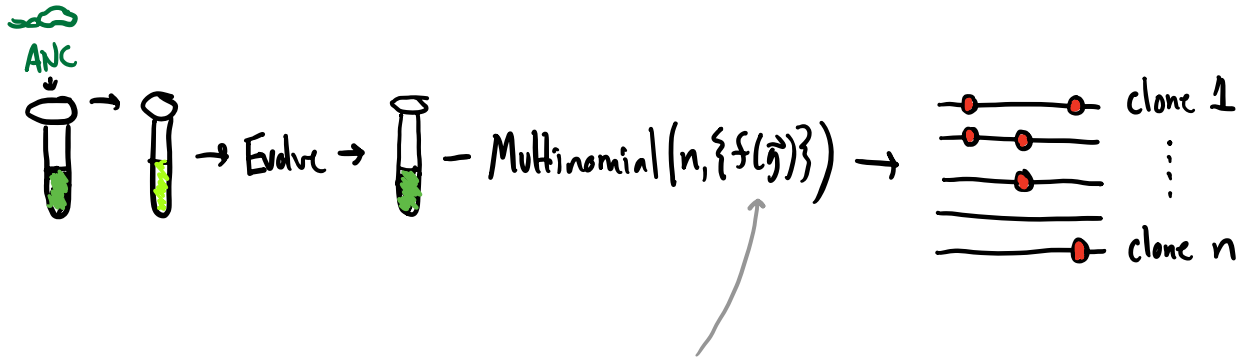# Chapter 9

# Multi-locus models of evolution

# Next Steps:

now that we have methods for measuring genomes (or amplicons)



$\Rightarrow$ need <u>models</u> to predict $f(\vec{g})$'s that arise during evolution

For genome of length $L \Rightarrow 2^L$ possible genotypes

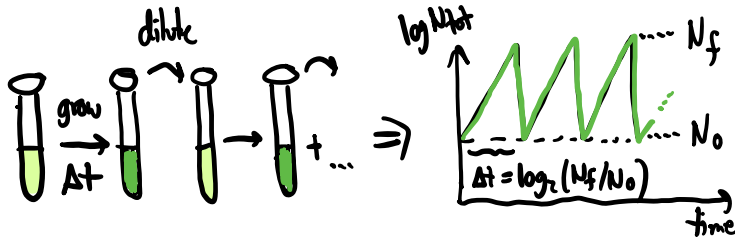e.g. $L=1$: $g=0,1 \Rightarrow \underbrace{f(1) \equiv f}_{\text{mutant}}, \underbrace{f(0)=1-f}_{\text{WT}}$

$L=2$: $\vec{g} = \underbrace{(0,0)}_{\text{WT}}, \underbrace{(1,0), (0,1)}_{\text{single mutants}}, \underbrace{(1,1)}_{\text{double mutant}}$

$L=3$: $\vec{g} = (0,0,0), \underbrace{(1,0,0)}_{\text{single}}, \ldots, \underbrace{(1,1,0)}_{\text{double}}, \ldots, \underbrace{(1,1,1)}_{\text{triple mutant}}$

$\vdots$

etc.

Can we generalize our <u>serial dilution</u> (& <u>diffusion</u>) models?



$$\frac{\partial f(\vec{g})}{\partial t} = \textcolor{red}{???}$$

① <u>Genetic drift</u>:  first assume no growth rate differences...
(& no mutations)

⇒ After 1 day of growth (before dilution):

$$f(\vec{g}) \xrightarrow{\Delta t} \frac{f(\vec{g})e^{r\Delta t}}{\sum\limits_{\vec{g}'} f(\vec{g}')e^{r\Delta t}} = \frac{f(\vec{g})}{\sum\limits_{\vec{g}'} f(\vec{g}')} = f(\vec{g}) \quad \left(\begin{array}{l}\text{i.e. no change}\\ \text{in freqs} \checkmark\end{array}\right)$$

⇒ After dilution step:

① $n(\vec{g}, t+\Delta t) \sim \text{Poisson}\left(\bar{N}_0 \cdot f(\vec{g})\right)$   (sampling)

② $f(\vec{g}, t+\Delta t) = \dfrac{n(\vec{g}, t+\Delta t)}{\sum\limits_{\vec{g}'} n(\vec{g}', t+\Delta t)}$   (re-normalize)

