



2021 EMBO Course in Population Genomics

Learning about evolution by building coalescent trees

Leo Speidel



How are we all related?

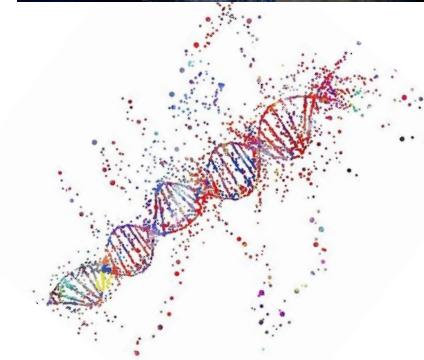
Everyone is different

- Yet people share lots of characteristics
(hair colour, blood group, disease susceptibility)
- Relatives tend to be more similar



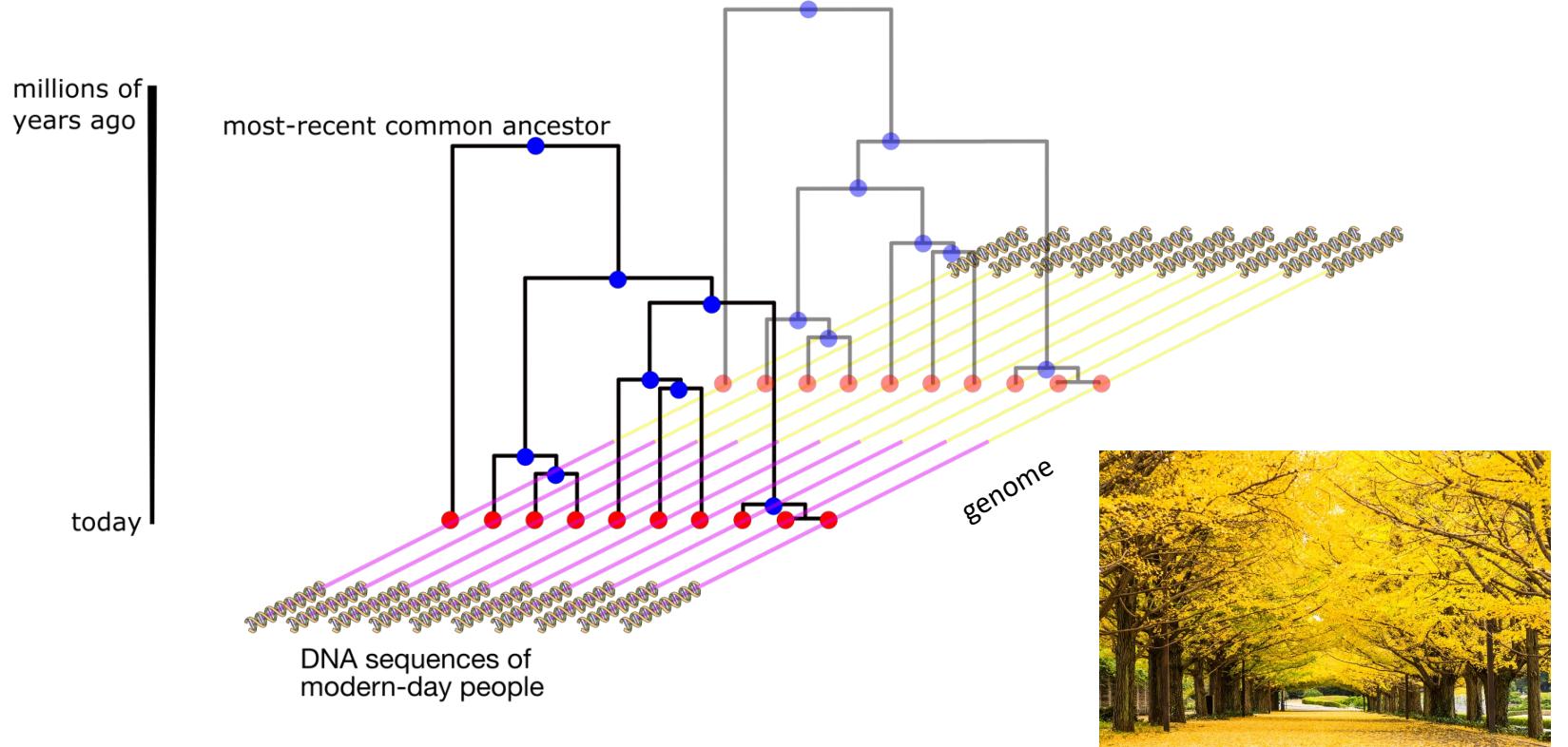
How does genetic variation arise and how is it maintained?

Big family tree?



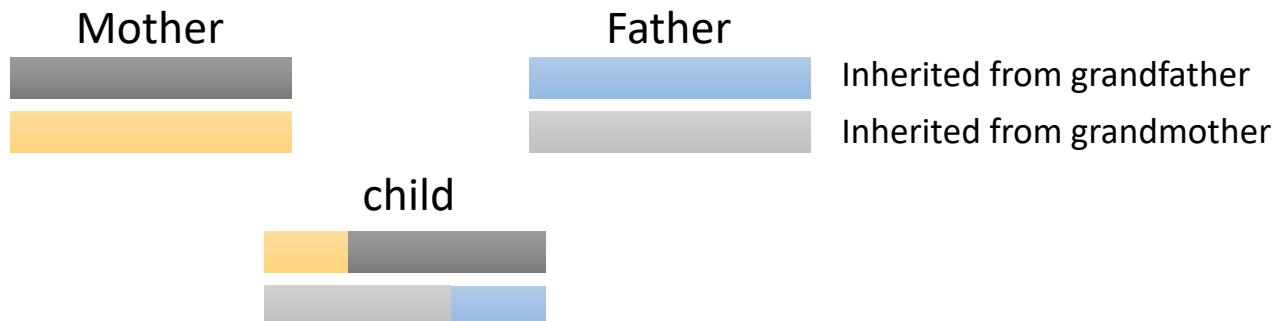
Key concept: Genealogies

We are genetically related through a sequence of trees



Reason for why trees change along the genome

Recombination:



Processes impacting tree shapes and genetic diversity

Genetic variation is shaped by

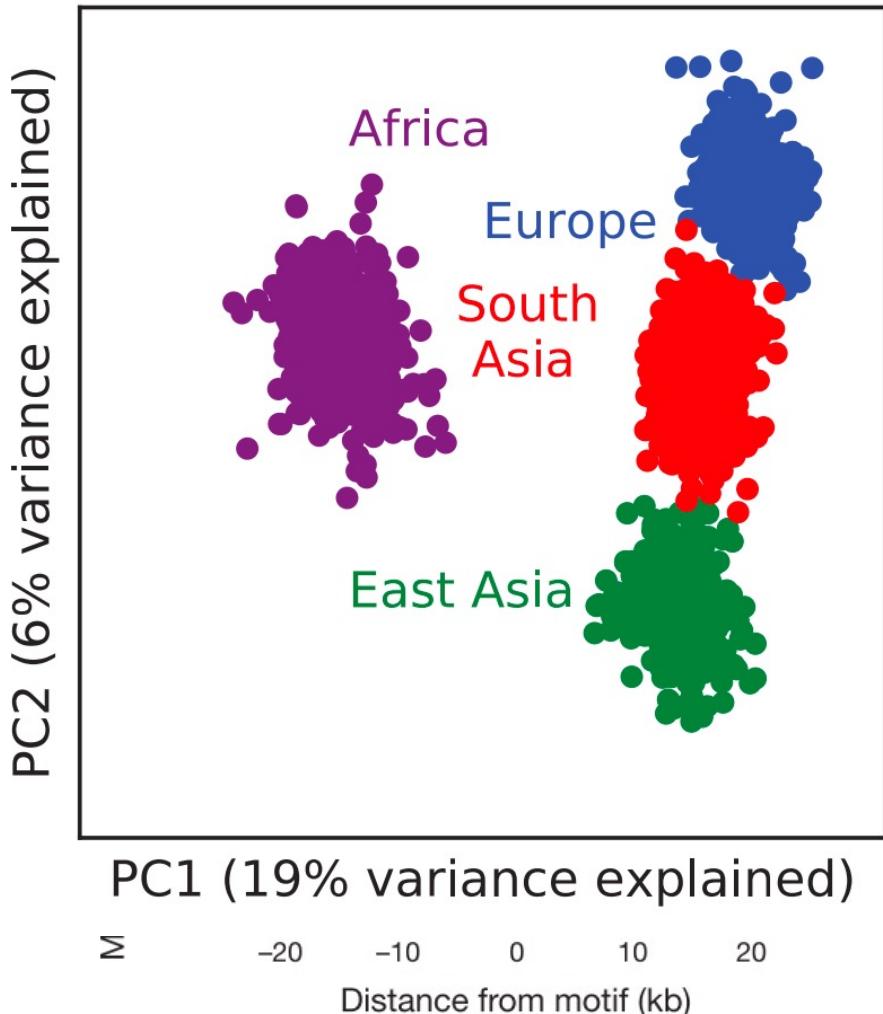
- Structure and migrations
- Population bottlenecks
- Admixture
- Mutation
- Recombination
- Selection
- Etc



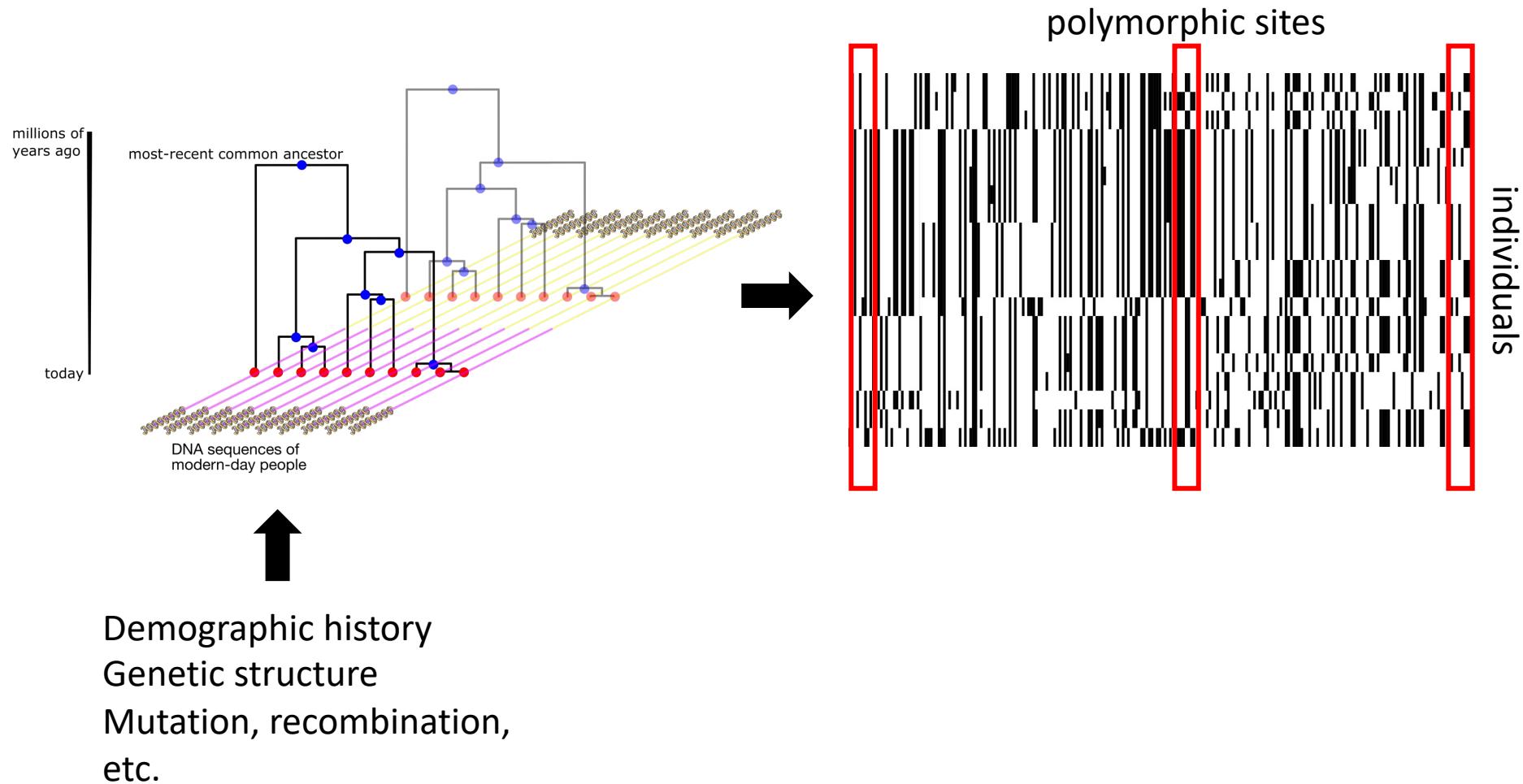
These evolve over time!
“evolution of evolution”

