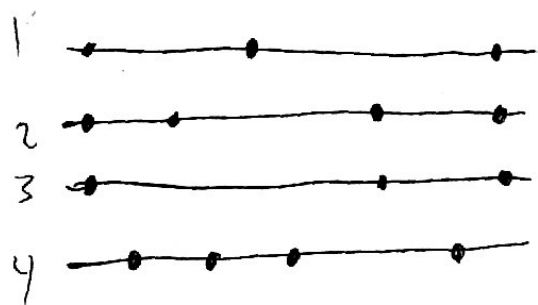
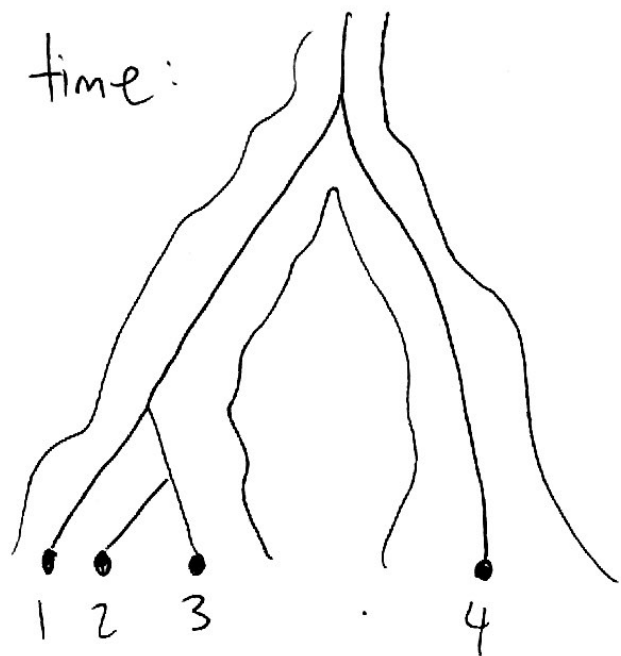


Coalescent theory II : selection & recombination

(1)

Last time:



mutations: poisson process w/ rate U .

we introduced a backward-in-time picture (coalescent theory) for describing sequences from a long neutral & non-recombining genome.

\Rightarrow powerful because it separated genealogy of sample from mutations that occur conditioned on genealogy.

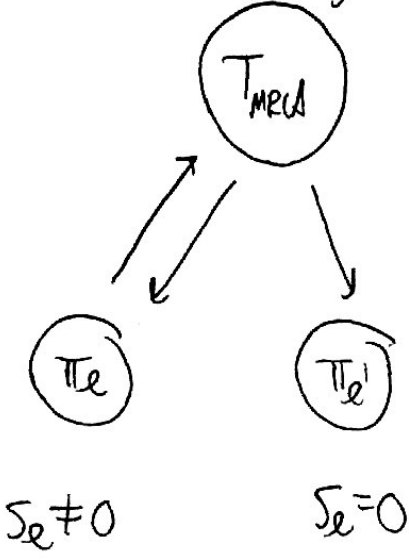
genealogy: Kingman coalescent
 $\rho_c \sim 1/N(t)$

\Rightarrow ~~downside~~ lots of use in pop. gen of higher organisms.

downside is that very hard to get selection back into picture.

\Rightarrow basic problem is that causation diagram gets reversed:

causation diagram



genealogy



* there still is a genealogy between individuals in sample, so can paint on neutral mutations if you know TMRCA.
 (i.e., neutral mutations \approx time)

\Rightarrow but can't paint on selected mutations in this way

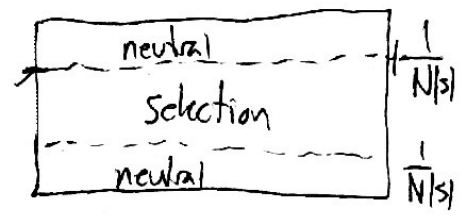
\Rightarrow even worse, no longer know ~~that~~ how to calculate TMRCA because it will depend on which mutations are there.

(catch-22)

\Rightarrow when is this likely to be a problem? let's try to estimate.