$\Rightarrow$ if repeat our Taylor expansions from Ch 4 ($\bar{N}_0$ large):

$$n(\vec{g}, t + \Delta t) \sim \bar{N}_0 f(\vec{g}, t) + \sqrt{\bar{N}_0 f(\vec{g}, t)} \cdot Z_{\vec{g}}$$

where $Z_{\vec{g}} \overset{iid}{\sim} \text{Gaussian}(0, 1)$

$\Rightarrow$ Taylor expand $f(\vec{g}) = \dfrac{n(\vec{g})}{\sum_{\vec{g}'} n(\vec{g}')}$ :

$$f(\vec{g}, t + \delta t) = f(\vec{g}, t) + \sqrt{\frac{f(\vec{g})\, \delta t}{N_e}} \, Z_{\vec{g}} - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')\, \delta t}{N_e}} \, Z_{\vec{g}'}$$
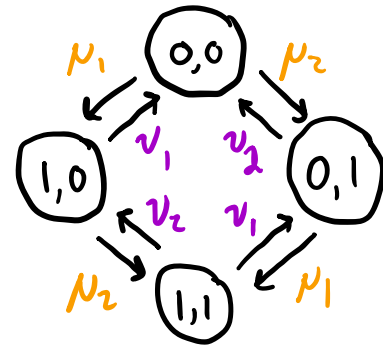
depends on $Z_{\vec{g}}$ @ other $\vec{g}'$!

$\Rightarrow$ $\underline{\text{correlations}}$ between $\delta f(\vec{g})$ & $\delta f(\vec{g}')$ $\Rightarrow$ keeps $f(\vec{g}, t)$
normalized!

$$\sum_{\vec{g}} f(\vec{g}, t + \delta t) = \sum_{\vec{g}} f(\vec{g}) + \sum_{\vec{g}} \sqrt{\frac{f(\vec{g})\delta t}{N_e}} \, Z_{\vec{g}} - \sum_{\vec{g}} f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')\delta t}{N_e}} \, Z_{\vec{g}'}$$

$$= 1 \quad \Rightarrow \text{stays normalized @ later times}$$

## ② Mutations:

=> easiest to start w/ L=2 :

$\mu_1$ (0,0) $\mu_2$
$\nu_1$ $\nu_2$
(1,0) $\nu_2$ $\nu_1$ (0,1)
$\mu_2$ (1,1) $\mu_1$

=> key feature: can only move by ~1 step @ a time

=> generalizing results from L=1 case, after 1 dilution:

mutations into genotype

$$n(1,0,t+\Delta t) \sim Poisson\left(N_0 f(1,0,t) + N_0 \Delta t\left[\overbrace{\mu_1 f(0,0,t) + \nu_2 f(1,1,t)}\right]\right.$$
$$\left. - N_0 \Delta t\left[\underbrace{\mu_2 f(1,0,t) + \nu_1 f(1,0,t)}\right]\right)$$

mutations out of genotype

=> __continuum limit__ (i.e. Taylor expansions):

$$\delta f(1,0)_{mut} = \left[\mu_1 f(0,0) + \nu_2 f(1,1) - \mu_2 f(1,0) - \nu_1 f(1,0)\right]\delta t$$

(+ noise from drift)

linear in ✳
genotype freqs

$\Rightarrow$ larger $L$'s are similar, but more work to write out...

$\Rightarrow$ one way is:

mutations **into** genotype

$$\left[\frac{\delta f(\vec{g})}{\delta t}\right]_{mut} \equiv \sum_{\substack{nearest \\ neighbors \\ \vec{g}'}} \sum_{\ell=1}^{L} \left[ \mu_\ell f(\vec{g}') g_\ell (1-g'_\ell) + \nu_\ell f(\vec{g}')(1-g_\ell) g'_\ell \right]$$

mutations **out of** genotype

$$- \sum_{\ell=1}^{L} \left[ \mu_\ell f(\vec{g})(1-g_\ell) + \nu_\ell f(\vec{g}) g_\ell \right]$$

$$\left(\frac{\delta f(\vec{g})}{\delta t}\right)_{mut} \equiv \sum_{\vec{g}'} \left[ \underbrace{M(\vec{g}' \to \vec{g})}_{} f(\vec{g}') - \underbrace{M(\vec{g} \to \vec{g}')}_{} f(\vec{g}) \right]$$

$2^L \times 2^L$ matrix of mut'n rates

<u>Note</u>: mutation matrix normalized s.t. $\sum_{\vec{g}} \left(\frac{\delta f(\vec{g})}{\delta t}\right)_{mut} = 0$

$$\Rightarrow \text{ ensures that } \sum_{\vec{g}} f(\vec{g}, t+\delta t) = \underbrace{\sum_{\vec{g}} f(\vec{g}, t)}_{1} + \underbrace{\sum_{\vec{g}} \delta f_{mut}(\vec{g})}_{0} = 1$$

