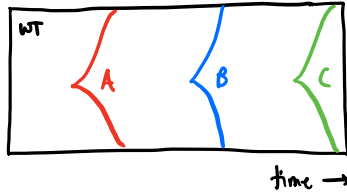


## **Chapter II**

# **Neutral theory and the coalescent**

# Neutral theory & the Coalescent

Successive mutations:



$$\frac{ds(\bar{g})}{dt} = \underbrace{\sim(x-\bar{x})}_{\text{blue}} + \underbrace{\sim L\mu}_{\text{orange}} \xrightarrow{\epsilon} + \underbrace{\sim \rho}_{\text{purple}} + \underbrace{\sim \frac{\pi}{\sqrt{\mu}}}_{\text{green}}$$

$\Rightarrow$   $\sim 1$  variant present @ high freqs

$\Rightarrow$  solved by reducing to  $L=1$  model

\* But genomes in data separated by multiple mut'n's  
(e.g. humans, 2 individuals differ by  $\sim 1$  mut / 1000 bp)

$\Rightarrow$  need to understand what's going on in these cases...

$$\frac{ds(\bar{g})}{dt} = \underbrace{\sim(x-\bar{x})}_{\text{blue}} + \underbrace{\sim L\mu}_{\text{orange}} + \underbrace{\sim \rho}_{\text{purple}} + \underbrace{\sim \frac{\pi}{\sqrt{\mu}}}_{\text{green}}$$

$\Rightarrow$  one other limit that's well understood:

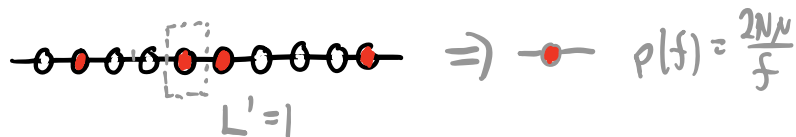
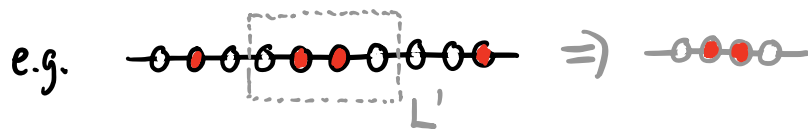
neutral evolution in nonrecombining genome

When  $X(\vec{g})=0 + e=0$ , left with:  $(\mu_e = \nu_e)$

$$\frac{dS(\vec{g})}{dt} = \underbrace{\sum_{|\vec{g}'-\vec{g}|=1} \sum_e \mu_e f(\vec{g}') \left[ g_e(1-g_e') + (1-g_e)g_e' \right]}_{\text{incoming mutations}} - \underbrace{\sum_e \mu_e f(\vec{g})}_{\text{outgoing mutations}}$$

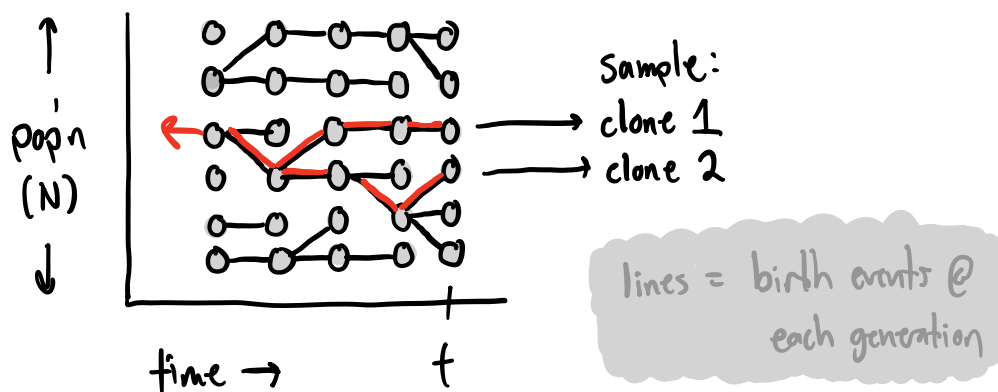
$$+ \sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}') \quad \text{genetic drift}$$

Key insight: sites don't actually influence each other (because neutral)

