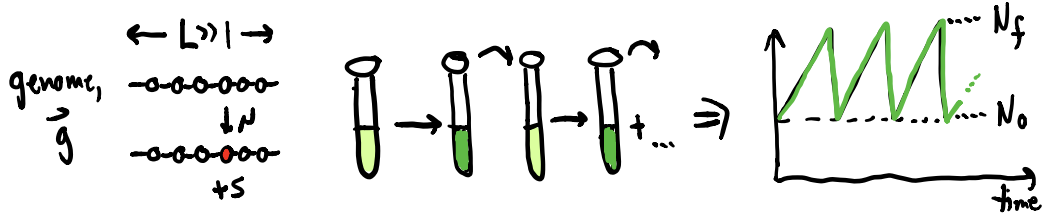


Announcements: Problem Set 4 Posted (Last one!)

DUE: 3/9/21

Last time: Intro to multi-locus models of evolution



$$\frac{\partial f(\vec{g}, t)}{\partial t} = \underbrace{(X(\vec{g}) - \bar{X}) f(\vec{g})}_{\text{natural selection}} + \underbrace{\sum_{\vec{g}'} M(\vec{g}' \rightarrow \vec{g}) f(\vec{g}') - M(\vec{g} \rightarrow \vec{g}') f(\vec{g})}_{\text{mutations}} + \underbrace{\left[\sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}') \right]}_{\text{genetic drift}}$$

$\bar{X}(t) = \int X(\vec{g}) f(\vec{g}) d\vec{g}$
fitness of genotype *mean fitness of pop'n*
incoming *outgoing*



$$\frac{df}{dt} = \underbrace{s f(1-f)}_{\text{natural selection}} + \underbrace{\mu(1-f) - \nu f}_{\text{mutations}} + \underbrace{\sqrt{\frac{f(1-f)}{N}} \eta(t)}_{\text{genetic drift}}$$

so far, similar to L=1 case but with more dimensions...

- Today:
- ① two new biological features that enter for $L \geq 2$.
 - ② How can we start to understand these models?

④ "Epistasis" : properties of $\vec{g} \rightarrow X(\vec{g})$ map
("fitness landscape")

\Rightarrow easiest to motivate w/ $L=2$ case (e.g. 2 gene deletions)

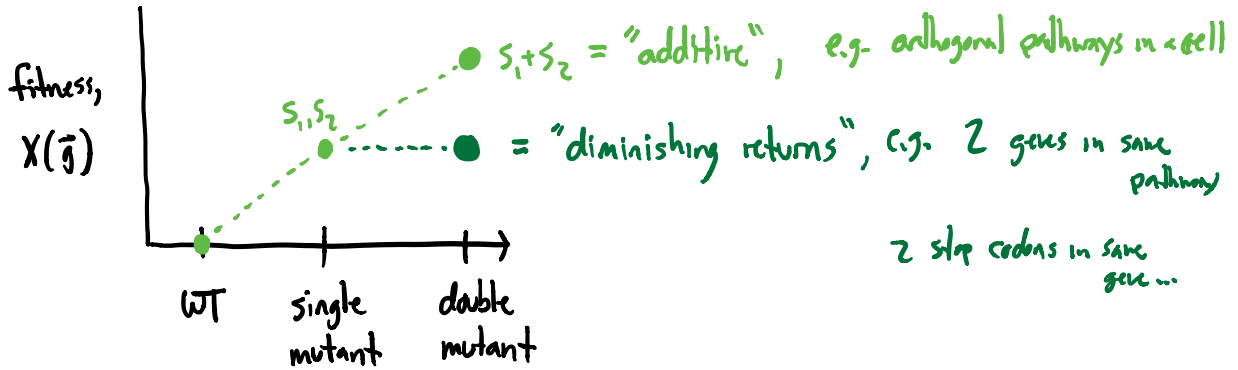
$$X(0,0) \equiv 0 \quad (\text{convention})$$

$$\left. \begin{array}{l} X(1,0) \equiv S_1 \\ X(0,1) \equiv S_2 \end{array} \right\} \text{could measure, e.g. gene deletion screen (HW2)}$$

