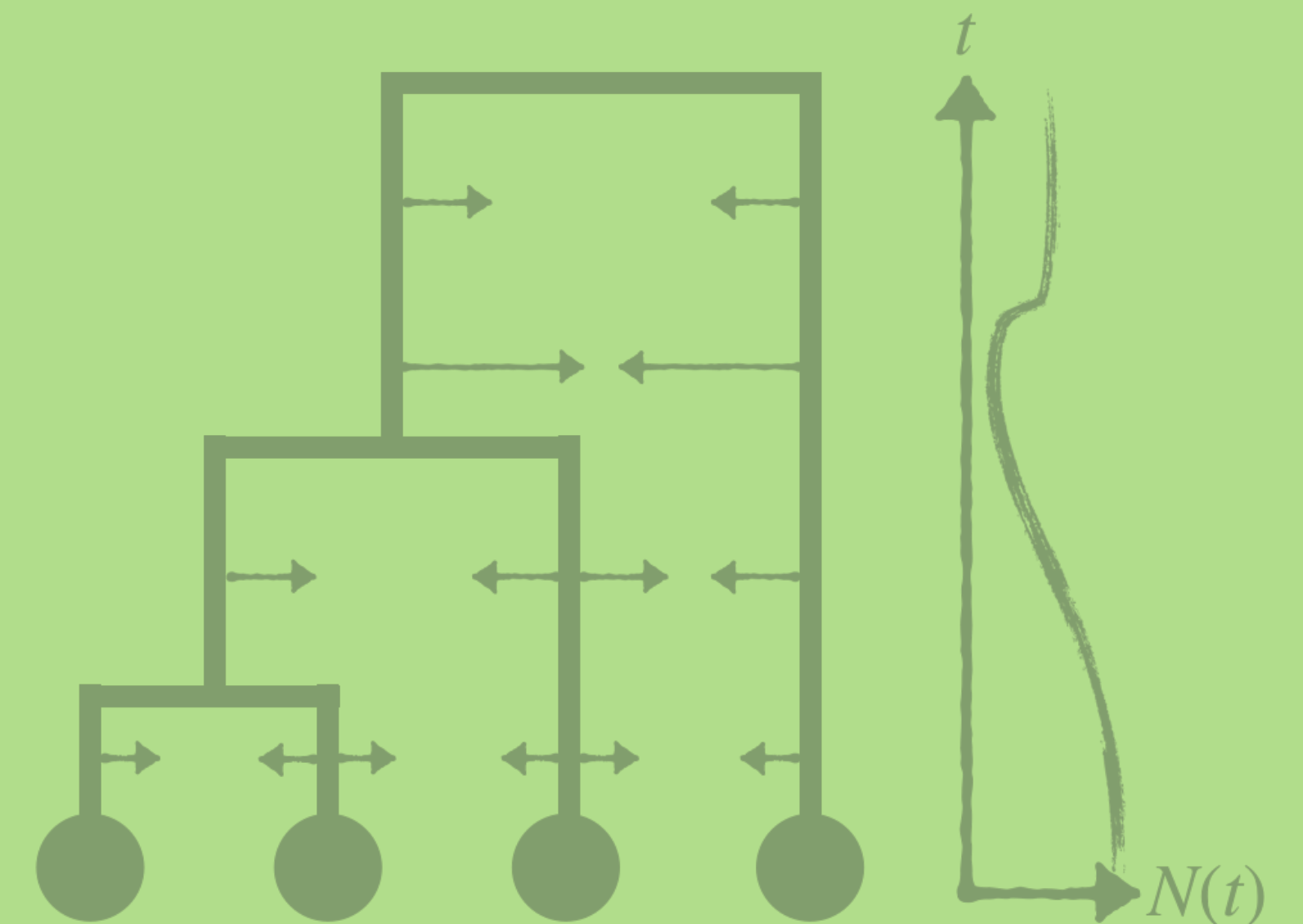
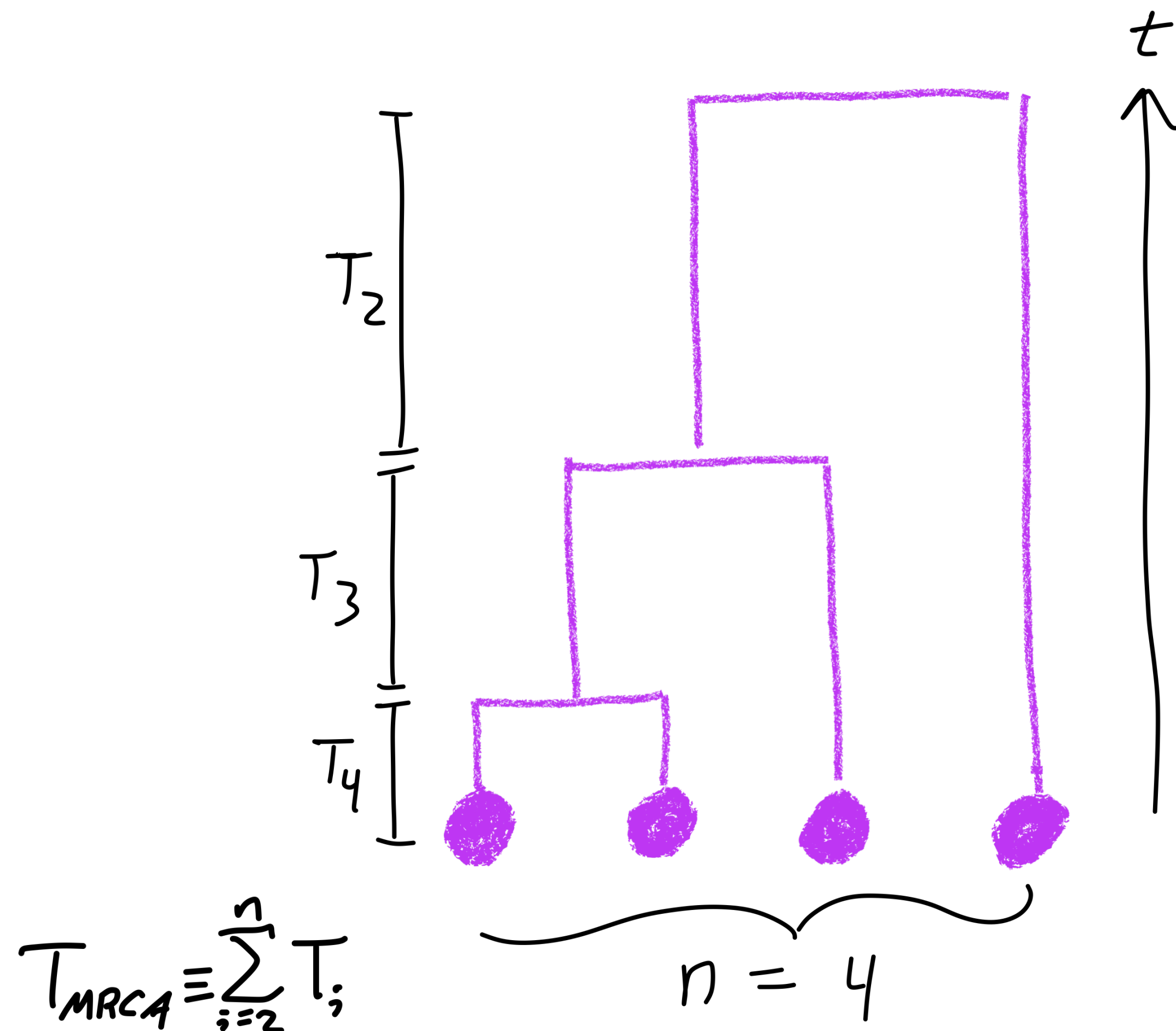


GENOME 541: population genetic inference

II. The coalescent with recombination



Previously on...



"intercoalescence"
times

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

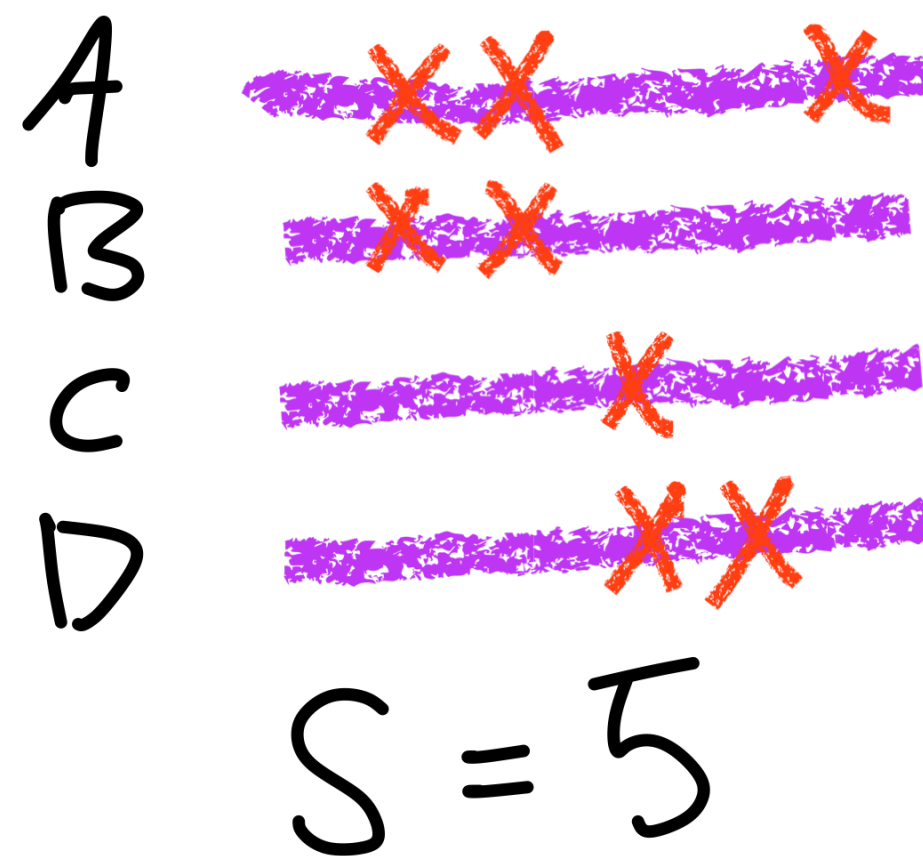
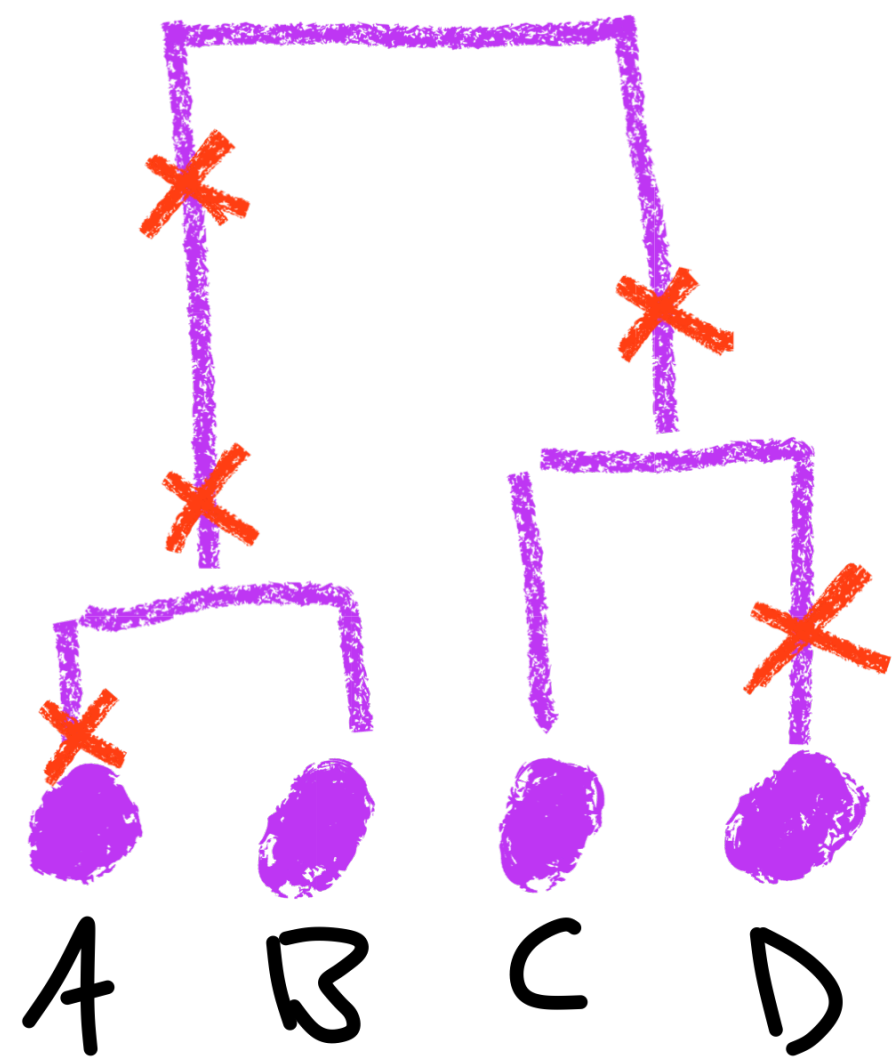
Exp. dist.