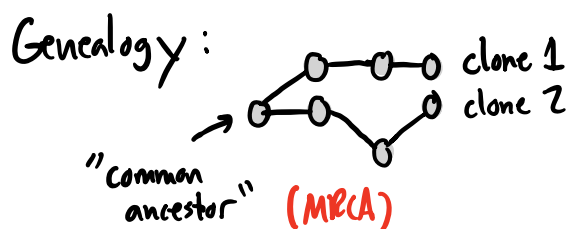


⇒ 2nd key insight: can take  $L'=0$  —

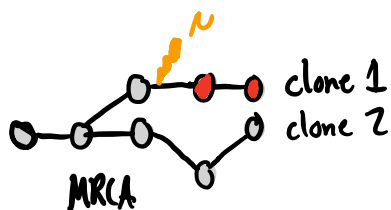
E.g. simulation of neutral pop'n in Wright-Fisher model:



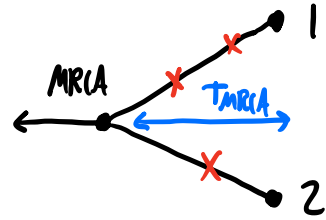
⇒ key insight: lines also = **genealogical relationships** backward in time!



⇓ differences between sampled individuals  
= mutations on genealogy



⇒ Mut's @ site  $l \approx$  Poisson Process  
w/ rate  $\mu_e$  on each branch



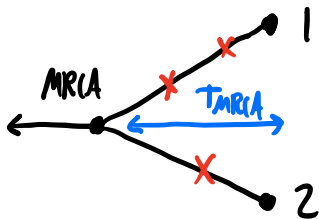
↳ total # mut's  
 $\sim$  Poisson( $2T_{MRCA}\mu_e$ )

⇒ 2 extreme limits:

(1)  $\mu_e T_{MRCA} \ll 1 \Rightarrow$  0 or 1 mutations on whole tree

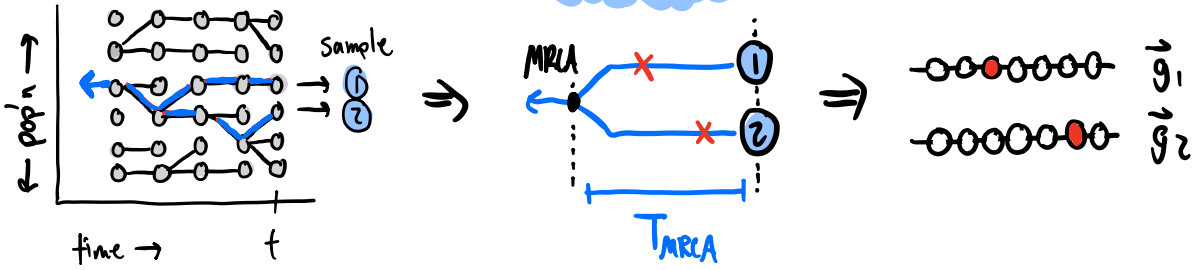
$$\Rightarrow \Pr(\text{genetic diff @ site } l \mid T_{MRCA}) \approx$$

(2)  $\mu_e T_{MRCA} \gg 1 \Rightarrow$  lots of forward & backward mutations along each branch.



$$\Rightarrow \Pr(\text{genetic diff @ site } l \mid T_{MRCA}) =$$

Recap:



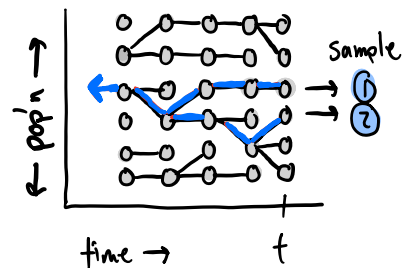
⇒ Given genealogy ( $T_{MRCA}$ ), mutations occur as Poisson Process along each branch ("mutation painting")

$$\Pr[\text{difference @ site } l \mid T_{MRCA}] \approx \begin{cases} 2\mu_e T_{MRCA} & \text{if } \mu T_{MRCA} \ll 1, \\ 1/2 & \text{else.} \end{cases}$$

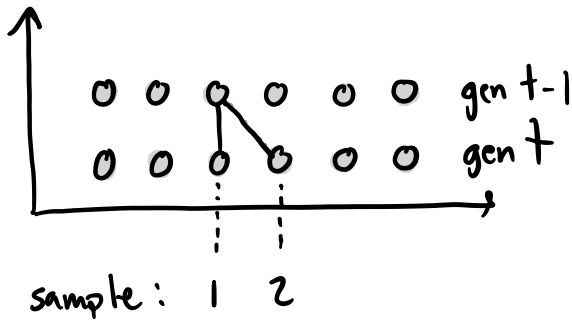
Question: what determines genealogy ( $T_{MRCA}$ )?

⇒ Note:  $T_{MRCA}$  is random quantity

(genealogy will vary from sample-to-sample & simulation-to-simulation...)