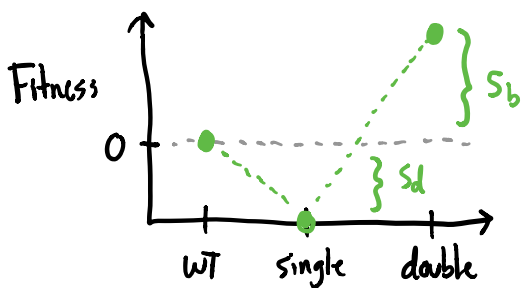
$$X(1,1) \equiv ? \equiv \underbrace{S_1 + S_2}_{\text{"additive part"}} + \underbrace{\epsilon}_{\text{"epistasis"}} \\ (\text{how much deviation from additivity})$$

e.g. " $\epsilon > 0$ " \Rightarrow "positive epistasis" \Rightarrow "sign epistasis"
" $\epsilon < 0$ " \Rightarrow "negative epistasis" etc. etc.

Often easiest to express w/ picture:

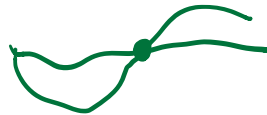


⇒ people often interested in scenarios like:

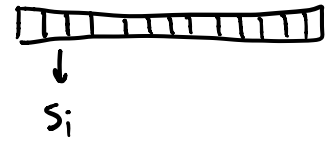


"fitness valley crossing"

e.g. initiation of cancer
contact residues in proteins



⇒ gets even more complicated for $L > 2$:



$$X(\vec{g}) \equiv \underbrace{\sum_{e=1}^L s_e g_e}_{\text{additive part ("coupon collecting")}} + \underbrace{\epsilon(\vec{g})}_{\text{epistatic part.}}$$

⇒ can write as Taylor expansion around WT:

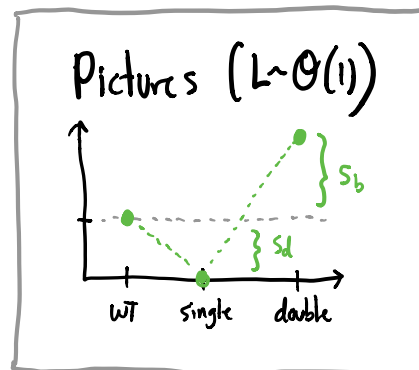
$$\epsilon(\vec{g}) = \underbrace{\sum_{e=1}^L \sum_{e'=1}^L \epsilon_{ee'} g_e g_{e'}}_{\text{"pairwise epistasis"}} + \underbrace{\sum_e \sum_{e'} \sum_{e''} \epsilon_{eee'} g_e g_{e'} g_{e''}}_{\text{"higher order epistasis"}} + \dots$$

⇒ hard to parametrize in general (active area of research!)

⇒ in practice, people often use:

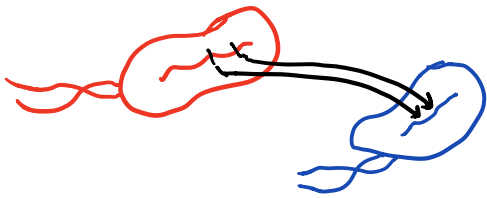
Additive model ($L \gg 1$)

$$X(\vec{g}) \approx \sum_{e=1}^L s_e g_e$$



⇒ other new bit of biology for $L \geq 2$:

⑤ Recombination (exchange of genetic material between different individuals)



Many different mechanisms!

