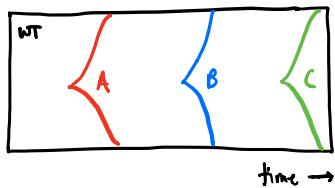


Chapter II

Neutral theory and the coalescent

Neutral theory & the Coalescent

Successive mutations:



$$\frac{dS(\vec{g})}{dt} = \sim(x-\bar{x}) + \cancel{\sim L\cdot\mu} \rightarrow \epsilon \\ + \sim\epsilon + \sim\frac{\pi}{JN}$$

⇒ ~1 variant present @ high freqs ⇒ solved by reducing to $L=1$ model

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

⇒ need to understand what's going on in these cases...

$$\frac{dS(\vec{g})}{dt} = \cancel{\sim(x-\bar{x})} + \sim L\cdot\mu \rightarrow 0 \\ + \sim\epsilon + \sim\frac{\pi}{JN} \rightarrow 0$$

⇒ one other limit that's well understood:
neutral evolution in nonrecombining genome

when $X(\vec{g})=0 \wedge \rho=0$, left with: $(N_e = \nu_e)$

$$\frac{\partial f(\vec{g})}{\partial t} = \sum_e \sum_{|\vec{g}' - \vec{g}|=1} N_e f(\vec{g}') \left[g_e(1-g_e) + (1-g_e)g_e' \right] - \sum_e N_e f(\vec{g})$$

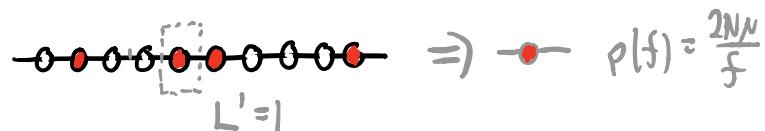
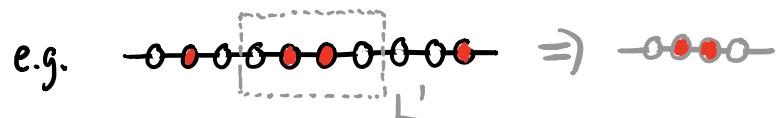
incoming mutations

outgoing mutations

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

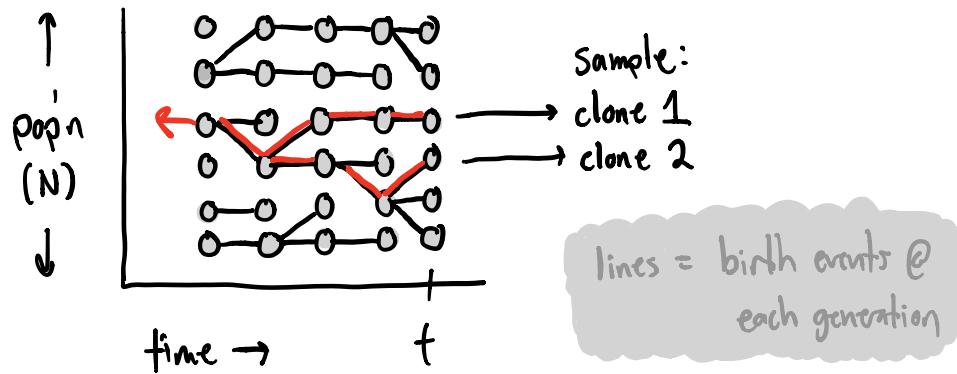
genetic drift

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

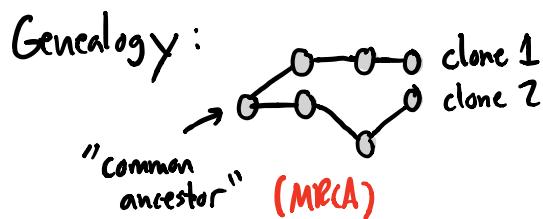


\Rightarrow 2nd key insight: can take $L'=0$ —

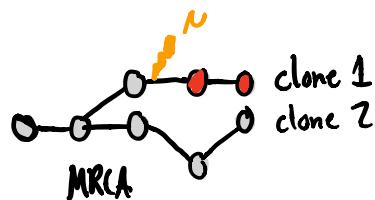
E.g. simulation of neutral pop'n in Wright-Fisher model:



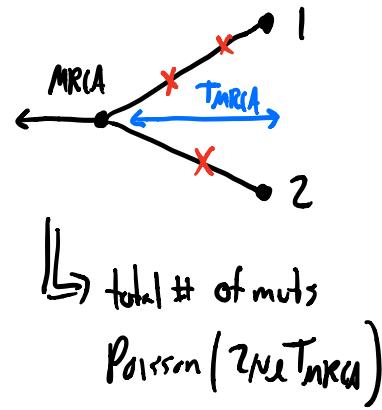
⇒ key insight: lines also = **genealogical relationships**
backward in time!



↓ differences between sampled individuals
 = mutations on genealogy



\Rightarrow Mutations occur as Poisson Process
w/ rate μ_e on each branch

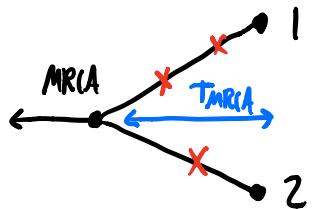


\Rightarrow 2 extreme limits:

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

$$\Rightarrow \Pr[\text{genetic diff} @ \text{site } e] =$$

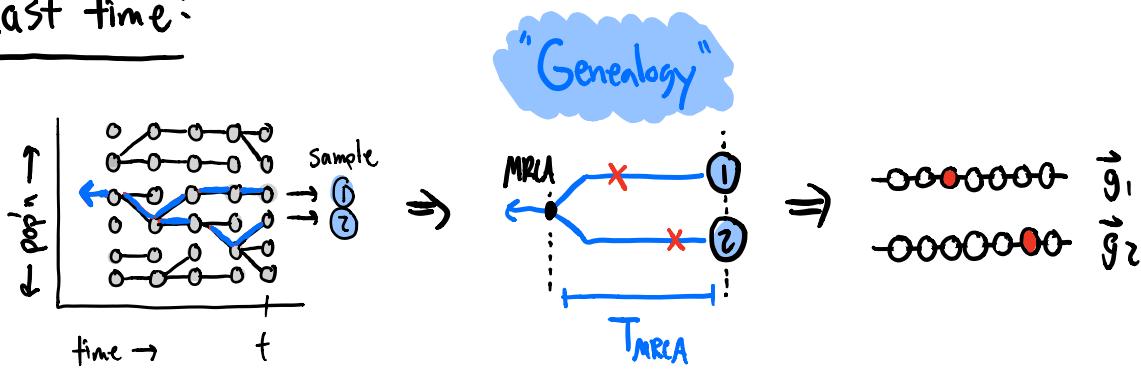
(2) $\mu_e T_{\text{MRCAs}} \gg 1 \Rightarrow$ lots of forward & backward mutations along each branch.



$$\Rightarrow \Pr[\text{genetic diff} @ \text{site } L] =$$

\Rightarrow Key Question: what sets T_{MRCAs} ?

Last time:



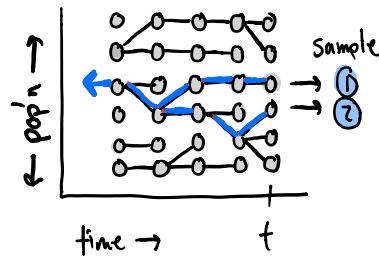
⇒ Given genealogy (T_{MRCA}), mutations occur as Poisson Process along each branch ("mutation painting")

