

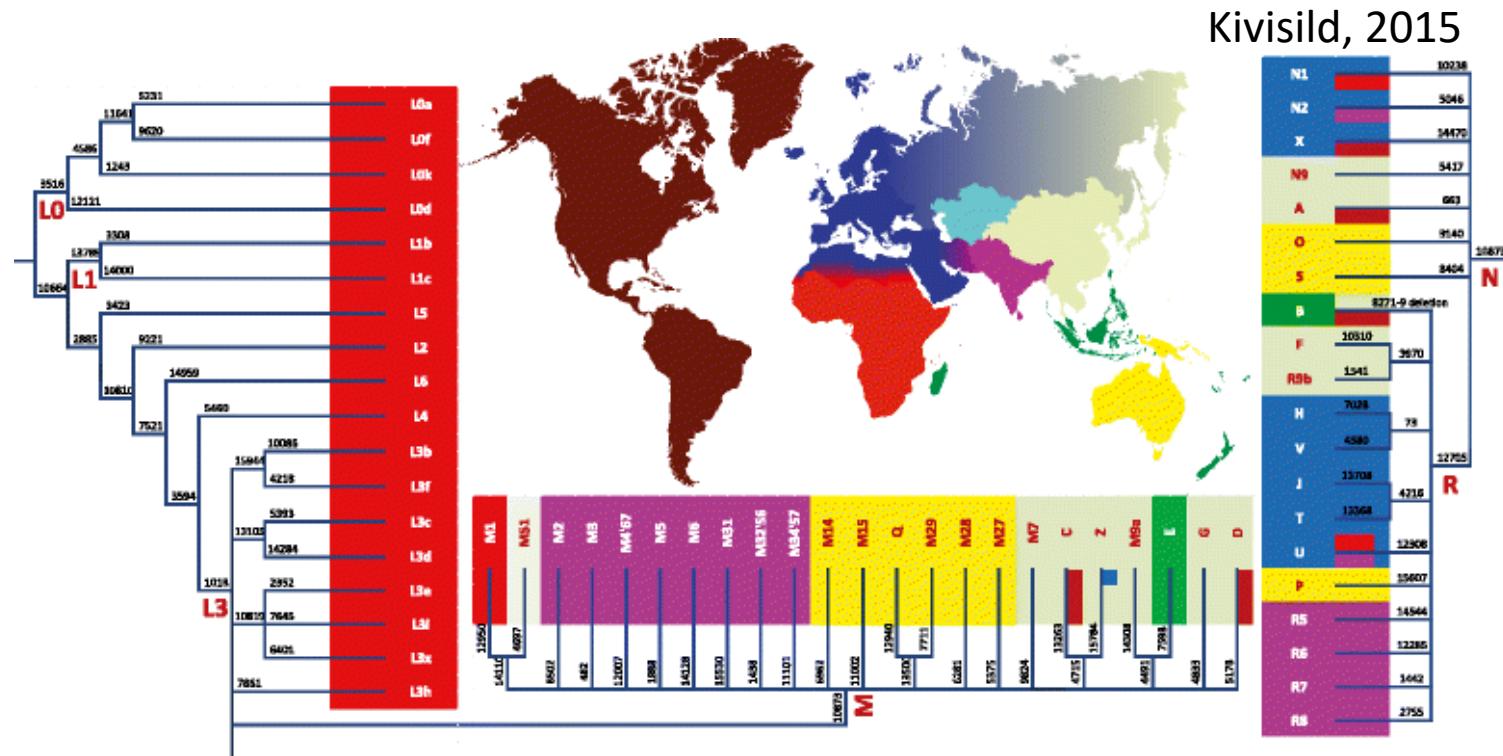
Learning about evolution by building coalescent trees

Leo Speidel

The genetic tree(s) relating humans (or any other species)

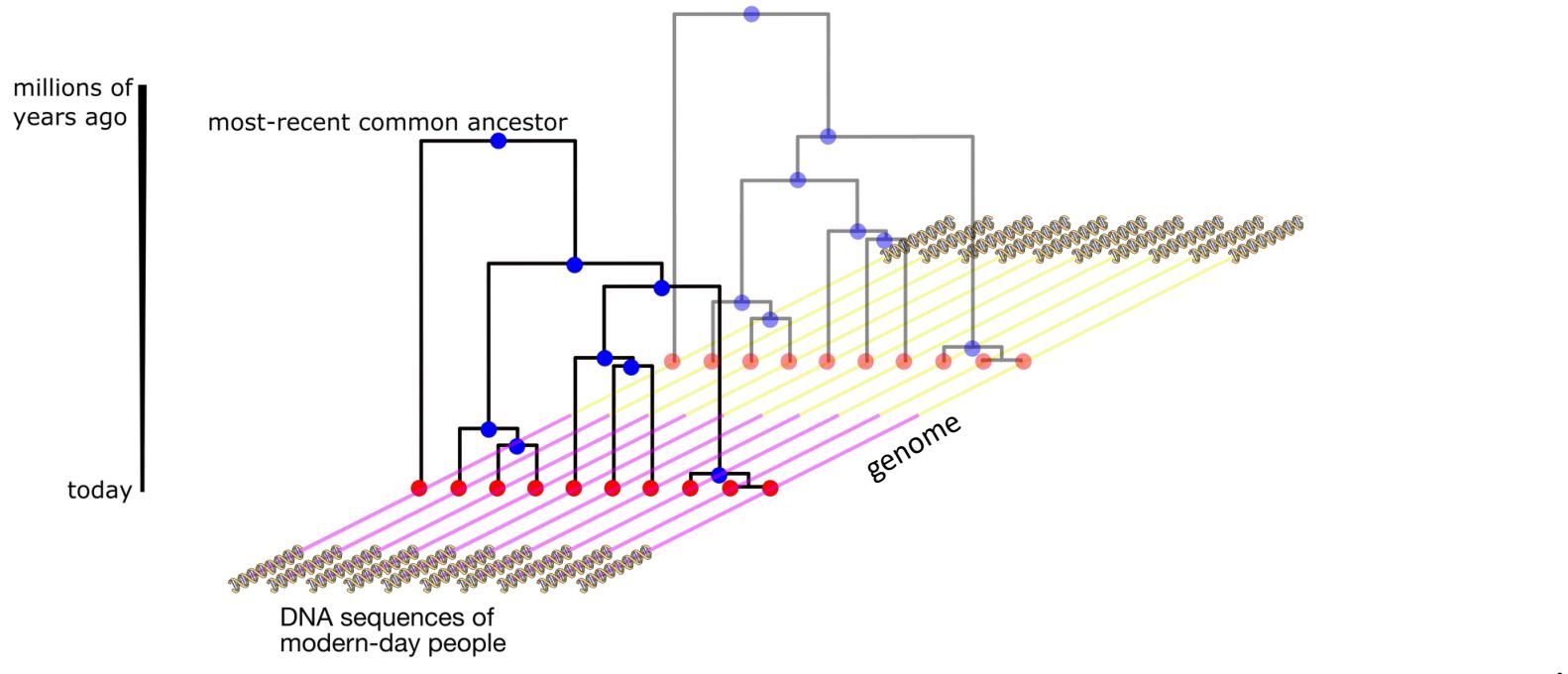
Reconstructing >100,000s years of evolution!

- 22 autosomal chromosomes
 - 2 sex chromosomes
 - X chromosome
 - Y chromosome
 - Mitochondrial genome
- } Different trees in different parts of the genome,
due to recombination
- } 1 tree each (maternal/paternal)



Key concept: Genealogies

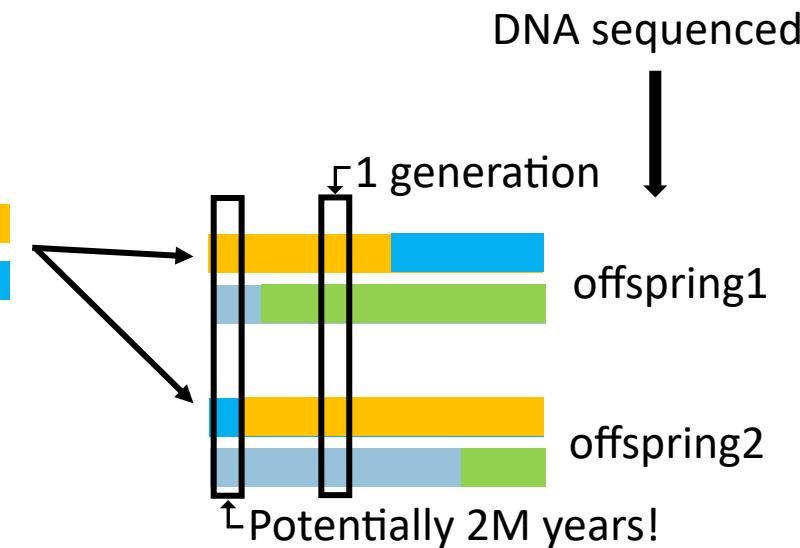
We are genetically related through a sequence of trees



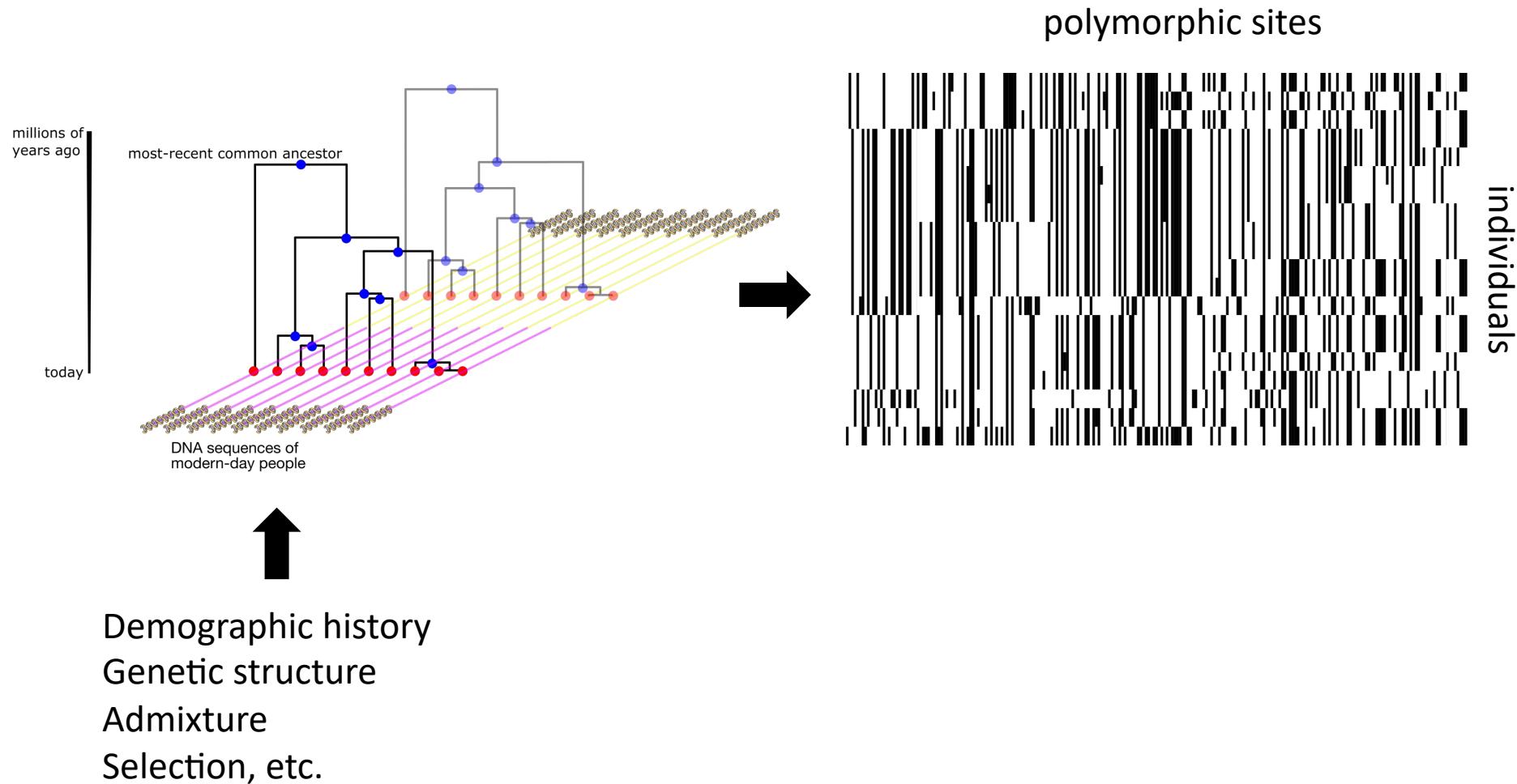
Reason for why trees change along the genome

Recombination:

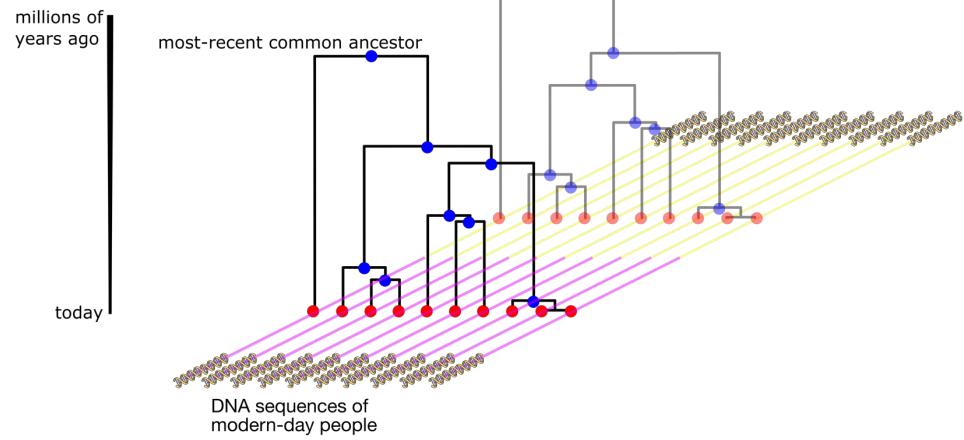
from grandmother
from grandfather



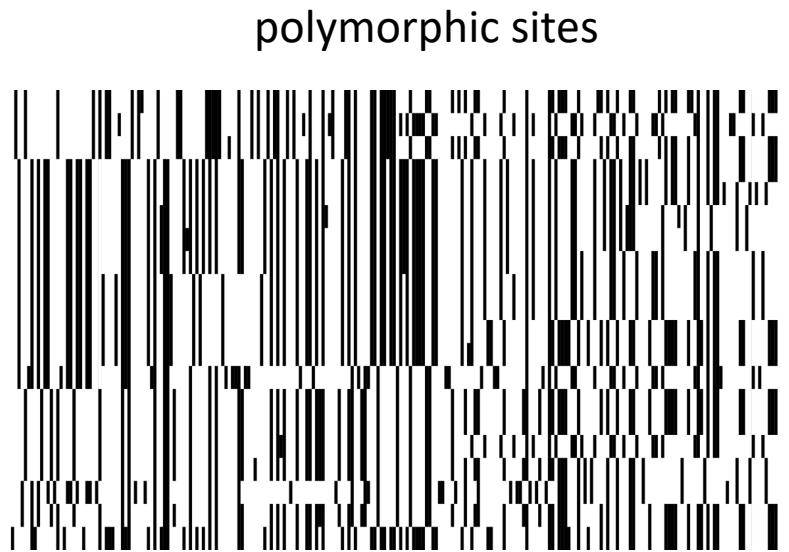
Fundamental forces impact data (only) through underlying genealogies



Many canonical approaches



Demographic history
Genetic structure
Admixture
Selection, etc.

