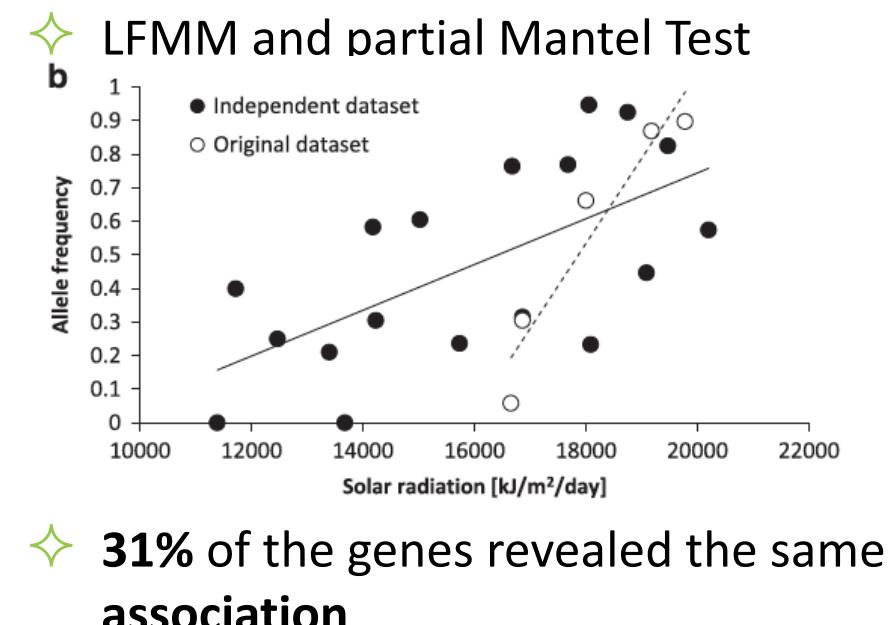
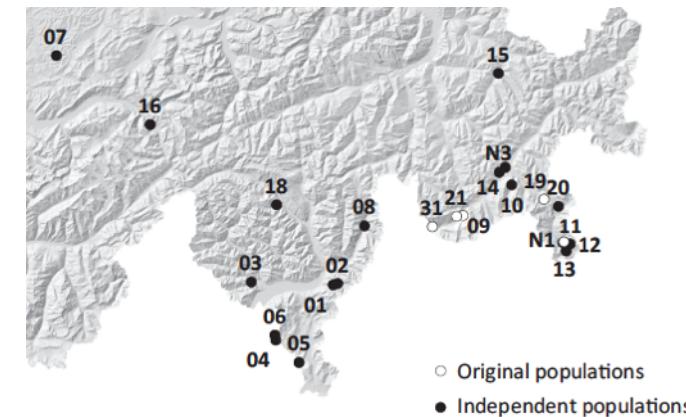
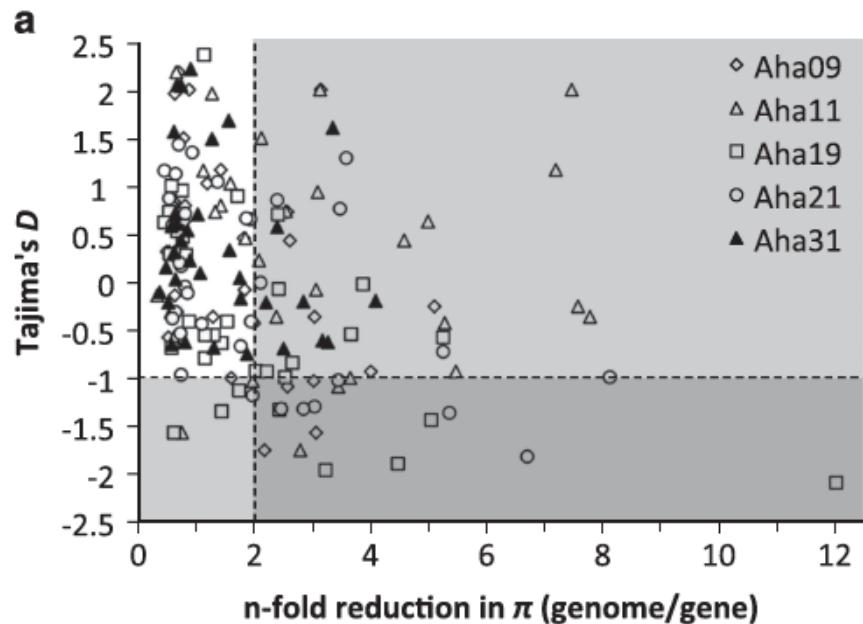
# Local remains local... ... in a highly heterogeneous environment

