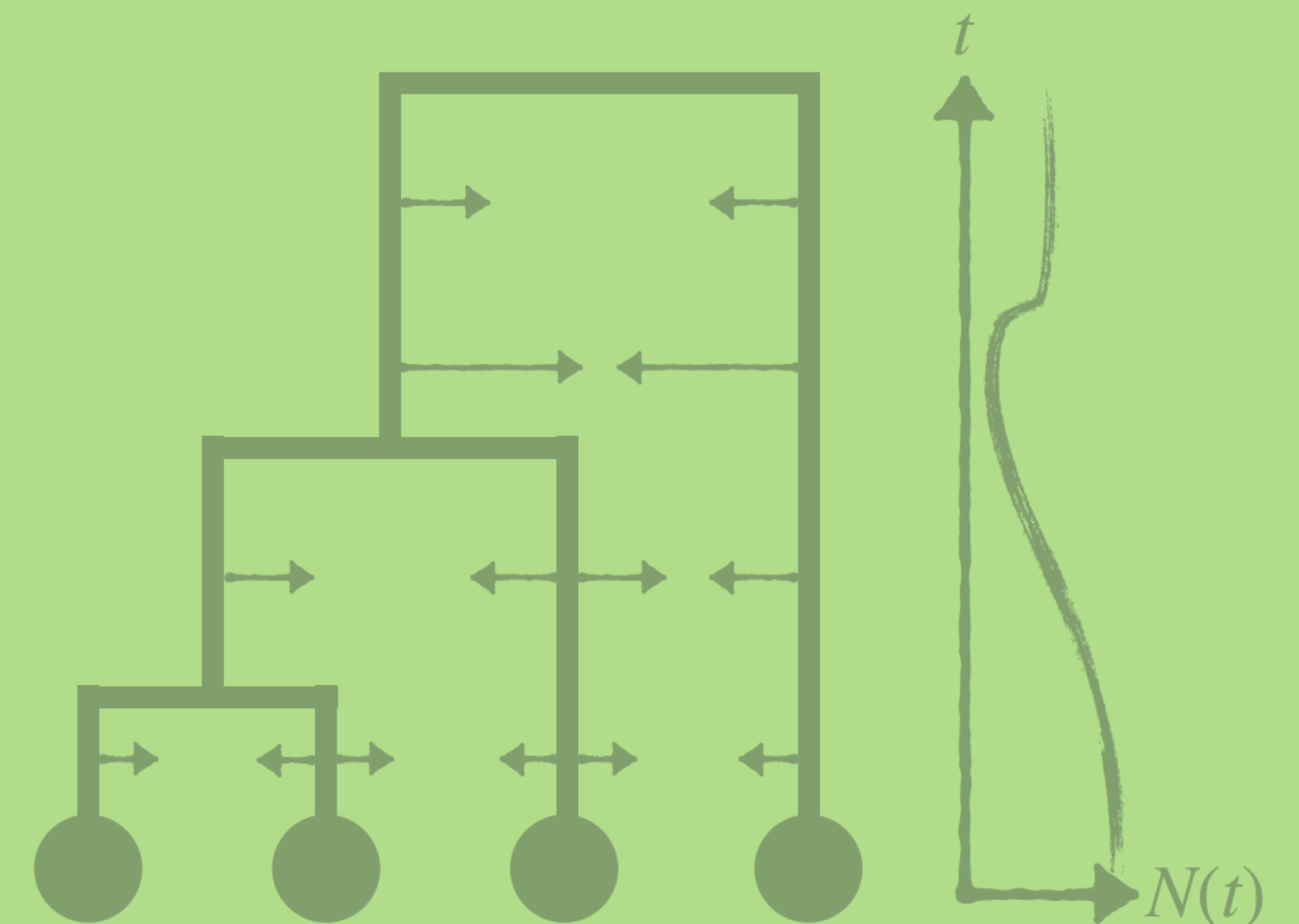# GENOME 541: population genetic inference
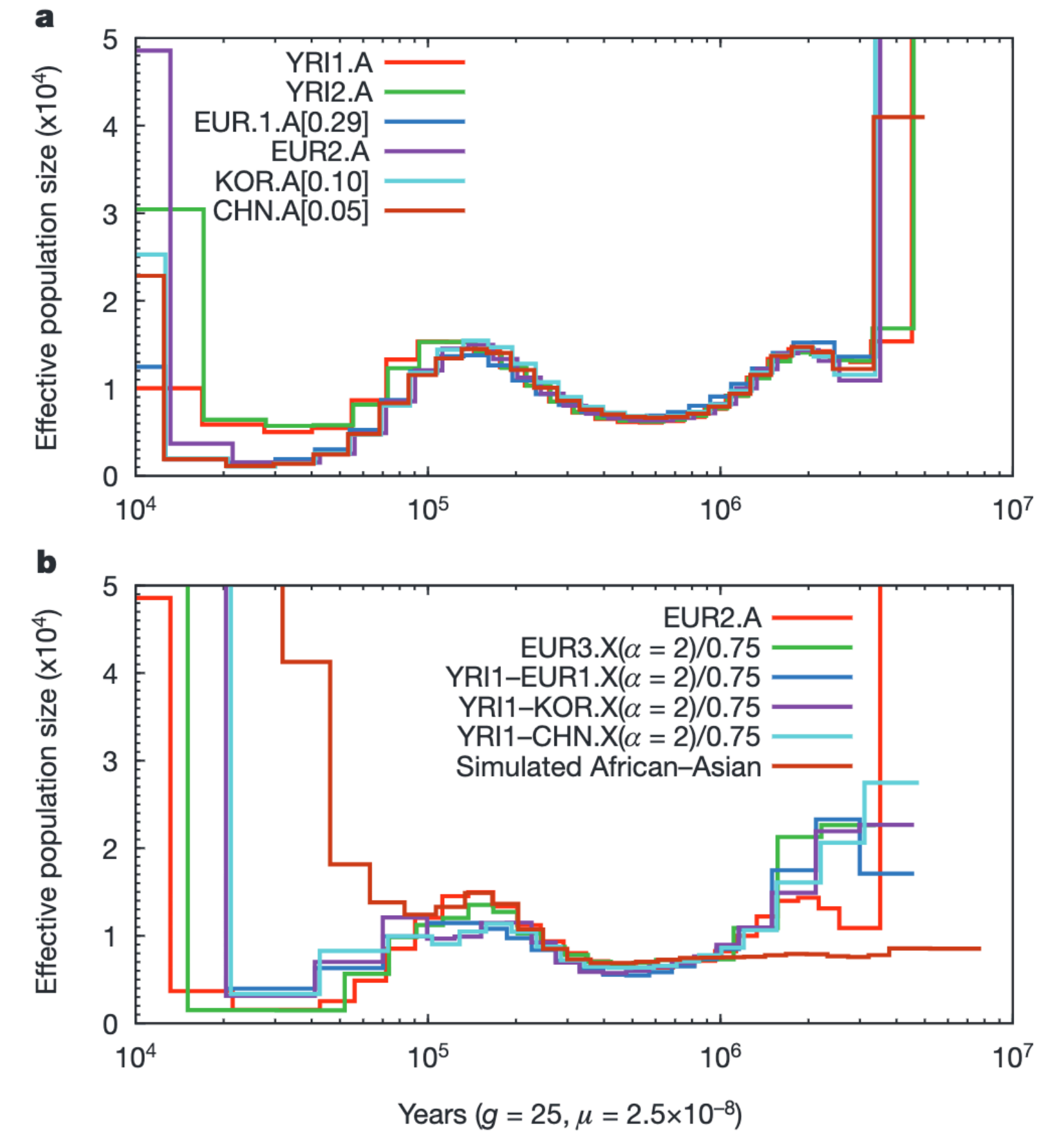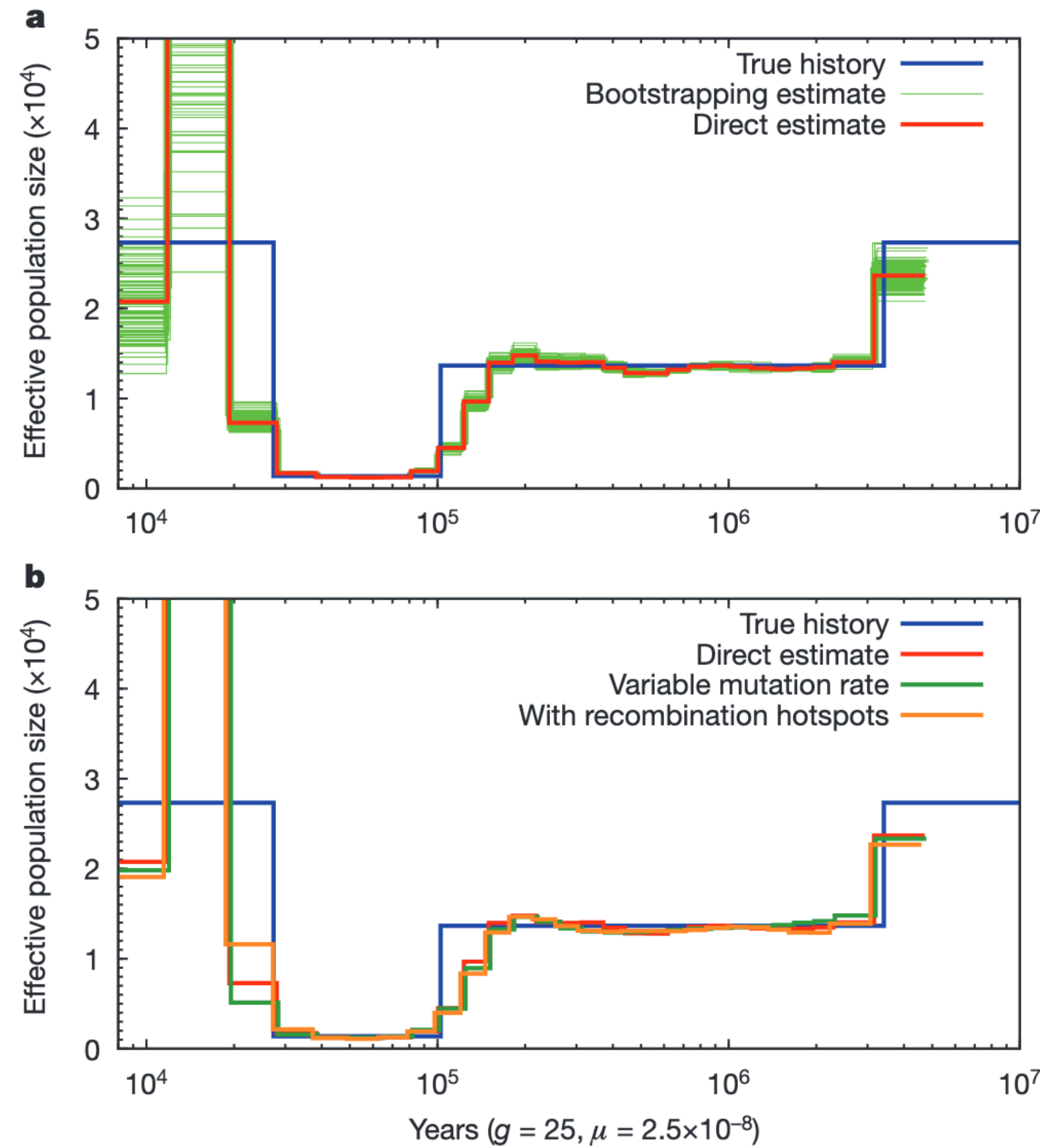