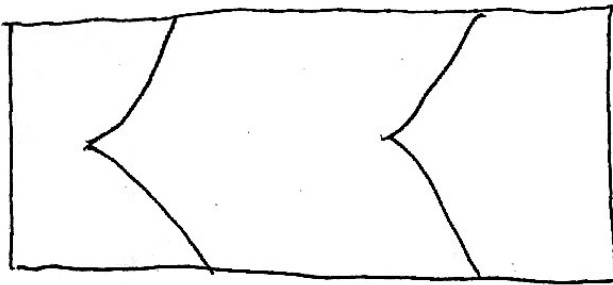


# Neutral Theory and the Coalescent

①

In the last lecture, we discussed the behavior of multi-locus models of evolution in the successive mutations regime



where only  $\sim 1$  genetic variant is present @ high frequencies @ any given time.

In this case, behavior was exactly solvable because it reduced to an effective single locus model.

In general, genomes in data are separated by multiple mutations (e.g. humans,  $\sim 1$  mutation every 1000bp between 2 individuals)

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

~~today~~ today, we'll focus on one limit that is very well understood: neutral evolution on ~~an~~ a nonrecombining genome.

it might seem unrealistic (and it might be unrealistic, @ least for bacteria), but the neutral model is extremely influential in how people have come to think about + talk about genetic data, so it's worth being aware of it.

(2)

one motivation for considering this limit is that selection + recombination contribute nonlinear terms to the full multi-locus SDE we wrote down last class. Setting  $X(\vec{g}) = \text{const}$  +  $\rho = 0$  ~~also~~ leads to a linear system, and we'd expect that linear systems should be possible to analyze (@ least in principle). At first glance, linear system is still pretty complicated:

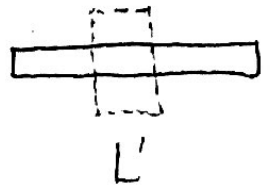
$$\frac{df(\vec{g})}{dt} = \left[ \sum_{\substack{\vec{g}' \text{ s.t.} \\ |\vec{g}' - \vec{g}| = 1}} \sum_{e=1}^L \mu_e f(\vec{g}') [g_e(1-g_e') + (1-g_e)g_e'] - f(\vec{g}) \sum_{e=1}^L \mu_e \right]$$

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

insight comes from realizing that the sites don't actually influence each other (because they are neutral).

3

$\Rightarrow$  i.e., we could focus on a subset of the genome  $L'$ , and write down a corresponding neutral model for the genotypes @ just these loci, and it would have the same form.



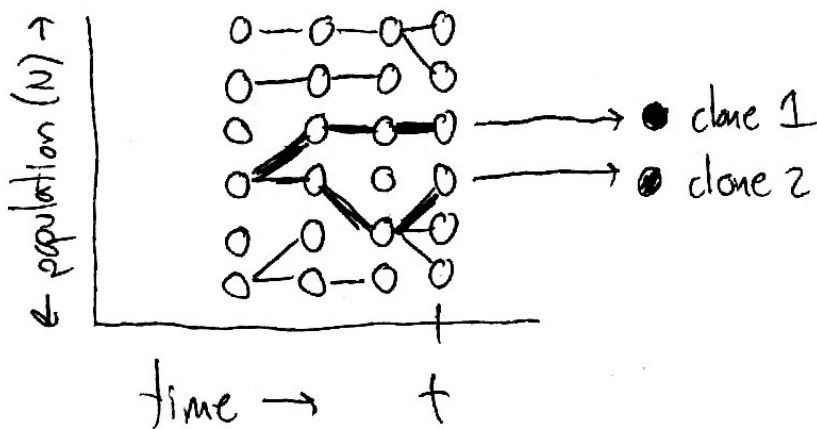
$\Rightarrow$  we saw an extreme form of this a few lectures ago where we took  $L'=1$ , and saw that the behavior of a single site in a long neutral genome is indistinguishable from a single-locus model.

