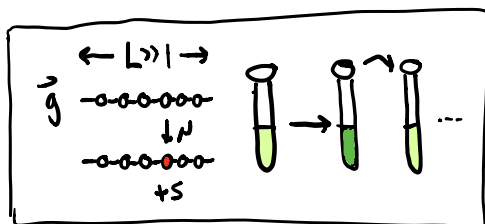


Announcements: PSET 4 DUE 3/9/21; Office hrs ~~12:30~~ → 1pm Thurs

Last time: Multi-locus models of evolution ($L \gg 1$)



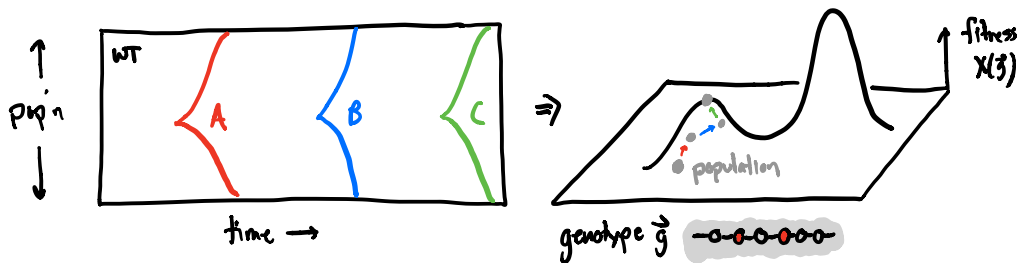
$$\frac{d\bar{f}(\vec{g})}{dt} = \underbrace{[\bar{X}(\vec{g}) - \bar{X}(t)] f(\vec{g})}_{\text{selection (nonlinear)}} + \underbrace{\sum_{\vec{g}'} M(\vec{g}' \rightarrow \vec{g}) f(\vec{g}') - M(\vec{g} \rightarrow \vec{g}') f(\vec{g})}_{\text{mutation (linear, "local")}}$$

$$+ \underbrace{\rho \sum_{\vec{g}_1, \vec{g}_2} T(\vec{g}_1, \vec{g}_2 \rightarrow \vec{g}) f(\vec{g}) - \rho f(\vec{g})}_{\text{recombination (nonlinear, non-local)}}$$

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

No closed-form sol'n in the general case...

① Successive mutations regime [$N \cdot L \cdot \mu \rightarrow 0$]



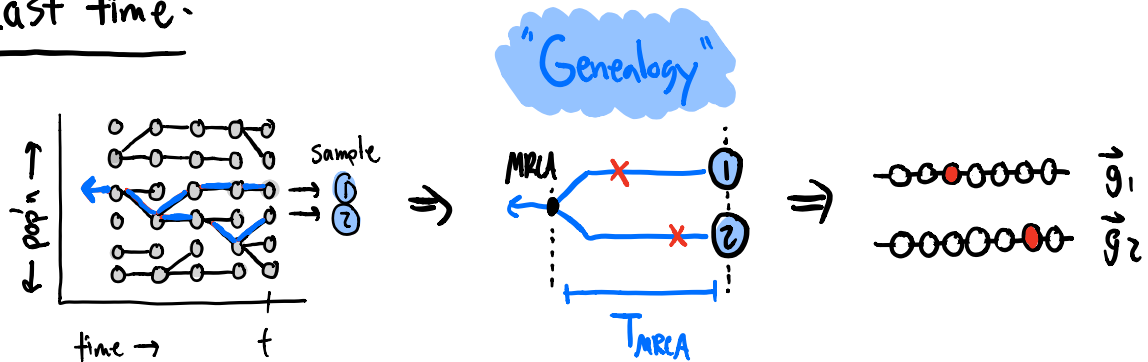
Today: ② The Neutral Limit (natural selection ≈ 0)

↳ new tool: "coalescent theory"

↓

"backward-time approaches"

Last time:



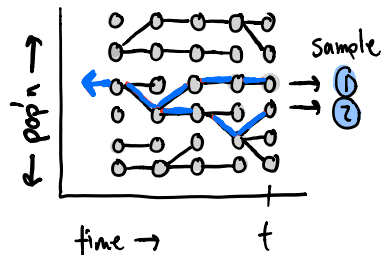
\Rightarrow Given genealogy (T_{MRCA}), mutations occur as Poisson Process along each branch ("mutation painting")

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

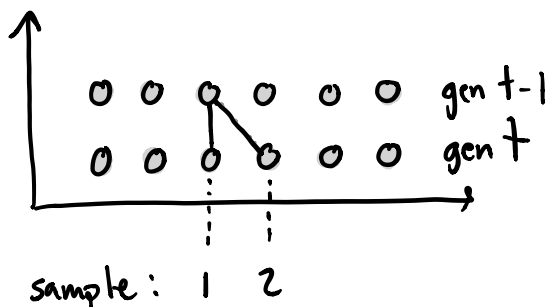
Today: what determines genealogy (T_{MRCA})?