ORIGINAL ARTICLE

## Local adaptation (mostly) remains local: reassessing environmental associations of climate-related candidate SNPs in *Arabidopsis halleri*

C Rellstab<sup>1</sup>, MC Fischer<sup>2</sup>, S Zoller<sup>3</sup>, R Graf<sup>4</sup>, A Tedder<sup>4</sup>, KK Shimizu<sup>4</sup>, A Widmer<sup>2</sup>, R Holderegger<sup>1,2</sup> and F Gugerli<sup>1</sup>

- ❖ 74 candidate SNPs tested for EAA on a large geographic scale
- ❖ 444 individuals in **23 populations**
- ❖ **Signature of selection** was mainly observed in single populations



# Last glacial maximum (18,000 years)

- ❖ The molecular signature of adaptation in Alpine populations of *A. halleri* is **highly complex** and **local**
- ❖  $F_{ST}$ -outlier approach detect **recent selection**
- ❖ Selective signals mostly found in low elevation populations
- ❖ Alpine plants are already adapted to harsh environments
- ❖ Adaptation to **warm** and **dry** environments



# Model based $F_{ST}$ -outlier approach



Estimates probability of each locus to be under selection (Island model)

❖ **BayeScan** (Foll & Gaggiotti 2008; Fischer et al. 2011)

➤ Bayesian model comparison

$$PO = \frac{\Pr(M_1 | \text{Data})}{\Pr(M_2 | \text{Data})}$$

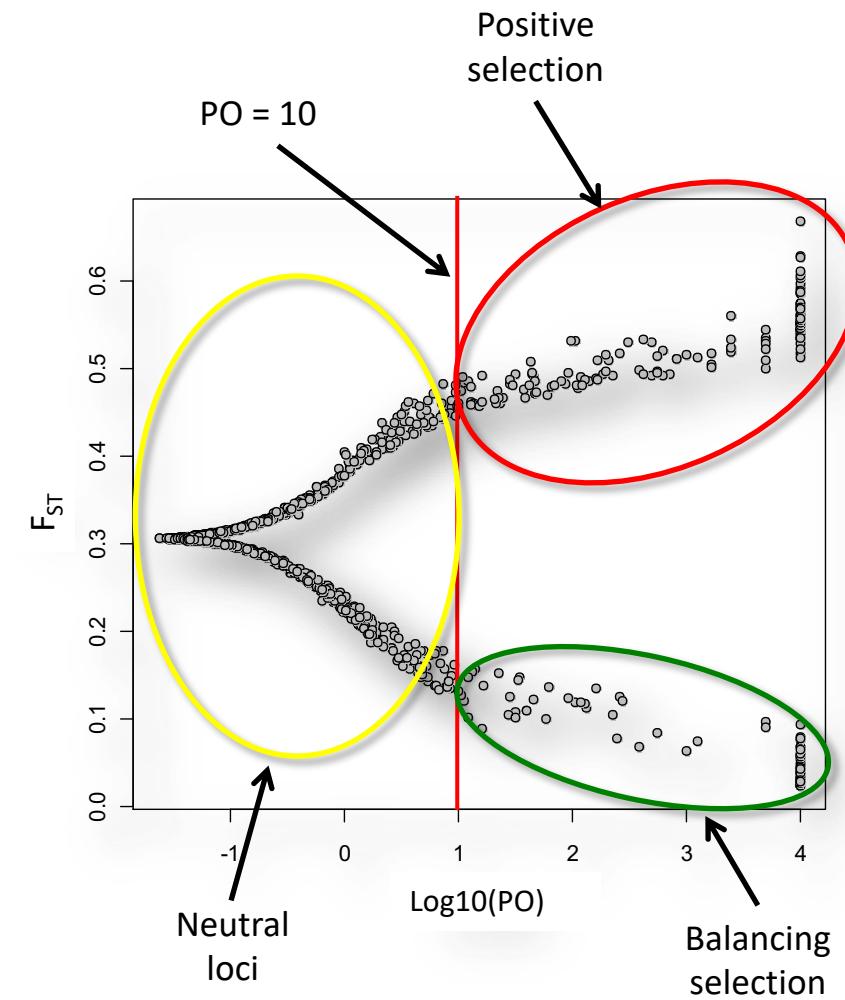
Posterior odds

Model with selection

Neutral model

A mathematical equation for Posterior Odds (PO) comparing two models. The numerator is the probability of Model 1 given the data, and the denominator is the probability of Model 2 given the data. Arrows point from the labels to the respective terms in the equation.