\Rightarrow remember for single site, ~~that~~

\Rightarrow needed $N|s| \ll 1$ for neutral.



for longer genome, selection term is now $[x(\vec{s}) - \bar{x}(t)] f(\vec{s})$

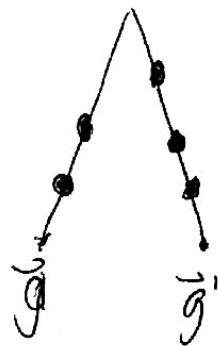
$$\hookrightarrow \bar{x}(t) = \sum_{\vec{s}} x(\vec{s}) f(\vec{s})$$

\Rightarrow suggests that need $|X(\vec{g}) - X(\vec{g}')| \ll \frac{1}{N}$ for most pairs of genomes (3)

can we estimate this difference w/ a self consistency argument?

\Rightarrow ~~if~~ if neutral limit is good approx,
then # mutations between individuals is

$$\sim \text{~~2U~~ } 2U \langle T_2 \rangle$$



but ~~mut~~ mutations will occur on branches equally, & tend to cancel out. ~~(~~ (Central limit theorem).

\Rightarrow if each mutation has effect $\pm s$, then

$$|X(\vec{g}) - X(\vec{g}')| \sim \sqrt{2U \langle T_2 \rangle s^2} \sim \sqrt{2NU s^2}$$

\Rightarrow neutral model is good approx if $\sqrt{2NU s^2} \ll \frac{1}{N}$

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

~~if~~ \Rightarrow if $NU \gg 1$, can be violated even if $Ns \ll 1$

e.g. if $Ns \approx 0.1$ (nearly neutral in single-locus setting)

and $NU = \langle \pi \rangle L \approx \begin{cases} 10^4 & \text{for bacteria in gut} \\ 10^6 & \text{for humans} \end{cases}$ (ignoring recombination, for now)

\Rightarrow then $|X(\vec{s}) - X(\vec{s}')|/N \sim \sqrt{NU(Ns)^2} = \sqrt{10^4 \times 10^{-2}} = 10 \gg 1$

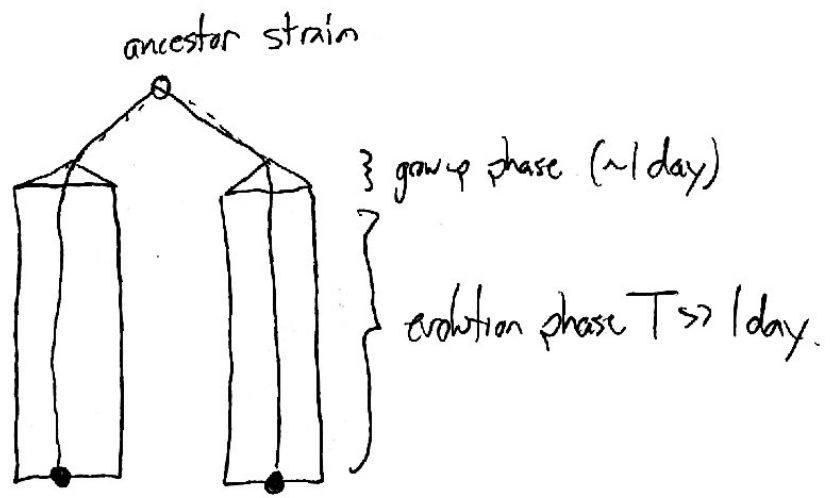
\Rightarrow in this case, neutral model would not be good approx. & we'll have to turn to different methods...

In some cases, neutral picture can still be salvaged if we only care about predicting neutral mutations (e.g. synonymous mutations) and we can find some other way to predict genealogy.

e.g. in evolution experiment:

\Rightarrow if pick 1 individual from each pop'n, we know what their genealogy looks like

(up to initial grow up period.)