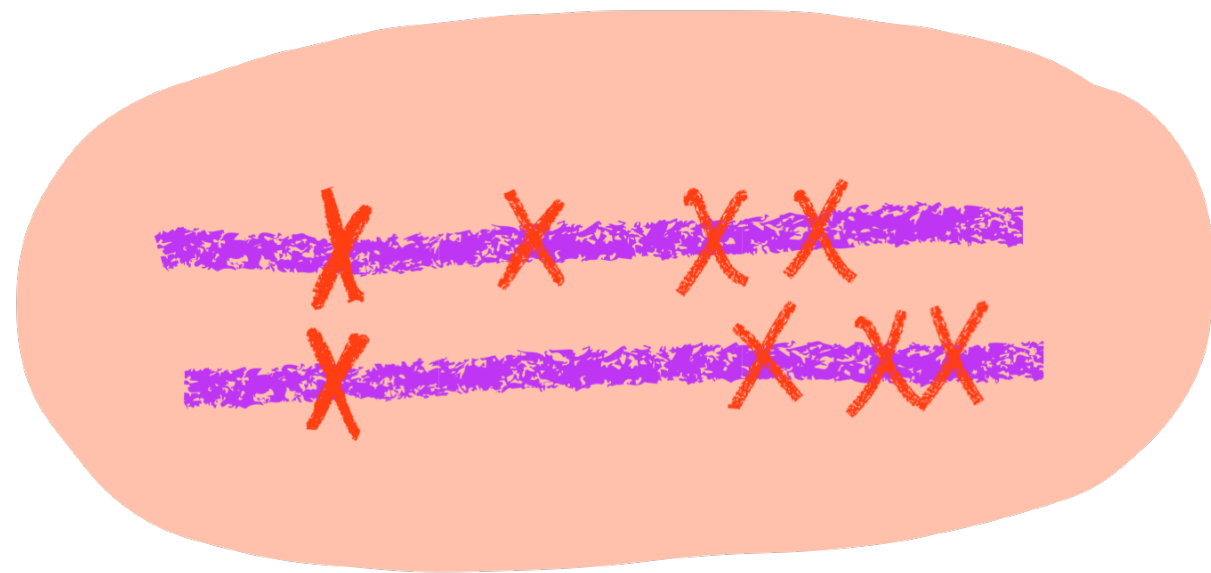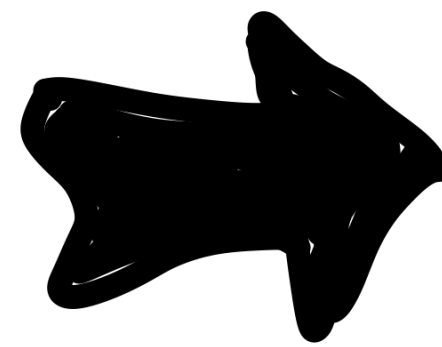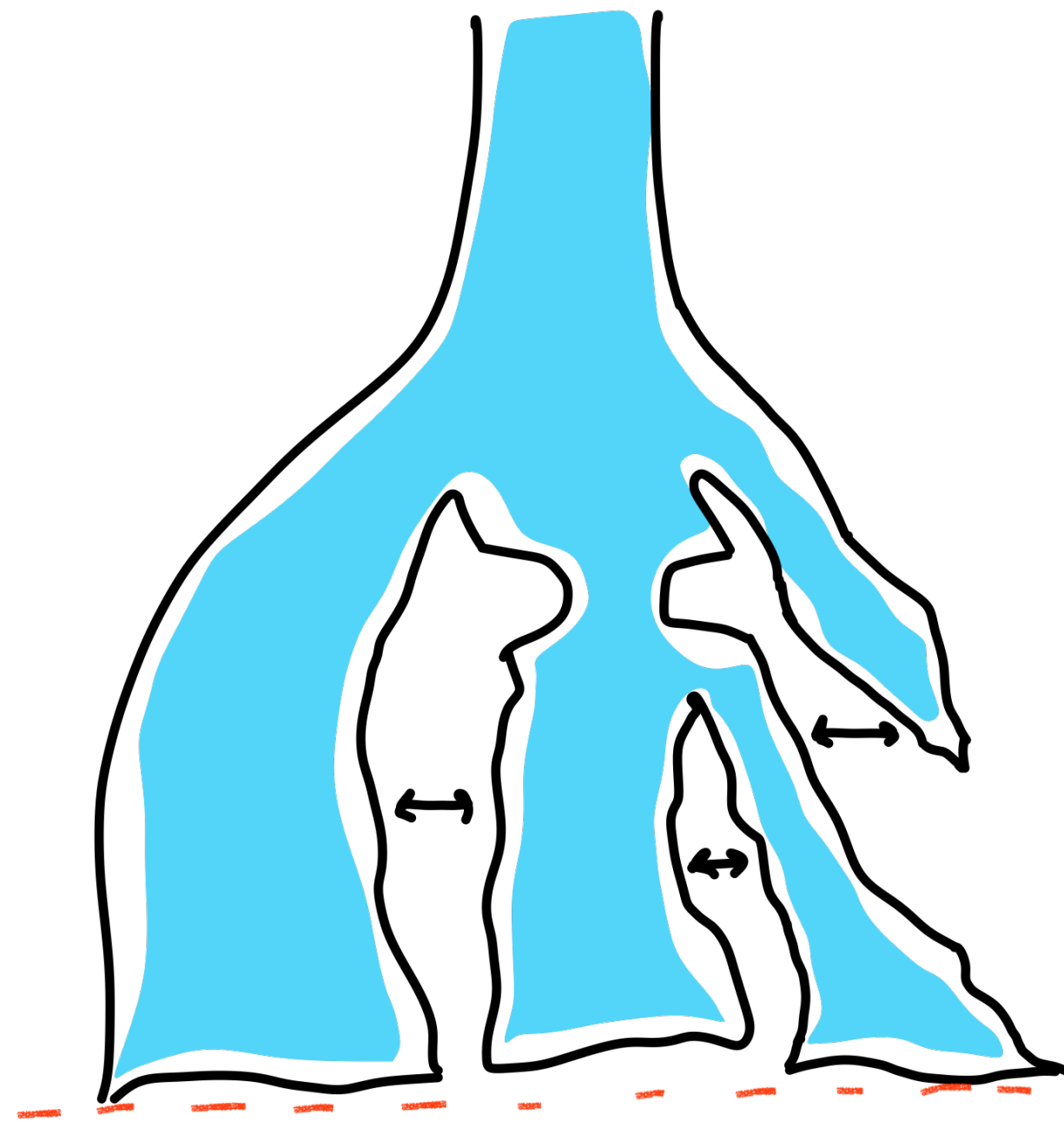
## I. The coalescent (again)