\Rightarrow 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:



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

"coalesced"

total # of ancestors

↑

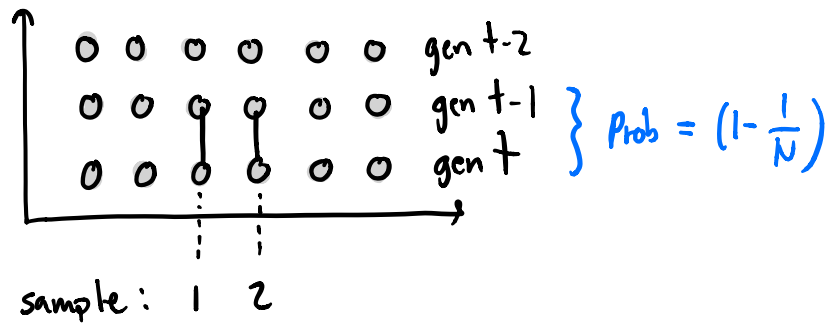
prob that 2 individuals draw it.

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

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

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

Process repeats itself w/ next gen:



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

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

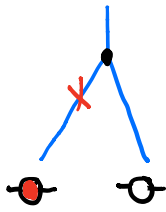
\Rightarrow coalescence is also a Poisson Process w/ rate $\frac{1}{N}$!

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

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

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

$$\Pr(\text{diff @ site } e) = \int \underbrace{\Pr(\text{diff} | T_{\text{MRC}})}_{\text{mutation painting}} \underbrace{P(T_{\text{MRC}})}_{\text{coalescent}} dT_{\text{MRC}}$$

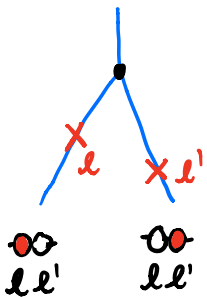


$$\approx \int 2N_e T_{\text{MRC}} P(T_{\text{MRC}}) dT_{\text{MRC}} = 2N_e \langle T_{\text{MRC}} \rangle = 2N_e \mu$$

\Rightarrow matches our previous result for $\langle \pi \rangle$, ↖ ✓
 since $\langle \pi \rangle \equiv \Pr(\text{diff @ site } e)$

\Rightarrow Distribution of T_{MRC} becomes more important when considering mutations @ multiple sites, e.g.

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



$$= \int \underbrace{\Pr[\pi_e=1 | T_{\text{MRC}}] \Pr[\pi_{e'}=1 | T_{\text{MRC}}]}_{\text{mutations are neutral so can't influence each other}} P(T_{\text{MRC}}) dT_{\text{MRC}}$$

$$\begin{aligned}
&= \int (z_{N_e} T_{MRCR}) \cdot (z_{N_e'} T_{MRCR}) \cdot p(T_{MRCR}) \cdot dT_{MRCR} \\
&= (z_{N_e}) \cdot (z_{N_e'}) \cdot \langle T_{MRCR}^2 \rangle = (2\mu_e) \cdot (2\mu_e') \cdot (2N^2) \\
&= 2 \cdot (2\mu_e N) \cdot (2\mu_e' N) \\
&= 2 \cdot \Pr(\pi_e) \cdot \Pr(\pi_{e'}) \geq \Pr(\pi_e) \Pr(\pi_{e'})
\end{aligned}$$

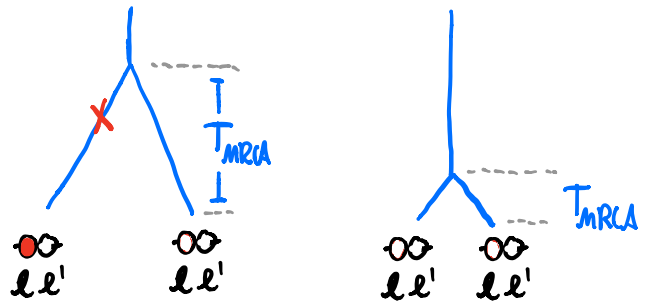
\Rightarrow joint prob of mutations is not independent

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

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

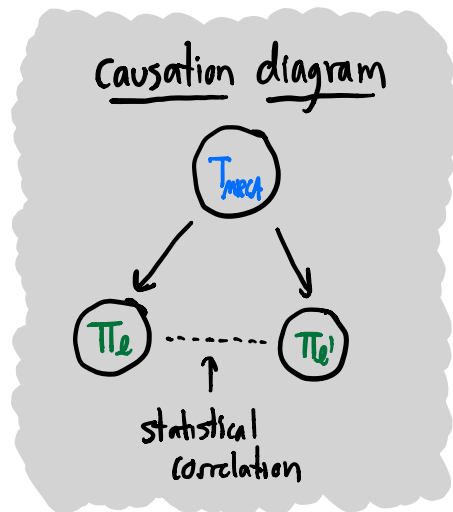
⇒ what's going on?