❖ **PO > 10** strong evidence for accepting a model  
(Jeffreys 1961)



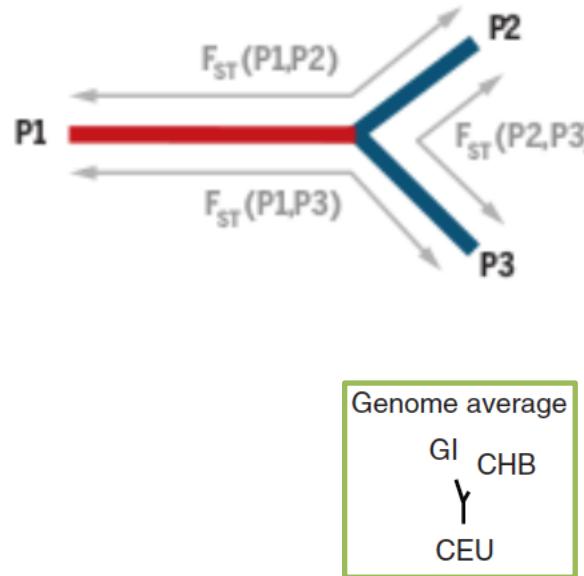
# Population Branch Statistics (PBS)

## Adaptation of Greenland Intuits to high altitude

- 191 Greenland Intuits (GI; P1)
- 44 Han Chines (CHB; P2)
- 40 European(CEU; P3)

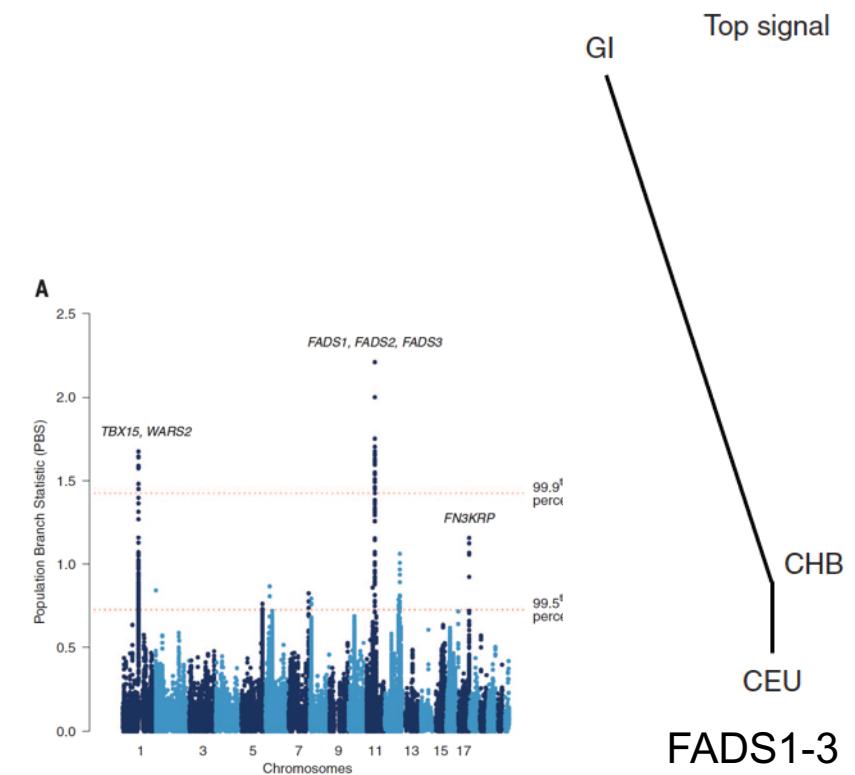
### ❖ Locus specific branch length

- $PBS = (T^{12} + T^{13} - T^{23})/2$
- $T^{12} = -\log(1 - F_{ST}^{12})$



### ❖ FADS1-3

- Fatty acid desaturase
- Mutation compensate for high-fat diet



# Genotype likelihoods (ANGSD)

- Medium/low coverage individual sequencing: **1x-10x**
- Ultra low coverage individual sequencing: **<1x**