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

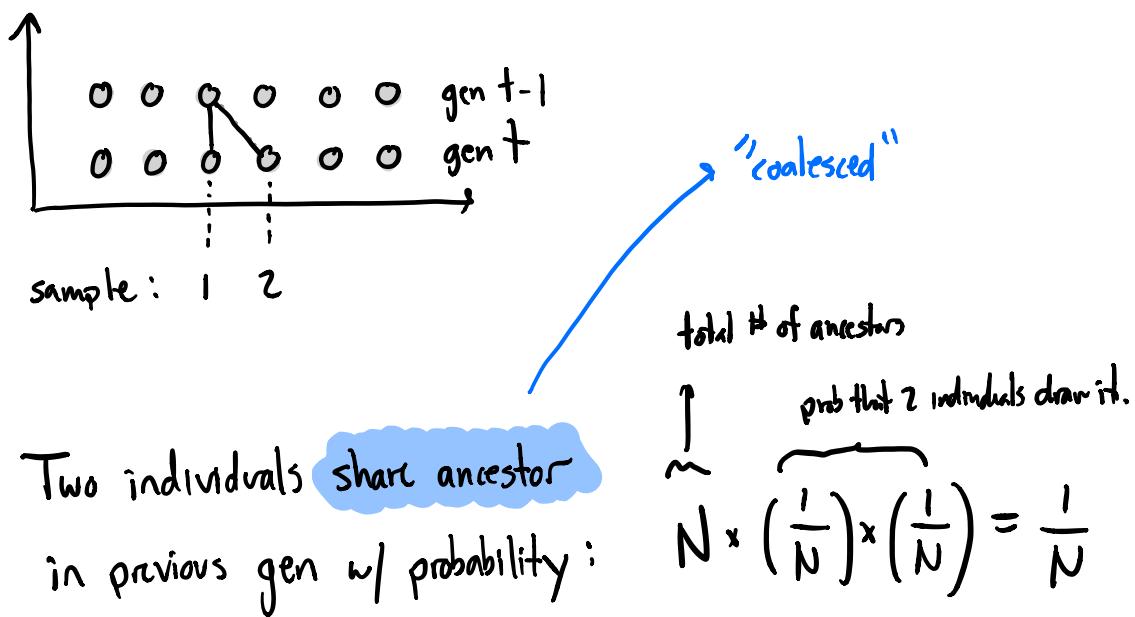
Today: what determines genealogy (T_{MRCA})?

⇒ Note: T_{MRCA} is random quantity

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



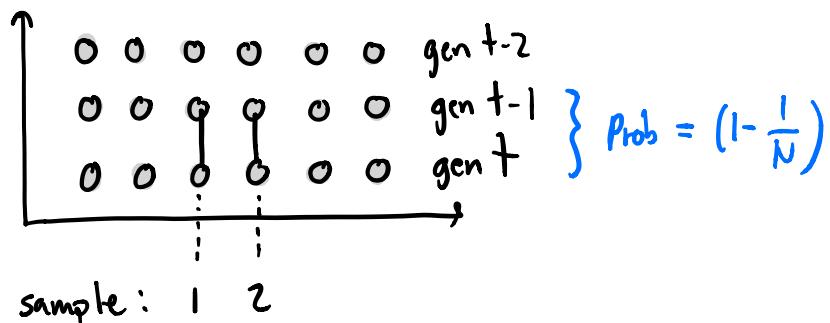
\Rightarrow key insight: start from present & work backward in time:



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

\Rightarrow otherwise, diff ancestors in gen $t-1 \Rightarrow$ repeat!

Process repeats itself w/ next gen:

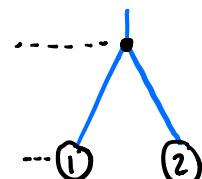


$$\Rightarrow \text{w/ prob } \frac{1}{N} \left(1 - \frac{1}{N}\right) \Rightarrow T_{\text{MRCA}} = 2$$

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

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

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



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

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

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

$$\approx \int_{(\mu T_{\text{MRCA}})} 2N_e T_{\text{MRCA}} p(T_{\text{MRCA}}) dT_{\text{MRCA}} = 2N_e \langle T_{\text{MRCA}} \rangle = 2N_e$$

\Rightarrow matches our previous result for $\langle \pi \rangle$, ✓.