⇒ consider 2 trees:

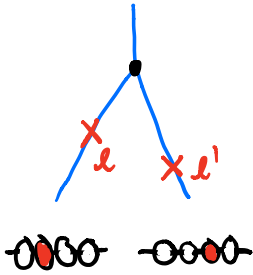
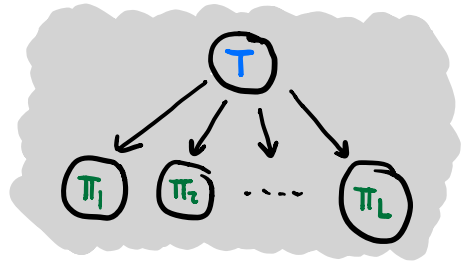


⇒ conditioned on $\pi_e = 1$, likely had **bigger-than-avg T_{MRCAs}**

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



⇒ can keep adding more sites this way...



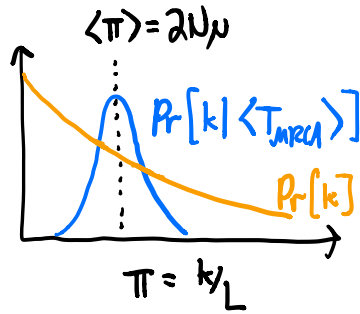
⇒ when $\mu_e T_{MRCA} \ll 1$, most mutations will occur @ **unique site in genome**
 "infinite-sites approximation"

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

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

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

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



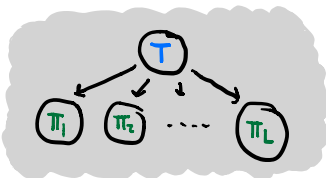
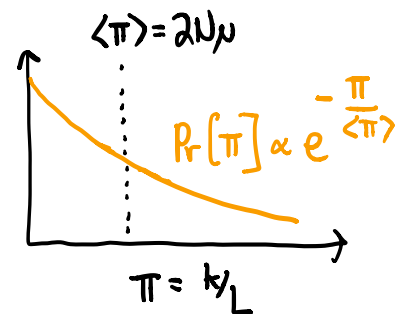
\Rightarrow one advantage of coalescent approach:

\Rightarrow simple predictions for uncertainty in π (not just avg)

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

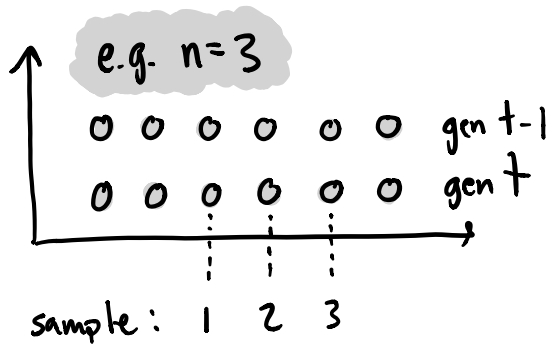
$$\Rightarrow \text{or } C_v^2 \equiv \frac{\text{Var}(\pi)}{\langle \pi \rangle^2} = \frac{1+2NU}{2NU} \geq 1$$

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



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

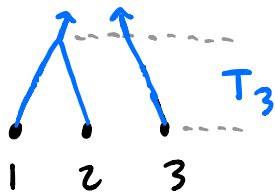
Larger sample sizes ($n > 2$)



\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 = $N \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) = \frac{1}{N^2}$

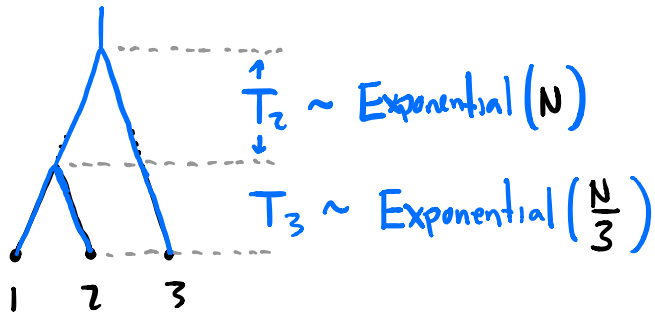
\Rightarrow when $N \gg 1 \rightarrow$ only need to worry about **pairwise coalescence**
(known as "Kingman's coalescent") (all pairs are equally likely to coalesce)



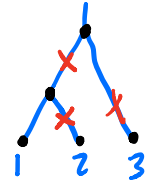
\Rightarrow total prob of coalescence = $\frac{3}{N}$ per gen

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

\Rightarrow now we have sample of $n=2 \dots \Rightarrow$ repeat!



