

Natural Selection

Ryan D. Hernandez

The Effect of Positive Selection

Adaptive

Neutral

Nearly Neutral

Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



The Effect of Positive Selection

Adaptive

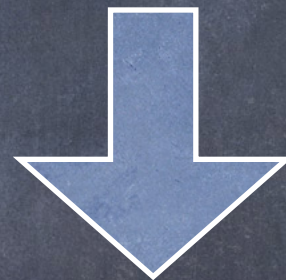
Neutral

Nearly Neutral

Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



Site-Frequency Spectrum

■ SNM

■ AfAm (Human)

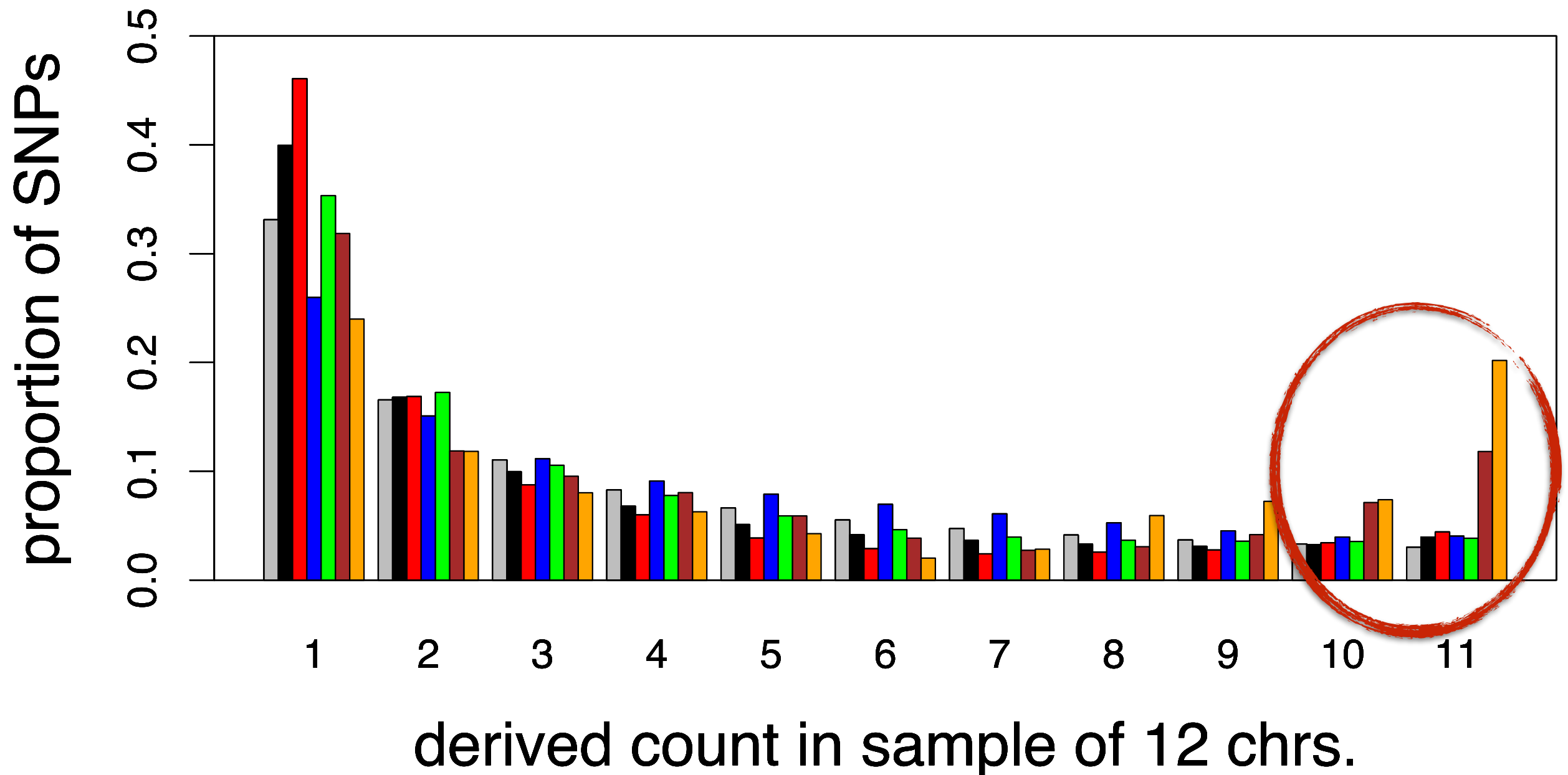
■ Ch (RheMac)

■ In (RheMac)

■ Rufi (rice)

■ Indica (rice)

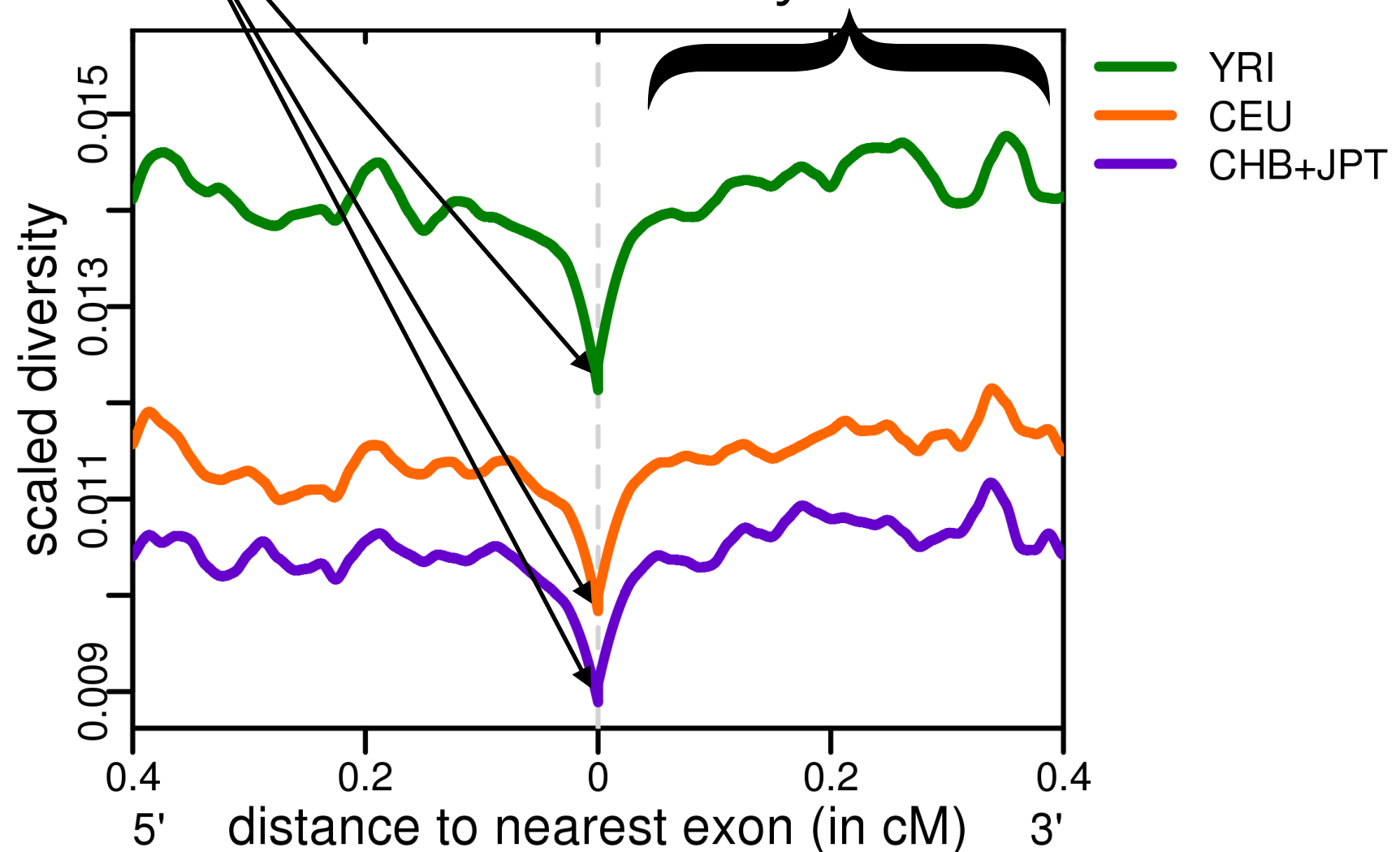
■ Japonica (rice)



Coding regions tend to have the lowest levels of diversity in the genome

Diversity at selected loci

Diversity at linked sites



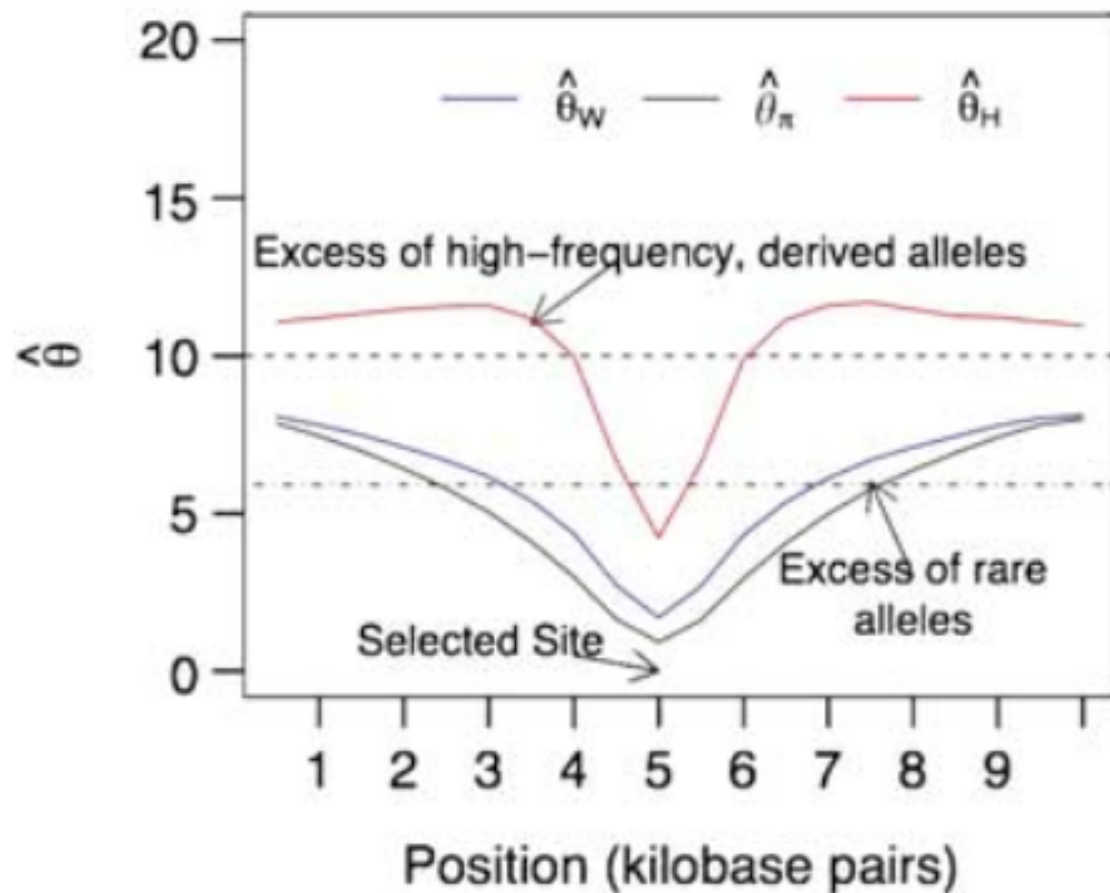
What are the predominant evolutionary forces driving human genomes?!

Eyre-Walker & Keightley (2009) ~**40%** of amino acid substitutions were **advantageous**

Boyko et al (2008) **10-20%** of amino acid substitutions were **advantageous**

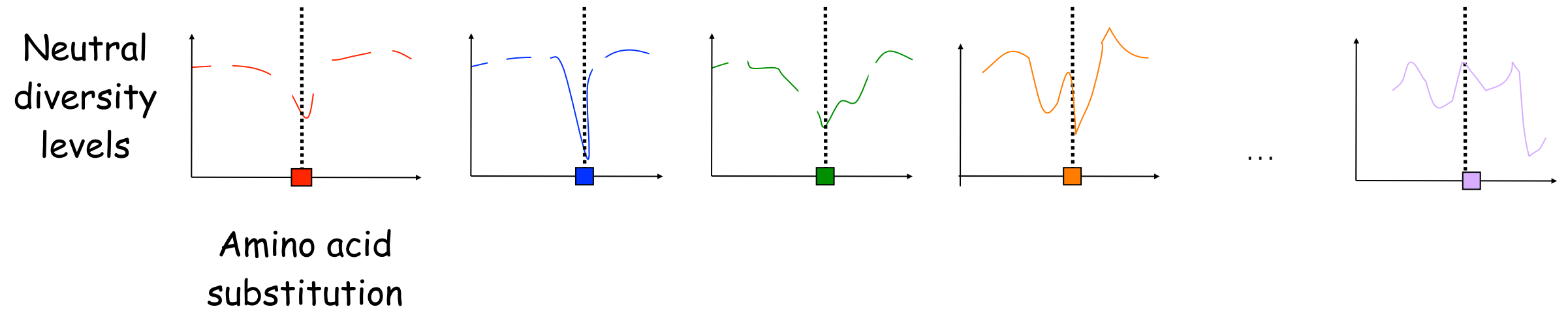
Williamson et al (2007) **10%** of the genome affected by **selective sweeps**

Diversity levels around a selective sweep

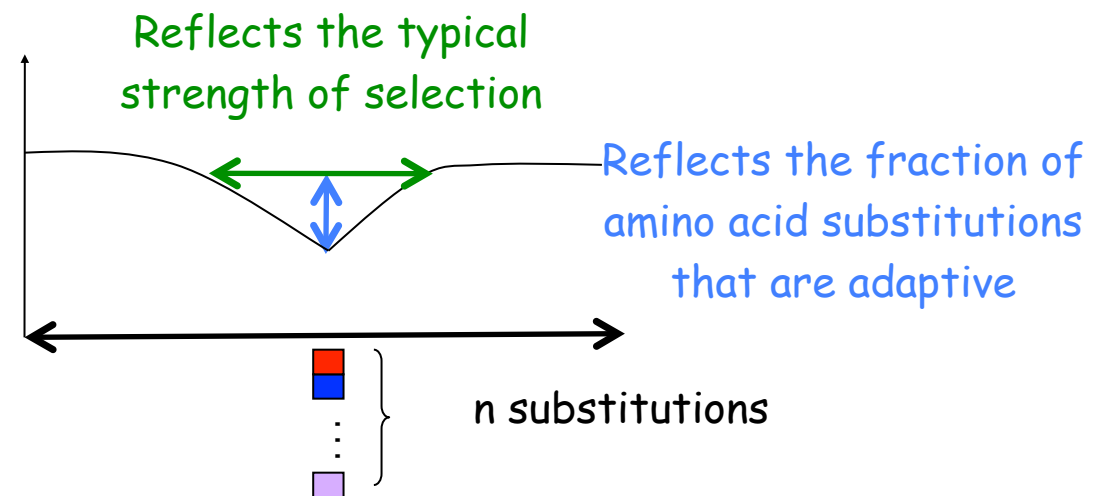


Thornton et al (2007): Simulation of patterns of **neutral** diversity around a **selective sweep**

The footprint of adaptive amino acid substitutions

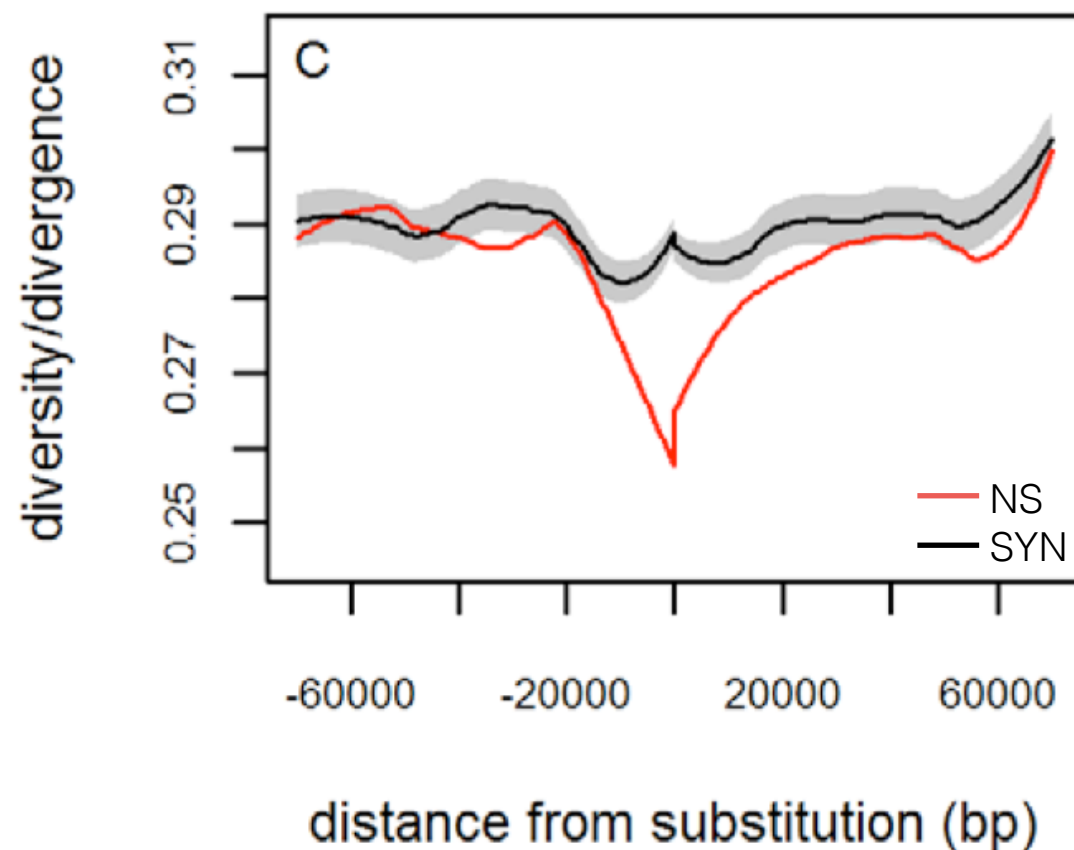


- Goal: compare the pattern around **amino acid substitutions** to the pattern around **synonymous substitutions**.