Medium depth sequencing

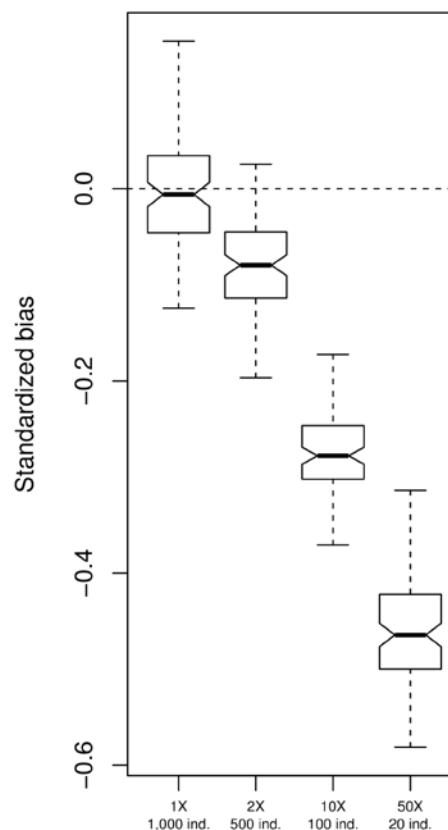


Ultra low depth sequencing

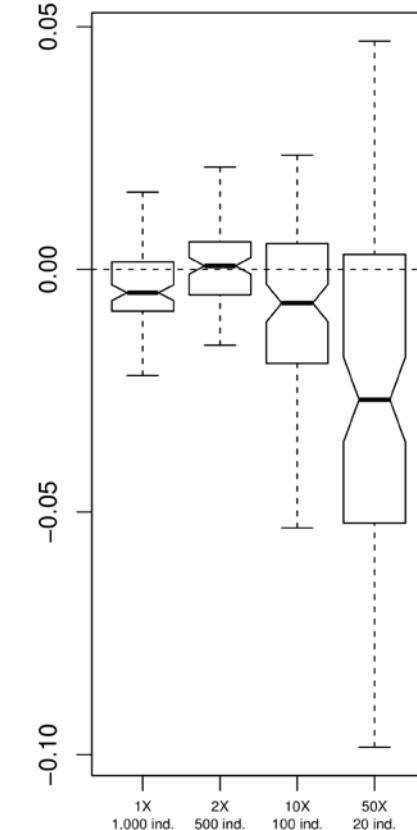


- Sequencing many individuals!
- Cheap library preparation, ~5.- CHF (Therkildsen & Palumbi 2017)

Number of segregating sites



Expected heterozygosity



# Genotype likelihoods

## ➤ Using genotype likelihoods (ANGSD) instead of genotypes for downstream analysis

Summarise the data in 10 genotype likelihoods

	A	C	G	T
A	1	2	3	4
C		5	6	7
G			8	9
T				10

The likelihood

$$P(\text{Data} | G = \{A_1, A_2\}) \propto P(X | G = \{A_1, A_2\}) = P(X | G)$$

where  $A \in \{A, C, G, T\}$

10 genotype likelihoods

	A	C	G	T
A	0.0	0.001	0.0	0.01
C		0.02	0.001	0.12
G			0.0	0.003
T				0.001

TCCTTTTTTT  
CCCTTTTTTT

## Base (b)

- Nr. of reads
- Base quality (Error)
- (Mapping quality)
- Pop. allele freq. as prior => **Genotype posterior** used for downstream analysis

# ANGSD (Genotype likelihoods)

ANGSD

**Table 1 Overview of analyses implemented in ANGSD**

Analysis	Basis	Reference
<b>Contamination estimates</b> based on the X-chromosomes	BC	[19] <sup>b</sup>
<b>Type specific error estimation</b> estimated by simultaneously estimating allele frequencies and genotype likelihoods	GL	[10]
<b>Type specific error estimation</b> based on an outgroup and a high quality genome	BC	[20] <sup>ab</sup>
<b>Genotype likelihoods (GL)</b> (diploids)	BC/Seq	[6,8,10,15]
<b>Allele frequencies</b> for a site	BC/GL/GP	[21] <sup>b</sup> [10]
<b>SNP discovery (LRT)</b> used for rejecting that the allele frequency is different from zero	GL	[10]
<b>Genotype posteriors (GP)</b> can be used for calling genotypes by specifying a cutoff	GL/SAF	[9,10]
<b>Sample allele frequencies (SAF)</b> the probability of all read data given the sample allele frequency	GL/GP	[9] <sup>b</sup>
Population differentiation statistics $F_{st}$	SAF	[14] <sup>ac</sup>
Population structure via principle components analysis ( <b>PCA</b> )	GP	[14] <sup>ac</sup>
<b>Admixture analysis (NGSadmix)</b> NGS data	GL	[22] <sup>ab</sup>
Detection of ancient admixture <b>ABBA-BABA/d-statistics</b>	BC	[20] <sup>b</sup>
Estimation of <b>SFS (1D)</b>	SAF	[9] <sup>ab</sup>
Estimation of <b>SFS (2D)</b>	SAF	
<b>Selection scans</b> , Neutrality tests (e.g $\theta$ 's and Tajima's D)	SAF	[12] <sup>ab</sup>
Estimation of individual and site-wise <b>Inbreeding</b> coefficients. Also MAF and GP estimation for inbred individuals	GL	[13] <sup>abc</sup>
<b>Allele frequency based association</b> for case/control data)	GL	[10]
<b>Association score test</b> in a generalized linear model framework for both quantitative and case/control data while allowing for additional covariates	GL-GP	[11] <sup>b</sup>

Table of the supported analyses in ANGSD. <sup>a</sup>indicates methods that require a secondary program in ANGSD package. <sup>b</sup>indicates methods for which ANGSD is the *de facto* implementation and <sup>c</sup>are user supplied extensions for ANGSD. The basis for each analysis is either the sequencing data (Seq), base counts (BC), genotype likelihood (GL), sample allele frequencies (SAF) or genotype probabilities (GP).