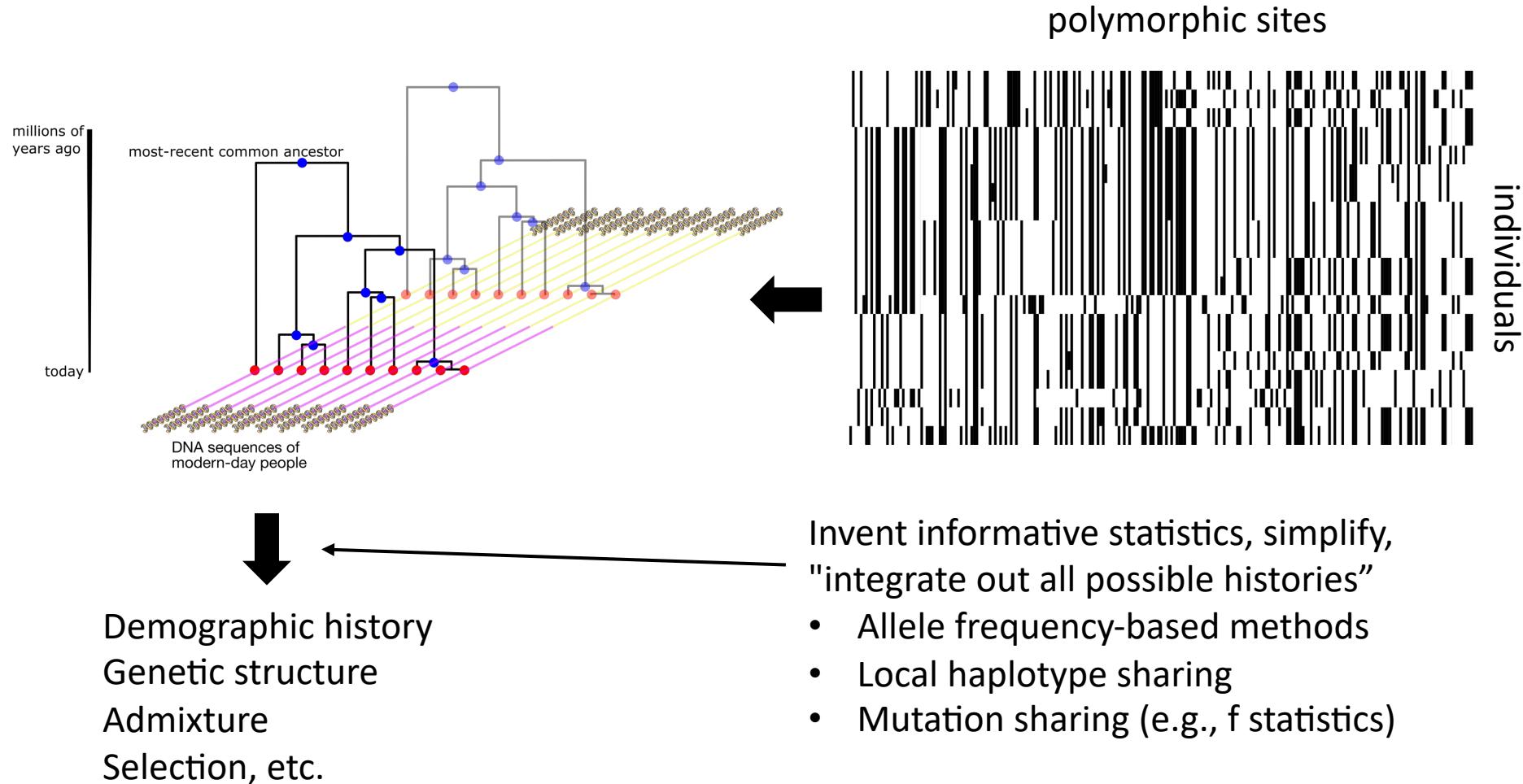


Invent informative statistics, simplify,
"integrate out all possible histories"

- Allele frequency-based methods
- Local haplotype sharing
- Mutation sharing (e.g., f statistics)



Genealogies are the “unobserved link” between evolutionary processes and genetic variation



Challenges: computationally very challenging to sample trees from the data, and modern datasets can contain >50,000 individuals and >100,000,000 mutations

We can now build these trees for many thousands of individuals and add a time dimension onto our genetic data

Old problem, lots of methods, but few can scale:

- ARGweaver
 - Rent+
 - Tsinfer + tsdate
 - ARG-Needle
 - Relate
- } Infers Ancestral Recombination Graphs
- } Published in since 2019,
scale to large sample sizes

We will talk about Relate, but principles of tree-based inference applies more generally!

Relate Home Getting Started Input data Add-on modules Parallelise Relate

Relate

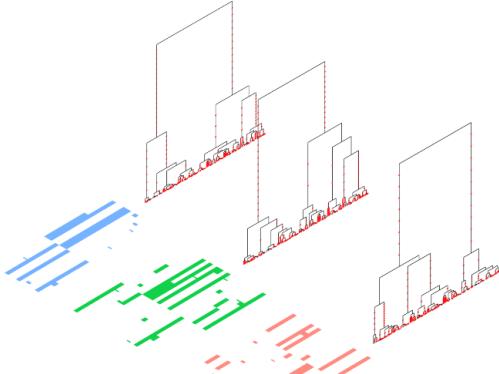
Software to estimate genome-wide genealogies for thousands of samples

Relate estimates genome-wide genealogies in the form of trees that adapt to changes in local ancestry caused by recombination. The method, which is scalable to thousands of samples, is described in the following paper. Please cite this paper if you use our software in your study.

Citations:

- (original Relate paper) Leo Speidel, Marie Forest, Sinan Shi, Simon Myers. A method for estimating genome-wide genealogies for thousands of samples. *Nature Genetics* 51: 1321-1329, 2019.
- (update, v1.1.) Leo Speidel, Lara Cassidy, Robert W. Davies, Garrett Hellenthal, Pontus Skoglund, Simon R. Myers. Inferring population histories for ancient genomes using genome-wide genealogies. *Molecular Biology and Evolution* 38: 3497-3511, 2021.

Contact: leo.speidel@outlook.com
Website: <https://leospeidel.com>



Download

Relate is available for academic use. To see rules for non-academic use, please read the [LICENCE](#) file, which is included with each software download.

Pre-compiled binaries (last updated: 7/11/2021)

I agree with the [terms and conditions](#)

Linux (x86_64, dynamic) - v1.1.8
Linux (x86_64, static) - v1.1.8
Mac OS X (Intel) - v1.1.8
Mac OS X (M1) - v1.1.8

Github repository

Alternatively, you can [compile your own version](#) by downloading the source code from this [github repository](#).

In the downloaded directory, we have included a toy data set. You can try out Relate using this toy data set by following the instructions on our [getting started](#) page.

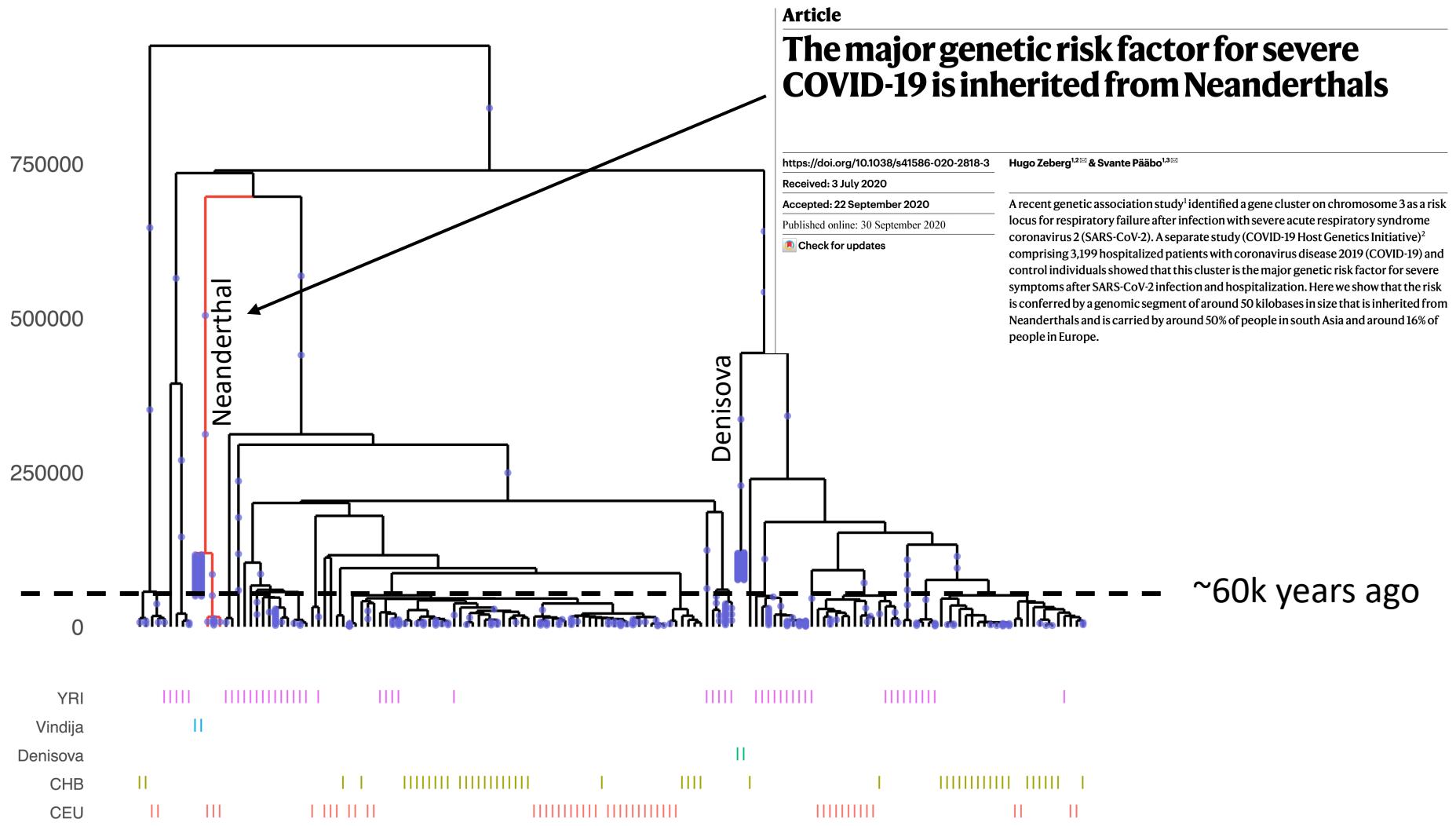
If you have any problems getting the program to work on your machine or would like to request an executable for a platform not shown here, please send a message to leo.speidel [at] outlook [dot] com.

<https://myersgroup.github.io/relate/>

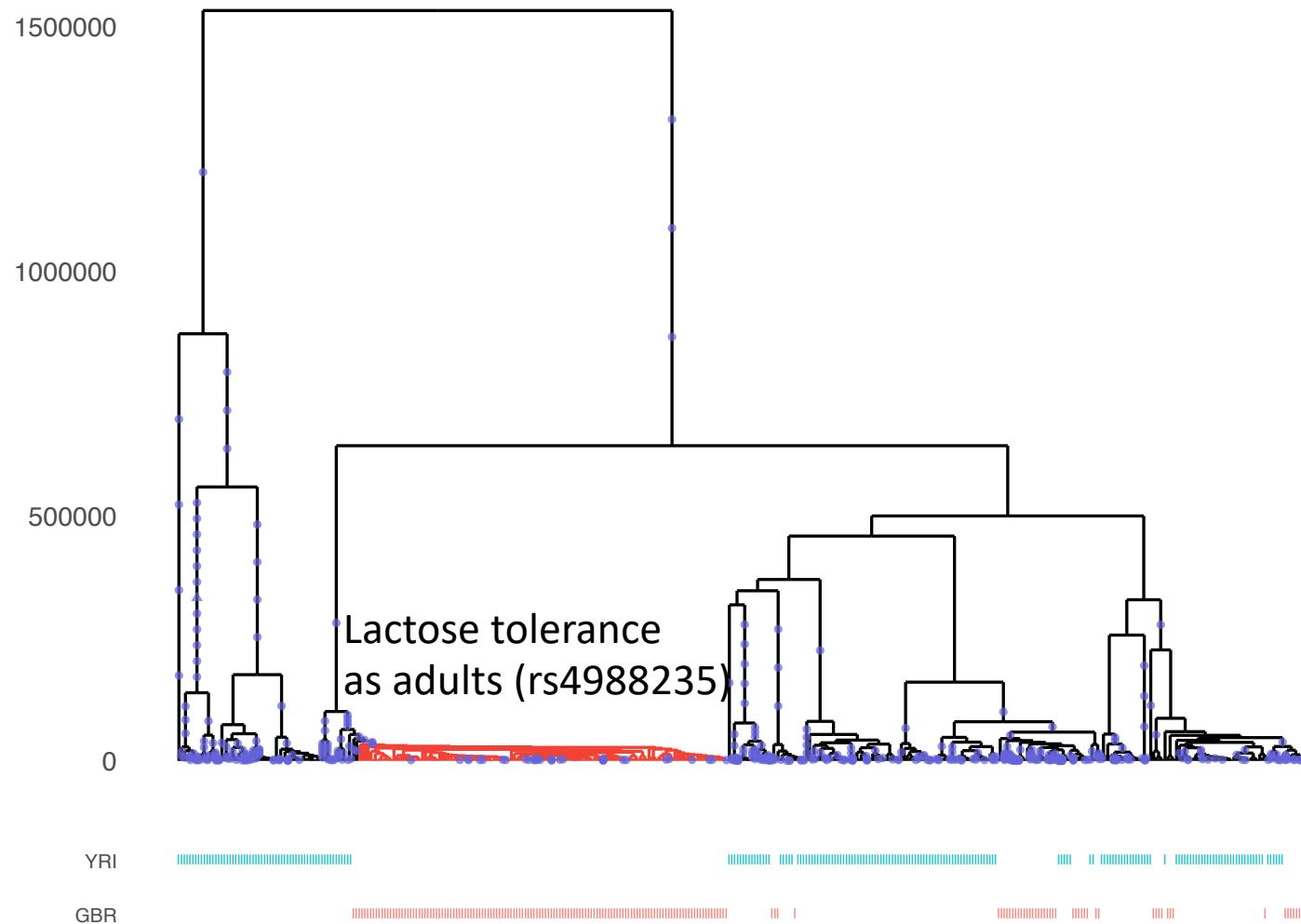
Key features:

- Fast & accurate
- Robust to errors!
- Jointly infers branch lengths and demographic history
- Moderns and ancients
- Lots of add-on tools for various types of analyses

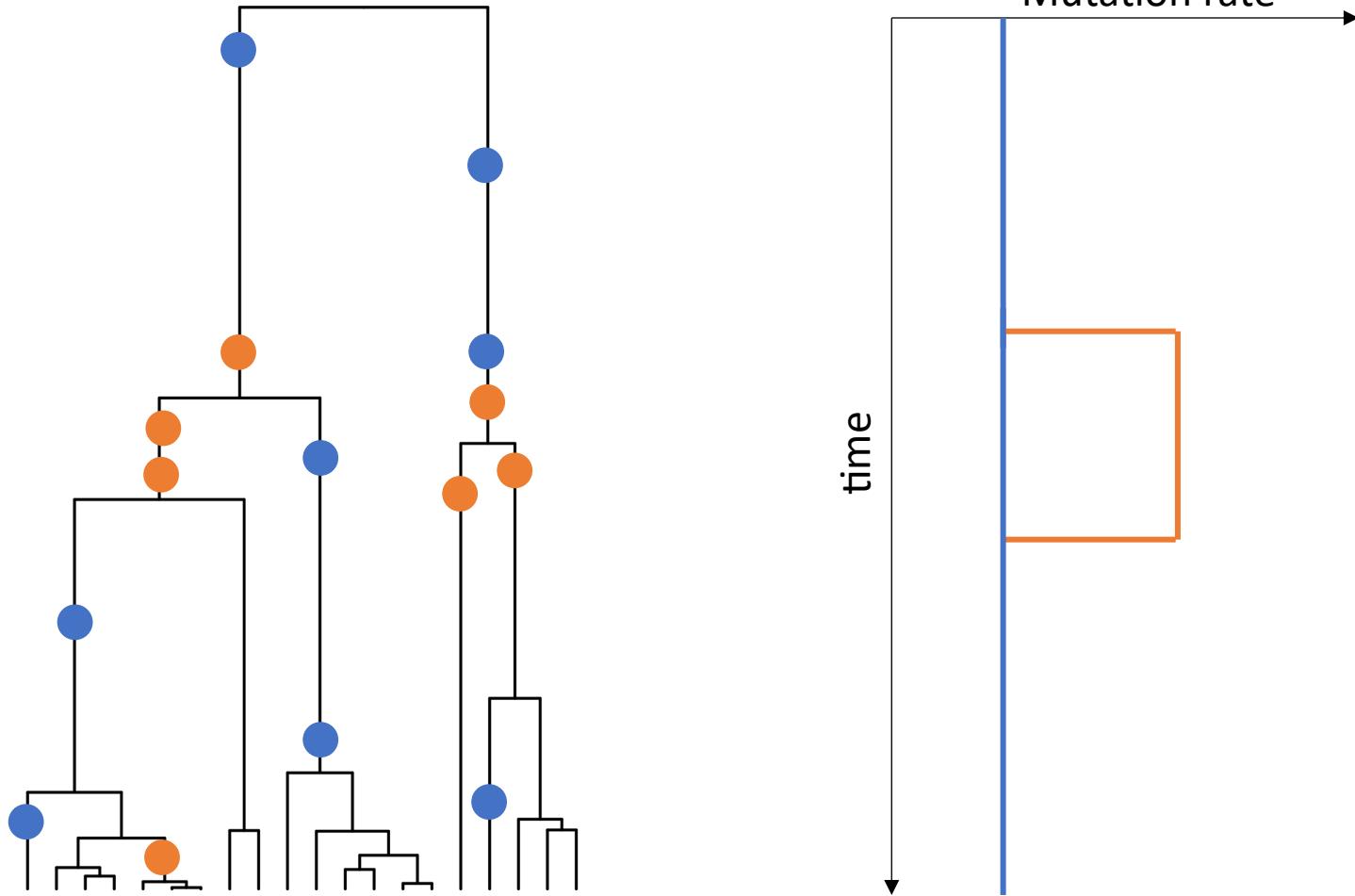
One locus can already tell us a lot about our history

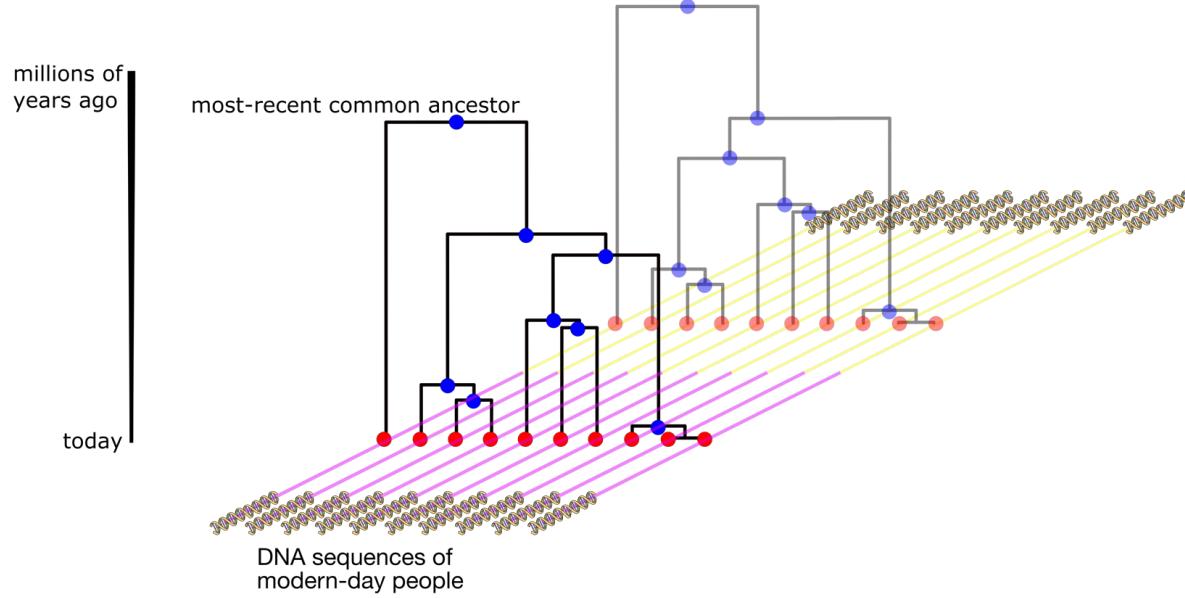


Positive selection: rapidly spreading lineage



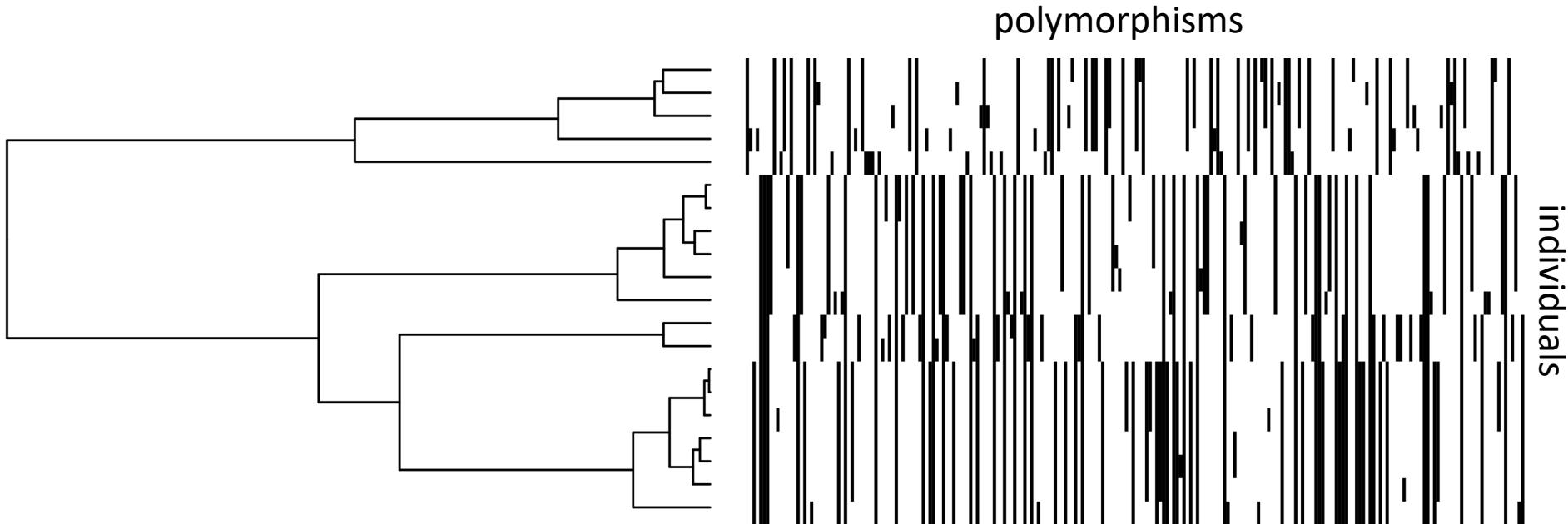
Clusters of mutations in time can capture changes in mutation rate





Inferring genealogies from genetic variation

Data and the underlying tree structure



- **Every mutation shows the existence of a branch**
- Mutations are “ordered by inclusion”
- No two branches (mutations) ever show only partial overlap

Count derived mutations to reconstruct tree topology

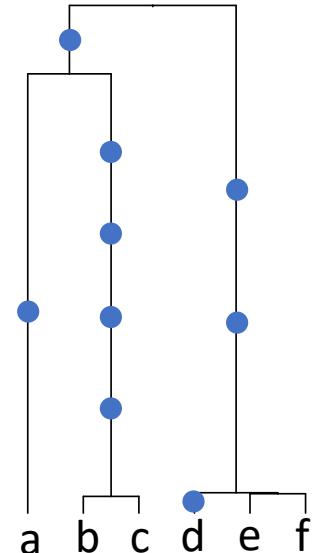
For tree topology, we want to quantify the order in which we are related to others

1. Count number of “derived mutations” to get order of coalescences

- E.g., sequence a has 1 derived mutation to (b,c)
2 derived mutation to (d,e,f)

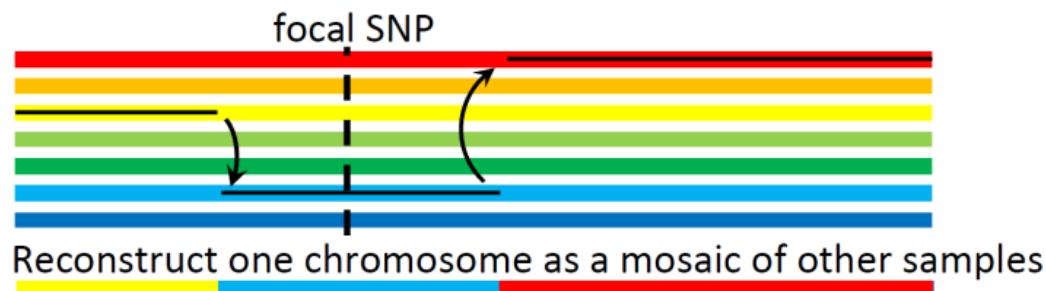
2. Coalesce mutually closest lineages

- No recombination: Guaranteed to build tree consistent with data
- This is not the case if we use “pairwise differences” (UPGMA)
 - a and (d,e,f) are closer than a and (b,c)
- Use **chromosome painting** to count derived mutations accounting for recombination



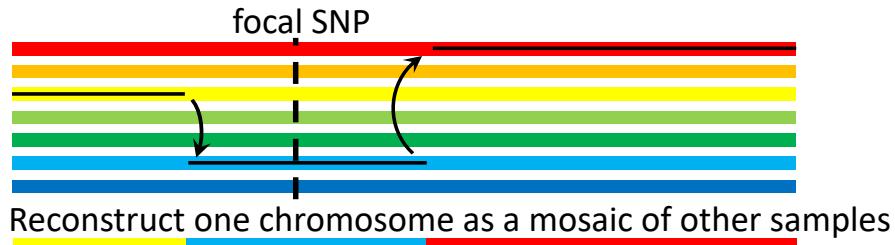
Hidden Markov model (HMM)