⇒ but many share same basic behavior:

① Focal individual f is chosen to undergo recombination

⇒ w/ probability ρ per individual per-gen → e.g. mating
viruses/phage
uptake of DNA
cellular DNA, d

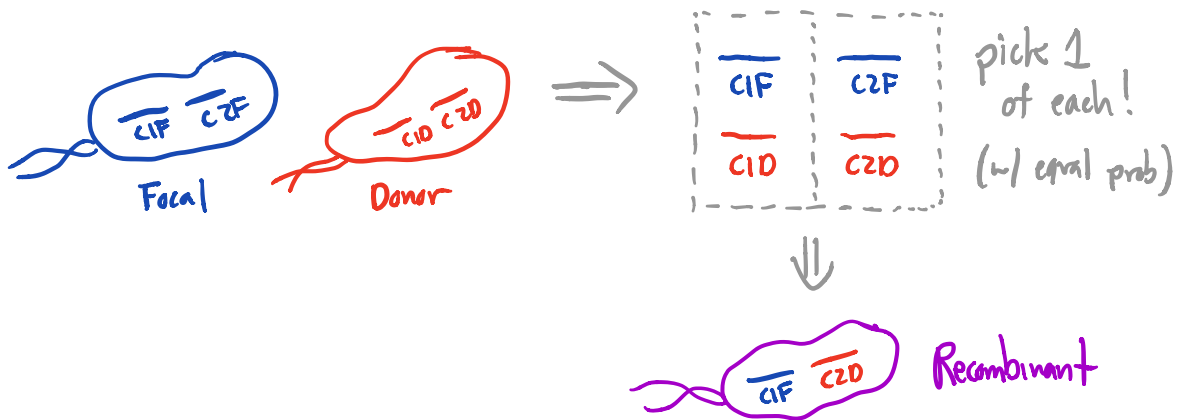
② Donor individual d is chosen to donate portion of genome

⇒ probability $\sim \frac{1}{N} \Rightarrow f(\vec{g})$ for any individual of that genotype.

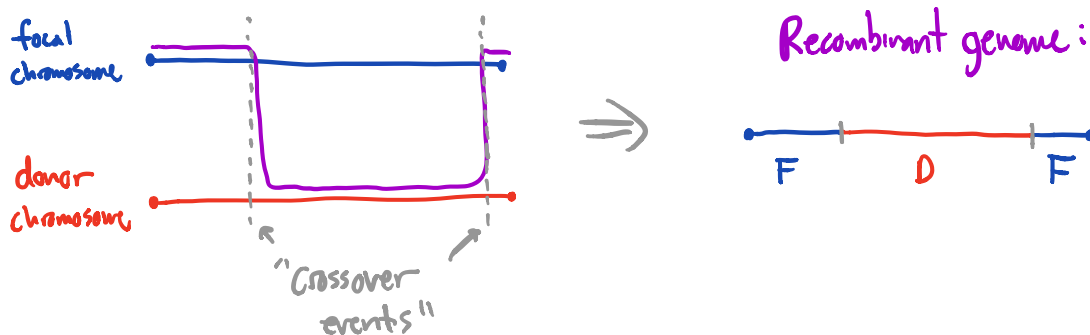
③ Some piece of donor's DNA is integrated into focal genome
⇒ producing "recombinant"

⇒ different mechanisms enter @ this step:

① Reassortment (e.g. different chromosomes, e.g. yeast, humans, influenza.)



② Crossover Recombination (e.g. w/in chromosomes in humans)



⇒ often modeled w/ ~ 1 crossover per recombination event

w/ location chosen uniformly across chromosome 

⇒ in practice, "hot spots" + "cold spots" ⇒ "recombination map"

⇒ effective recombination rates vary over many orders of magnitude for different sites in same genome!