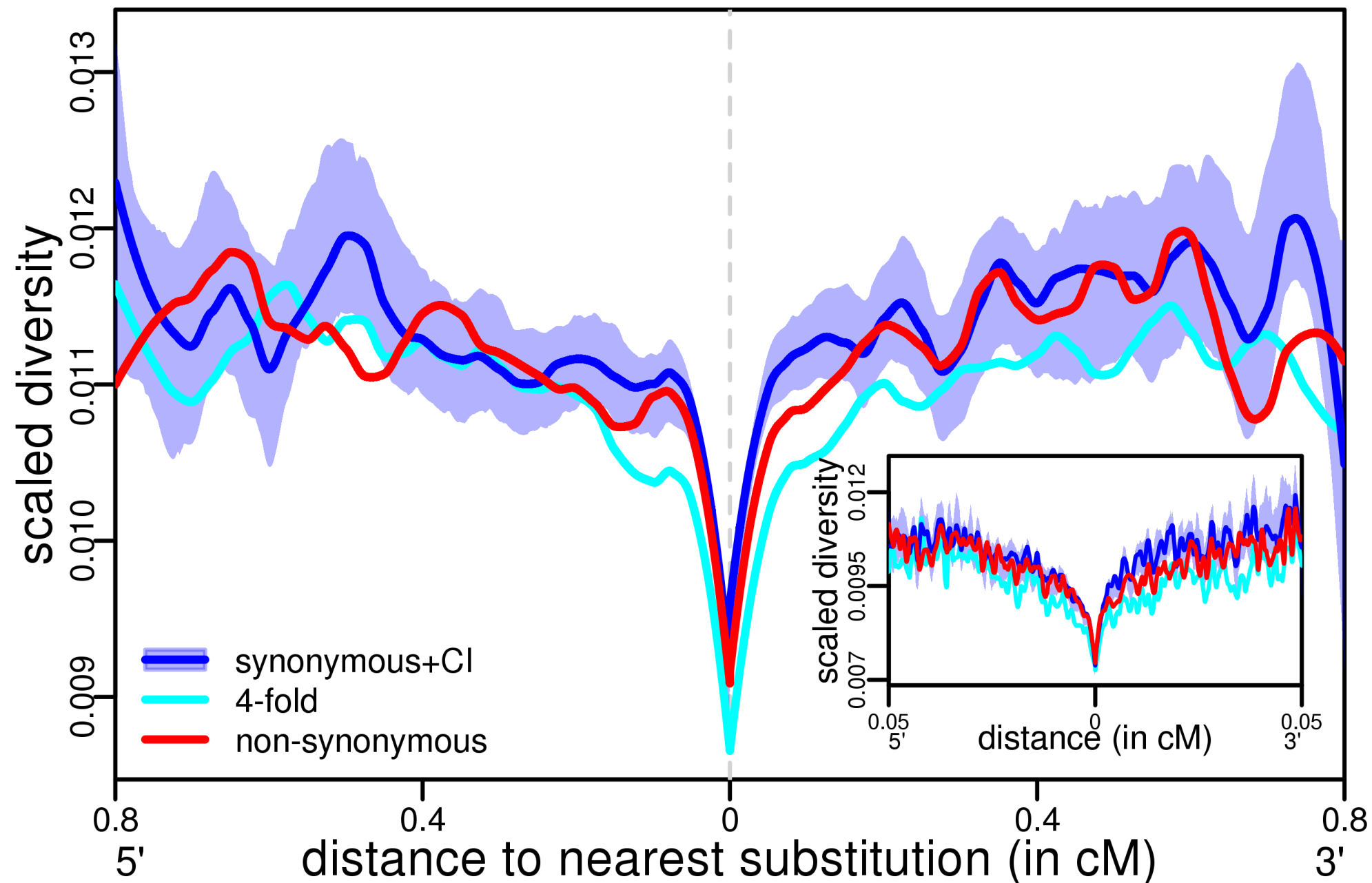
Other organisms...

Drosophila



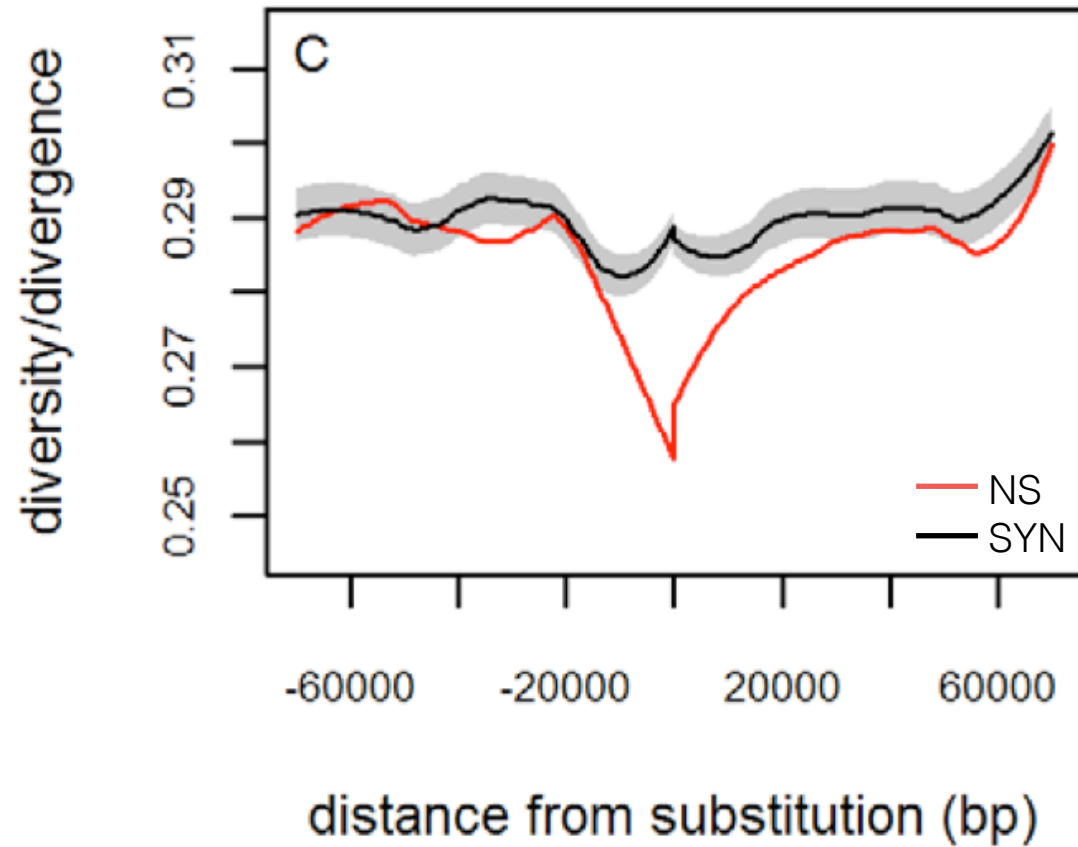
Sattath et al (2011) estimate
~13% of amino acid
substitutions were adaptive.

Observed Patterns of Diversity Around Human Substitutions

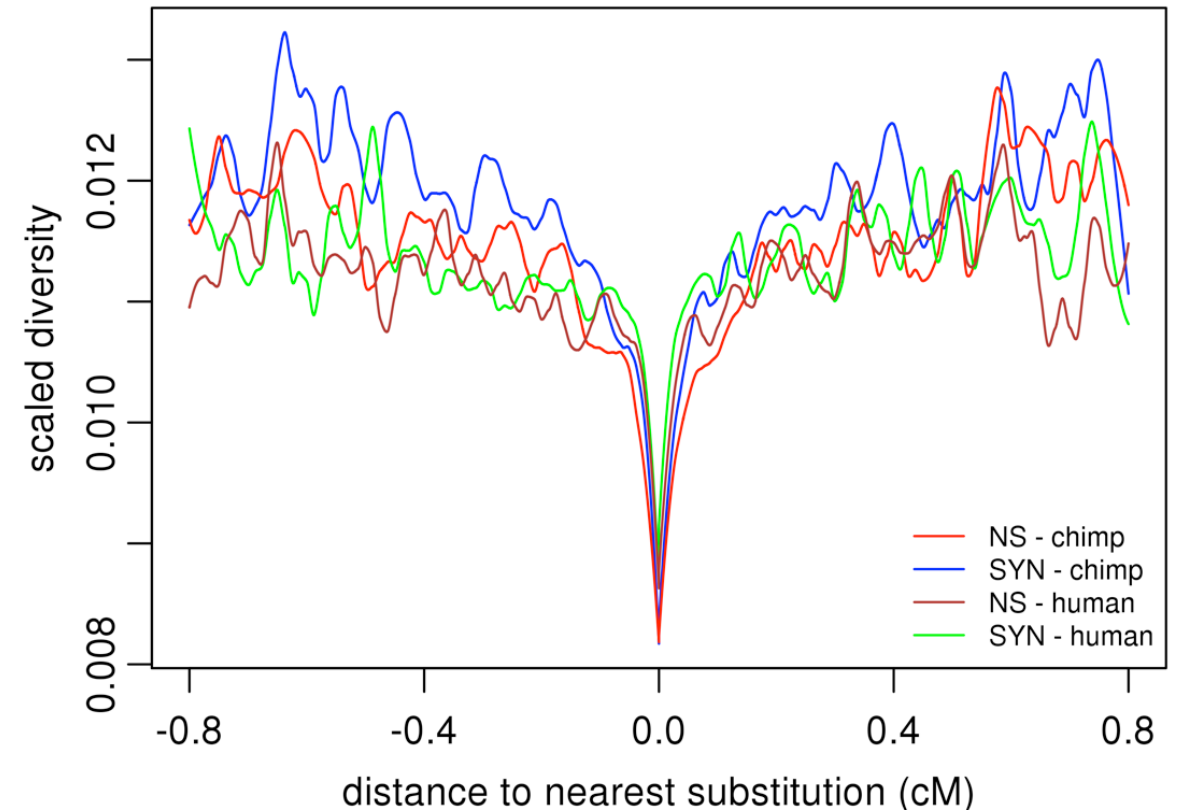


Other organisms...

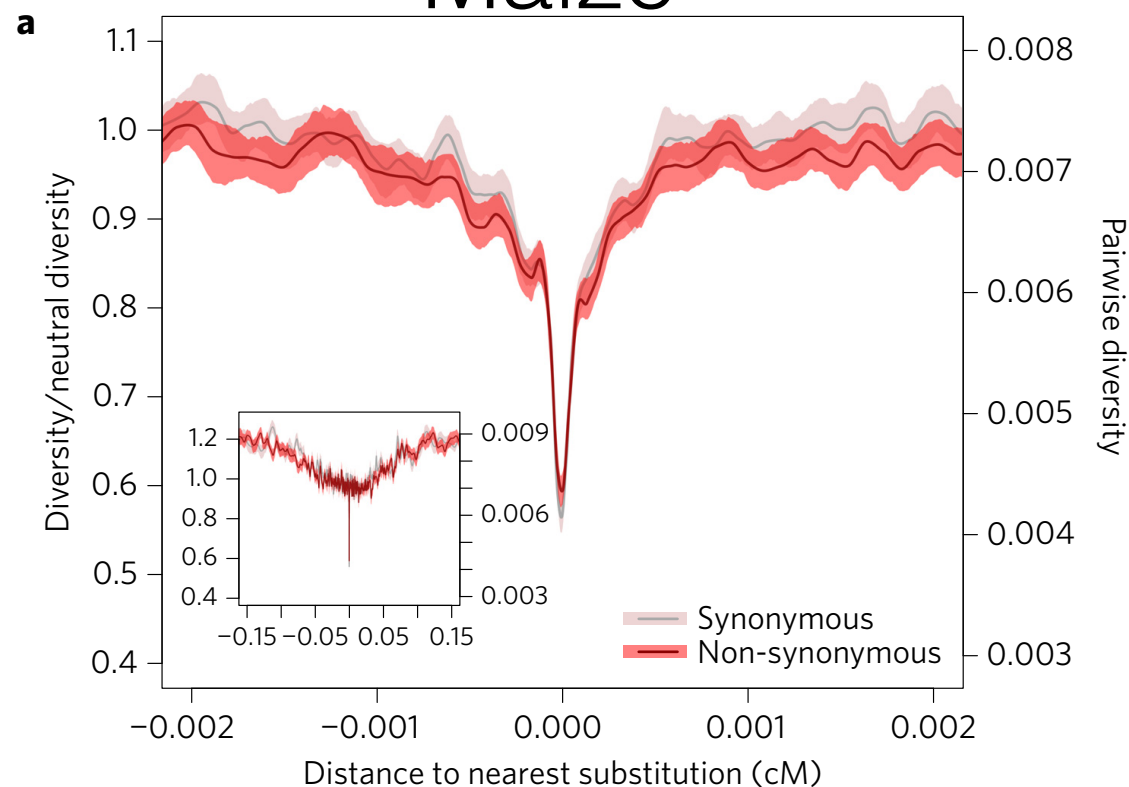
Drosophila



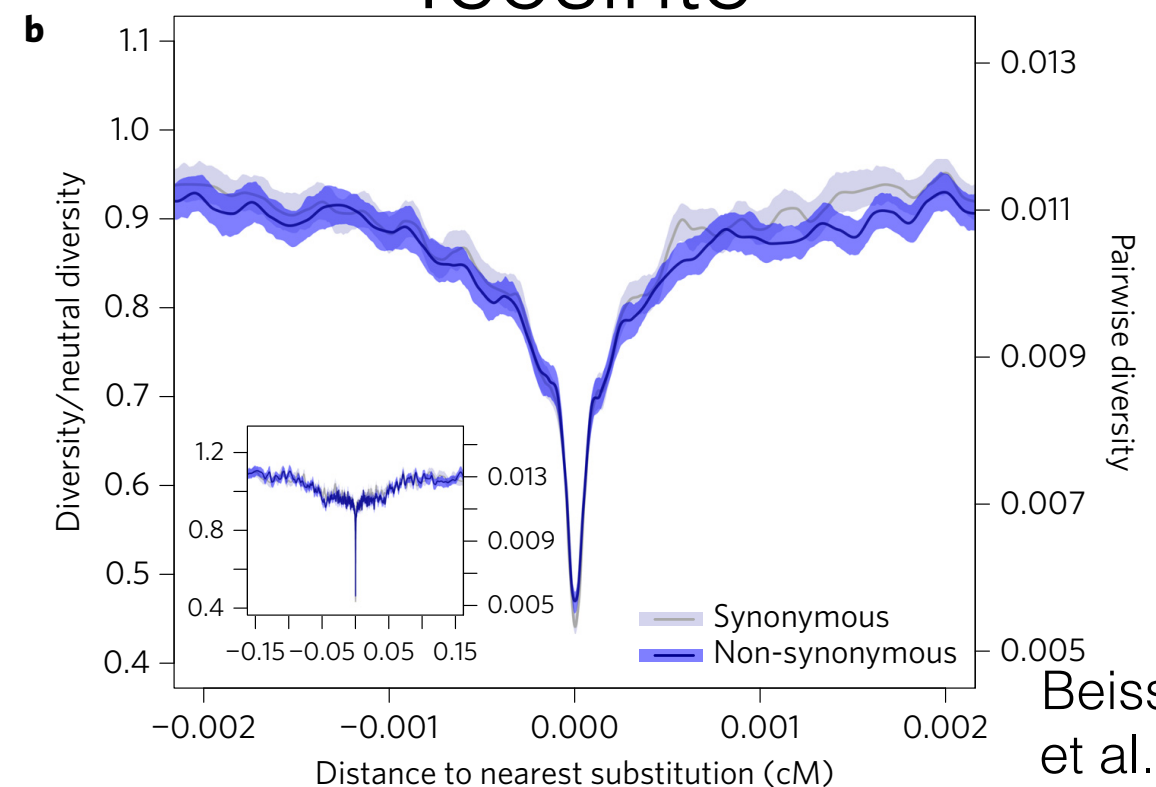
Chimpanzee



Maize



Teosinte



Beissinger ,
et al. (2016)