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

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

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

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

$$\begin{aligned}
&= \int (2N_e T_{\text{match}}) \cdot (2N_{e'} T_{\text{match}}) \cdot \rho(T_{\text{match}}) \cdot dT_{\text{match}} \\
&= (2N_e) \cdot (2N_{e'}) \cdot \langle T_{\text{match}}^2 \rangle = (2N_e) \cdot (2N_{e'}) \cdot (2N^2) \\
&= 2 \cdot (2N_e N) \cdot (2N_{e'} N) \\
&= 2 \cdot \Pr(\pi_e = 1) \cdot \Pr(\pi_{e'} = 1) \geq \Pr(\pi_e = 1) \Pr(\pi_{e'} = 1)
\end{aligned}$$

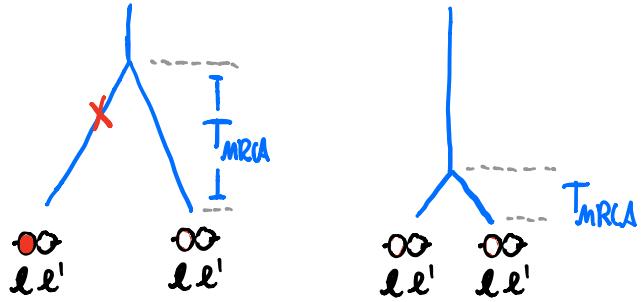
\Rightarrow joint prob of mutations is not independent

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

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

⇒ what's going on?

⇒ consider 2 trees:

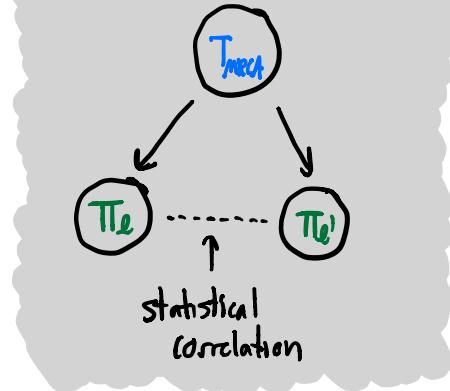


⇒ conditioned on $\pi_{le} = 1$, likely had bigger-than-avg T_{MRCA}

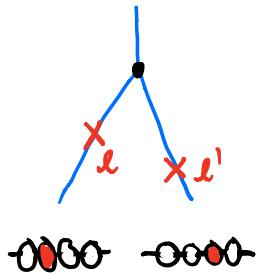
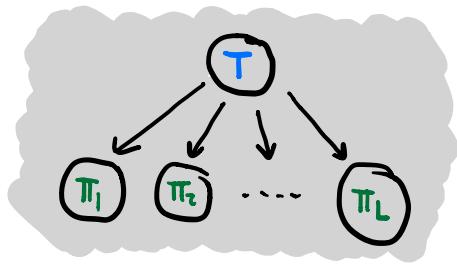
⇒ i.e. mutations don't interact, but are still coupled

by shared genealogy

Causation diagram



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



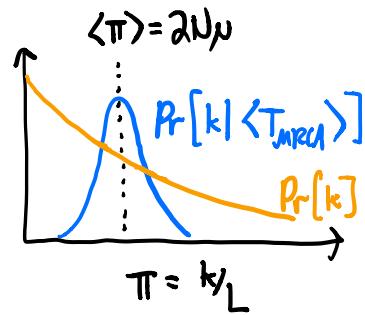
\Rightarrow when $\mu e T_{\text{MRCA}} \ll 1$, most mutations
will occur @ unique site in genome
"infinite-sites approximation"

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

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

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

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



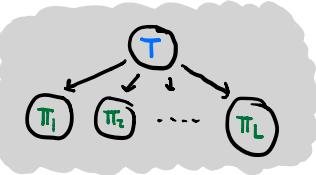
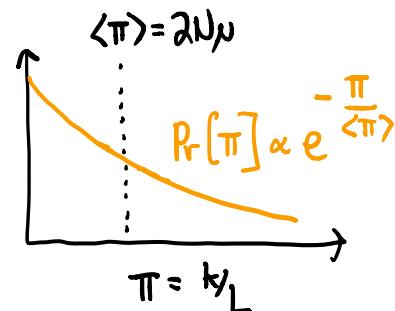
\Rightarrow one advantage of coalescent approach :