⇒ key insight: start from present & work backward in time:



→ "coalesced"

⇒ Two individuals share ancestor in previous gen w/ probability:

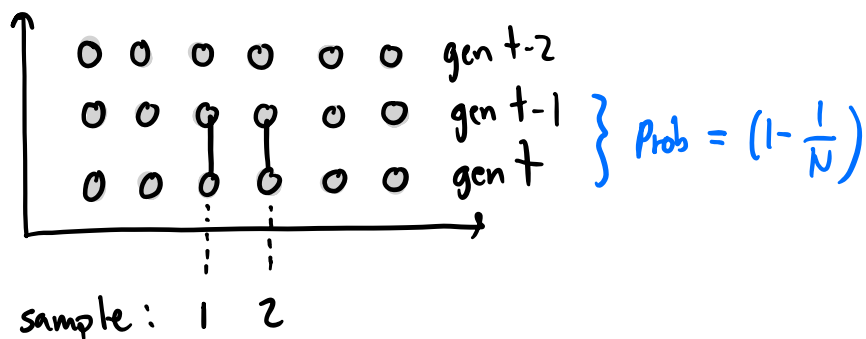
prob that both draw same

$$N \times \left(\frac{1}{N}\right) \times \left(\frac{1}{N}\right) = \frac{1}{N}$$

↳ # possible ancestors

⇒ w/ probability  $\frac{1}{N}$  ⇒  $T_{MRC} = 1$

⇒ otherwise, diff ancestors in gen t-1 ⇒ repeat!

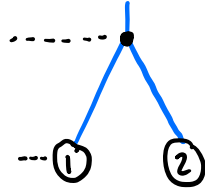


Process repeats itself w/ next gen:

$$\Rightarrow \text{w/ prob } \frac{1}{N} \left(1 - \frac{1}{N}\right) \Rightarrow T_{MRCA} = 2$$

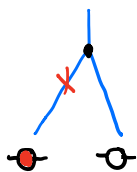
$$\Rightarrow \text{w/ prob } \frac{1}{N} \left(1 - \frac{1}{N}\right)^2 \Rightarrow T_{MRCA} = 3$$

$\Rightarrow$  coalescence is also a Poisson Process w/ rate  $\frac{1}{N}$ !

$$\Rightarrow T_{MRCA} \sim \text{Exponential}(N)$$


$$\Rightarrow \langle T_{MRCA} \rangle = N \quad \sqrt{\text{Var}(T_{MRCA})} = N$$

$\Rightarrow$  total probability of mutation @ site  $l$  is integral over  $T_{MRCA}$ :



$$\Pr[\text{difference @ site } l] = \int \underbrace{\Pr(\text{diff} | T_{MRCA})}_{\text{mutation painting}} \underbrace{p(T_{MRCA})}_{\text{coalescent}} dT_{MRCA}$$

$$\approx_{(N \ll 1)} \int 2N_e T_{MRCA} \cdot p(T_{MRCA}) dT_{MRCA}$$

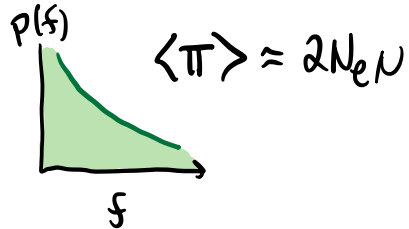


$$= 2N_e \langle T_{MRCA} \rangle$$

$$= 2N\mu_e$$

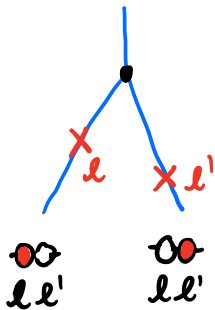
⇒ matches our previous result

Since  $\langle \pi \rangle \equiv \Pr(\text{diff @ site } e)$



⇒ Distribution of  $T_{MRCA}$  becomes more important when considering mutations @ multiple sites, e.g.

$$\Pr(\text{diff @ site } e \text{ \& } e') = \int \Pr[\pi_e=1, \pi_{e'}=1 | T_{MRCA}] p(T_{MRCA}) dT_{MRCA}$$



$$= \int \underbrace{\Pr[\pi_e=1 | T_{MRCA}] \Pr[\pi_{e'}=1 | T_{MRCA}]} p(T_{MRCA}) dT_{MRCA}$$

mut's are neutral, so can't affect each other!