# LETTER

# Inference of human population history from individual whole-genome sequences
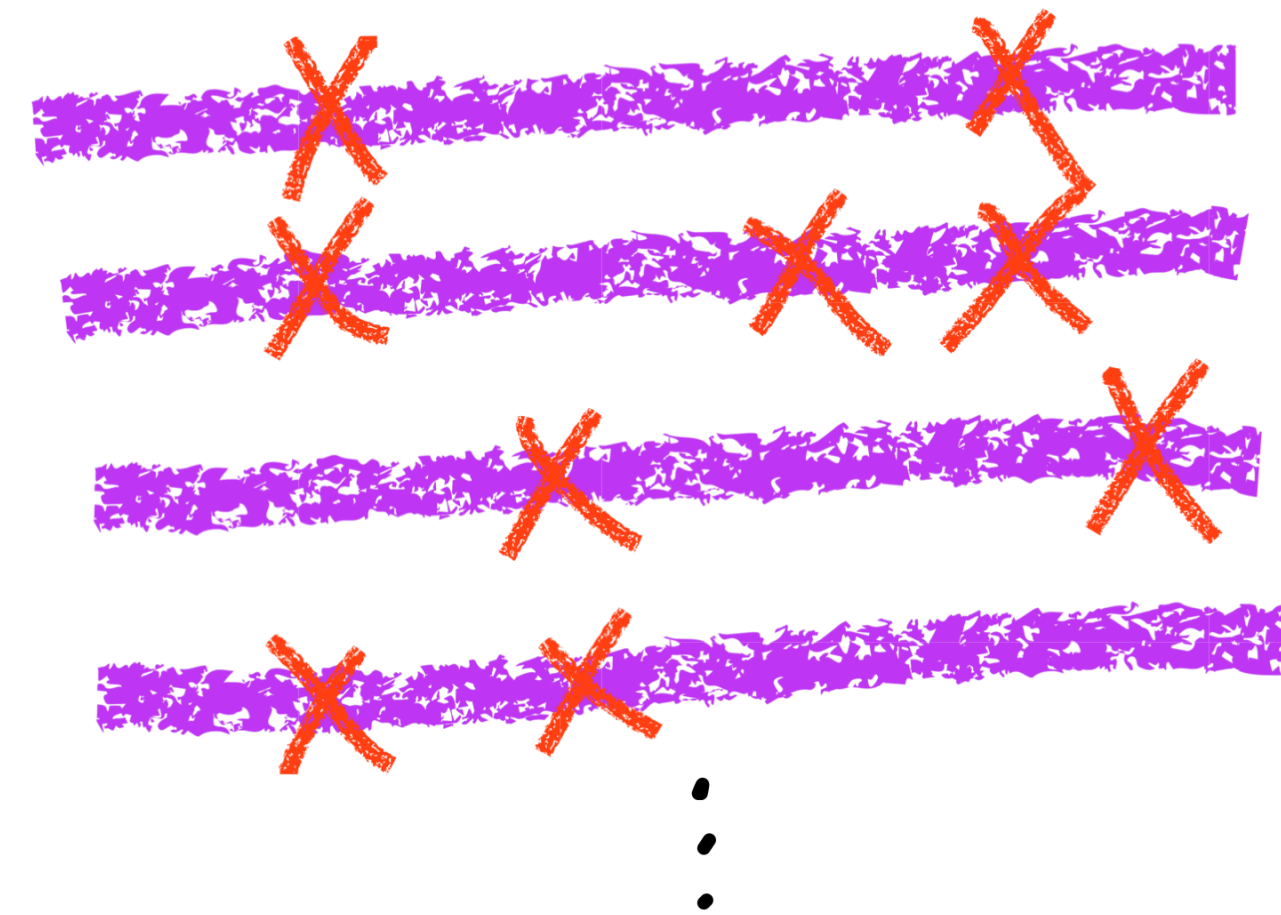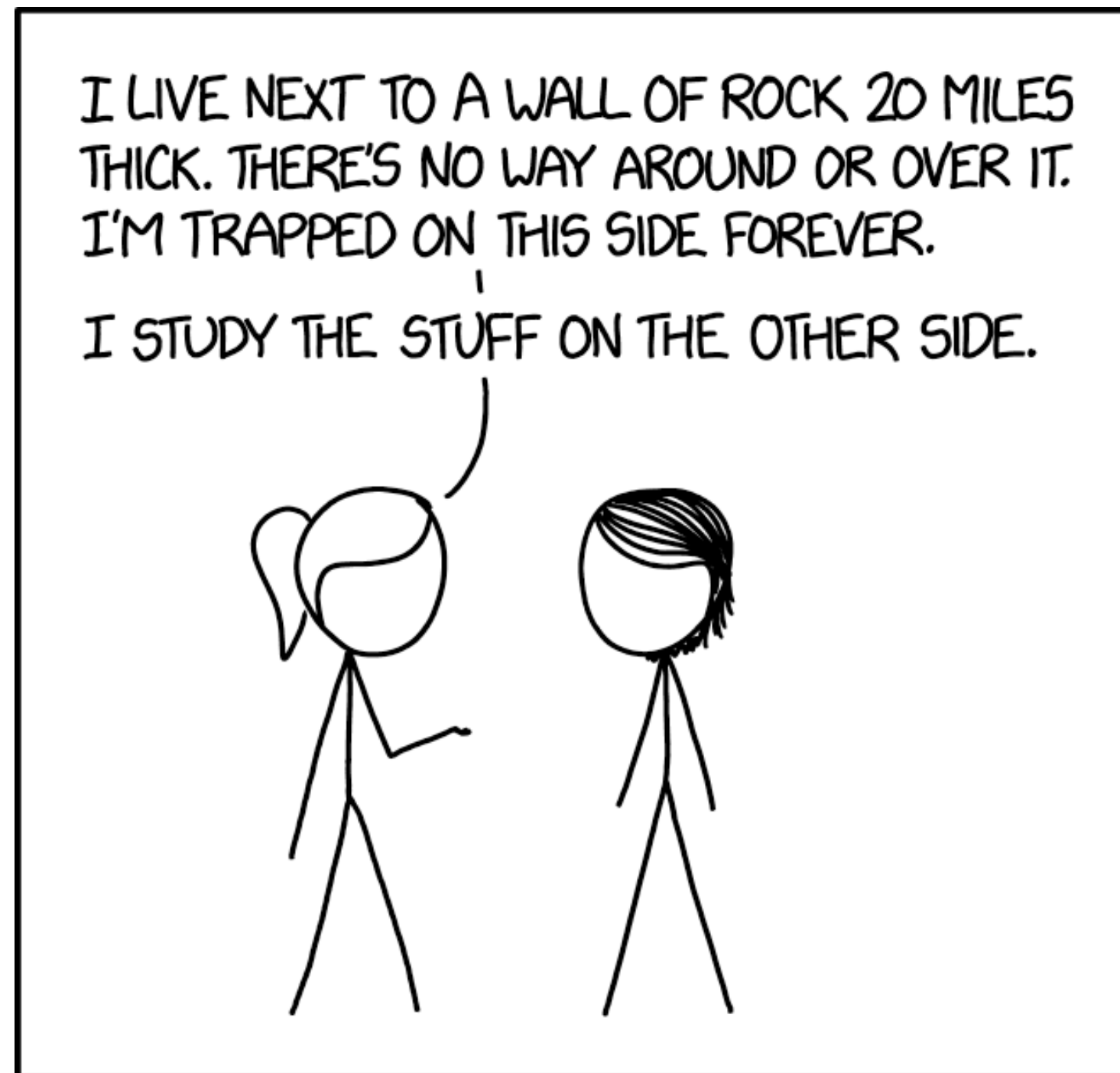
Heng Li[1,2] & Richard Durbin[1]