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

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

Previously on...

segregating sites, S , equals # mutations in the sample's history (infinite sites approximation)



Constant N case:

$$E[S] = \mu E[T_{\text{total}}]$$

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

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

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

total branch length

Genetic diversity stats

segregating sites, S , equals # mutations in the sample's history (infinite sites approximation):

$$E[S] = 4Nm \sum_{i=1}^{n-1} \frac{1}{i}$$

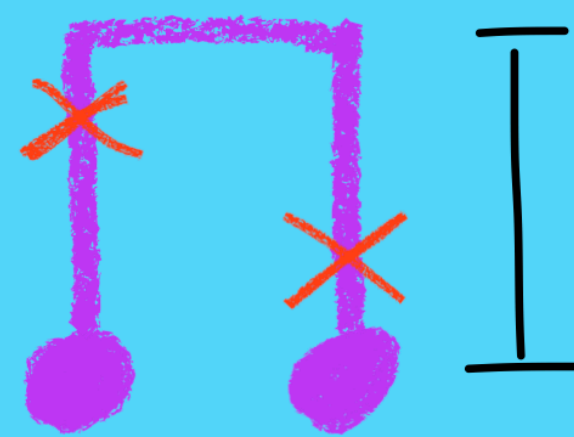
Pairwise divergence, π , # mutations in the history of two sampled haplotypes:

$$E[\pi] = 4Nm$$

Tajima's D: null hypothesis, standard neutral coalescent (constant N)

$$\frac{E[S]}{\sum_{i=1}^{n-1} \frac{1}{i}} - E[\pi] = 0$$

Derivation:

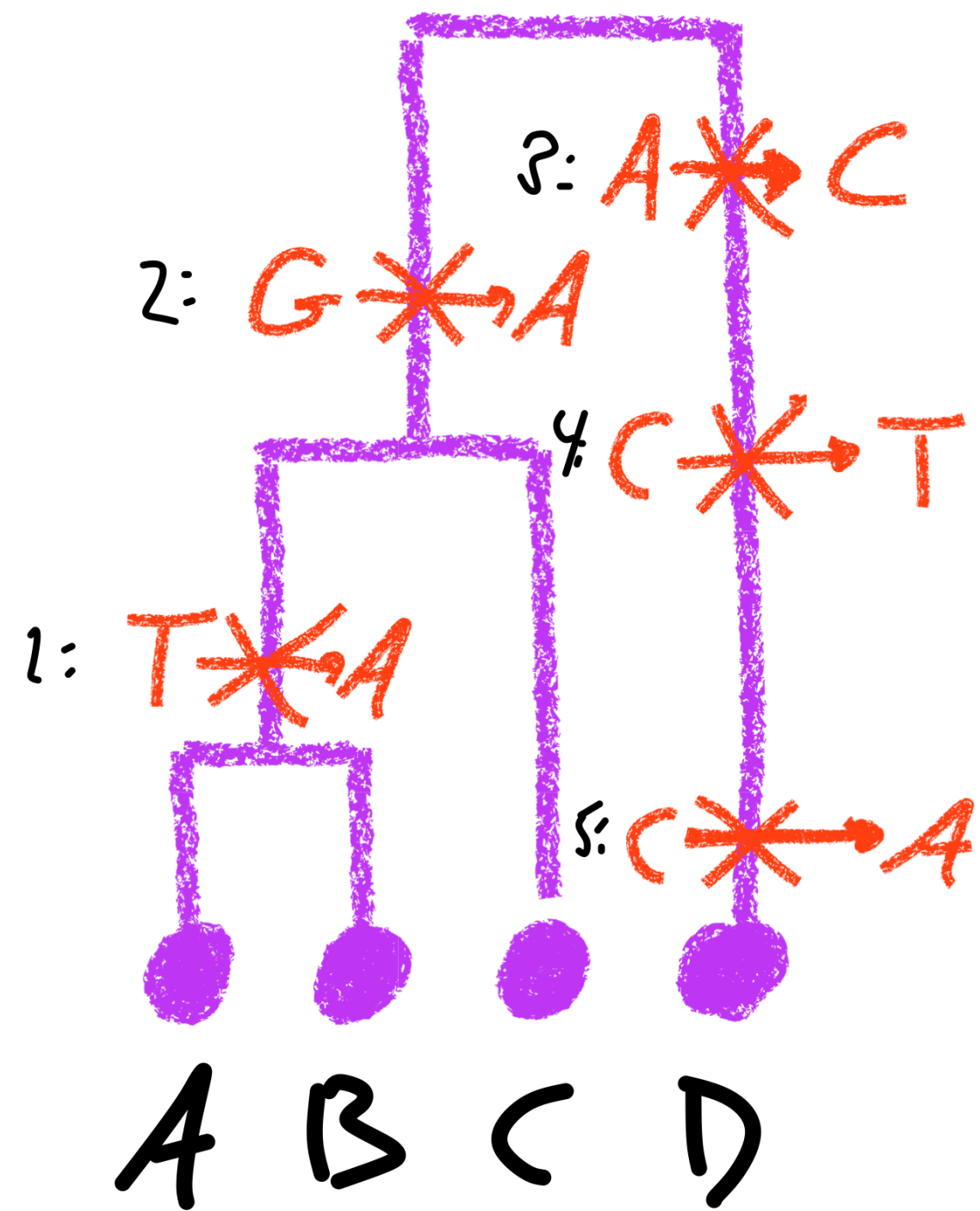


$$E[T_2] = \frac{1}{\frac{1}{2N}} = 2N$$

$$\rightarrow E[\pi] = 4Nm$$

Question: why 4?

Example



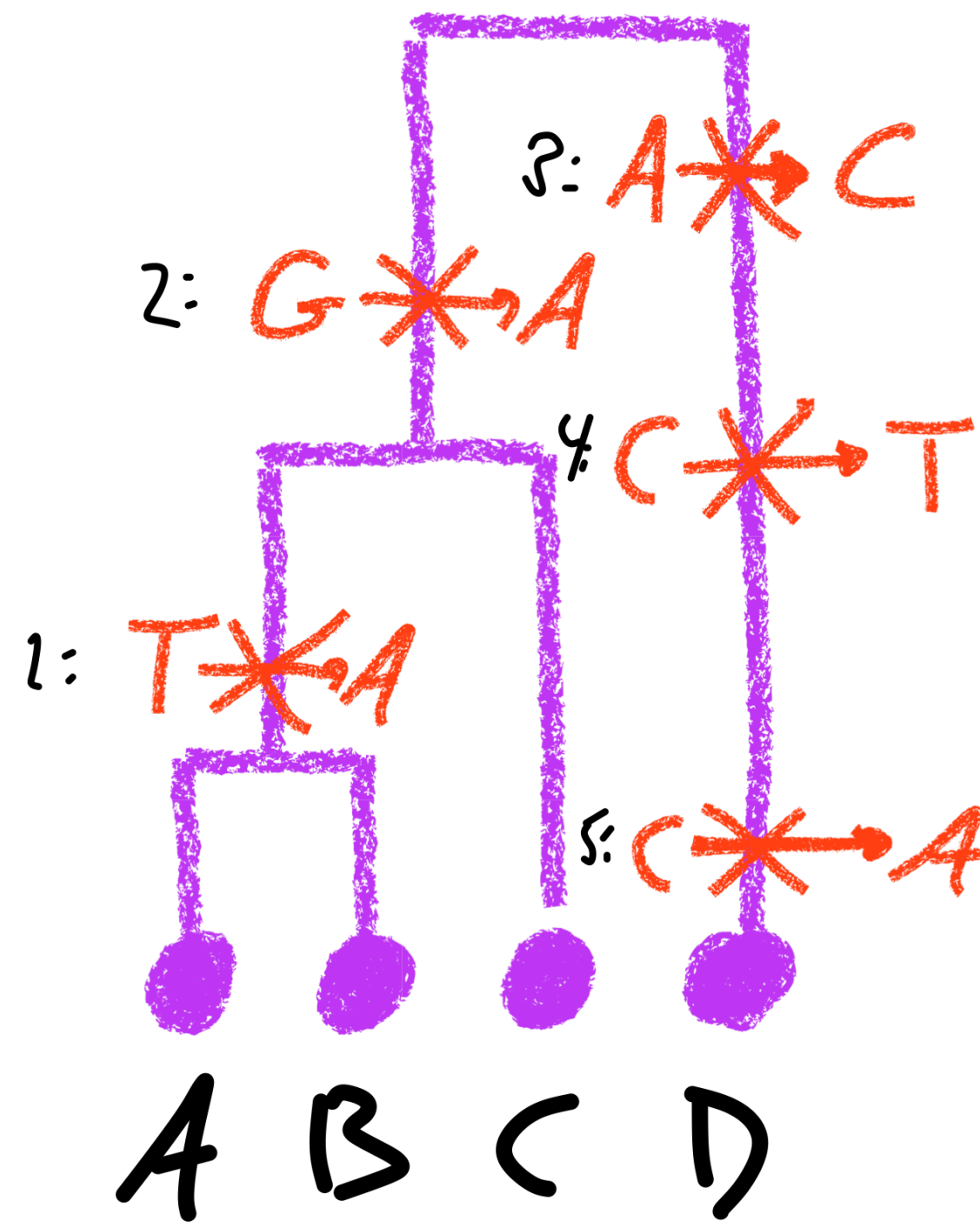
		1	2	3	4	5			
A	•••	A	•••	A	•••	C	•••	C	•••
B	•••	A	•••	A	•••	C	•••	C	•••
C	•••	T	•••	A	•••	C	•••	C	•••
D	•••	T	•••	G	•••	C	•••	T	•••

$$S = 5$$

$$\pi = \frac{1}{\binom{4}{2}} (0 + 1 + 5 + 2 + 5 + 4) = \frac{8}{3}$$

$$D = \frac{5}{1 + \frac{1}{2} + \frac{2}{3}} - \frac{8}{3} = \dots = \frac{2}{33}$$

Example: sample frequency spectrum (SFS)



	1	2	3	4	5
A	A	A	A	C	C
B	A	A	A	C	C
C	T	A	A	C	C
D	T	G	C	T	A



SFS: histogram of mutant allele frequencies

$\vec{z} = [z_2, z_3, \dots, z_{n-2}]$, $z_i = \# \text{mutations with frequency } i \text{ in sample}$

Theory: sample frequency spectrum (SFS)