# The genomic signature of selection

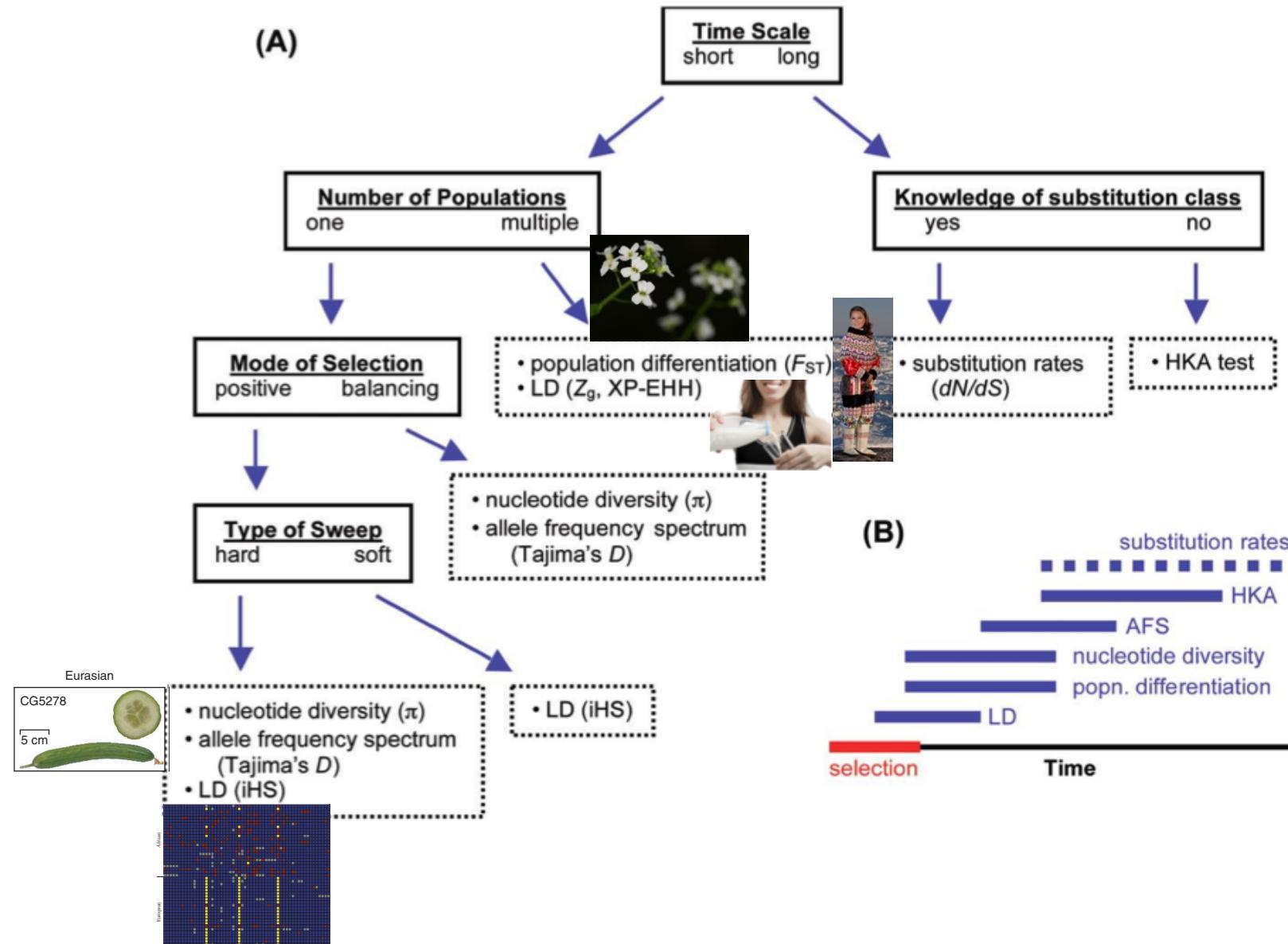


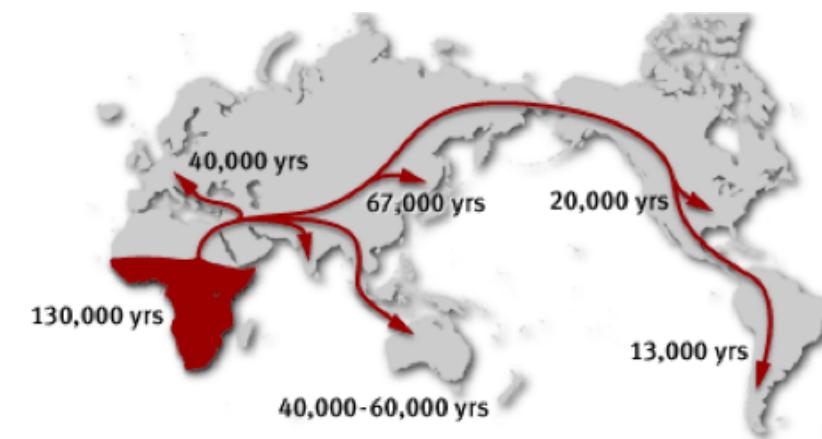
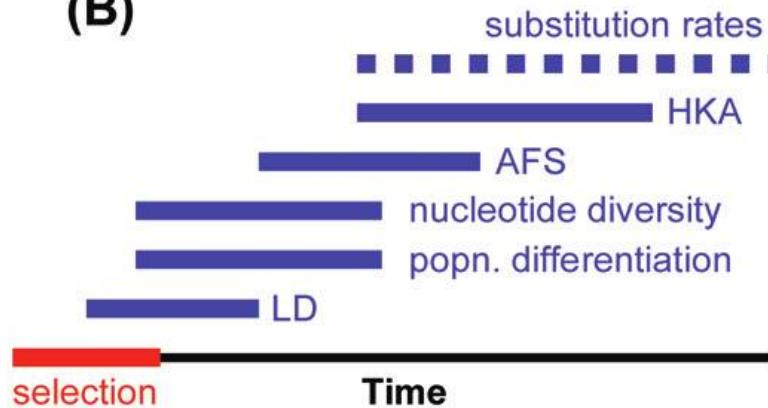
Figure 1 of Hohenlohe et al. 2010

# Different time scale of selection tests

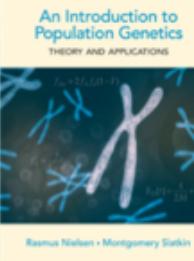
Example in humans

❖ Substitution rates ( $dN/dS$ )	>1,000,000 years	~ >50,000 generations
❖ Skew of allele frequency spectra	<200,000 years	~ <10,000 generations
❖ Reduced level of genetic variation ( $\pi$ )	<200,000 years	~ <10,000 generations
❖ Population differentiation ( $F_{ST}$ -outlier)	<80,000 years	~ <4,000 generations
❖ Linkage disequilibrium (LD)	<30,000 years	~ <1,500 generations

(B)