$\Rightarrow$  thus, in a similar way as in the SSWM limit, we might be able to understand what's going on due to some effective reduction to a single locus behavior that we already understand

$\Rightarrow$  in this case, the true insight comes from going one step further and considering a zero locus model, ~~that~~ (that is, taking out mutations entirely.)

to explain what I mean, suppose we simulate a neutral population in a Wright fisher model, and we sample 2 random individuals from the population @ the end.

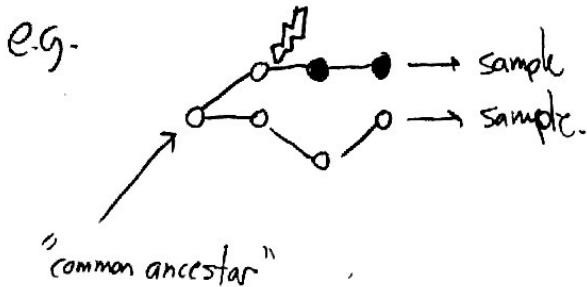
If we examined the guts of the simulation, it would look like:



where lines illustrate birth events that form next generation.

going backward in time, lines also illustrate genealogical relationships between the sampled (or unsampled individuals).

In this case, ~~the~~ the 2 clones would have a genotic difference @ a ~~particular~~ particular site if one of the divisions ~~along the genealogy~~ along the genealogy resulted in a mutation, before the common ancestor

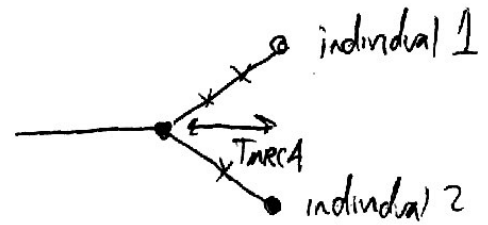


\* Since mutations are neutral, they can have no influence on this hypothetical genealogy itself.

⇒ we are free to "paint" them on after the fact.

5

E.g. if 2 individuals shared a common ancestor  $T_{MRCA}$  generations ago, then mutations occur as a Poisson process @ rate  $\mu_e$  on each branch (length  $T_{MRCA}$ )



⇒ 2 extreme limits:

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

$$\Pr[\text{genetic difference @ site } e] = 1 - (1 - \mu_e)^{2T_{MRCA}} \approx 2\mu_e T_{MRCA}$$

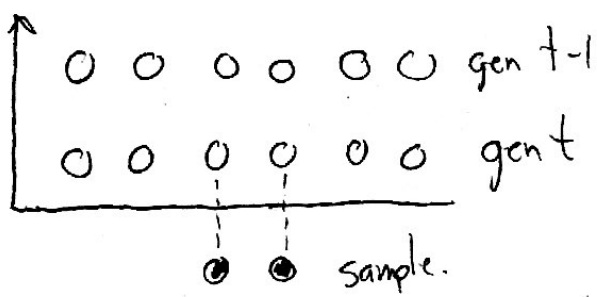
(2)  $\mu_e T_{MRCA} \Rightarrow$  lots of forward & reverse mutations along each branch of tree.

$$\Pr[\text{genetic difference @ site } e] = \frac{1}{2}$$

⇒ seems pretty straight forward. question is what sets  $T_{MRCA}$ ?

In principle, this is a random quantity, since genealogies in W.F. simulation will vary if you re-run the simulation.

⇒ Suppose we start from present day population & work back in time:



all ~~individuals~~ individuals in gen t must have some ancestor in previous generation. ⇒ in WF model, these ancestors are chosen uniformly @ random from previous generation (w/ replacement.)