\Rightarrow Done! can now paint on mutations...

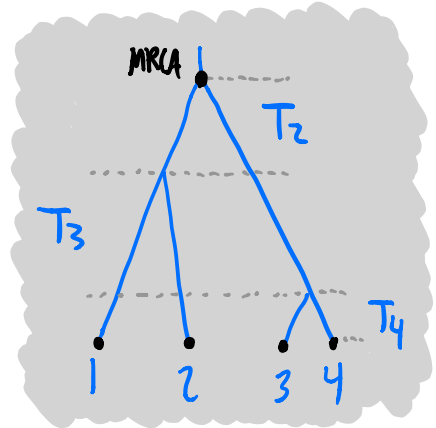


Easily generalizes to sample of size n:

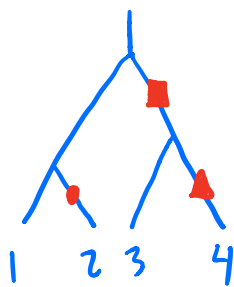
① @ each step, only consider coalescence between pairs of lineages

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

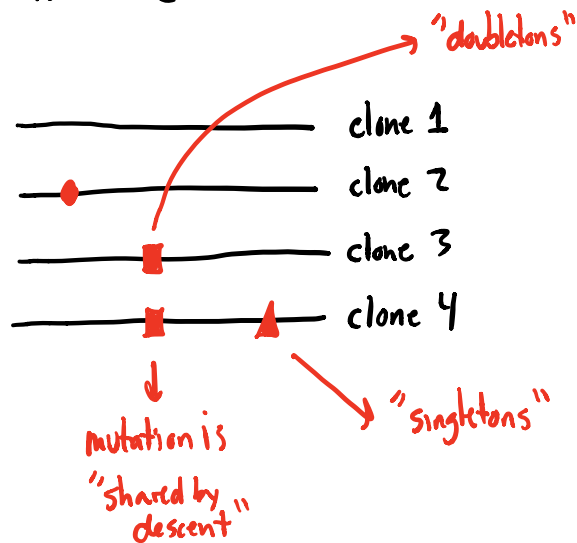
③ choose random pair to coalesce repeat!



④ then can paint mutations on @ end:



\Rightarrow



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

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

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

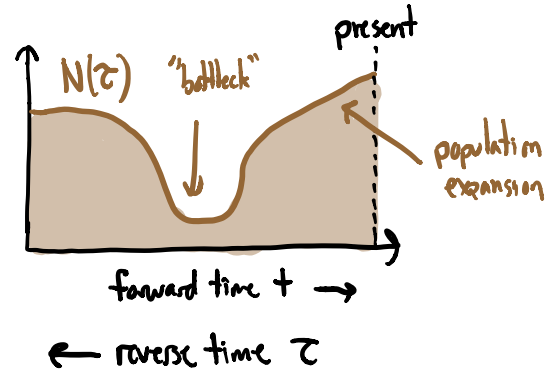
\Rightarrow compare to single-locus prediction (easy!)

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

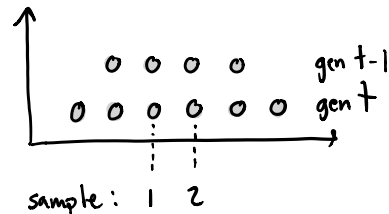
\Rightarrow why use coalescent picture then??