# The big picture

Evolutionary Process

Genetic Variation

measure stuff here

# Inverse problems



$$g(x) = \int K(x,y)f(y)\mathrm{d}y$$

# Inference frameworks
## Wright-Fisher process at the core

Summary statistics and simple calculations (small data, simple population)

$\pi$ : average pairwise divergence

$S$ : # segregating sites

$T_D$ : deviation from neutrality

Probabalistic models (big data, more complex populations)

Forward time:
- PDEs / diffusion
- selective sweeps

Reverse time:
- coalescent geneology
- coalescent HMM

Machine learning (big data, arbitrarily complex population)

Supervised learning

alignment

NN

# Problem
## The Wright-Fisher process is wrong!

Question: what's missing w/ WF model?

It turns out...

WF
overlapping generations
diploidy 2 sexes
assortative mating
migration
selection?

Universe of models
(WF just one)

big N

The coalescent limit

Universality
PV = nRT
Oh nos, where are the molecules?

# *Cole*-escent theory
## Inspired by GENOME 551 with Cole Trapnell

- 20 students in a class, numbered 1, 2, ..., 20

- Each day of class, professor rolls 20-sided die to choose a student to call on

**Question:** How many classes do I expect it to take me to get called on?

Each class I get called w/ Prob. $\frac{1}{20}$

$$\implies \mathbb{E}[T] = \frac{1}{\frac{1}{20}} = 20 \text{ classes}$$

# *Cole*-escent theory

When there are $i$ students left to call on, the prob of rolling one of these remaining $i$ is $\frac{i}{20}$

$$\mathbb{E}[T_i] = \frac{1}{i/20} = \frac{20}{i}$$

$T_4$

$i = 20, 19, \ldots, 1$

So

$$\mathbb{E}[T_{20 \to 0}] = \mathbb{E}[T_{20}] + \mathbb{E}[T_{19}] + \ldots \mathbb{E}[T_1]$$

$$= 1 + \frac{20}{19} + \ldots + \frac{20}{1} = 1 + 20 \sum_{i=1}^{19} \frac{1}{i} \approx 72$$

# *Cole*-escent theory

Generalize: Class size $N$ (and $N$-sided die)
sample of $n \leq N$ students

$$\mathbb{E}[T_{n \to 0}] = \sum_{i=1}^{n} \mathbb{E}[T_i] = \sum_{i=1}^{n} \frac{1}{i/N} = N \sum_{i=1}^{n} \frac{1}{i}$$

$$P = \frac{i}{N}$$

Each interval is geometrically distributed

$$\mathbb{P}(T_i = t_i) = \frac{i}{N}\left(1 - \frac{i}{N}\right)^{t_i - 1}$$

$$p(t_i) \simeq \frac{i}{N} e^{-\frac{i}{N} t_i}$$

continuous

for large $N$:
$\to$ approx. exponential dist.