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

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

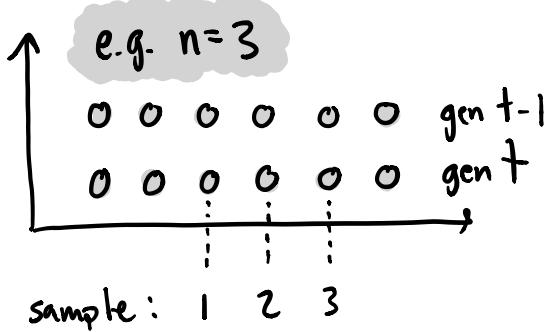
$$\Rightarrow \text{or } C_V^2 \equiv \frac{\text{Var}(\pi)}{\langle \pi \rangle^2} = \frac{1+2NU}{2NU} \geq 1$$

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



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

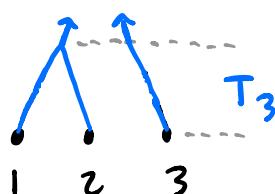
Larger Sample Sizes ($n > 2$)



\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 = $N \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) \cdot \left(\frac{1}{N}\right) = \frac{1}{N^3}$

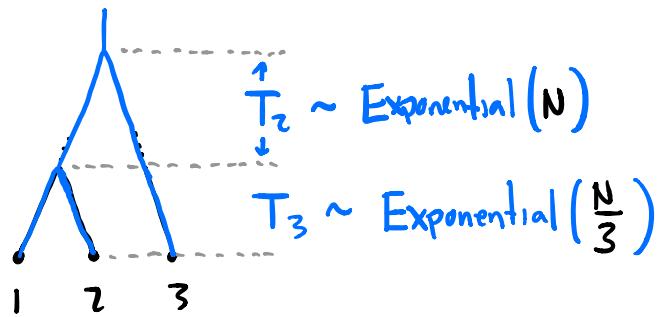
\Rightarrow when $N \gg 1 \rightarrow$ only need to worry about **pairwise coalescence**
 (known as "Kingman's coalescent") **(all pairs are equally likely to coalesce)**



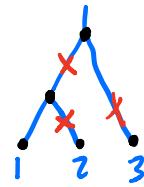
\Rightarrow total prob of coalescence = $\frac{3}{N}$ per gen

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

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

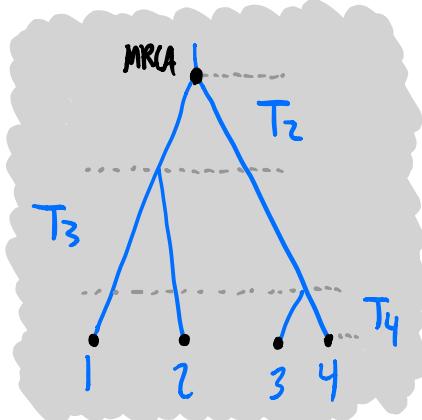


\Rightarrow Done! can now paint on mutations...