Harris et al, eLife 2017

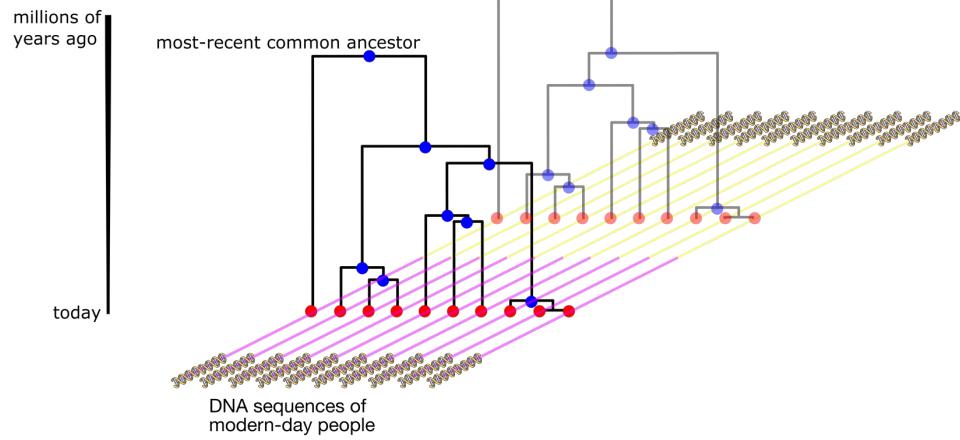
A. PCA of human mutation spectra



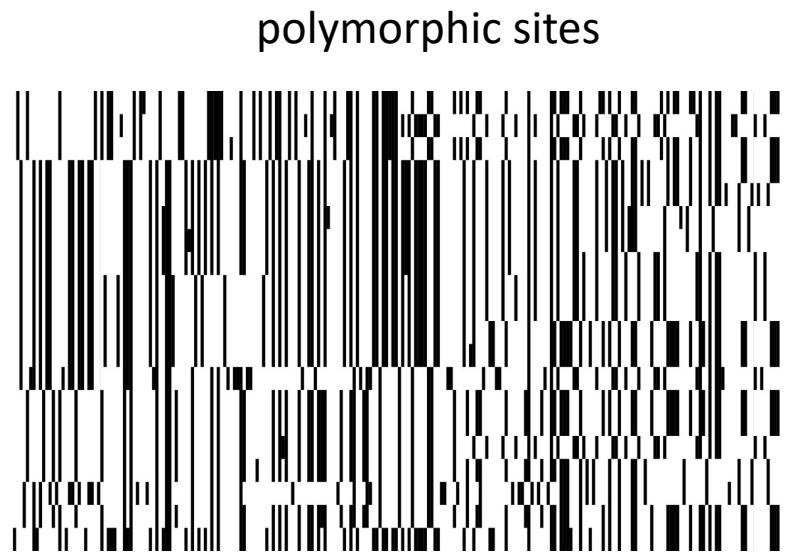
Fundamental forces impact data (only) through underlying genealogies



Many canonical approaches

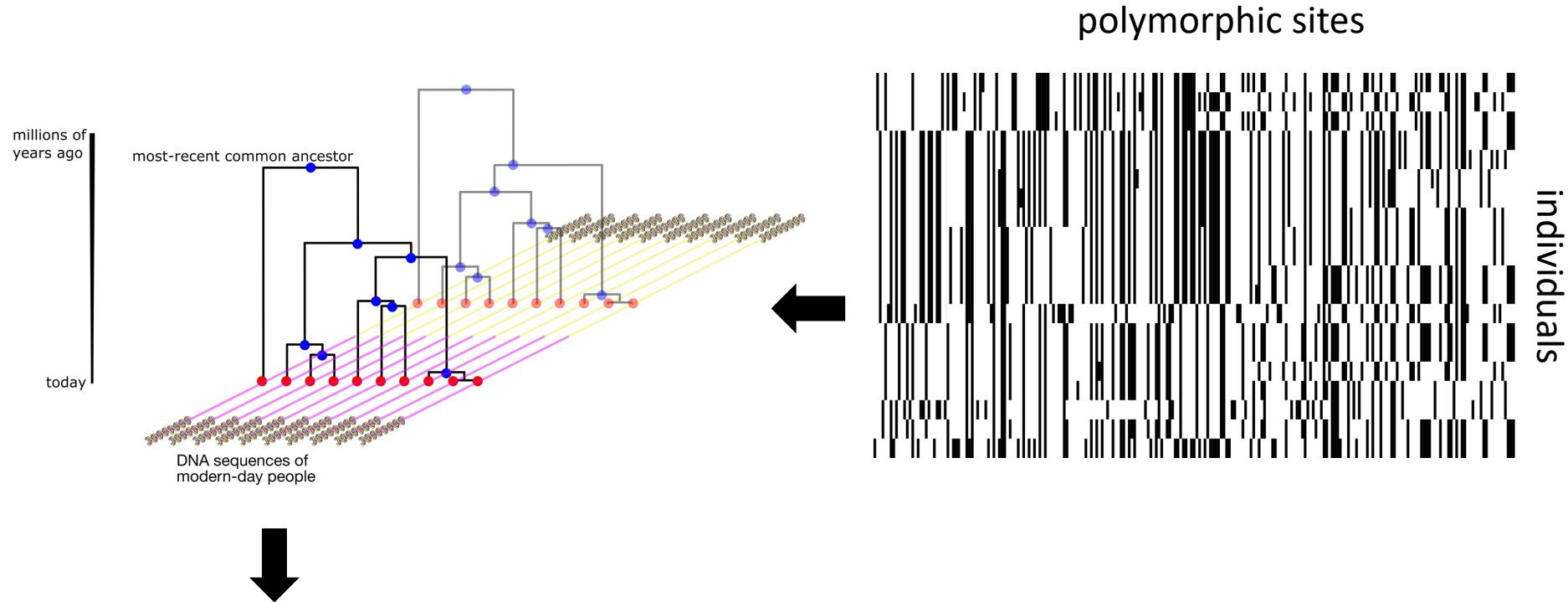


Demographic history
Genetic structure
Mutation, recombination,
etc.



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

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



Demographic history
Genetic structure
Mutation, recombination,
etc.

Inferred trees reorganise all the information available from the data about these processes in a very accessible & powerful way

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

Inferring genealogies

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

- ARGweaver] Infers Ancestral Recombination Graphs
- Rent+
- Tsinfer + tsdate] Published in 2019/2021,
• Relate] scale to large sample sizes

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

Relate

L. Speidel, M. Forest, S. Shi, S. Myers. Nature Genetics 2019

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



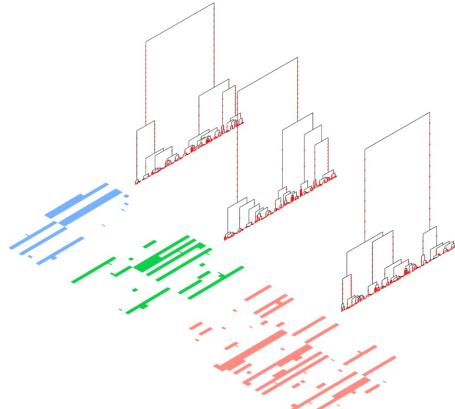
Software to estimate genome-wide genealogies for thousands of samples

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.

Citation: 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.

Contact: leo.speidel@outlook.com

Website: <https://leospeidel.wordpress.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: 10/01/2020)

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

Linux (x86_64, dynamic) - v1.0.17

Linux (x86_64, static) - v1.0.17

Mac OSX - v1.0.17

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.

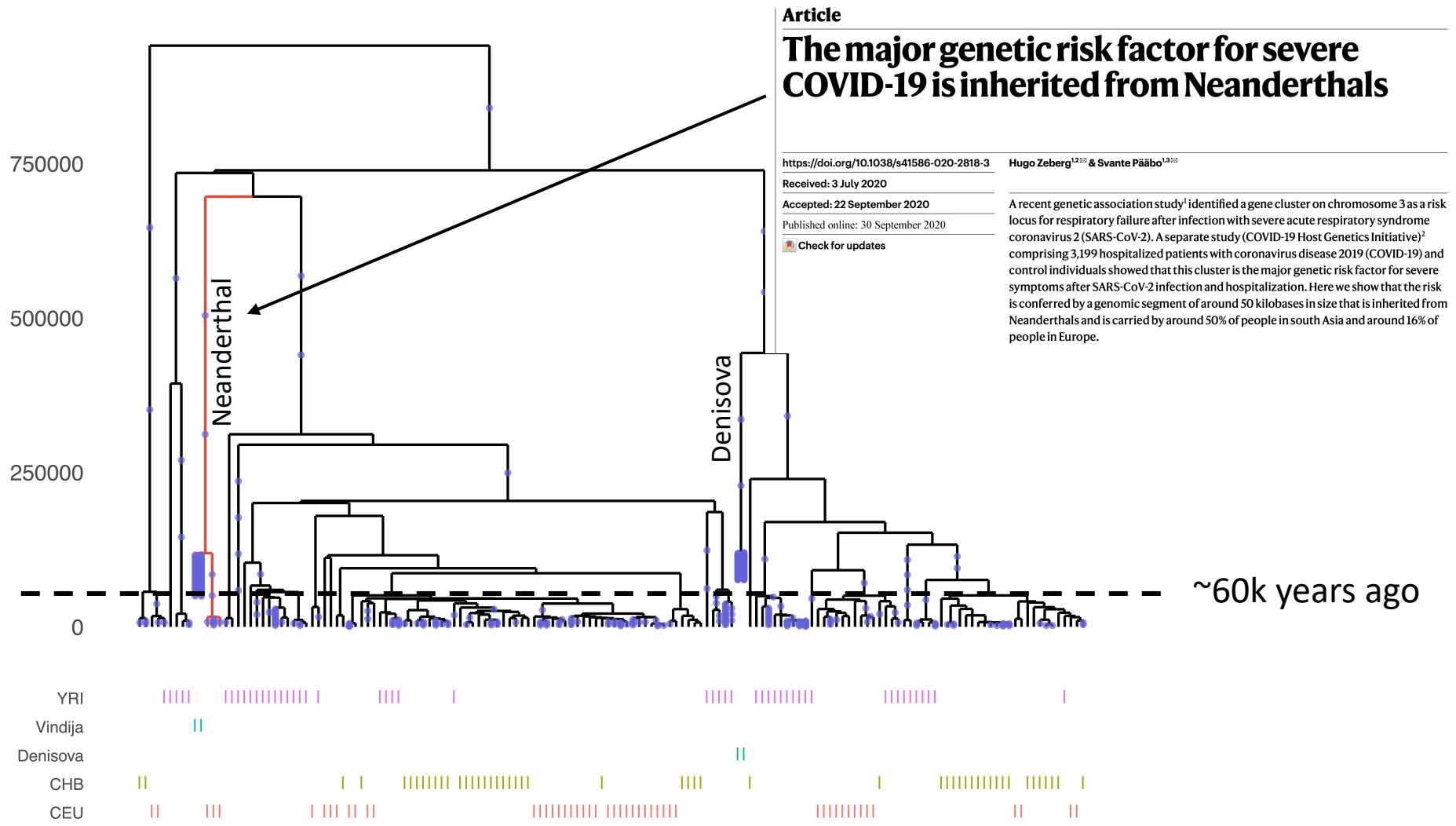
If you have any problems getting the program to work on your machine or would like to request an executable for a platform not shown here, please send a message to leo.speidel [at] outlook [dot] com.

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

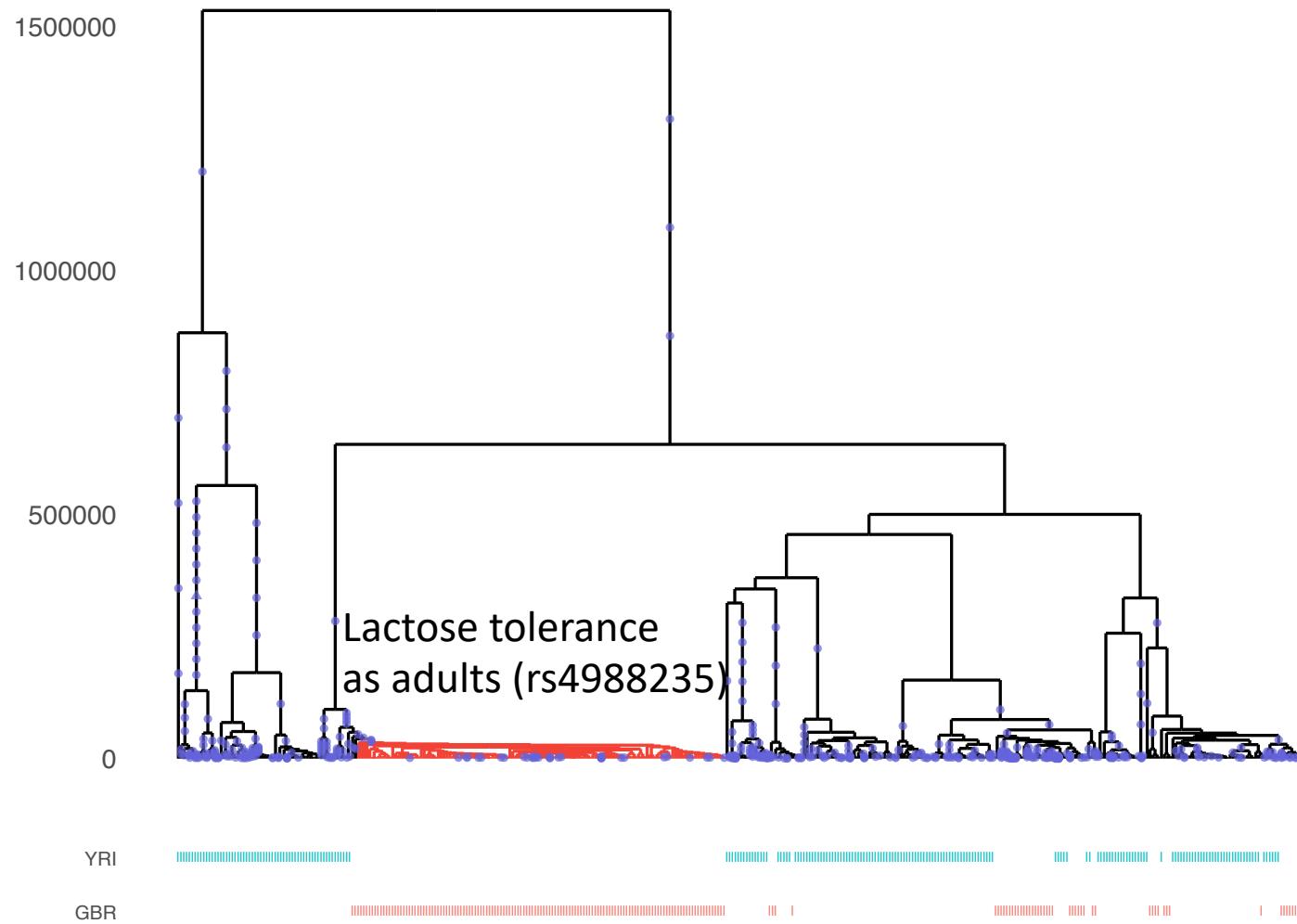
Key features:

- 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

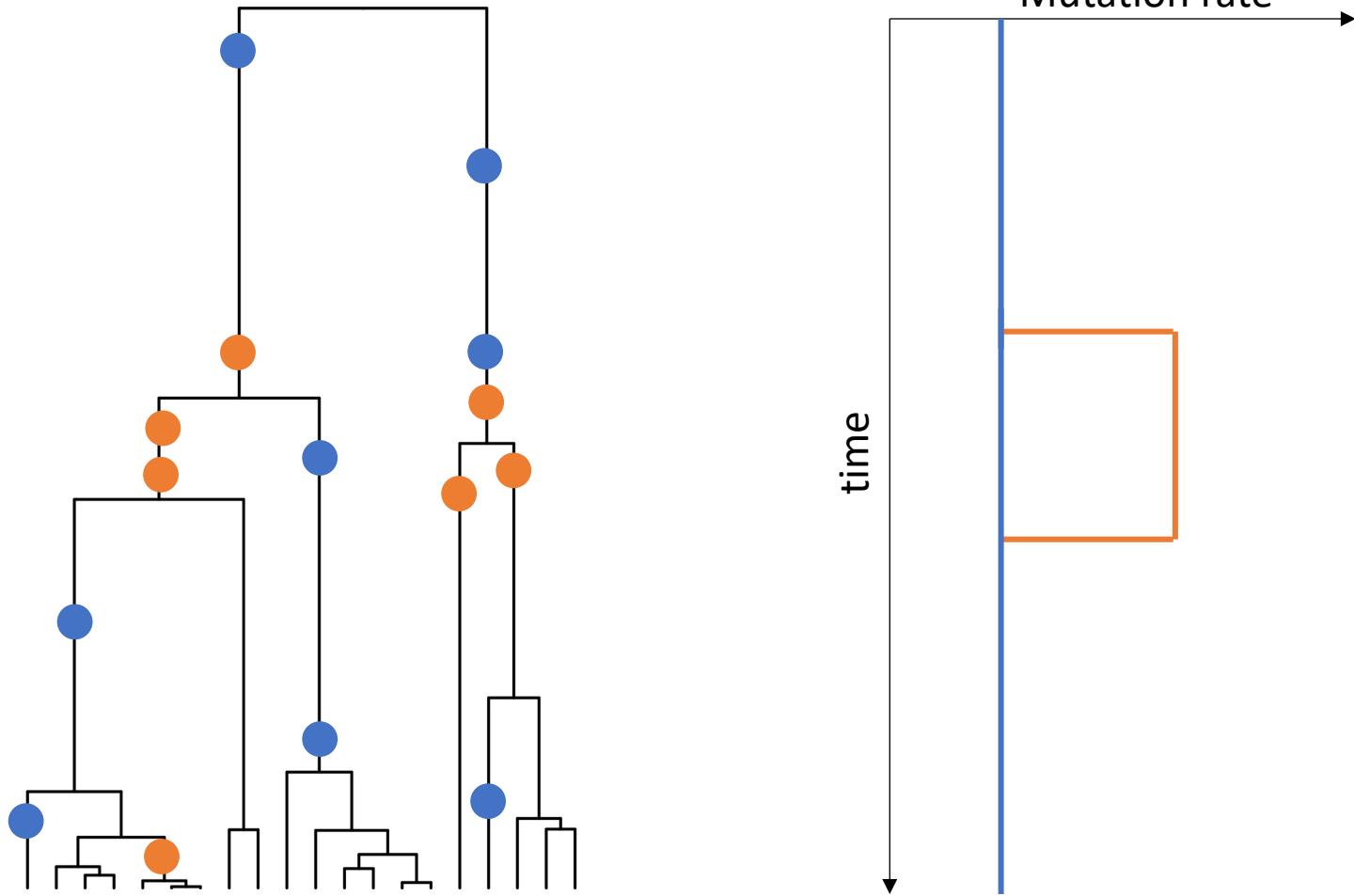
One locus can already tell us a lot about our history

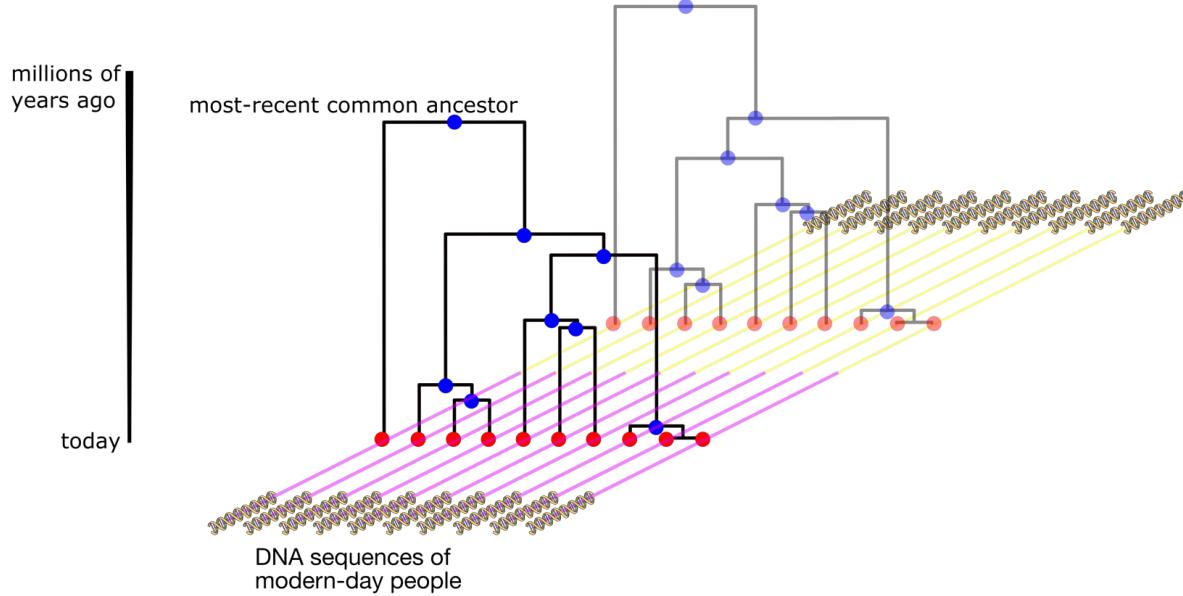


Positive selection: rapidly spreading lineage



Clusters of mutations in time can capture changes in mutation rate





Inferring genealogies

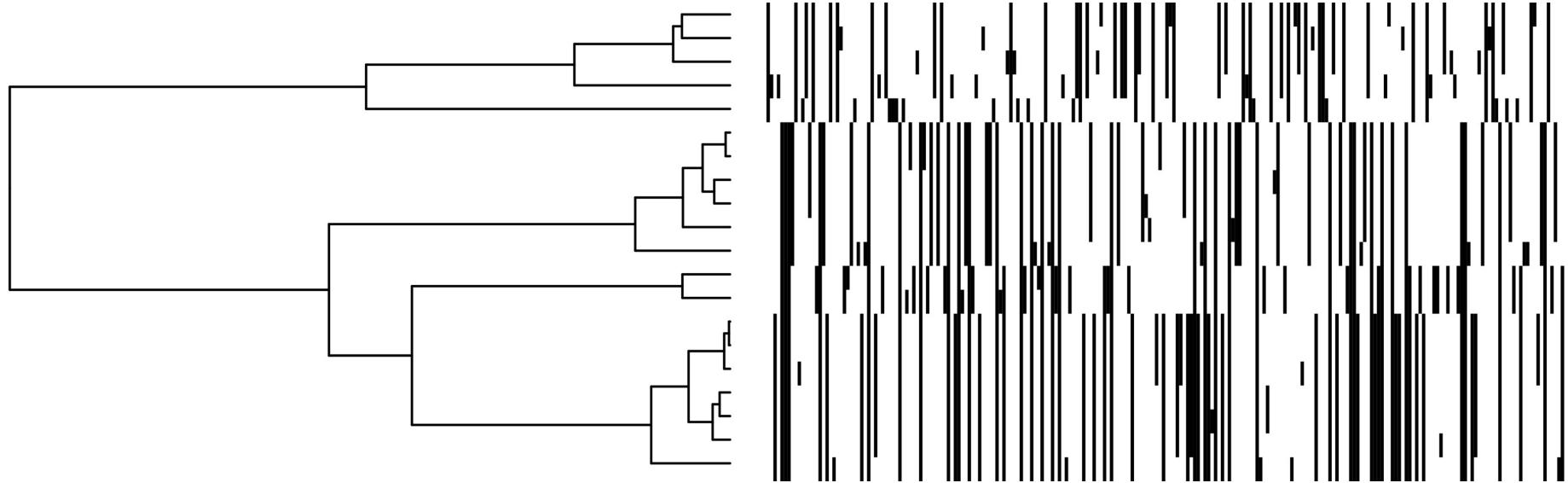
Old problem

- 90s – 2000s: Use the coalescent with recombination to infer genealogies, e.g. MCMC

$$P(\text{tree} \mid \text{Data}) \propto P(\text{Data} \mid \text{tree}) * \underbrace{P(\text{tree})}_{\text{coalescent}}$$
- Relate and tsinfer (2019):
 - Separate inference of tree topology and branch lengths
 - Use fast but approximate approach for tree topology
 - Use coalescent for branch lengths

Scalable to many thousands of samples

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

Use hidden Markov model to identify stretches where relationships to other are unchanged

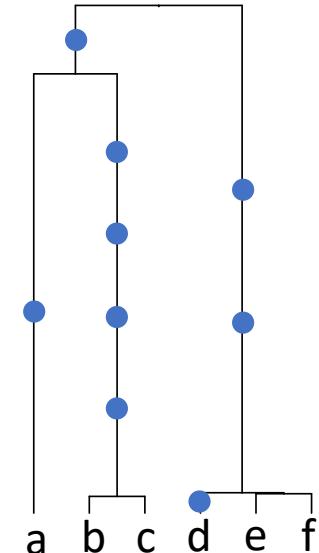
For tree topology, we want to quantify the order in which we are related to others

1. Count number of “derived mutations”

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

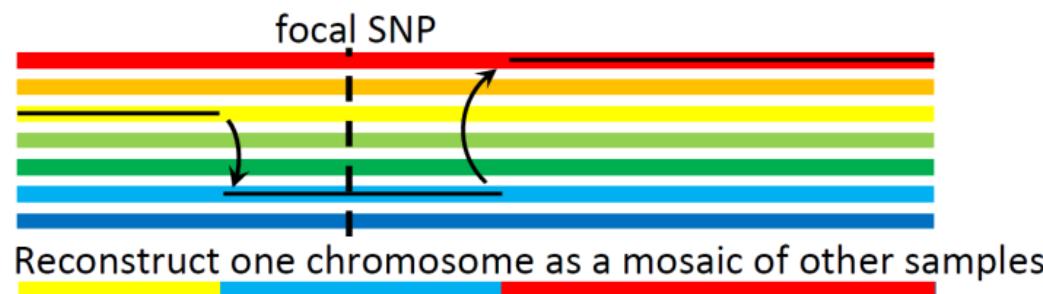
2. Coalesce mutually closest lineages

- 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 HMM to count derived mutations accounting for recombination



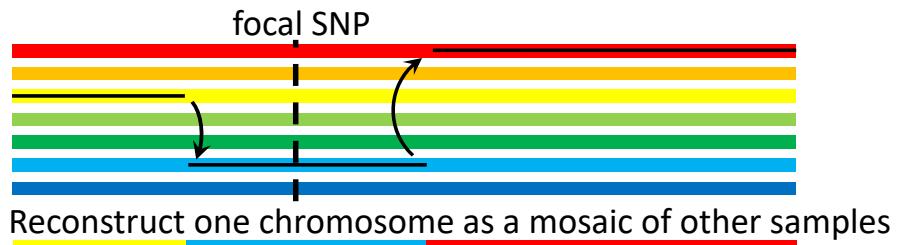
Hidden Markov model (HMM)