The Effect of Negative Selection

Adaptive

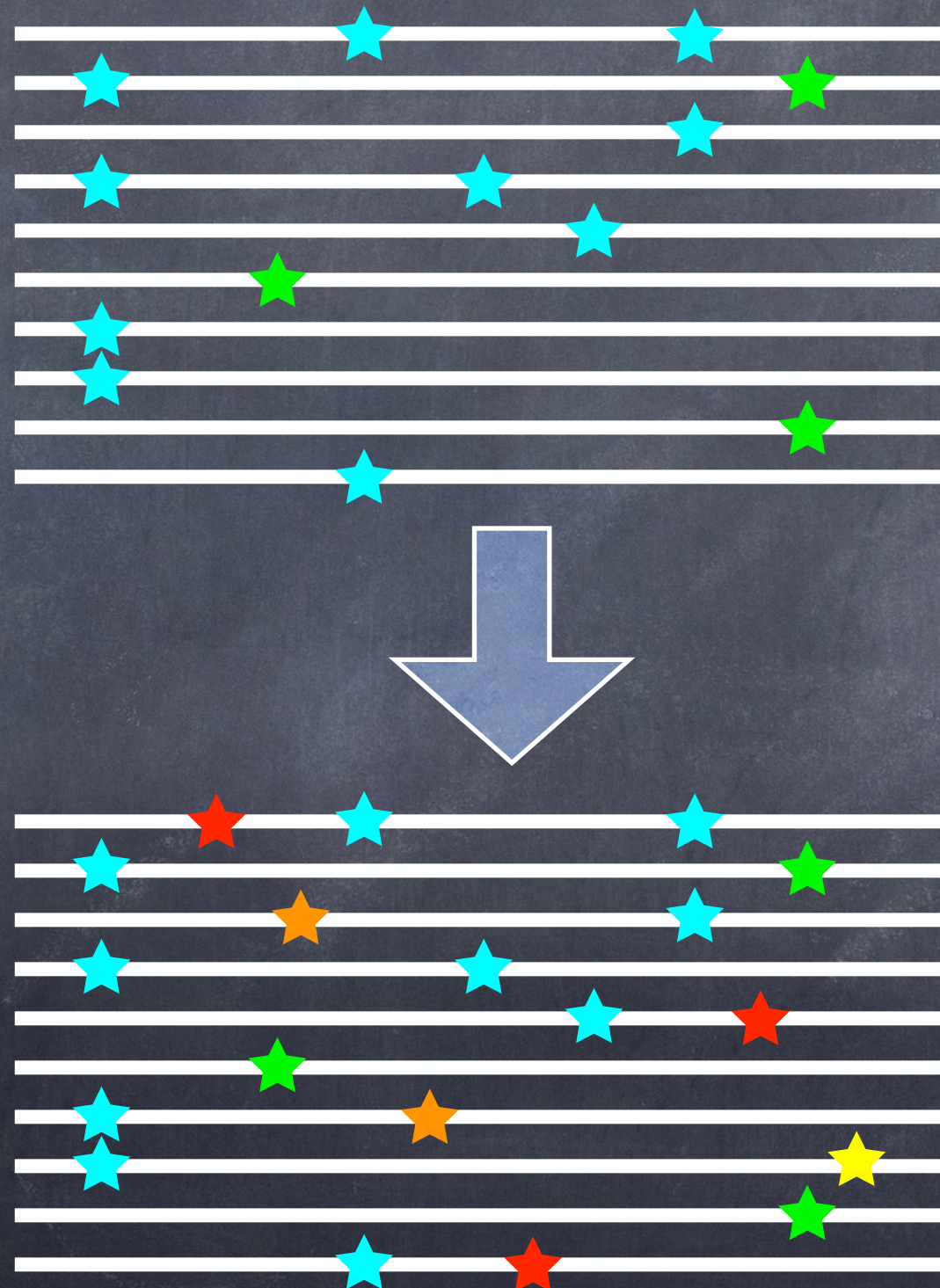
Neutral

Nearly Neutral

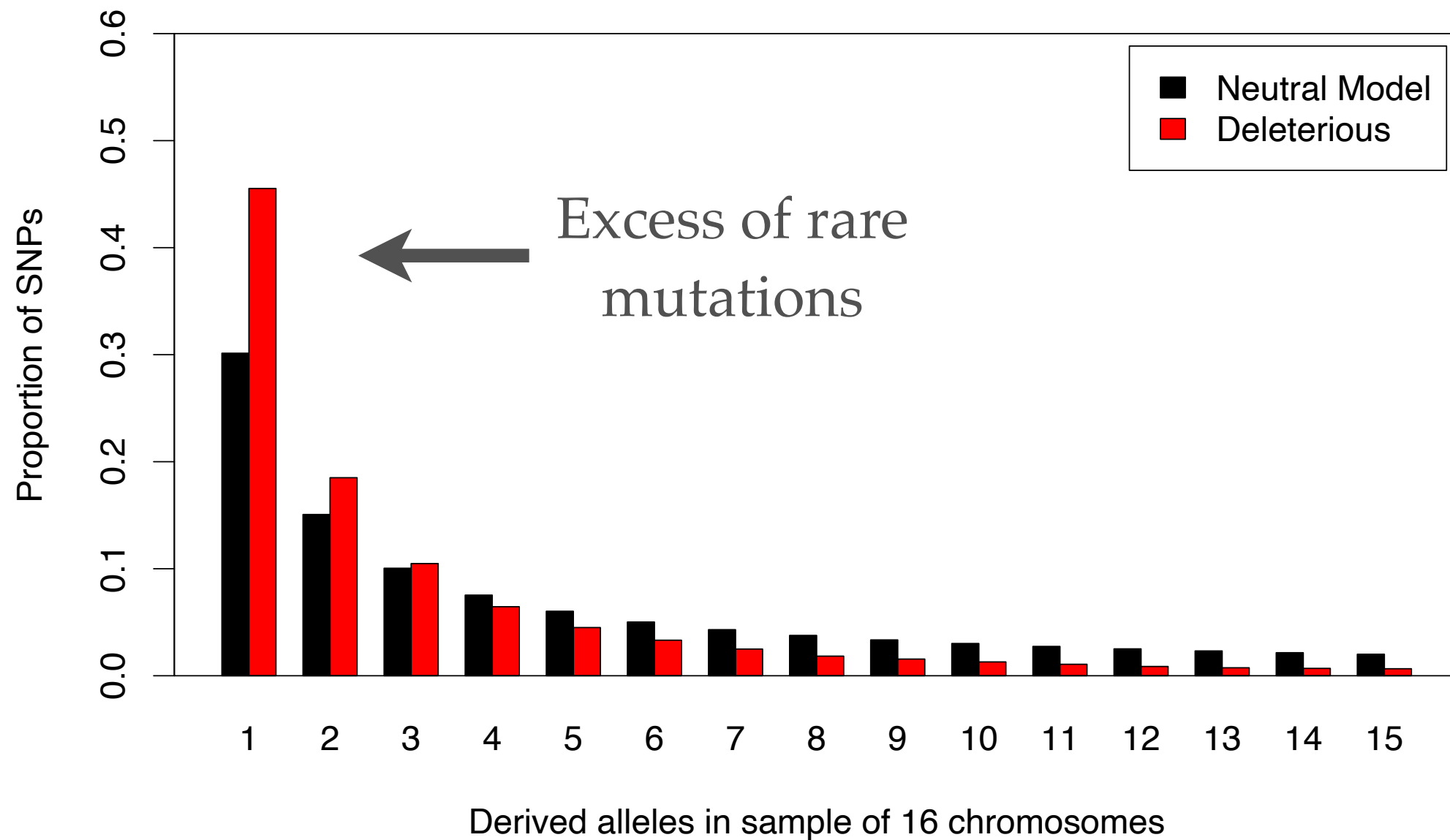
Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



Site-Frequency Spectrum



The Effect of Negative Selection

Consequences:

- Some proportion of chromosomes eliminated each generation
 - ➡ Decreased effective population size ($f_0 N_e$)
 - ➡ Decreased neutral variation ($f_0 \pi$)
- While neutral variation can be lost, some neutral mutations may increase in frequency

Background
selection

Background selection (BGS)

- Definition: The reduction of diversity at a **neutral** locus due to the effects of purifying selection at linked deleterious loci
- Can estimate the effect of BGS by comparing **observed** diversity at neutral sites compared to the level of diversity you would **expect** under neutrality!
- π/π_0

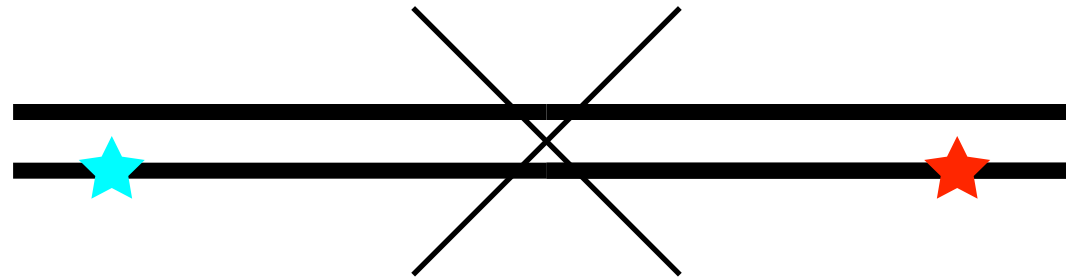
Earlier Theoretical Work

Hudson & Kaplan (1995)

$$f_0 = \exp \left(-\frac{U}{s + R} \right)$$

- U = deleterious mutation rate
- s = selection coefficient
- R = recombination rate

Effect of Recombination



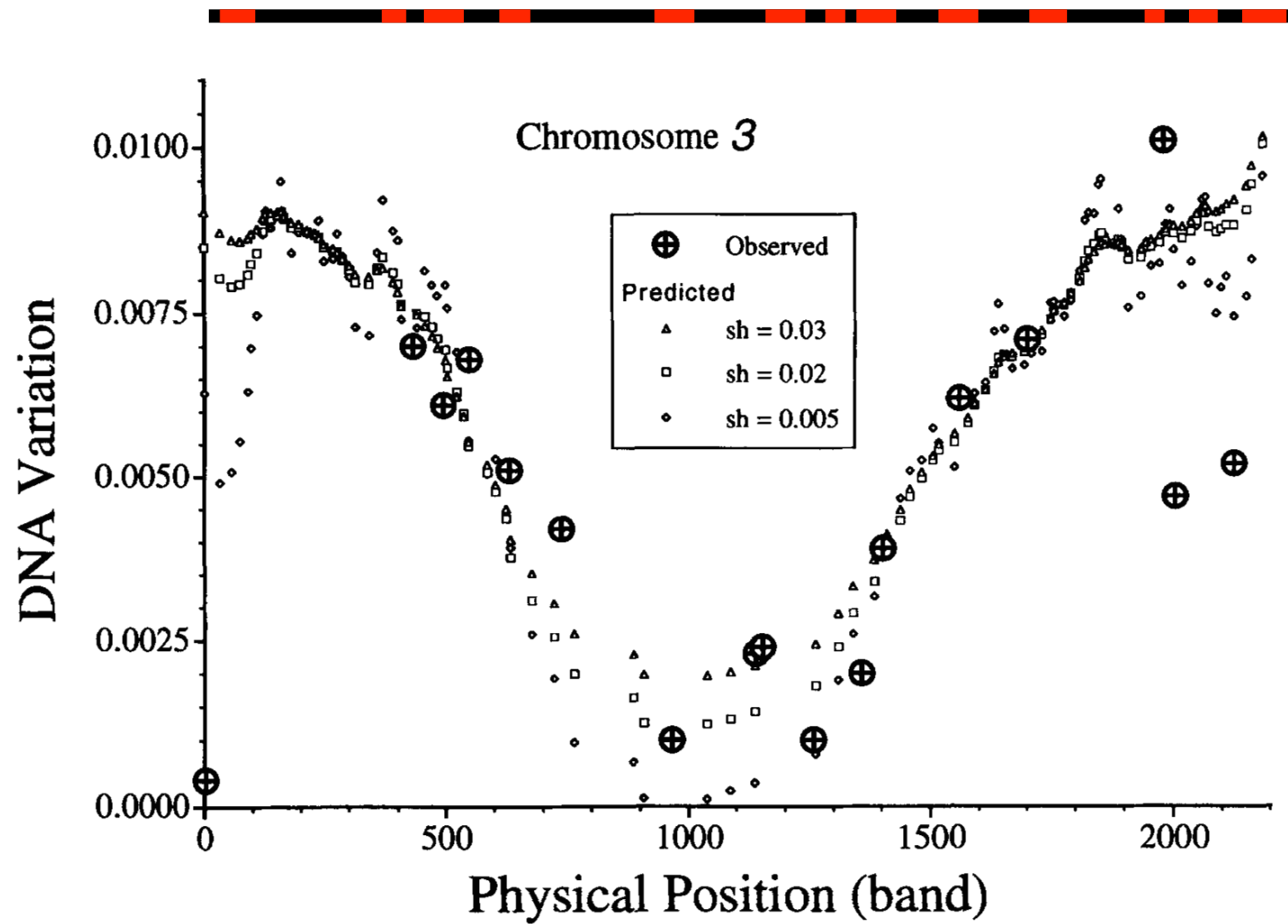
With recombination, neutral mutations can escape the grip of deleterious mutations.

Multiple Targets of Deleterious Mutations



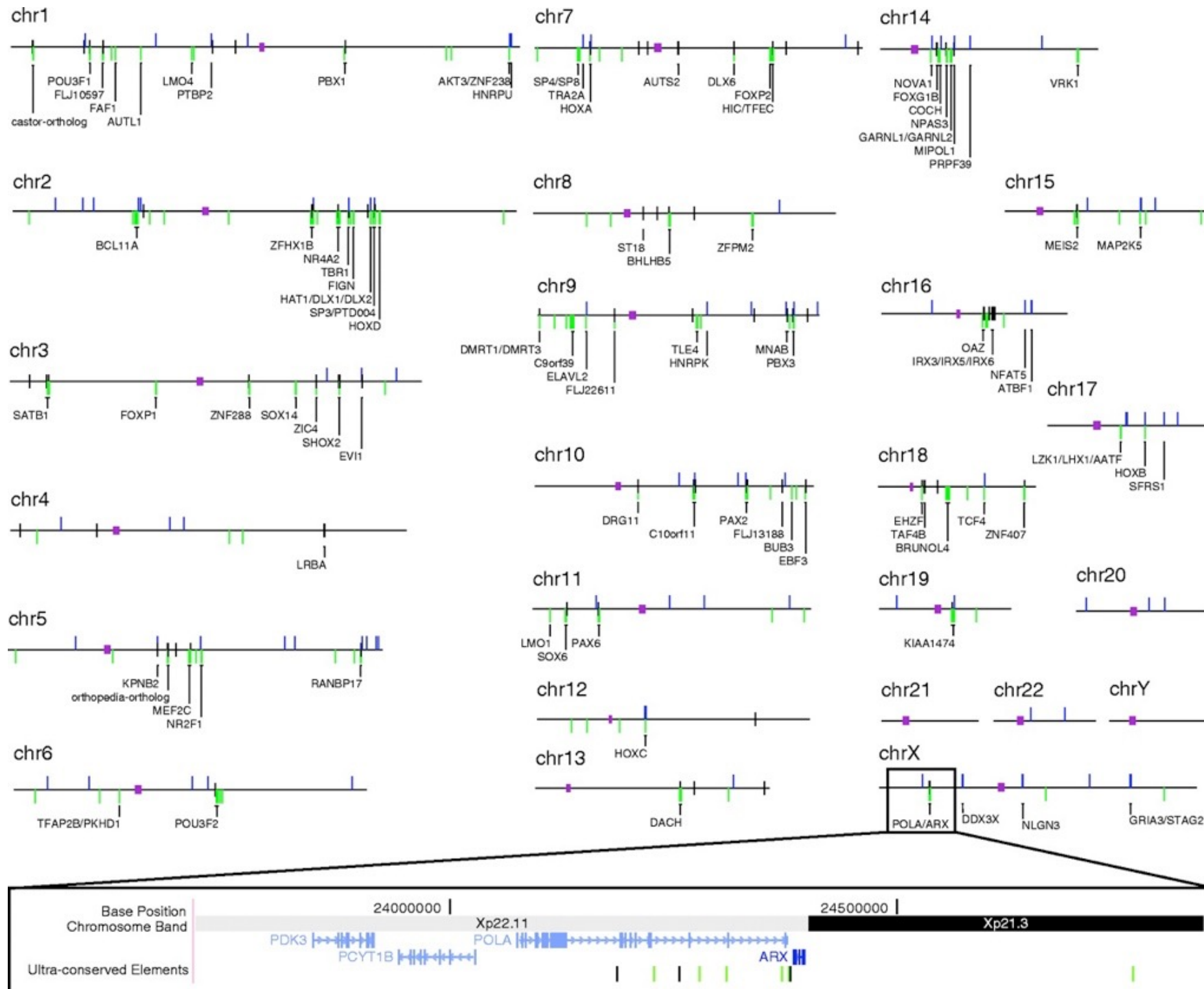
Consider a chromosome composed of neutral loci and deleterious loci

Drosophila



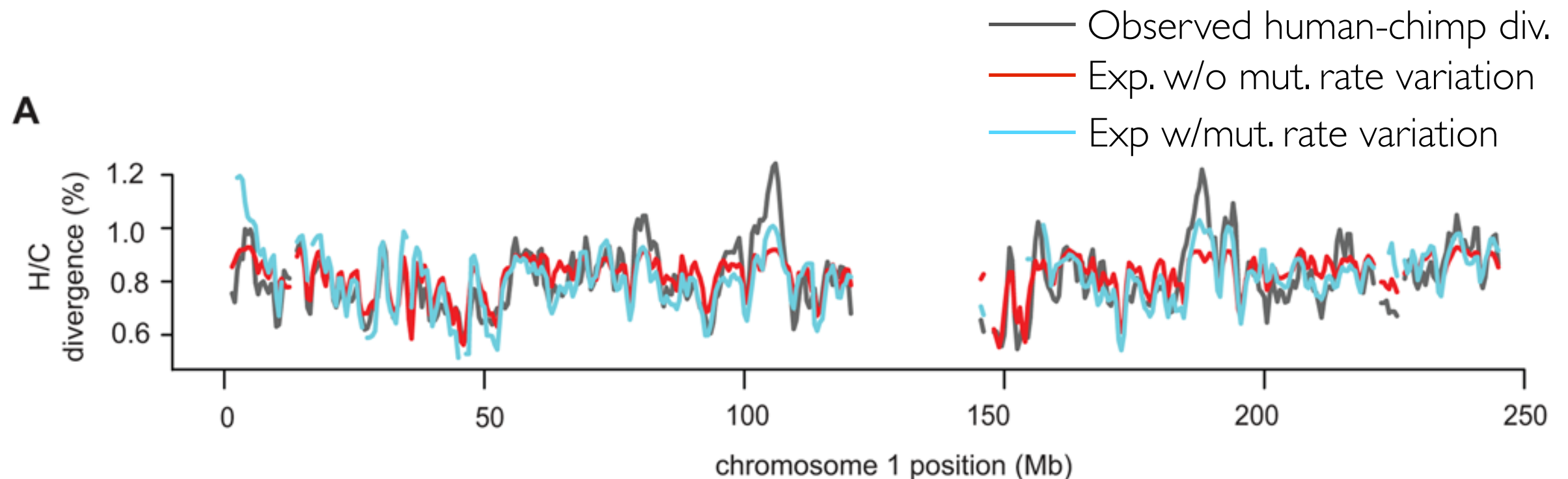
Hudson & Kaplan (1995)

Distribution of Ultraconserved Elements in the Human Genome

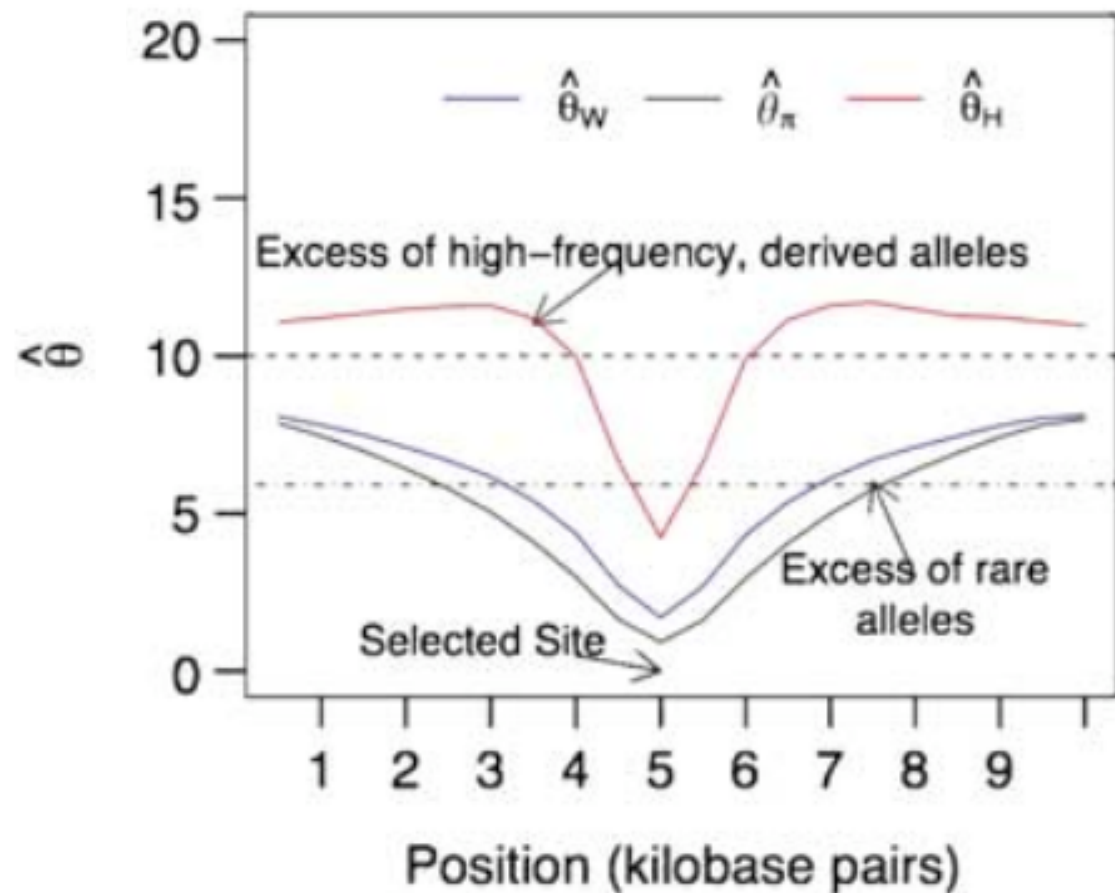


Background Selection

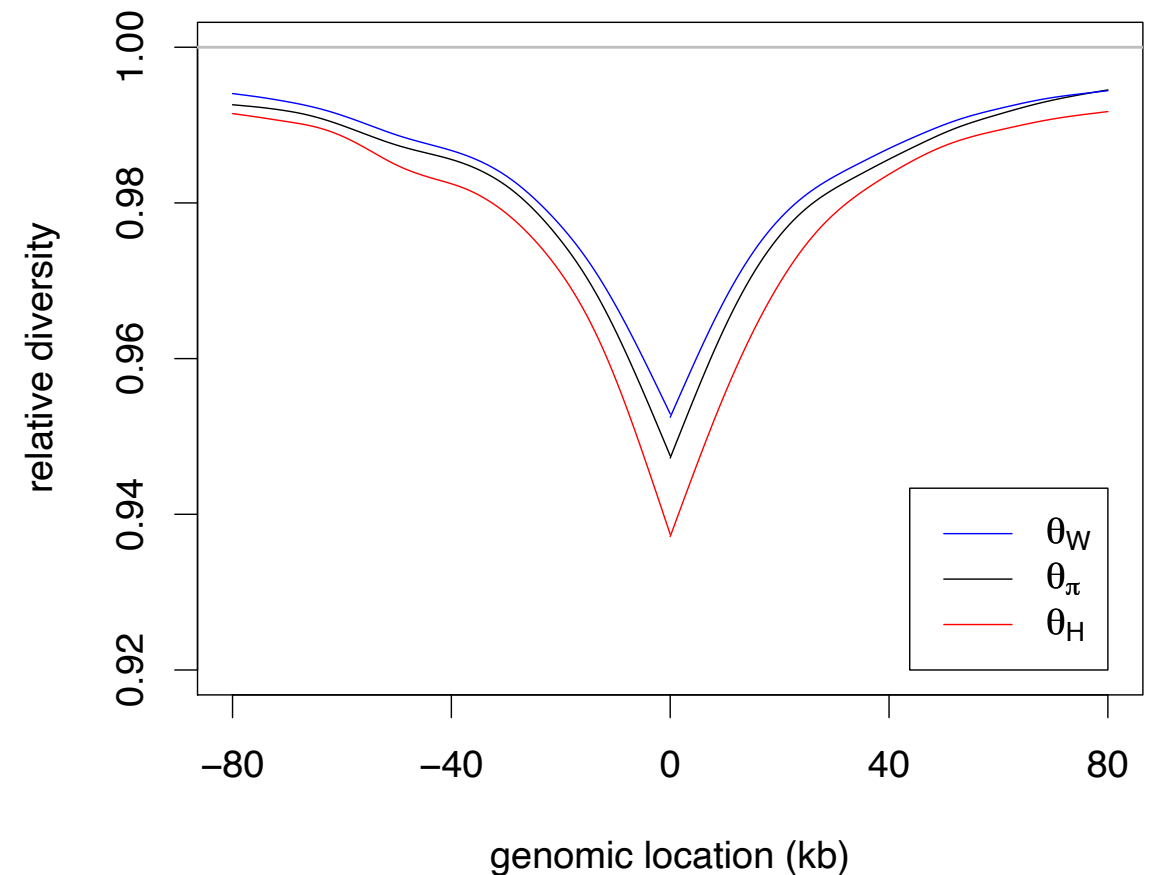
- The effects of the continual removal of deleterious mutations by natural selection on variability at linked sites.



