

Searching for Signatures of Selection with Selscan

Ryan Hernandez

The Effect of Positive Selection

Adaptive

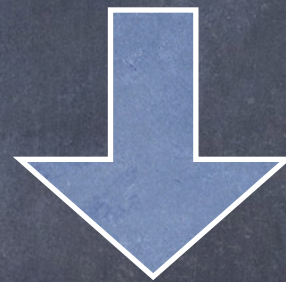
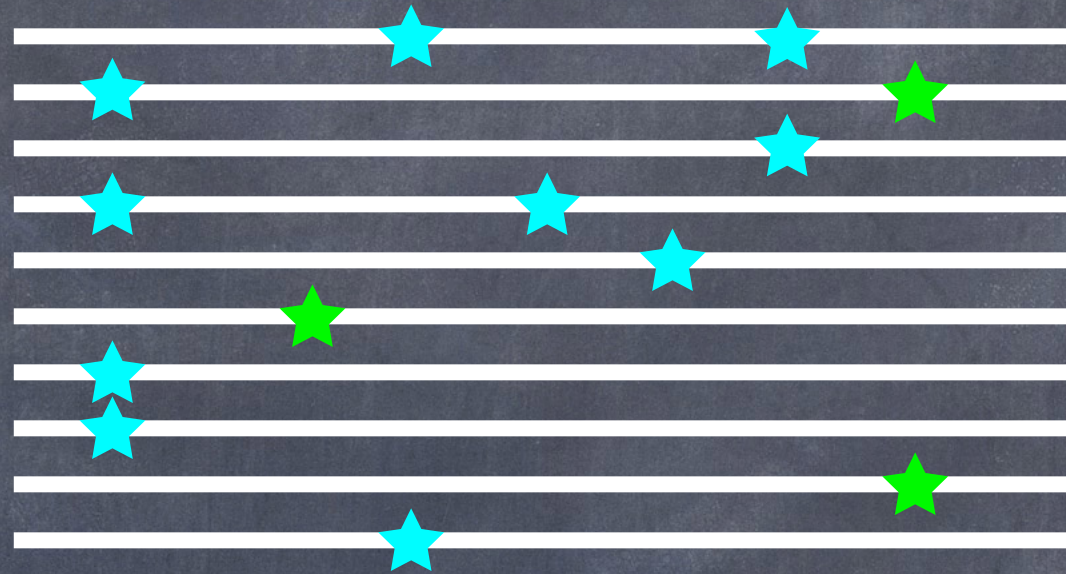
Neutral

Nearly Neutral

Mildly Deleterious

Fairly Deleterious

Strongly Deleterious



How do we capture this process in a statistic?

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

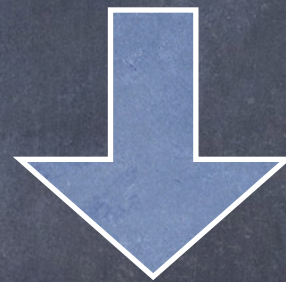
Neutral

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

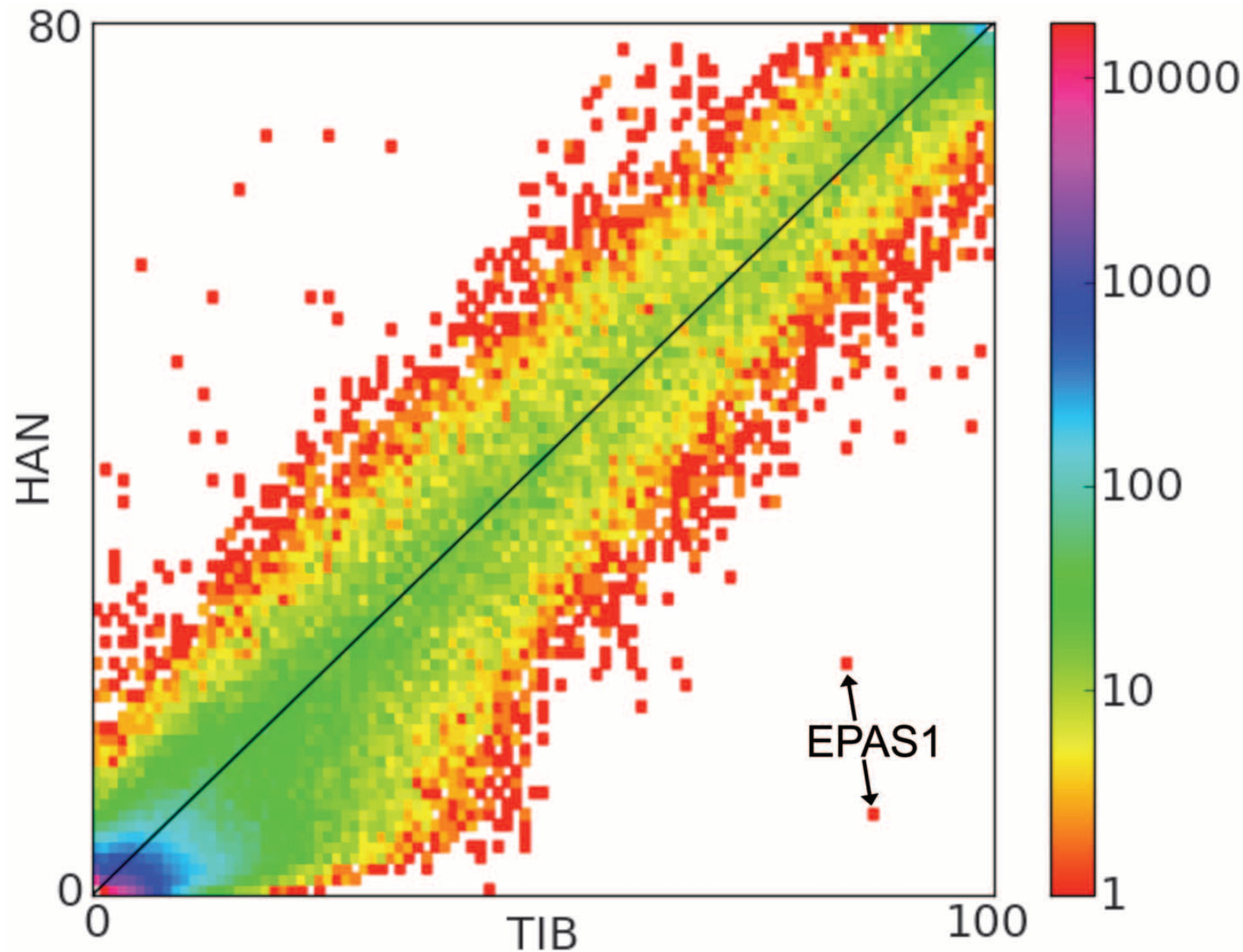
- Imagine two populations diverged several thousand years ago.
- One population stayed where it was, but the other migrated up a mountain to the Tibetan Plateau.
 - Many environmental changes...
 - Not obvious where in the genome to look for adaptations
 - Try exome sequencing

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 Weihan 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 †

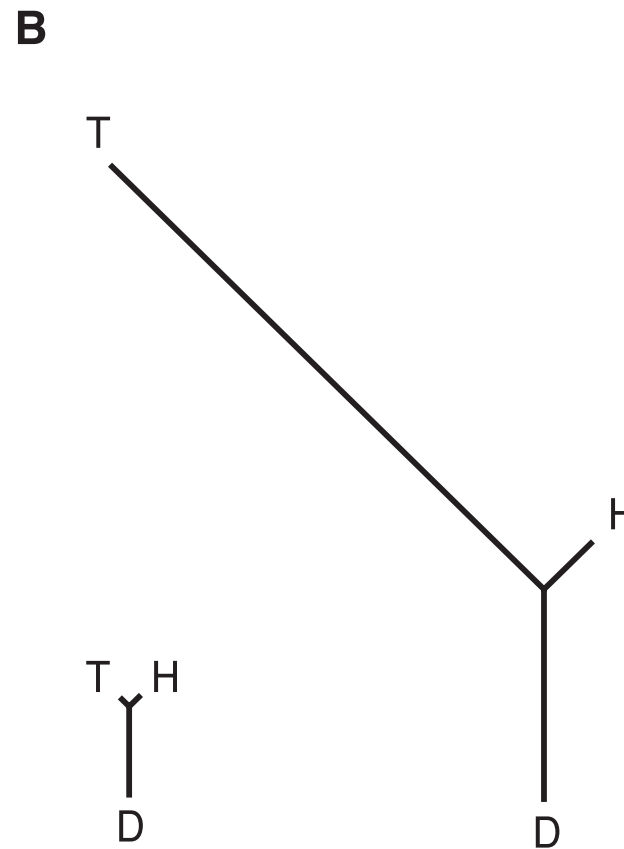
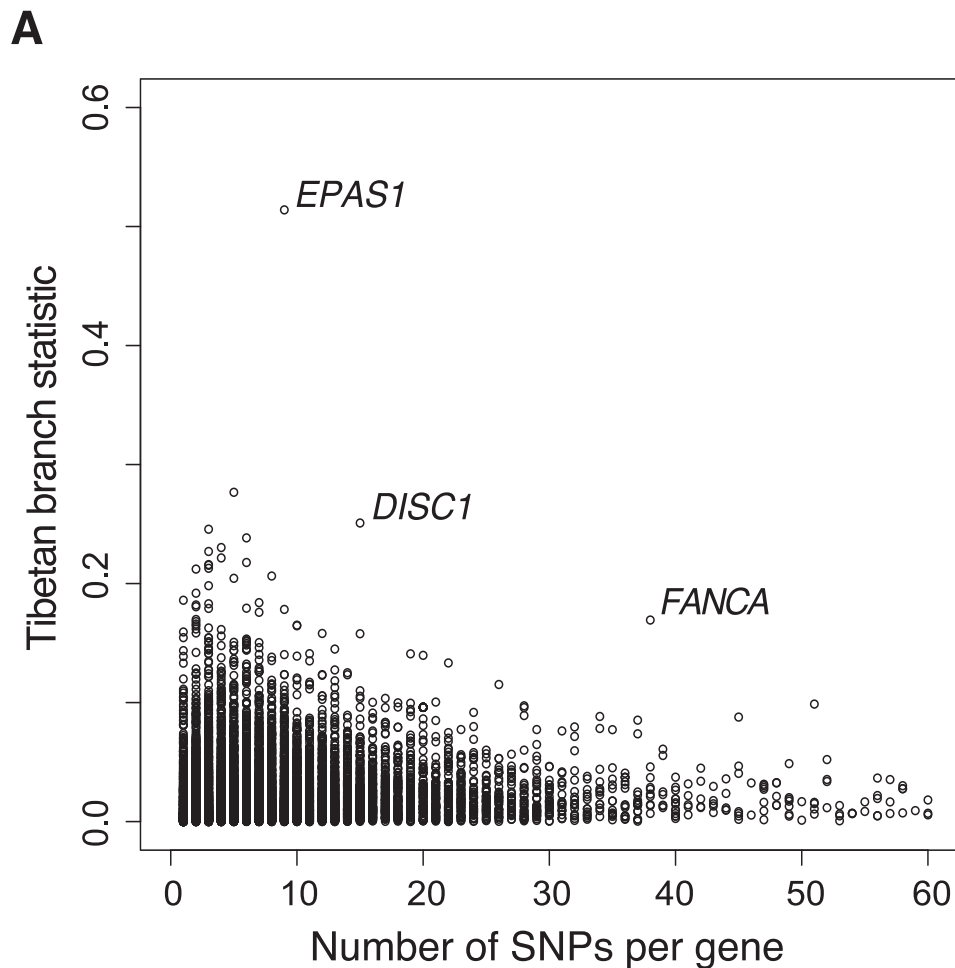
Testing for Population Divergence