$$= \int (2N_e T_{MRCA}) \cdot (2N_{e'} T_{MRCA}) \cdot p(T_{MRCA}) \cdot dT_{MRCA}$$

$$= (2N_e) \cdot (2N_{e'}) \cdot \langle T_{MRCA}^2 \rangle = (2N_e) \cdot (2N_{e'}) \cdot (2N^2)$$

$$= 2 \cdot (2N_e N) \cdot (2N_{e'} N)$$

$$= 2 \cdot \Pr(\pi_e) \cdot \Pr(\pi_{e'}) \geq \Pr(\pi_e) \Pr(\pi_{e'})$$

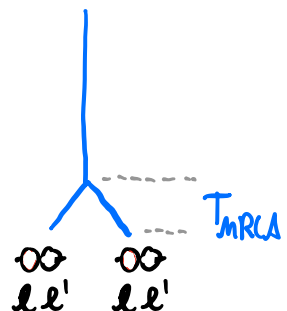
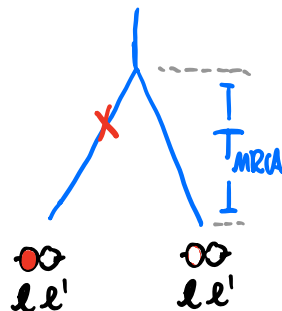
Upshot: joint prob of mut'ns is not independent:

$$\Pr(\pi_{e'}=1 | \pi_e=1) = \frac{\Pr(\pi_e=1, \pi_{e'}=1)}{\Pr(\pi_e)} = 2 \Pr(\pi_{e'}=1)$$

But previously said that neutral mutations can't influence each other directly...

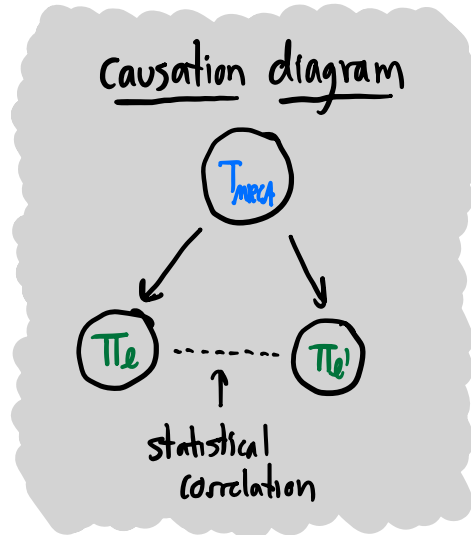
$\Rightarrow$  what's going on?

$\Rightarrow$  consider 2 trees:

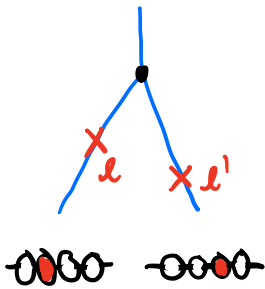
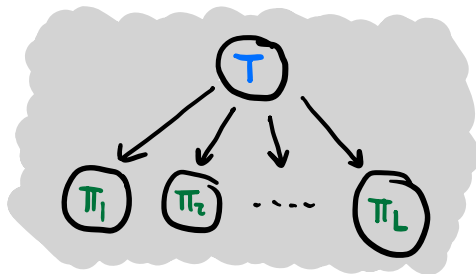


$\Rightarrow$  conditioned on  $\pi_e=1$ , likely had **bigger-than-avg  $T_{MRCA}$**

$\Rightarrow$  i.e. mutations don't interact,  
but are still coupled  
by **shared genealogy**



$\Rightarrow$  can keep adding  
more sites this way...



$\Rightarrow$  when  $\mu_e T_{MRCA} \ll 1$ , most mutations  
will occur @ **unique site in genome**  
"infinite-sites approximation"

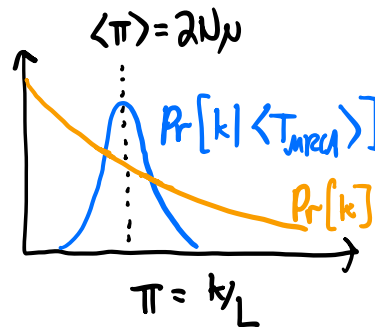
$\Rightarrow$  total # mut'ns ( $k$ ) is Poisson Process w/ rate  $U \equiv \sum_{e=1}^L \mu_e$

$$\Rightarrow \Pr[k | T_{\text{MRCA}}] = \frac{(2UT_{\text{MRCA}})^k}{k!} e^{-2UT_{\text{MRCA}}}$$

$$\begin{aligned} \Rightarrow \Pr[k] &= \int \Pr[k | T_{\text{MRCA}}] p(T_{\text{MRCA}}) dT_{\text{MRCA}} \\ &= \int \frac{(2UT)^k}{k!} e^{-2UT} \frac{1}{N} e^{-T/N} dT \end{aligned}$$

$$\Rightarrow \Pr[k] = \frac{(2NU)^k}{(2NU+1)^{k+1}}$$

total # diffs  $\text{---} \circ \bullet \circ \circ \text{---}$   
 btw 2 genomes  $\text{---} \circ \circ \circ \bullet \text{---}$