SFS: histogram of mutant allele frequencies

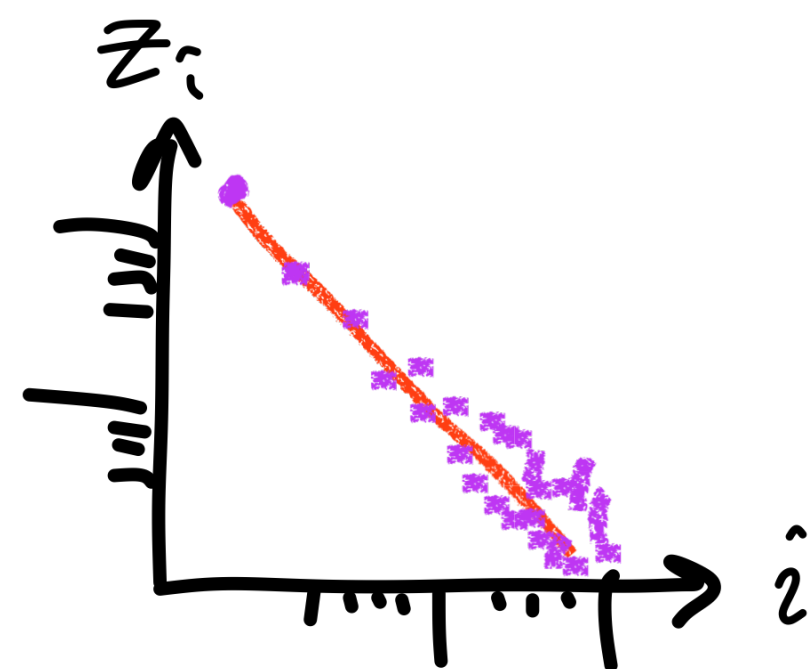
$\vec{z} = [z_2, z_2, \dots, z_{n-2}]$, $z_i = \# \text{mutations with frequency } i \text{ in sample}$

With a bit of work ...

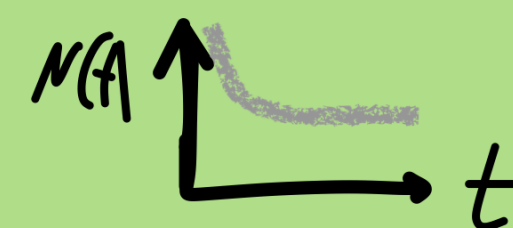
$$E[z_i] = \frac{4N\mu}{i}$$

(for constant N) $i=1, 2, \dots, n-1$

$$\begin{aligned} \log z_i &= \log \left(\frac{4N\mu}{i} \right) \\ &= \log 4N\mu - \log i \end{aligned}$$

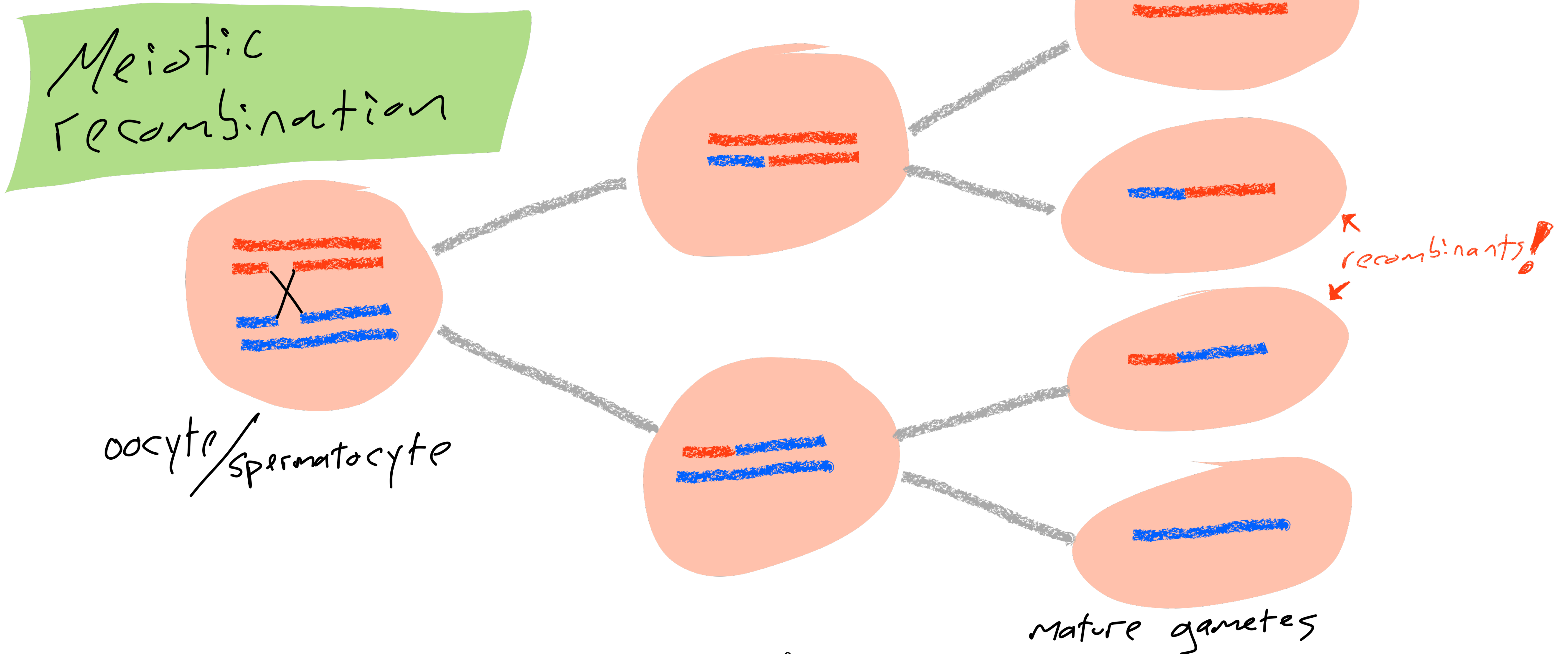


Question: what happens to SFS if $N(t)$ is changing?