Li and Stephens, Genetics, 2003; Lawson et al., PLOS Genetics, 2012

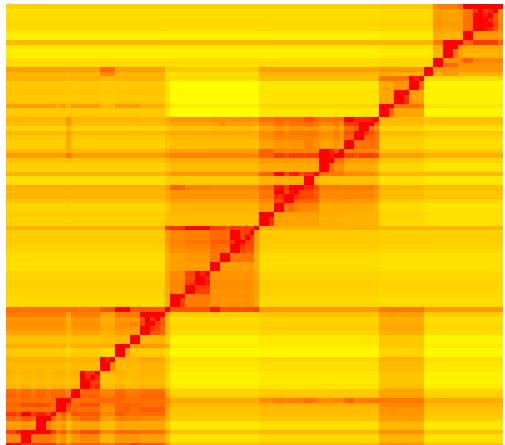


Summary of Relate

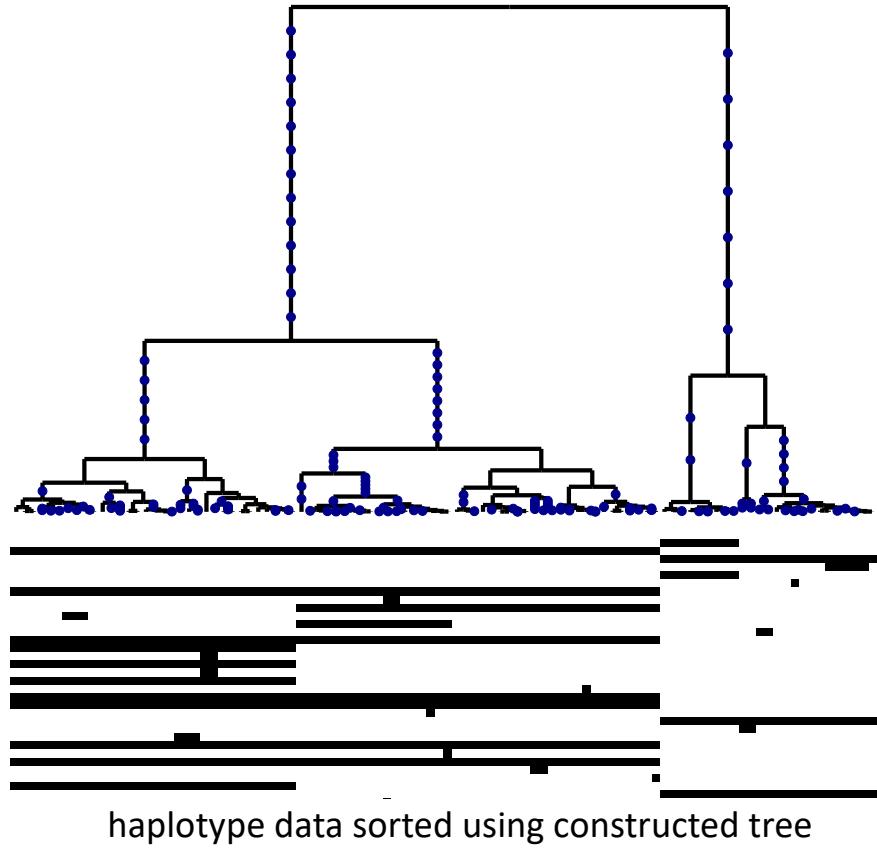
Hidden Markov model (HMM)



Distance matrix for focal SNP

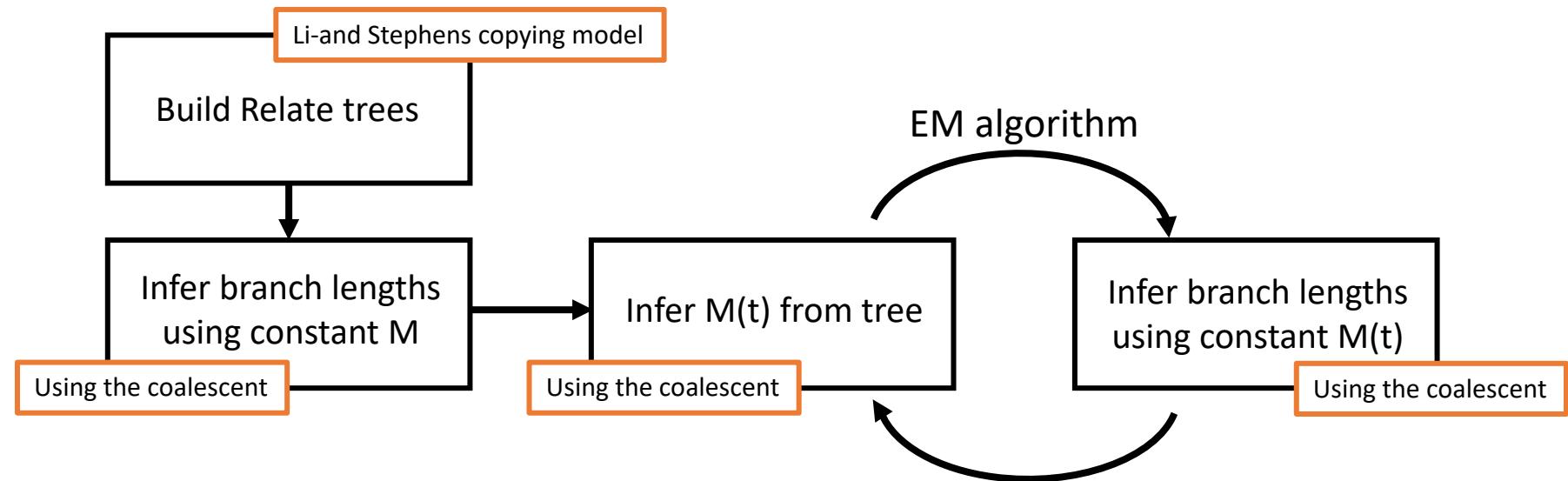
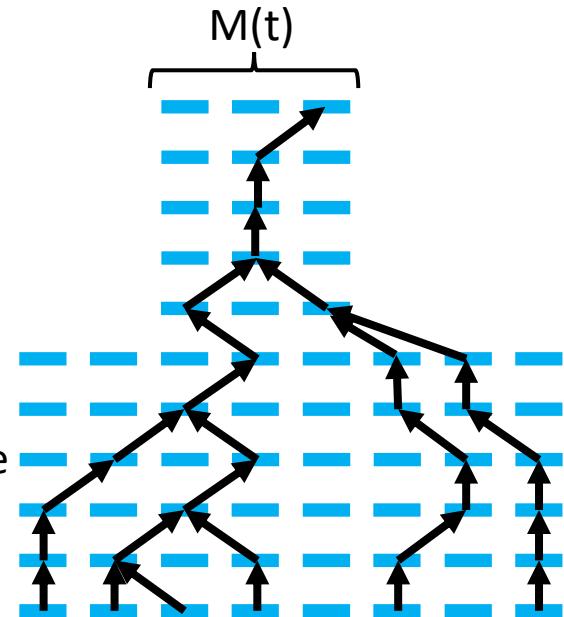


Hierarchical clustering
&
MCMC for branch lengths



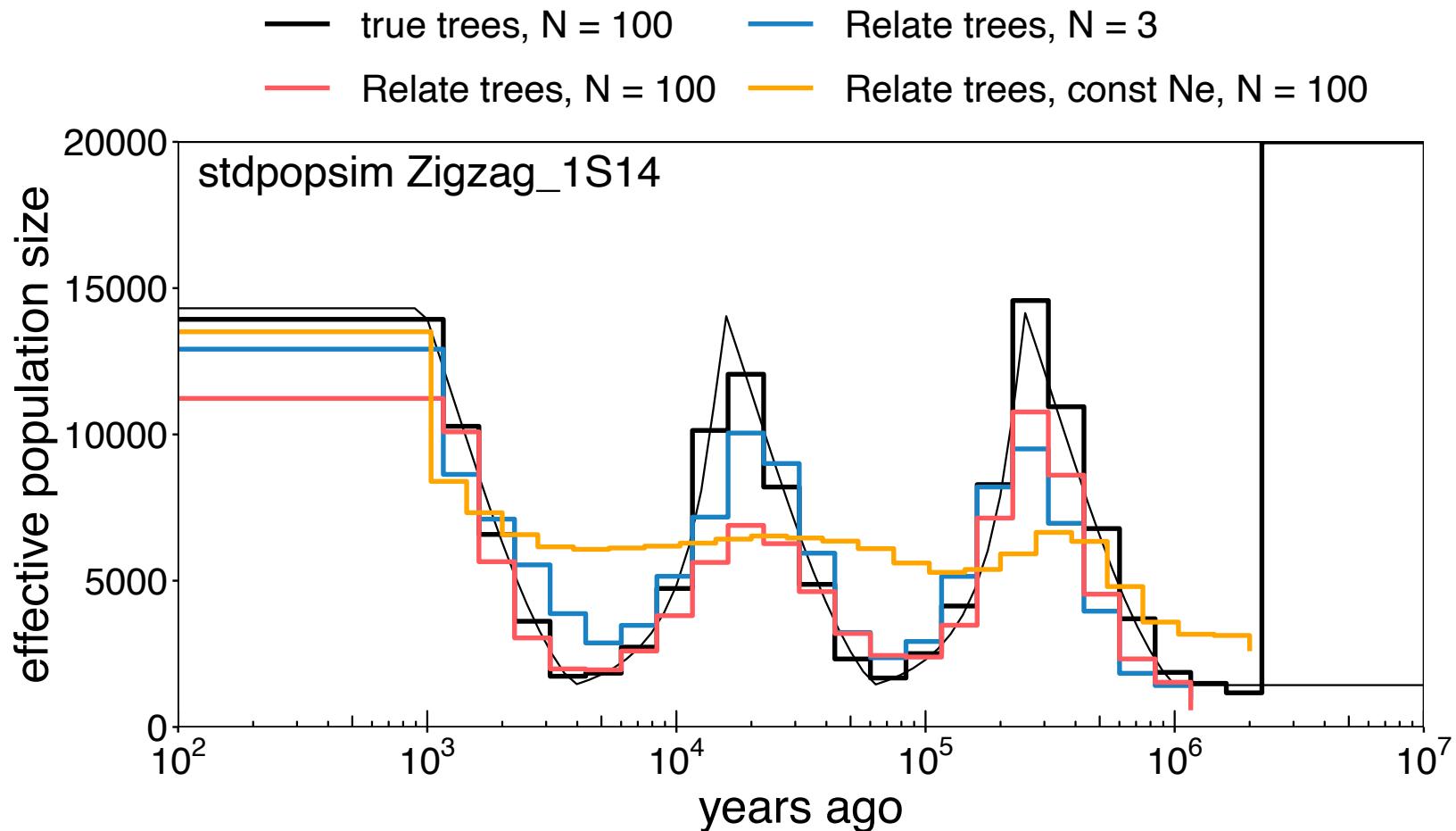
Branch lengths and population size is estimated jointly in an EM algorithm

- Expected branch lengths depend on population size $M(t)$ (or coalescence rates $1/M(t)$)
- While there are j lineages, the rate at which a coalescence happens is $\binom{j}{2}/M(t)$ a time t ago
- Demography is shared genome-wide, so we average across trees
- So within a time interval, scaled fraction of trees where coalescence occurs is inversely proportional to $M(t)$



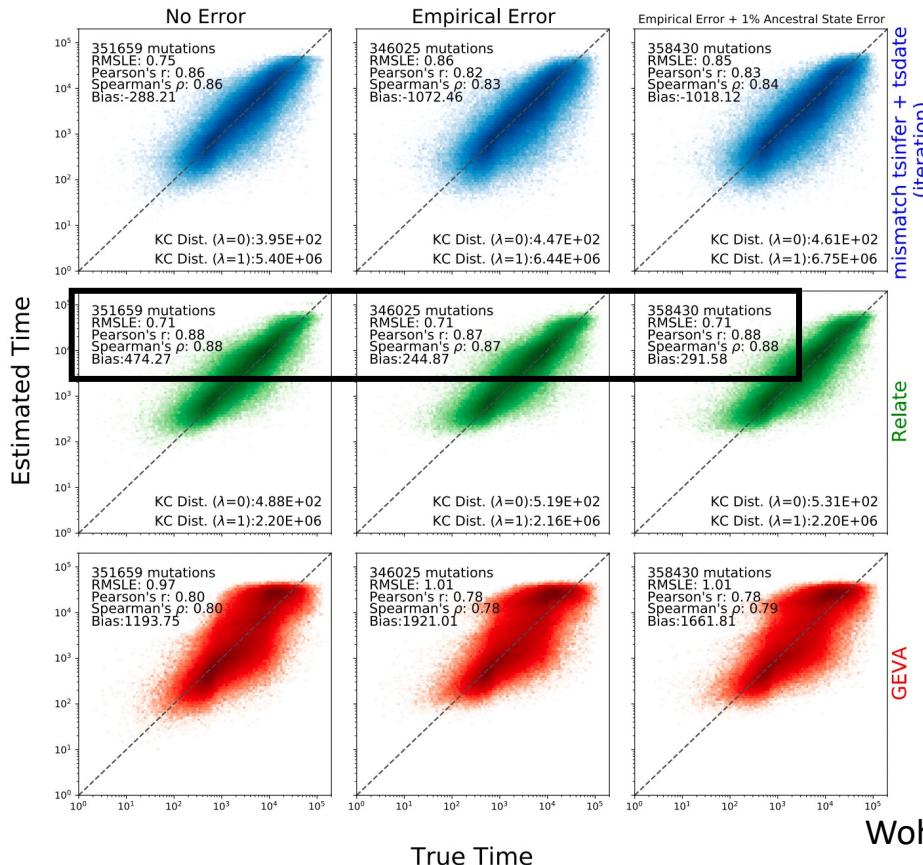
Population size changes through time are jointly inferred in Relate

- Effective population size = inverse coalescence rate
- N: number of diploid samples



Speed and accuracy of Relate

- About 14,000 times faster than previous method, ARGWEAVER (1 min. vs. 200 hours), slower than tsinfer + tsdate
- Builds “correct” tree if no recombination
- Accurate, robust to data errors
- Can sample posterior branch lengths

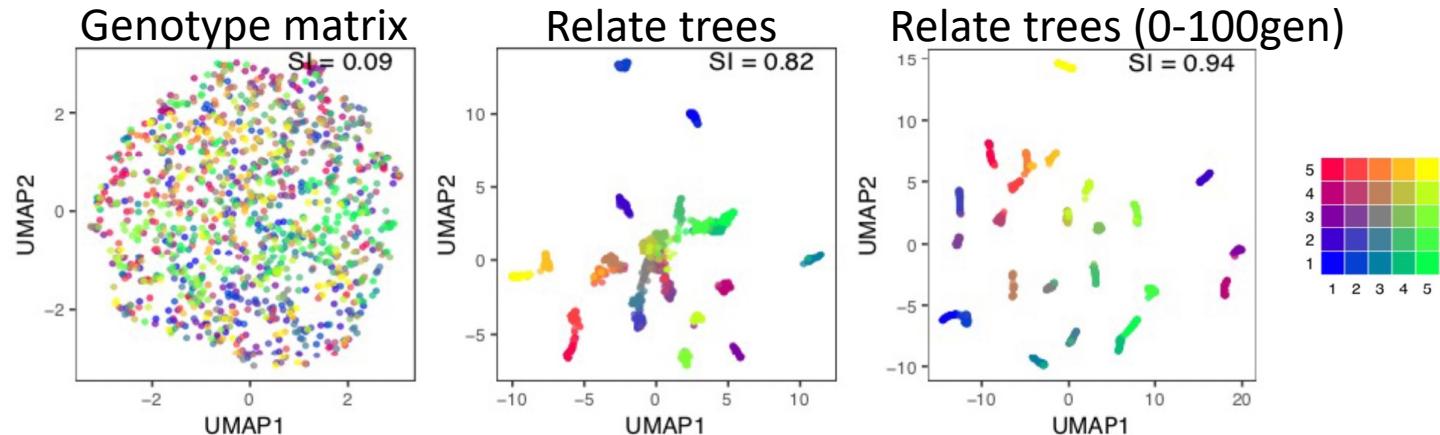
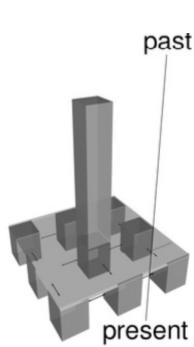