~~the model is a Poisson process~~

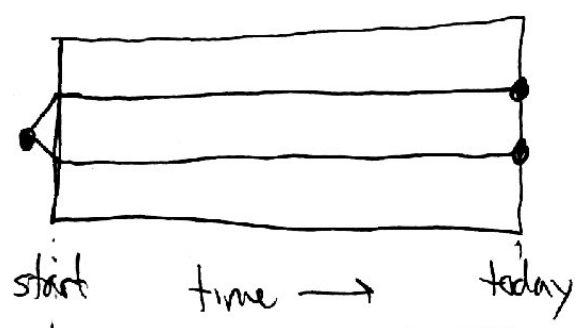
neutral mutations btw \approx Poisson ($2UT$)
1 & 2

regardless of any kind of selection (or anything else) w/in pop'n.

\Rightarrow why doesn't this work for larger sample sizes?

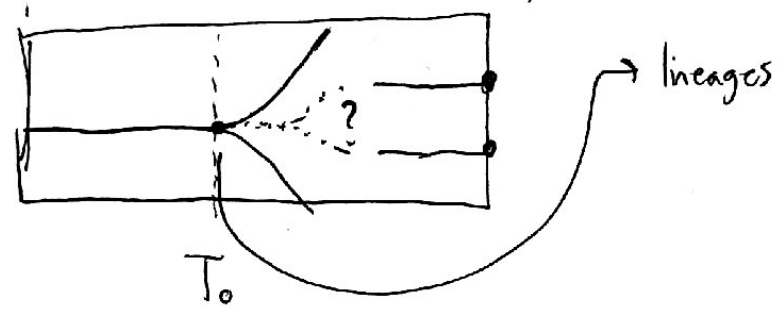
\Rightarrow let's consider 2 scenarios:

(a) Neutral



if $N \gg T$
 \Rightarrow no coalescence until start of experiment.

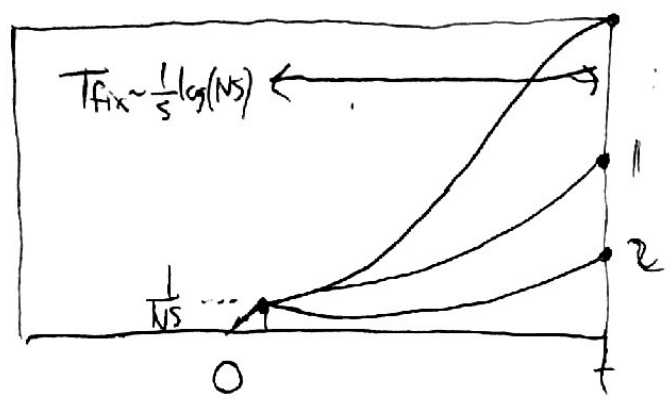
(b) Sweep



\Rightarrow mutations w/in population very different in 2 ~~sc~~ scenarios.

\Rightarrow in case where selected mutation is from SSWM regime can still make progress (common trend \Rightarrow reduce to single locus selection) & things will be easier

in this case:



⇒ ~~the population size is N~~ • sweep acts like effective population w/ size $Nf(t)$

⇒ since $Ns \gg 1 \Rightarrow T_{fix} \ll N \Rightarrow$ no coalescence until $f(t)$ gets small!

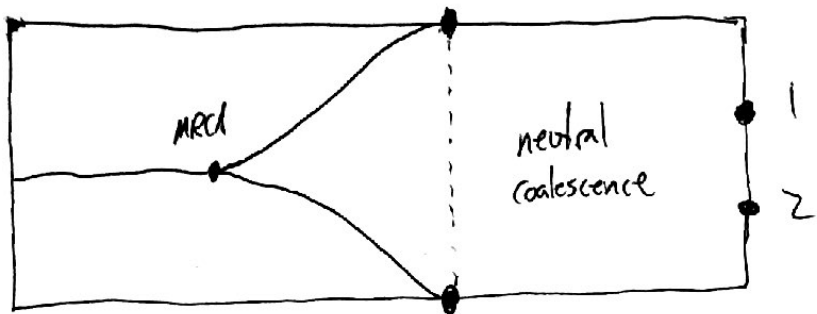
How small? Need $\frac{t}{Nf(t)} \approx 1$ (decent chance of coalescing in that amount of time.)