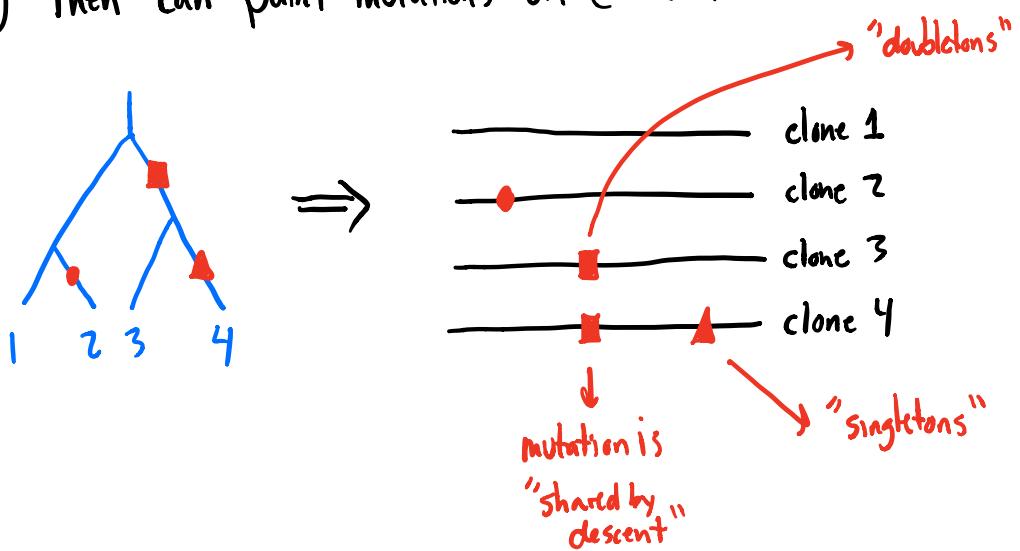


Easily generalizes to sample of size n:

- ① @ each step, only consider coalescence between pairs of lineages ↗
- ② Time until next coalescence event is $T_n \sim \text{Exponential}(N/(n))$
- 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. $\langle \# \text{ doubletons in sample } n=4 \rangle = \left\langle \begin{array}{c} \text{Diagram 1} \\ + \\ \text{Diagram 2} \end{array} \right\rangle$

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

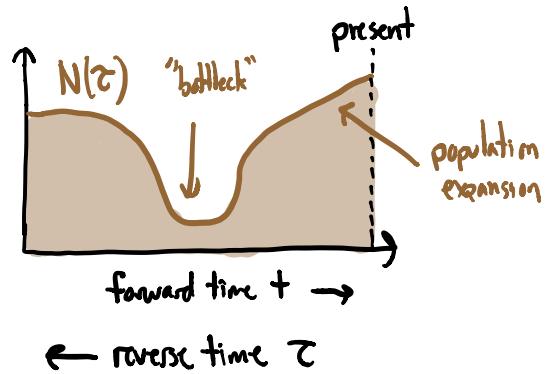
\Rightarrow compare to single-locus prediction (easy!)

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

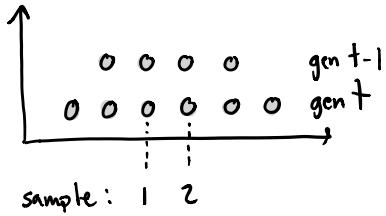
\Rightarrow why use coalescent picture then??

Answer: coalescent picture makes it easy to model demography!

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



⇒ coalescent picture still works, but coalescent prob $\rightarrow Y_{N(\tau)}$

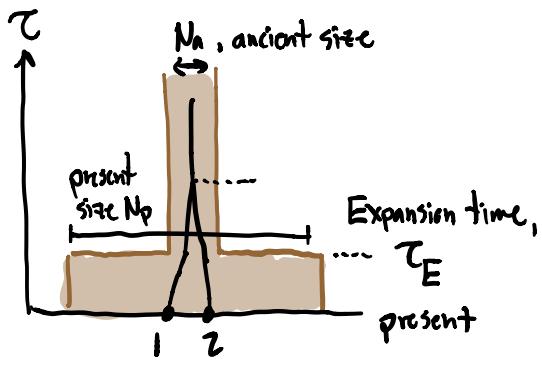


⇒ coalescence = "inhomogeneous" Poisson process:

$$\Rightarrow \Pr[T_2 > \tau] = \prod_{i=2}^2 \left[1 - \frac{1}{N(\tau)} \right] \approx e^{- \int_0^\tau \frac{dp}{N(\tau')} d\tau'}$$

$$\Rightarrow \Pr[T_2 = \tau] = \frac{1}{(2N)} e^{- \int_0^\tau \frac{dp}{N(\tau')} d\tau'}$$

Simple example: rapid expansion in recent past



$$\textcircled{1} \text{ if } N_p \gg \infty \quad (\tau_E \ll N_p)$$

\Rightarrow no coalescence until τ_E

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

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

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

if $N_p N \sim 100 \Rightarrow \pi \sim e^{-0.3}$ in humans ??