Genealogy-based inference of human evolutionary history

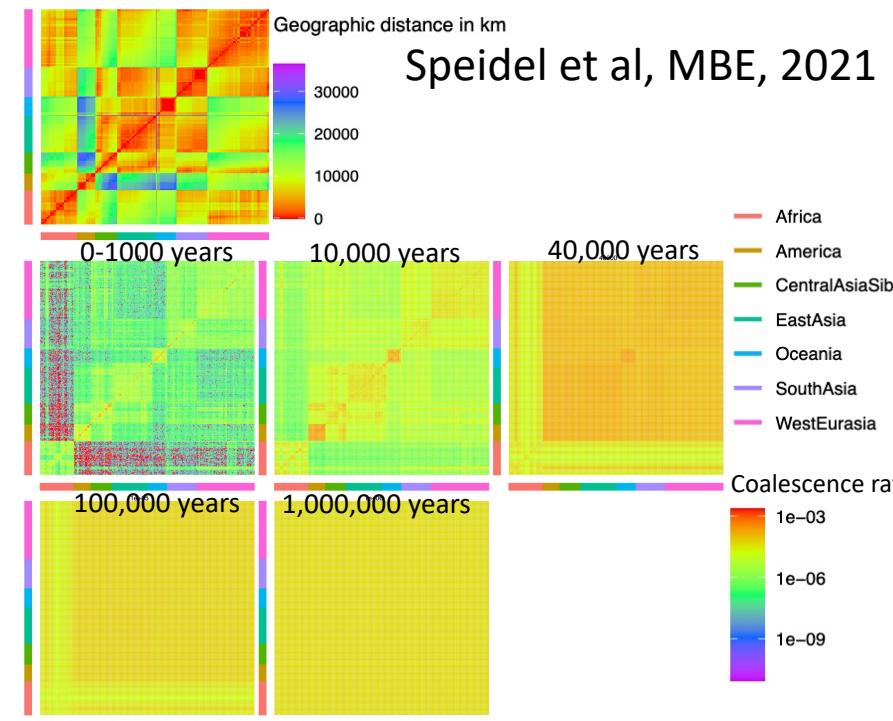
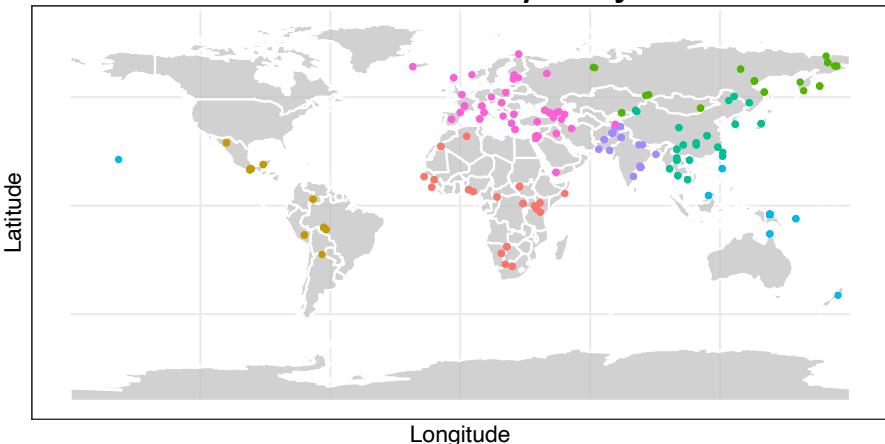
One reconstruction of history, many applications that are self consistent

Inferring fine-scale population structure and how it changed through time

Fan, Mancuso, Chiang, AJHG, 2022

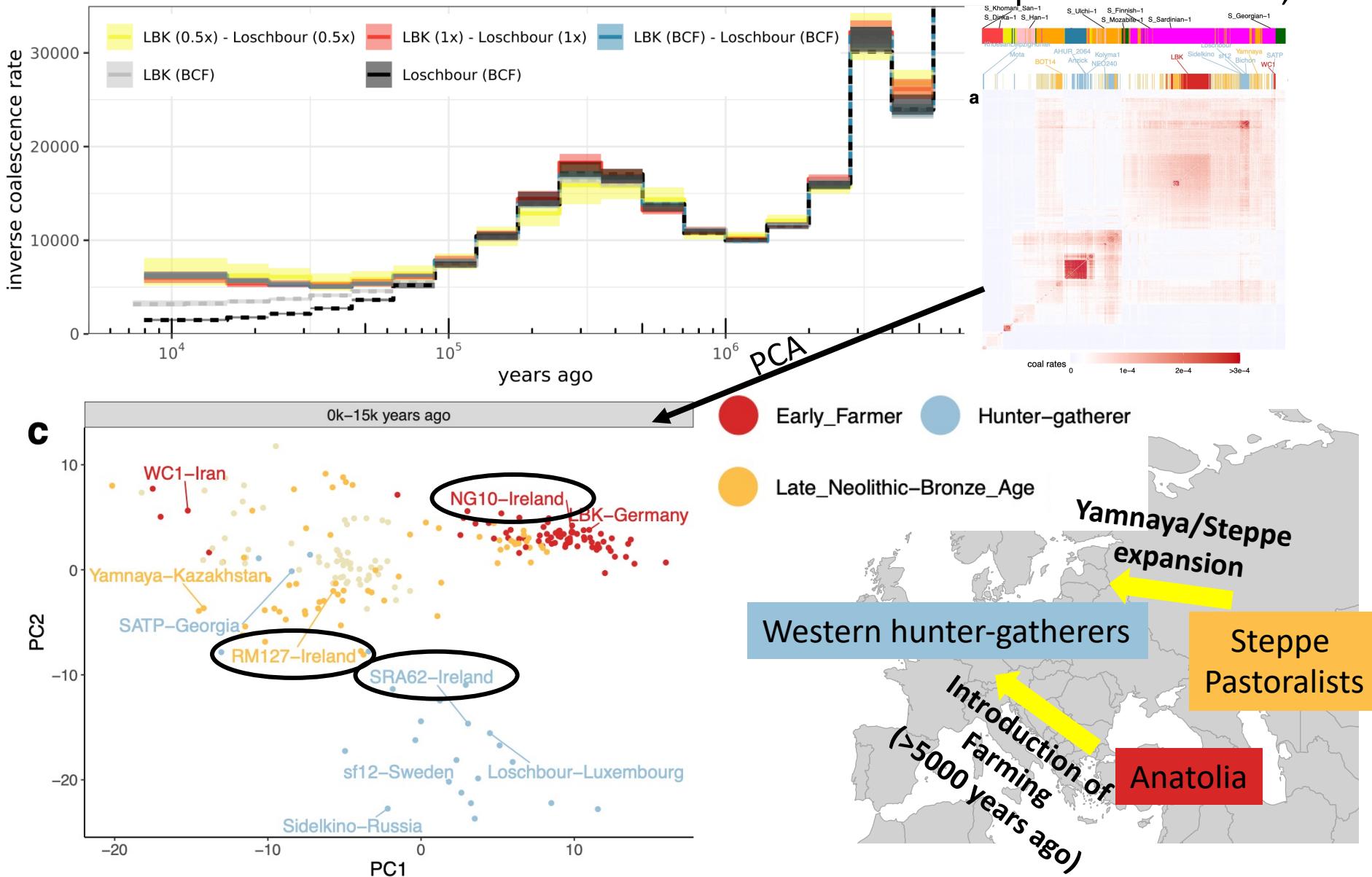


Simons Genome Diversity Project

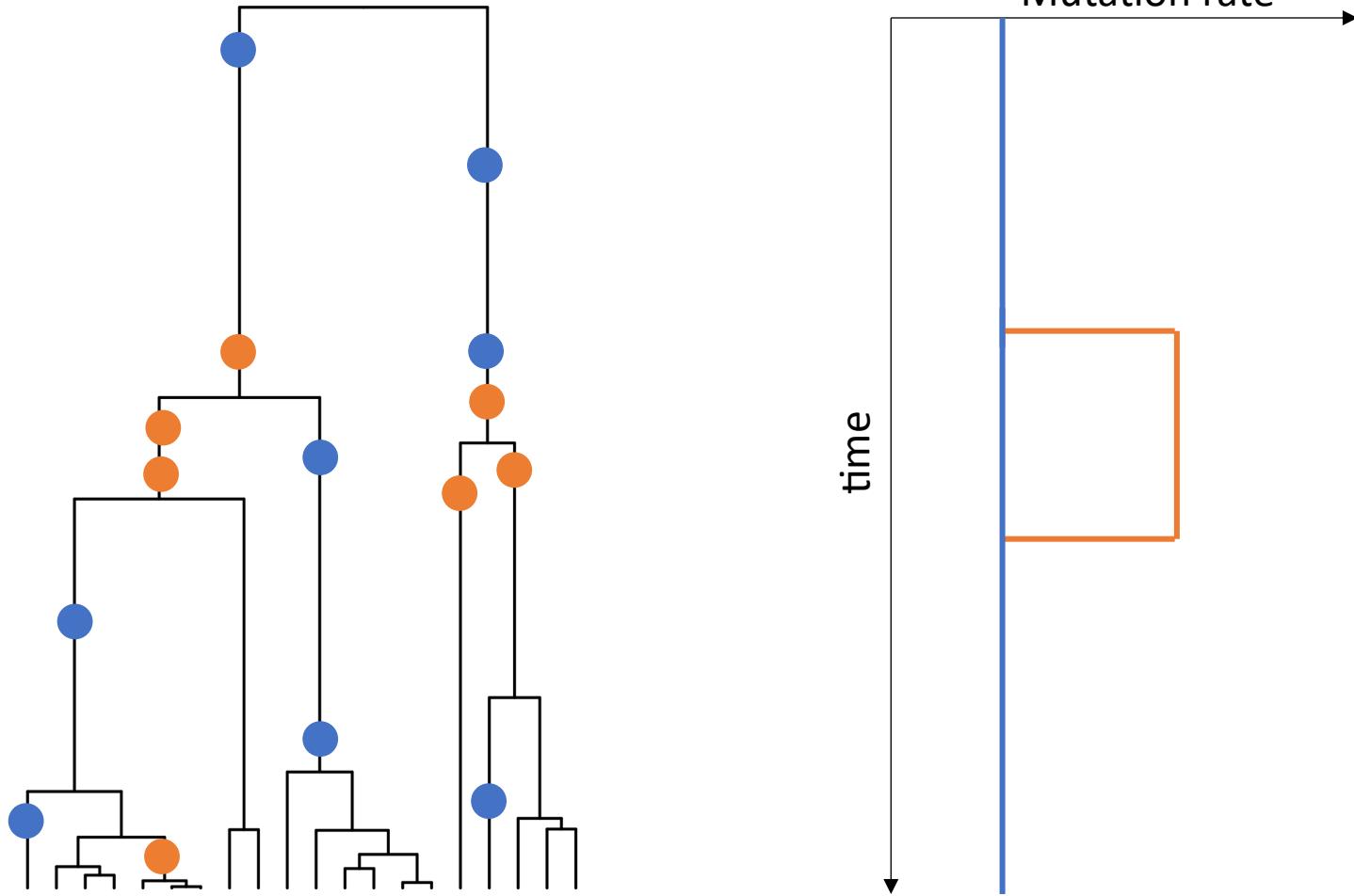


Colate: Inferring coalescence rates for low-coverage, unphased (ancient) genomes

Speidel et al. MBE, 2021

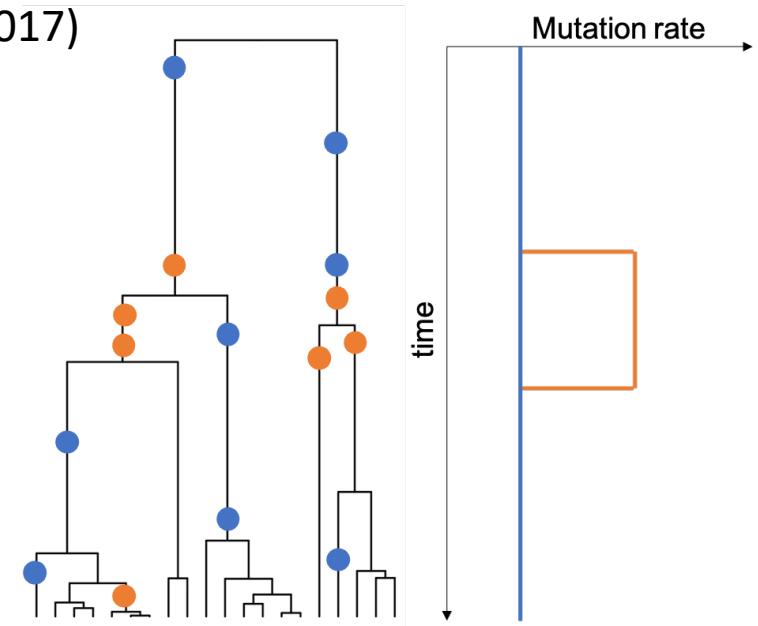
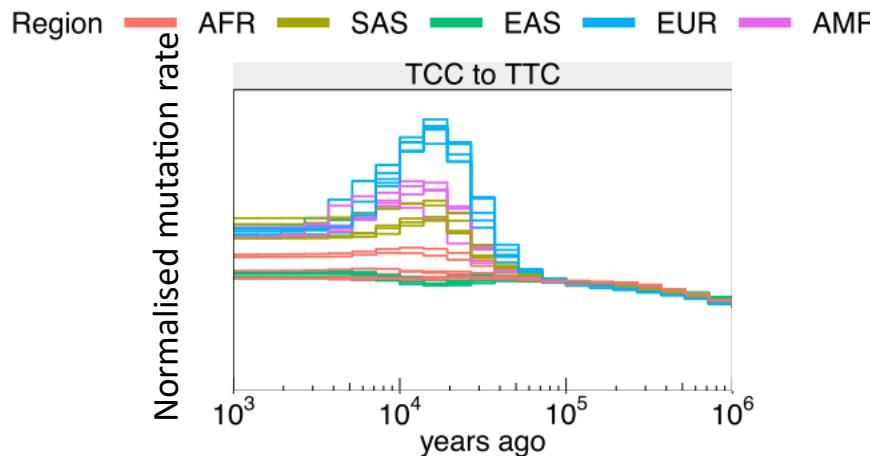