Diversity levels around sites subject to natural selection

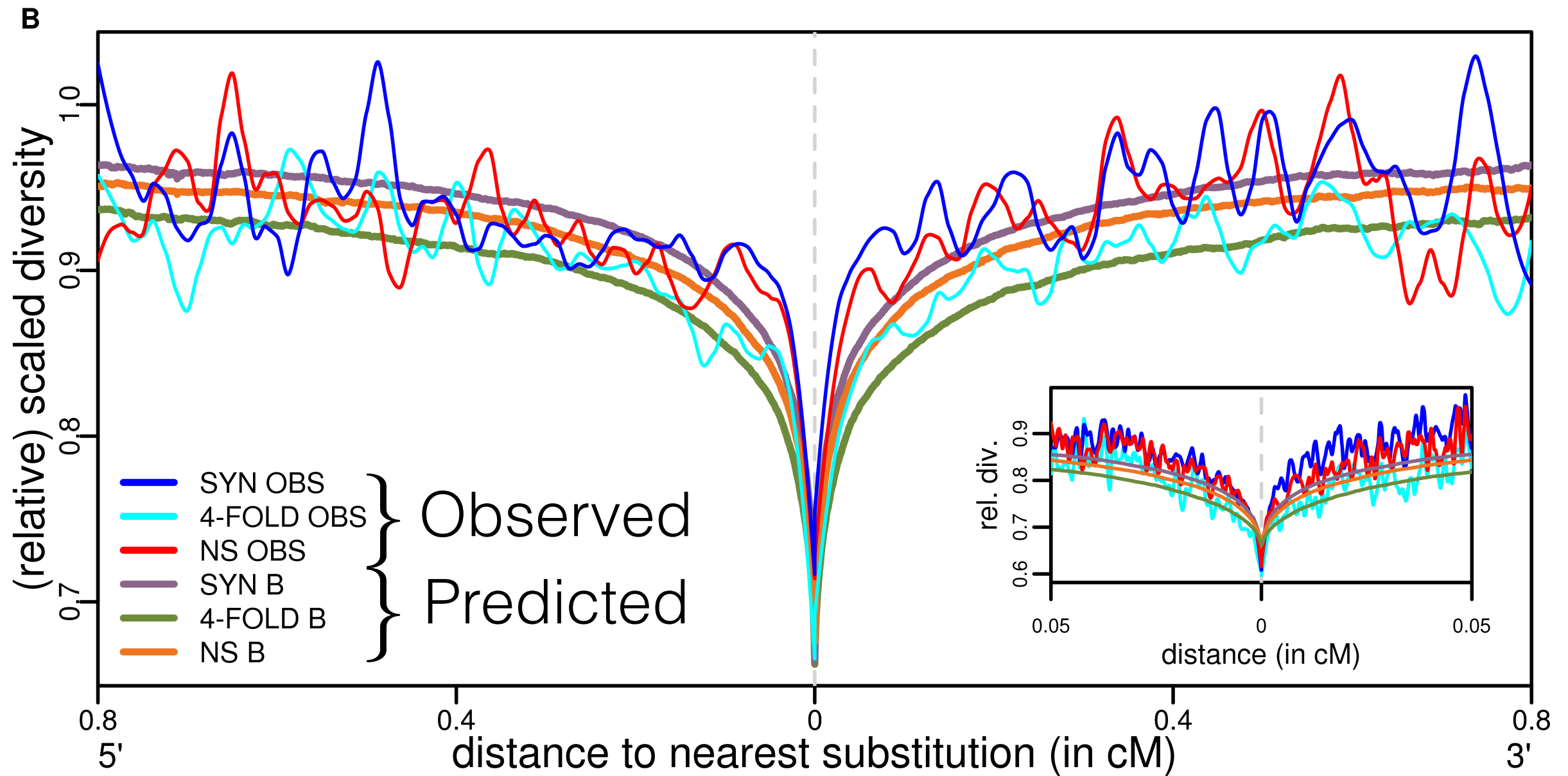


Thornton et al (2007): Simulation of patterns of **neutral** diversity around a **selective sweep**

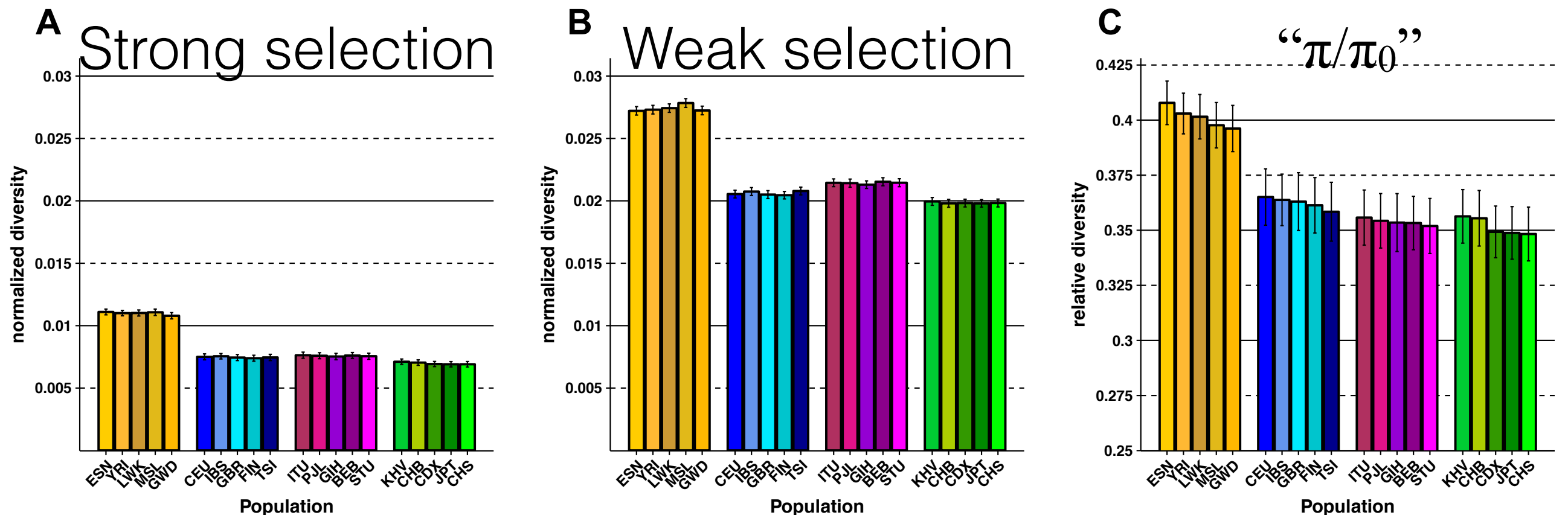


Simulation of patterns of **neutral** diversity around a 700bp **deleterious locus** with $\gamma = -5$.

Modeling the data



BGS Features



- Neutral sites in 1000 Genomes Project data: 20 non-admixed populations
- The strength of background selection varies across populations!
 - Stronger effects in bottlenecked Out-Of-Africa populations

BGS Features

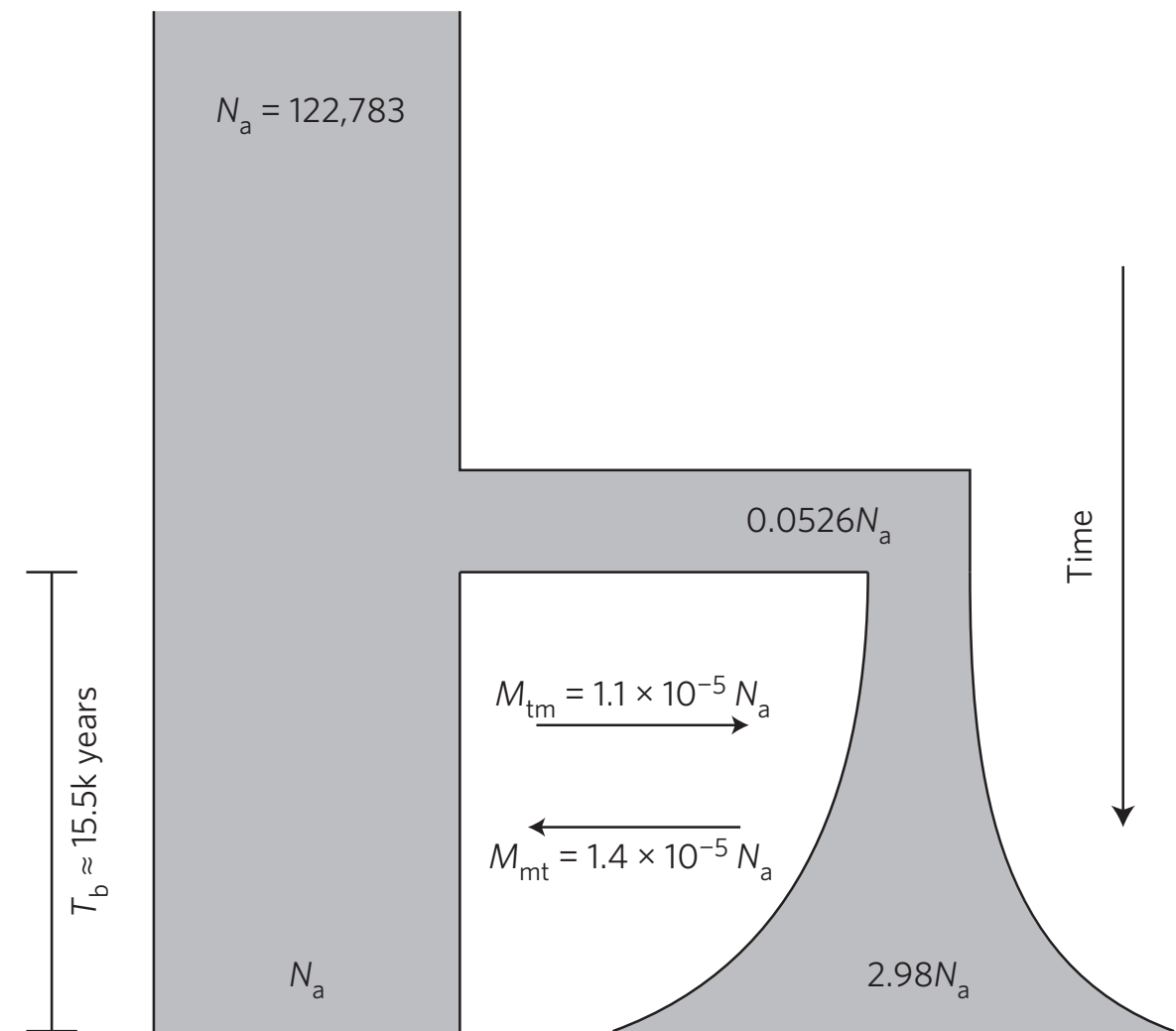
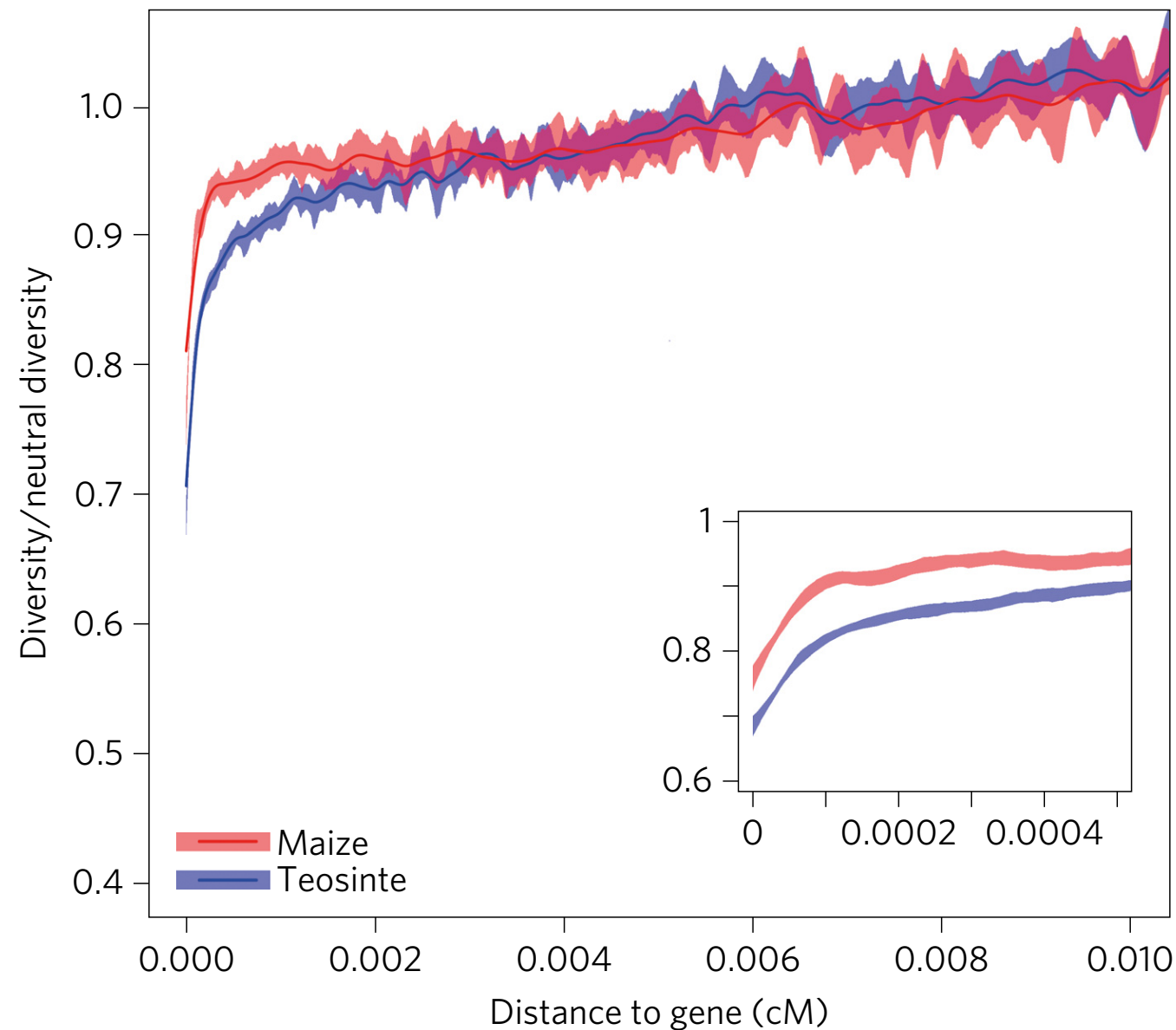
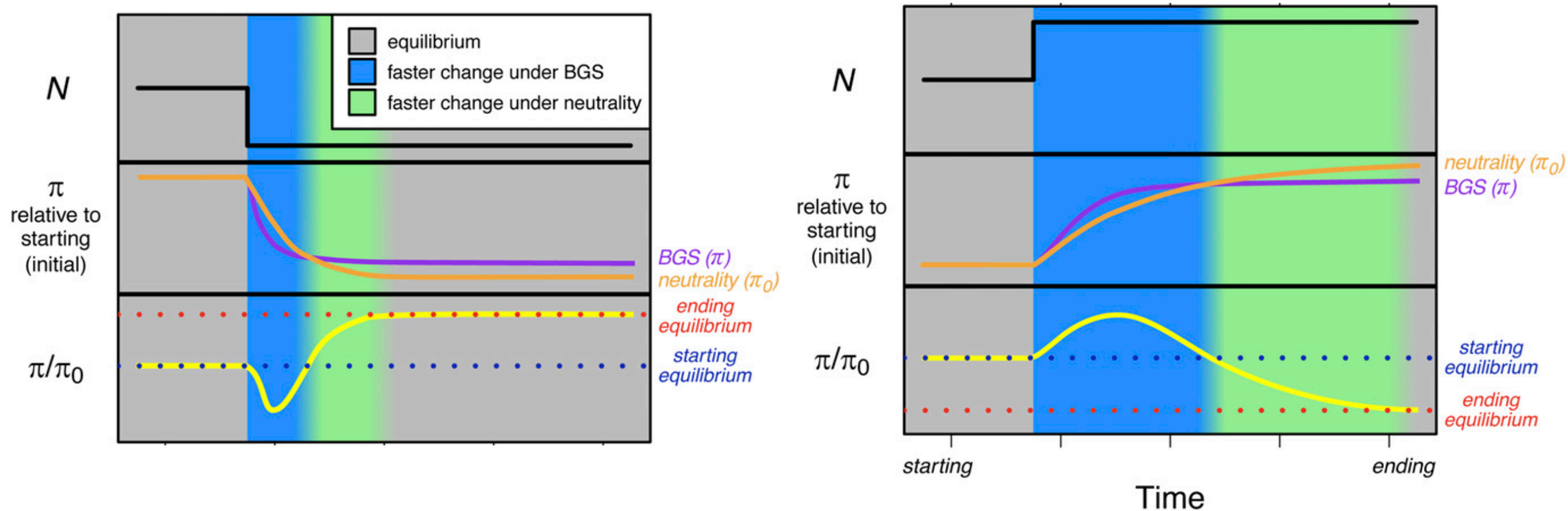


Figure 2 | Estimated demographic history of maize and teosinte.

Parameter estimates for a basic bottleneck model of maize domestication. See Methods for details.

- Strength of BGS varies between Maize and Teosinte
- Stronger in ancestral Teosinte population!

Demographic Models Matter!



- For **both** contraction and expansion models:
 - π/π_0 can be greater than or less than the ancestral population depending on time!

Background Selection & Disease?