⇒ probability that two individuals share a common ancestor is  $\frac{1}{N}$  ⇒ in this case, we say they have "coalesced"

⇒ w/ probability  $\frac{1}{N}$ ,  $T_{MRCA} = 1$  ← also known as "coalescence time"

⇒ w/ probability  $(1 - \frac{1}{N})$ , individuals descend from different parents in previous generation.

then process repeats itself:

(7)

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

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

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

$$\text{i.e., } \langle T_{\text{MRCA}} \rangle = N, \sqrt{\text{Var}(T_{\text{MRCA}})} = N$$

$\Rightarrow$  the total probability of observing a mutation @ site  $l$  can then be obtained by integrating over  $T_{\text{MRCA}}$ :

$$\text{Pr} \left[ \begin{array}{c} \text{difference} \\ \text{@ site } l \end{array} \right] = \int \text{Pr}(\text{diff} | T_{\text{MRCA}}) p(T_{\text{MRCA}}) dt_{\text{MRCA}}$$

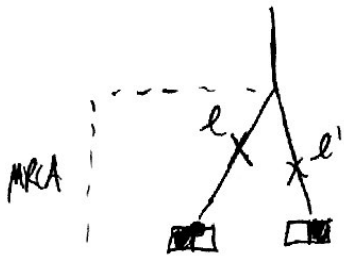
$$(\text{when } N \gg 1) \approx \int 2\mu T_{\text{MRCA}} p(T_{\text{MRCA}}) dt_{\text{MRCA}} = 2\mu \langle T_{\text{MRCA}} \rangle = 2\mu N$$

$\hookrightarrow$  this matches w/ what we derived for  $\langle \pi \rangle$  before, since

$$\langle \pi \rangle \equiv \text{Pr} \left( \begin{array}{c} \text{difference} \\ \text{@ site } l \end{array} \right)$$

distribution of  $T_{MRCR}$  becomes important when we consider mutations @ multiple sites. e.g.

$$Pr \left[ \begin{array}{c} \text{diff @ site} \\ \underline{l} \text{ and } \underline{l'} \end{array} \right] = \int Pr[\pi_{e=l}, \pi_{e'=l'} | T_{MRCR}] P(T_{MRCR}) dT_{MRCR}$$



$$= \int \underbrace{Pr[\pi_{e=l} | T_{MRCR}] Pr[\pi_{e'=l'} | T_{MRCR}]}_{\text{mutations are neutral, so can't affect each other.}} P(T_{MRCR}) dT_{MRCR}$$

$$= \int (2\mu_e T_{MRCR}) (2\mu_{e'} T_{MRCR}) P(T_{MRCR}) dT_{MRCR}$$

$$= 2 \cdot (2\mu_e N) \cdot (2\mu_{e'} N) \quad \left[ \text{since } \langle T^2 \rangle = 2N \right]$$

$$= 2 \cdot Pr[\pi_{e=l}] Pr[\pi_{e'=l'}] \geq Pr[\pi_{e=l}] Pr[\pi_{e'=l'}]$$

joint  
=> probability of having a mutation at both sites is not independent.

$$\Rightarrow Cov(\pi_{e=l}, \pi_{e'=l'}) = \frac{\langle \pi_{e=l} \pi_{e'=l'} \rangle - \langle \pi_{e=l} \rangle \langle \pi_{e'=l'} \rangle}{\langle \pi_{e=l} \rangle \langle \pi_{e'=l'} \rangle} = 1$$

twice as likely!

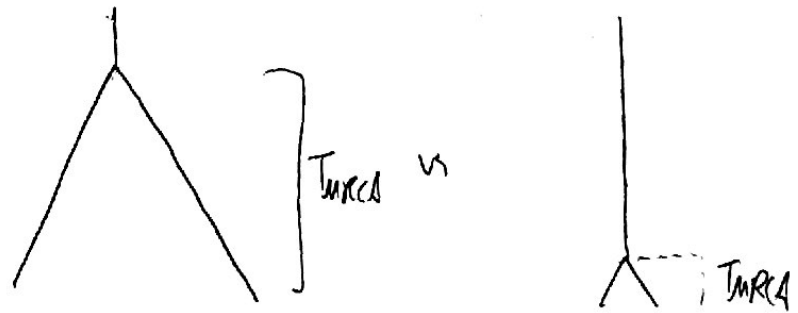
=> knowing that you have a mutation @ site  $l$  makes it more likely to observe a mutation @ site  $l'$ !



But previously said that mutations don't influence other directly...