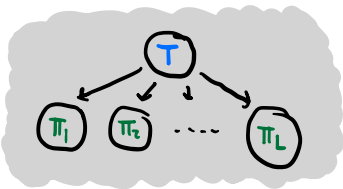
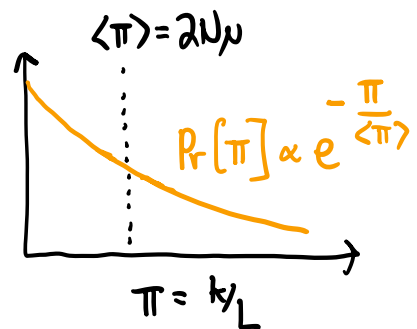
⇒ one advantage of coalescent approach:

⇒ simple predictions for uncertainty in  $\pi$  (not just avg)

$$\text{e.g. } \text{Var}(\pi) = \frac{\text{Var}(k)}{L^2} = \frac{(1+2Nu)2Nu}{L^2}$$

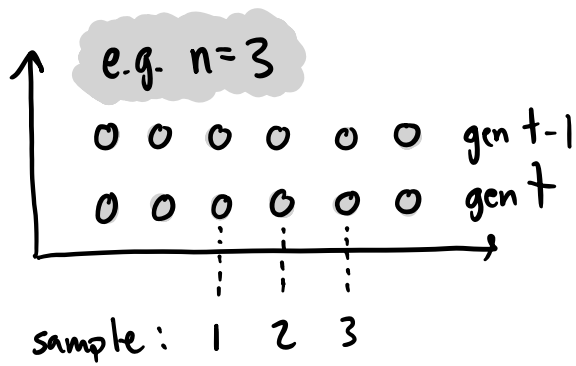
$$\Rightarrow \text{or } C_V^2 \equiv \frac{\text{Var}(\pi)}{\langle \pi \rangle^2} = \frac{1+2Nu}{2Nu} \geq 1$$

$\Rightarrow$  i.e.  $\pi$  does not self-average on a long asexual genome!



$\Rightarrow$  fluct'ns in  $T_{MRCA}$  affect many sites!

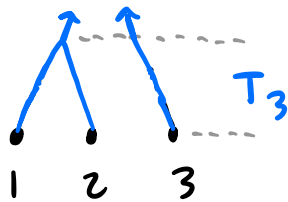
## Larger sample sizes ( $n > 2$ )



$\Rightarrow$  Prob that any 2 share ancestor is  $\frac{1}{N} \left[ \times \binom{3}{2} \text{ pairs} \right]$

$\Rightarrow$  Prob that all 3 share ancestor =  $N \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) = \frac{1}{N^2}$

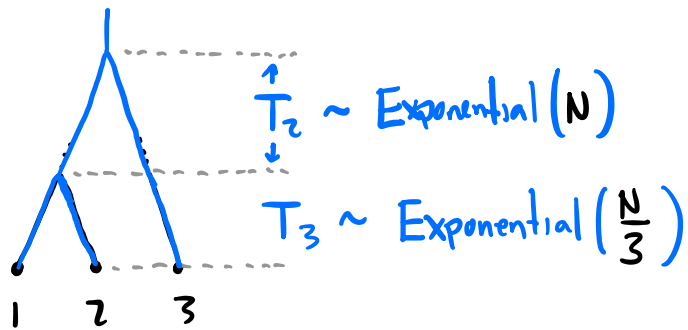
$\Rightarrow$  when  $N \gg 1 \rightarrow$  only need to worry about **pairwise coalescence**  
 (known as "Kingman's coalescent") (all pairs are equally likely to coalesce)



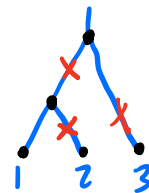
$\Rightarrow$  total prob of coalescence =  $\frac{3}{N}$  per gen

$\Rightarrow T_3 \sim \text{Exponential} \left( \frac{N}{3} \right)$

⇒ now we have sample of  $n=2...$  ⇒ repeat!



⇒ Done! can now paint on mutations...