Background selection drives patterns of genetic variation.

- But does it matter?
- Does it have implications for studying complex traits?

To find out, we looked at the NHGRI GWAS database:

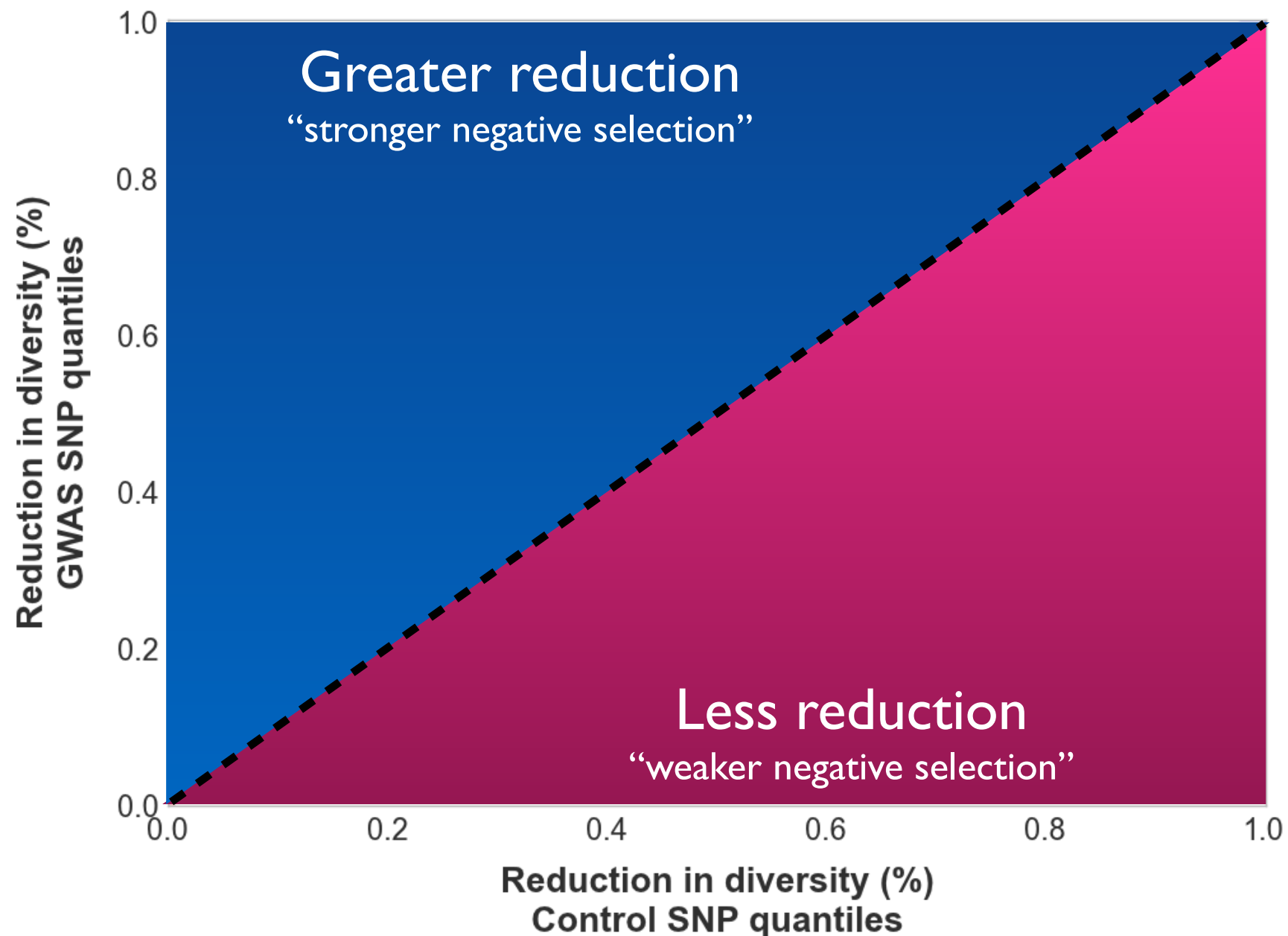
www.genome.gov/gwastudies/

Published Genome-Wide Associations through 07/2020

Published GWA at $p \leq 5 \times 10^{-8}$

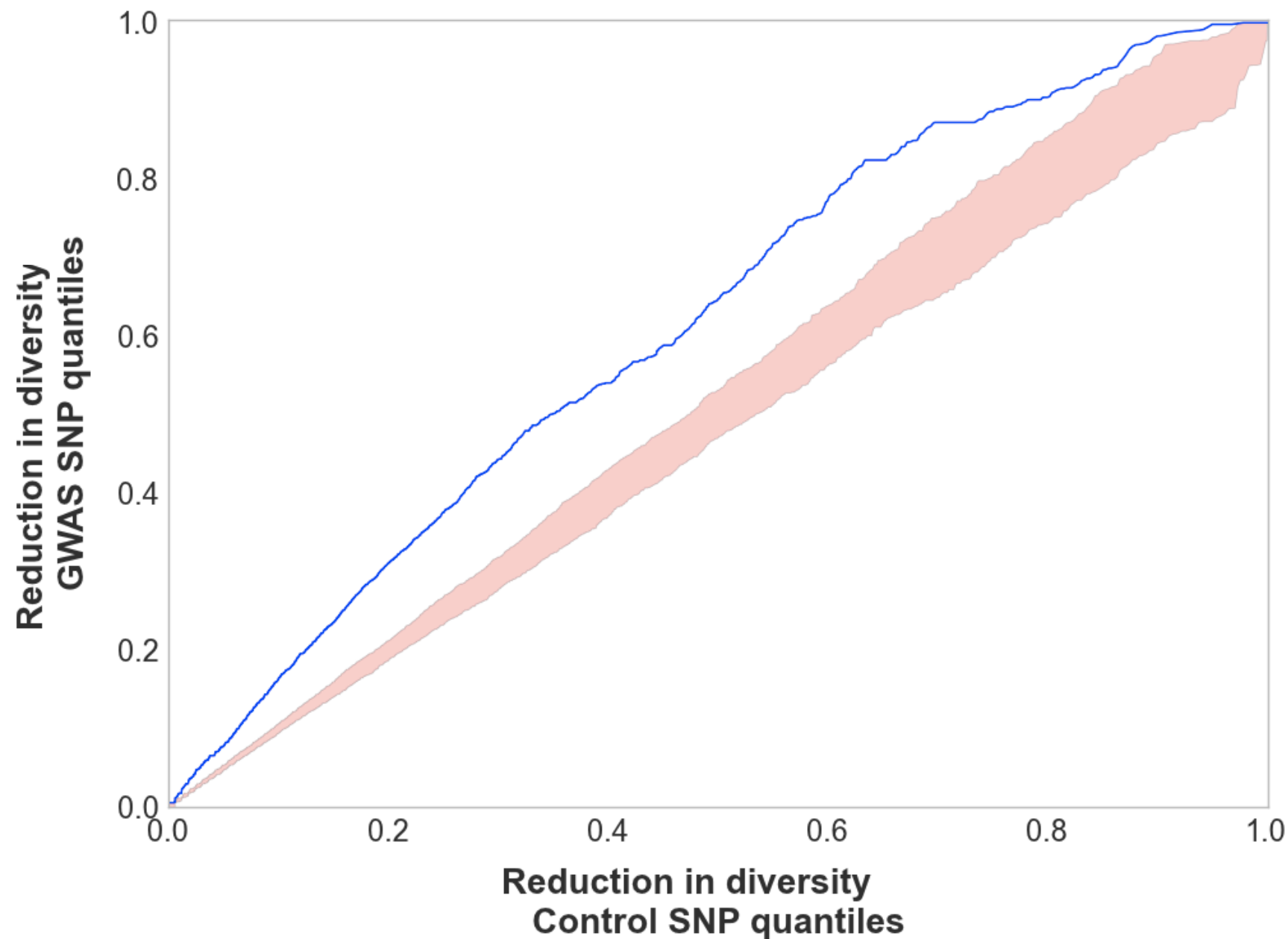


Effects of Linked Selection



- QQ-plot of the reduction in diversity around GWAS hits compared to background.

Effects of Linked Selection



- Greater reduction in diversity around GWAS hits indicates a strong, local burden of negative selection.

Genetic Load

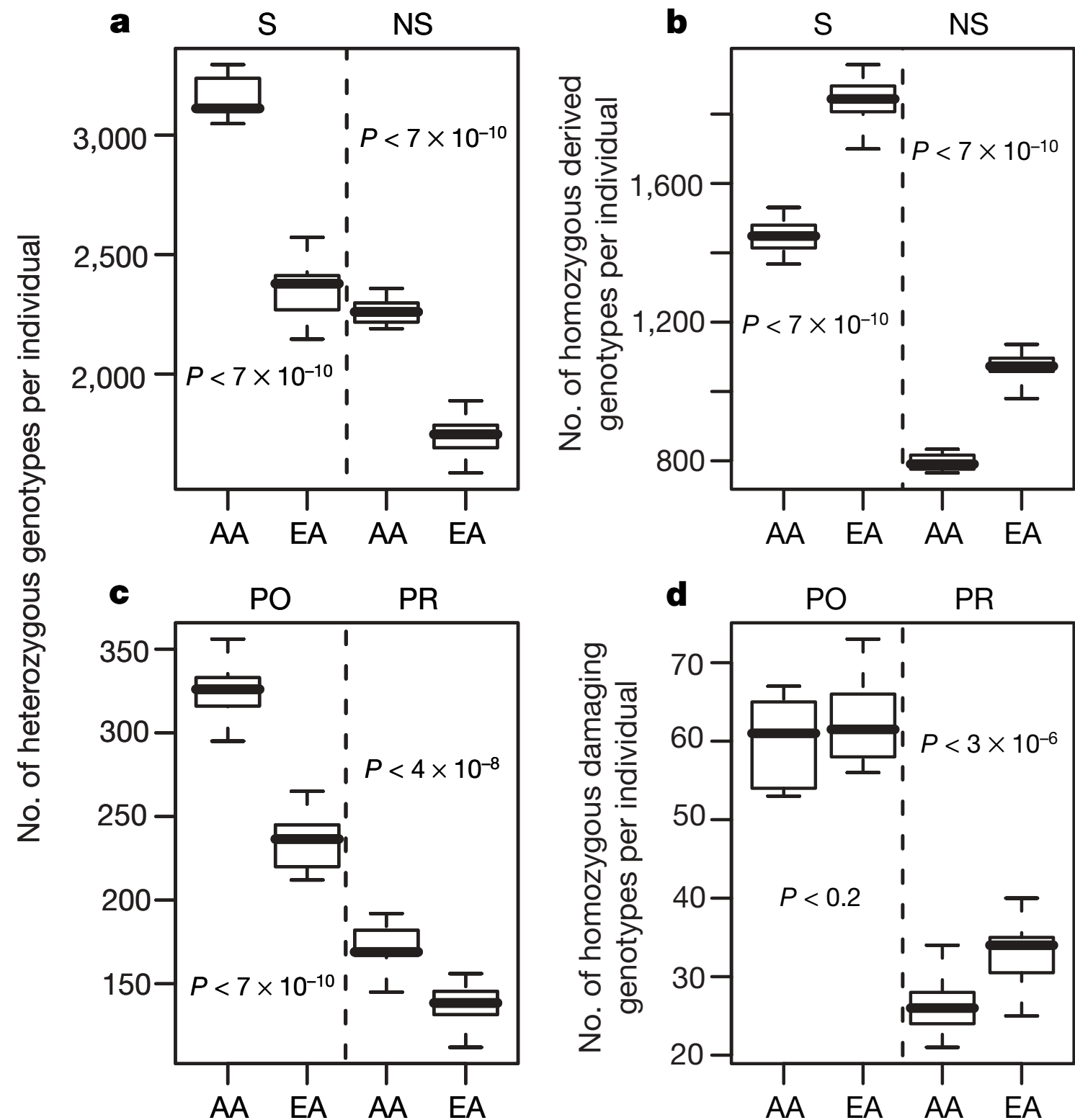
- Genetic load is the reduction in population mean fitness due to deleterious mutations compared to a (hypothetical) mutation-free population.
- Load is the outcome of the evolutionary process of a population.
- But, unlike other features of genetic variation, it cannot be directly observed.
- Must be indirectly inferred.

Inferring Genetic Load

- Empirical counting approaches:
 - Under an additive model, the number of derived deleterious alleles will be proportional to genetic load
 - Under a recessive model, the number of homozygous derived genotypes will be proportional to load

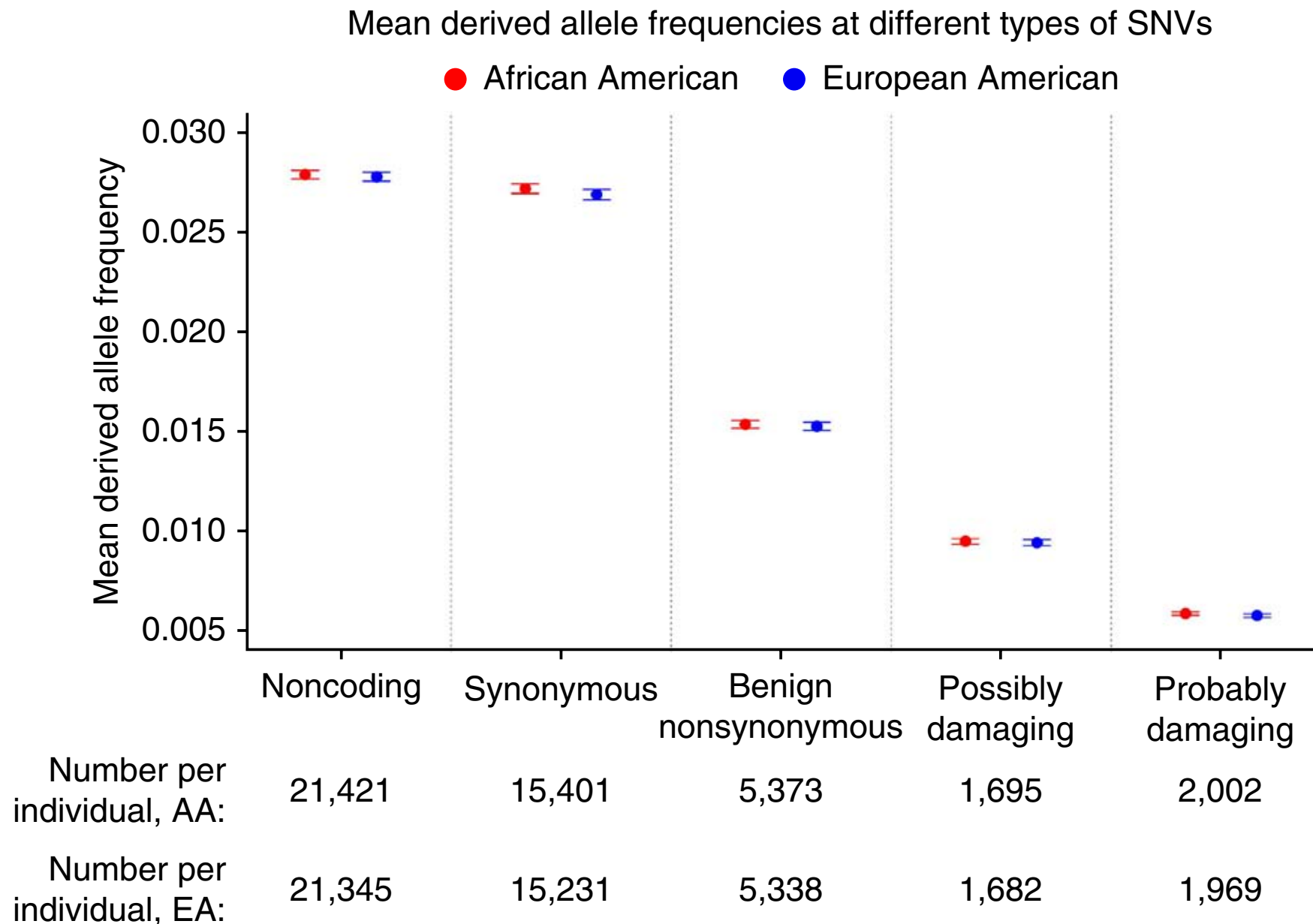
Inferring Genetic Load

- It is widely appreciated that African ancestry individuals have more variation overall than individuals with European ancestry.
- However, European individuals have more **homozygous** variation.
- Increased load?