⇒ no coalescence until $f(t) \sim \frac{1}{Ns}$ ($t \sim \frac{1}{5}$)

⇒ $T_2 \approx \frac{1}{5} \log(Ns) \pm O(\frac{1}{5})$ → difference is small.

~~what if mutation had fixed before time of sampling?~~

what if mutation had fixed before time of sampling?



$T_{\text{test}} \sim \frac{1}{N_b S}$ (exponential distn)

$T_{\text{fix}} \sim \frac{1}{3} \log(N_s) \ll T_{\text{test}}$ (SSWM)

⇒ 2 regimes:

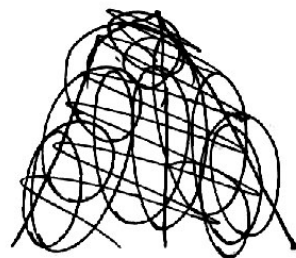
① if $N \ll T_{\text{test}}$ ⇒ coalescence before sweep.
 ($(N_b/N_s) \ll 1$) (looks neutral)

② if $N \gg T_{\text{test}}$ ⇒ no coalescence until $\pm O(\frac{1}{3})$
 of establishment time.
 ($(N_b/N_s) \gg 1$)

$T_2 \approx T_{\text{test}} \sim \text{Exponential} \left(\frac{1}{N_b S} \right)$

⇒ in this case $\pi \approx \langle T_2 \rangle U_0 \approx \frac{U_0}{U_b} \frac{1}{N_s} \rightarrow$ anti-correlated w/ pop size.

⇒ works for larger sample sizes:



Simple selective sweep:



"Star like genealogy"
(mostly singletons)

neutral

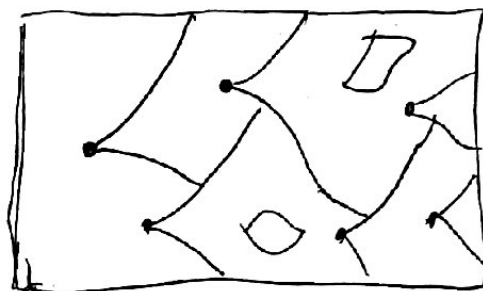


regular Kingman genealogy
(last coalescence is longest)

↳ similar to population expansion

w/ more than 1 selected mutation,
things get complicated. (clonal interference)

⇒ will revisit in a few lectures



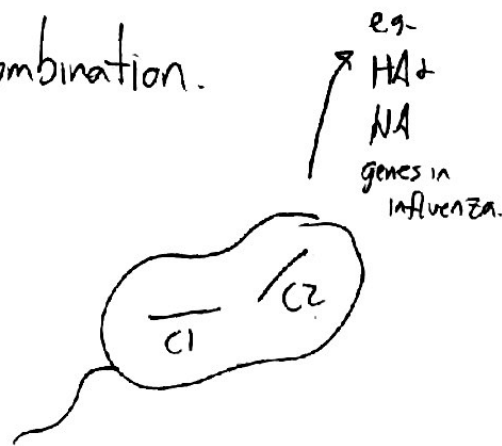
so far, have been talking about asexual populations

⇒ what about recombination? does coalescent picture work there too?

let's go back to neutral case and consider recombination.

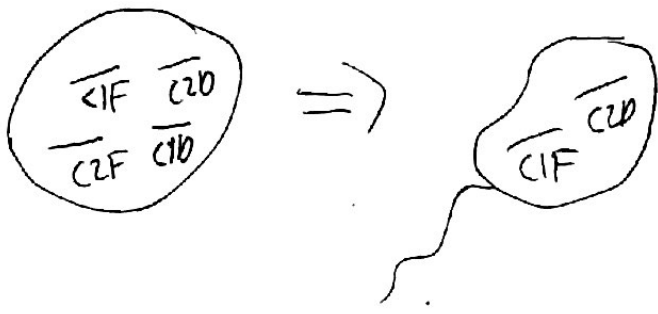
⇒ to keep things simple, let's consider a reassortment model of recombination w/