$p(t_i)$ ← rate $\frac{i}{N}$

$t_i$

Memoryless:

$$p(t_i \mid t_{i+1}) = p(t_i)$$

# Coalescent theory
## Only slightly fancier

| | Cole-escent | coalescent |
|---|---|---|
| time | classes into the future | generations into the past |
| events | dice rolls, one student | coalescences, pairs of individuals |
| rate of events w/ $i$ individuals | $\dfrac{i}{N}$ ← students uncalled ← class size | $\dfrac{\binom{i}{2}}{2N}$ ← $=\frac{1}{2}i(i-1)$ , #pairs ← population size ↑ if diploid |

# Coalescent theory

## Watch those factors of 2!

Haploid population $N = 6$

$T_2$

two haploid individuals

$n = 2$ sampled haplotypes

Diploid Population $2N = 6$ ($N = 3$)

one diploid individual

$n = 2$ sampled haplotypes

# Coalescent theory
## Updating the previous results



$t$

$T_2$

$T_3$

$T_4$

$$T_{MRCA} = \sum_{i=2}^{n} T_i$$

$n = 4$

"intercoalescence" times

$$\mathbb{E}[T_i] = \frac{2N}{\binom{i}{2}}$$

Exp. dist.

$$p(t_i) = \frac{\binom{i}{2}}{2N} e^{-\frac{\binom{i}{2}}{2N} t_i}$$

- Each **pair** is a process w/ rate $\frac{1}{2N}$
- The pairs race to coalesce
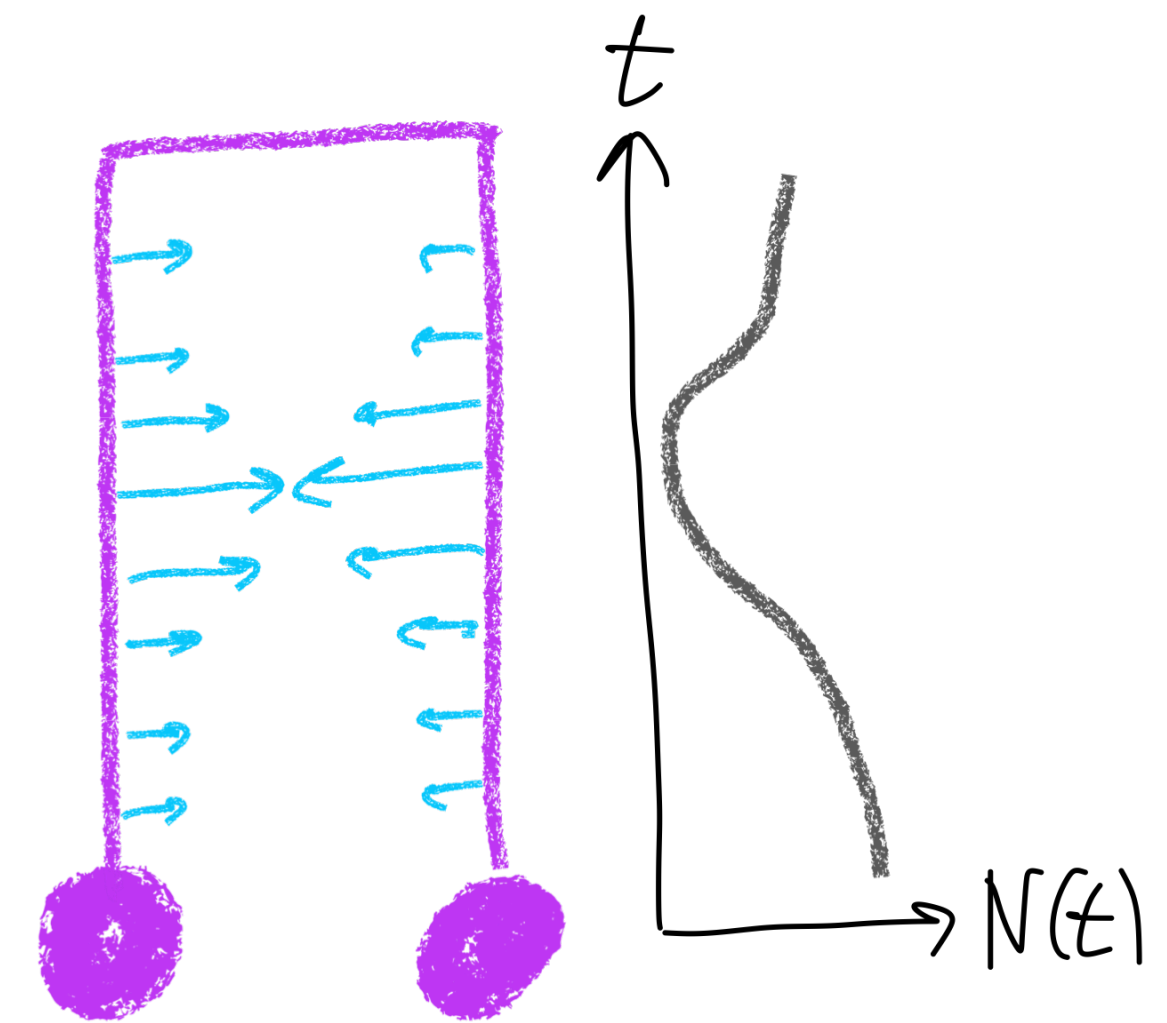
12

# Coalescent theory

## Population size determines coalescence rate

> What if population size varies over time ? $N(t)$

[ Like students adding/dropping mid-quarter in Cole-escent theory ]

$N(t)$ distorts time scale from the standard coalescent

- time compressed when $N(t)$ is small
- time stretched when $N(t)$ is large

$t$

$N(t)$

**The details :**

$$\mathbb{P}\left(T_i = t_i\right) = \frac{\binom{i}{2}}{2 N_{t_i}} \prod_{j=1}^{t_i - 1}\left(1 - \frac{\binom{i}{2}}{2 N_j}\right)$$

$$\xrightarrow{\text{big } N} \; p\left(t_i\right) = \frac{\binom{i}{2}}{2 N(t)} e^{-\binom{i}{2}\int_0^t \frac{ds}{2 N(s)}}$$
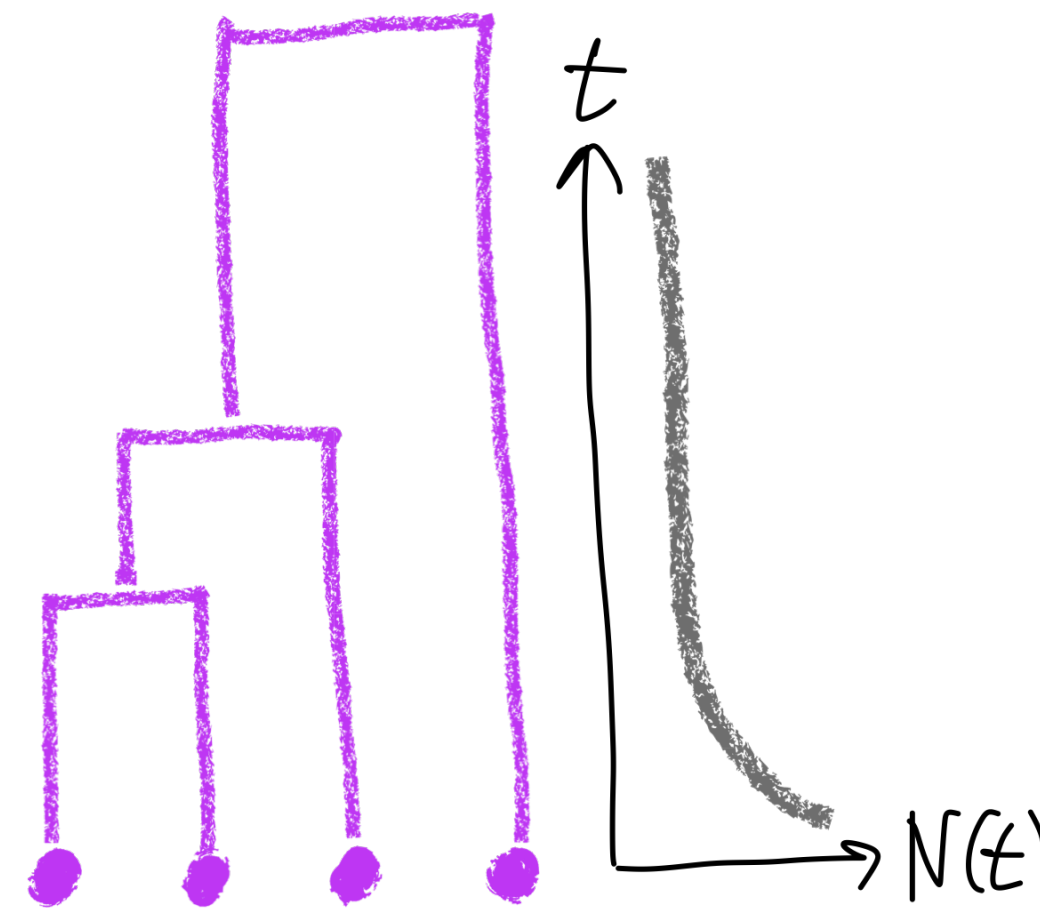
inhomogeneous Poisson process