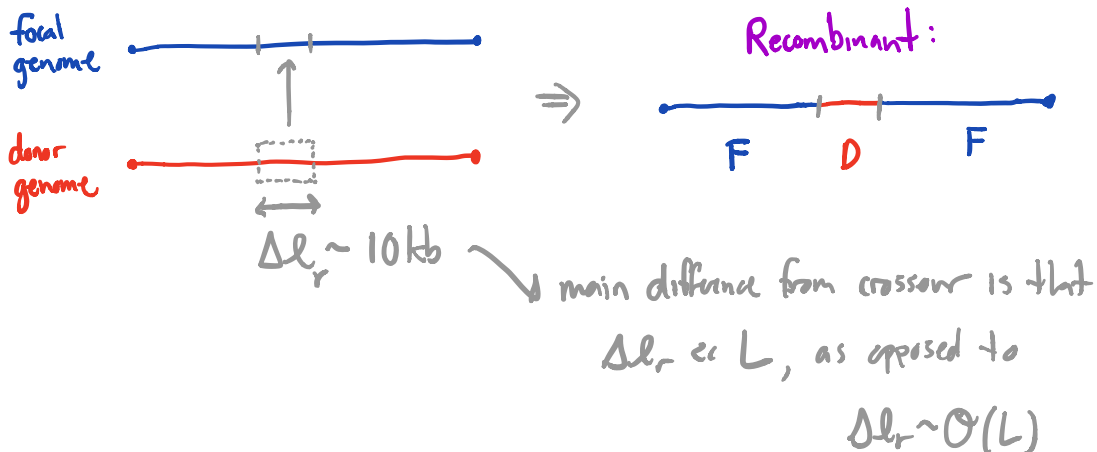
⇒ e.g. in humans ⇒ $L_{\text{chrom}} \sim 10^8$ bp

⇒ $P(\text{recomb}) \sim 100\%$ if 2 ends of same chrom

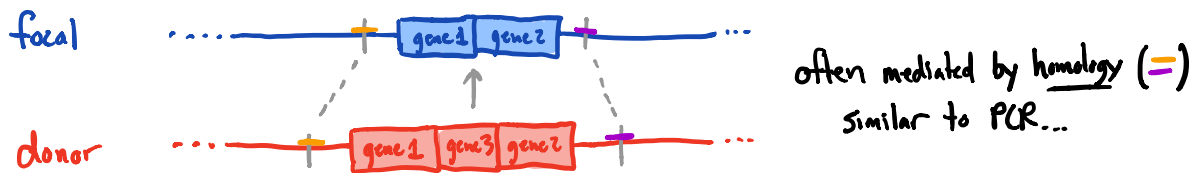
⇒ $P(\text{recomb}) \sim 10^{-8}$ if neighboring base pairs

Ⓒ "Horizontal gene transfer" / "gene conversion"

⇒ lingo is a little controversial, but basic idea pretty simple:



⇒ also a mechanism for gaining + losing genes ("accessory genome")

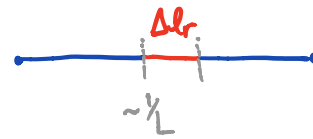


⇒ active area of research!

⇒ but in this class, will mostly focus on "core genome"

⇒ simplest HGT model:

$\Delta L_r = \text{const}$, location ~ uniform



So far: individual-based picture...

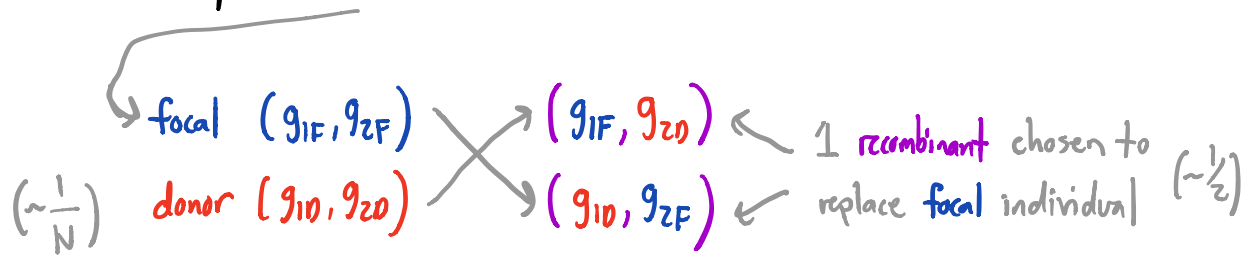
⇒ can we translate to continuum limit?

$$\left(\frac{\delta \mathcal{F}(\vec{g})}{\delta t} \right)_{\text{rec}} = ???$$