# Further Reading

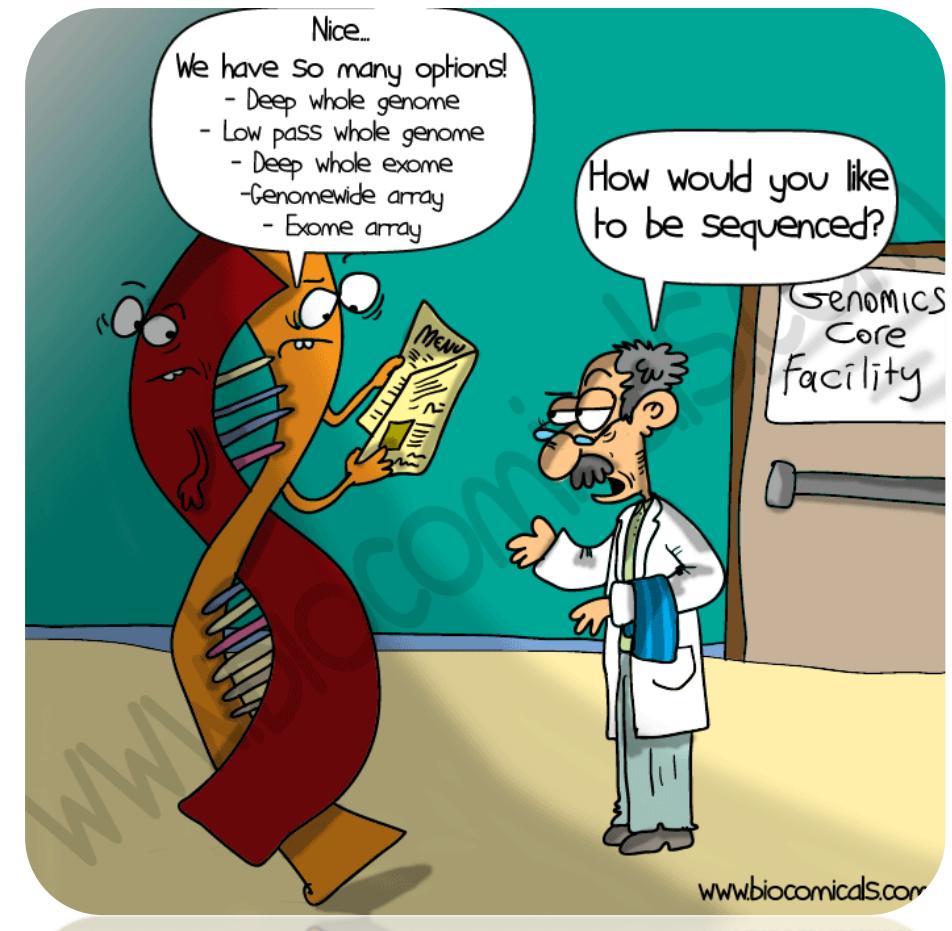


- Nielsen R & Slatkin M (2013) An Introduction to Population Genetics: Theory and applications. *Sinauer*
- Cutter AD, Payseur BA (2013) Genomic signatures of selection at linked sites: unifying the disparity among species. *Nature Reviews Genetics* **14**, 262-274.
- Fischer MC, Rellstab C, Leuzinger M, et al. (2017) Estimating genomic diversity and population differentiation – an empirical comparison of microsatellite and SNP variation in *Arabidopsis halleri*. *BMC Genomics* **18**, 69.
- Fischer MC, Foll M, Heckel G, Excoffier L (2014) Continental-scale footprint of balancing and positive selection in a small rodent (*Microtus arvalis*). *PLoS One* **9**, e112332.
- Fischer MC, Rellstab C, Tedder A, et al. (2013) Population genomic footprints of selection and associations with climate in natural populations of *Arabidopsis halleri* from the Alps. *Molecular Ecology* **22**, 5594-5607.
- Fischer MC, Foll M, Excoffier L, Heckel G (2011) Enhanced AFLP genome scans detect local adaptation in high-altitude populations of a small rodent (*Microtus arvalis*). *Molecular Ecology* **20**, 1450-1462.
- Fumagalli M (2013) Assessing the effect of sequencing depth and sample size in population genetics inferences. *PLoS One* **8**, e79667.
- Fumagalli M, Moltke I, Grarup N, et al. (2015) Greenlandic Inuit show genetic signatures of diet and climate adaptation. *Science* **349**, 1343-1347.
- Hohenlohe PA, Phillips PC, Cresko WA (2010) Using population genomics to detect selection in natural populations: key concepts and methodological considerations. *International Journal of Plant Sciences* **171**, 1059-1071.
- Korneliussen TS, Albrechtsen A, Nielsen R (2014) ANGSD: analysis of next generation sequencing data. *BMC Bioinformatics* **15**, 1-13.
- Oleksyk TK, Smith MW, O'Brien SJ (2010) Genome-wide scans for footprints of natural selection. *Philosophical Transactions of the Royal Society B: Biological Sciences* **365**, 185-205.
- Storz JF (2005) Using genome scans of DNA polymorphism to infer adaptive population divergence. *Molecular Ecology* **14**, 671-688.