# Coalescent theory
## Population size determines coalescence rate

Constant
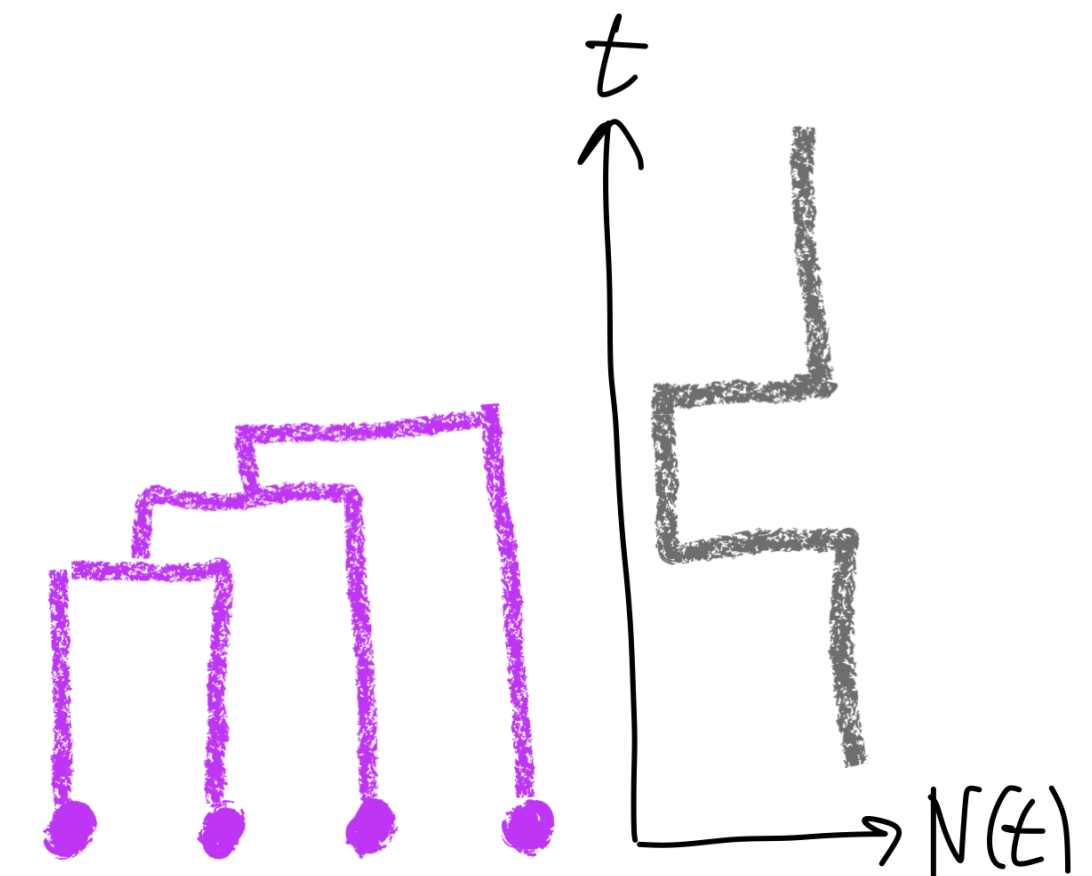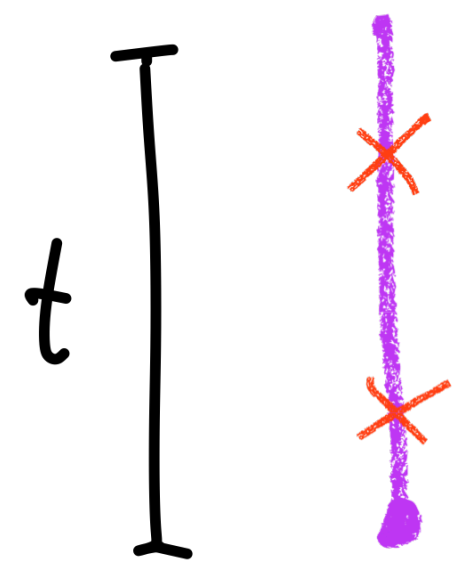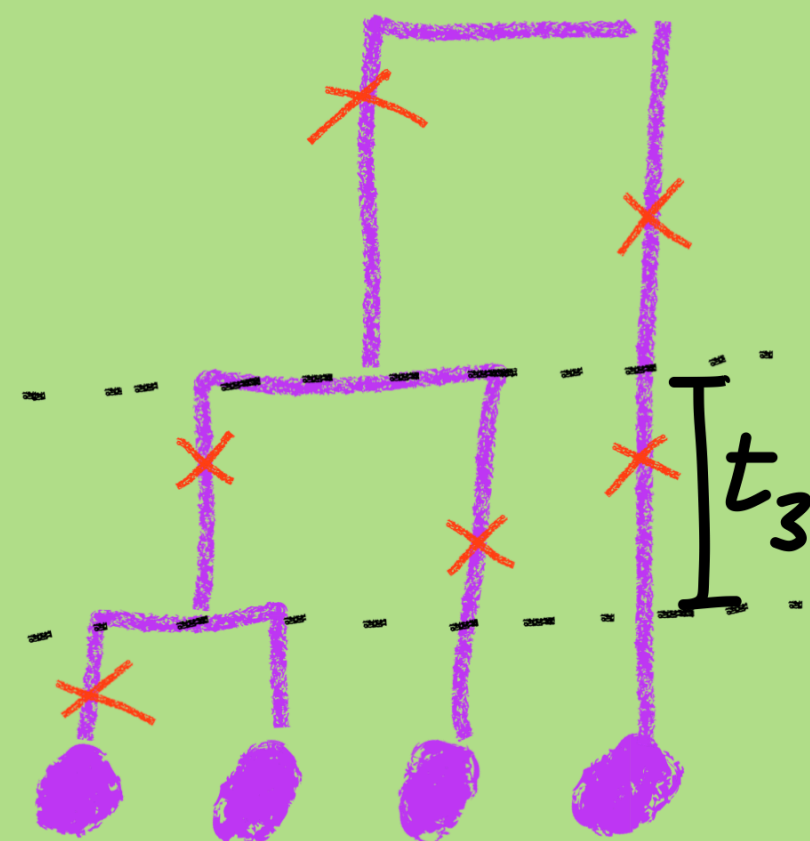
Exponential Growth