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

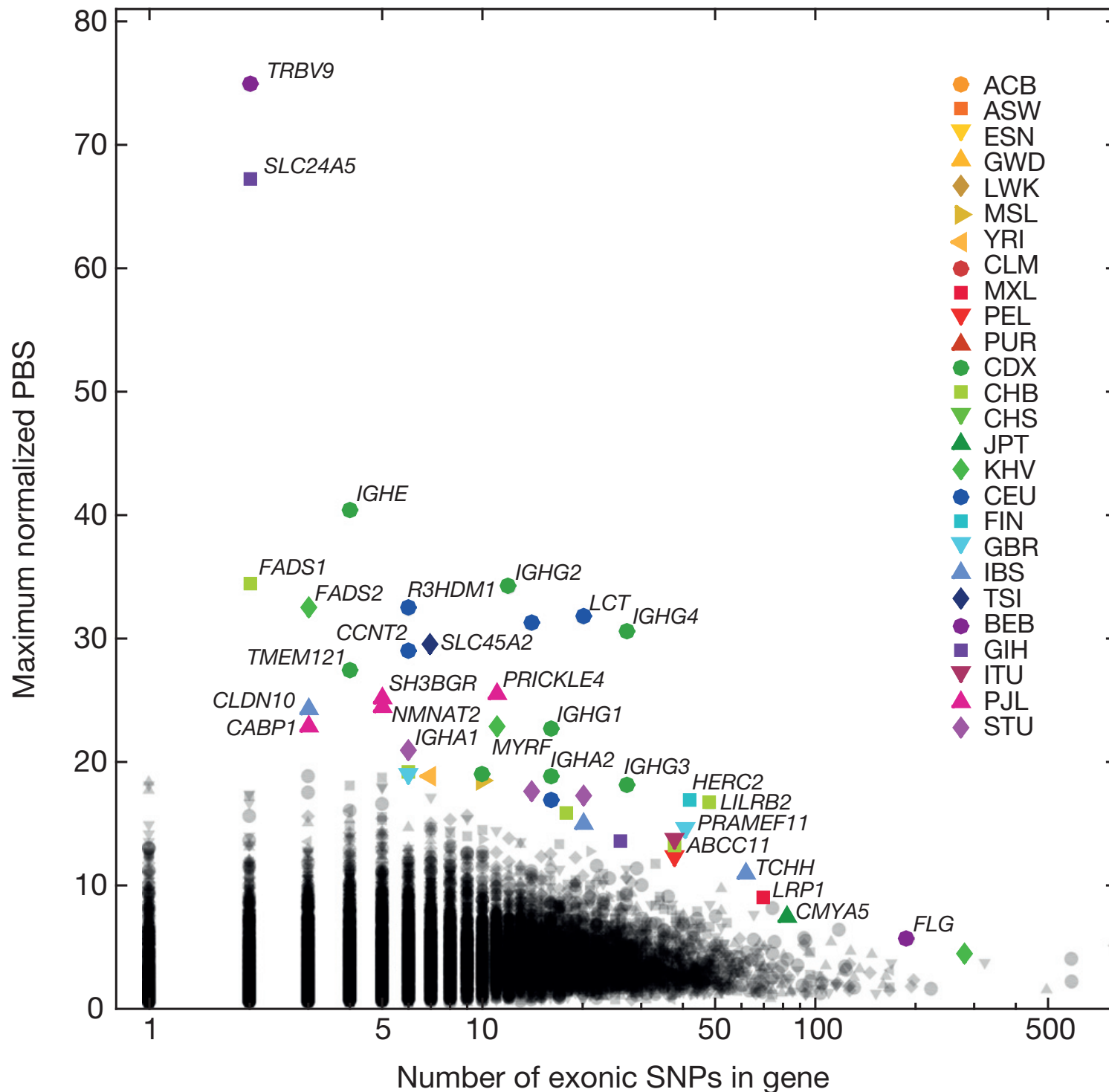


EPAS1: a transcription factor involved in response to hypoxia

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

Testing for Population Divergence

b



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

Types of Positive Selection

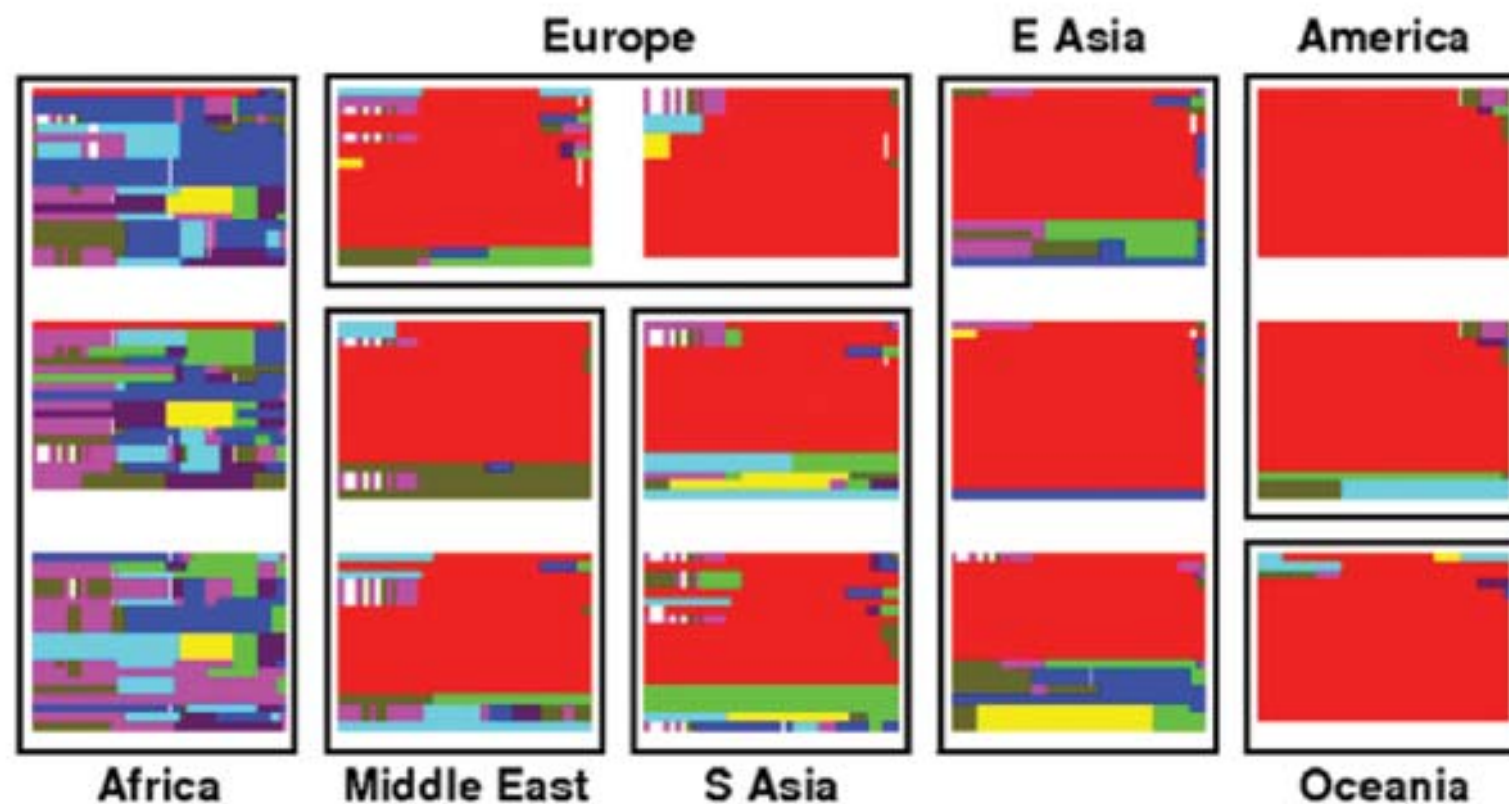
- Selection acts in one population but not another
- Selection operates on a new mutation
 - Selection will act to increase the frequency of the allele
 - Results in a young allele at relatively high frequency
 - The test is simple:
 - Are there young alleles at unusually high frequency?

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

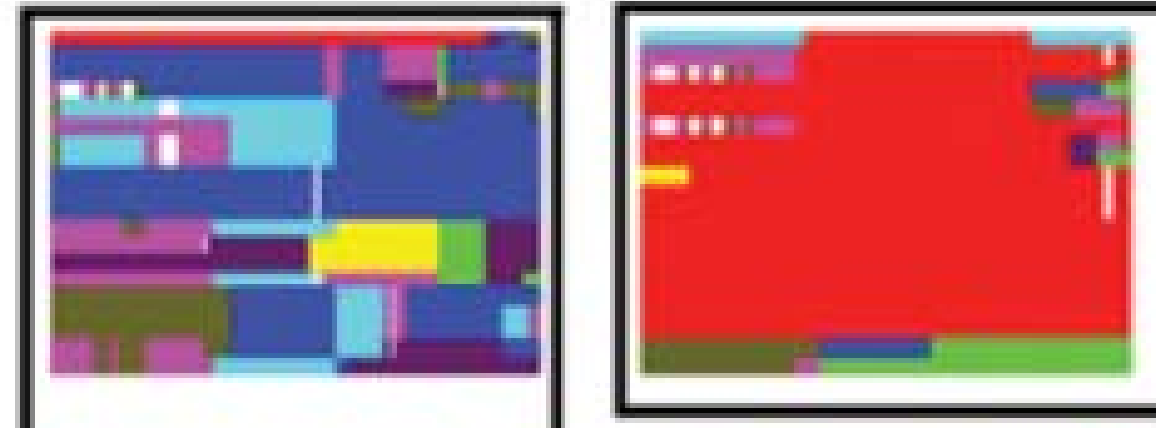
Selection acts in one population but not another

Selection acts on a new mutation

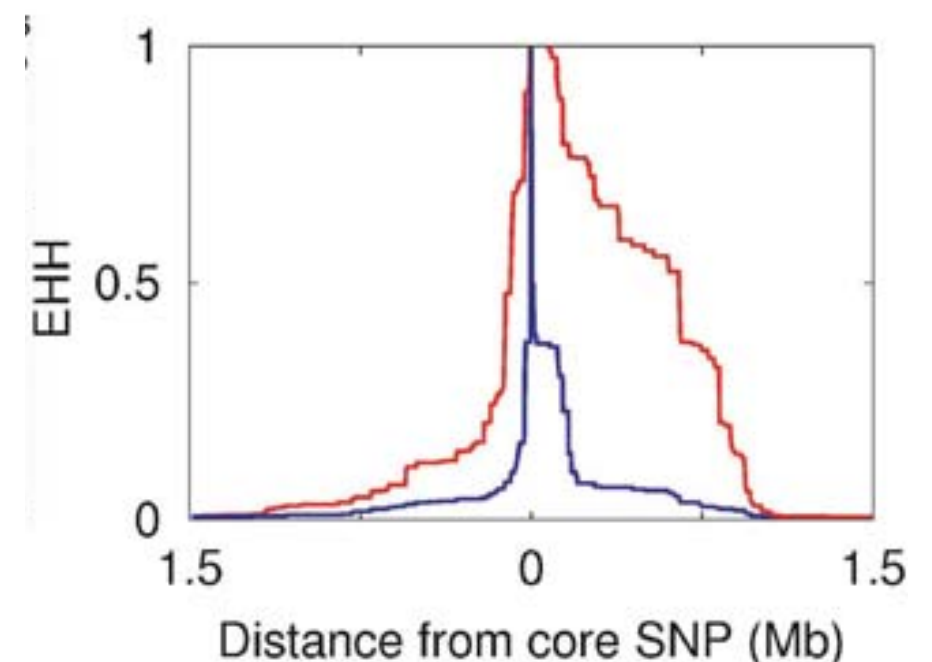
- Selection acts on new mutations primarily in one population
 - In this case, we expect high divergence and long haplotypes in one population

Divergence of a Young Allele

- Recall the haplotype patterns before for just two populations:

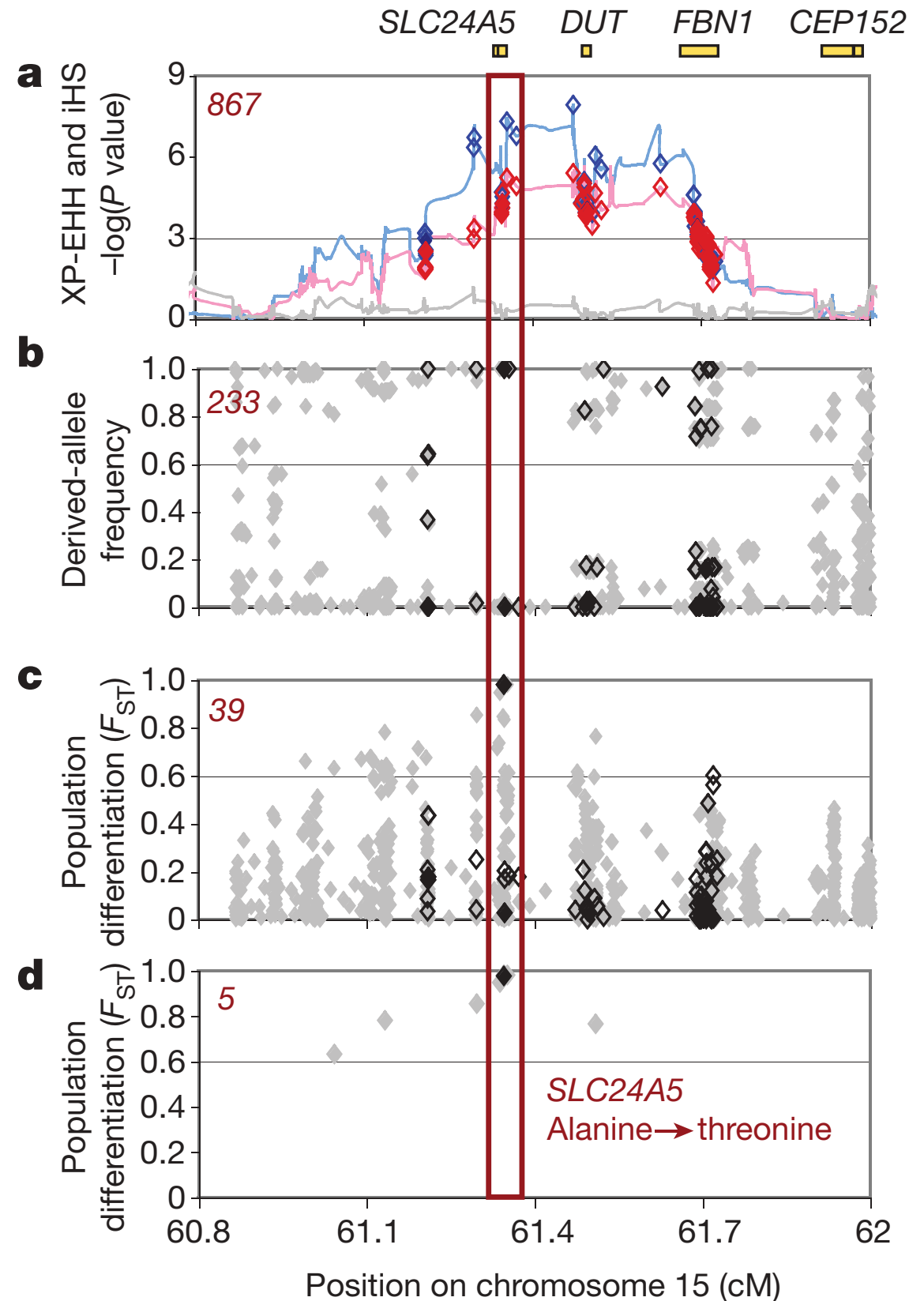


- These can be plotted as the probability that two randomly chosen individuals have an identical haplotype as a function of distance from the core SNP:
- Comparing the area under these two curves is the basis for XP-EHH



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

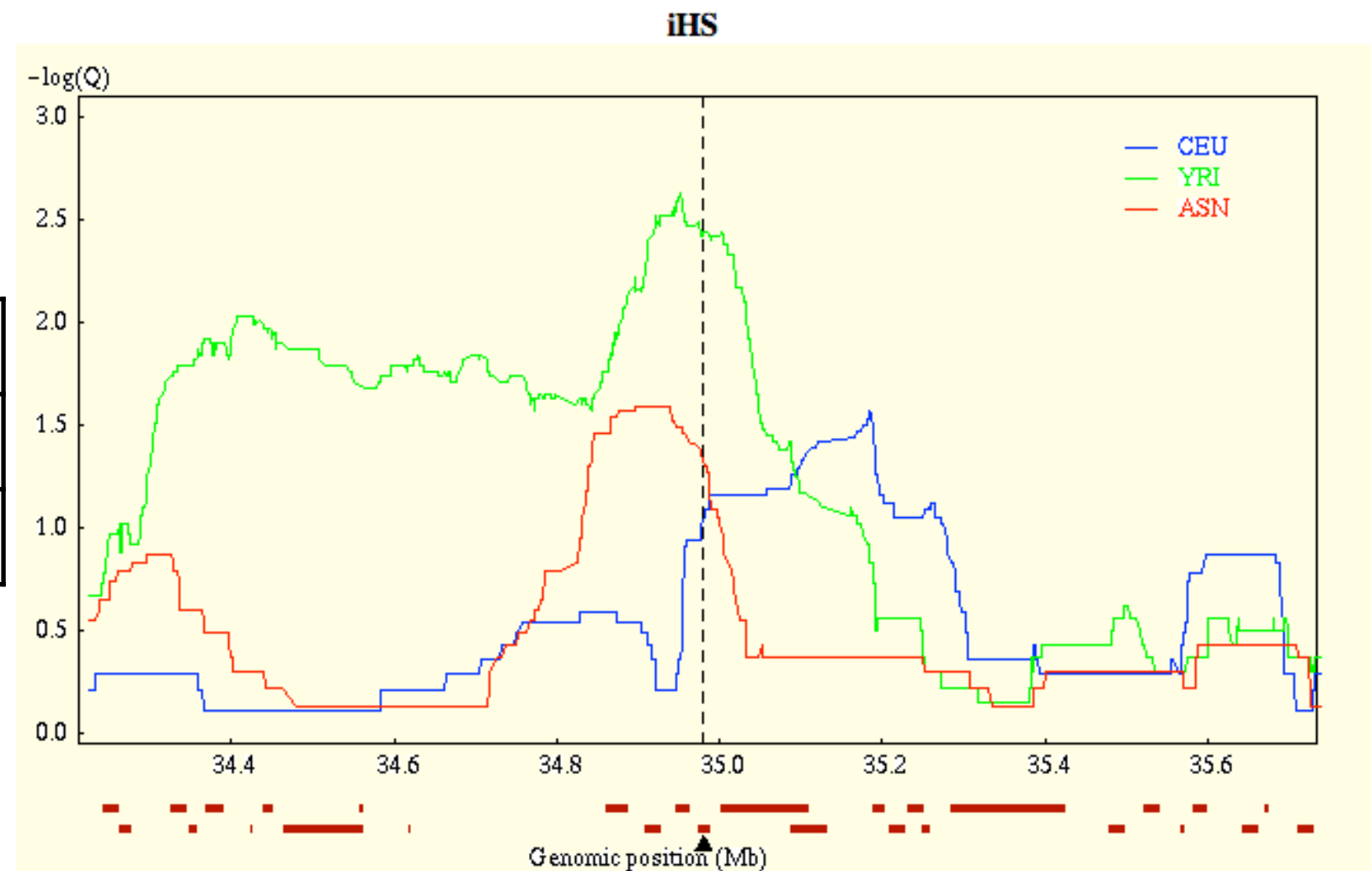
Case Study: Kidney Disease in African Americans

- 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

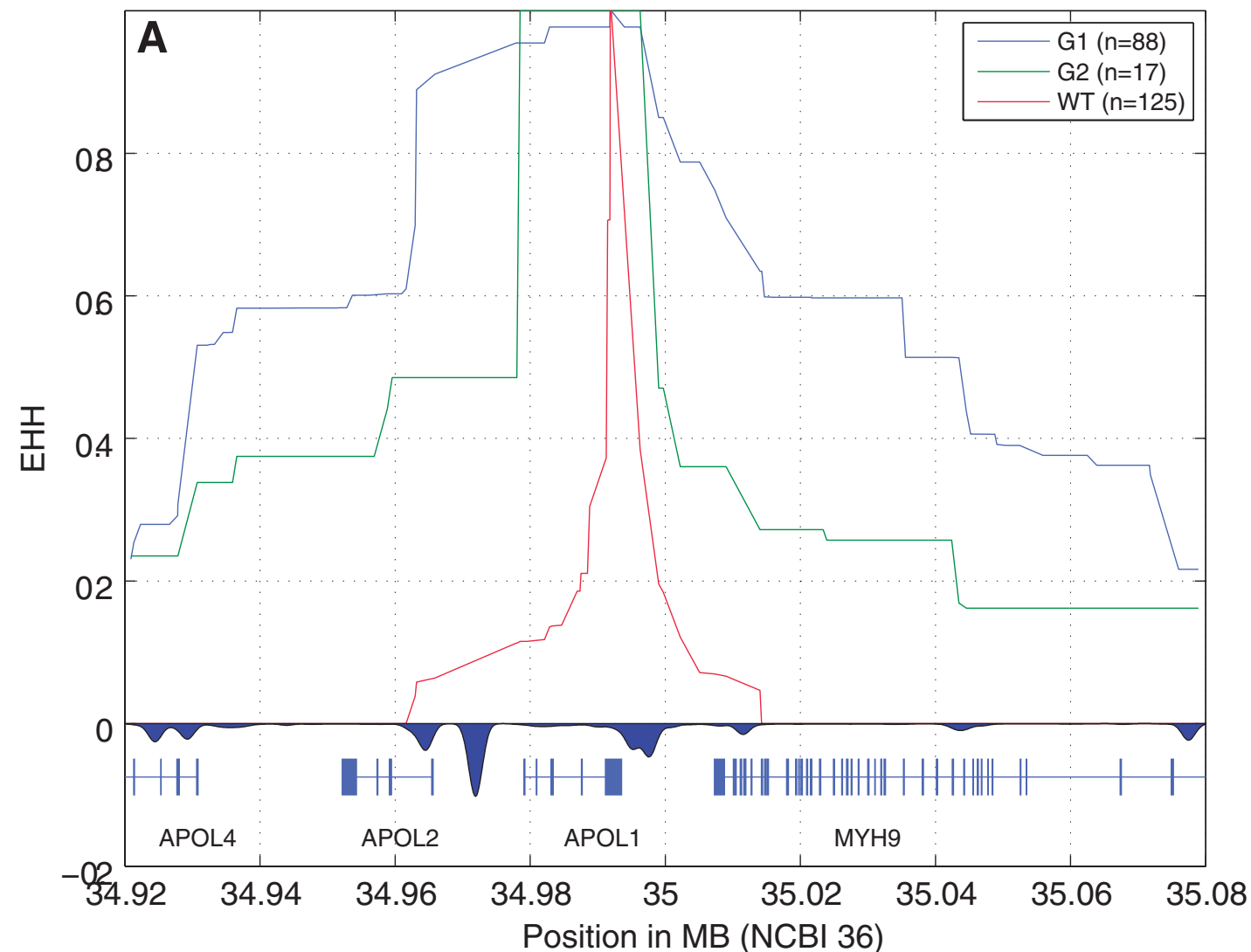
- Looking at signatures of selection adds valuable insight.
- Consider iHS from haplotter.uchicago.edu (more on this later):