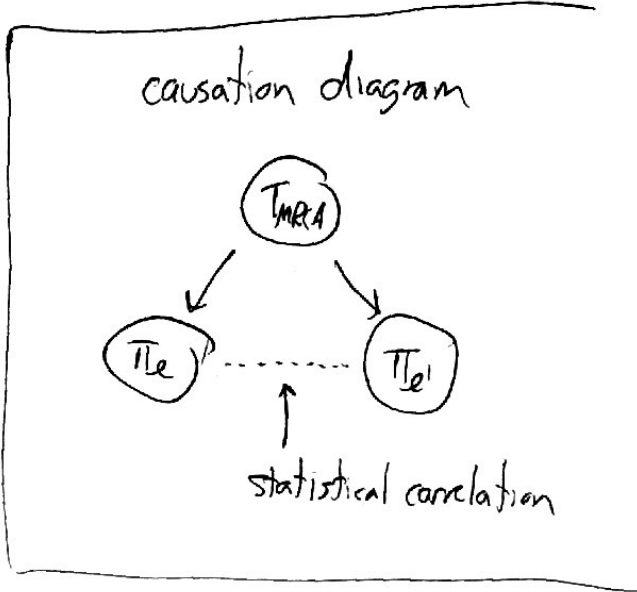
⇒ what's going on?

⇒ mutations more likely to occur when  $T_{MRCA}$  is bigger



⇒ conditioned on  $\pi_{l'}=1$ , probably had a bigger than avg  $T_{MRCA}$  so more likely to have mutation @ site  $l'$  too.

i.e., mutation processes do not interact, but are still coupled together through random genealogy that happens to be present in a given population



⇒ can keep adding more sites in this way. In limit that  $N_e T_{MRCA} \ll 1$ , then most mutations will occur @ a unique site in the genome. ("infinite sites" assumption), then conditioned on  $T_{MRCA}$ ,

~~the total # of mutations is a Poisson process with rate U~~

~~with rate U~~ the total # of mutations occurs as a Poisson process w/ rate  $U \equiv \sum_{l=1}^L N_e$  per gen.

i.e., if  $k$  = total # of mutational differences between

the two individuals, then  $Pr[k | T_{MRCA}] = \frac{(2\mu T_{MRCA})^k}{k!} e^{-2\mu T_{MRCA}}$

$$Pr[k] = \int Pr[k | T_{MRCA}] P(T_{MRCA}) dT_{MRCA} = \int \frac{(2\mu T_{MRCA})^k}{k!} \frac{1}{N} e^{-(2\mu + \frac{1}{N})T_{MRCA}} dT_{MRCA}$$

$$= \frac{(2\mu N)^k}{(2\mu N + 1)^{k+1}} \Rightarrow \text{geometric distribution w/ prob } \frac{2\mu N}{1 + 2\mu N}$$

(again, broader than  $Pr[k | T_{MRCA}]$ )

$\Rightarrow$  so one advantage of coalescent approach is that it provides simple predictions for ~~variance~~ ~~uncertainty~~ ~~in~~  $\pi$  rather than just  $\langle \pi \rangle$ .

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

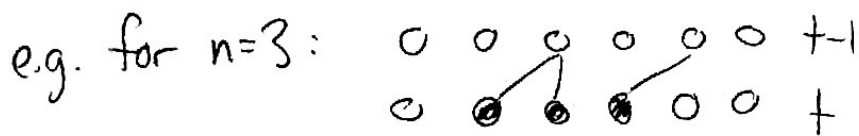
$$\text{or } CV \equiv \frac{Var(\pi)}{\langle \pi \rangle^2} = \frac{(1 + 2\mu N)}{2\mu N} \geq 1$$

i.e.,  $\pi$  does not self average on a long a sexual genome.

$\Rightarrow$  fluctuations in  $T_{MRCA}$  influence large #s of sites simultaneously.

can also consider sample sizes larger than 2.

(11)



prob that any 2 share a common ancestor in previous generation is  $\frac{1}{N}$   
 $\Rightarrow \binom{3}{2} = 3$  total pairs.

$\Rightarrow$  probability that all 3 share an ancestor in previous generation is

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

$\swarrow$  # of ancestors possible       $\searrow$  prob that each one chooses that ancestor.