③ **Selection** (growth rate differences)

If growth rate of genotype $\vec{g}$ is $\equiv r + X(\vec{g})$

$\Rightarrow$ then after 1 cycle of growth:

$$f(\vec{g}) \longrightarrow \frac{f(\vec{g})\, e^{[r+X(\vec{g})]\Delta t}}{\sum_{\vec{g}'} f(\vec{g}')\, e^{[r+X(\vec{g}')]\Delta t}} = \frac{f(\vec{g})\, e^{X(\vec{g})\Delta t}}{\sum_{\vec{g}'} f(\vec{g}')\, e^{X(\vec{g}')\Delta t}}$$

$\Rightarrow$ if $X(\vec{g})\Delta t \ll 1$ (continuum limit) $\Rightarrow$ Taylor expand:

$$f(\vec{g}, t+\delta t) \approx f(\vec{g}, t) + \left[X(\vec{g}) - \bar{X}(t)\right] f(\vec{g}, t)\, \delta t$$

where $\bar{X}(t) \equiv \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t)$ (population mean fitness)

✱ note: <u>not</u> an ensemble avg! $\langle \bar{X}(t) f(\vec{g}, t)\rangle \neq \langle \bar{X}(t)\rangle \langle f(\vec{g}, t)\rangle$

$\Rightarrow$ stays normalized: $\sum_{\vec{g}} f(\vec{g}, t+\delta t) = \sum_{\vec{g}} f(\vec{g}, t) + \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t) - \sum_{\vec{g}} f(\vec{g}) \sum_{\vec{g}'} X(\vec{g}') f(\vec{g}') = 1$

$\Rightarrow$ 2 new biological features that enter for $L \geq 2$:

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)

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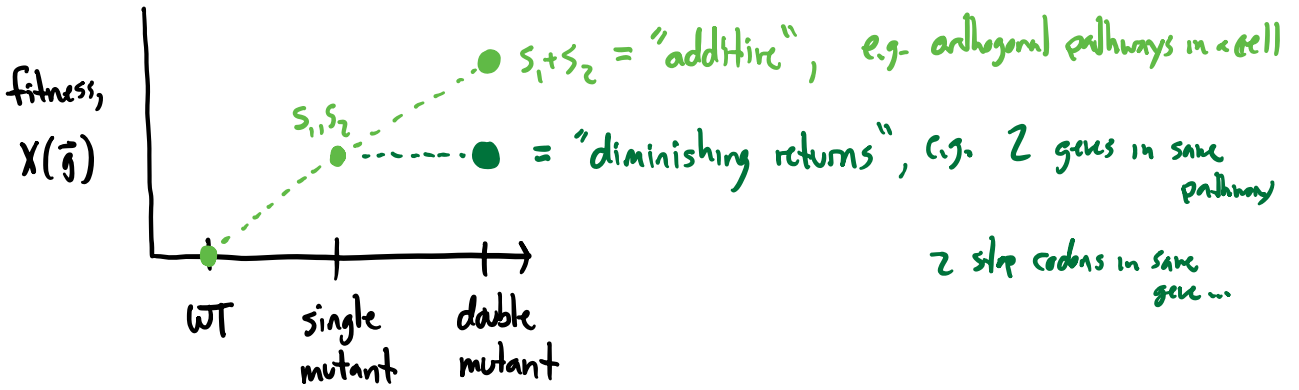
$\left. \begin{array}{l} X(1,0) \equiv S_1 \\ X(0,1) \equiv S_2 \end{array} \right\}$ could measure, e.g. gene deletion screen (HW2)

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

"epistasis"
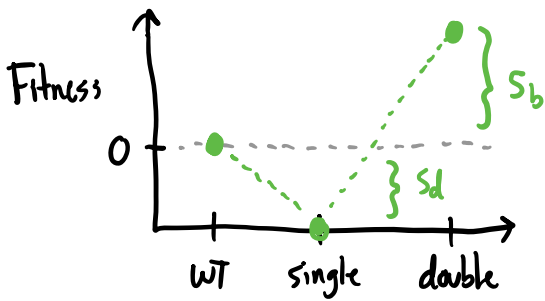(how much deviation from additivity)

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

# Often easiest to express w/ picture:



fitness, $X(\bar{g})$

$S_1 + S_2 =$ "additive", e.g. orthogonal pathways in a cell

$S_1, S_2$

$=$ "diminishing returns", e.g. 2 genes in same pathway

2 stop codons in same gene ...