Gene	iHS p-value
APOLI	0.0033
MYH9	0.014



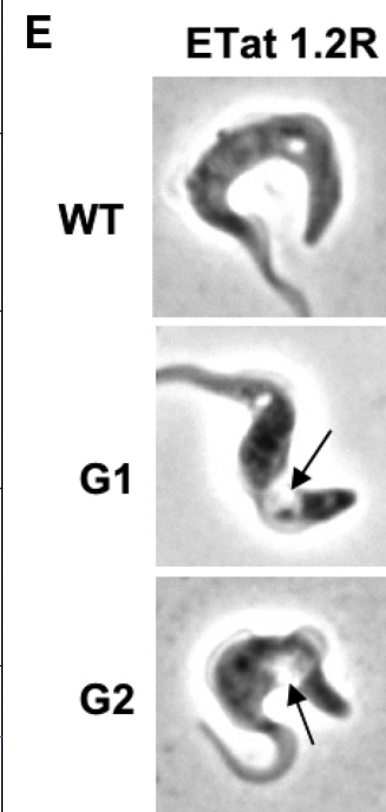
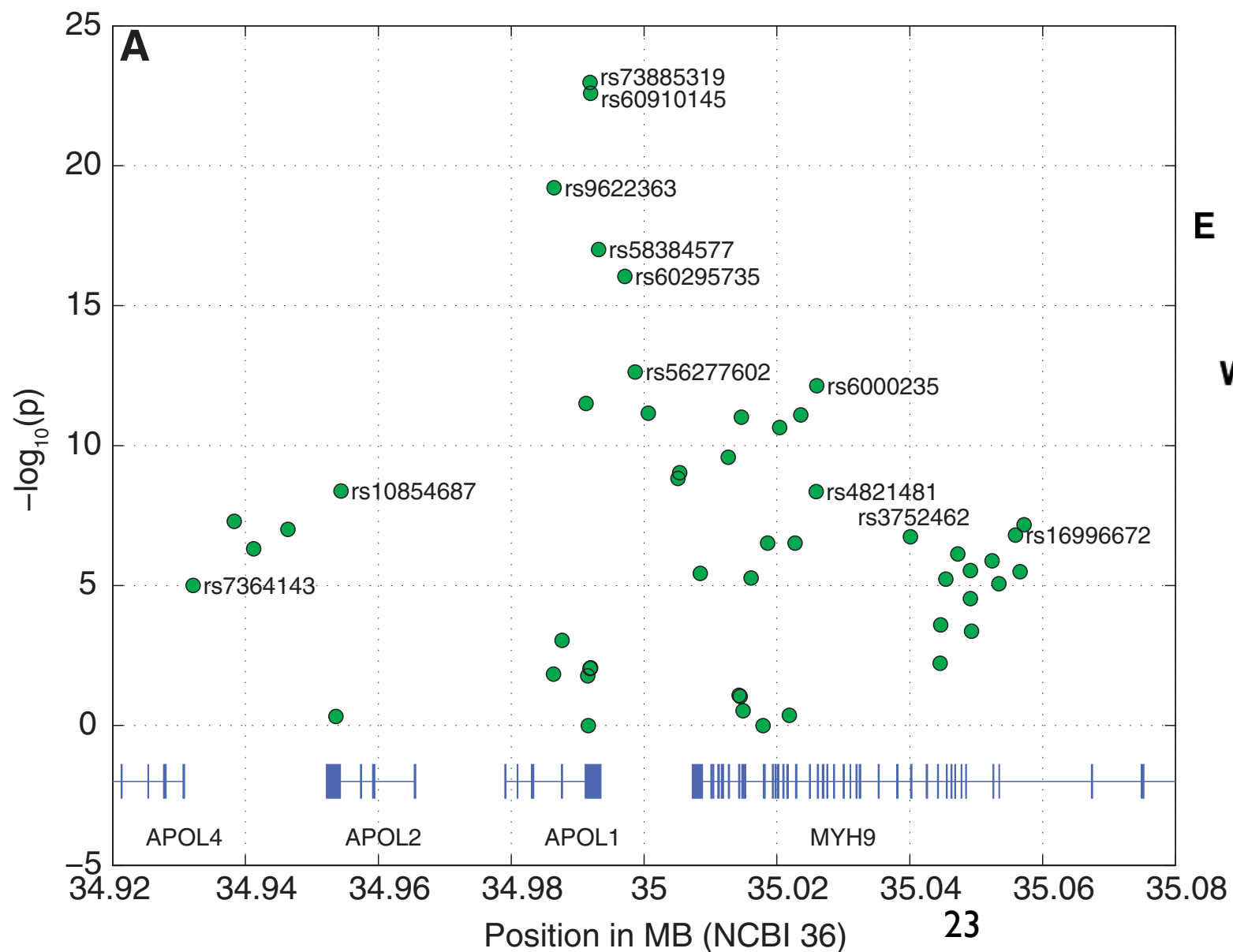
Case Study: Kidney Disease in African Americans

- 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 APOLI!



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

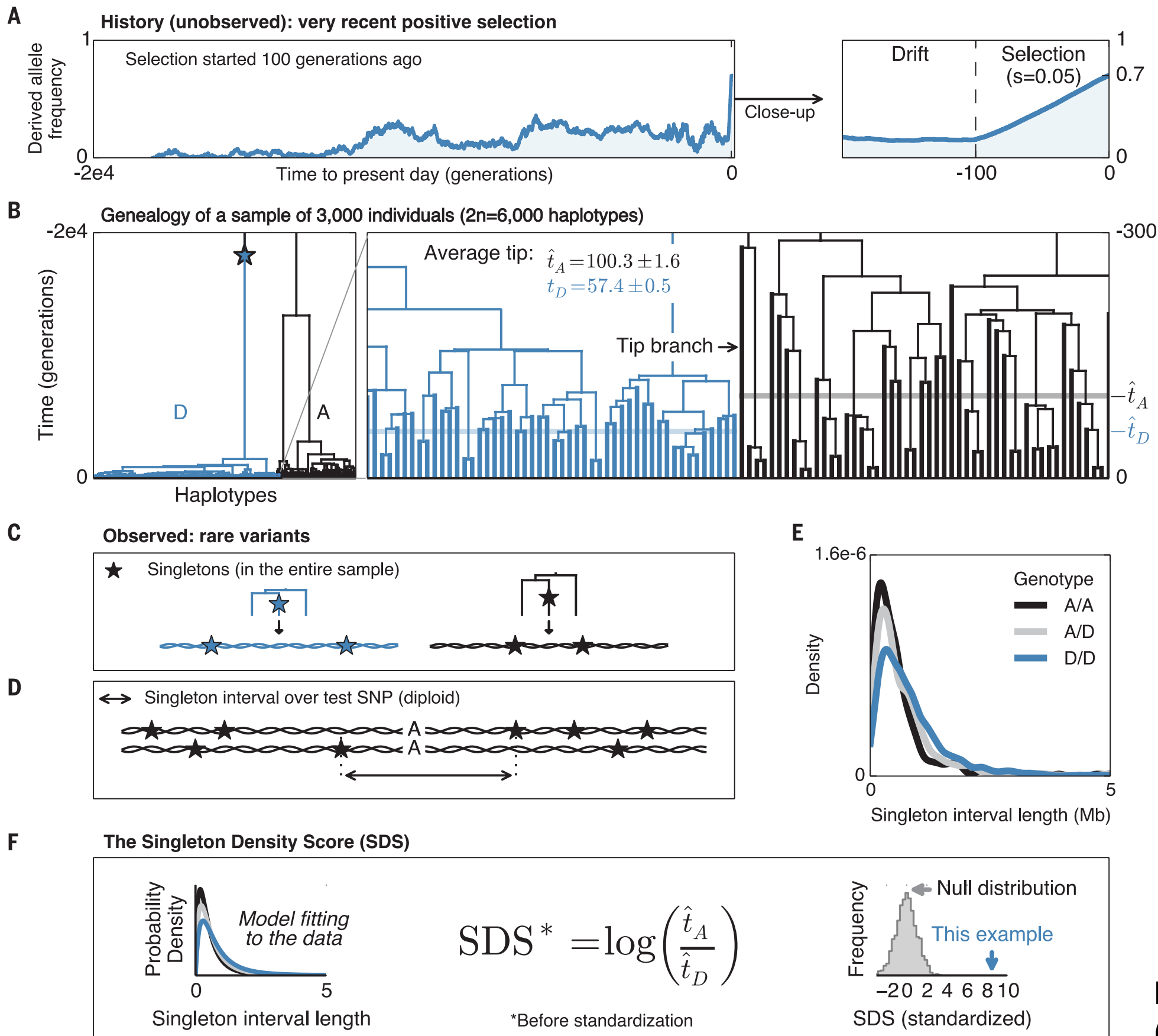
WGS

- 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

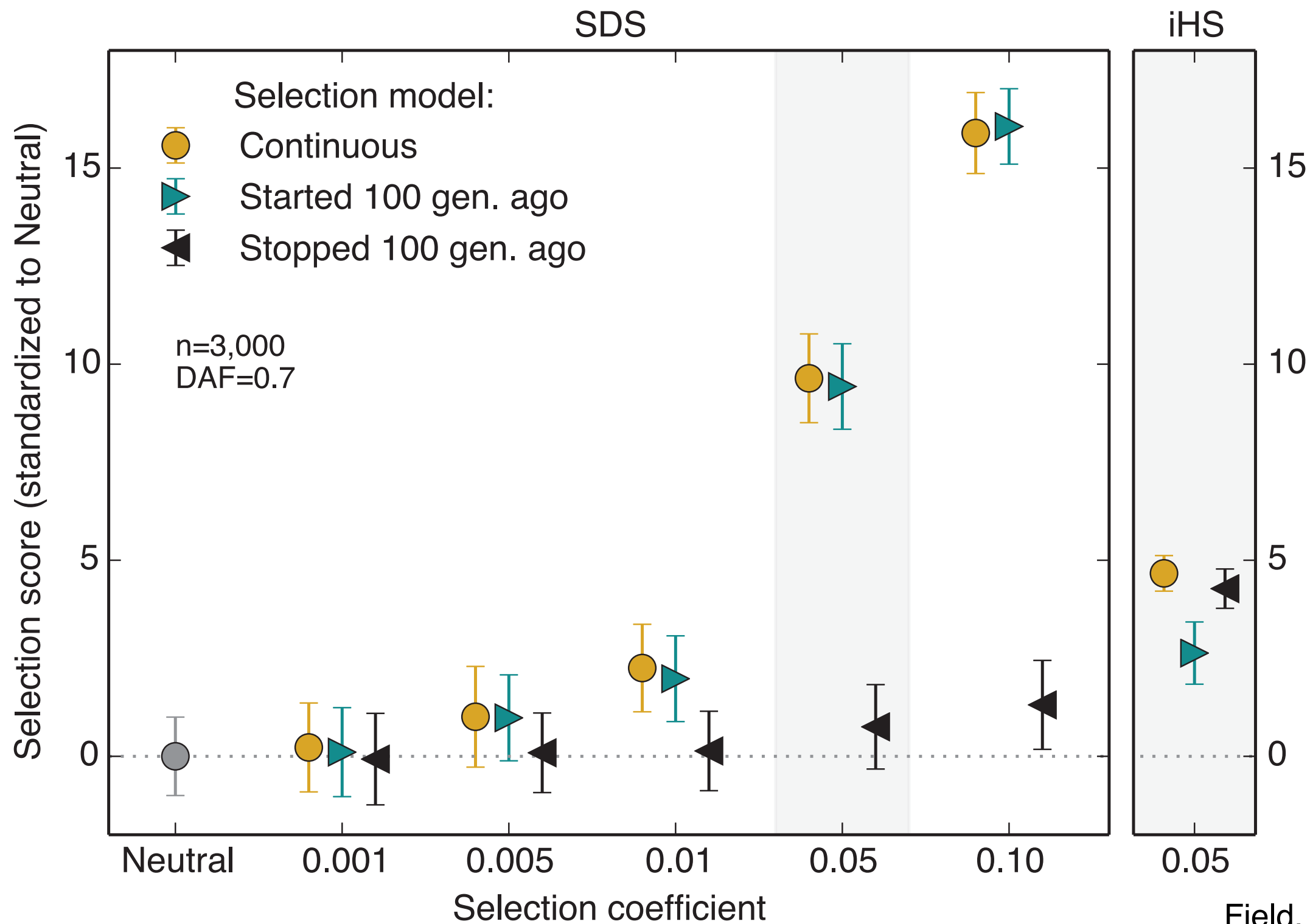


Field, et al.
(Science, 2016)