(when  $N \gg 1$ )

$\Rightarrow$  only have to worry about pairwise coalescence (all pairs equally likely to coalesce)  
 (known as "kingman coalescent")

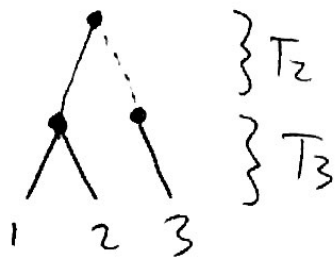
~~total # of pairs is 3~~ total # of pairs is 3, so total probability of getting a coalescence b/w a pair is  $\frac{3}{N}$ .

$\Rightarrow$  time until this happens is  $T_3 \equiv \text{Exponential}\left(\frac{N}{3}\right)$

$\Rightarrow$  once this happens, choose random pair to coalesce. then have effective sample of  $n=2$ .

$\Rightarrow$  time till coalescence of these is

$T_2 \equiv \text{Exponential}(N)$  as before



at this point, all share common ancestor, so done!

⇒ can then paint mutations on as before

⇒ easily generalizes to sample of size n:

@ each step, need only consider coalescence btw pairs of individuals ⇒ time until next coalescence event =  $T_n = \text{Exponential}\left(\frac{N}{\binom{n}{2}}\right)$

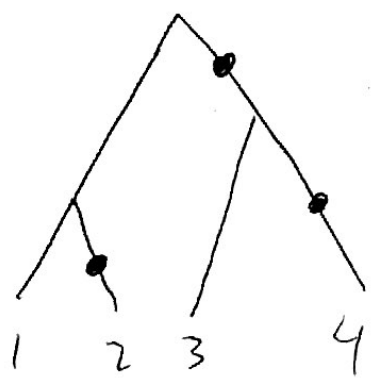
⇒ once you draw  $T_n$ , choose a random pair to coalesce. then repeat w/ sample of size  $n-1$ . →  $n=2$ .



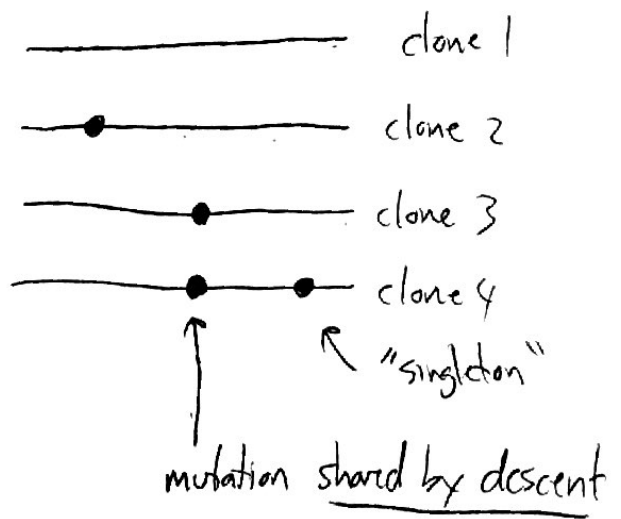
⇒ then can paint mutations on @ end.

⇒ mutations that occur higher up in tree are present in all individuals that inherit from that branch.

e.g.



⇔



(only one mutation event)

easy to simulate this process, for any  $n$  ("coalescent simulations") but much harder to calculate anything analytically for ~~any~~  $n \geq 4$ .

e.g.  $\langle \# \text{ of doubletons in sample of size } 4 \rangle = \langle \text{all trees that look like this} \rangle$

i.e. avg over mutations occurring + branch lengths + topologies

$\Rightarrow$  gets very hard for  $n \geq 4$ , even when  $n \rightarrow \infty$ .