⇒ easiest to start w/ $L=2$ case ⇒ $\vec{g} = (g_1, g_2)$

⇒ all mechanisms have same net effect:

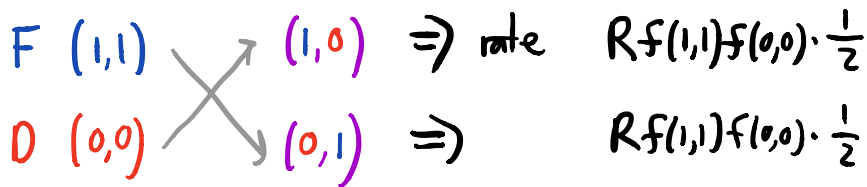
⇒ w/ rate R [function of $\rho, L, \Delta t, \dots$ etc.]



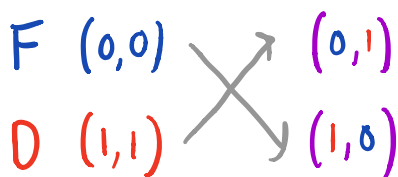
⇒ total outflow from recombination: $-Rf(\vec{g})$

⇒ total inflow? $2^2 \times 2^2 = 16$ possible focal/donor combos

case 1 (of 16):



case 2 (of 16):



same!

Case 3 (of 16):

$$\begin{array}{l} (1,1) \xrightarrow{\quad} (1,0) \Rightarrow Rf(1,1)f(1,0) \frac{1}{2} \\ (1,0) \xrightarrow{\quad} (1,1) \Rightarrow Rf(1,1)f(1,0) \frac{1}{2} \end{array}$$

\Rightarrow after tabulating all 16 combinations (all 32 recombinants)
can add them up to obtain:

$$\left(\frac{\delta f(1,1)}{\delta t} \right)_{rec} = Rf(1,0)f(0,1) - Rf(1,1)f(0,0) \quad \leftarrow \text{same!}$$

$$\left(\frac{\delta f(0,0)}{\delta t} \right)_{rec} = Rf(1,0)f(0,1) - Rf(1,1)f(0,0)$$

$$\left(\frac{\delta f(1,0)}{\delta t} \right)_{rec} = Rf(1,1)f(0,0) - Rf(1,0)f(0,1)$$

$$\left(\frac{\delta f(0,1)}{\delta t} \right)_{rec} = \text{same.}$$

\Rightarrow normalized so that $\sum_{\vec{g}} \delta f(\vec{g})_{rec} = 0 \quad \checkmark$

⇒ harder to write down explicitly for $L > 2$

but will have general form:

$$\left(\frac{\delta f(\vec{g})}{\delta t} \right)_{\text{rec}} = \rho \sum_{\vec{g}_F, \vec{g}_0} \underbrace{T(\vec{g}_F, \vec{g}_0 \rightarrow \vec{g})}_{\substack{\text{"recombination"} \\ \text{kernel}} \rightarrow \text{"tensor"}} \underbrace{f(\vec{g}_F) f(\vec{g}_0)}_{\text{incoming recombinants}} - \underbrace{\rho f(\vec{g})}_{\text{outgoing recombinants.}}$$

nonlinear!

⇒ unlike mutation, can create genotypes far from \vec{g} !

Putting everything together, general multilocus model looks like:

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

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

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

Problem: No exact solution for stationary dist'n, p_{fix} , etc.
 - even for $L=2$!

\Rightarrow What do we do instead?!? \Rightarrow asymptotic approx's

Question: Given parameters ("knobs") $L, N, X(\vec{g}), M, \rho, T$

\Rightarrow what are some limits where we might understand understand this SDE?