WT          single mutant          double mutant

$\Rightarrow$ people often interested in scenarios like:



Fitness

0

$\}\, S_b$

$\}\, S_d$

WT     single     double

"fitness valley crossing"

e.g. initiation of cancer
contact btw AA's in proteins
$\hookrightarrow$ (Problem 6 of HW 3)

$\Rightarrow$ gets even more complicated for $L > 2$:



$$X(\vec{g}) \equiv \sum_{\ell=1}^{L} s_\ell g_\ell + \epsilon(\vec{g})$$

$\underbrace{\phantom{\sum_{\ell=1}^{L} s_\ell g_\ell}}$ additive part ("coupon collecting")

$\underbrace{\phantom{\epsilon(\vec{g})}}$ epistatic part.

$\Rightarrow$ can write as Taylor expansion around WT:

$$\epsilon(\vec{g}) = \sum_{\ell=1}^{L} \sum_{\ell'=1}^{L} \epsilon_{\ell\ell'} g_\ell g_{\ell'} + \sum_{\ell} \sum_{\ell'} \sum_{\ell''} \epsilon_{\ell\ell'\ell''} g_\ell g_{\ell'} g_{\ell''} + \ldots$$

$\underbrace{\phantom{\sum \sum}}$ "pairwise epistasis"

$\underbrace{\phantom{\sum \sum \sum}}$ "higher order epistasis"
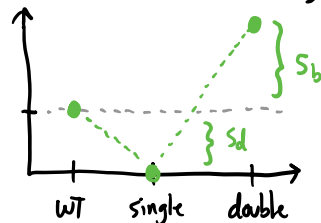
$\Rightarrow$ hard to parameterize in general (active area of research!)

$\quad \Rightarrow$ in practice, people often use:

Additive model ($L \gg 1$)

$$X(\vec{g}) \approx \sum_{\ell=1}^{L} s_\ell g_\ell$$

Pictures ($L \sim \mathcal{O}(1)$)

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

(5) <u>Recombination</u>    ( exchange of genetic material
                              between different individuals )



Many different mechanisms!

   $\Rightarrow$ but many share same basic behavior:

(1) **Focal individual** (f) is chosen to undergo recombination

   $\Rightarrow$ w/ probability $\rho$ per individual per-gen    $\nearrow$ e.g. mating
                                                                      viruses/phage
                                                                      uptake of DNA
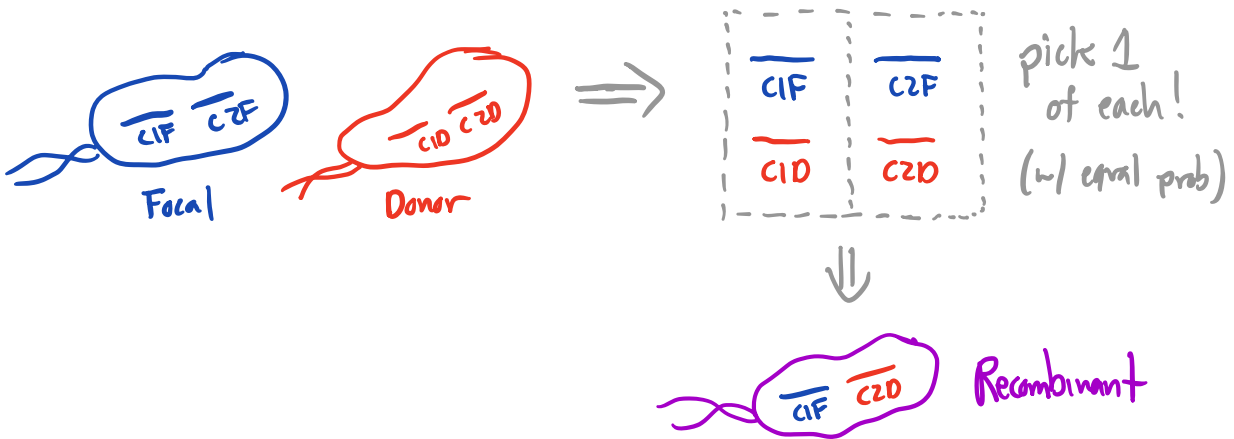                                                                      cellular DNA, d