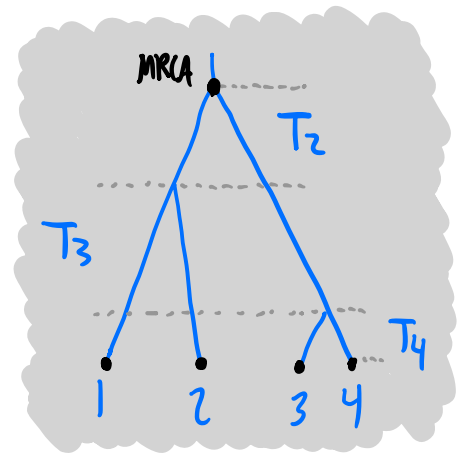


Easily generalizes to sample of size n:

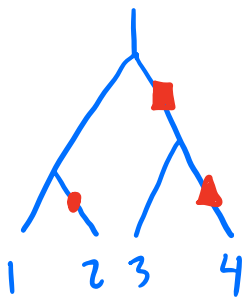
① @ each step, only consider coalescence between pairs of lineages

② Time until next coalescence event is  $T_n \sim \text{Exponential}(N/\binom{N}{2})$

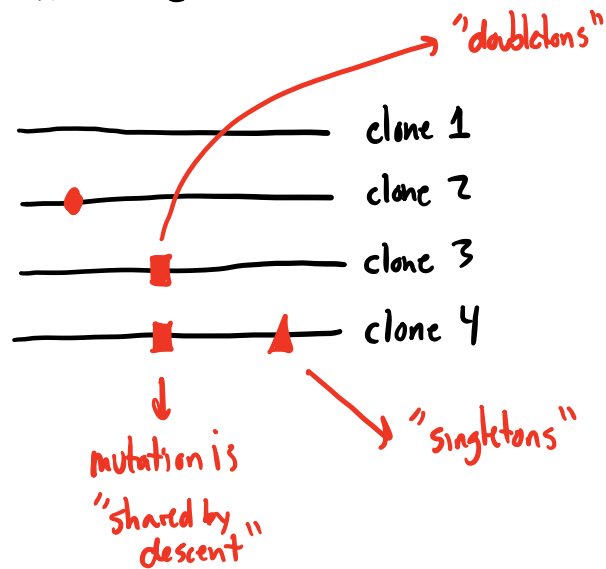
③ choose random pair to coalesce repeat!



④ then can paint mutations on @ end:



⇒



⇒ easy to simulate for  $n > 2$ , but hard to calculate...



e.g.  $\langle \# \text{ doubletons in sample } n=4 \rangle = \langle \text{ [tree with 2 doubletons]} + \text{ [tree with 1 doubleton]} \rangle$

- $\Rightarrow$  must avg over:
- ① tree topologies
  - ② branch lengths | topology
  - ③ mutation painting | branch lengths

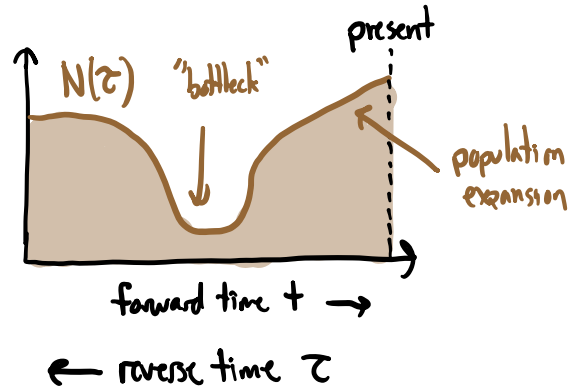
$\Rightarrow$  compare to single-locus prediction (easy!)

$$\langle \# \text{ doubletons in } n=4 \rangle = \int \binom{4}{2} f^2 (1-f)^2 \cdot \left( \frac{2N\mu}{f} \right) \cdot df = N\mu$$

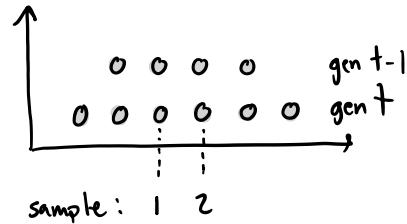
$\Rightarrow$  why use coalescent picture then??