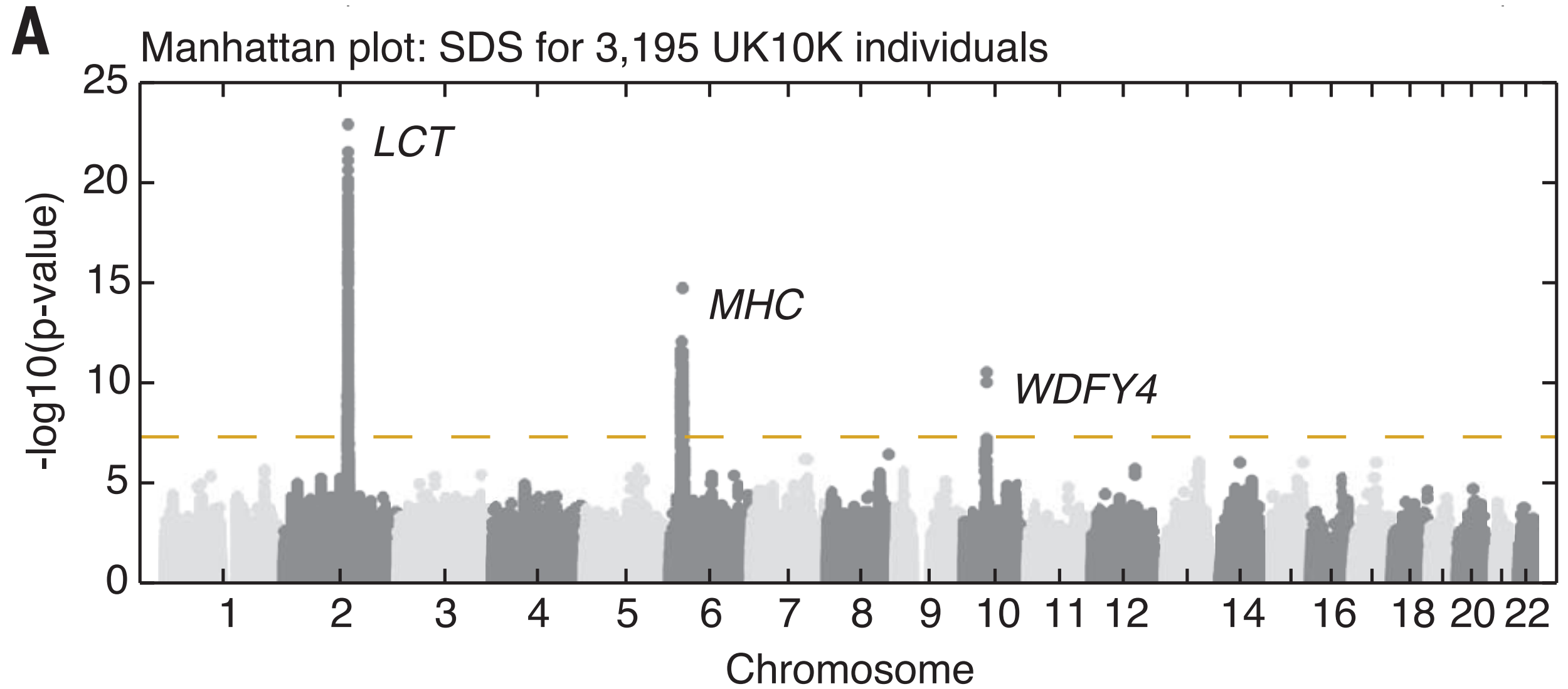
SDS

B

Simulations: signal and specificity of our method to recent history



SDS



Conclusions

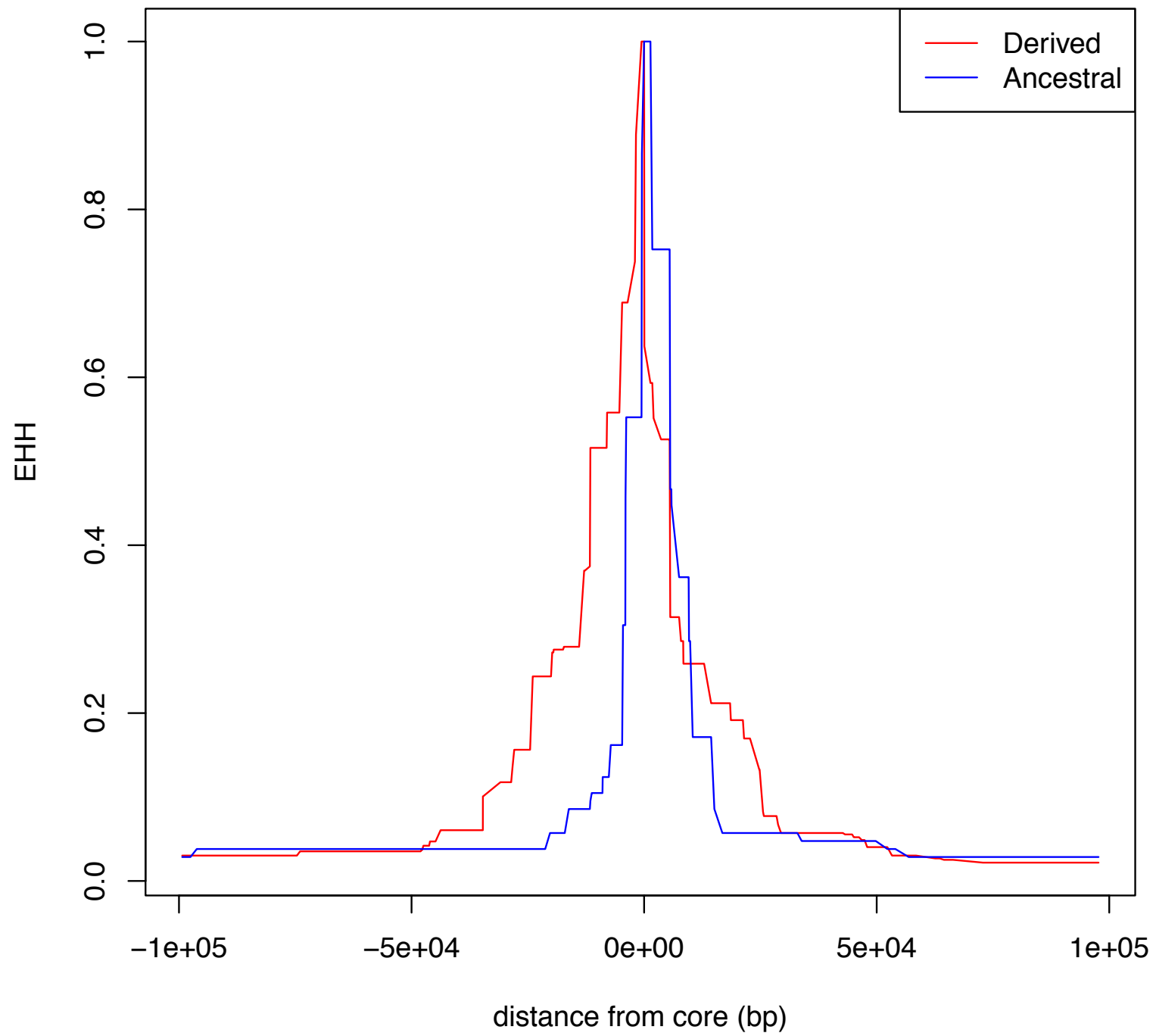
- Natural selection leaves distinctive footprints within patterns of genetic variation.
- This occurs because alleles driven by natural selection tend to be younger than neutral alleles at the same frequency.
- Characterizing signatures of natural selection around disease associated loci can sometimes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

Extended Haplotype Homozygosity

- Sabeti, et al. (*Nature*, 2002) proposed EHH
- Designed to track the decay of haplotype identity away from a locus of interest
- If selection acts quickly enough
- Originally derives from ideas in Hudson, et al. (*Genetics*, 1994).