(2) **Donor individual** (d) is chosen to donate portion of genome

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

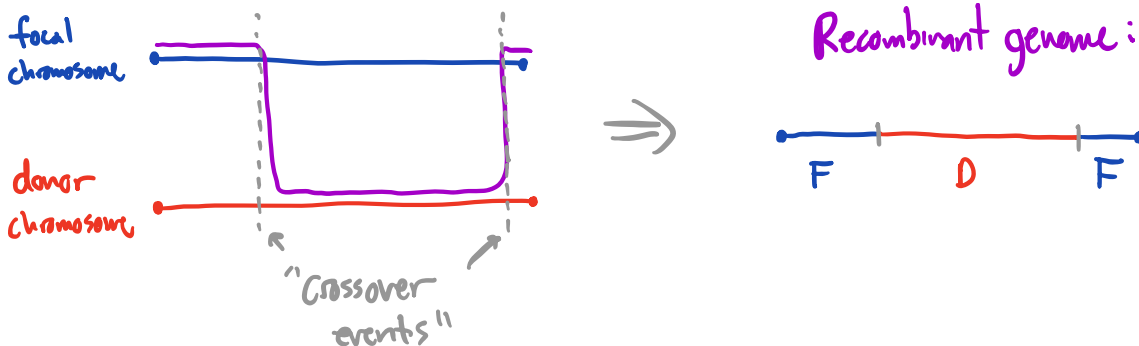③ Some piece of __donor's__ DNA is integrated into __focal__ genome

⇒ producing "__recombinant__"

⇒ different mechanisms enter @ this step:

ⓐ __Reassortment__ (e.g. different chromosomes, e.g. yeast, humans, influenza.)



Focal      Donor

| C1F | C2F |
| C1D | C2D |

pick 1 of each! (w/ equal prob)

C1F C2D   Recombinant

ⓑ __Crossover Recombination__ (e.g. w/in chromosomes in humans)

focal chromosome

donor chromosome

"Crossover events"

Recombinant genome:

F   D   F

$\Rightarrow$ often modeled w/ ~1 crossover per recombination event

w/ location chosen <u>uniformly</u> across chromosome



$\Rightarrow$ in practice, "hot spots" & "cold spots" $\Rightarrow$ "recombination map"

$\Rightarrow$ effective recombination rates vary over <u>many</u> <u>orders-of-magnitude</u> for different pairs of sites in <u>same genome</u>!
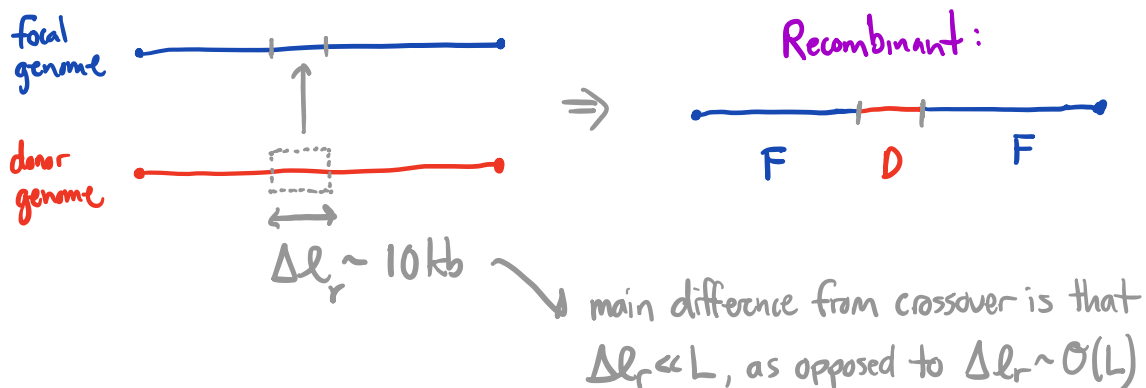
$\Rightarrow$ e.g. in humans $\Rightarrow L_{chrom} \sim 10^8$ bp (x 23 chromosomes)

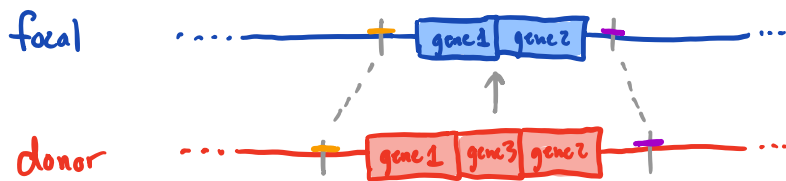$\Rightarrow P(recomb) \sim 100\%$ if opp. ends of same chrom (or diff chroms)