Answer: coalescent picture makes it easy to model demography!

e.g. what if  $N$  was not constant, but varied historically in time:



$\Rightarrow$  coalescent picture still works, but coalescent prob  $\rightarrow 1/N(\tau)$

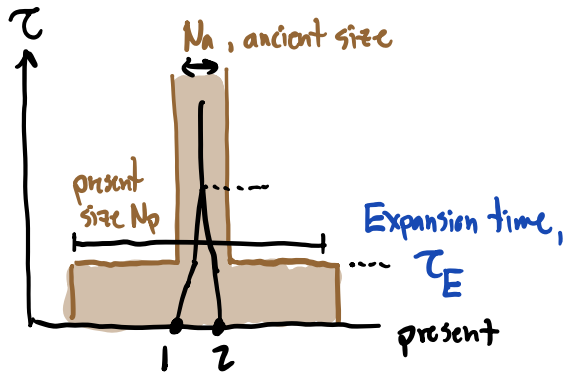


$\Rightarrow$  coalescence = "inhomogeneous" Poisson process:

$$\Rightarrow \Pr[T_2 > \tau] = \prod_{\tau'=1}^{\tau} \left[ 1 - \frac{1}{N(\tau')} \right] \approx e^{-\int_0^{\tau} \frac{dz'}{N(z')}}$$

$$\Rightarrow \Pr[T_2 = \tau] = \frac{1}{N(\tau)} e^{-\int_0^{\tau} \frac{dz}{N(z)}}$$

Simple example: rapid expansion in recent past



If  $N_p \gg \infty$  &  $\tau_E \ll N_p$ :

- ① no coalescence until  $\tau_E$
- ② coalescence @ rate  $\frac{1}{N_a}$  after

$$\Rightarrow \langle T_2 \rangle = \tau_E + N_a$$

$$\Rightarrow \langle \pi \rangle = 2\mu \langle T_2 \rangle = 2\mu(\tau_E + N_a) \approx 2\mu N_a \quad \left( \text{if } \tau_E \ll N_a \right)$$

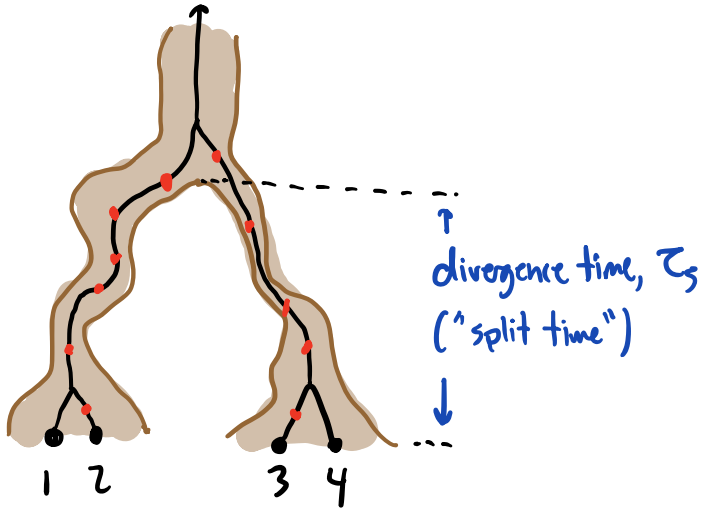
( compare to forward time calc...  $\frac{df}{dt} = \mu(1-f) - \nu f + \sqrt{\frac{f(1-f)}{N(t)}} \eta(t) \Rightarrow p(f,t)$  )

can revisit our earlier puzzle: if  $N_p \cdot \mu \sim 100$  in humans  
why  $\langle \pi \rangle \sim 10^{-3}$ ?

$\Rightarrow$  one answer:  $N(t)$  was smaller backward in time!

$$\Rightarrow N_a \approx 10^5 \quad (\tau_E \ll 10^5 \text{ gens})$$

can also easily add population structure



$\Rightarrow$  Pr(coalescence)  
 between pop'ns = 0  
 until time  $\tau = \tau_s$

$\Rightarrow$  much of pop gen is about inferring these demographic models

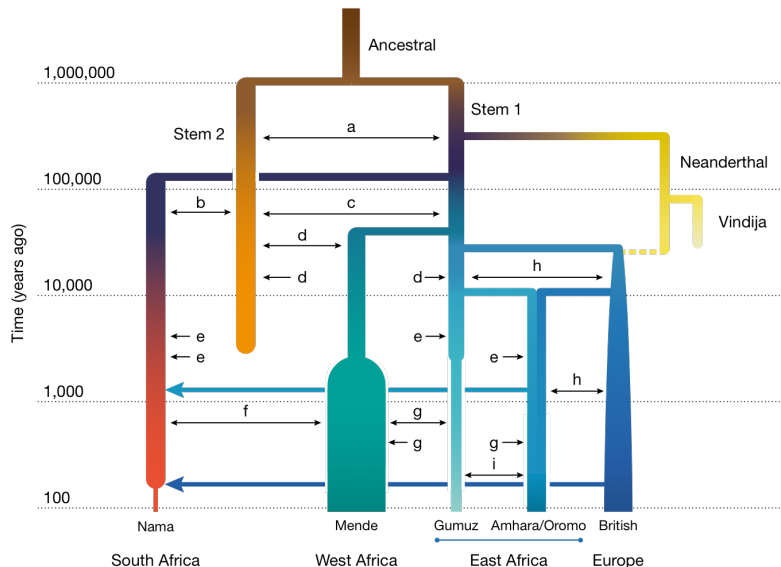
e.g. :