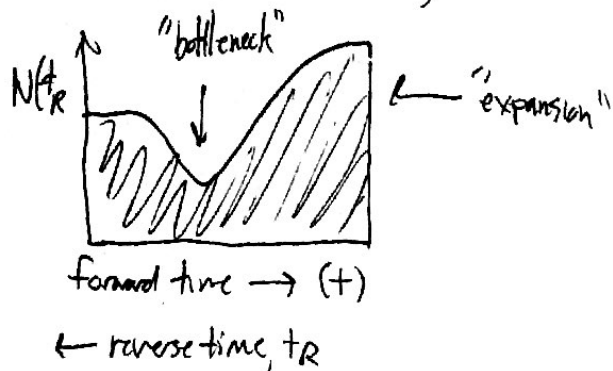
$\Rightarrow$  compare to single locus prediction:

$$\langle \# \text{ doubletons in sample of size } 4 \rangle = \int \text{[scribbled out]} \cdot \binom{4}{2} f^2 (1-f)^{4-2} \cdot \frac{2N\mu}{f} df$$

$$= \frac{4N\mu}{5} \text{ (easy!)}$$

why would we ever use coalescent picture then?

⇒ answer is that coalescent picture makes it very easy to model demography. e.g. if  $N$  was not constant, but ~~fixed~~ varied historically back in time,  $N(t)$



then coalescence picture still works, except that now coalescence probability changes in each generation,  $\frac{1}{N(t)}$

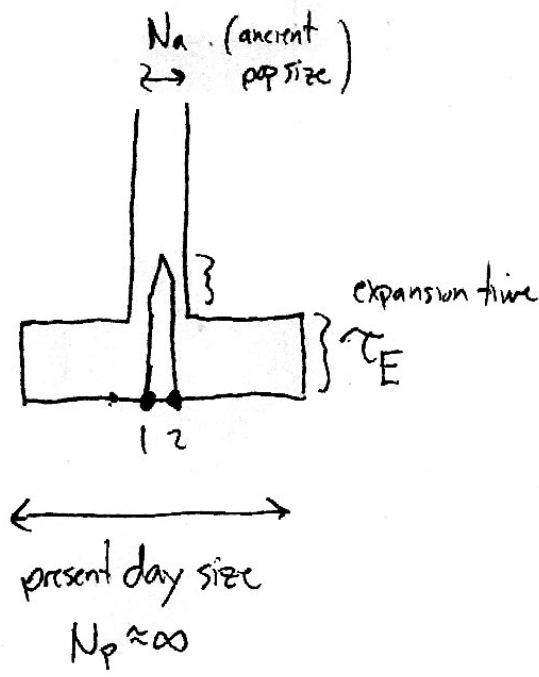
⇒ time to coalescence is inhomogeneous poisson process:

e.g. ~~scribble~~ 
$$p(T_{n,2} = t) = \frac{\binom{n}{2}}{N(t)} e^{-\int_0^t \frac{\binom{n}{2}}{N(t')} dt'}$$

⇒ if  $N$  is larger ⇒ less likely to coalesce.  
if  $N$  smaller ⇒ more likely to coalesce.

⇒ otherwise, everything is the same (same topologies, same mutation painting)

e.g. simple case: rapid expansion from  $N=N_a$  to ~~scribble~~  $T_e$  generations ago.  
 $N=N_p \approx \infty$



$\Rightarrow$  no coalescence for first  $\tau_E$  generations (assuming  $\tau_E \ll N_p$ )

$\Rightarrow$  then regular coalescence @ rate  $\frac{1}{N_a}$  afterwards.

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

$\Rightarrow \langle \Pi \rangle = 2\mu \langle T_2 \rangle = 2\mu (\tau_E + N_a) \approx 2\mu N_a$   
 (if  $\tau_E \ll N_a$ )

Hence, in this case  $\langle \Pi \rangle$  is dominated by historical (ancestral) population size  $N_a$ , and has nothing to do w/ population size now.

*sometimes people will call this effective pop size ( $N_e$ ).  $\rightarrow$  don't!*

$\Rightarrow$  addresses key paradox we had w/ human data @ ~~beginning~~ beginning of class. if  $N_p \gg 1$  why  $\Pi \sim 10^{-3}$ ?