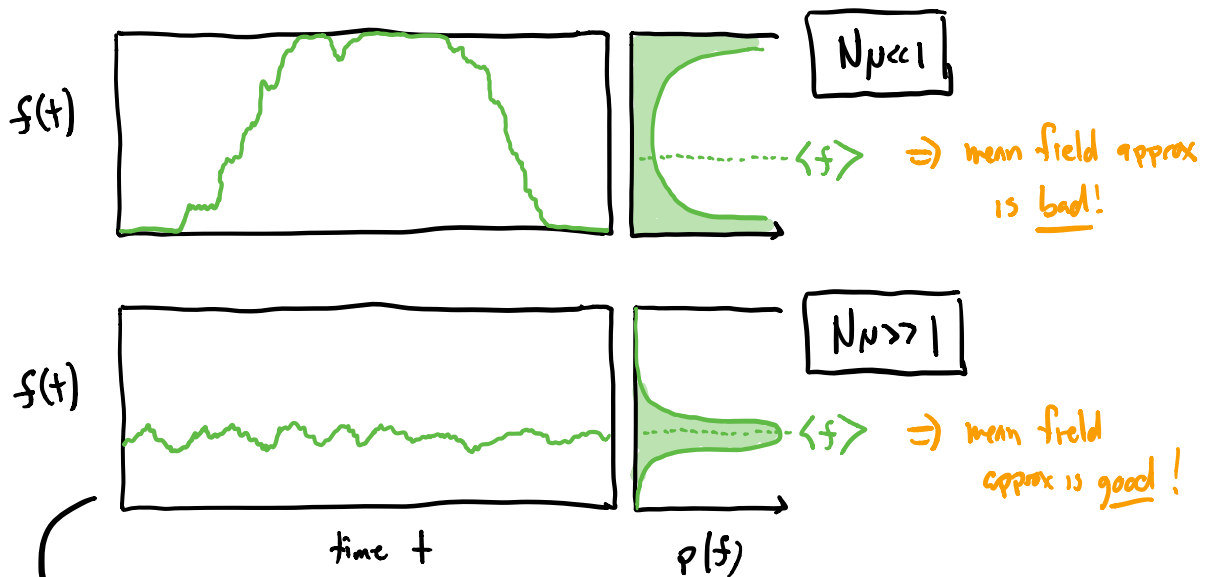
$$\frac{df(\vec{g})}{dt} = \underbrace{\sim (X - \bar{X})}_{\text{blue}} + \underbrace{\sim L^2 \mu}_{\text{orange}}$$

$$+ \underbrace{\sim \rho}_{\text{purple}} + \underbrace{\sim \frac{\rho}{\sqrt{N}}}_{\text{green}}$$

① Obvious answer: $L=1 \Rightarrow$ cheating! *

② in physics, might be primed to take $N \rightarrow \infty$ limit ...
("mean field approx") since @ least noise goes away ...
 \Rightarrow is this a good approx here?

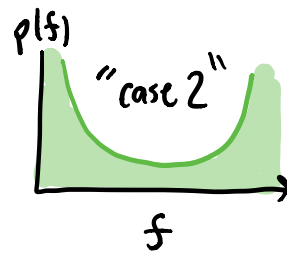
\Rightarrow Recall for $L=1$ case, 2 different regimes when $t \rightarrow \infty$:



key feature: large # of individuals in both genotypes @ same time
 \Rightarrow so fluctuations are small.

\Rightarrow e.g. for $L=2$, might be ok \Rightarrow but for $L \gg 1 \Rightarrow 2^L \gg N!$
e.g. $L=1000 \text{ bp} \Rightarrow 2^L \sim 10^{300}!$

⇒ large L will always look like
 (@ least in some dimensions)



⇒ noise always relevant!

Need to look for other
 approximations of SDE...

$$\frac{ds(\vec{q})}{dt} = \sim (x - \bar{x}) + \sim L \times \mu + \sim \rho + \sim \frac{\sigma}{\sqrt{L}}$$

Let's revisit our first idea ($L=1$)

⇒ even if $L \gg 1$, if behavior "looks like" $L=1$ case,
 ⇒ can use what we already know...

③ Successive mutations regime (i.e. treat **mutation** as small correction)

⇒ what if mutation rates are low enough that
 only 1 or 2 genotypes are present @ a time?