Bottleneck

# Coalescent theory

## Mutations

# mutations on branch of length $t$ is Poisson rv w/ mean $\mu t$

$$\mathbb{P}(k \mid t) = \frac{(\mu t)^k e^{-\mu t}}{k!}$$
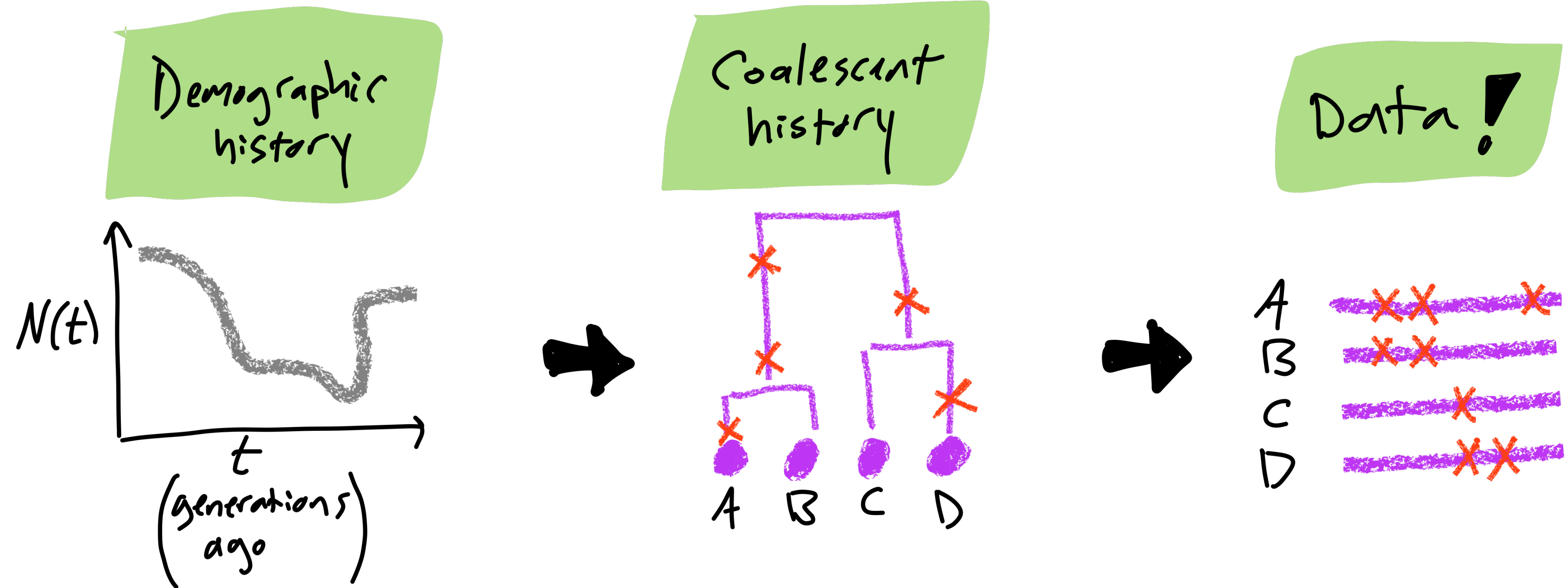
mutations per genome per generation