Clusters of mutations in time can capture changes in mutation rate



TCC/TTC mutation rates have experienced a strong pulse in the Upper Paleolithic

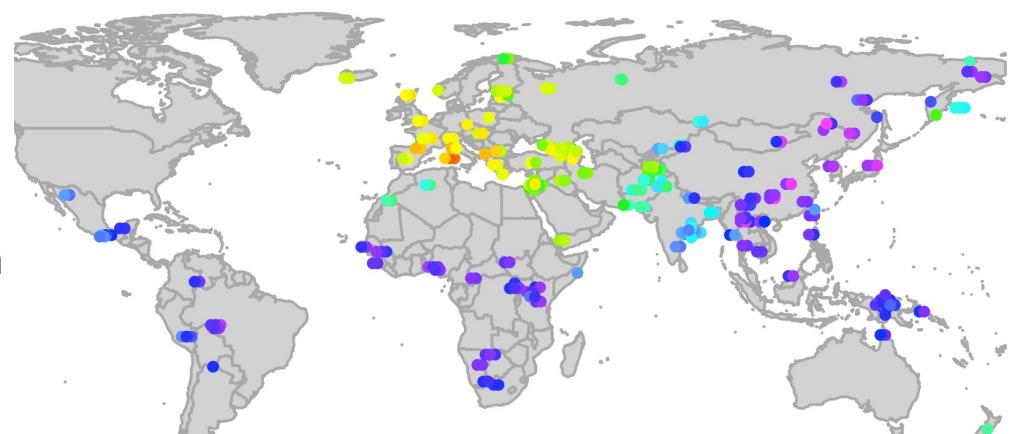
- First reported by Kelley Harris (PNAS 2015, eLife 2017)
- Unknown cause (genetic?, environmental?)
- Previously mainly studied in modern groups



Speidel, Nature Genetics, 2019

How did this spread to all West Eurasians today?

Colour shows strength of elevation in TCC/TTC mutation rates



TCC/TTC mutation rate increase happened and spread before farming, >15k years ago!

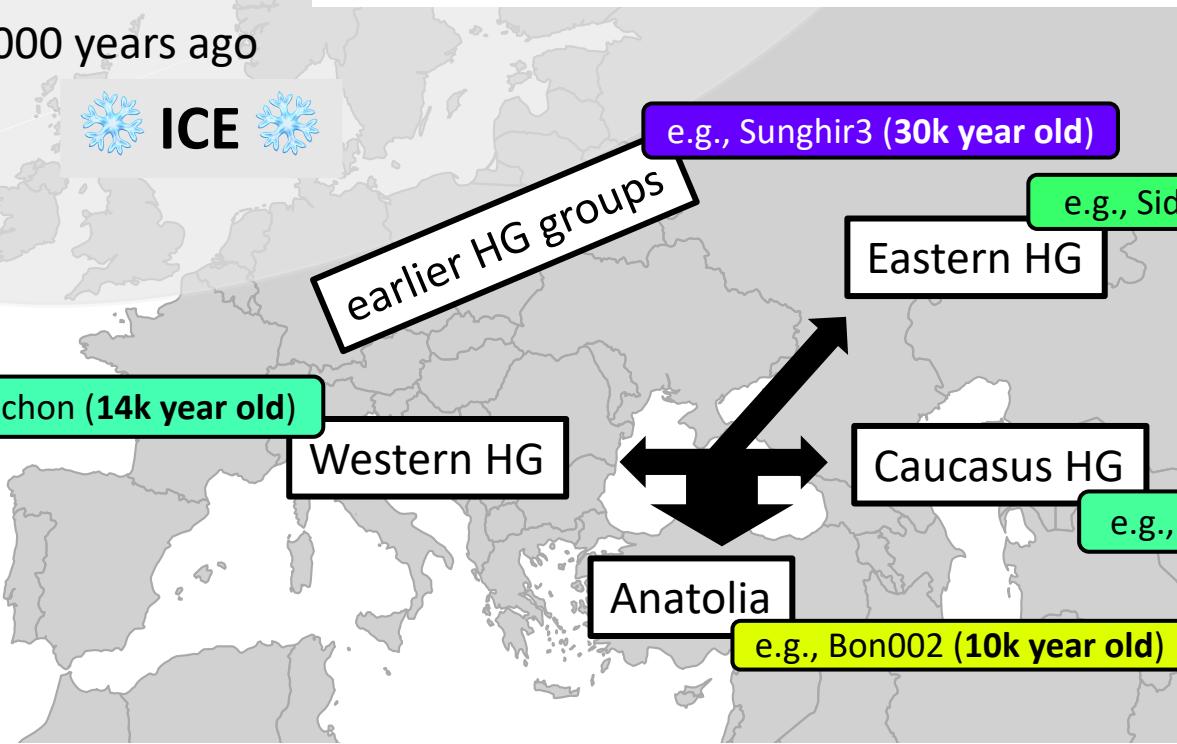
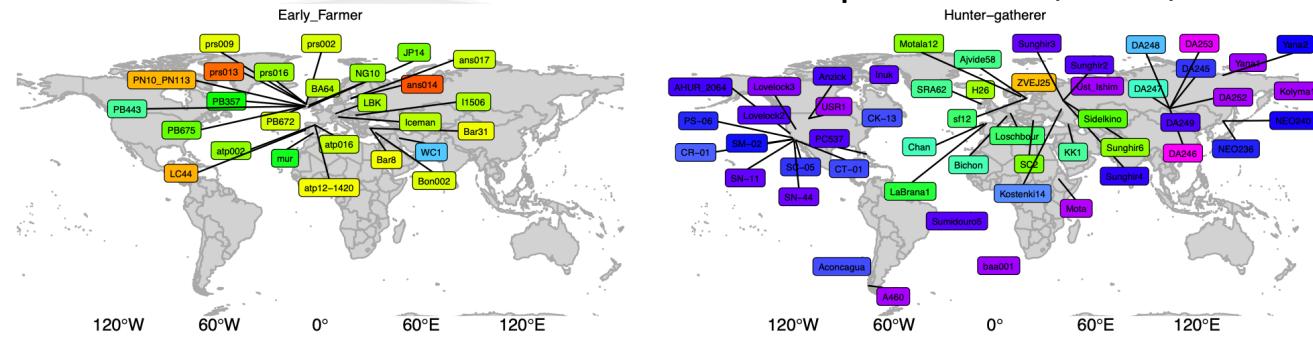
Ice Age Europe



>20,000 years ago

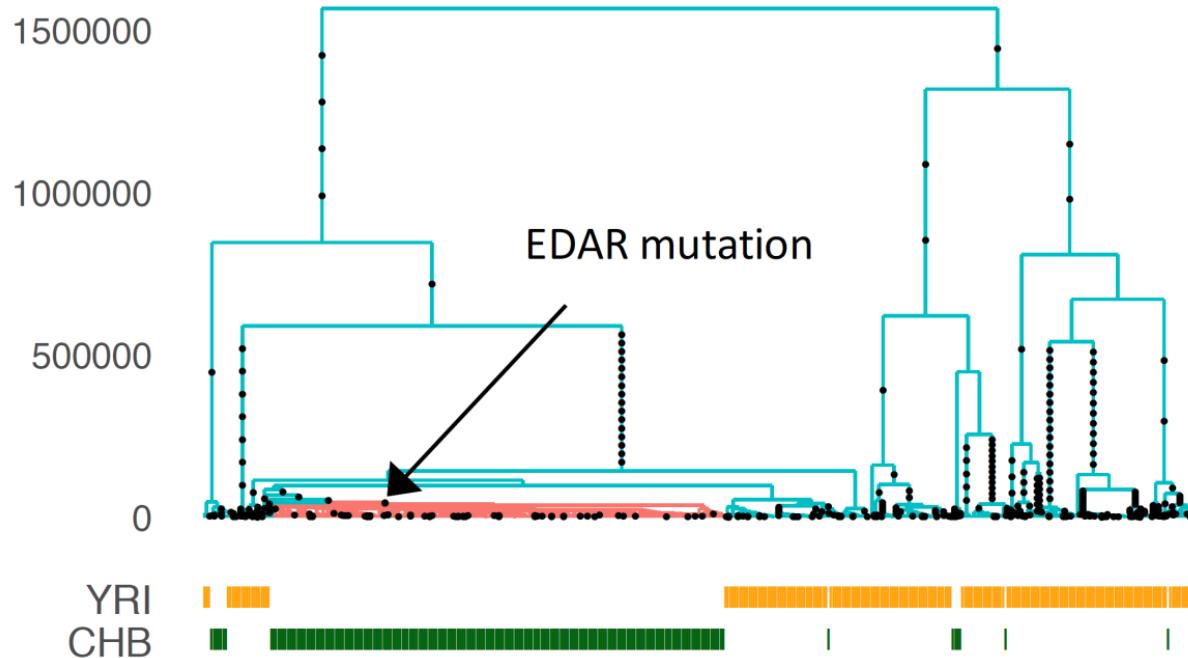


ICE



Colour shows strength of elevation in TCC/TTC mutation rates

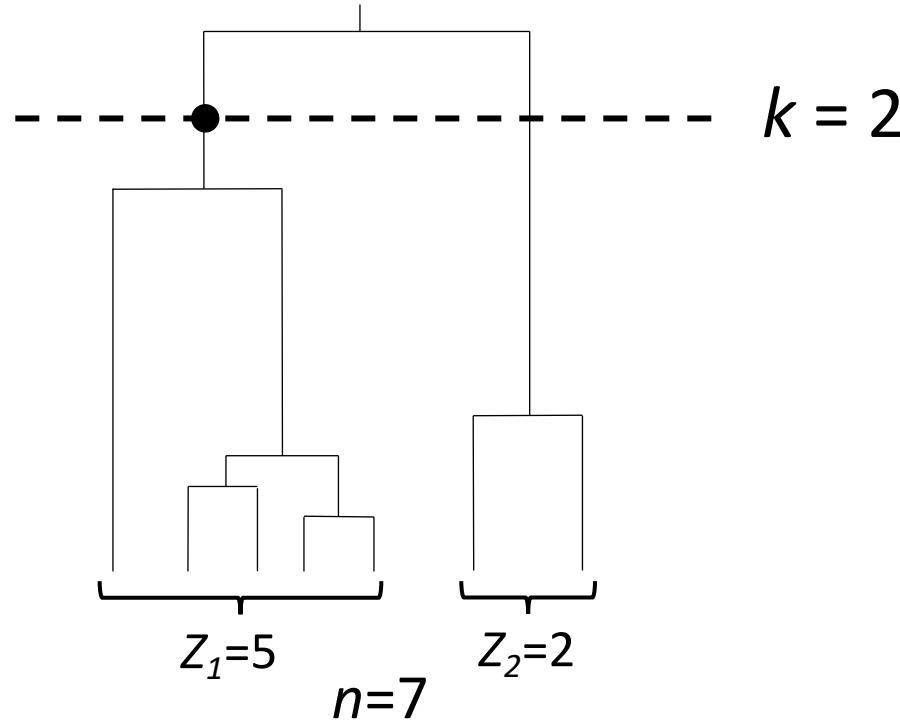




Quantifying positive natural selection on a single mutation

- Genetic adaptations to changing environment, diet, lifestyles,...
- Use trees incorporating demographic history

How quickly does a mutation spread in the neutral case?



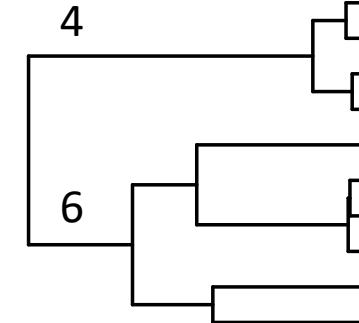
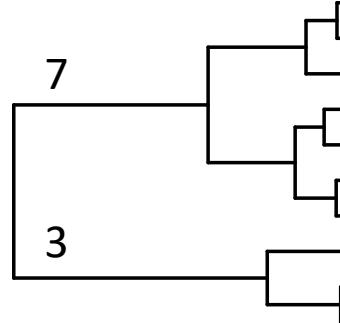
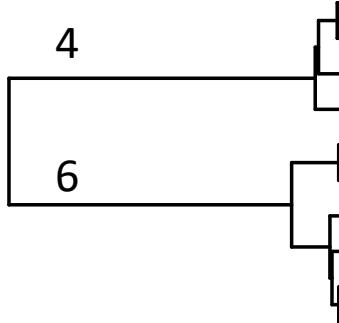
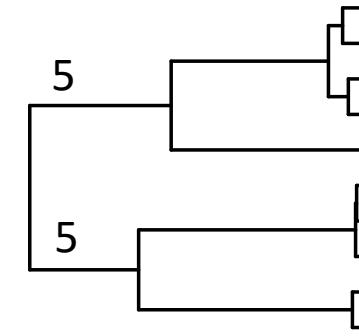
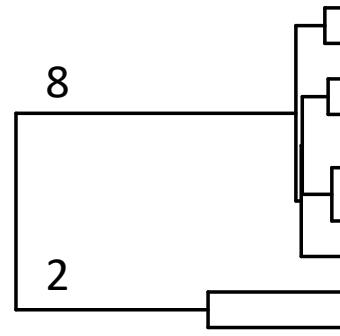
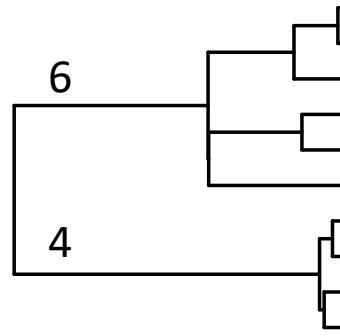
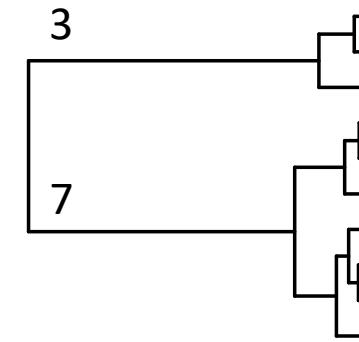
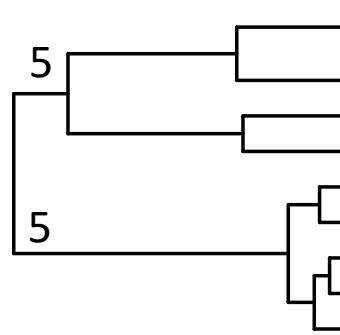
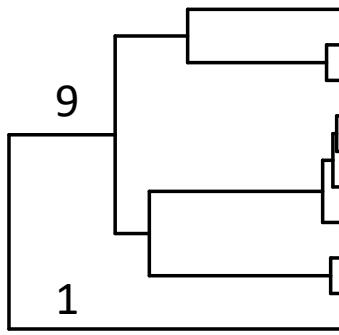
We can write down the analytical distribution for the number of descendants of a mutation arising while k lineages remain

Example: if $k=2$, this is just a **uniform distribution**

$$P(5 \text{ descendants}) = 1/6$$

The $k = 2$ case

We expect “unbalanced” tree shapes!