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

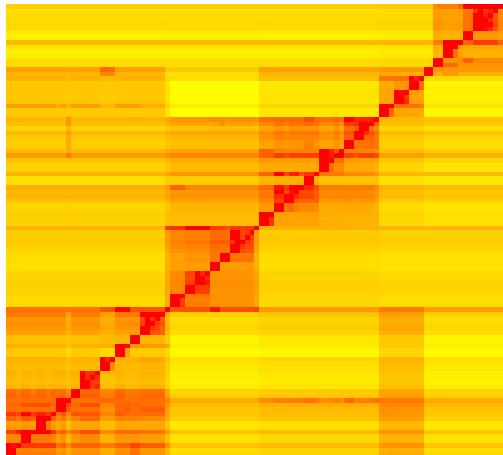


Summary of Relate

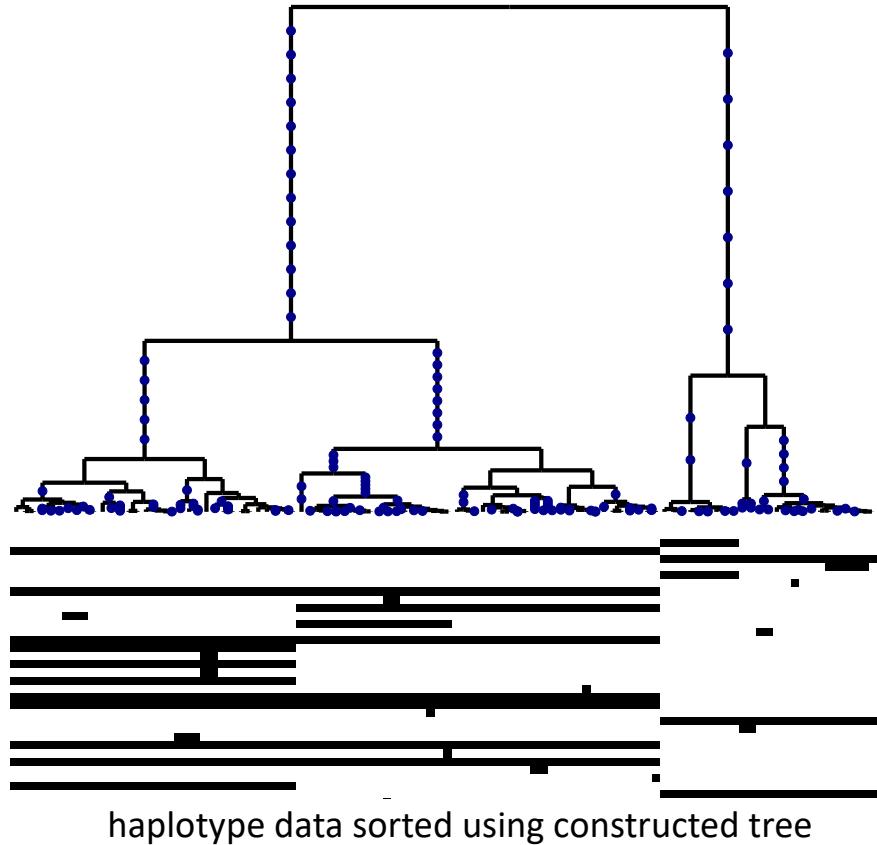
Hidden Markov model (HMM)



Distance matrix for focal SNP

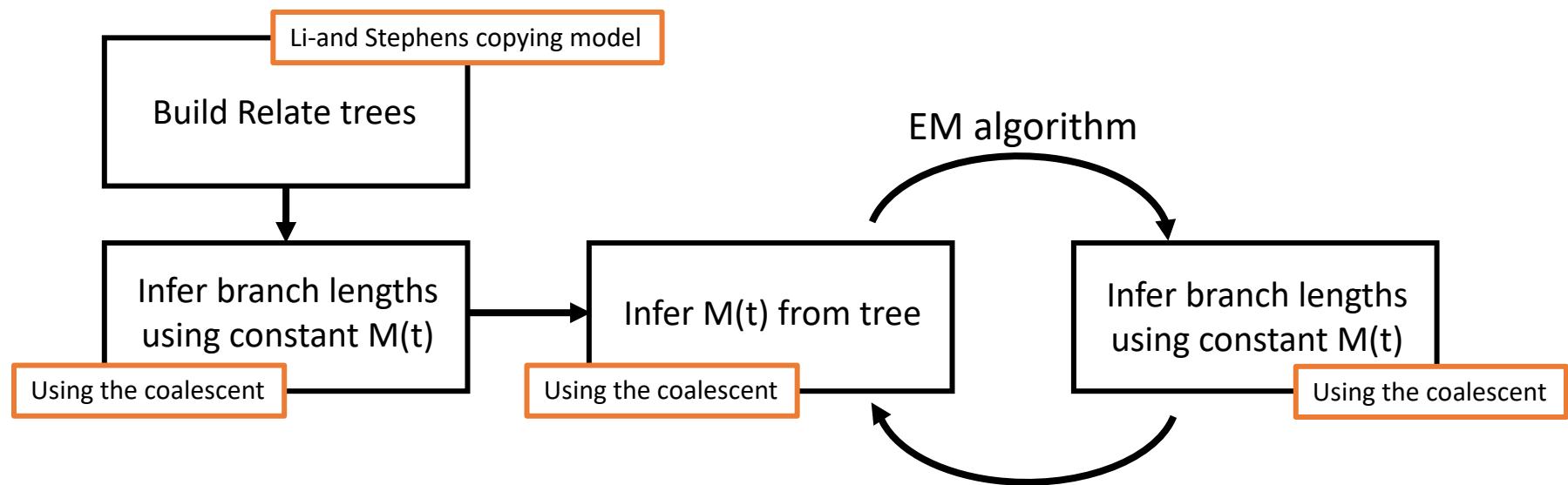
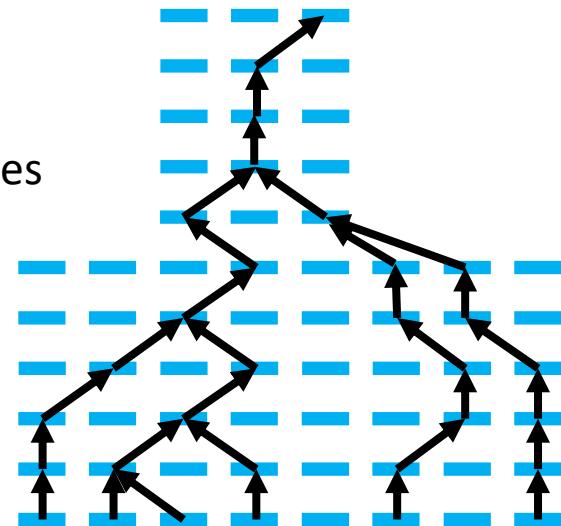


Hierarchical clustering
&
MCMC for branch lengths



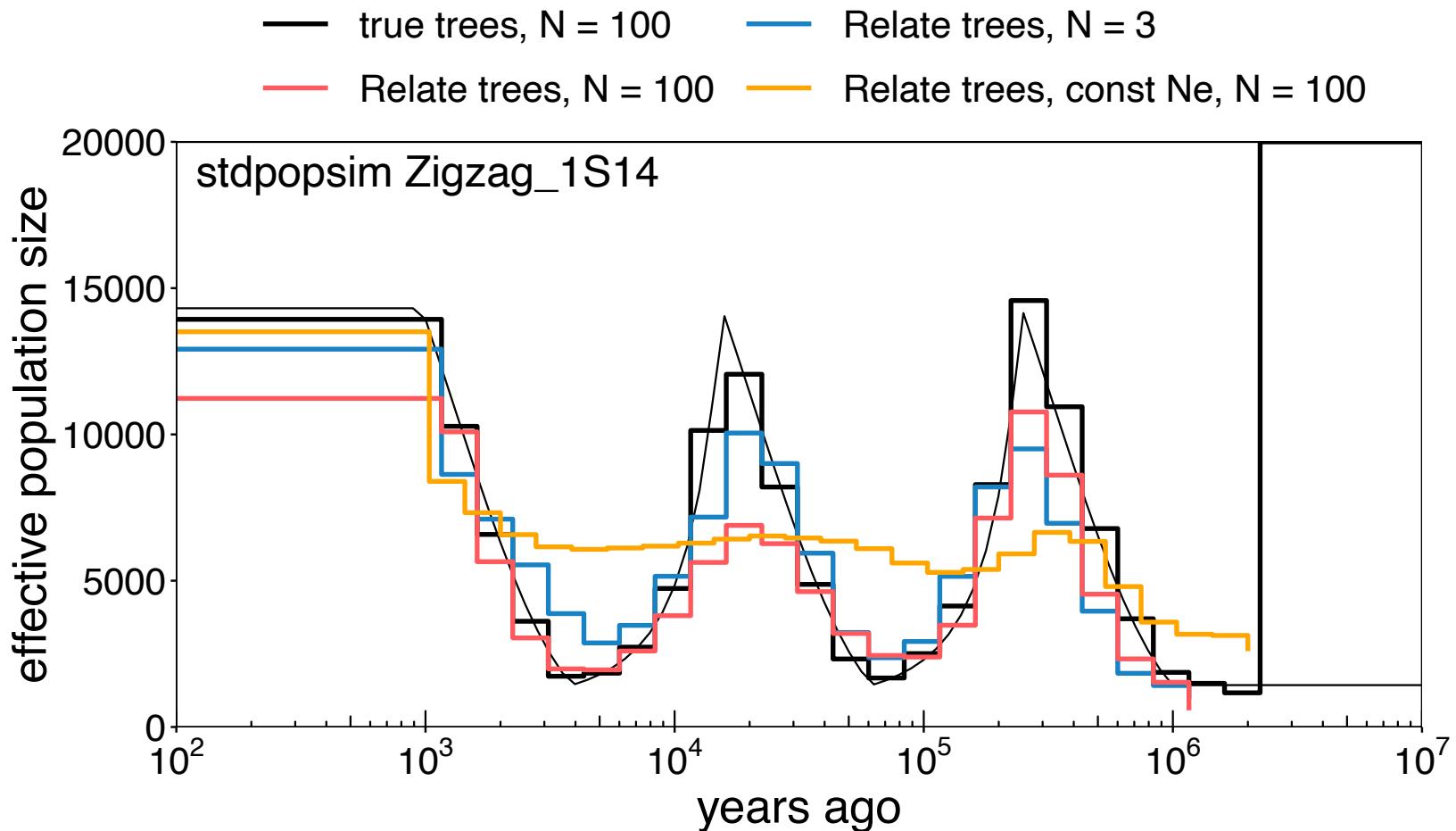
Branch lengths and population size is estimated jointly in an EM algorithm

- 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 $M(t)$

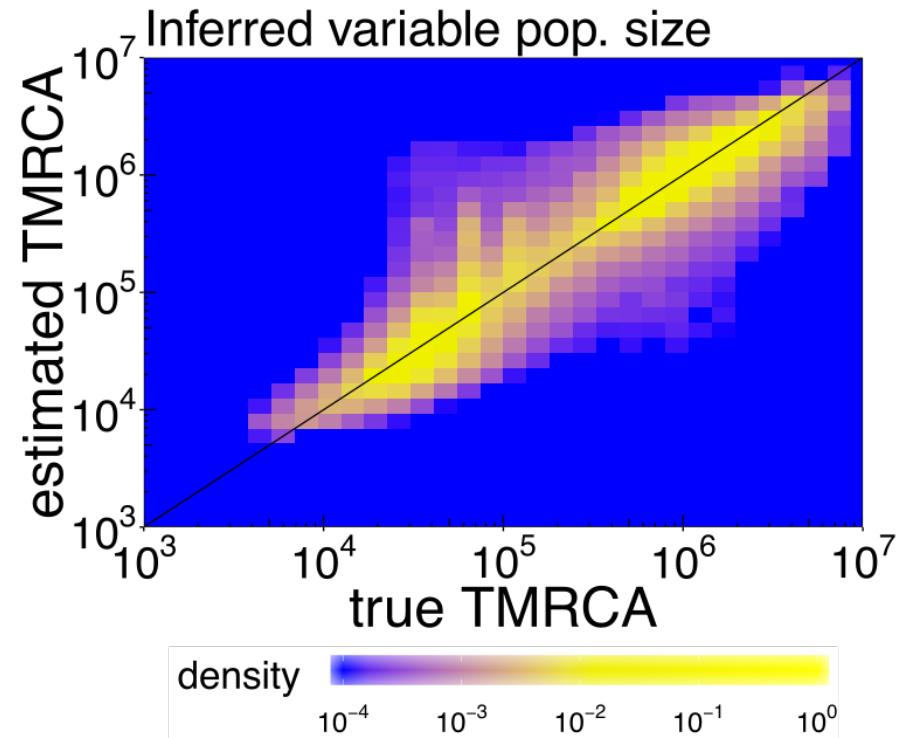
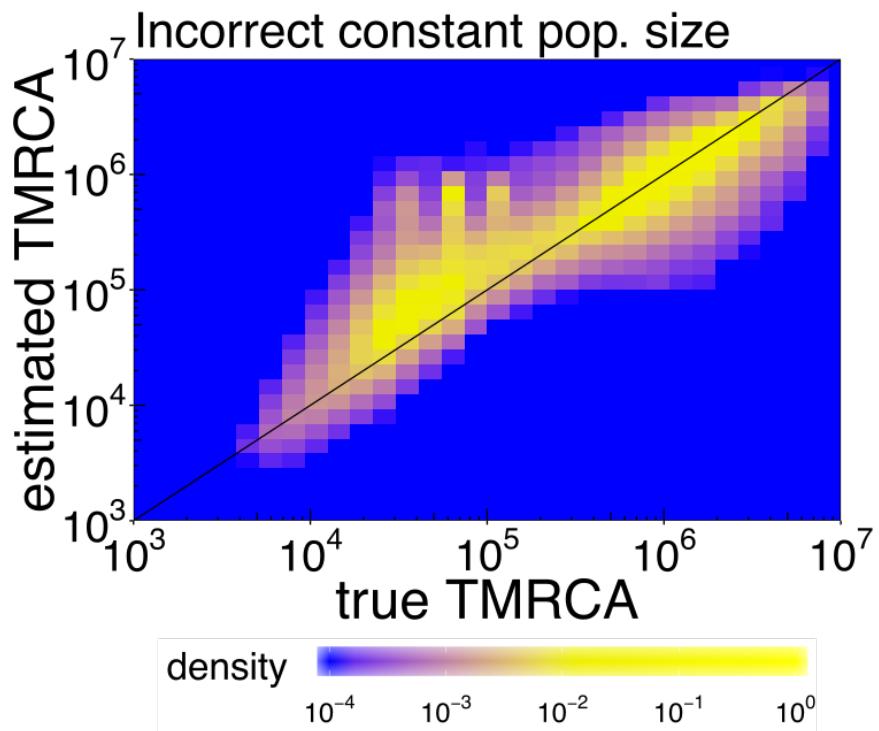