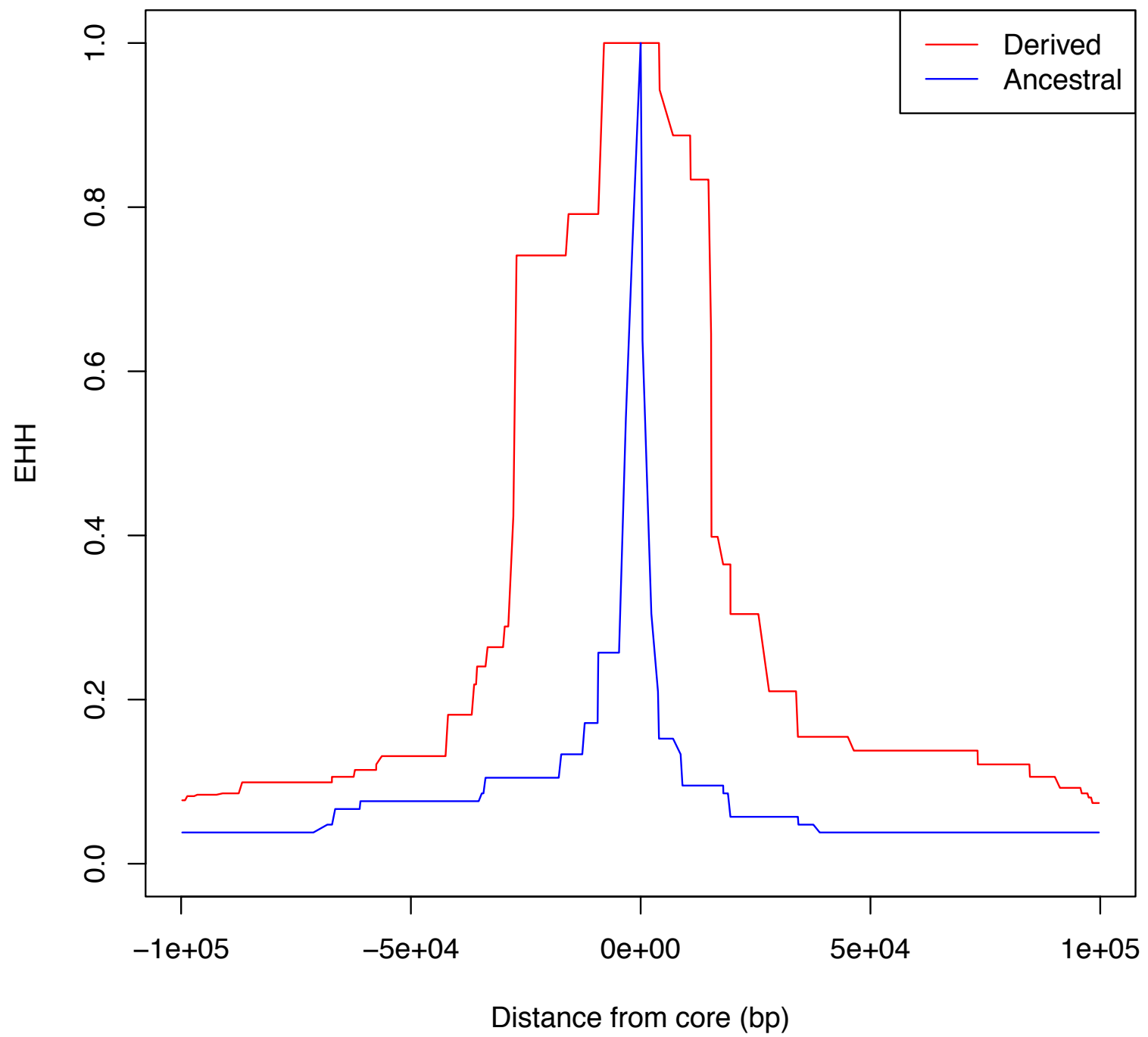
EHH

$s = 0.01, N_e = 10,000$



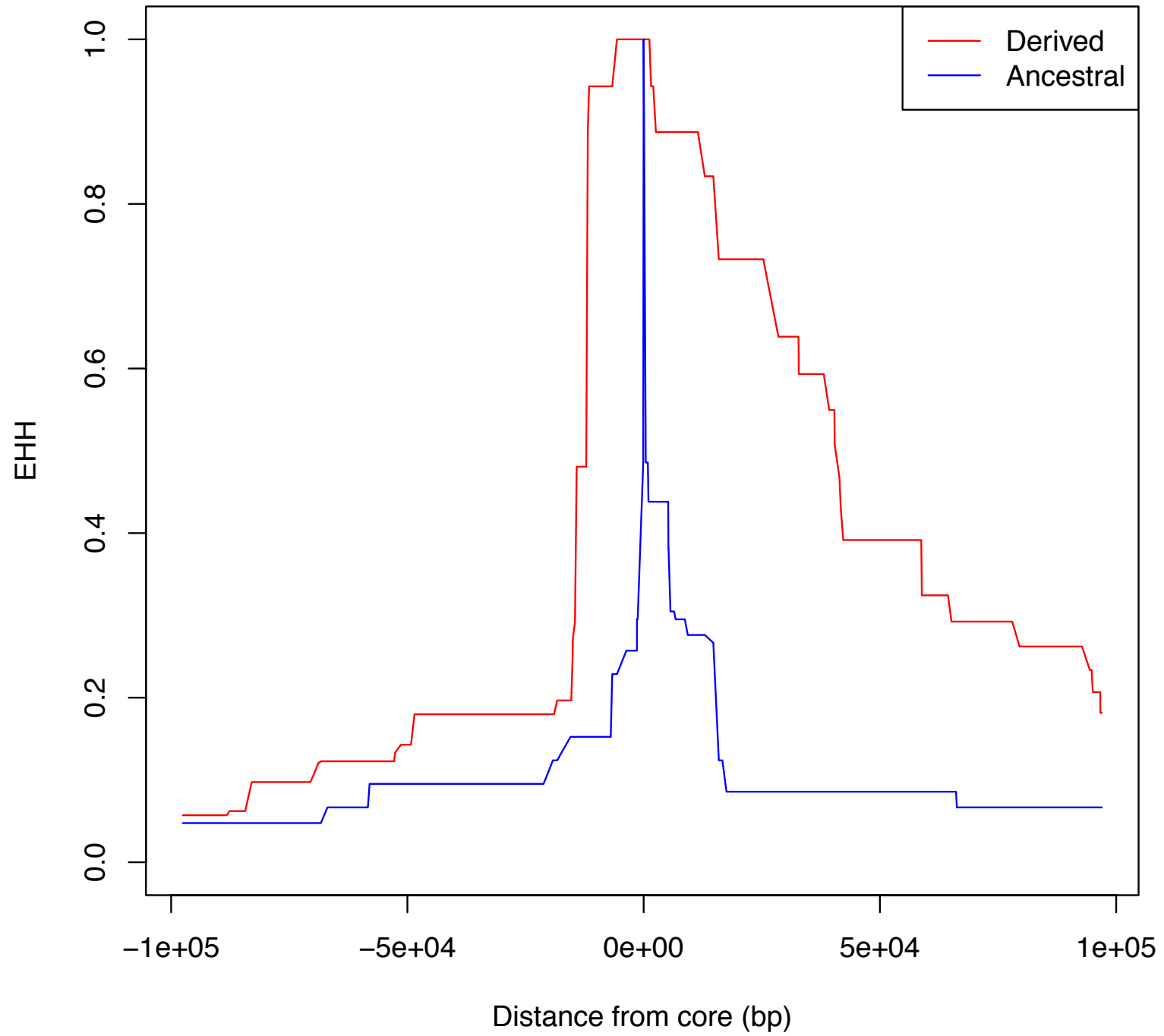
EHH

$s = 0.02, N_e = 10,000$



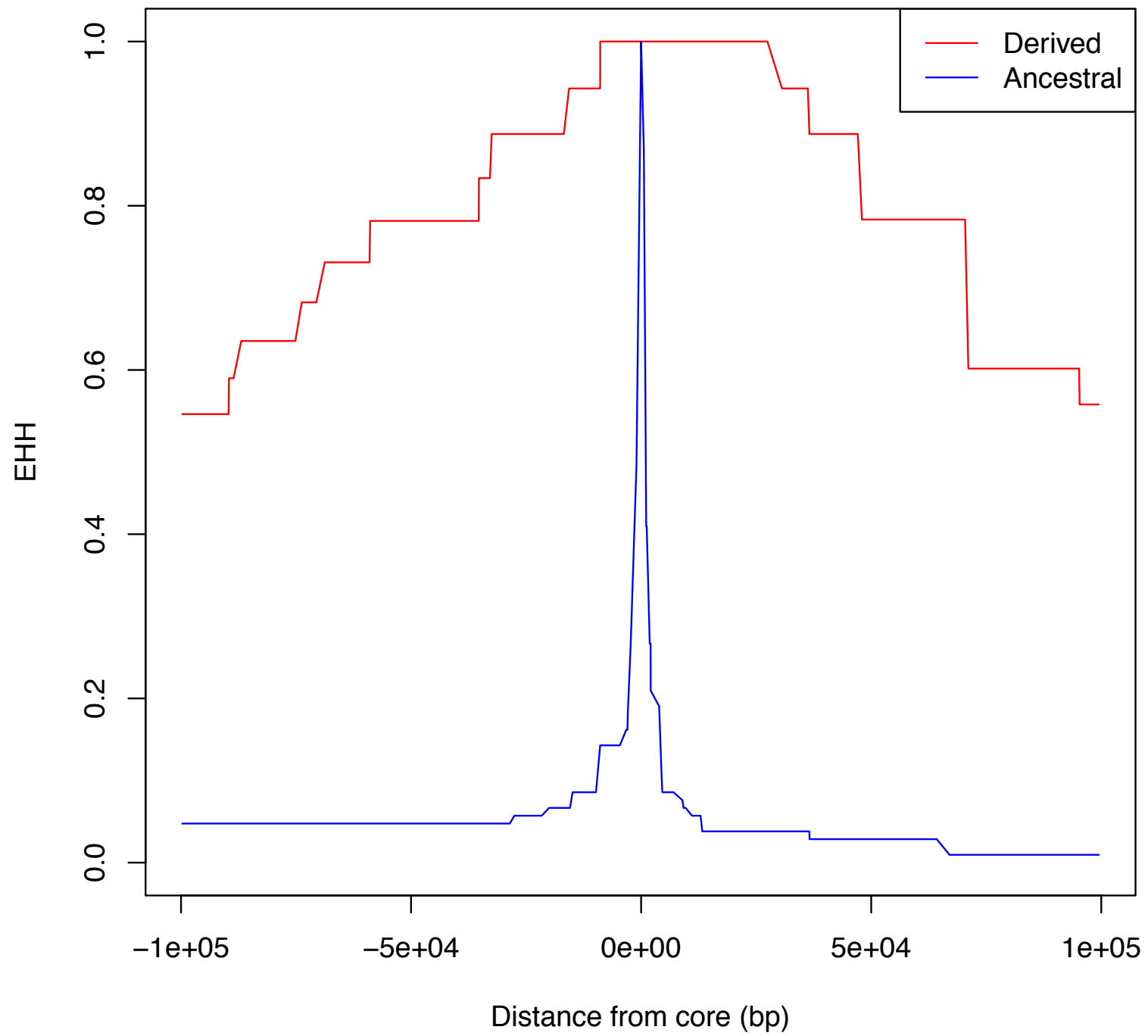
EHH

$s = 0.05, N_e = 10,000$



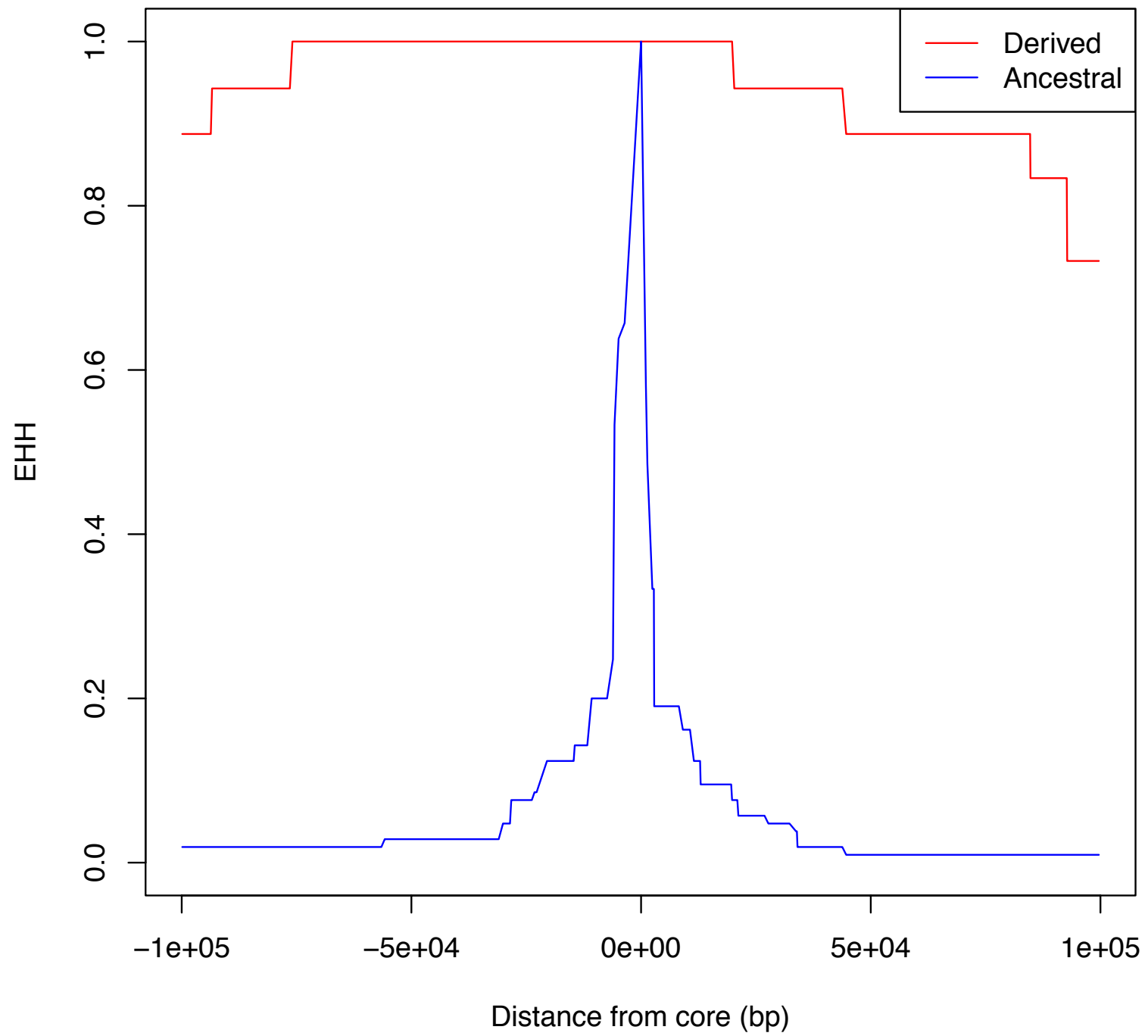
EHH

$s = 0.10, N_e = 10,000$



EHH

$s = 0.50, N_e = 10,000$



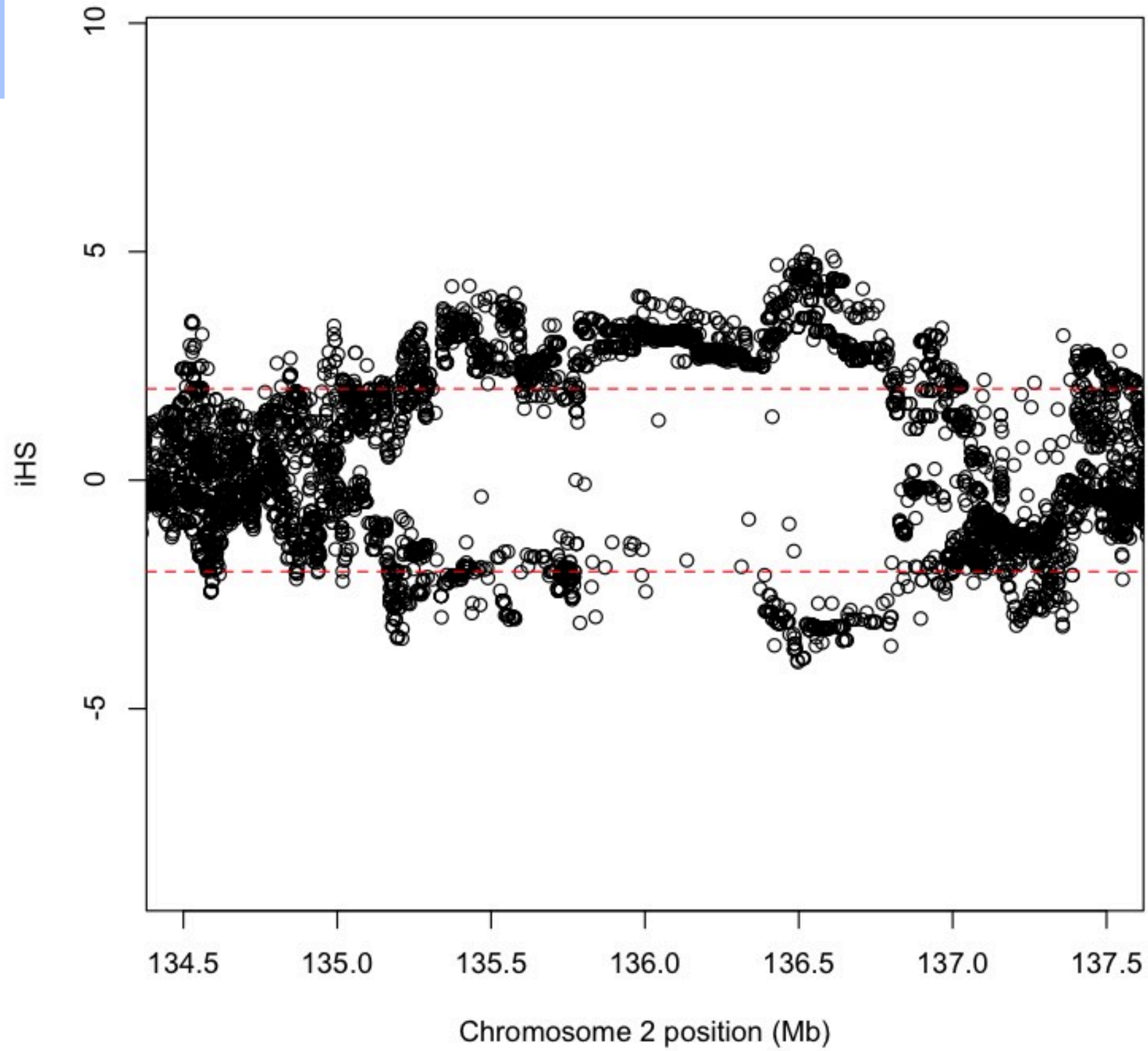
EHH

- When querying a specific region of the genome, for each core haplotype, calculate EHH for successively longer surrounding haplotypes.
- Statistical significance is determined by comparing EHH scores to neutral simulations and random control regions of the genome.