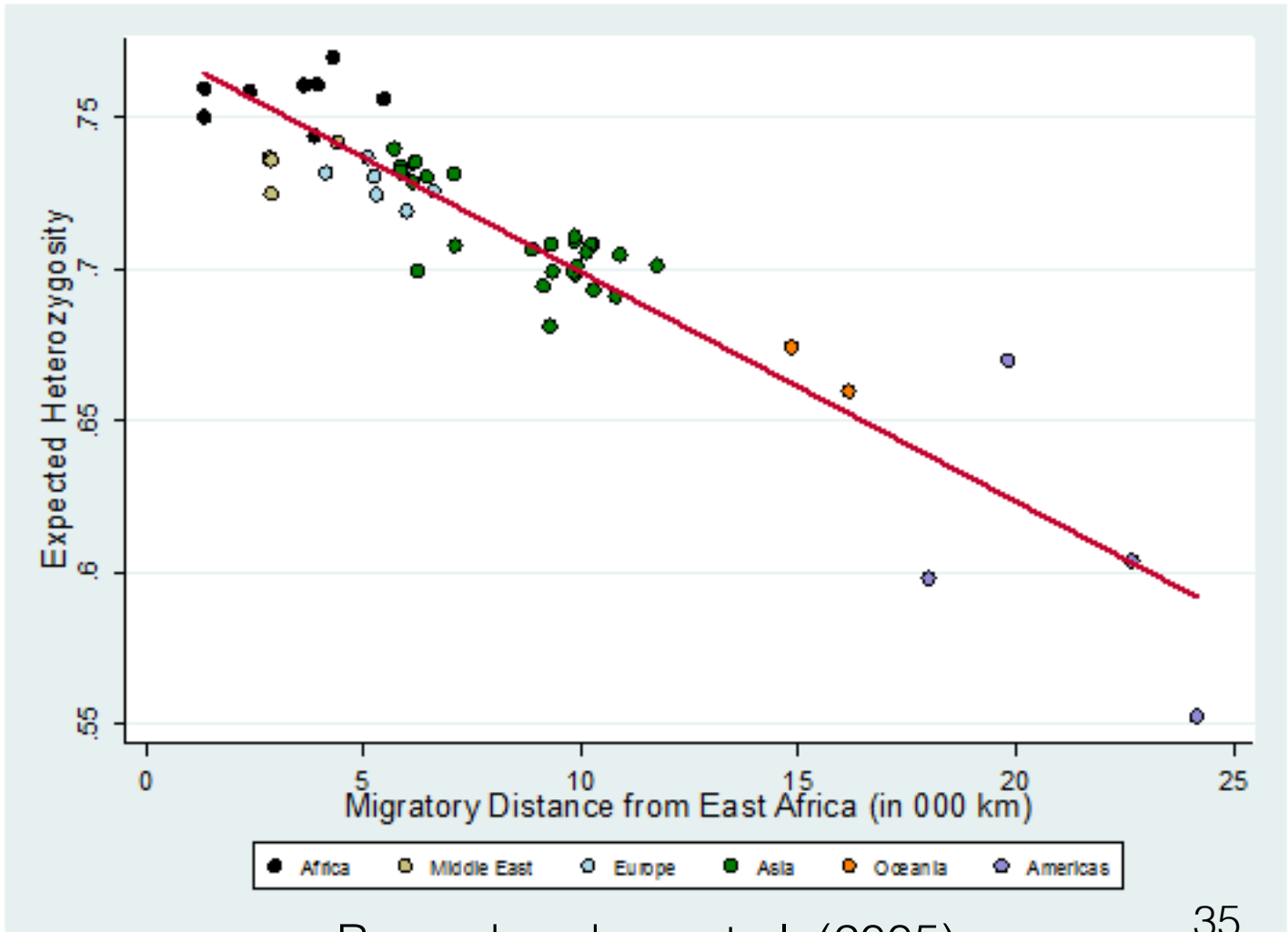
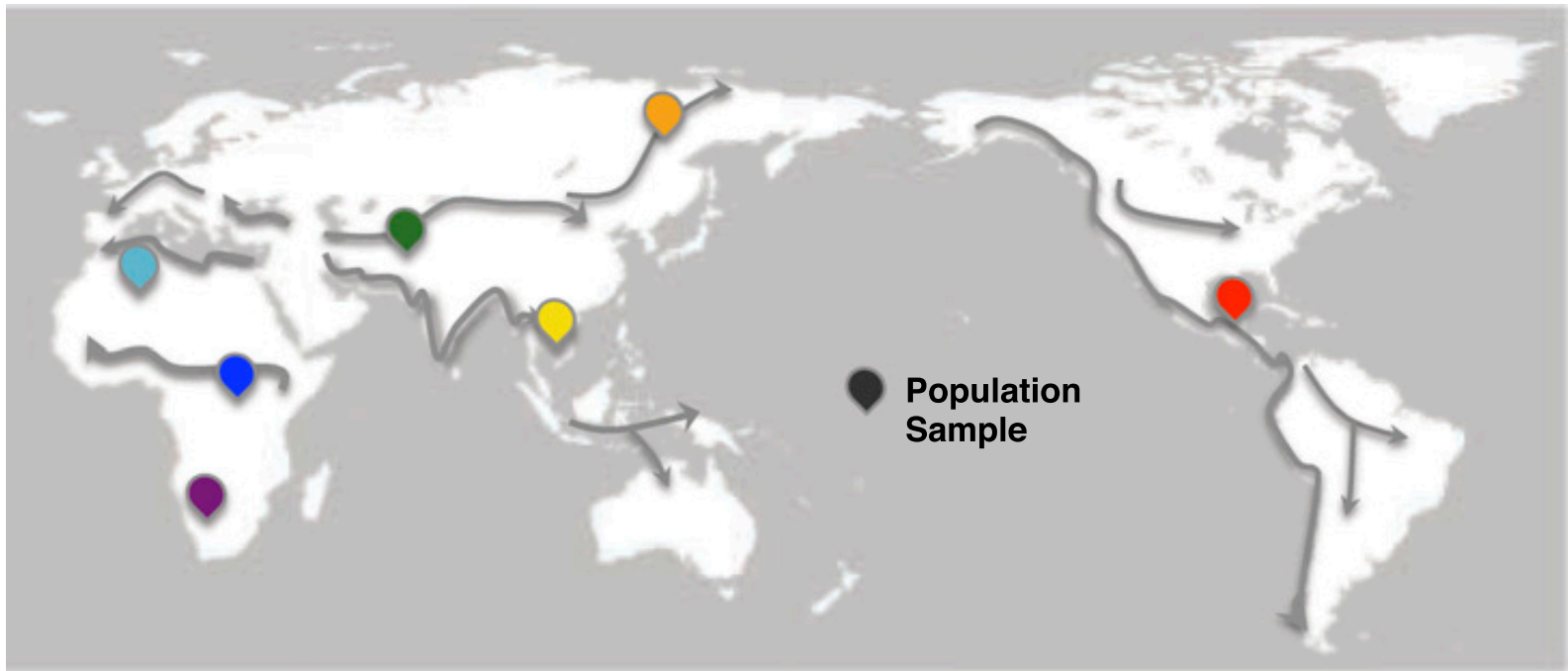


Inferring Genetic Load

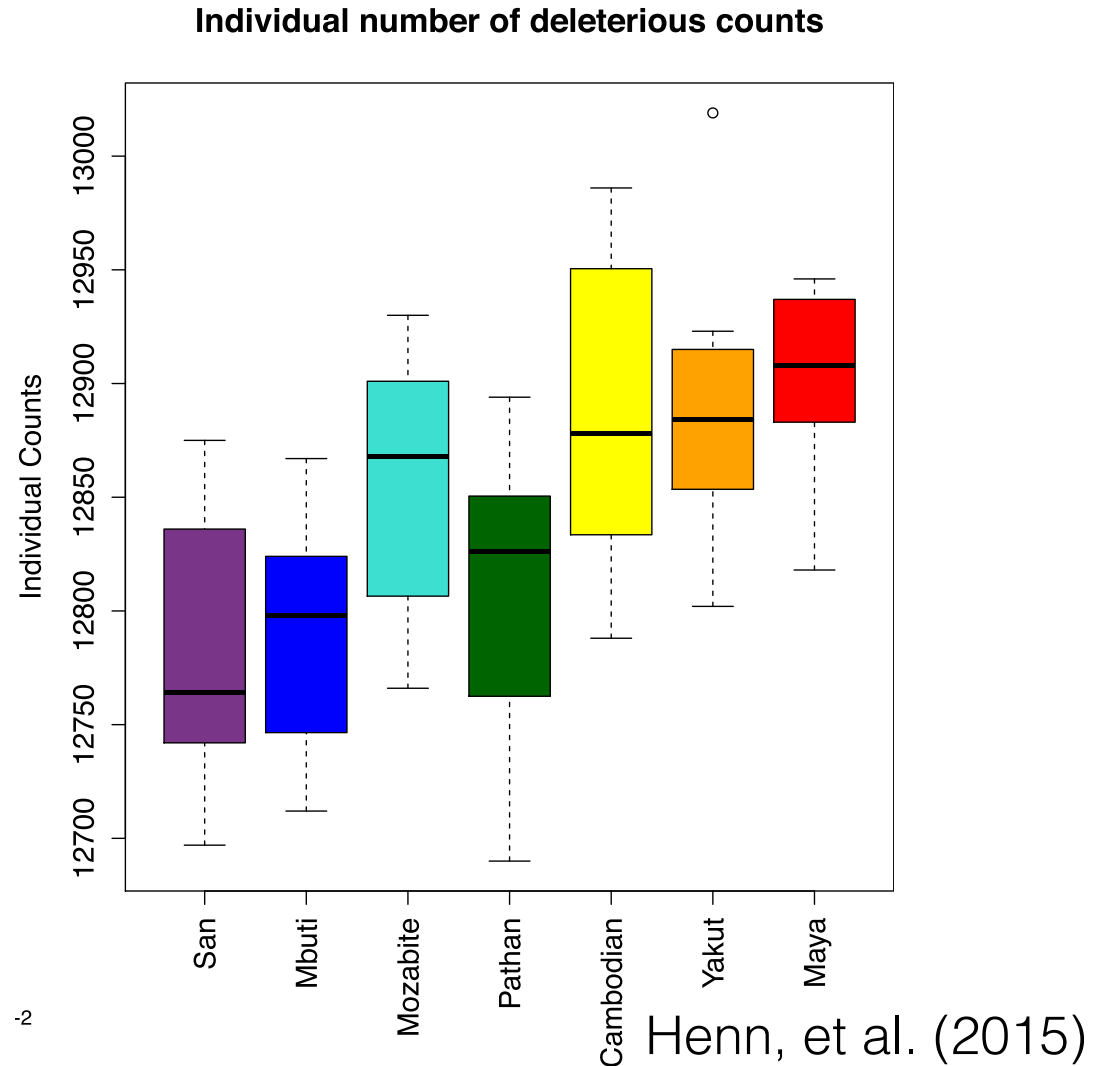
- Het. and hom. derived alleles appear to balance between African and European Americans!
- All individuals have ~same number of derived alleles!



Serial Founder Effects on Genetic Load



35



-2

Key Feature of Natural Selection

- Alleles change frequency unusually fast
 - Positive selection tends to increase frequency
 - Negative selection tends to decrease frequency
- All tests for natural selection seek to identify this feature using different aspects of the data.
- While negative selection shapes majority of patterns of variation in many species, positive selection may drive patterns of local variation.

The Effect of Positive Selection

Adaptive

Neutral

Nearly Neutral

Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



The Effect of Positive Selection

Adaptive

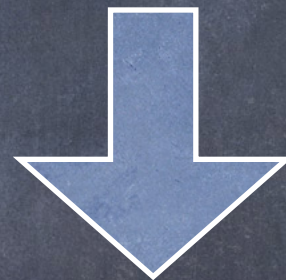
Neutral

Nearly Neutral

Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



Types of Positive Selection

- Selection acts in one population but not another
 - Frequencies of the selected alleles in one population will go up relatively quickly compared to the frequencies of those same alleles in the other population.
 - The test is simple:
 - Are there alleles that have unusually large allele frequency differences between two populations?

Testing for Population Divergence

- Imagine two populations diverged several thousand years ago.
- One population stayed where it was, but the other migrated up a mountain to the Tibetan Plateau.
 - Many environmental changes...
 - Not obvious where in the genome to look for adaptations
 - Try exome sequencing

Testing for Population Divergence

Sequencing of 50 Human Exomes Reveals Adaptation to High Altitude

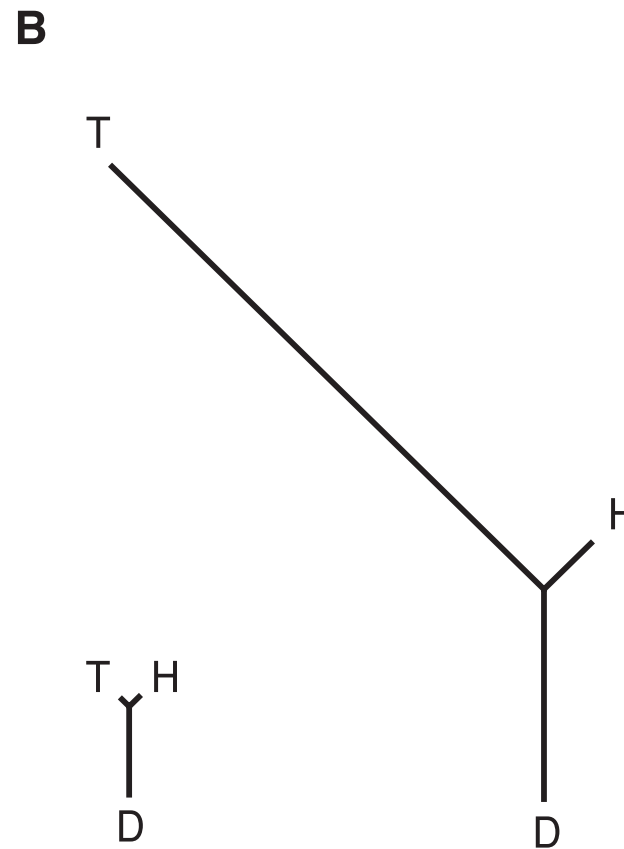
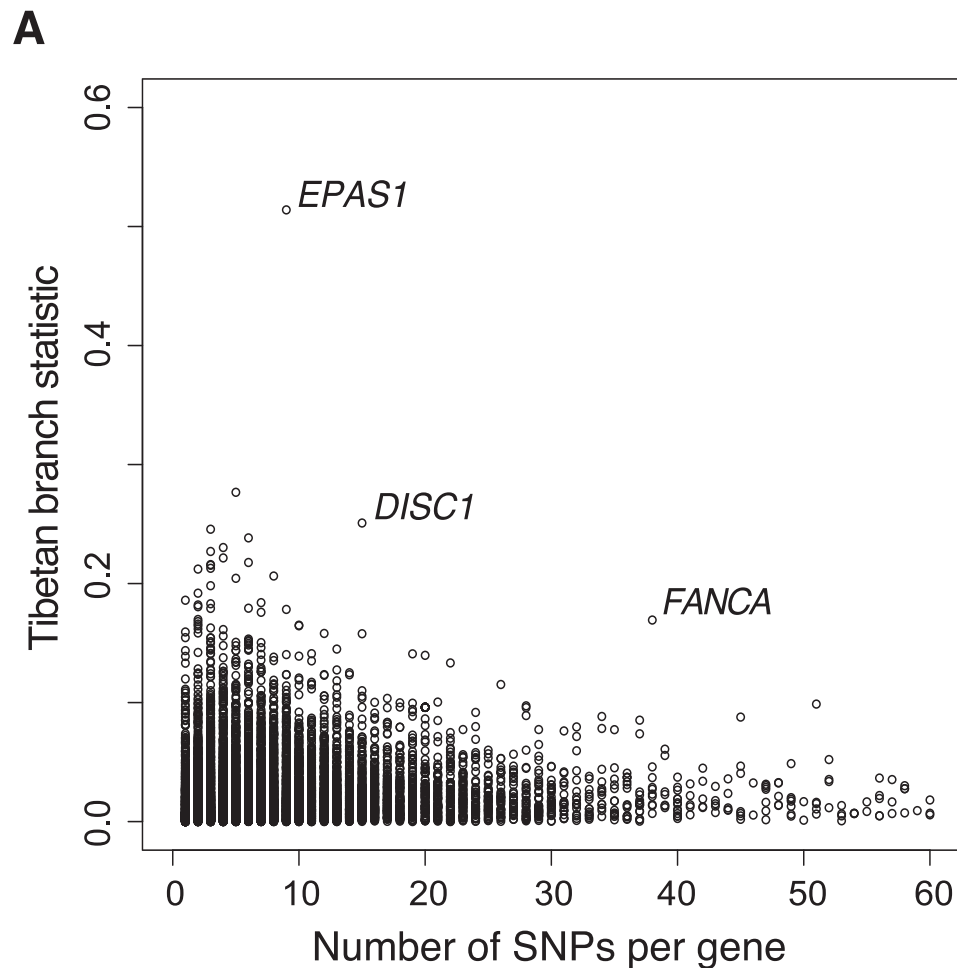
Xin Yi,^{1,2*} Yu Liang,^{1,2*} Emilia Huerta-Sanchez,^{3*} Xin Jin,^{1,4*} Zha Xi Ping Cuo,^{2,5*} John E. Pool,^{3,6*} Xun Xu,¹ Hui Jiang,¹ Nicolas Vinckenbosch,³ Thorfinn Sand Korneliussen,⁷ Hancheng Zheng,^{1,4} Tao Liu,¹ Weiming He,^{1,8} Kui Li,^{2,5} Ruibang Luo,^{1,4} Xifang Nie,¹ Honglong Wu,^{1,9} Meiru Zhao,¹ Hongzhi Cao,^{1,9} Jing Zou,¹ Ying Shan,^{1,4} Shuzheng Li,¹ Qi Yang,¹ Asan,^{1,2} Peixiang Ni,¹ Geng Tian,^{1,2} Junming Xu,¹ Xiao Liu,¹ Tao Jiang,^{1,9} Renhua Wu,¹ Guangyu Zhou,¹ Meifang Tang,¹ Junjie Qin,¹ Tong Wang,¹ Shuijian Feng,¹ Guohong Li,¹ Huasang,¹ Jiangbai Luosang,¹ Wei Wang,¹ Fang Chen,¹ Yading Wang,¹ Xiaoguang Zheng,^{1,2} Zhuo Li,¹ Zhuoma Bianba,¹⁰ Ge Yang,¹⁰ Xinpeng Wang,¹¹ Shuhui Tang,¹¹ Guoyi Gao,¹² Yong Chen,⁵ Zhen Luo,⁵ Lamu Gusang,⁵ Zheng Cao,¹ Qinghui Zhang,¹ Weihai Ouyang,¹ Xiaoli Ren,¹ Huiqing Liang,¹ Huisong Zheng,¹ Yebo Huang,¹ Jingxiang Li,¹ Lars Bolund,¹ Karsten Kristiansen,^{1,7} Yingrui Li,¹ Yong Zhang,¹ Xiuqing Zhang,¹ Ruiqiang Li,^{1,7} Songgang Li,¹ Huanming Yang,¹ Rasmus Nielsen,^{1,3,7}† Jun Wang,^{1,7}† Jian Wang¹†

Testing for Population Divergence

EPAS1: a transcription factor involved in response to hypoxia

- To find these types of signatures:
 - Compare allele frequencies using F_{st}

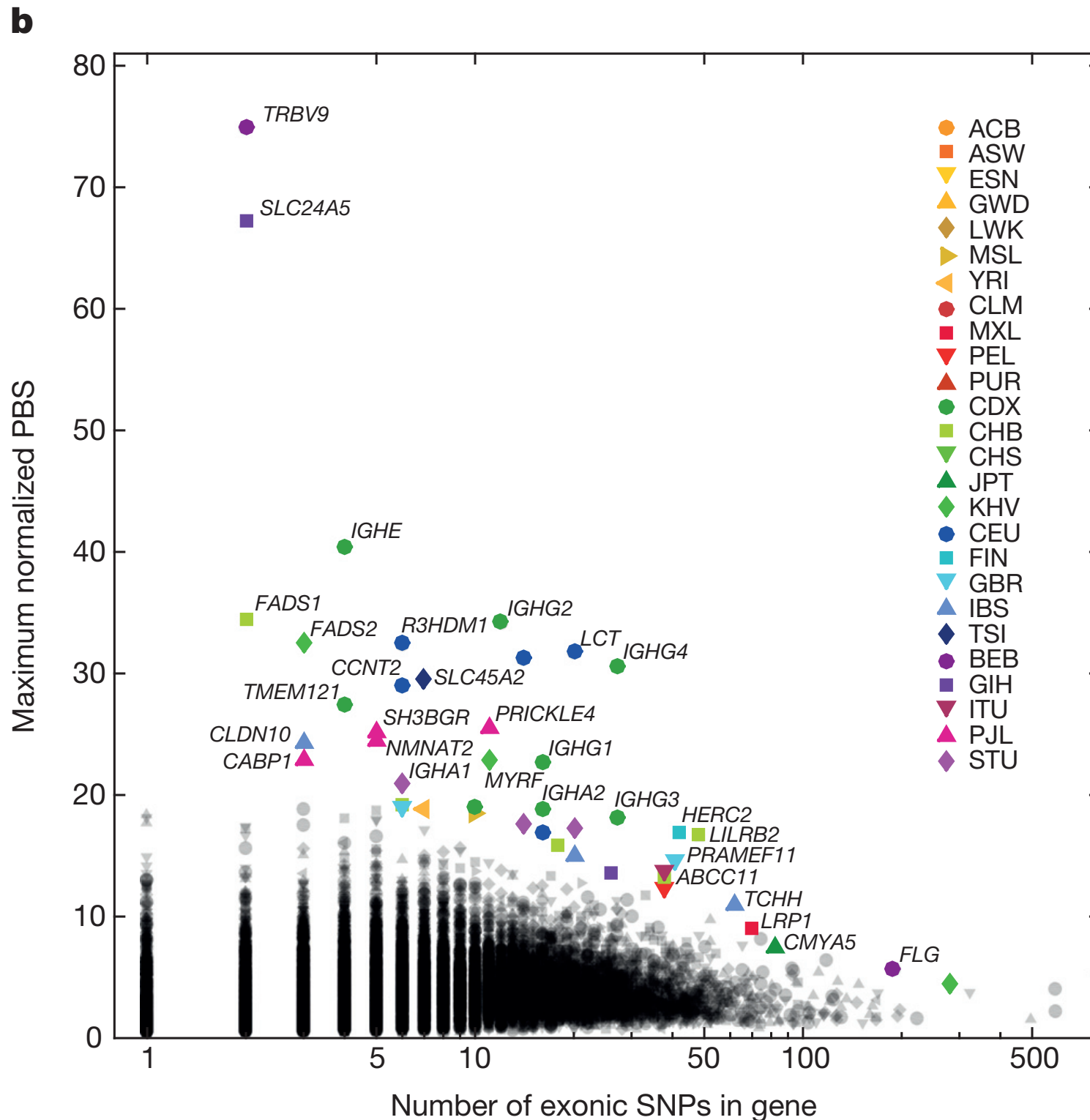
Testing for Population Divergence



EPAS1: a transcription factor involved in response to hypoxia


- To find these types of signatures:
 - Compare allele frequencies using F_{st}

Testing for Population Divergence



- Applying this statistic to 26 human populations
- Several known genes
- Several novel ones