Population size changes through time are jointly inferred in Relate

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

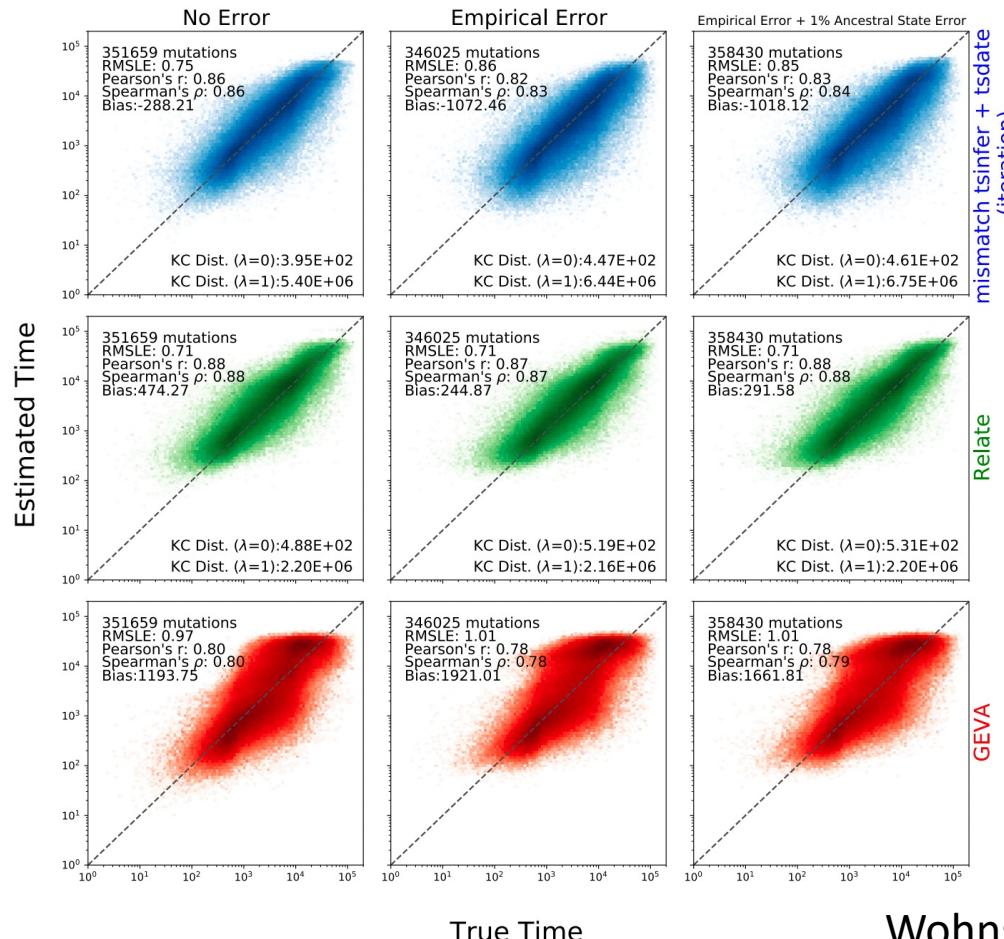


Simulated data with variable population size (European-like demographic history)



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



What Our DNA Can Tell Us About the History of Humans

Authors



Leo Speidel



Clare Bycroft

Young Reviewers



Mariana



Anna-Marie



Zara



Luckily not our reviews...

This seems important, but the way it is written is so boring. I can't even get to the end. Could the authors maybe sound excited about what they are doing?

Reviewer, Age 12

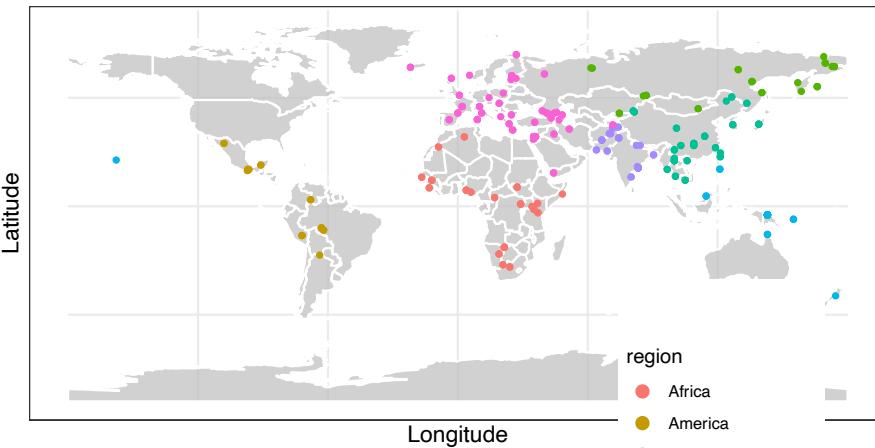
The writers of the article did not make it clear why such an expensive and involved research project was done to begin with ... It seemed like a fruitless task.

Reviewer, Age 14

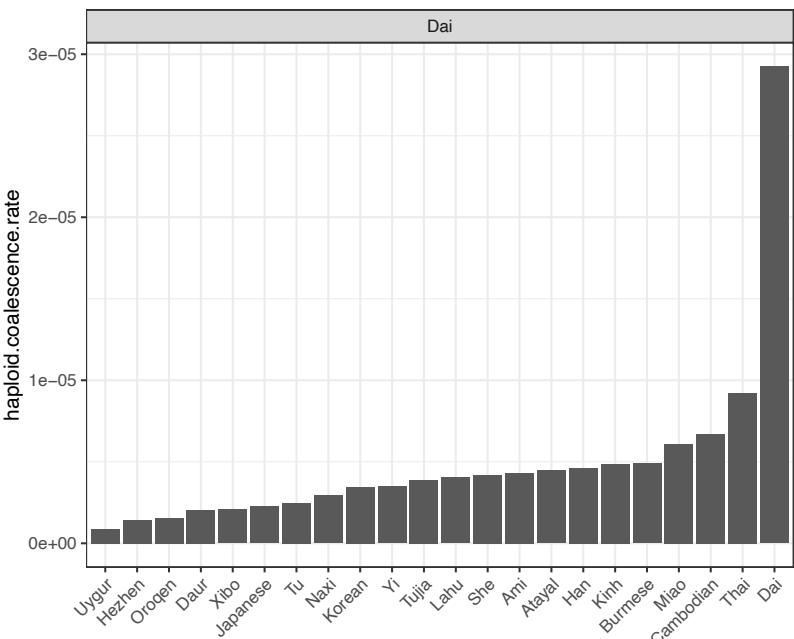
Genealogy-based inference of human evolutionary history

Coalescence rates track relatedness through time

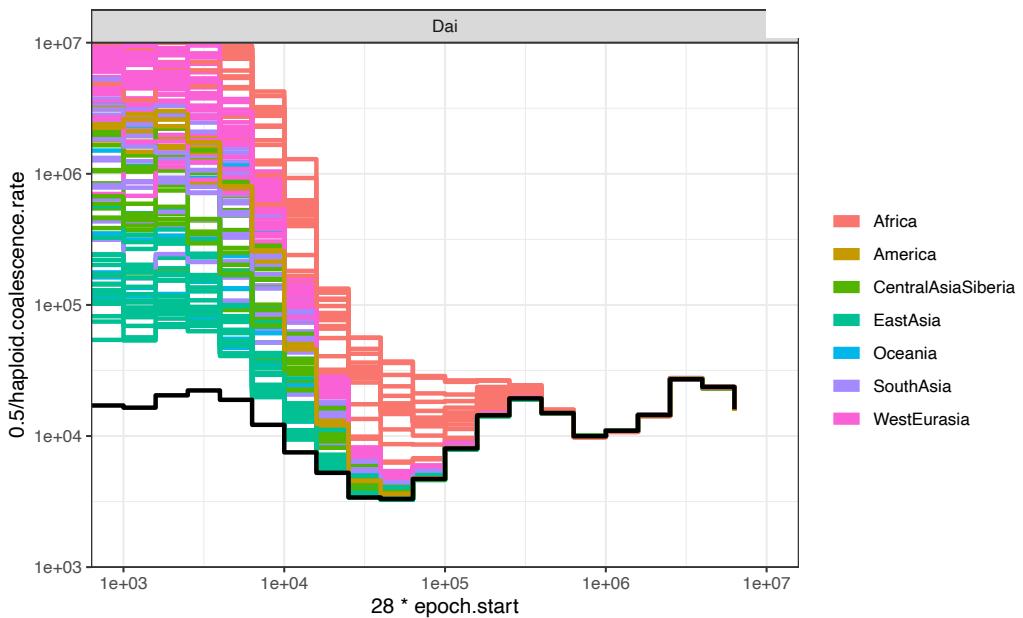
Coalescence rates = $1/M(t)$
 $M(t)$: effective population size



Relatedness < 1000 years in East Asia

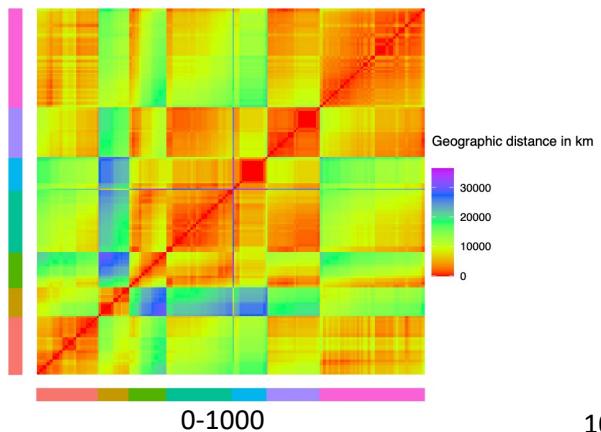


Relatedness through time globally



Structure through time

Geographic distance

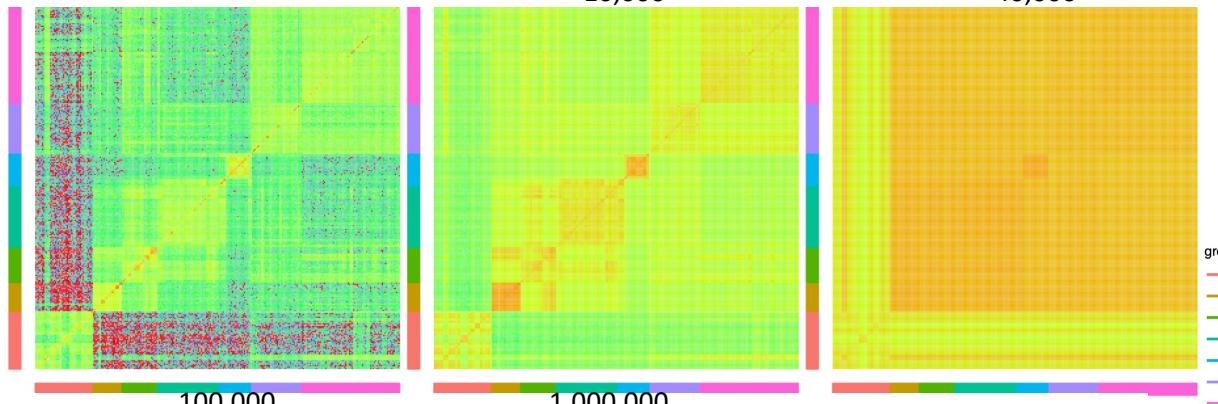


Latitude

Longitude

- region
- Africa
 - America
 - CentralAsiaSiberia
 - EastAsia
 - Oceania
 - SouthAsia
 - WestEurasia