Genome-wide scans

- The EHH approach does not lend itself to a genome-wide scan.
- Voight, et al. (2006) create a genome-wide scan statistic based on EHH called integrated Haplotype Score (iHS).

CEU TGP Phase 3, lactase (LCT) region



Caveats

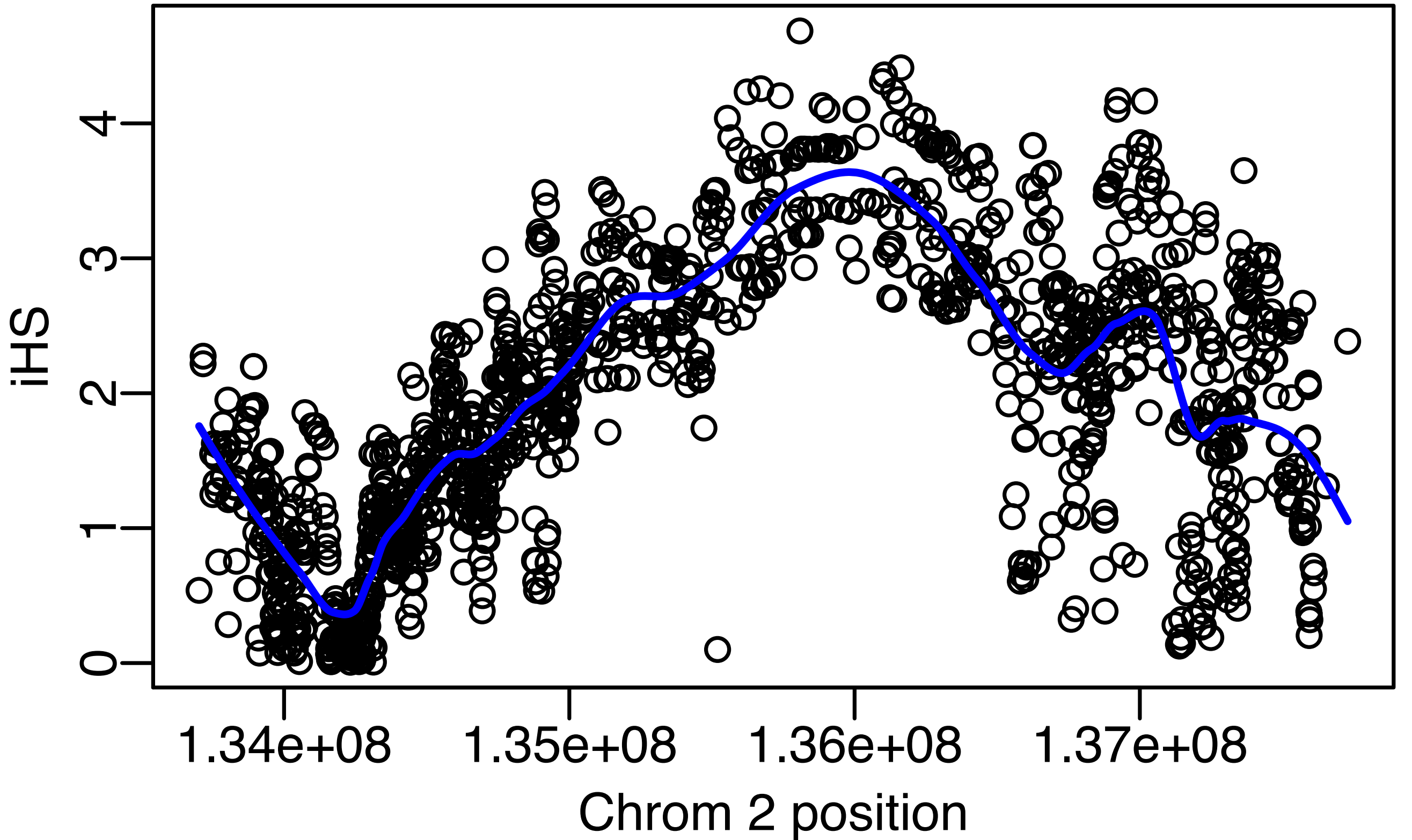
- Power may be overstated.
 - If a large proportion of the genome is non-neutral, we lose power to detect the weakest selected variants because of genome-wide normalization.
- iHS no formal test to decide significance.
 - Take top 1% of signals
- XP-EHH more sensitive to demographics
 - i.e. comparing populations with serial bottlenecks separating them
- Important to combine *multiple lines* of evidence!

Breakout Groups!!

Running iHS with `selscan`

- Open up your command prompt (i.e., rev your engines)
- Let's give iHS a go!
- Let's consider the LCT gene.
- Follow commands in `Day3_AM2_CMD.txt`
 - You will need `selscan.zip`
 - In terminal run:
 - `selscan ...`
 - `norm ...`
 - Plot it in R!

iHS



Other populations??

- Now run selscan on the YRI population
- YRI is a sample of individuals from Yoruba, Nigeria, where they do not have a long tradition of domesticating cows.
- Update the selscan commands by replacing “CEU” with “YRI”
- Breakout Groups!!
 -

Do you think there is selection in this region in the YRI population?

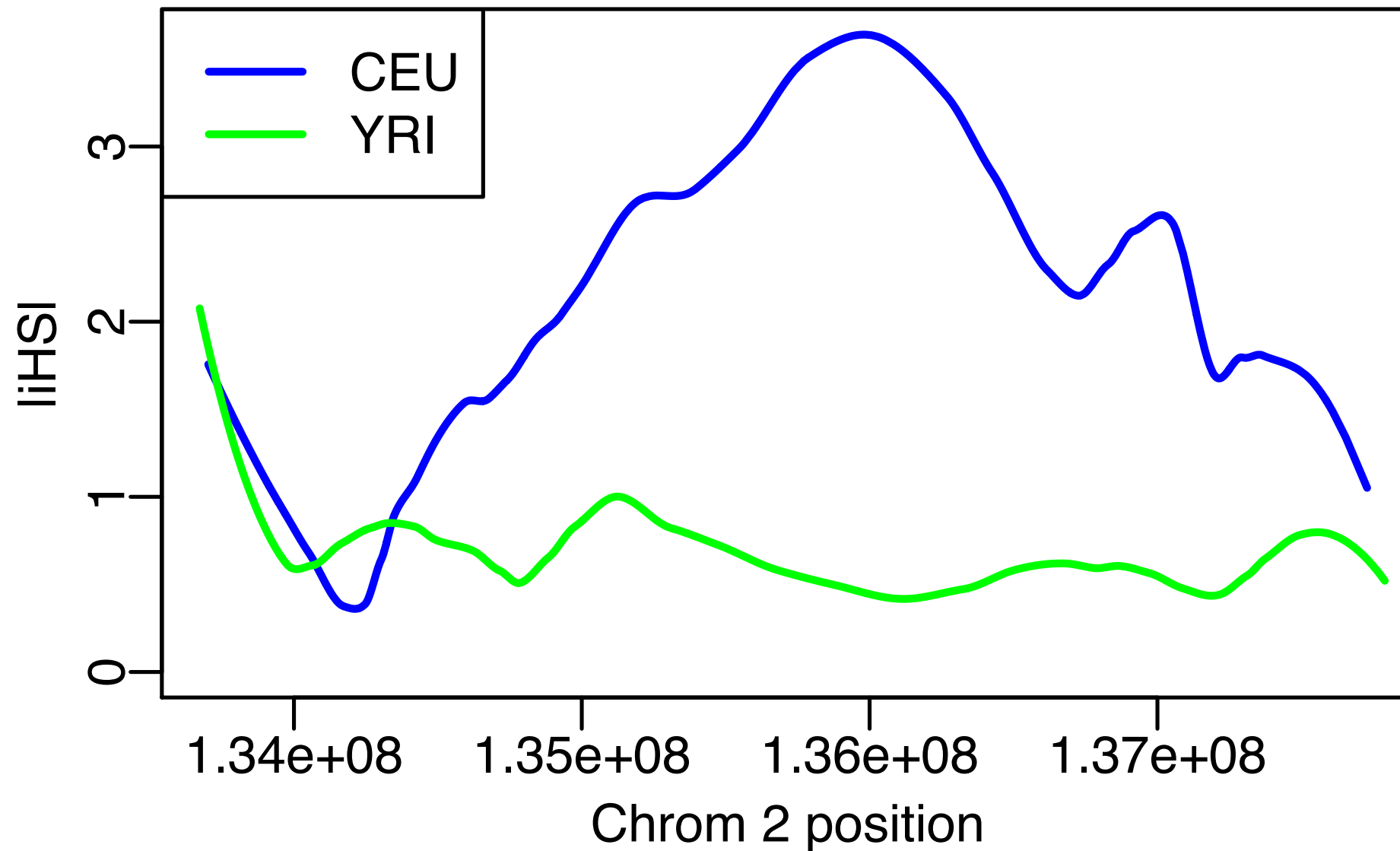
Yes!

No!

Unclear

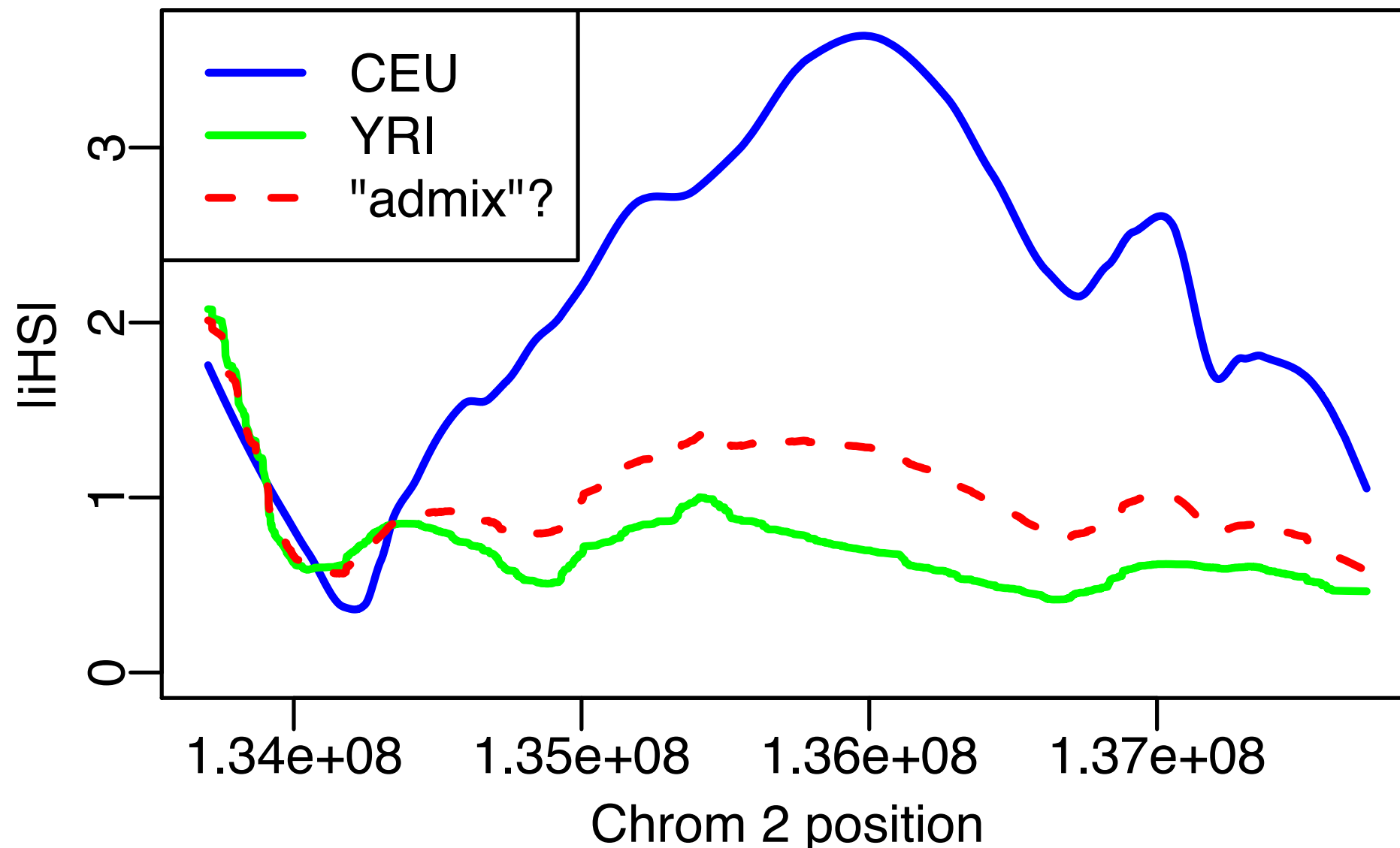
Other populations??