Non tree-like ancestry

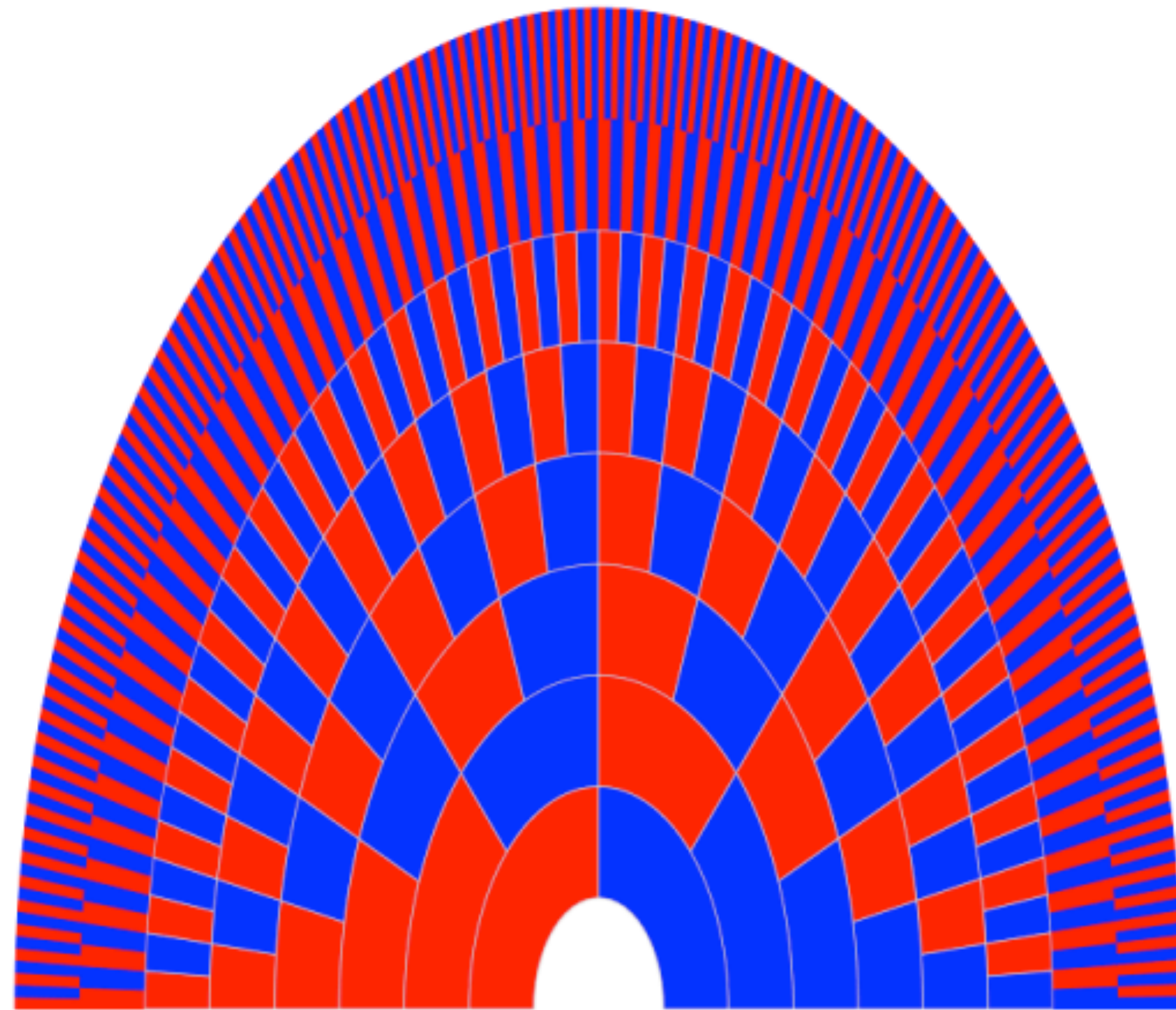
Recombining genomes as mosaics



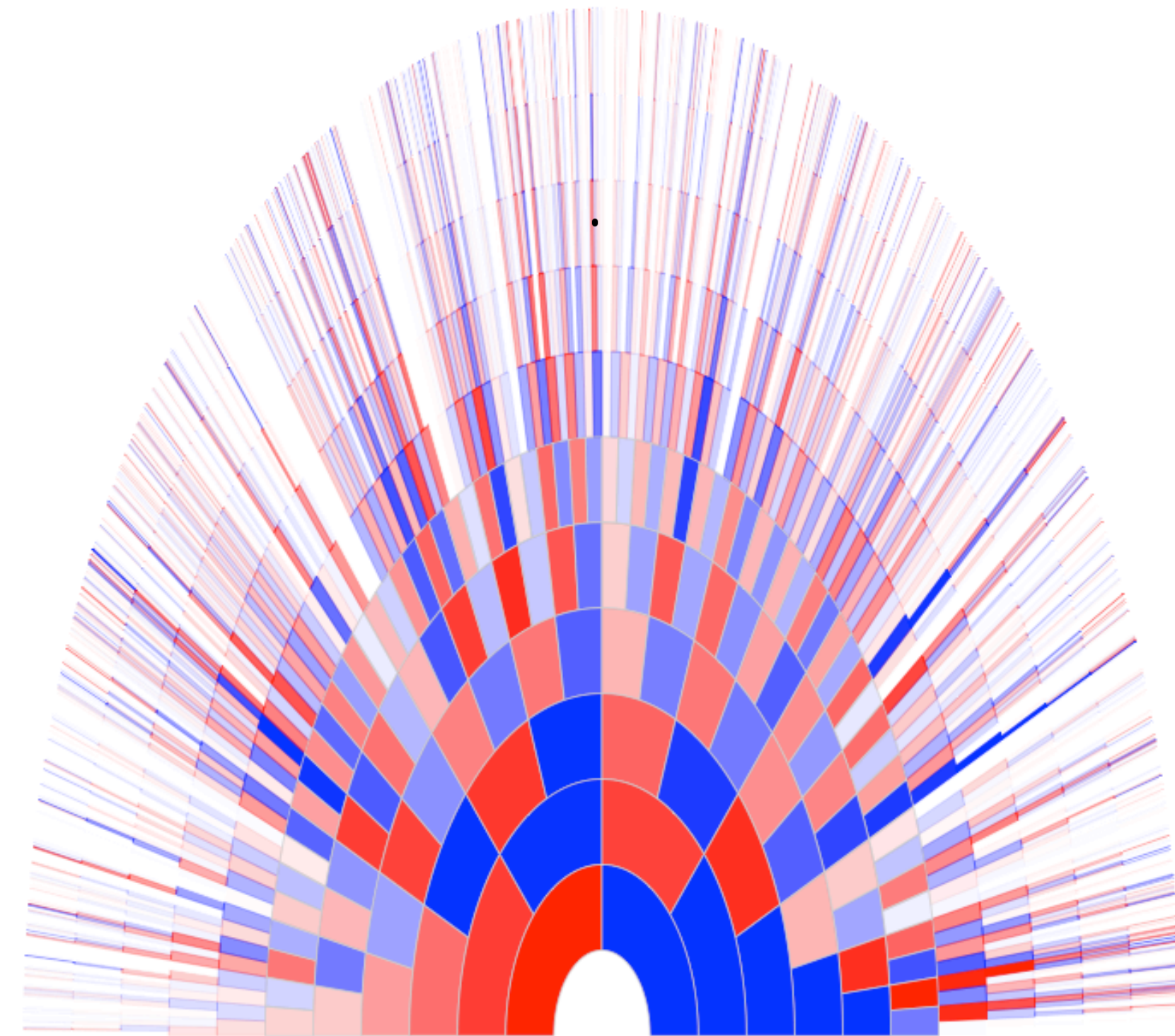
Non tree-like ancestry

Recombining genomes as mosaics

genealogical ancestry



genetic ancestry

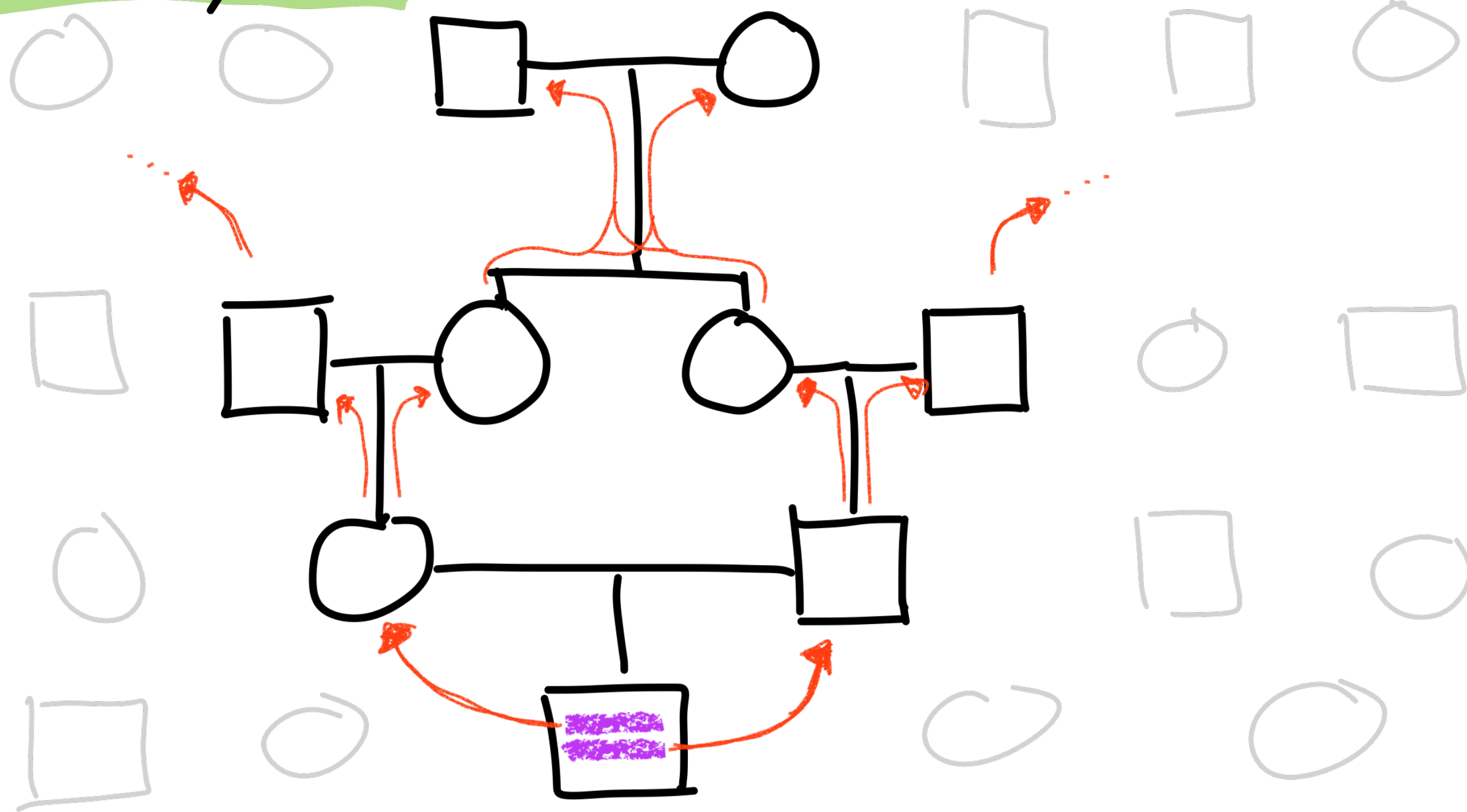


<https://gcbias.org/2013/11/11/how-does-your-number-of-genetic-ancestors-grow-back-over-time/>

Non tree-like ancestry

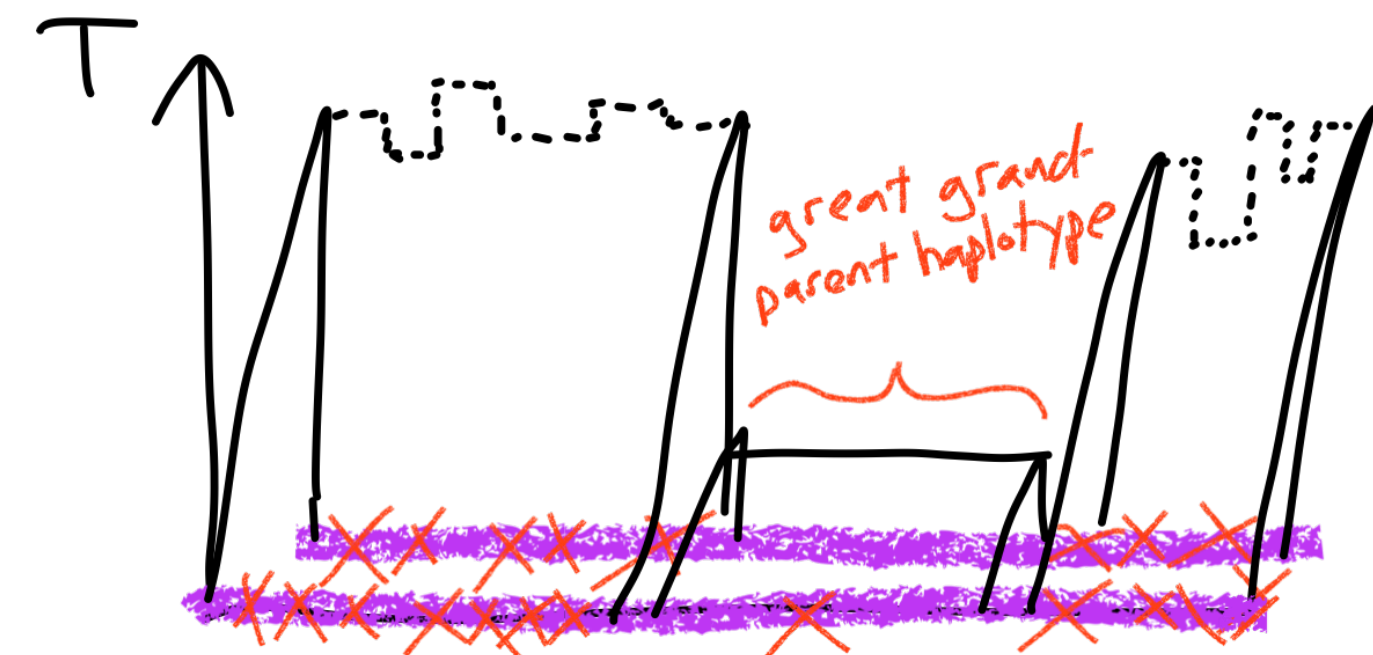
Recombining genomes as mosaics

1st - cousin mating



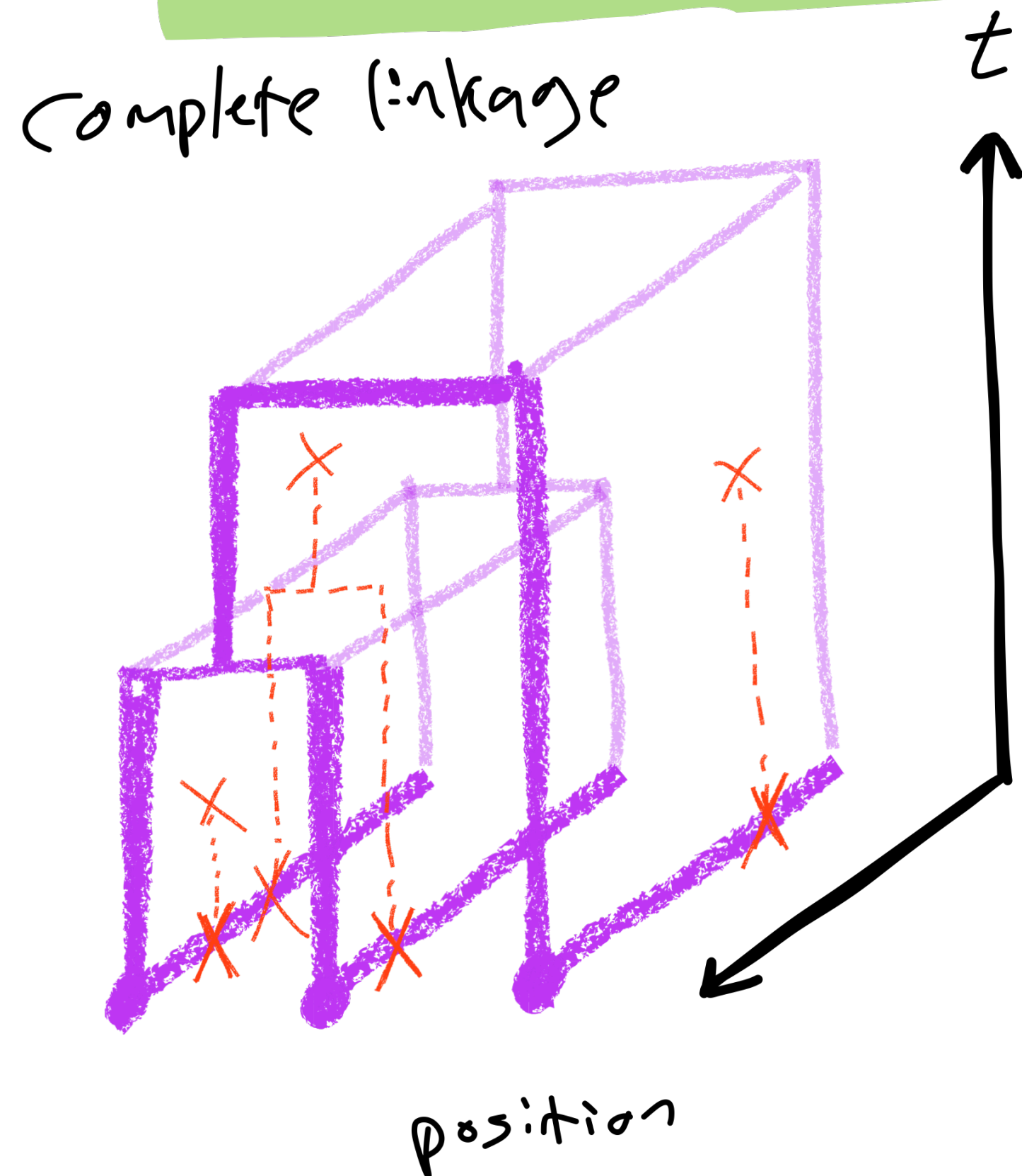
Question:

- probability $\boxed{?}$ to coalesce in great grandparent?
- Otherwise much deeper coalescence



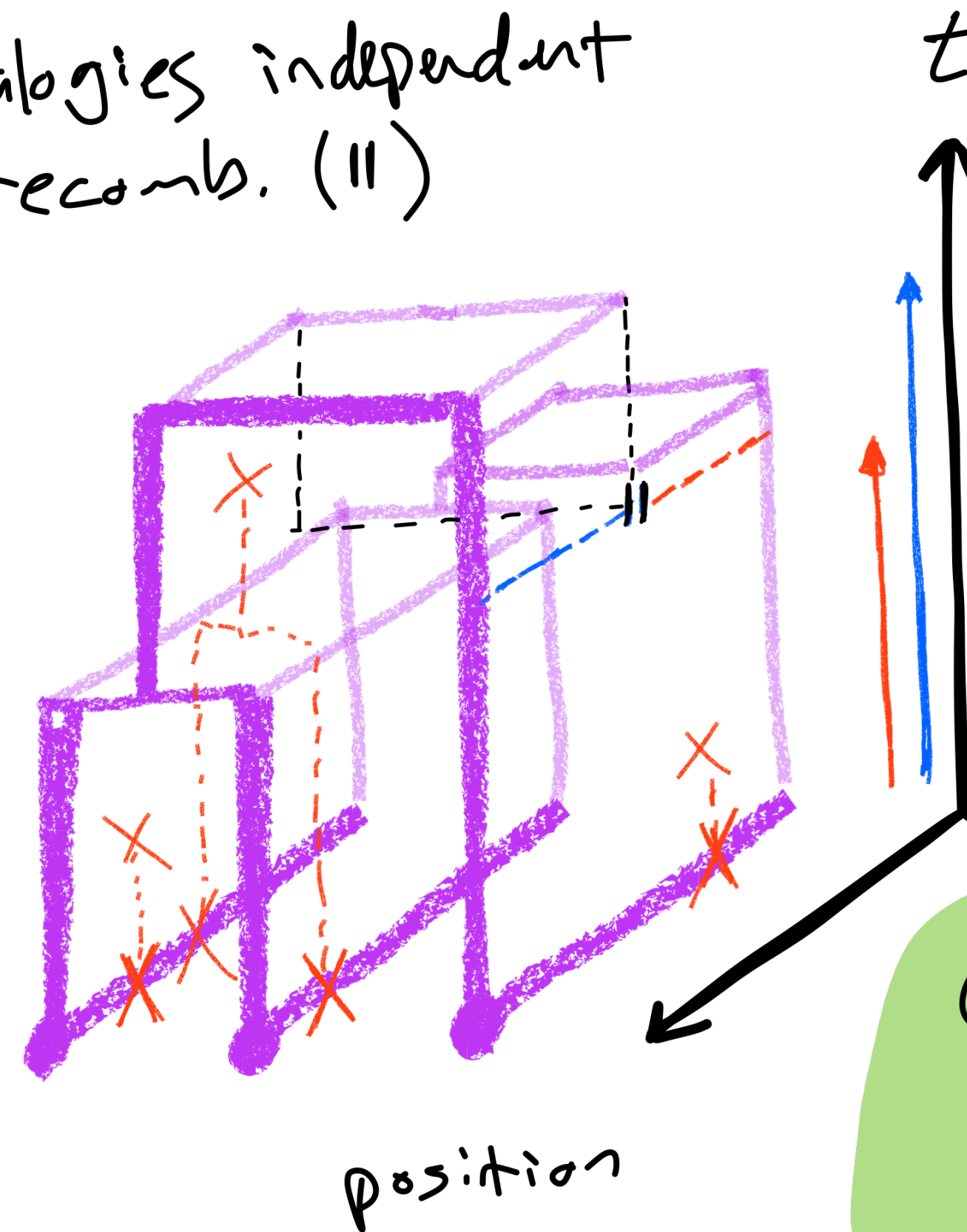
The coalescent with recombination

Without recombination:
coalescences and mutations



With recombination:
add Poisson recombination events!

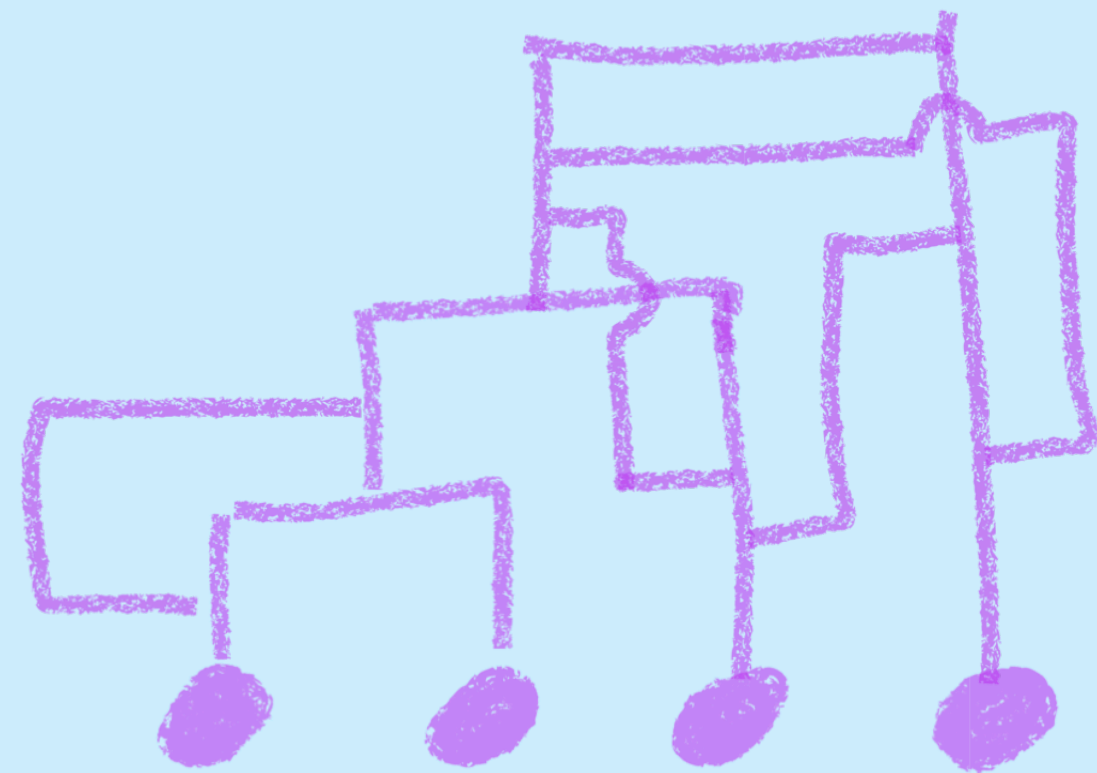
Genealogies independent
above recomb. (II)



Genealogies
decorrelate
(sample different
TMRCA)

Question: does this
change $E[S]$, $E[\pi]$,
 $E[Z_i]$?

THE ANCESTRAL



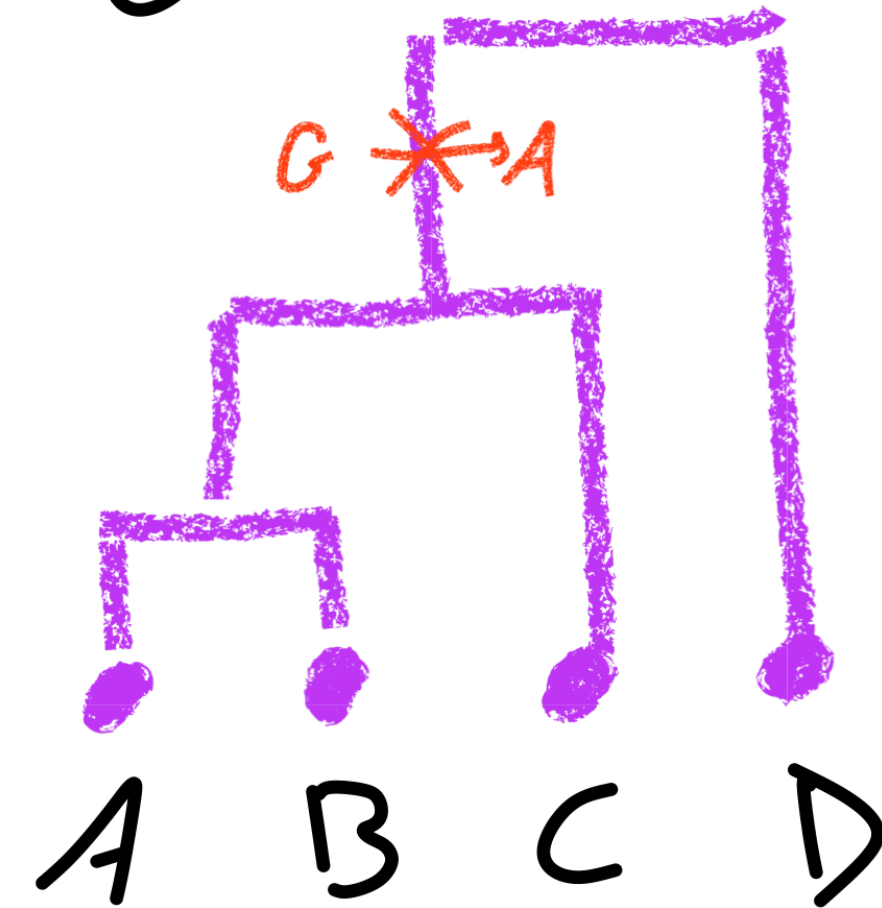
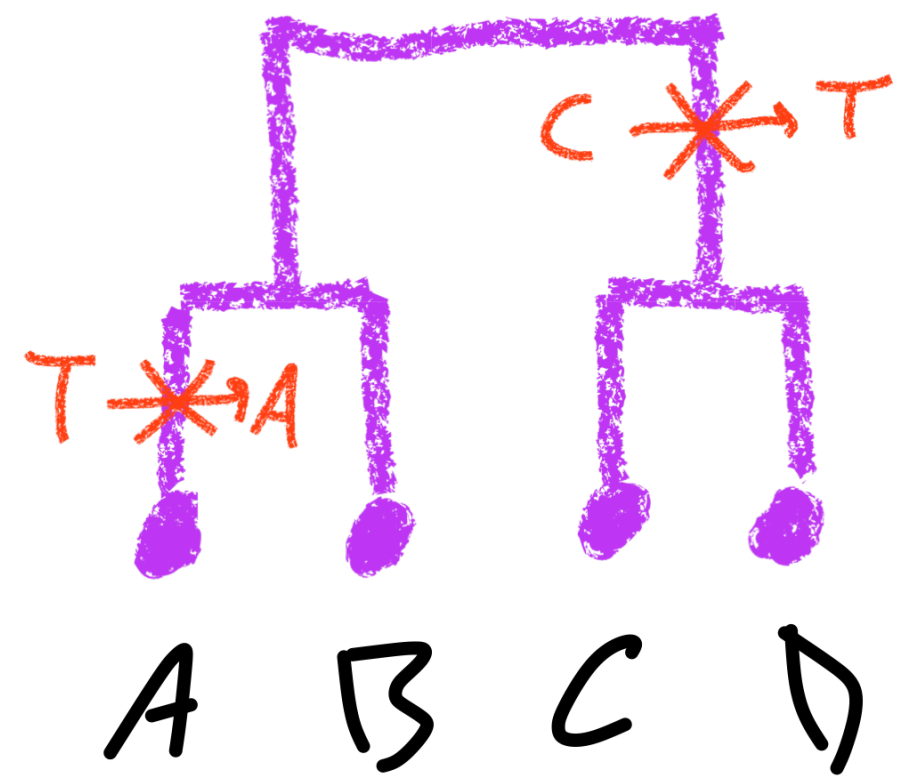
RECOMBINATION