Coalescence rates

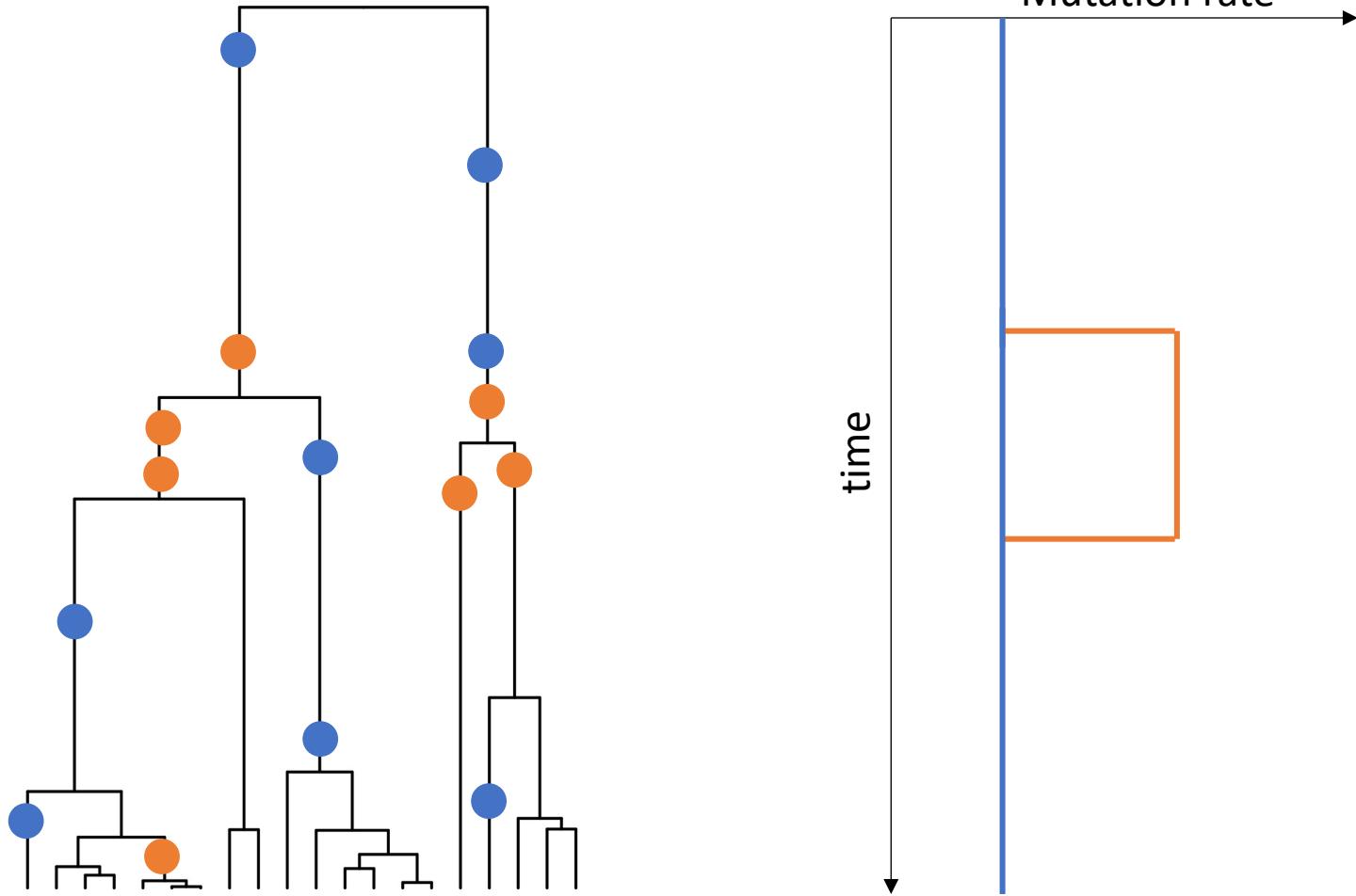


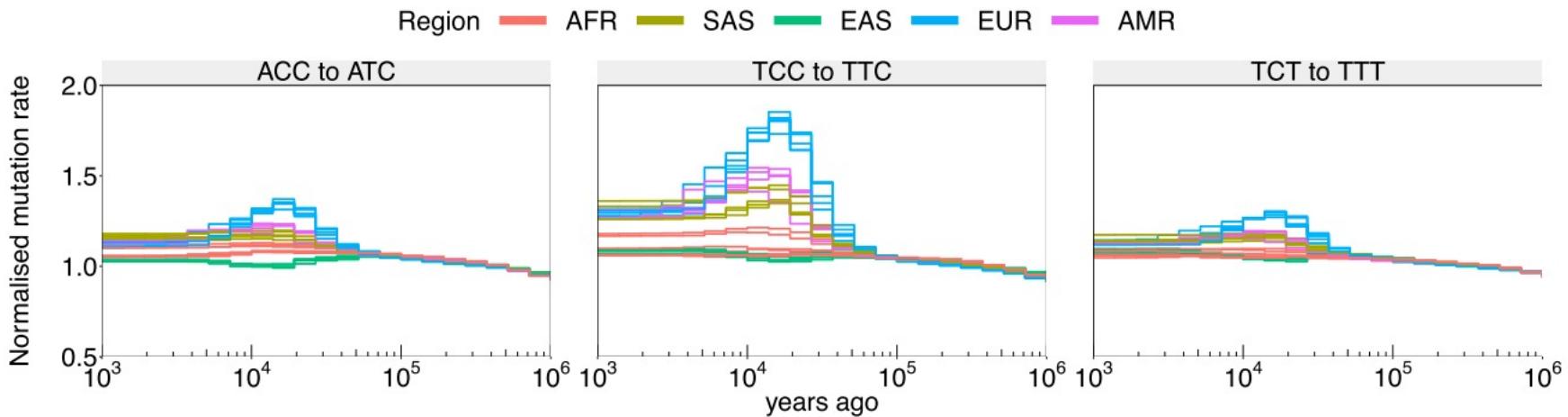
- group
- Africa
 - America
 - CentralAsiaSiberia
 - EastAsia
 - Oceania
 - SouthAsia
 - WestEurasia

- coal rate
- 1e-03
 - 1e-06
 - 1e-09

$$\text{Coal rate} = 1/M(t)$$

Clusters of mutations in time can capture changes in mutation rate





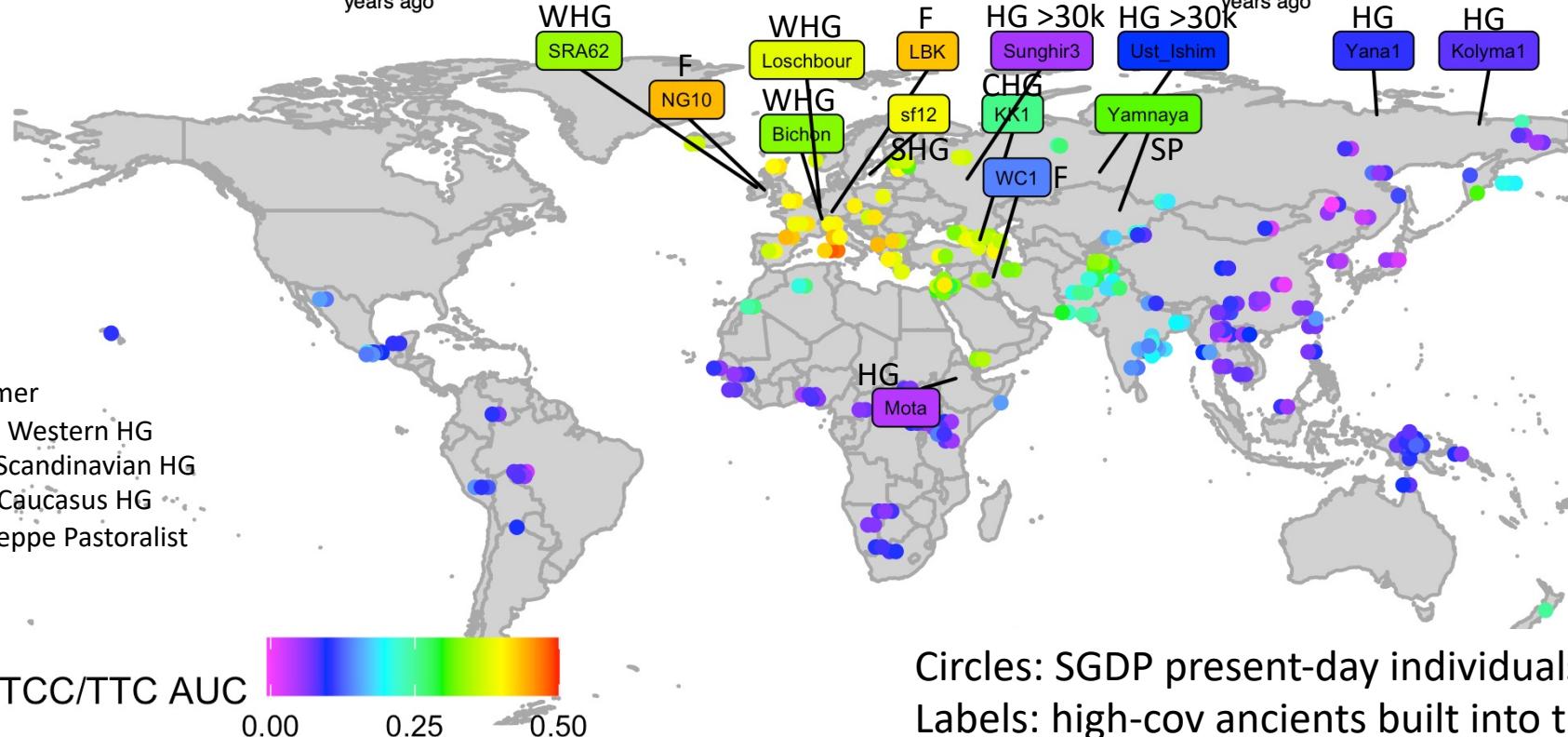
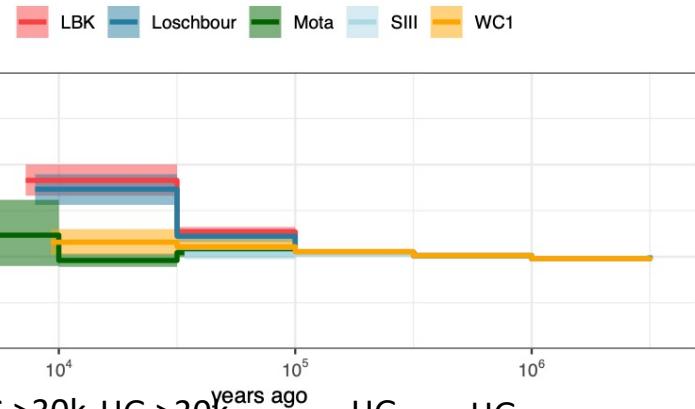
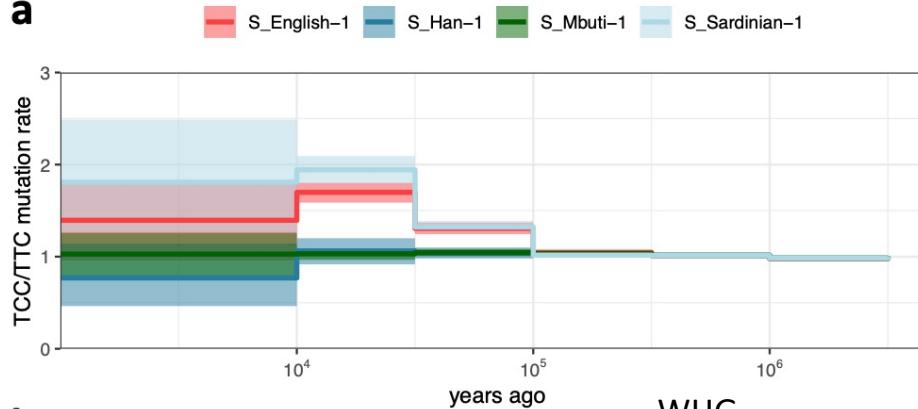
“Mysterious” TCC/TTC mutation rate in West Eurasia

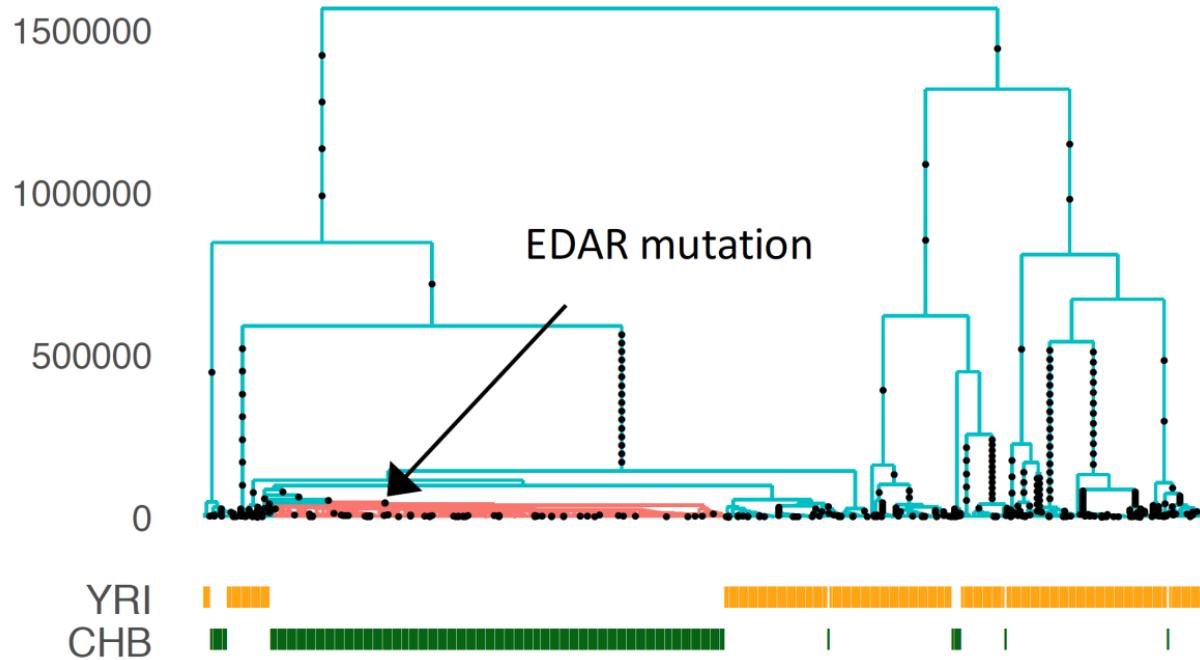
- First reported by Kelley Harris (PNAS 2015, eLife 2017)
- Strongest signal of all triplets and strongest in Europeans among modern-day people
- Unknown cause (genetic?, environmental?)
- Date of onset unclear (dates range from 10k – 80k)
- Studied mainly in moderns, handful of ancients (Mathieson & Reich 2018),
but detailed geographic & temporal pattern was unexplored