- “CEU” vs “YRI”



What about admixture?

- African American genomes contain admixture with African ancestry (~80%) and European ancestry (~20%).
- ASW is one sample of African Americans (from the Southwest)
- One guess might be that it should be intermediate

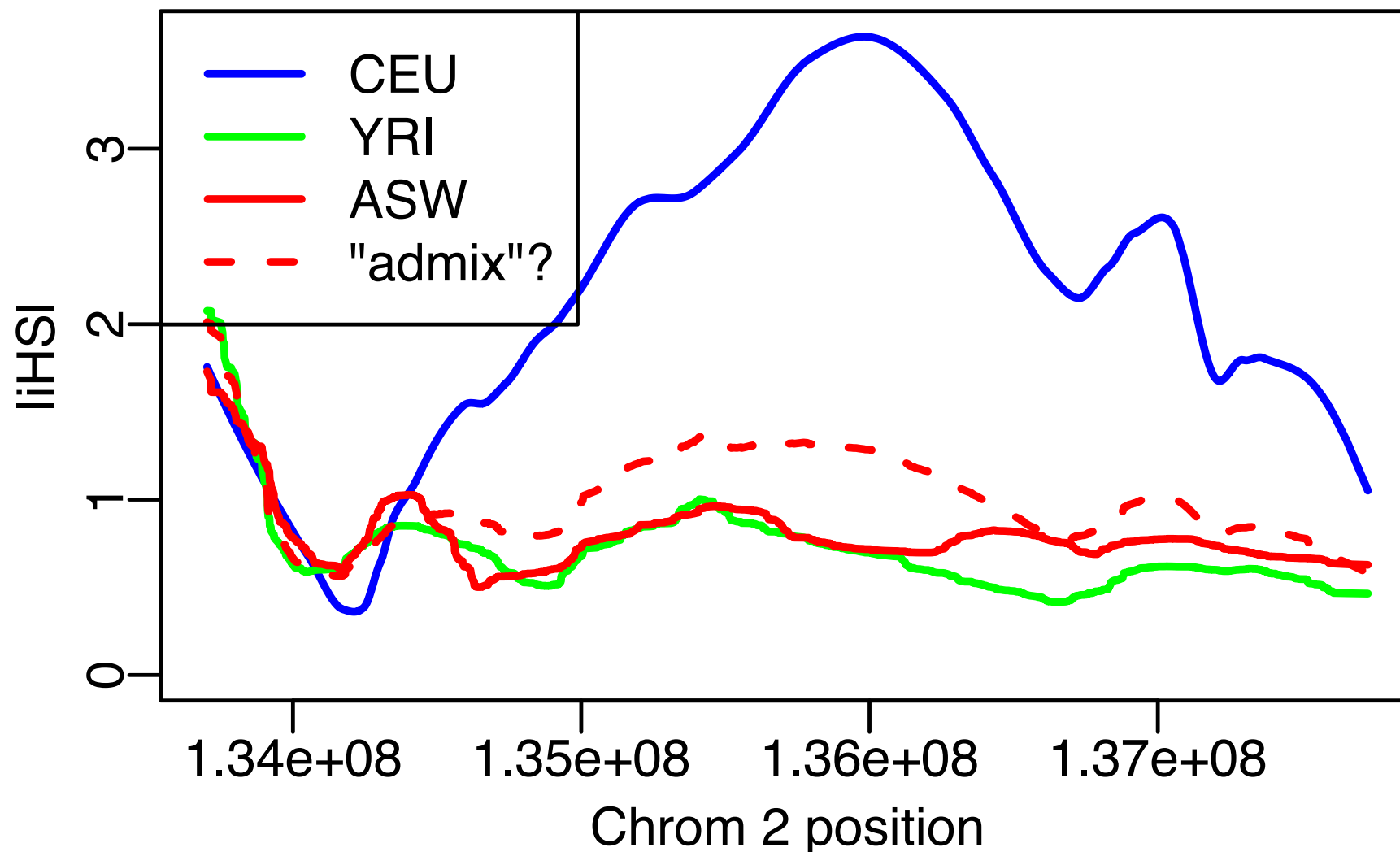


Other populations??

- Now run selscan on the ASW population
- Update the selscan command by replacing “CEU” with “ASW”
- Breakout groups!!

Other populations??

- Now run selscan on the ASW population
- Update the selscan command by replacing “CEU” with “ASW”
- In these data, ASW is much more similar to YRI than “expected”.



Summary

- iHS is one example of a statistic geared toward detecting a “classic sweep”.
- It is based on the idea that a new mutation has been selected, and quickly spread through the population.
- selscan is one piece of software that can run many different selection statistics in an efficient manner.

selscan

- Open your terminal/command prompt!
- Change to the new selscan directory
- For example:
 - `cd ~/Desktop/selscan/`
- There should 4 subdirectories:
 - `rhernandez$ ls`
`data linux osx win`
- Change Directory to where the data are:
 - `cd data`

selscan

- All the commands we are running can be found in the `selscan_CMD.txt` file.
- Copy the appropriate executable to the data directory:
- **osx:**
 - `cp ../osx/selscan .`
- **linux:**
 - `cp ../linux/selscan .`
- **Windows:**
 - `copy ..\win\selscan.exe .`

selscan

- Test that it works:
 - **osx/linux:** ./selscan **(Win: selscan.exe)**
selscan v1.1.0b
ERROR: Must specify one and only one of
EHH (-ehh)
iHS (--ihs)
XP-EHH (--xpehh)
PI (--pi)
nSL (--nsl)

selscan

- iHS requires 2 files, a **map** file and a **hap** file.
- `--map <string>`: A mapfile with one row per variant site.
 - Formatted with 4 columns:
 - `<chr#> <locusID> <genetic pos>`
`<physical pos>`
- `--hap <string>`: A hapfile with one row per haplotype, and one column per variant. Variants should be coded 0/1.

selscan

- Now run it!
- All in one line type:
 - `./selscan` (Win: `selscan.exe`)
`--ihs`
`--map CEU.chr2.map`
`--hap CEU.chr2.ihshap`
`--out CEU.chr2`

```
selscan v1.1.0b
Opening ../data/CEU.chr2.hap...
Loading 224 haplotypes and 1971 loci...
Opening ../data/CEU.chr2.map...
Loading map data for 1971 loci
--skip-low-freq set. Removing all variants < 0.05.
Removed 359 low frequency variants.
Starting iHS calculations with alt flag not set.
|=====>|
```

Normalize

- All in one line type:

- `./norm`

- `--ihs`

- `--files CEU.chr2.ihs.out bg.ihs.out`

```
norm v1.1.0aYou have provided 2 output files for joint
normalization.
```

```
Opened ../data/CEU.chr2.ihs.out
```

```
Opened ../data/bg.ihs.out
```

```
Total loci: 666285
```

```
Reading all frequency and iHS data.
```

```
Calculating mean and variance per frequency bin:
```

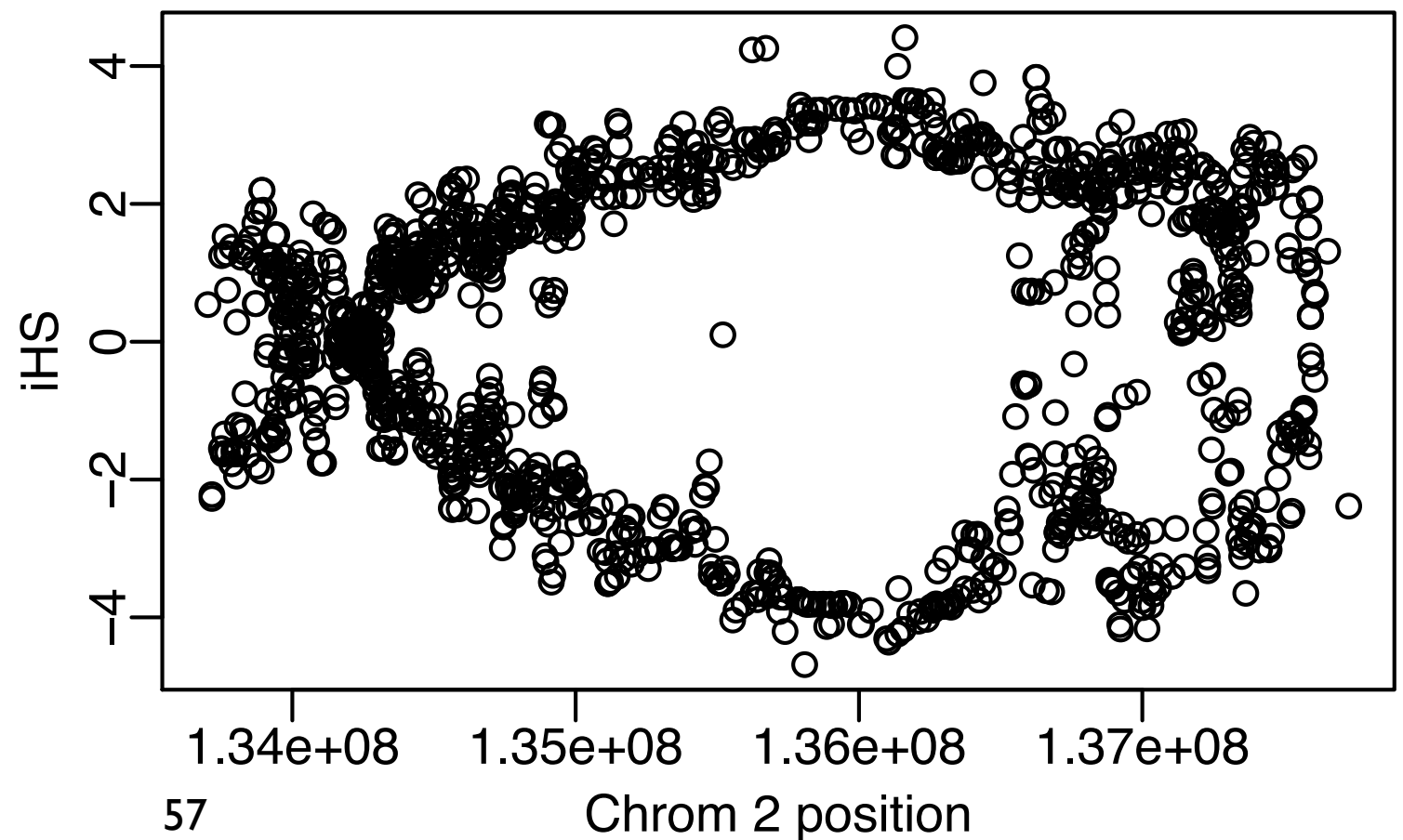
iHS

- Now let's plot it!
- Open R.
- Read in data for CEU:

```
setwd("cd ~/Desktop/selscan/data")
```

```
CEU=read.table("CEU.chr2.ihs.out.100bins.norm")
```

```
plot(CEU[,2], CEU[,7])
```



iHS

- Often analyze absolute value, and smooth it out.
- My preferred method for smoothing is using loess

```
SP=0.2 #this is the span, a parameter you can change (higher = more smoothing)
```

```
CEU.x=CEU[,2]; #the x-coordinates in Mb
```

```
y=abs(CEU[,7]) #iHS is actually the absolute value
```

```
CEU.loess=loess(y~CEU.x,span=SP,data.frame(x=CEU.x,y=y)); #step 1
```

```
CEU.predict=predict(CEU.loess,data.frame(x=CEU.x)); #step 2
```

```
plot(CEU[,2], abs(CEU[,7]))
```

```
lines(CEU.x, CEU.predict, lwd=2, col='blue')
```