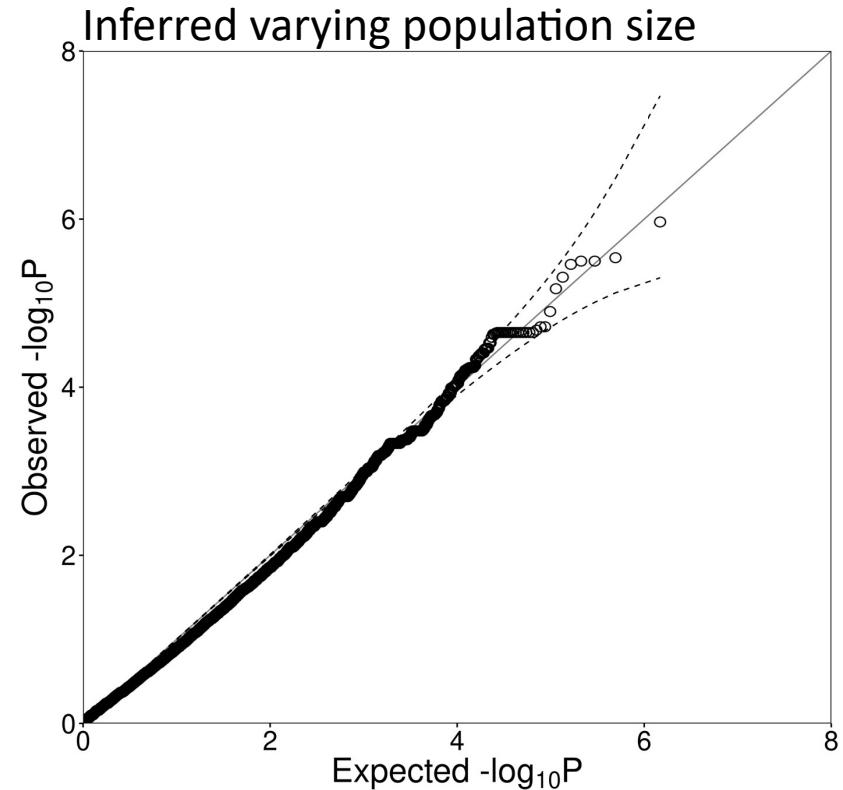
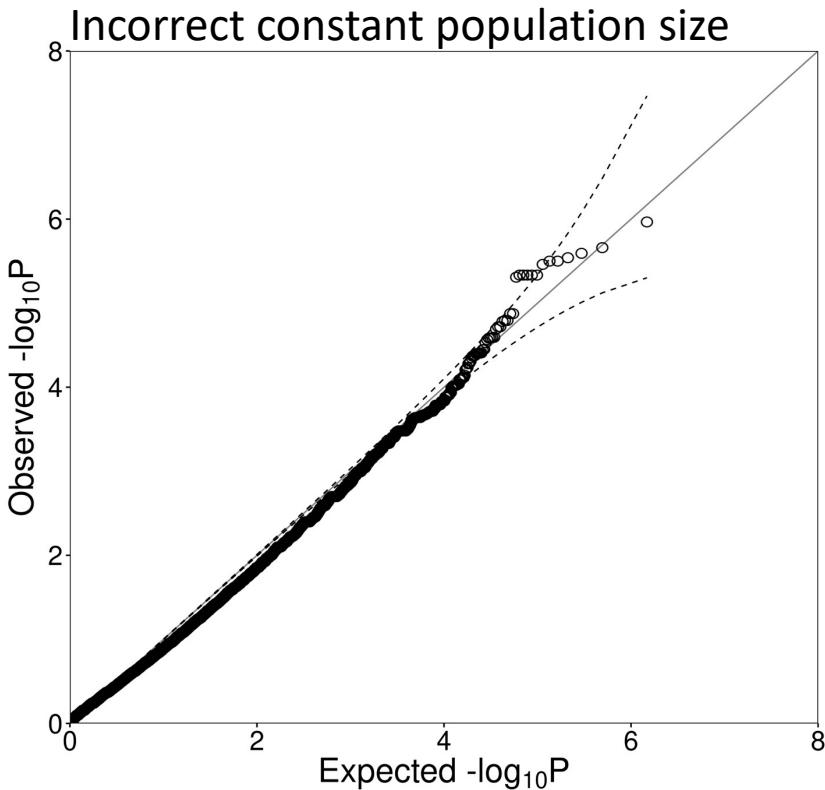


P-values: very well calibrated under null simulations of no selection

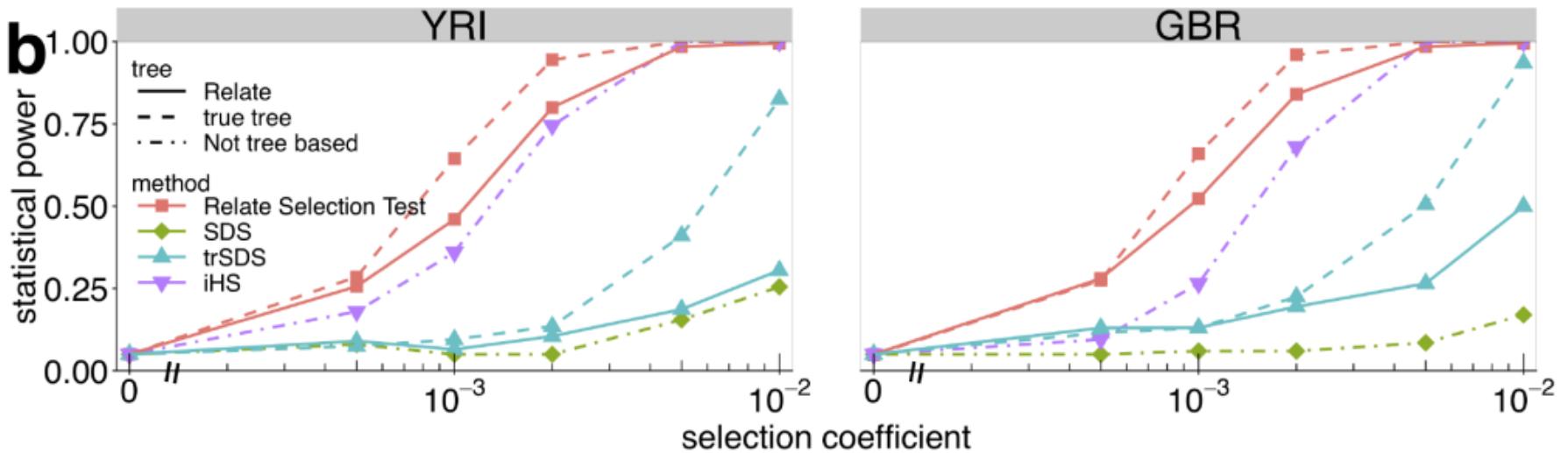
N=1000, 250Mb

Bottleneck population size

Quantile-quantile plot:

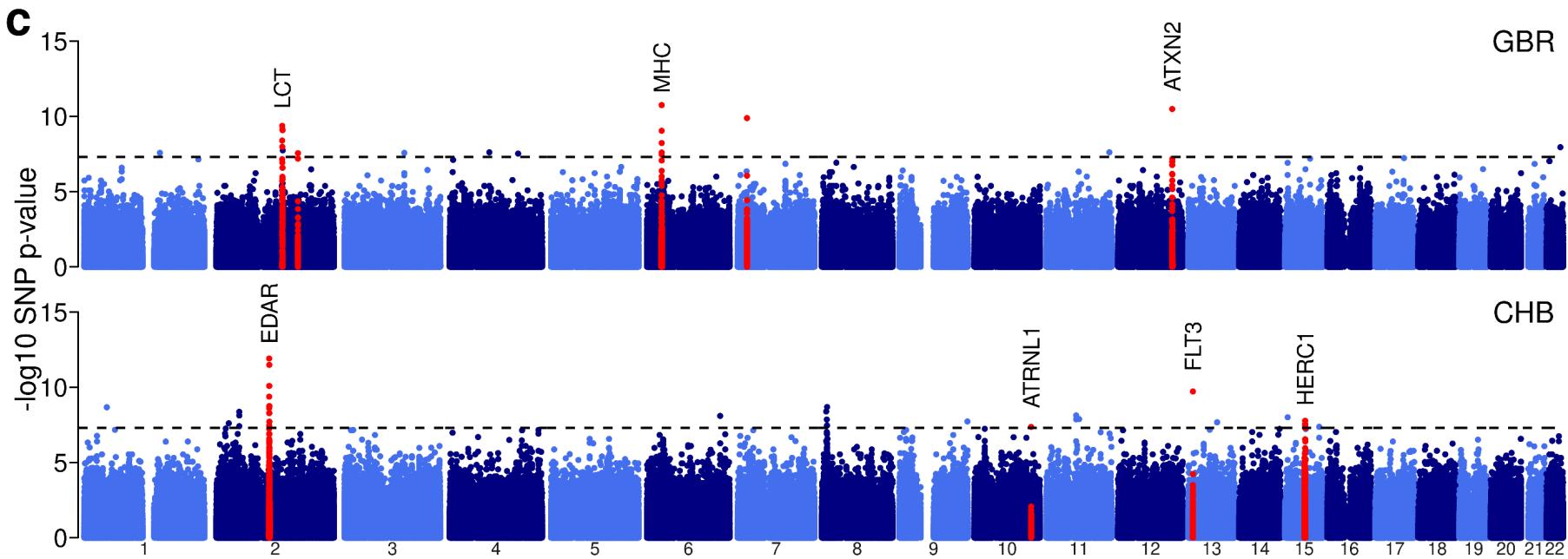


Improved power to see weak selection



Genome-wide selection p-values

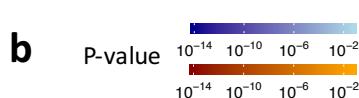
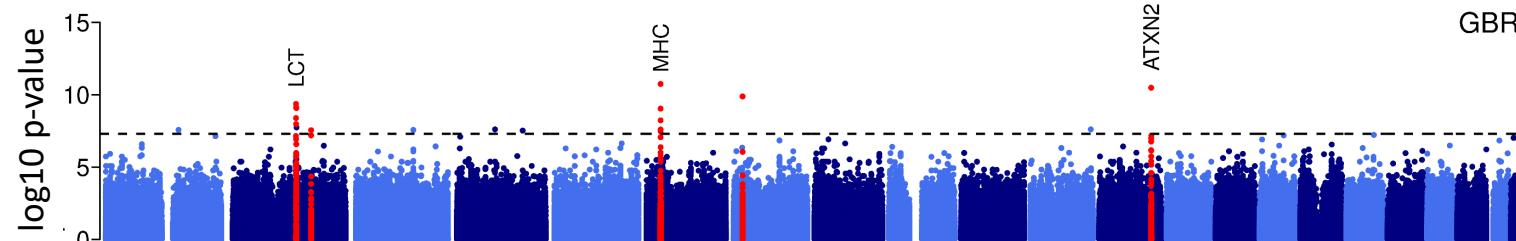
Given most traits are highly polygenic, expect mainly weak, polygenic selection



How does weak selection evidence vary by trait?