TCC/TTC mutation rate inferred from Relate genealogies

Speidel et al. bioRxiv, 2021

a



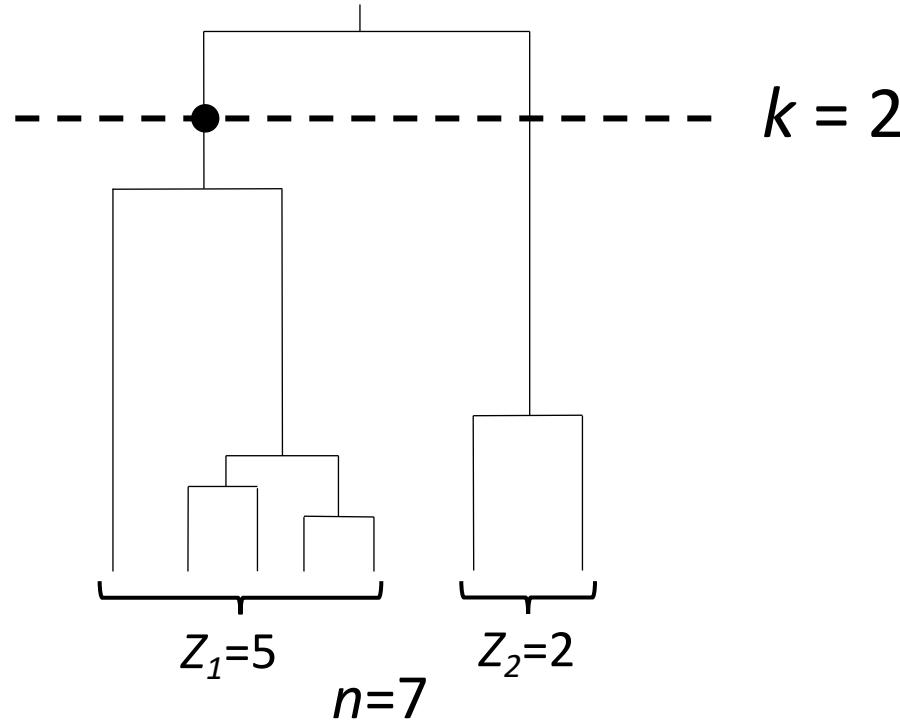


Quantifying positive natural selection on a single mutation

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

Reject a null model

How quickly does a mutation spread in the neutral case?



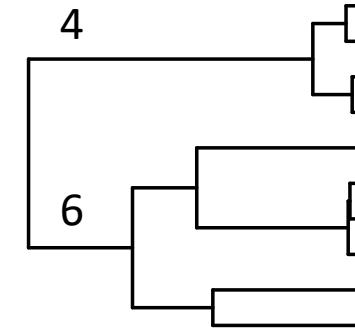
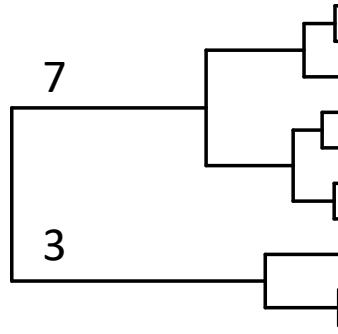
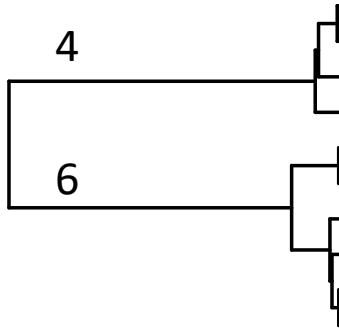
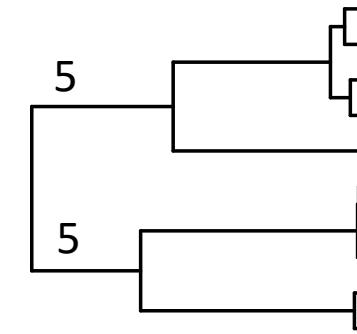
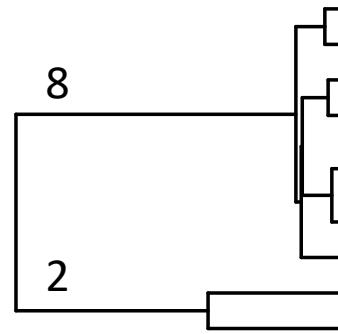
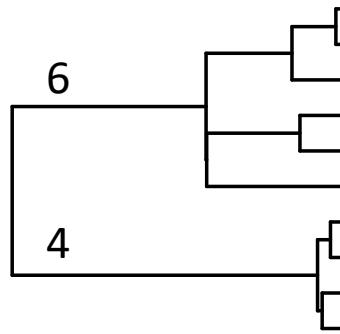
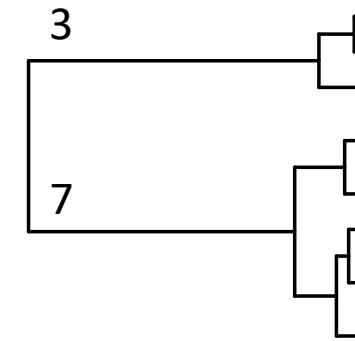
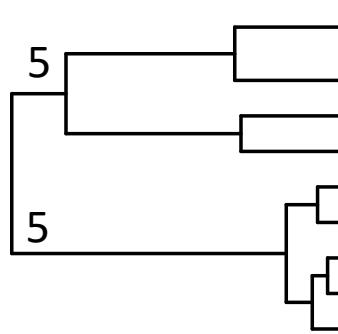
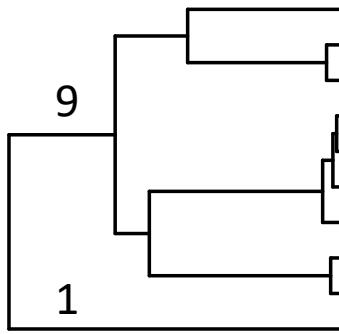
We can write down the analytical distribution for the number of descendants of a mutation arising while k lineages remain

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!

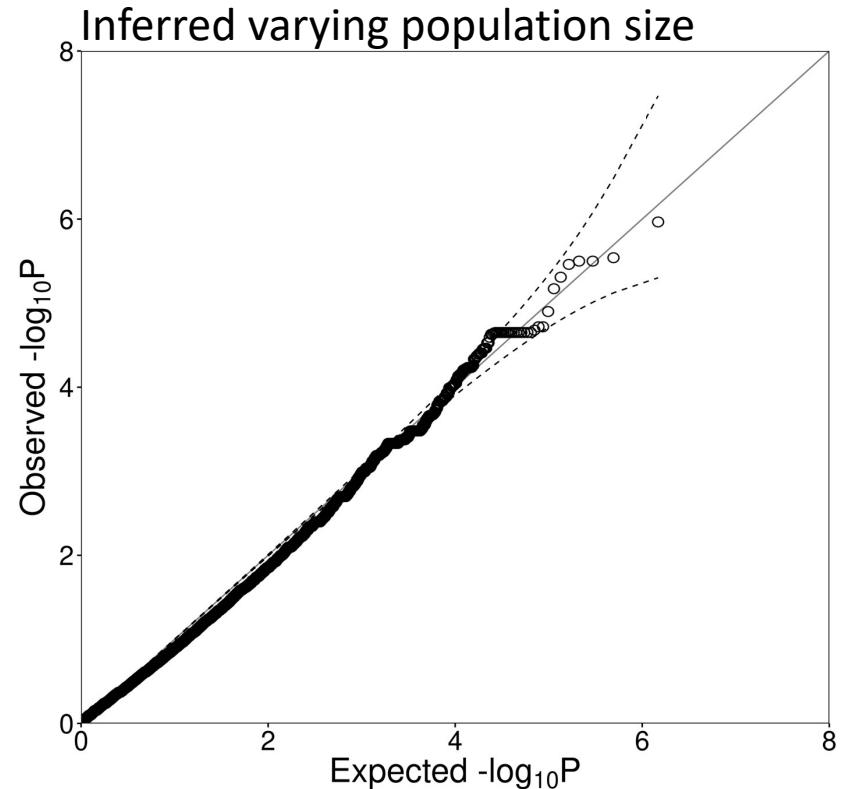
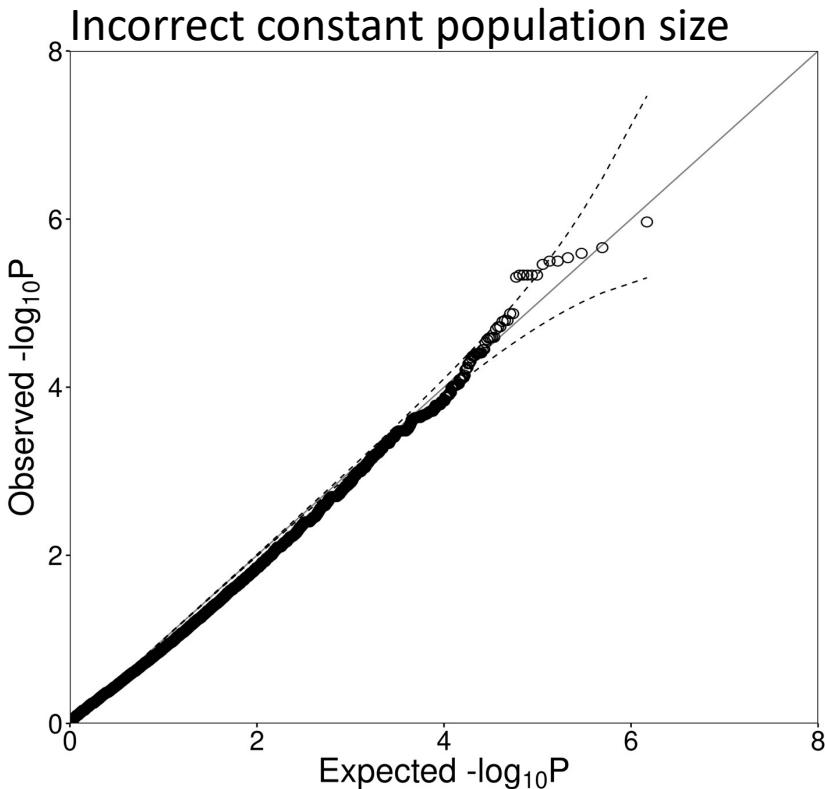


P-values: very well calibrated under null simulations of no selection

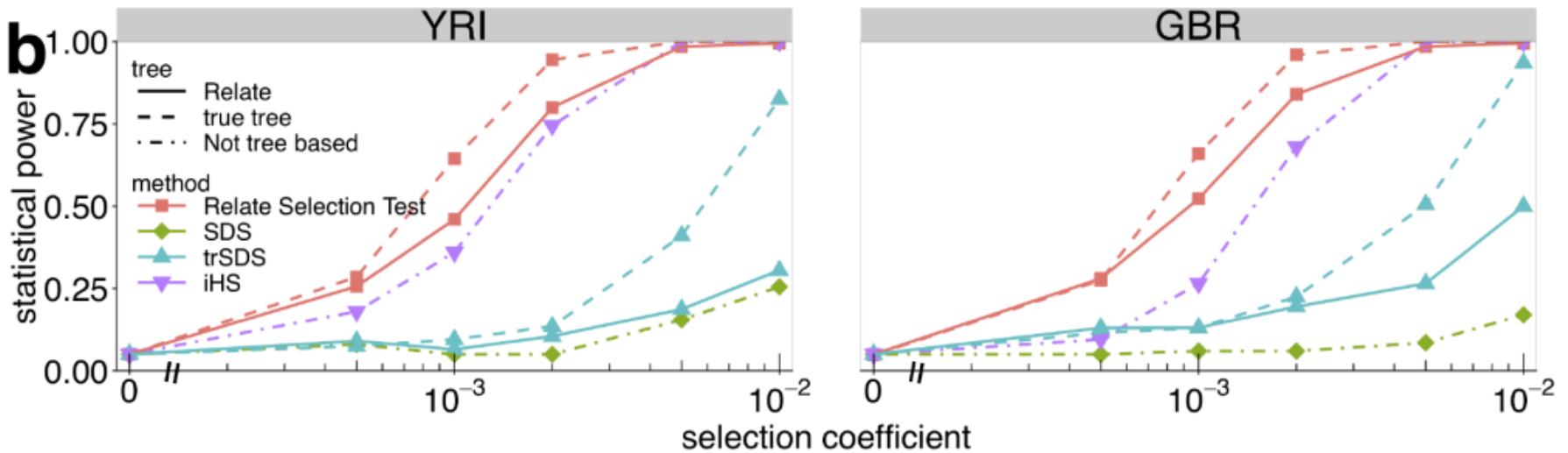
N=1000, 250Mb

Bottleneck population size

Quantile-quantile plot:

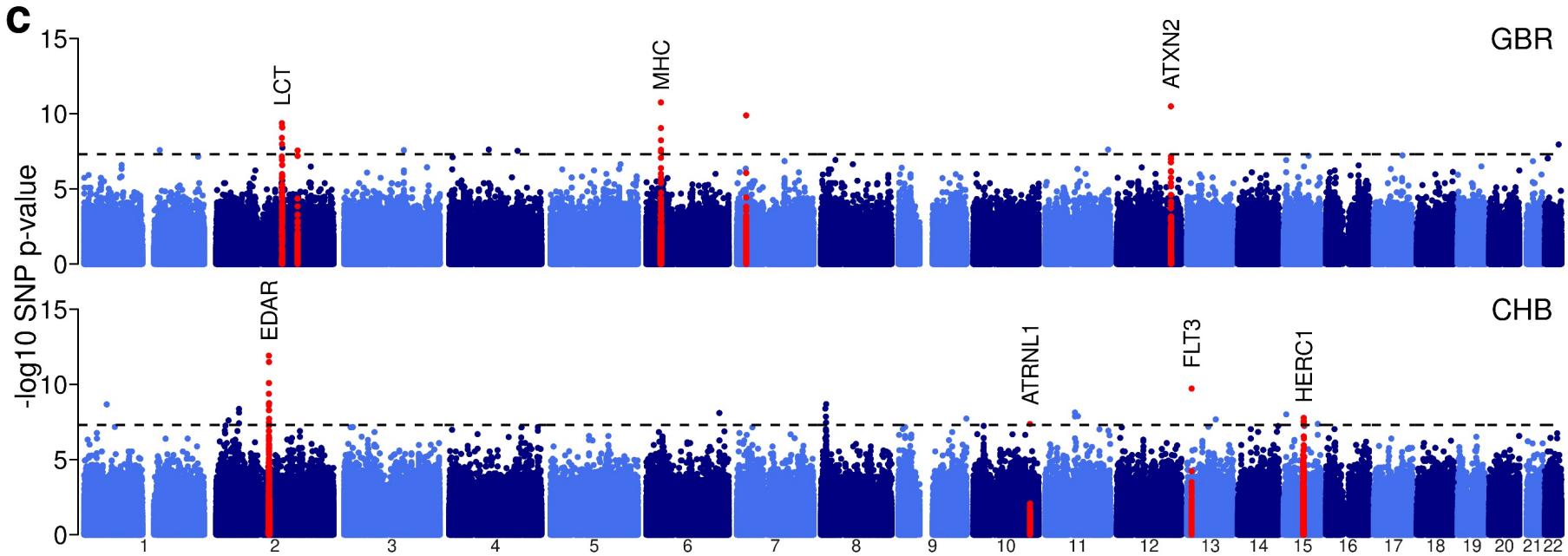


Improved power to see weak selection



Genome-wide selection p-values

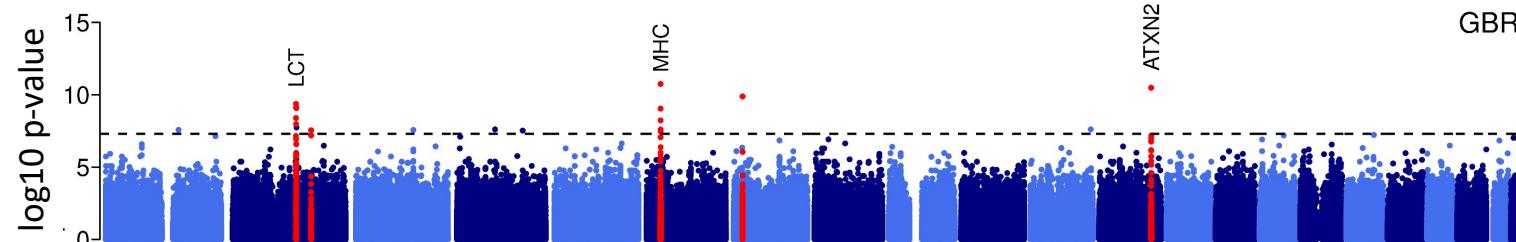
Given most traits are highly polygenic, expect mainly weak, polygenic selection



How does weak selection evidence vary by trait?

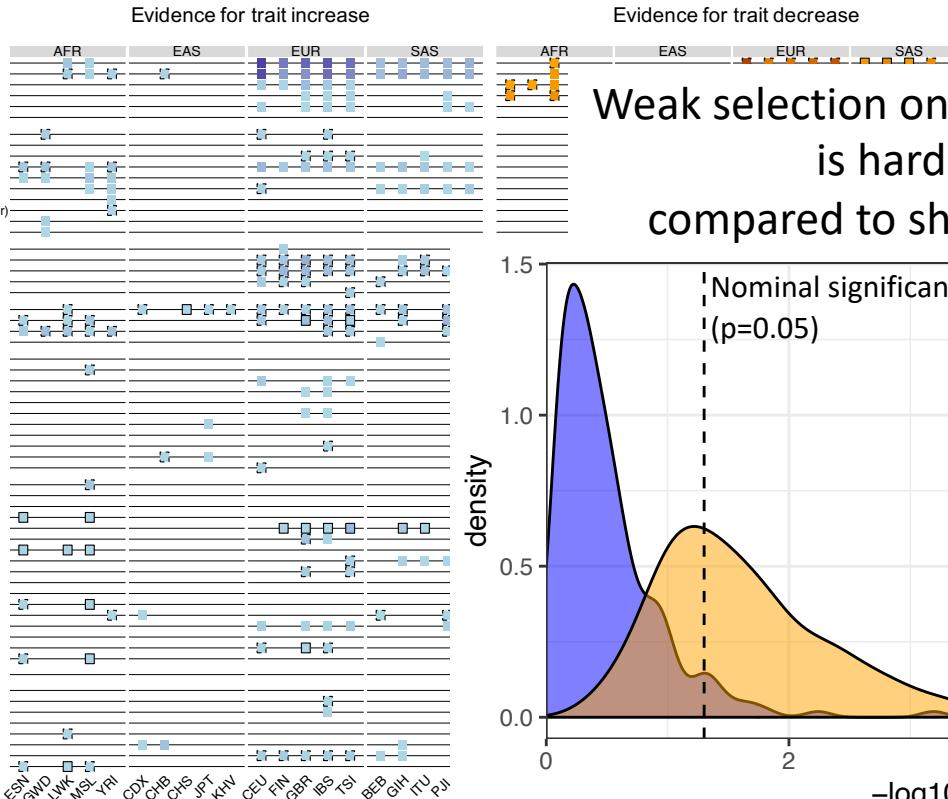
Many key events in our evolutionary history are only implicated as subtle effects in our genomes

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

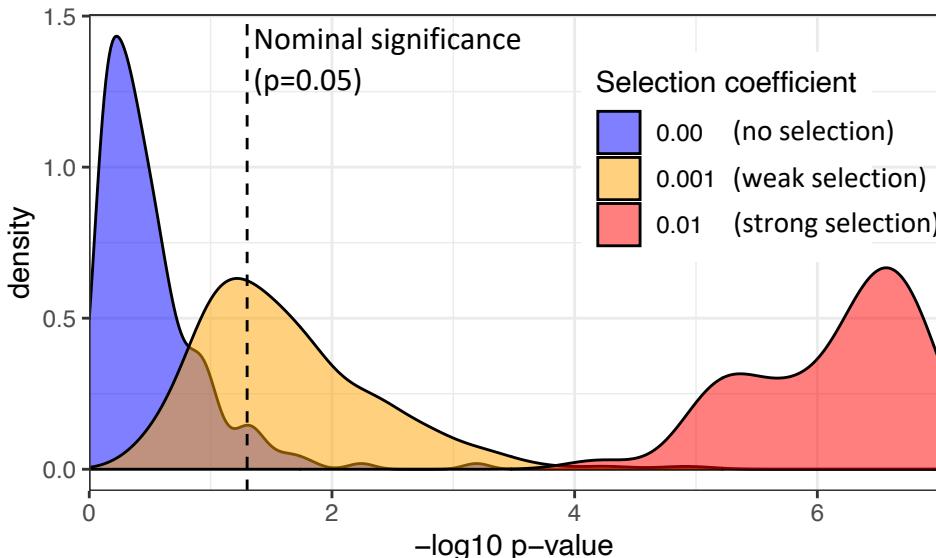


Physical traits	Sitting height (UKB) Standing height (UKB) BMI (UKB) Hip circumference (UKB) Waist circumference (UKB) Skin colour (UKB)
	Height BMI BMI (adj. for smoking behaviour) Hip circumference Hip circumference adj. for BMI Waist circumference Waist circumference adj. for BMI in active individuals 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 count Granulocyte count Basophil counts Neutrophil count
Lipids	Cholesterol, total Blood metabolite ratios Blood metabolite levels Blood glucose Glomerular filtration rate (creatinine) Glomerular filtration rate in non diabetics (creatinine) Gut microbiota (bacterial taxa) Educational attainment (years of education) Schizophrenia (PGC)
Other traits	

Lots of Polygenic selection signals

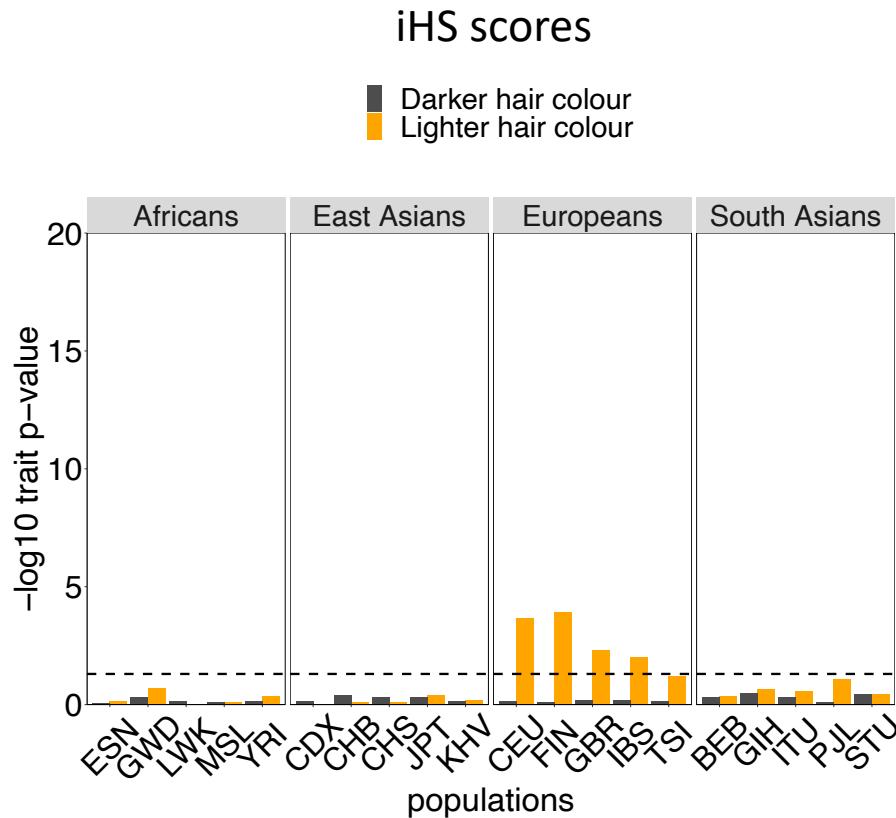
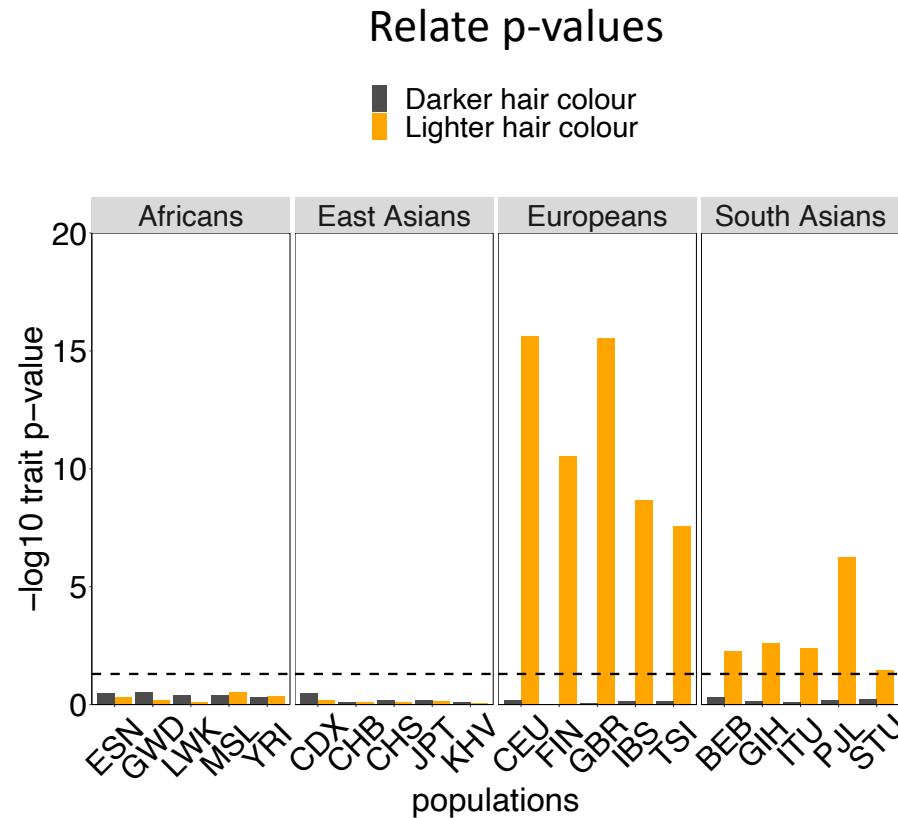


Weak selection on individual mutations
is hard to detect,
compared to shifts in distributions



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)



Summary & future work

- It is now possible to build genealogical trees for huge datasets, in humans and other species (currently 10,000 people or more)
- 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:
 - Many existing approaches are based on idea of coalescent trees
 - These can often be adapted work directly with trees!
 - Deeper analysis of varying types of selection
 - Ghost ancestries, directional migration etc
 - Recombination evolution
 - Etc

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

Thanks!

Simon Myers

Garrett Hellenthal

Pontus Skoglund

Aaron Stern

Lara Cassidy

Robbie Davies

Marie Forest

Sinan Shi

Sile Hu

Anubha Mahajan

Andrew Morris

Mark McCarthy

Daniel Falush

++