2 chromosomes of length ~~scribble~~ L



then @ rate e , individual mates w/ another

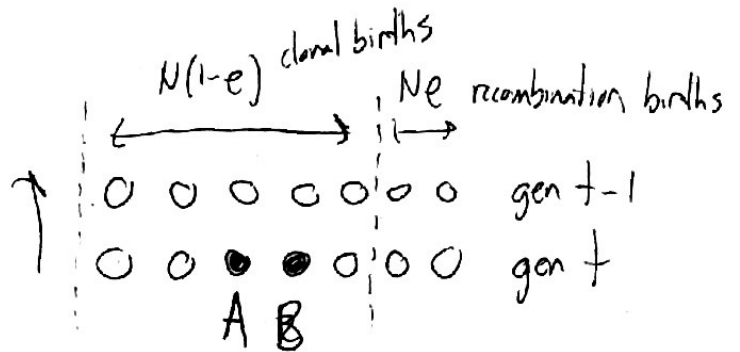
9



+ replaces chromosome w/ one from mate.

(i.e., e = successful recombination)

\Rightarrow now let's see if we can look back in time:



\Rightarrow in last generation, $N(1-e)$ individuals reproduced asexually
~~clonal~~ N_e individuals were result of recombination event.

\Rightarrow probability of coalescence $\approx \frac{1}{N}$

\Rightarrow probability that individual is recombinant $= \frac{N_e}{N(1-e)} \approx e$

\Rightarrow so w/ prob $e^{-T_{MRCA}(2e)}$, there were no recombination events in genealogy of sample.

in this case, T_{MRC} is same as asexual population,

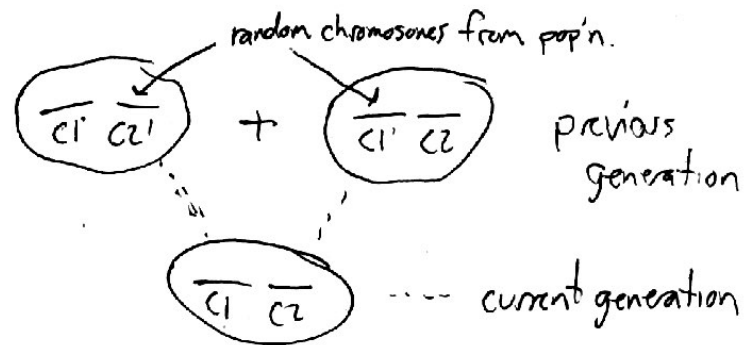
(10)

$$T_{MRC} \sim \text{Exponential}(N), \text{ so } P(\text{no recomb}) = \int_0^{\infty} e^{-T_{MRC}(2e) - T_{MRC}/N} \frac{1}{N} dT$$
$$\approx \frac{1}{1 + 2Ne}$$

\Rightarrow hence if $N_e \ll 1$, genomes behave as if effectively asexual.

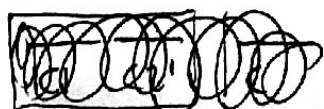
\Rightarrow however, if $N_e \gg 1$, very good chance that in ~~ancestral~~ history of sample, one of ~~ancestral~~ ancestors will be result of recombination event.

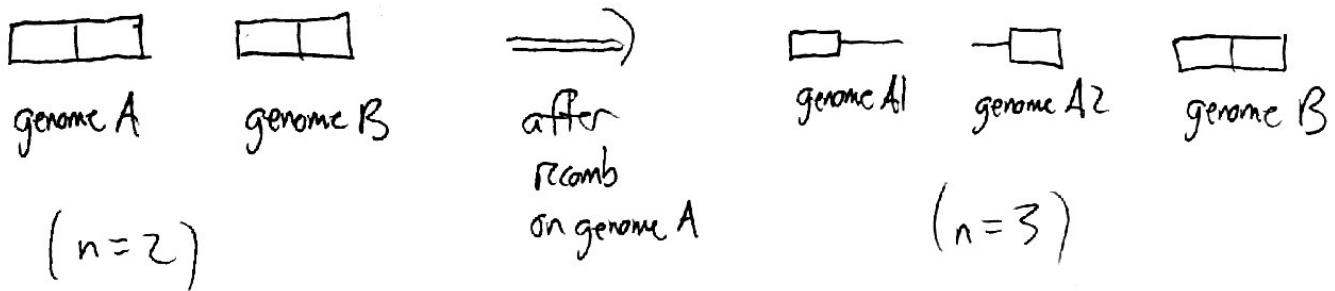
\Rightarrow what does this look like?