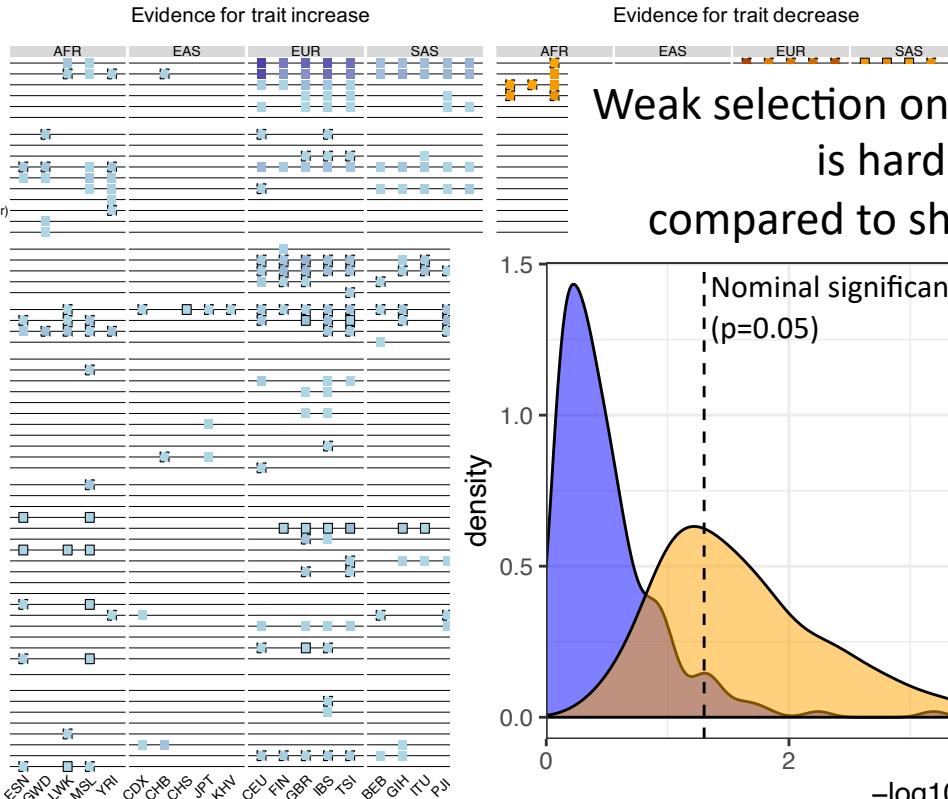
Many key events in our evolutionary history are only implicated as subtle effects in our genomes

Selection p-values: only a handful of “genome-wide significant” loci

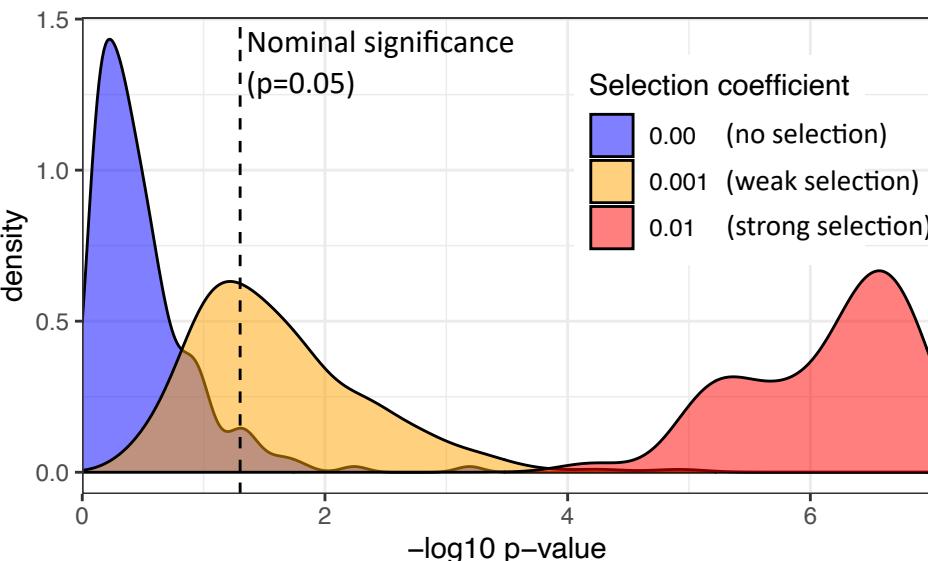


Physical traits	Sitting height (UKB) Standing height (UKB) BMI (UKB) Hip circumference (UKB) Waist circumference (UKB) Skin colour (UKB)
	Height BMI BMI (adj. for smoking behaviour) Hip circumference Hip circumference adj. for BMI Waist circumference Waist circumference adj. for BMI in active individuals Waist circumference adj. for BMI (adj. for smoking behaviour) Waist-to-hip ratio Waist-to-hip ratio adj. for BMI
Blood pressure	Blood pressure Diastolic blood pressure Systolic blood pressure Pulse pressure Resting heart rate
Platelets	Plateletcrit Mean platelet volume Platelet distribution width Platelet count Red blood cell count Hematocrit Hemoglobin concentration Mean corpuscular hemoglobin Mean corpuscular hemoglobin concentration Mean corpuscular volume
Red blood cells	Reticulocyte fraction of red cells Reticulocyte count Immature fraction of reticulocytes High light scatter reticulocyte count High light scatter reticulocyte percentage of red cells White blood cell count Lymphocyte percentage of white cells Lymphocyte count
White blood cells	Myeloid white cell count Monocyte percentage of white cells Monocyte count Eosinophil percentage of white cells Eosinophil percentage of granulocytes Neutrophil count Neutrophil percentage of white cells Neutrophil percentage of granulocytes Sum neutrophil eosinophil counts Sum eosinophil basophil counts Sum basophil neutrophil counts
Lipids	Cholesterol, total Blood metabolite ratios Blood metabolite levels Blood glucose Glomerular filtration rate (creatinine) Glomerular filtration rate in non diabetics (creatinine)
Other traits	Gut microbiota (bacterial taxa) Educational attainment (years of education) Schizophrenia (PGC)

Lots of Polygenic selection signals

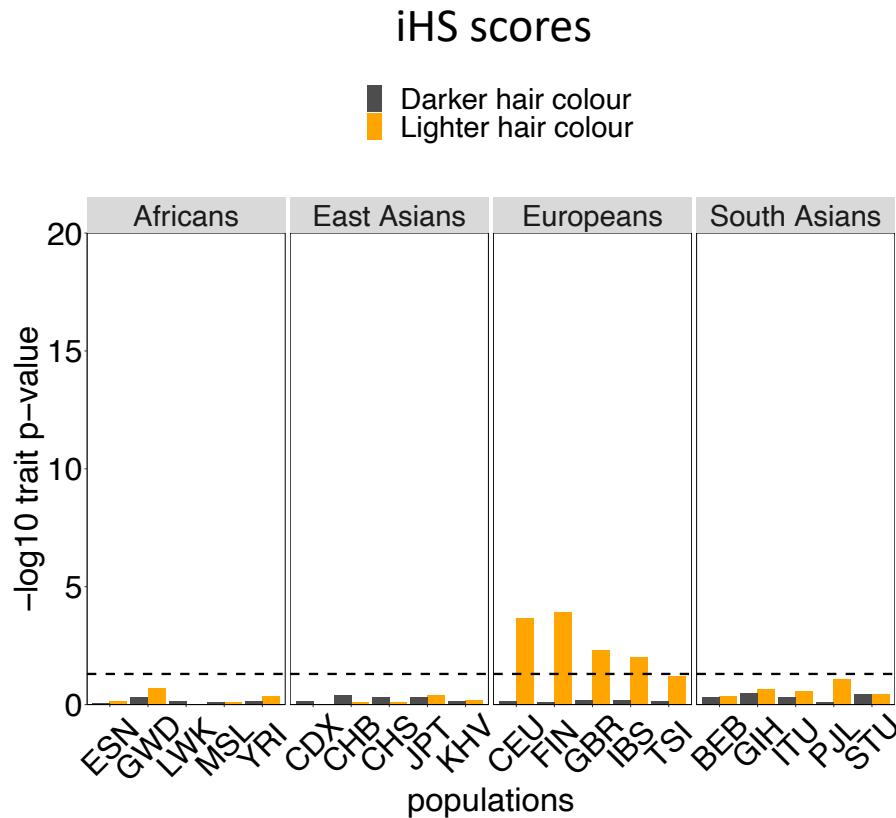
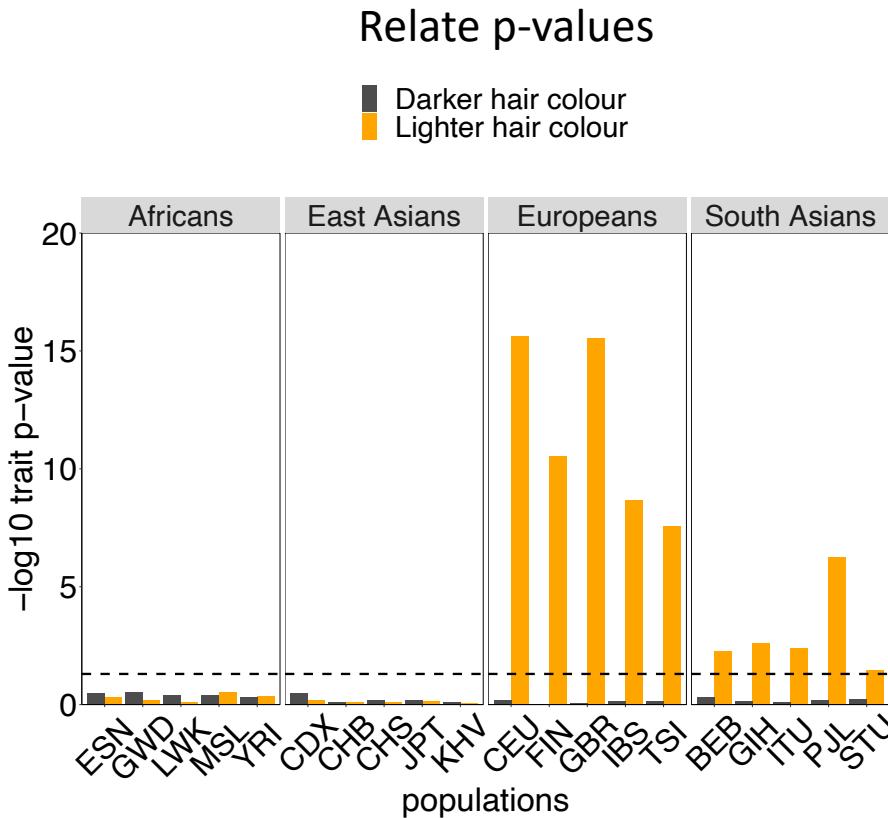


Weak selection on individual mutations
is hard to detect,
compared to shifts in distributions



Evidence of selection on a trait: hair colour

1. Use **effect direction** of "genome-wide significant" associations
2. Compare selection p-values to frequency matched random SNPs (Wilcoxon rank-sum test)



CLUES: Importance-sampling based method for inferring selection coefficients

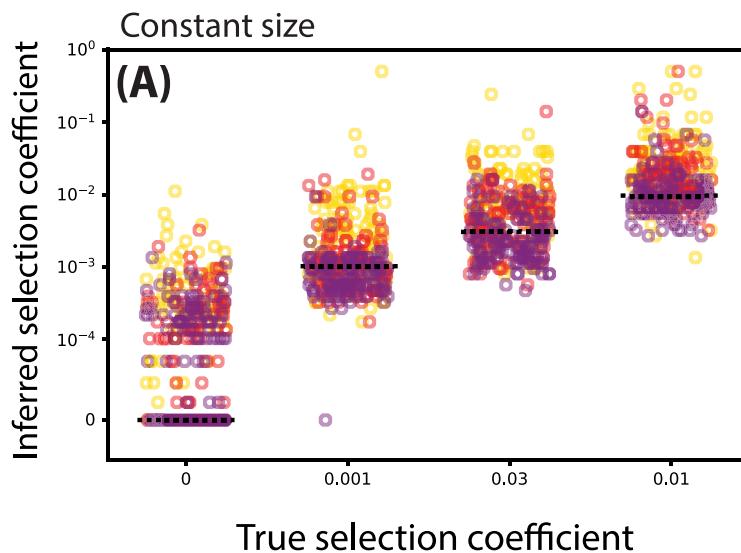
Aaron Stern



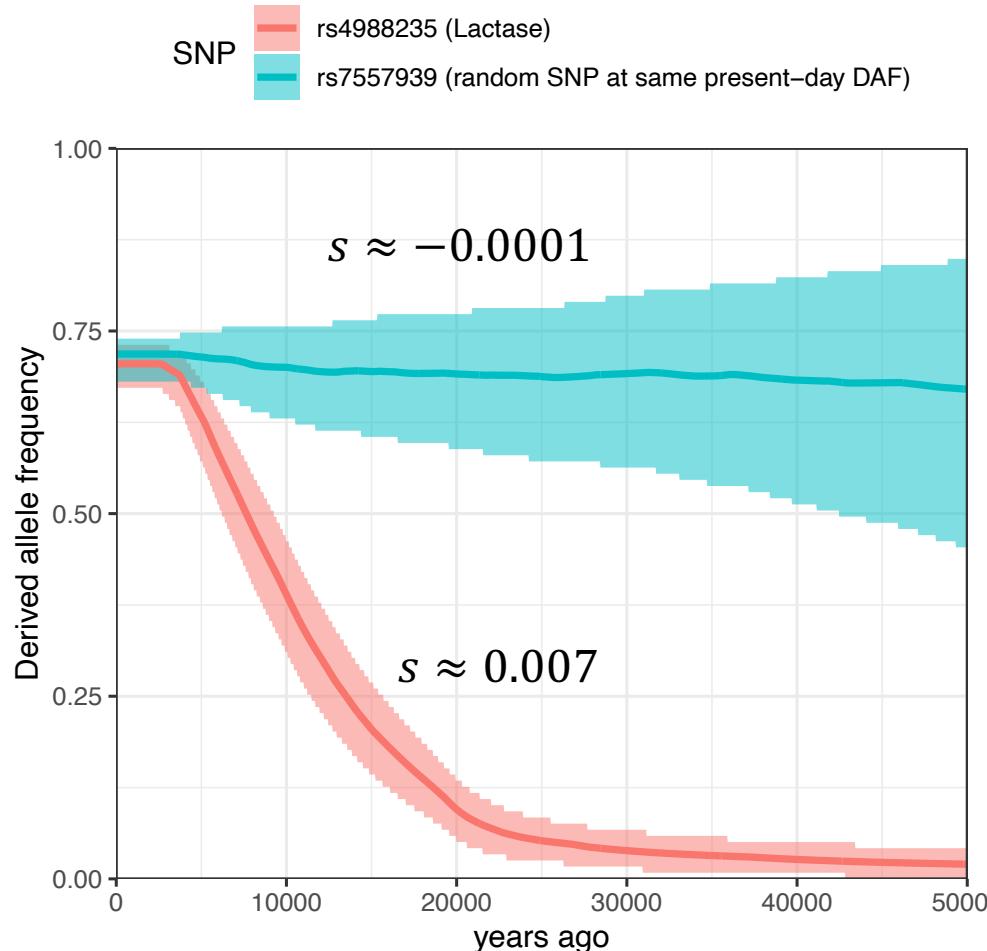
Aaron J. Stern, Peter R. Wilton, Rasmus Nielsen. **PLOS Genetics, 2019.**

Aaron J. Stern, Leo Speidel, Noah A. Zaitlen, Rasmus Nielsen. **AJHG 2021**

Simulations:



1000 Genomes Project British:



Summary & outlook



- It is now possible to build genealogical trees for huge datasets, in humans and other species (currently 10,000 individuals or more)
 - Humans (ancient and present)
 - Dogs and wolves
 - Mice
 - Bacteria
 - Atlantic cod, Cichlids
 - Waterhemp, Arabidopsis
- These trees capture information about many processes including
 - Migrations and ancient introgression
 - Mutation rate evolution
 - Trait evolution
 - (and many more things)
- Lots of scope for more methods using inferred genealogies under development

....creative approaches to leverage trees to answer biological questions!