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 $\rightarrow \frac{1}{N(\tau)}$

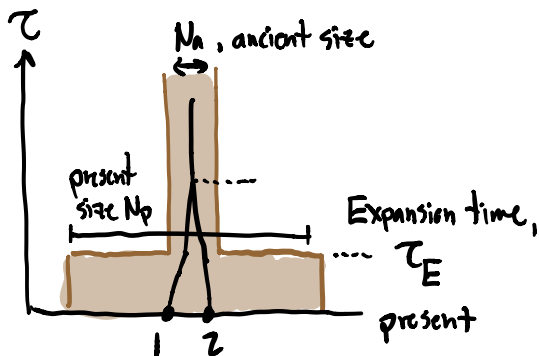


\Rightarrow coalescence = "inhomogeneous" Poisson process:

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

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

Simple example: rapid expansion in recent past



① if $N_p \gg \infty$ ($\tau_E \ll N_p$)

\Rightarrow no coalescence until τ_E

\Rightarrow coalescence @ rate $\frac{1}{N_a}$ after

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

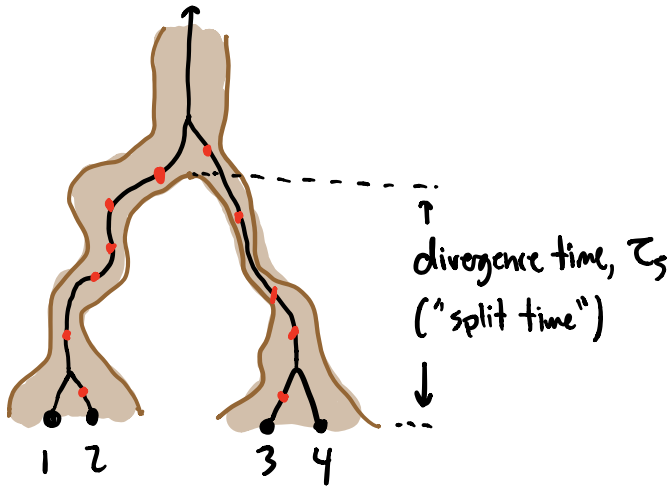
$\Rightarrow \langle \pi \rangle = 2N \langle T_2 \rangle = 2N(\tau_E + N_a) \approx 2NN_a$
(if $\tau_E \ll N_a$)

if $N_p \mu \sim 100 \Rightarrow$ why $\pi \sim 1e^{-03}$ in humans??

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

compare to $\frac{df}{dt} = \mu(1-f) - 2f + \sqrt{\frac{f(1-f)}{N(t)}} \eta(t)$

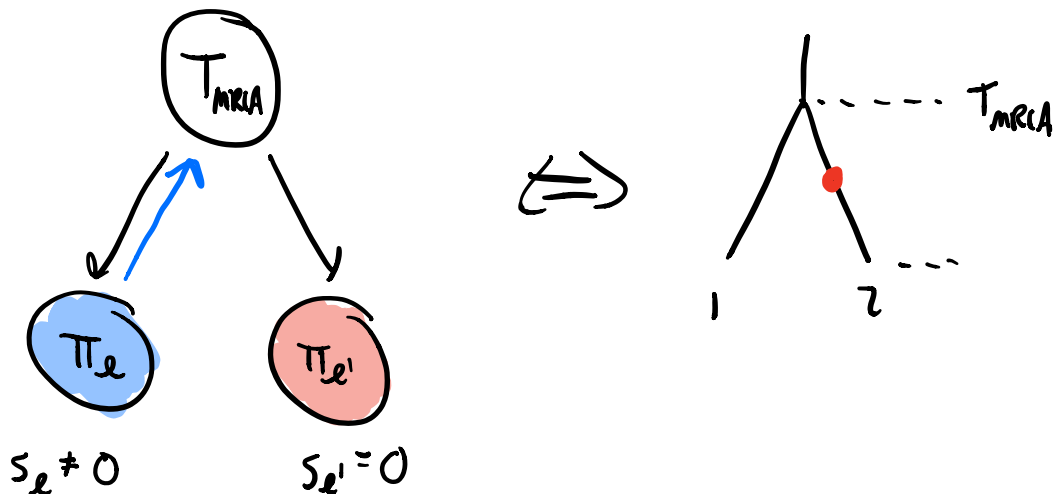
can also easily add population structure



\Rightarrow prob of coalescence
between pop's = 0
until time $\tau = \tau_s$

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

\Rightarrow downside: hard to add selection back in to picture...

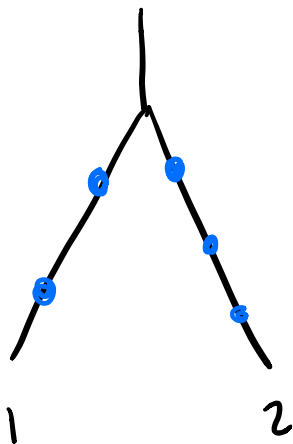


⇒ when is this going to be an issue?

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

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

⇒ suggests: $N|X(\vec{s}) - \bar{X}| \ll 1$ for neutrality



① assume effective neutrality:

⇒ total # mutations $\approx NU$

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

⇒ self consistent:

$$\boxed{(NU)(Ns)^2 \ll 1}$$

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

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



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