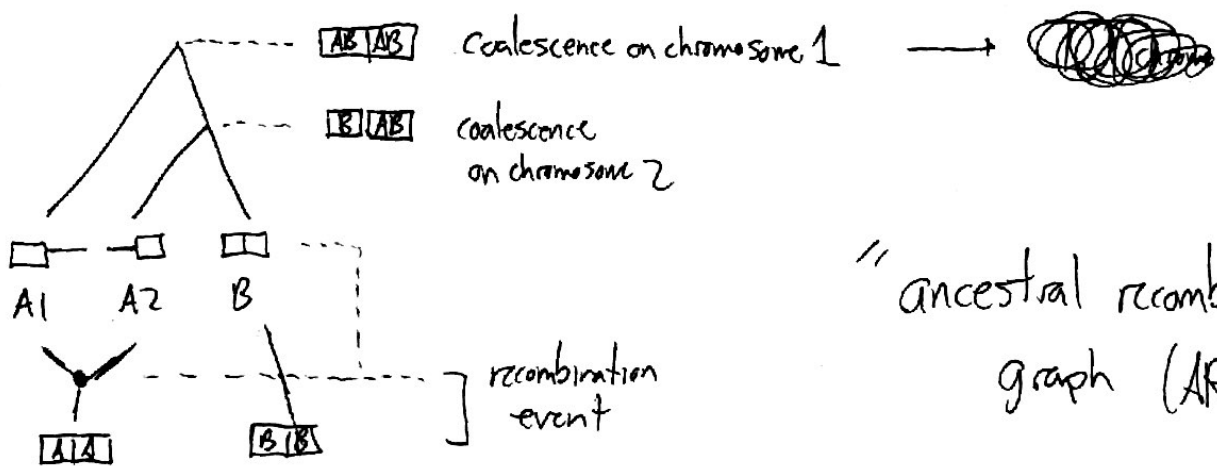
\Rightarrow in this case, ancestors of 2 chromosomes are different i.e., genealogies of 2 chromosomes separate.

\Rightarrow in coalescent picture, effectively increases our sample size:

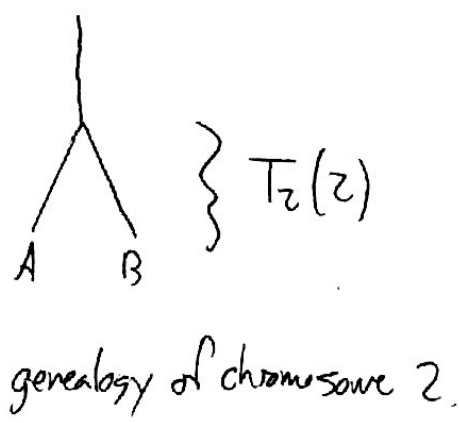
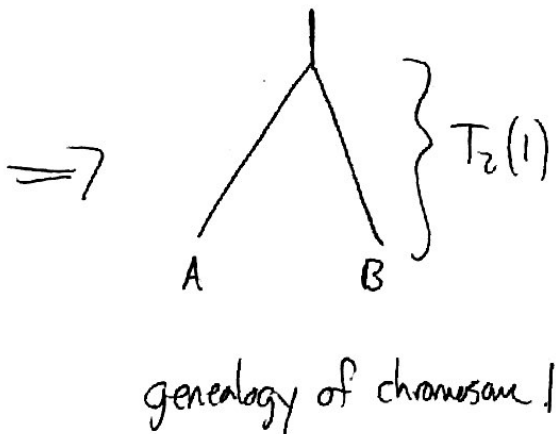