# mutations in intercoalescent interval $i$, of length $t_i$, is Poisson rv w/mean

$$i \mu t_i$$
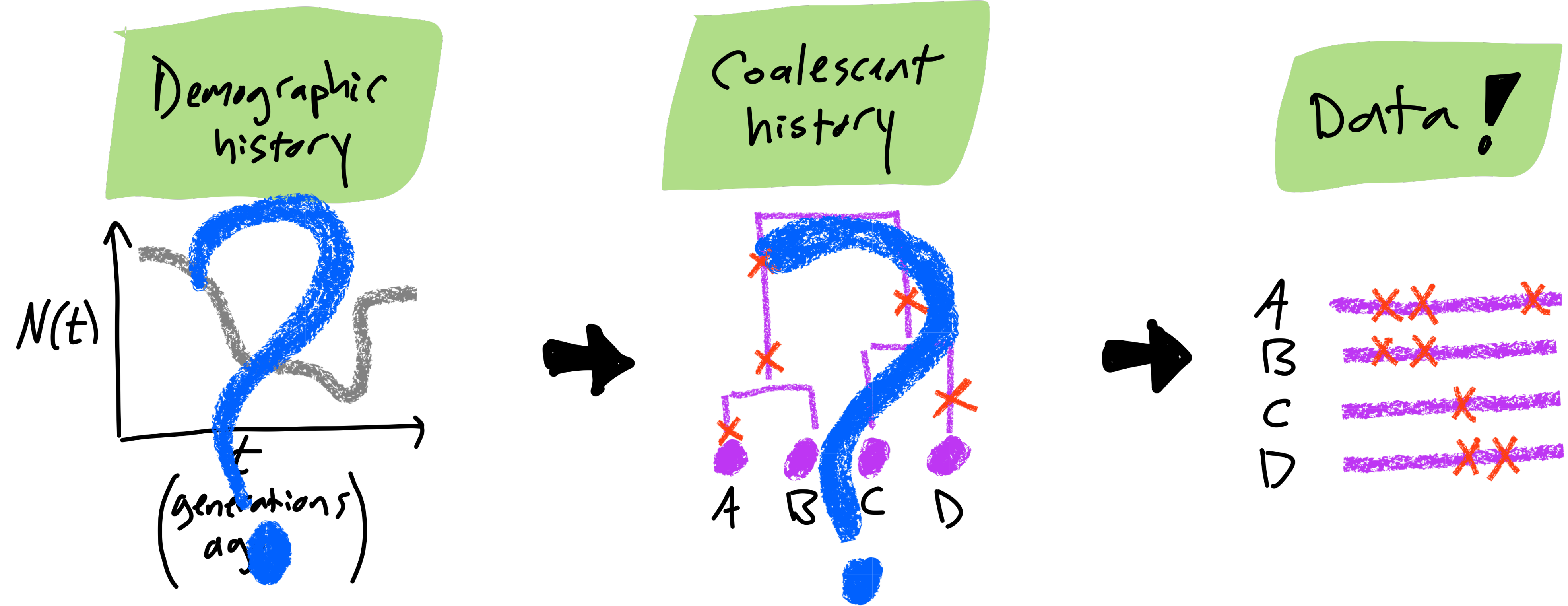
# lines in interval $i$

duration of interval $i$

$t_3$

# Coalescent theory
## Genetic diversity

# Coalescent theory
## Genetic diversity



Demographic history

$N(t)$

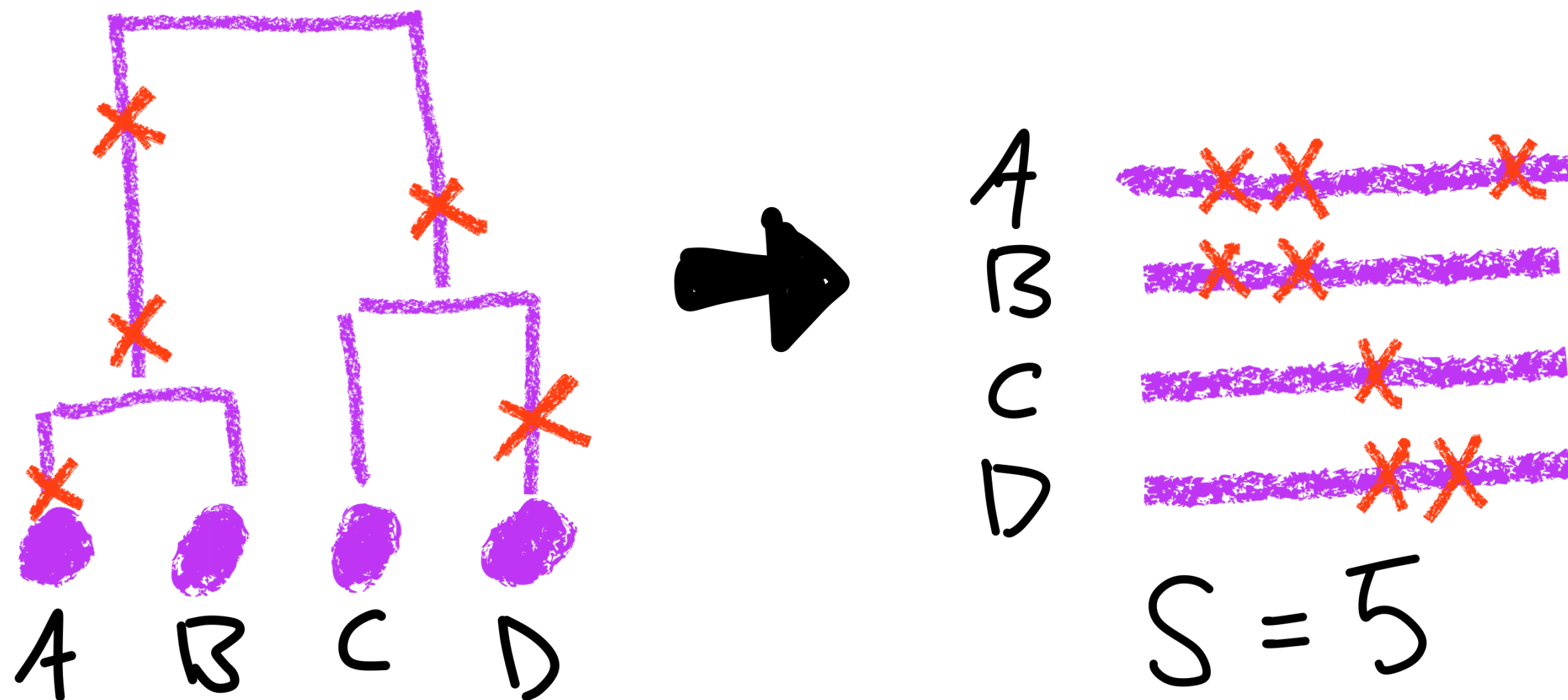$t$ (generations ago)

Coalescent history

A  B  C  D

Data!

A
B
C
D

# Coalescent theory
## Genetic diversity

# segregating sites, $S$, equals # mutations in the samples history (infinite sites approximation)



$S = 5$

Constant $N$ case:

$\mathbb{E}[S] = \mu \, \mathbb{E}[T_{total}]$ ← total branch length

$= \mu \sum_{i=2}^{n} i \, \mathbb{E}[T_i]$

$= \mu \sum_{i=2}^{n} i \, \frac{2N}{\binom{i}{2}}$

$= 4 \mu N \sum_{i=1}^{n-1} \frac{1}{i}$

# Recap
## Good to know for homework

Intercoalescent times are indep. exponential rvs

Intercoalescent interval $i$  diploid pop size  #sampled haplotypes

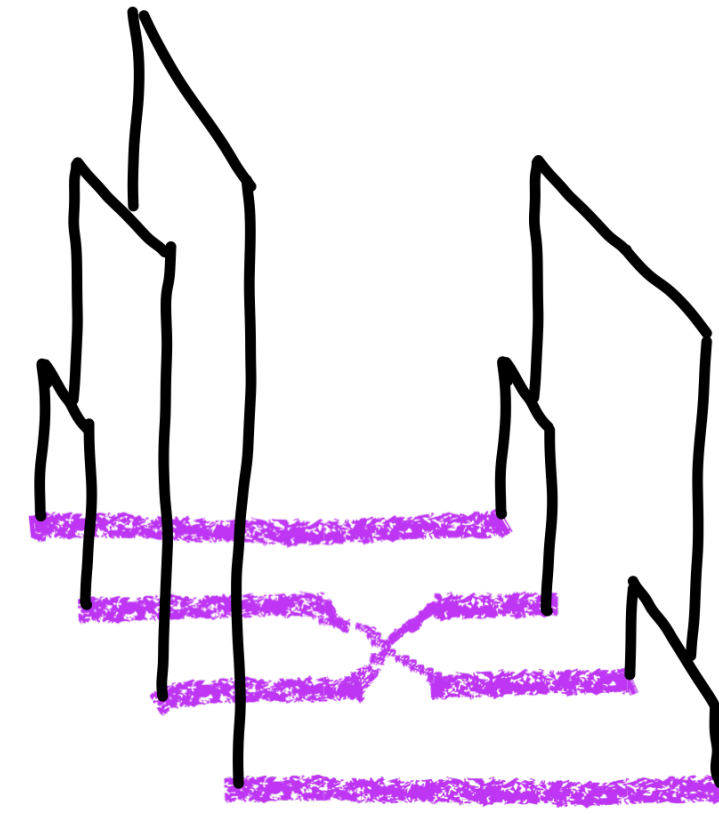$$T_i \sim \exp\left(\frac{2N}{\binom{i}{2}}\right), \text{ for } i = n, n-1, \dots, 2$$

#mutations on branches are indep. Poisson rvs

$$k \sim \text{Pois}\left(\mu t\right)$$
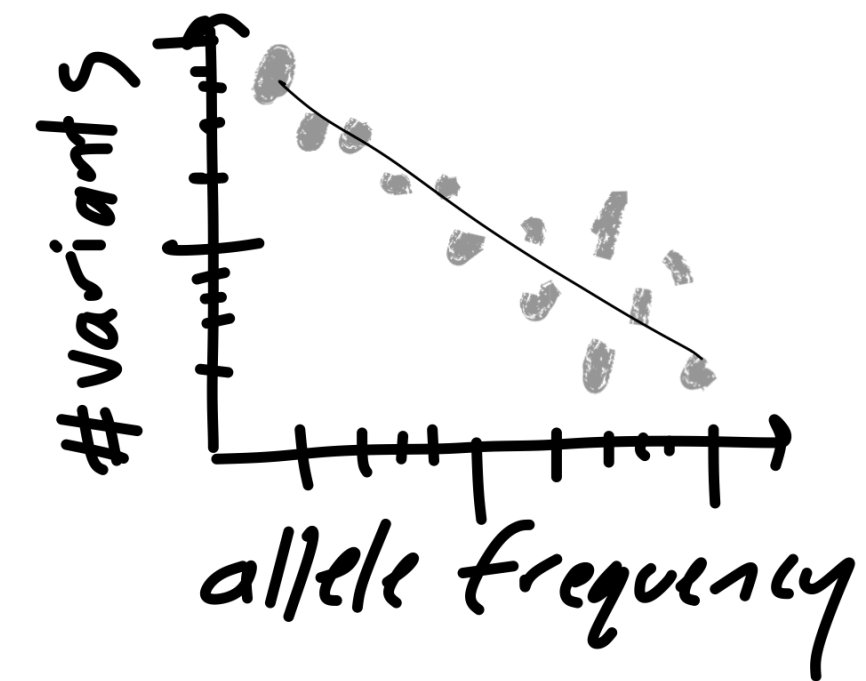
# mutations   mutation rate   branch length

# Next time

- The coalescent with recombination



- Sample Frequency spectrum (SFS)



#variants vs allele frequency

- Coalescent hidden Markov model



$T_{MRCA}$