- Rellstab C, Fischer MC, Zoller S, et al. (2017) Local adaptation (mostly) remains local: reassessing environmental associations of climate-related candidate SNPs in *Arabidopsis halleri*. *Heredity* **118**, 193-201.
- Rellstab C, Zoller S, Tedder A, Gugerli F, Fischer MC (2013) Validation of SNP allele frequencies determined by pooled next-generation sequencing in natural populations of a non-model plant species. *PLoS One* **8**, e80422.
- Rellstab C, Gugerli F, Eckert AJ, Hancock AM, Holderegger R (2015) A practical guide to environmental association analysis in landscape genomics. *Molecular Ecology* **24**, 4348-4370.
- Yi X, Liang Y, Huerta-Sanchez E, et al. (2010) Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* **329**, 75-78.
- Winter School: **Bioinformatics for Adaption Genomics (B@G)**
- Bioinformatics resources: <http://www.adaptation.ethz.ch/education/bag-winter-school-2018/teaching-resources.html>

# Thanks!

© Original Artist

Reproduction rights obtainable from  
[www.CartoonStock.com](http://www.CartoonStock.com)



**Box 1****Critical Population Genetic Concepts and Statistical Measures Used to Detect Selection in Population Genomics**

**Allele frequency spectrum (AFS):** The distribution of frequencies across alleles in a sample. Tests based on AFS using DNA sequence data rely on a few related statistics, all of which are comparisons between estimates of the population genetic parameter  $\theta = 4N\mu$ . The statistics are calculated as the difference of two such estimates, normalized by the expected variance of the difference under a neutral model, so that values below -2 or greater than 2 roughly exceed the 95% confidence limits about the neutral expectation of 0. However, the actual mean may frequently deviate from 0 (Thornton 2005; Wares 2010). Simonsen *et al.* (1995) compared three measures and found Tajima's  $D$  to have the most statistical power:

Tajima's  $D$ : Normalized difference between  $\pi$  and  $S$ , the number of segregating sites (Tajima 1989).