GRAPH

"ARG"

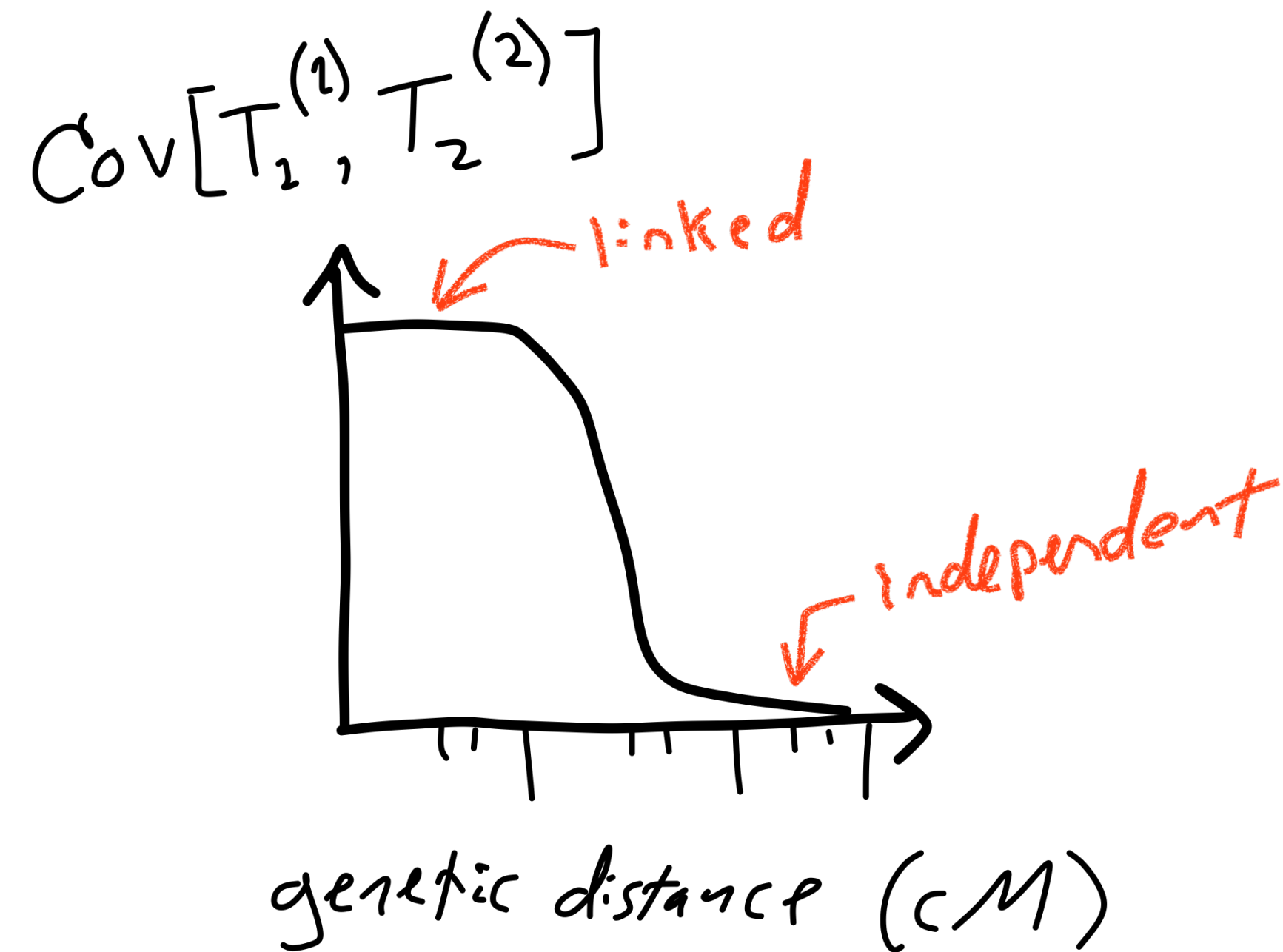
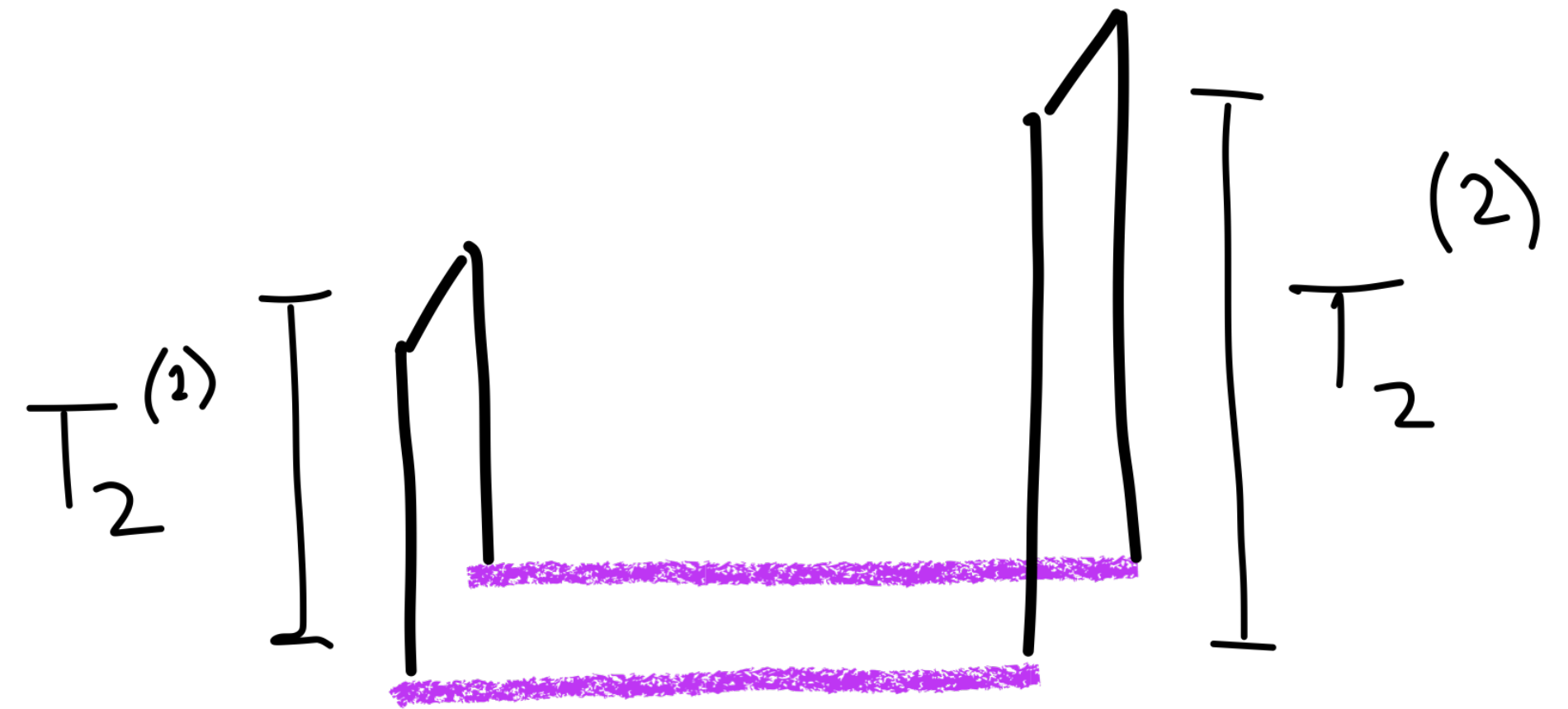
A sequence of trees

A	...	C	...	A	...	//	...	A	...
B	...	C	...	T	...	//	...	A	...
C	...	T	...	T	...	//	...	A	...
D	...	T	...	T	...	//	...	G	...



Decay of linkage

Larger genetic distance
 \Rightarrow recombination more likely



Deeper coalescence
 \Rightarrow shorter span

$X/t \sim \exp(4Nr t)$

genetic dist. to recomb. \uparrow recombination rate

$T_{MRC A}$

The equation shows the relationship between the genetic distance to recombination (X/t) and the recombination rate (4Nr t). The graph shows the time to the most recent common ancestor (T_{MRC A}) as a function of genetic distance to recombination. The curve starts at a high value for small distances and decays to zero for large distances. The horizontal bar at the base of the tree represents the time to the most recent common ancestor (T_{MRC A}).

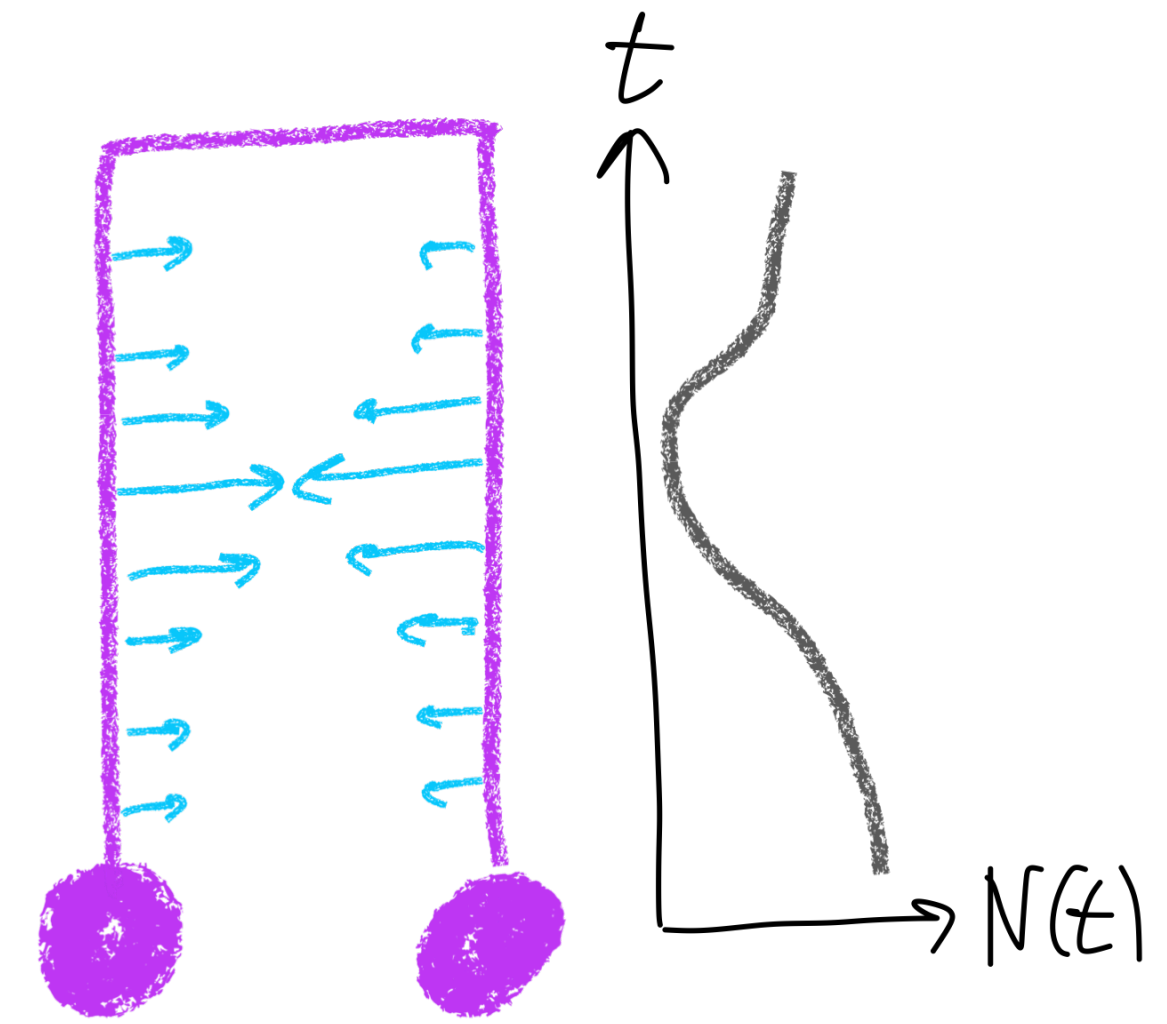
Previously on...