Article

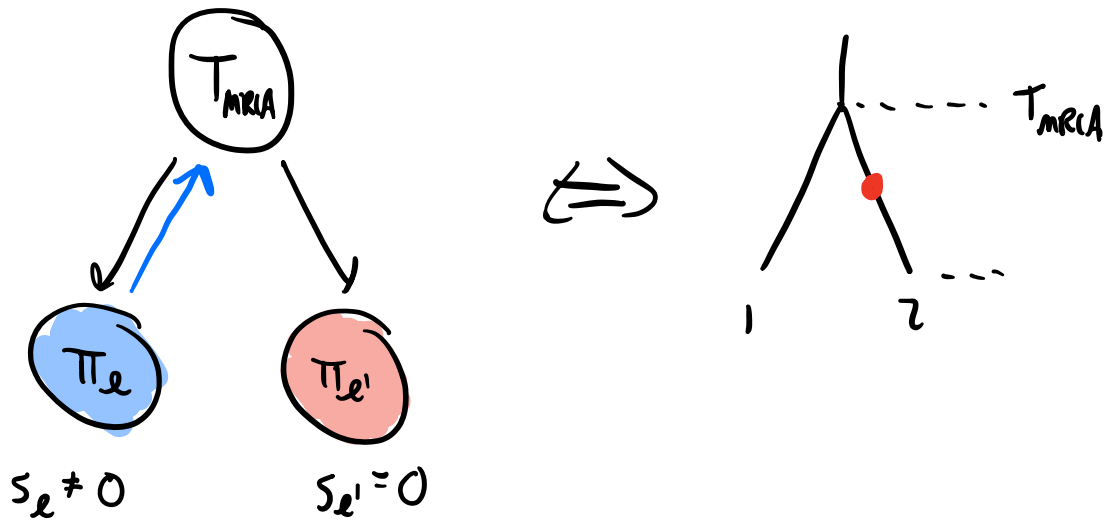
## A weakly structured stem for human origins in Africa

Aaron P. Ragsdale<sup>1</sup>, Timothy D. Weaver<sup>2</sup>, Elizabeth G. Atkinson<sup>3</sup>, Eileen G. Hoal<sup>4,5,6</sup>, Marlo Möller<sup>4,5,6</sup>, Brenna M. Henn<sup>2,7,8,9</sup> & Simon Gravel<sup>8,9,10</sup>

Published online: 17 May 2023



⇒ downside: hard to add **selection** back in to picture...

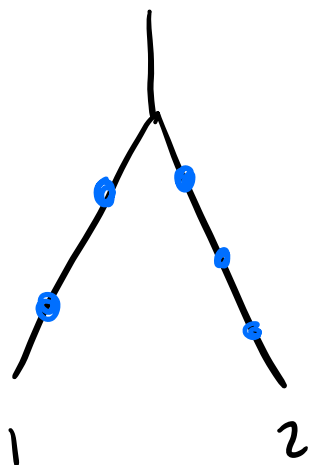


⇒ when is this going to be an issue?

⇒ for  $L=1$  case, needed  $N|s| \ll 1$  for effectively neutral.

⇒ for  $L \gg 1$ , selection looks like  $(x(\frac{1}{2}) - \bar{x}(t)) f(\frac{1}{2})$   
vs  
 $sf(1-s)$  in  $L=1$

⇒ suggests:  $N|x(\frac{1}{2}) - \bar{x}| \ll 1$  for neutrality



① assume effective neutrality:

$\Rightarrow$  total # mutations  $\approx NU$

$$|X(\vec{z}) - X(\vec{z}_0)| = \sqrt{NU s^2}$$

$\Rightarrow$  self consistent:

$$(NU)(Ns)^2 \ll 1$$

e.g.  $Ns \sim 0.1$  (neutral in single locus setting)

$$NU = \langle \pi \rangle L = \begin{cases} 10^4 & \text{for bacteria in a gut} \\ 10^6 & \text{for humans.} \end{cases}$$

$\downarrow$

$$\sqrt{10^4 \cdot (10^{-1})^2} = 10 \gg 1$$