Types of Positive Selection

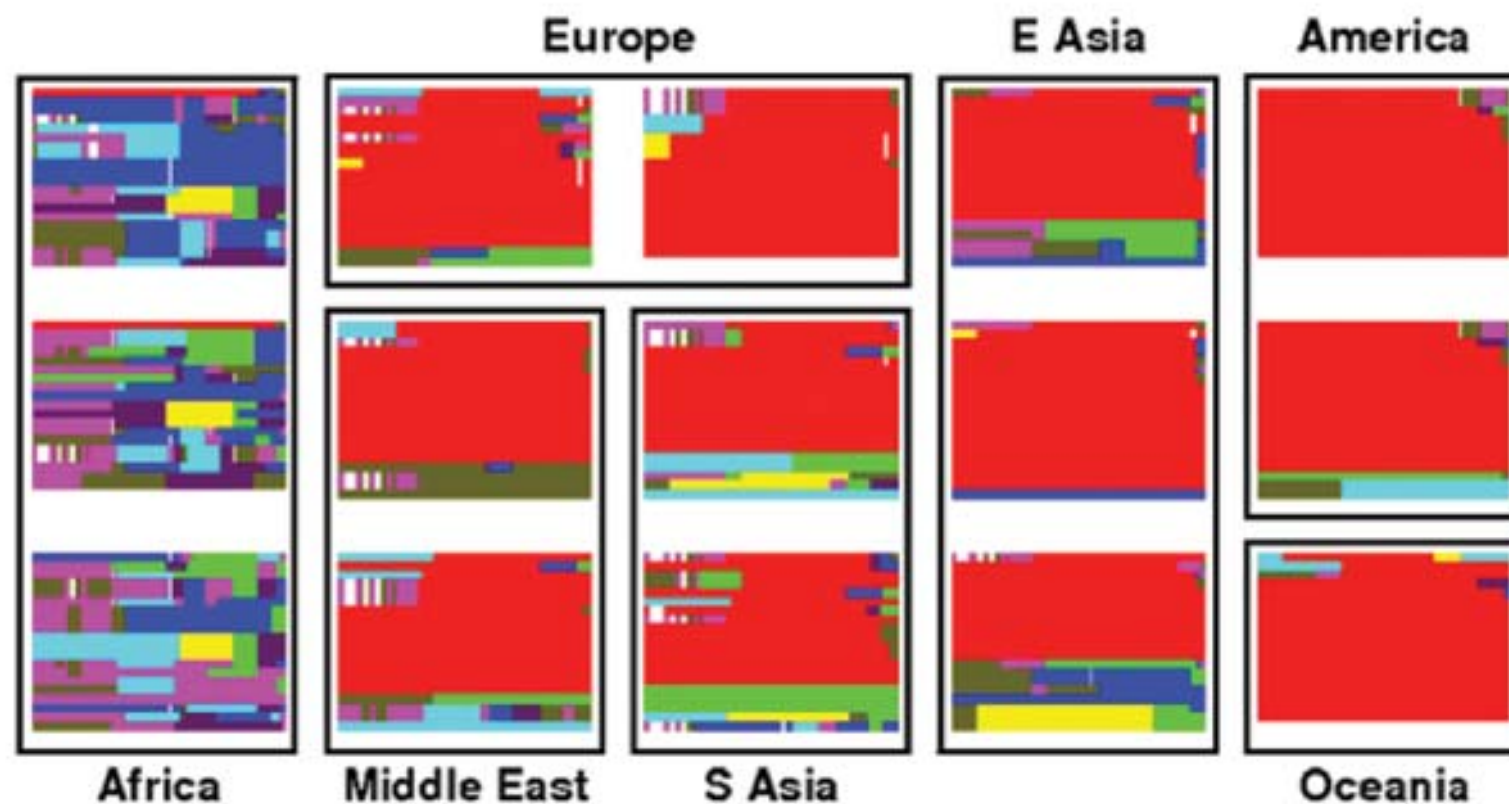
-  Selection acts in one population but not another
- Selection operates on a new mutation
 - Selection will act to increase the frequency of the allele
 - Results in a young allele at relatively high frequency
 - The test is simple:
 - Are there young alleles at unusually high frequency?

Testing for High Freq. Young Alleles

- The age of an allele can be assessed by measuring the amount of genetic variation around the allele.
 - As time passes:
 - Mutations occur nearby
 - Recombination breaks down the correlation between the allele and others nearby

Testing for High Freq. Young Alleles

- Example: Skin pigmentation
 - KITLG is a gene known to contribute to lighter skin in non-African populations.



- Each plot is a population.
- Each row is an individual's haplotype.
- Identical haplotypes have the same color.
- Large red blocks indicate long haplotypes with very little variation (i.e., young).

Testing for High Freq. Young Alleles

- Detecting these types of signatures:
 - Long Range Haplotype (**LRH**) or Extended Haplotype Homozygosity (**EHH**) {Sabeti, P. C. et al. Nature 419, 832-837 (2002)}.
 - integrated Haplotype Score (**iHS**) {Voight, B. F. et al. PLoS Biol 4, e72 (2006)}.
 - Composite Likelihood Ratio (**CLR**) {Williamson, S. H. et al. PLoS Genet 3, e90 (2007)}.

Types of Positive Selection

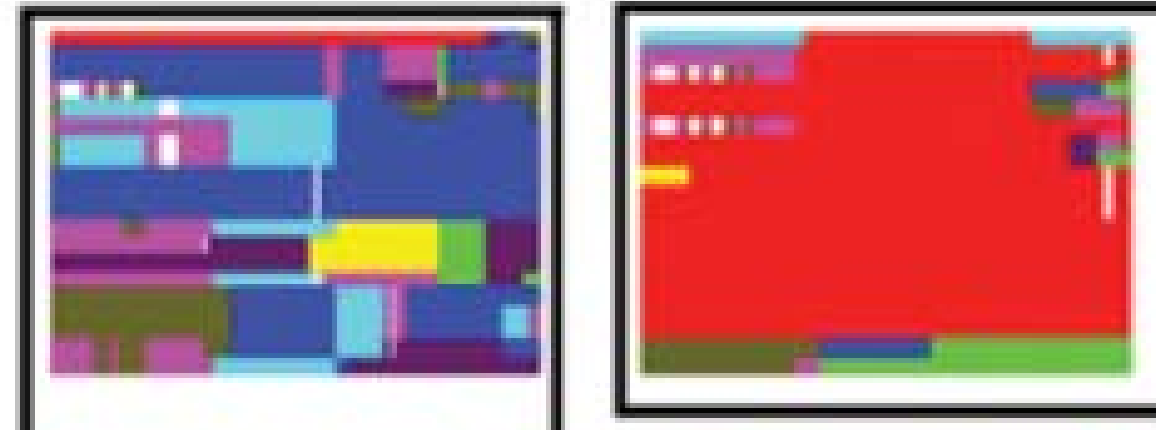
 Selection acts in one population but not another

 Selection acts on a new mutation

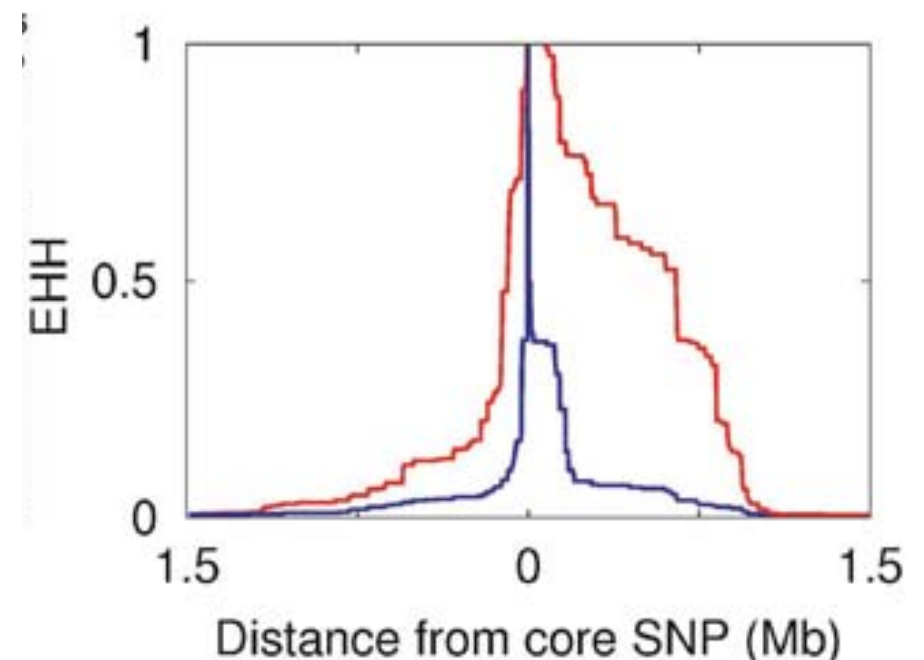
- Selection acts on new mutations primarily in one population
 - In this case, we expect high divergence and long haplotypes in one population

Divergence of a Young Allele

- Recall the haplotype patterns before for just two populations:

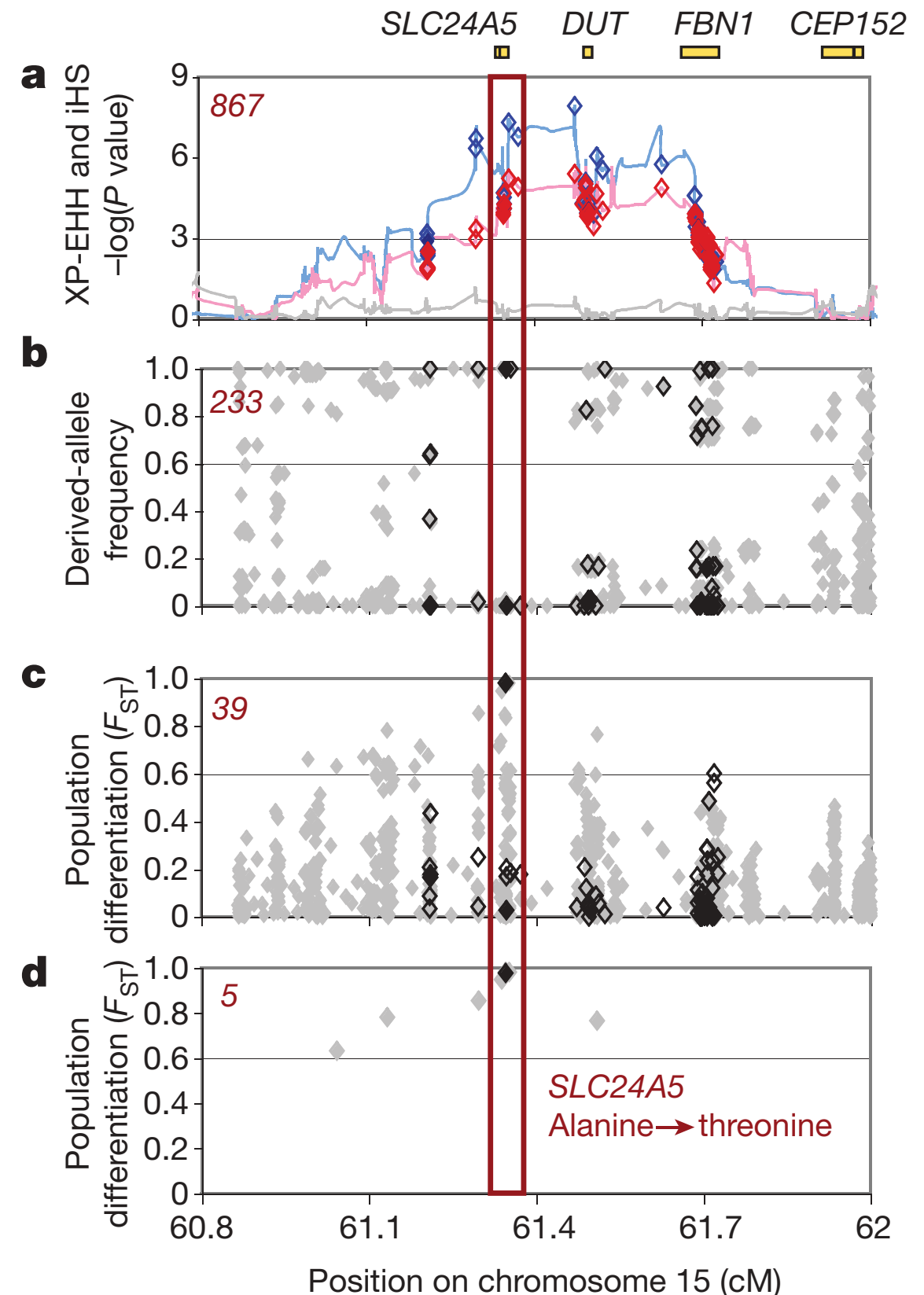


- These can be plotted as the probability that two randomly chosen individuals have an identical haplotype as a function of distance from the core SNP:
- Comparing the area under these two curves is the basis for XP-EHH



Divergence of a Young Allele

- XP-EHH rediscovers a nonsynonymous variant in *SLC24A5* contributing to lighter skin outside Africa.



Motivation

- Why should we care about finding signatures of natural selection?
 - It's cool... It's what often drives speciation
 - Understanding disease/complex traits

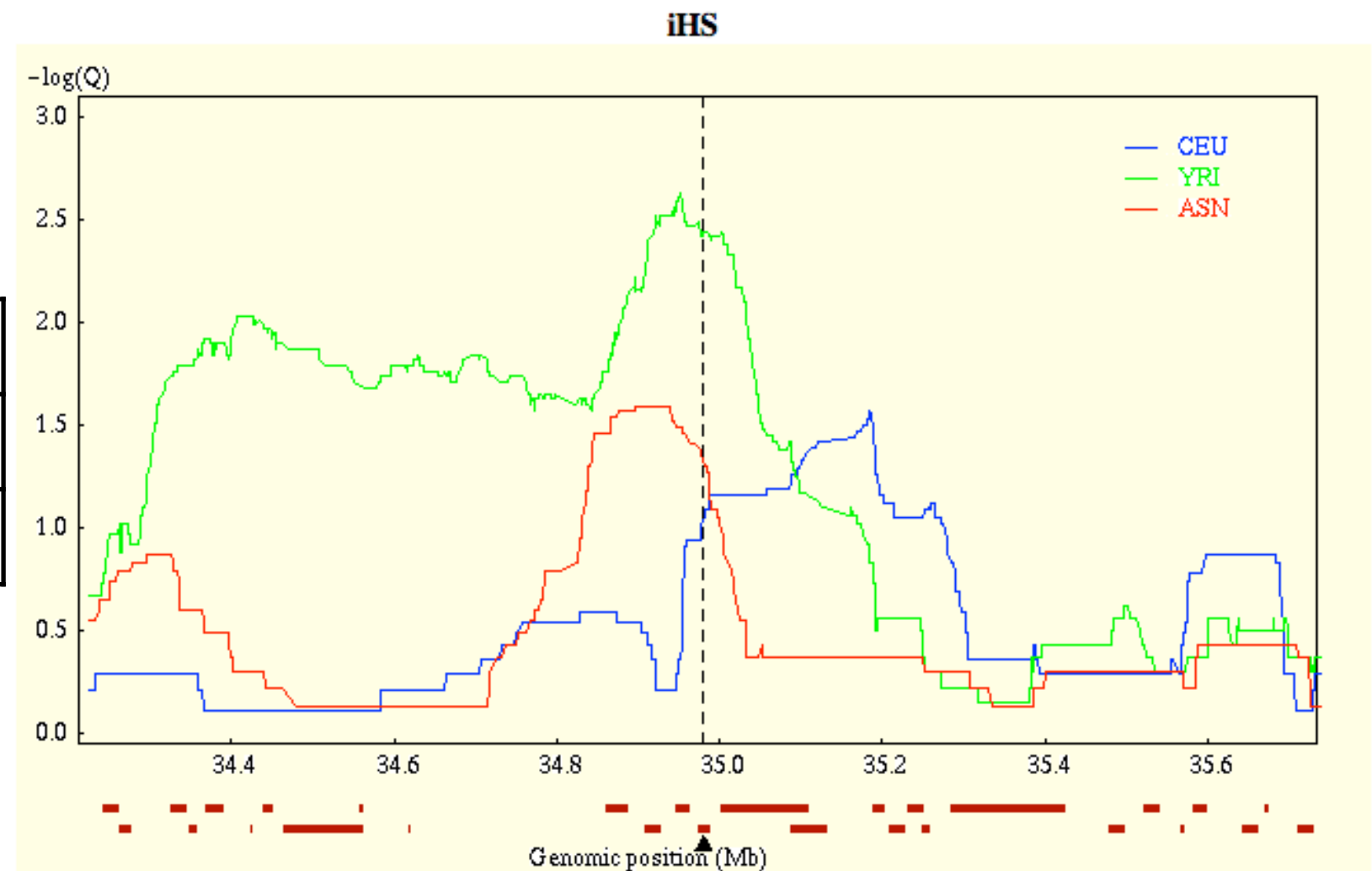
Case Study: Kidney Disease in African Americans

- Individuals of African descent have much higher incidence of kidney disease than individuals of European descent.
- GWAS had previously implicated the gene MYH9 with moderate effects ($p < 10^{-8}$)
- But there was no clear biological story.