Fu and Li's  $D^*$ : Normalized difference between  $S$  and the number of singletons  $\eta$  (alleles observed only once in a sample; Fu and Li 1993).

$F^*$ : Normalized difference between  $\pi$  and  $\eta$  (Fu and Li 1993).

**Background selection:** Ongoing selection against deleterious mutations that can result in the loss of linked neutral variation (Charlesworth *et al.* 1993).

**Balancing selection:** Here we define balancing selection broadly as the class of selective forces that maintain polymorphism over time. This can include, for example, frequency-dependent selection or heterozygote advantage (Charlesworth 2006).

**Coalescent theory:** A theoretical framework for understanding genetic variation based on the retrospective pattern of shared ancestry among alleles in a sample (Wakeley 2009).

**Divergent selection:** Positive selection acting differentially between separate populations.

**$dN/dS$ :** Ratio of nonsynonymous (amino acid-changing) to synonymous substitutions in a nucleotide sequence. Testing for selection based on this ratio typically uses aligned sequence data among populations or taxa and can detect selection over long timescales, although it requires multiple amino acid substitutions (i.e., recurrent selective sweeps).

**$F_{ST}$ :** A statistic describing the partitioning of allelic variance within versus among populations;  $F_{ST}$  ranges from 0 (no population differentiation) to 1 (complete population differentiation). There are multiple ways of calculating  $F_{ST}$  that can occasionally have substantial effects on its value but rarely its relative magnitude among loci (Charlesworth 1998; Holsinger and Weir 2009). Commonly used population genomic tests for selection based on identifying outliers in  $F_{ST}$  are as follows:

LOSITAN (Antao *et al.* 2008) computer software implements the method of Beaumont and Nichols (1996) to identify  $F_{ST}$  outliers based on heterozygosity, which affects the predicted neutral distribution of  $F_{ST}$ .

ARLEQUIN (Excoffier *et al.* 2009) software performs the same analysis, accounting for hierarchical population structure.

BAYESFST (Beaumont and Balding 2004) assesses the significance of a locus-specific parameter that indicates selection in a model of  $F_{ST}$ .

BAYESCAN (Foll and Gaggiotti 2008) modifies the approach of Beaumont and Balding (2004) to estimate the posterior probability of a locus being subject to selection.

DETSEL (Vitalis *et al.* 2003) uses coalescent simulations in a simple two-population model to identify  $F_{ST}$  outliers.

**Genetic draft:** The loss of genetic diversity and changes in AFS at loci linked to a selected locus during a selective sweep (Gillespie 2000).

**HKA test:** A test of the neutral prediction for the relationship between within-population diversity and among-population divergence (Hudson *et al.* 1987).

**Linkage disequilibrium (LD):** The correlation between alleles across loci. Traditionally, LD has been calculated as a function of a pair of loci, regardless of their physical position (Slackin 2008). This aspect of LD can be partitioned among populations in the statistic  $Z_g$  as a test of selection (Storz and Kelly 2008). Genome scans for selection also apply several of the following statistics that describe the decay of LD as a function of physical distance, also known as haplotype structure:

Extended haplotype homozygosity (EHH) measures the probability that any two randomly chosen haplotypes are identical over a given distance from a focal site (Sabeti *et al.* 2002).

Integrated haplotype score (iHS) integrates the area under the EHH curve (Voight *et al.* 2006). Huff *et al.* (2010) found this measure to have greater statistical power and to be more robust to complex demographics than two related alternatives.

Cross-population extended haplotype homozygosity (XP-EHH) compares EHH between two populations to test for interpopulation differences in the extent of LD (Sabeti *et al.* 2007).

**$\pi$ :** A measure of nucleotide diversity, calculated as the proportion of pairwise differences in a sample;  $\pi$  can be estimated either within or between populations and is directly used in some calculations of  $F_{ST}$  (Charlesworth 1998).

**Positive (directional) selection:** Selection in which one or a class of alleles is favored.

**Selective sweeps:** The increase in frequency of one or a class of alleles favored by selection. Hard sweeps result from selection on a single allele, typically a new mutation that is favored immediately on its appearance in a population. Soft sweeps are selection on standing genetic variation or on variants supplied by recurrent mutation or migration during the selective phase, so that a number of different alleles are collectively favored and increase in frequency. These alleles are typically considered to be neutral or even deleterious before a shift in selective regime (Hermisson and Pennings 2005).