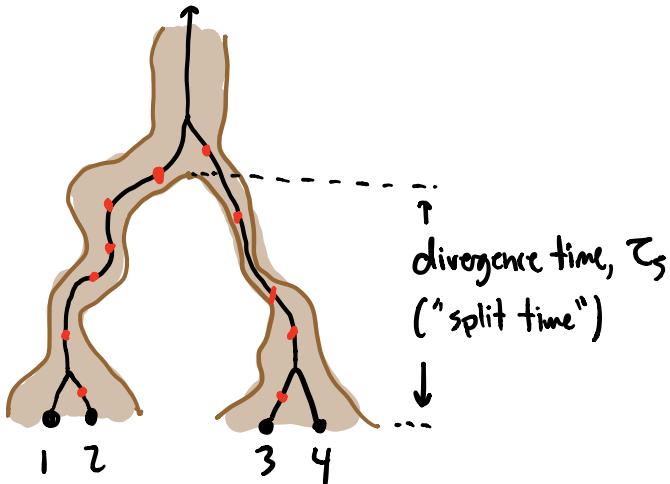
\Rightarrow answer: $N(t)$ was smaller backward in time.



Compare to

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

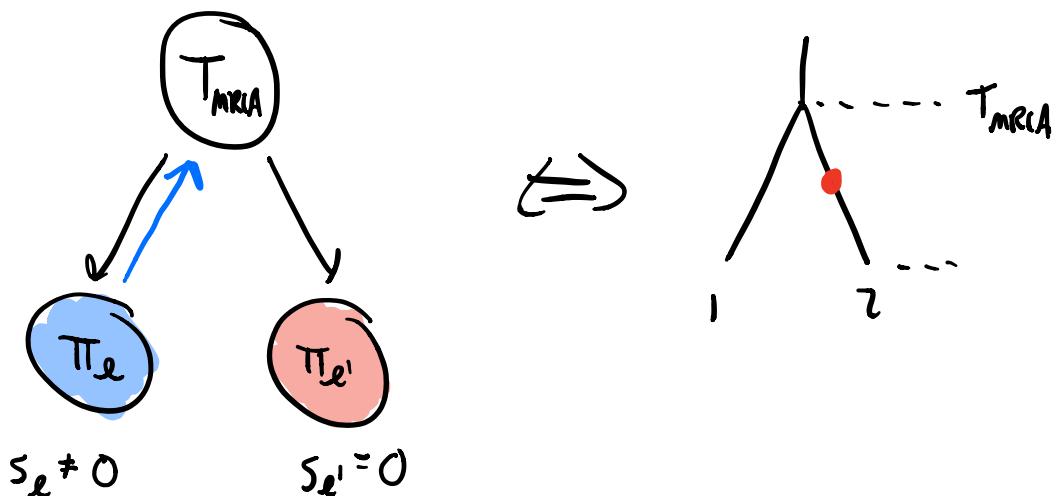
can also easily add population structure



\Rightarrow prob of coalescence
between pop'n's = 0
until time $T=T_s$

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

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

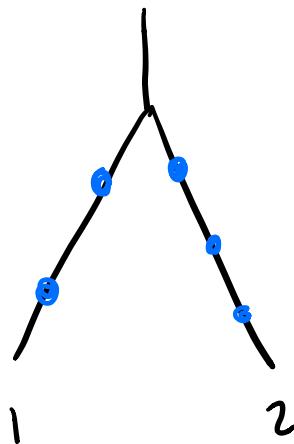


\Rightarrow when is this going to be an issue?

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

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

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



① assume effective neutrality:

\Rightarrow total # mutations $\approx N\bar{U}$

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

\Rightarrow self consistent:

$$(N\bar{U})(Ns)^2 \ll 1$$

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

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

↓

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