Population size determines coalescence rate

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

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

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



The details:

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

big $N \rightarrow P(t_i) = \frac{\binom{i}{2}}{2N(t_i)} \exp\left(-\binom{i}{2} \int_0^{t_i} \frac{ds}{2N(s)}\right)$

inhomogeneous Poisson process

Coalescent HMMs

Pairwise sequential Markov coalescent (PSMC)

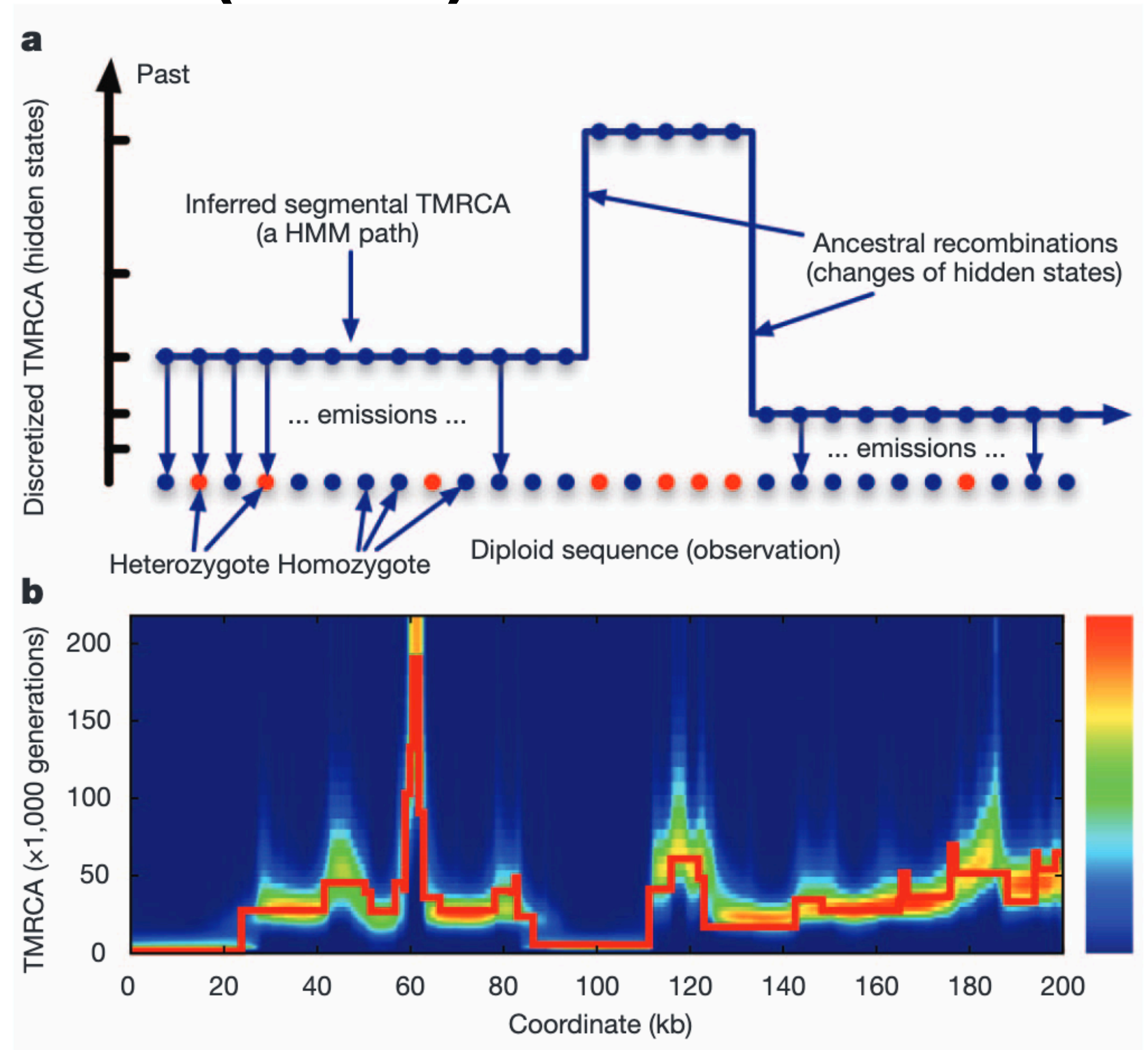
Hidden state: T_{MRCA}

Emission: heterozygosity $00100110...$

Transitions:

$$P(t|s) = \begin{cases} e^{-\rho s}, & s = t \\ (1 - e^{-\rho s})q(t|s), & s \neq t \end{cases}$$

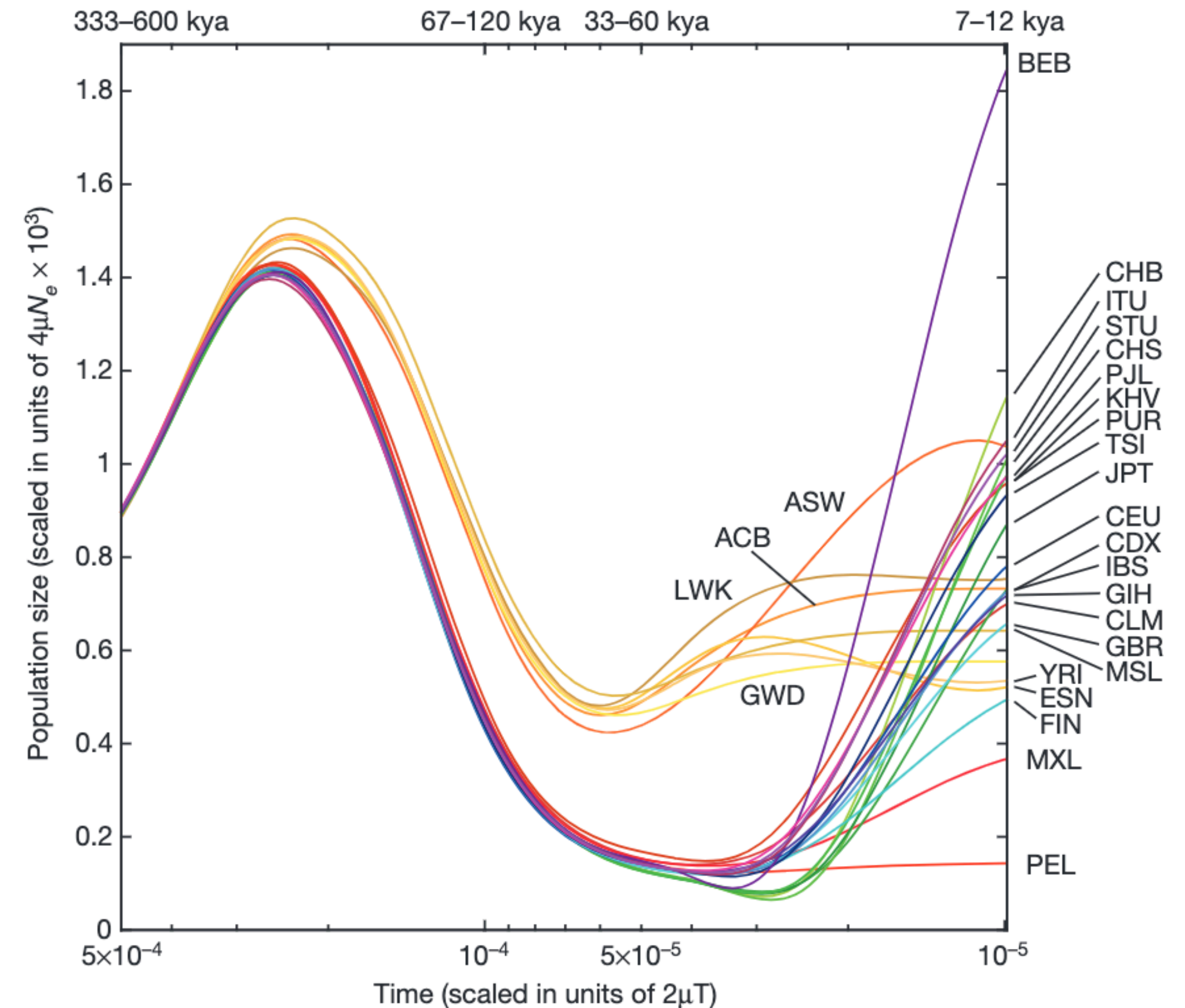
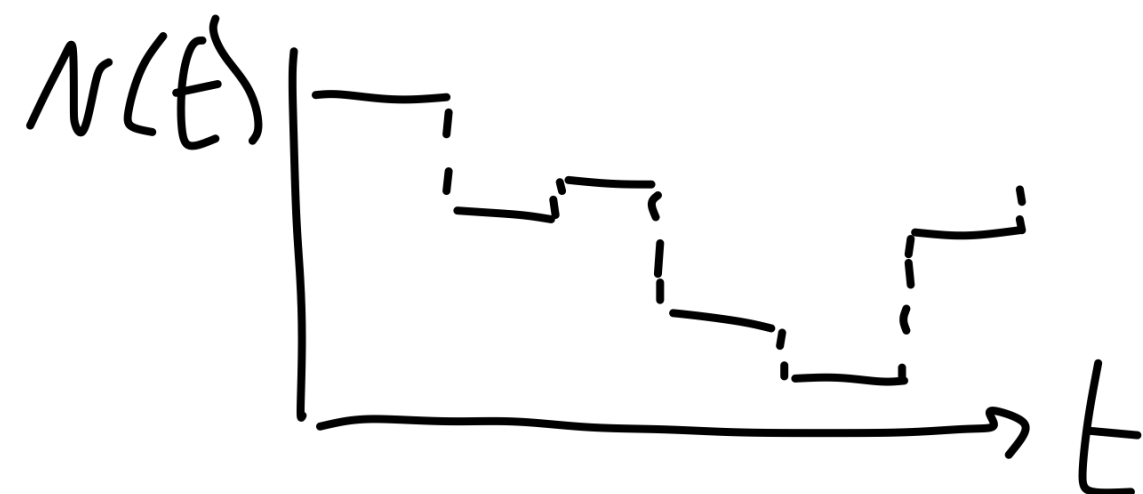
↑
complicated integral
of history $N(t)$