Now coalescent process continues w/ larger sample size:
 eg. if no other recombination events, could get:



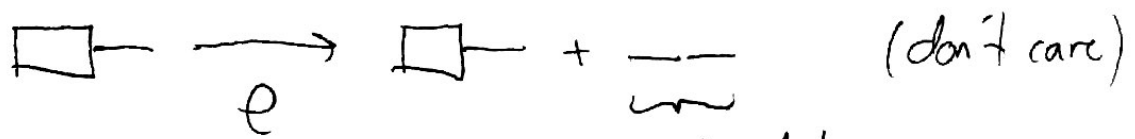
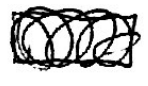
"ancestral recombination"
 graph (ARG)



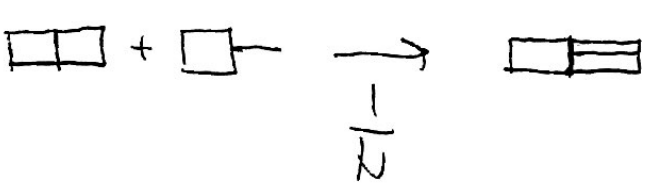
⇒ i.e., recombination allows genealogies to be different @ different sites in genome. (in asexual case, $T_2(1) = T_2(2)$)

of course, this is just one possible arga.

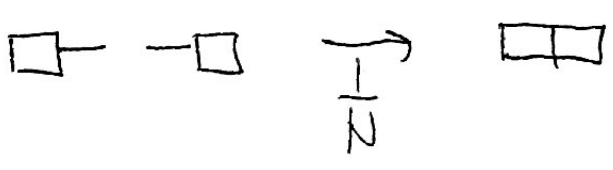
@ each stage, can have events:



no descendants
in sample, so can ignore.



just as likely.



let's see if we can simulate this in our heads in limit that $N_e \gg 1$

① start w/ $\boxed{AA} \quad \boxed{BB}$ coalescence $\frac{1}{N}$, recomb $2e$

② recomb happens first w/ high prob. ($\frac{1}{N} \ll 2e$), ~~edit~~

$\boxed{A} - - \boxed{A} \quad \boxed{B} \boxed{B}$ time $T_{rec} \sim \frac{1}{2e}$

Now coalescence $\frac{\binom{3}{2}}{N}$, recomb ϕ

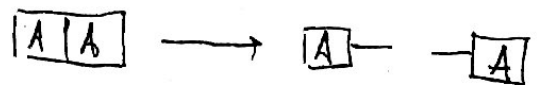
\Rightarrow recomb happens ~~at~~ first w/ high prob. $(\frac{\binom{3}{2}}{N} \ll e)$

\Rightarrow must now happen in other lineage: (Time $\sim \frac{1}{e}$)

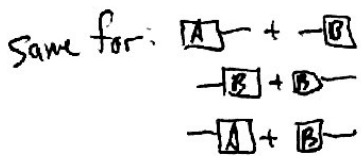


Now coalescence = $\frac{\binom{4}{2}}{N}$, recomb = 0. \Rightarrow next event must be coalescence. (T ~ N)

\Rightarrow 4 possible pairs. if $A-A \rightarrow AA$ then next event will be recomb.

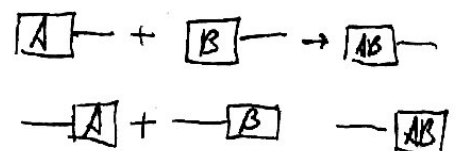


(only time $\sim \frac{1}{e} \ll N$)



\leftarrow so it's almost as if this coalescence event never happened

\Rightarrow can instead keep track of the 2 coalescence events that won't immediately get undone by recomb:



each of the two events happens independently w/ rate $\frac{1}{N}$

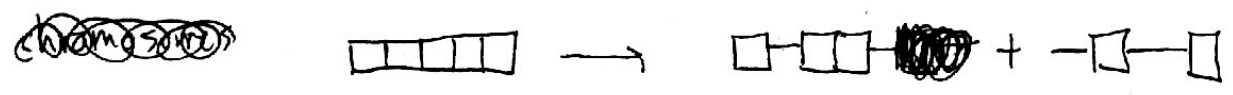
$\Rightarrow T_2(1), T_2(2) \overset{\text{independent}}{\sim} \text{Exponential}(N)$