Case Study: Kidney Disease in African Americans

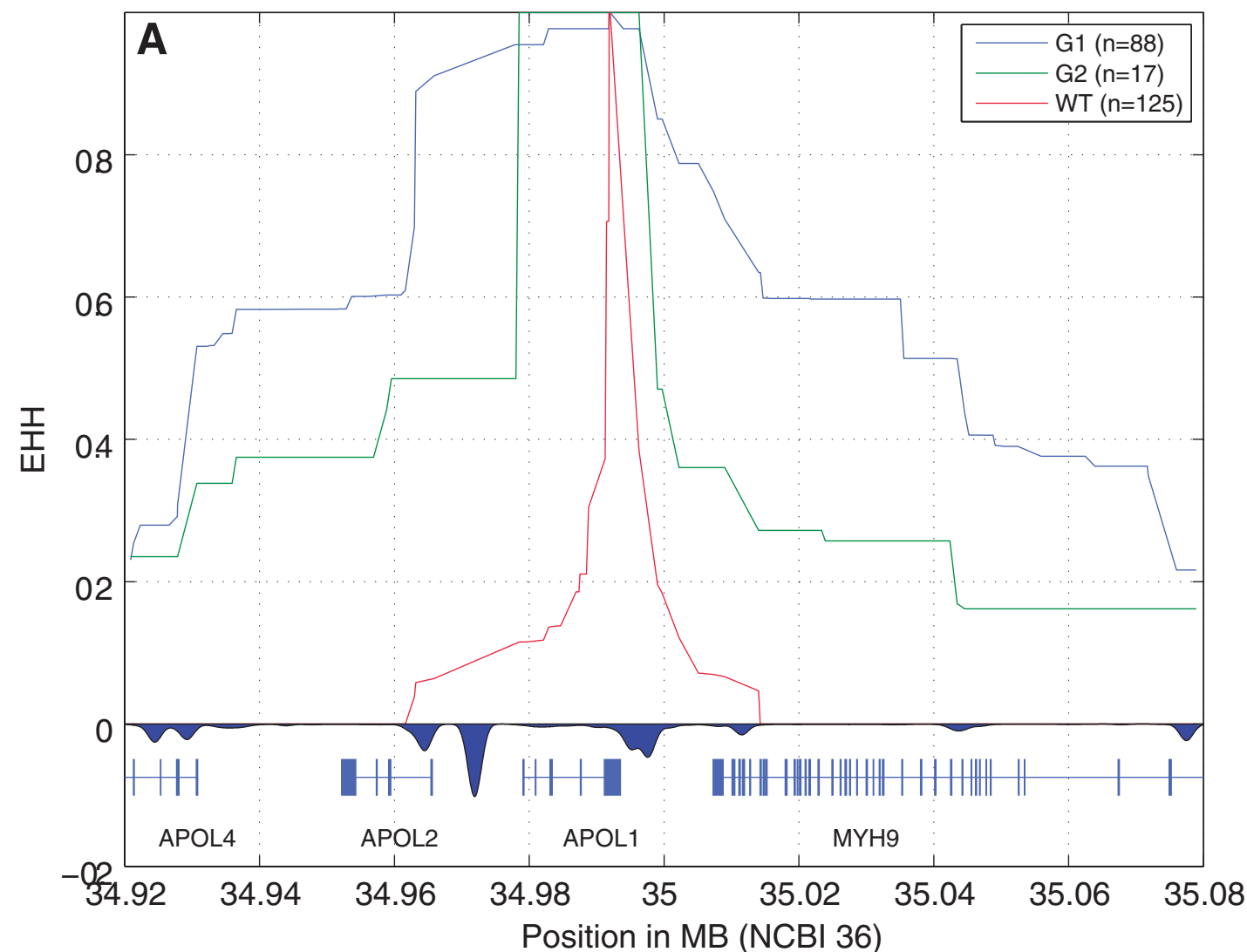
- Looking at signatures of selection adds valuable insight.
- Consider iHS from haplotter.uchicago.edu (more on this later):

Gene	iHS p-value
APOLI	0.0033
MYH9	0.014



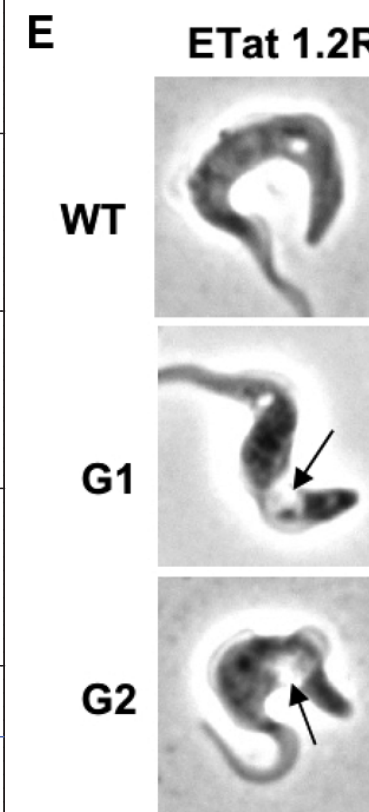
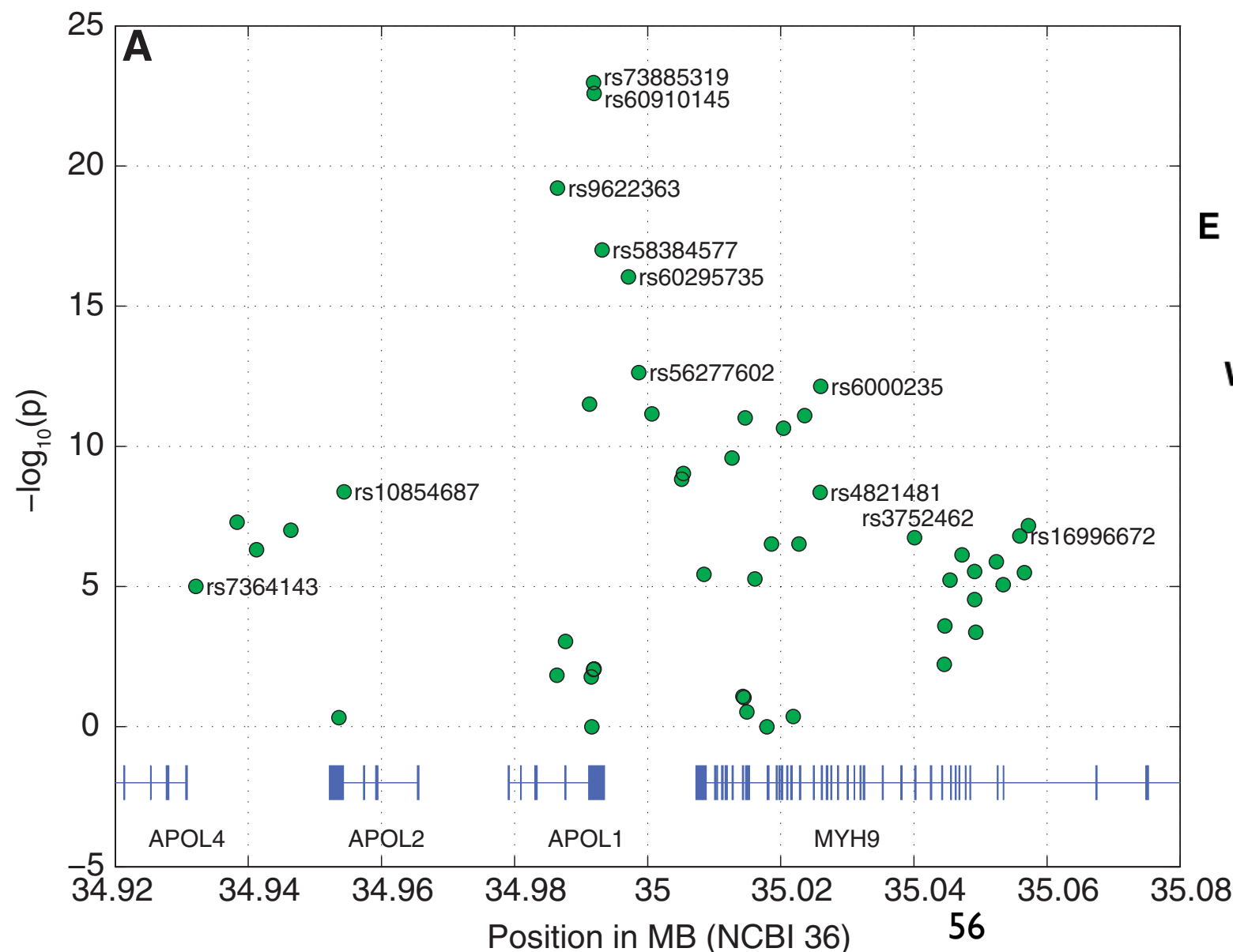
Case Study: Kidney Disease in African Americans

- Tag SNPs chosen across a broader region, and calculated EHH based on higher resolution data



Case Study: Kidney Disease in African Americans

- Subset of SNPs chosen based on signatures of selection genotyped on a larger panel strongly implicates APOL1!



Risk alleles confer resistance to trypanosomes (swelling of the lysosome).

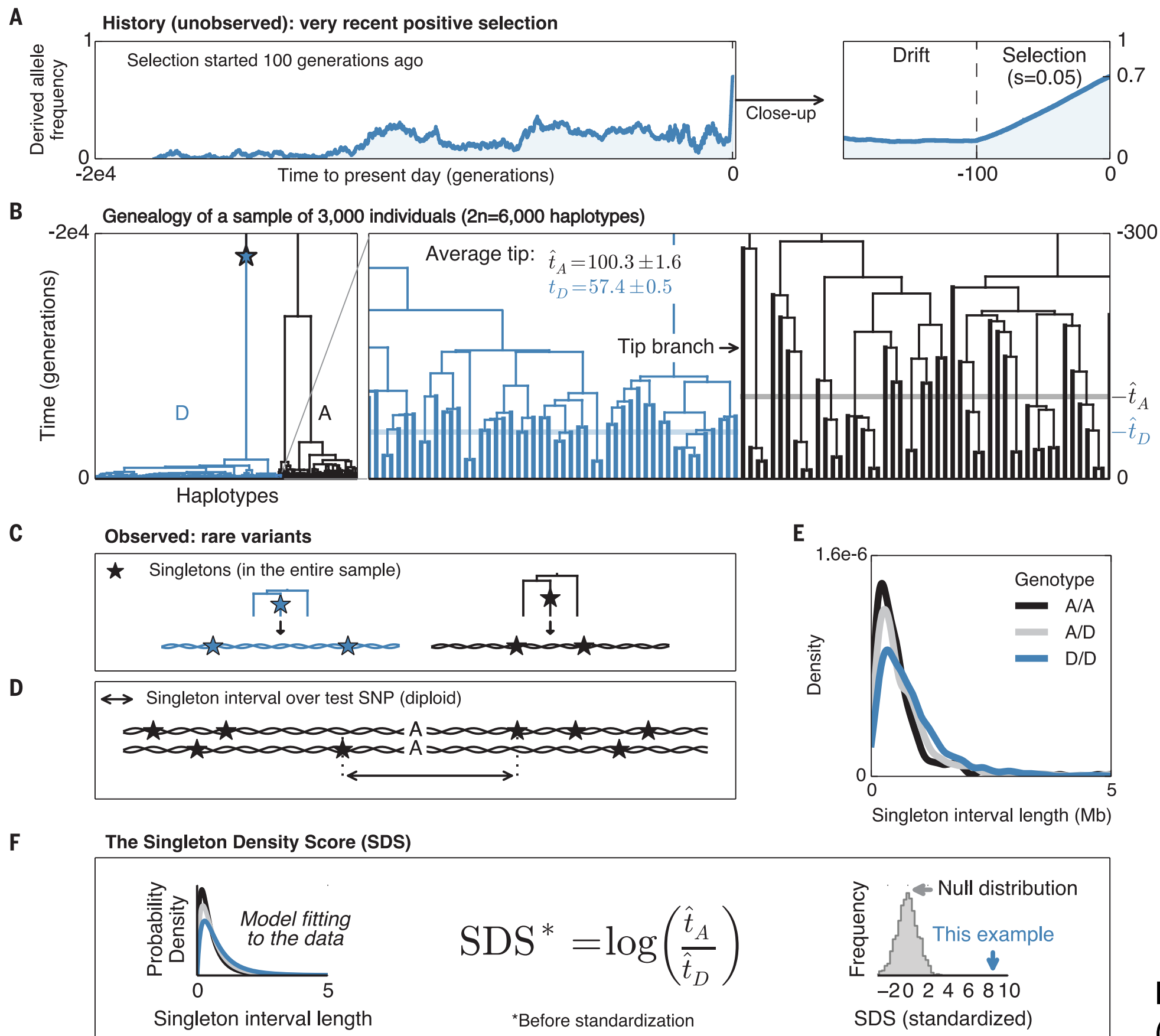
WGS

- The statistics described do not really handle whole genome sequencing data (WGS).
- Further, the timescale for when selection acted is not very well specified.
- With an abundance of rare variants, WGS should be informative about recent selection.
- Enter the Singleton Density Score (SDS).

SDS

- Field, et al. (*Science*, 2016) introduced the Singleton Density Score (SDS) to capitalize on WGS data with very large samples.
- In the presence of a sweep, the distribution of distances (across individuals) to the nearest singleton will be skewed towards longer distances.

SDS

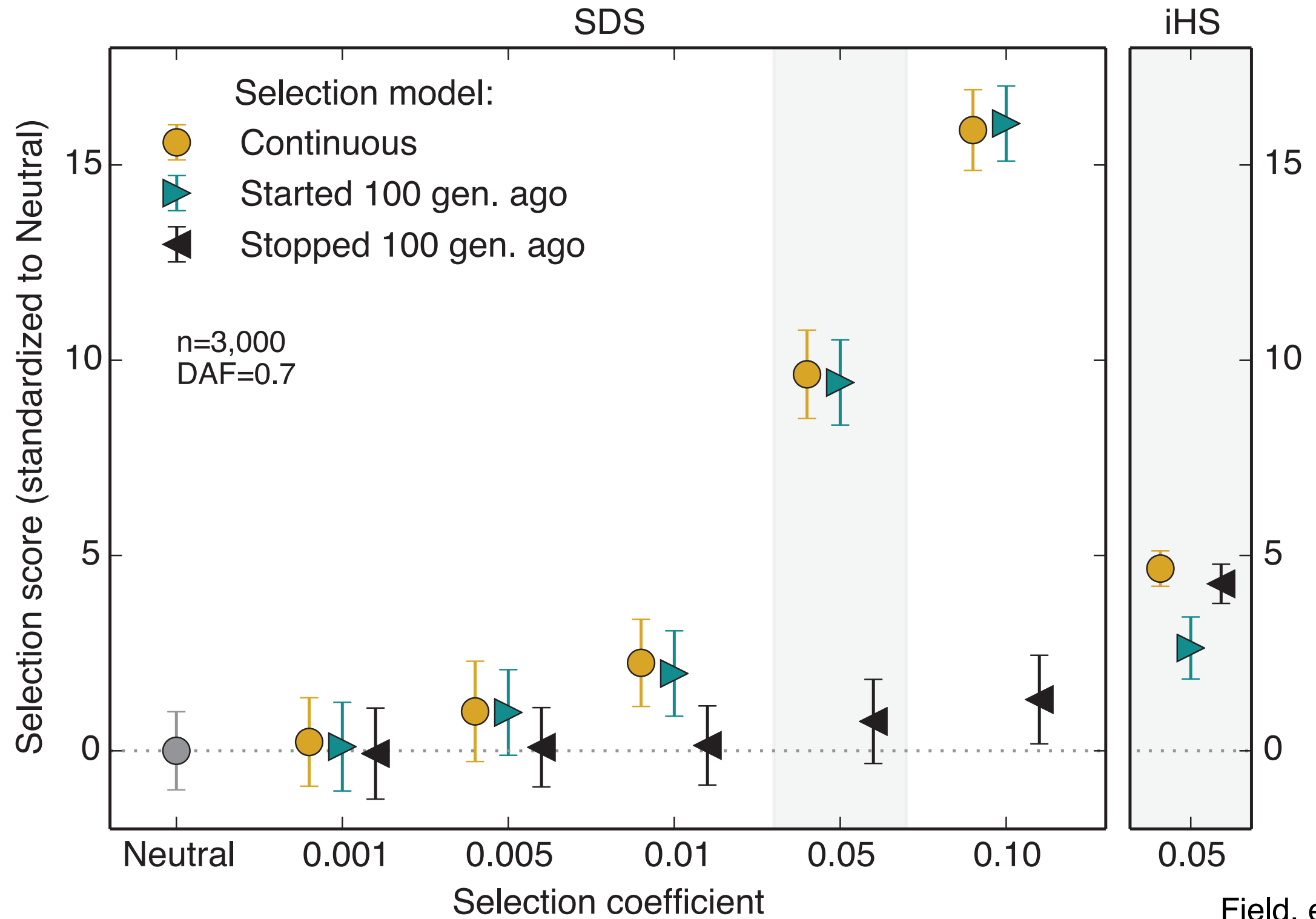


Field, et al.
(Science, 2016)

SDS

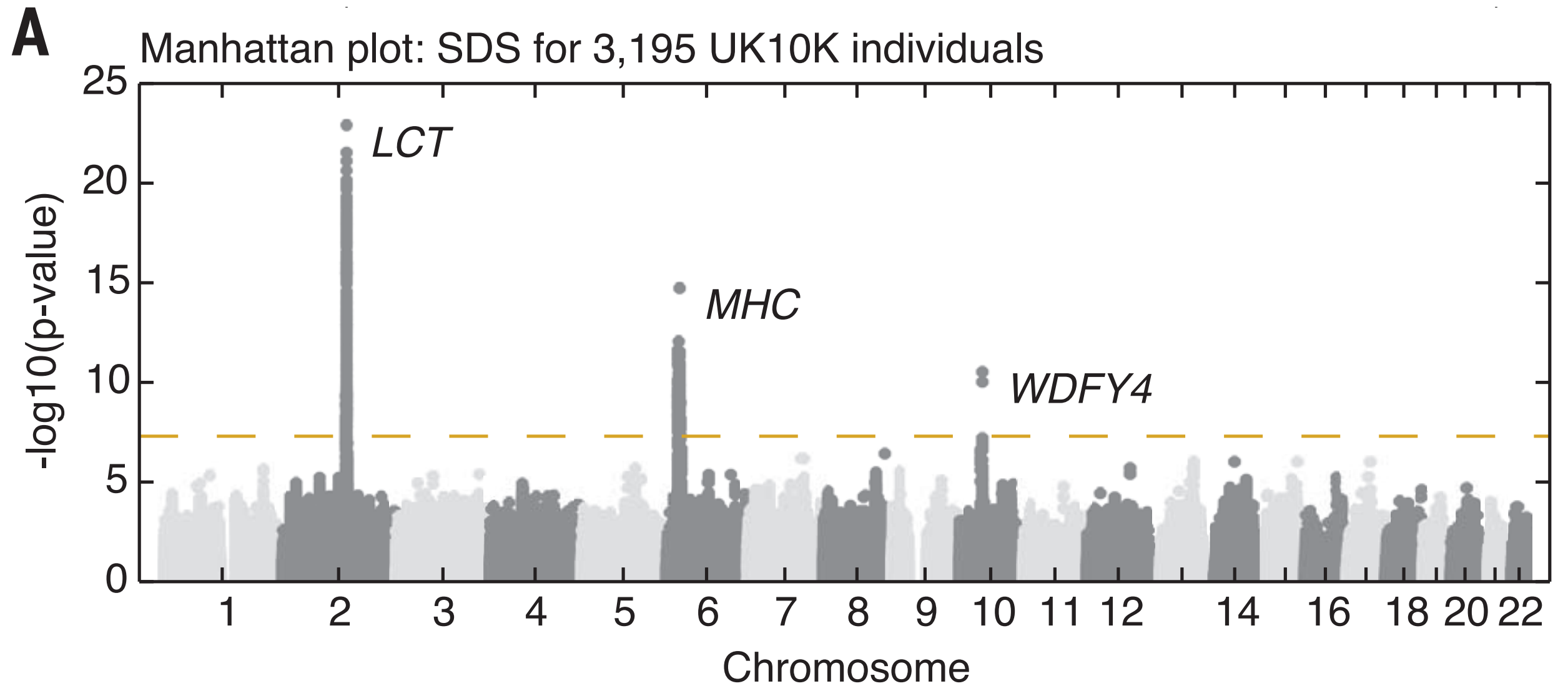
B

Simulations: signal and specificity of our method to recent history



Field, et al.
(Science, 2016)

SDS



Conclusions

- Natural selection leaves distinctive footprints within patterns of genetic variation.
- This occurs because alleles driven by natural selection tend to be younger than neutral alleles at the same frequency.
- Characterizing signatures of natural selection around disease associated loci can sometimes illuminate mechanistic relationships.