A global reference for human genetic variation

The 1000 Genomes Project Consortium*

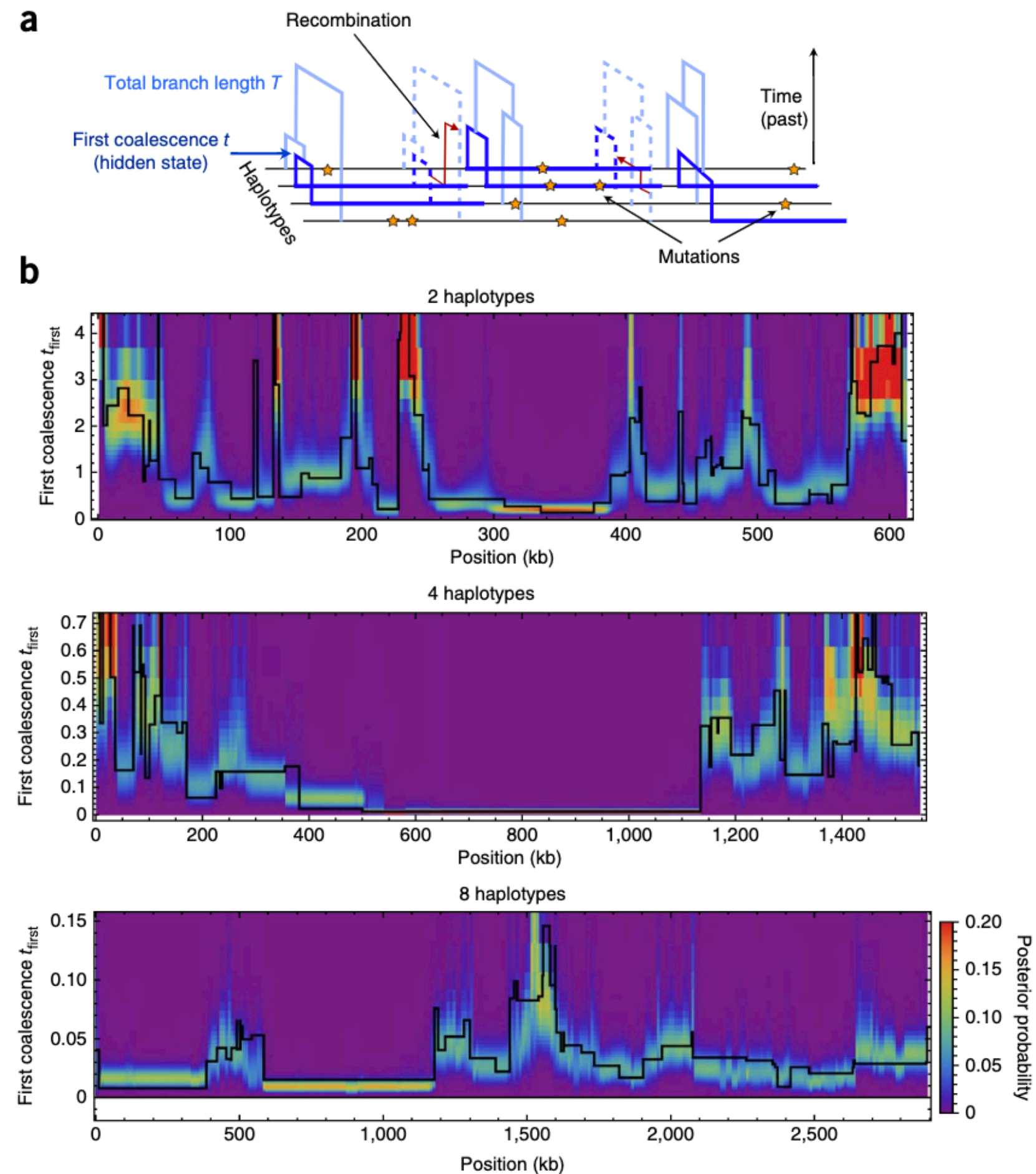
Averaging indiv. estimates
to get pop estimate



Generalizing beyond diploids

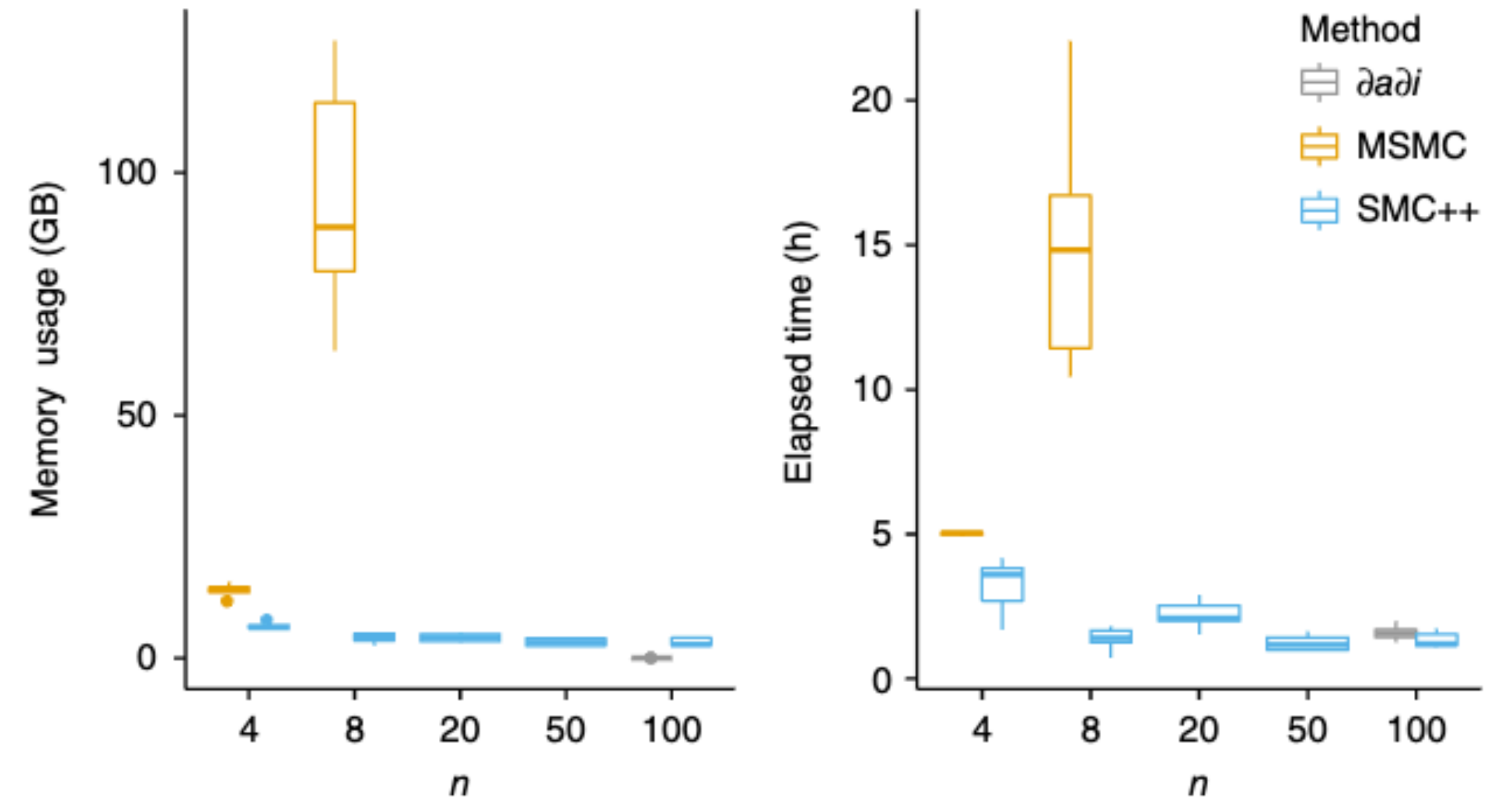
Inferring human population size and separation history from multiple genome sequences

Stephan Schiffels & Richard Durbin



Robust and scalable inference of population history from hundreds of unphased whole genomes

Jonathan Terhorst¹, John A Kamm^{1,2} & Yun S Song¹⁻⁴



The future: tree sequence inference

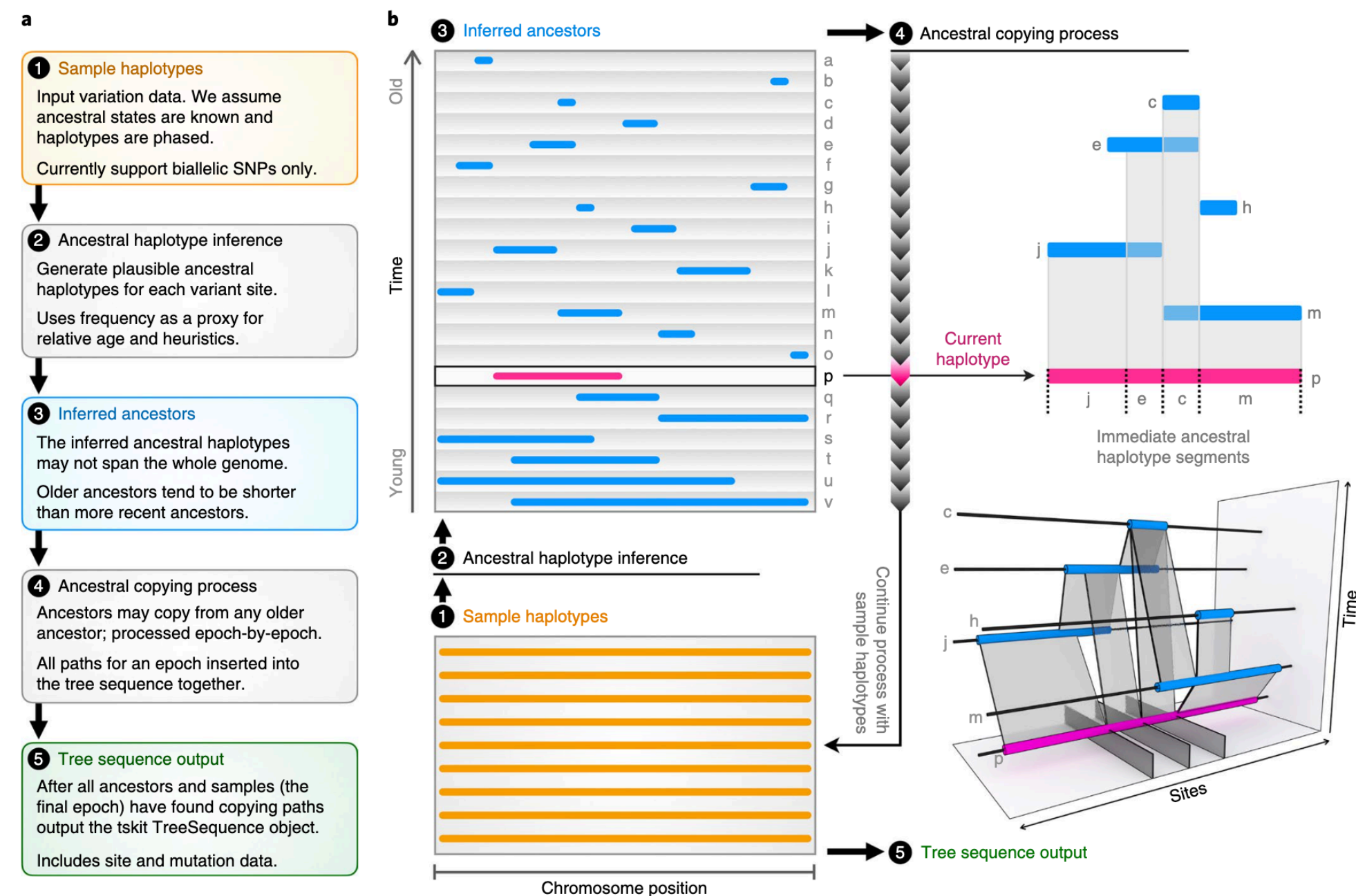
ARTICLES

<https://doi.org/10.1038/s41588-019-0483-y>

nature
genetics

Inferring whole-genome histories in large population datasets

Jerome Kelleher¹*, Yan Wong, Anthony W. Wohns¹, Chaimaa Fadil¹, Patrick K. Albers¹ and Gil McVean¹



nature
genetics

ARTICLES

<https://doi.org/10.1038/s41588-019-0484-x>

A method for genome-wide genealogy estimation for thousands of samples

Leo Speidel¹, Marie Forest², Sinan Shi¹ and Simon R. Myers^{1,3*}