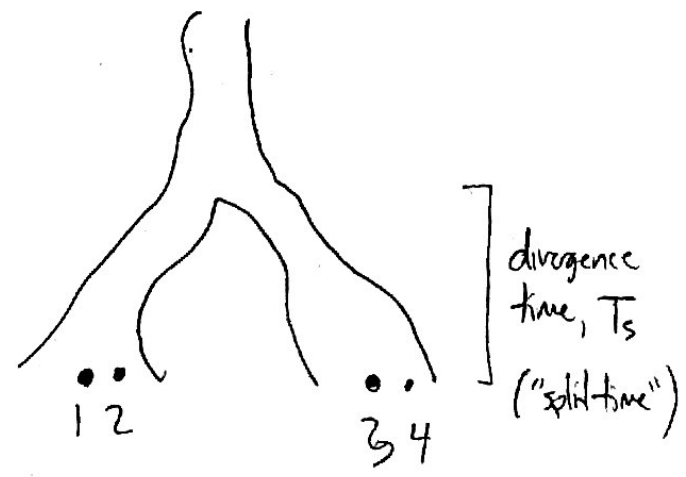
$\Rightarrow$  commonly accepted answer is that  $N(t)$  was smaller ~~in~~ back in time (e.g. out of africa).

$\Rightarrow$  genetic data let's us measure  $N(t)$ .

coalescent picture makes this much easier than corresponding forward time calculation:

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

⇒ can also easily add in population structure, e.g. isolated subpopulations:



- ⇒ prob coalescence between pop's x 0 until T<sub>S</sub>.
- ⇒ much of human pop gen is about inferring these population demographic models (much fancier of course)

~~\* when we use the neutral model as a good approximation~~

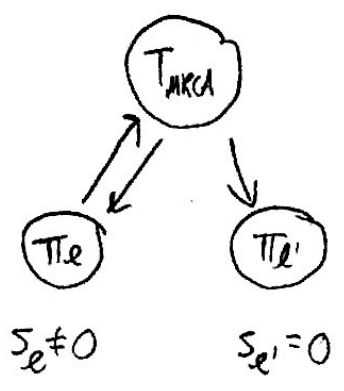
downside is that it is very hard to add selection back in to this picture.



⇒ basic problem is causation diagram gets reversed:

~~evolutionary model is painting backwards.~~





⇒ there still is a genealogy... so  
 can still paint on neutral mutations  
 (if you know T\_MRCA)

⇒ but can't paint on selected mutations  
 & worse... don't know what T\_MRCA dist'n  
 is any more!

⇒ this scenario is called "linked selection". will revisit  
 in a later lecture.

~~When is this a problem?~~ when is this a problem?

i.e. when is neutral limit a good approx?

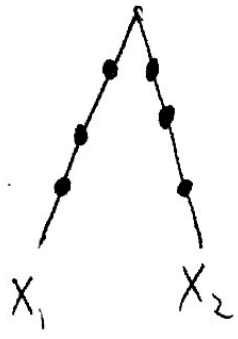
⇒ @ least need  $Ns_e \ll 1$  (as in single locus case)

⇒ but since looking @ multi-locus problem, need

$$|x(\vec{j}) - x(\vec{j}')}| \cdot N \ll 1$$

⇒ can we estimate these typical fitness differences  
 using a self consistency argument?

⇒ if neutral model is good approx, then  
 # mutations btw individuals is  $\sim 2N \langle T_2 \rangle L$



⇒ mutations occur on different branches equally  
 so tend to cancel each other out.

⇒ if each one has fitness effect  $\pm s$ .

then  $X_1 - X_2 \sim \pm \sqrt{\frac{2\mu \langle T_2 \rangle}{\pi} \cdot L \cdot s^2}$  (from C.L.T.)

so neutral model is good approx if  $(NU)(Ns)^2 \ll 1$

⇒ if  $NU \gg 1$ , can be bad even if  $Ns$  is very small!

e.g. if  $Ns = 0.1$  (hard to select on individually)

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

then genome wide  $\Delta X$ 's ~~are~~ are big!  
 ( $10^4 \times 10^{-2} = 10^2 \gg 1!$ )

⇒ in this case, we'll have to turn to alternative methods  
 to understand evolution in a long genome...