Putting everything together, we have:

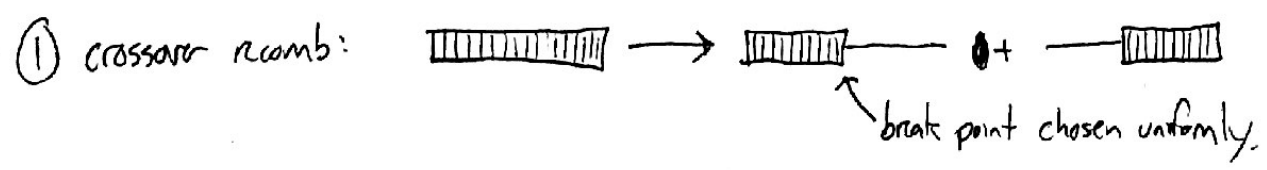
$$p(T_2(1), T_2(2)) \approx \begin{cases} \text{~~circled~~ } \frac{1}{N} e^{-\frac{T_2(1)}{N}} \delta(T_2(1) - T_2(2)) & \text{if } Ne \ll 1 \\ \left[\frac{1}{N} e^{-\frac{T_2(1)}{N}} \right] \cdot \left[\frac{1}{N} e^{-\frac{T_2(2)}{N}} \right] & \text{if } Ne \gg 1 \end{cases}$$

i.e., effectively asexual vs effectively independent.

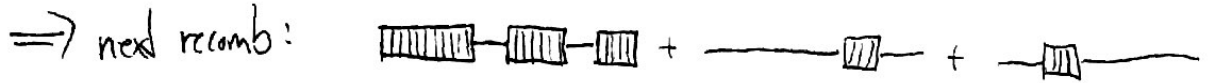
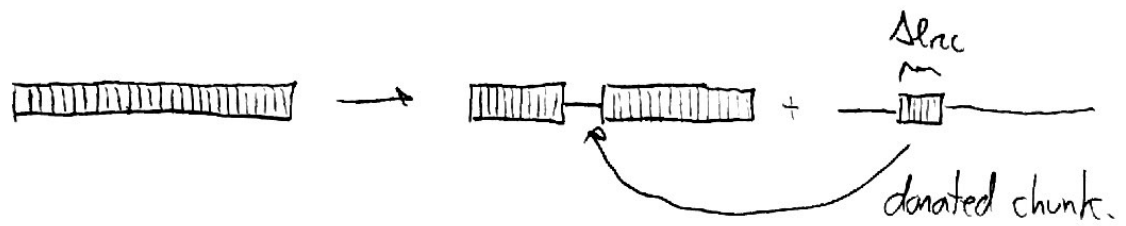
\Rightarrow works for more than 2 chromosomes:



\Rightarrow also works for other modes of recombination:



② HGT

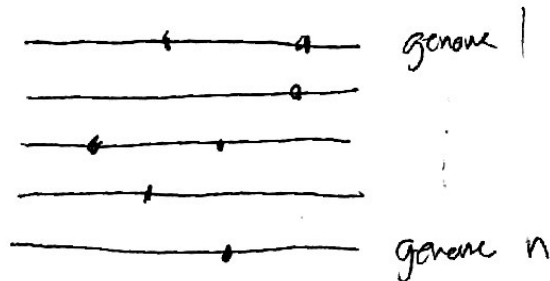


⇒ need $\sim \frac{L}{\Delta_{rec}}$ events to mix up whole genome.

⇒ but actually calculating statistics of ARG is very hard
 (because like coalescent tree w/ 2^c samples, where $c = \#$ of chromosomes.)
 ⇒ essentially intractable (lots of approx models out there...)

⇒ then you have to paint mutations on top of this.

⇒ hard to directly know what data looks like:



⇒ next time: back to forward time approach to calculate some summary statistics.