$\Rightarrow P(recomb) \sim 10^{-8}$ if neighboring bp

Ⓒ "Horizontal gene transfer (HGT)" / "gene conversion"

$\Rightarrow$ lingo is a little controversial, but basic idea pretty simple:



focal genome

donor genome

$\Delta \ell_r \sim 10$ kb

Recombinant:

F    D    F

main difference from crossover is that $\Delta \ell_r \ll L$, as opposed to $\Delta \ell_r \sim O(L)$

$\Rightarrow$ also a mechanism for gaining & losing genes ("accessory genome")



focal

donor

| gene1 | gene2 |

| gene1 | gene3 | gene2 |

often mediated by homology (—) similar to PCR...

$\Rightarrow$ active area of research!

$\Rightarrow$ but in this class, will mostly focus on "core genome"

$\Rightarrow$ simplest HGT model:

$\Delta \ell_r$ = const, location ~ uniform



$\Delta \ell_r$
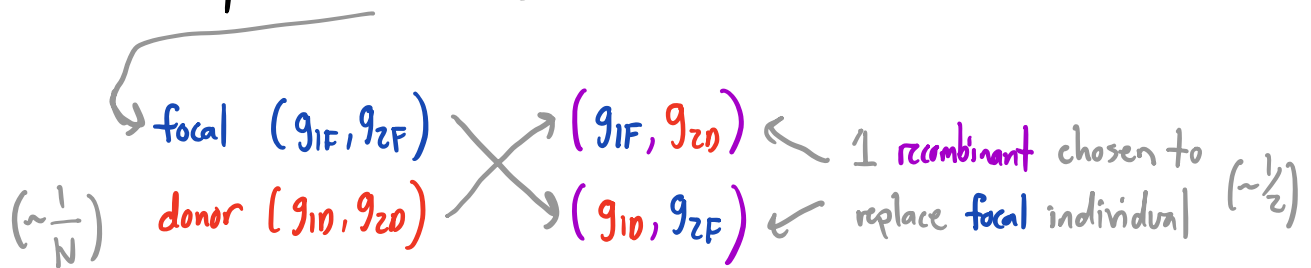
~½L

So far: individual-based picture...

$\Rightarrow$ can we translate to continuum limit?

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

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

$\Rightarrow$ all mechanisms have same net effect:

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

focal $(g_{1F}, g_{2F})$ $\times$ $(g_{1F}, g_{2D})$ $\leftarrow$ 1 recombinant chosen to

$(\sim \frac{1}{N})$ donor $(g_{1D}, g_{2D})$ $(g_{1D}, g_{2F})$ $\leftarrow$ replace focal individual $(\sim \frac{1}{2})$
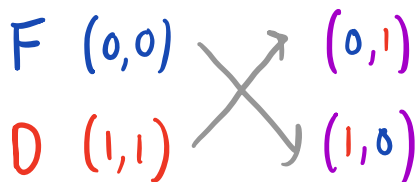
$\Rightarrow$ total <u>outflow</u> from recombination: $-R f(\vec{g})$

$\Rightarrow$ total <u>inflow</u>? $2^2 \times 2^2 = 16$ possible focal/donor combos

## case <u>1</u> (of 16):

F $(1,1)$ $\times$ $(1,0)$ $\Rightarrow$ rate $R f(1,1) f(0,0) \cdot \frac{1}{2}$

D $(0,0)$ $(0,1)$ $\Rightarrow$ $R f(1,1) f(0,0) \cdot \frac{1}{2}$

## case <u>2</u> (of 16): F $(0,0)$ $\times$ $(0,1)$ same!

D $(1,1)$ $(1,0)$

Case 3 (of 16):   (1,1)　　(1,0)  $\Rightarrow$  $R f(1,1) f(1,0) \frac{1}{2}$

(1,0)　　(1,1)  $\Rightarrow$  $R f(1,1) f(1,0) \frac{1}{2}$

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

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

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

same!

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

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

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

$\Rightarrow$ harder to write down explicitly for $L>2$ ....

but will have general form:

$$\left(\frac{\delta f(\vec{g})}{\delta t}\right)_{rec} = \rho \underbrace{\sum_{\vec{g}_F, \vec{g}_0} \overbrace{T(\vec{g}_F, \vec{g}_0 \to \vec{g})}^{} f(\vec{g}_F) f(\vec{g}_0)}_{\substack{\text{"recombination}\\ \text{kernel"} \Rightarrow \text{"tensor"}}} \Big) - \rho f(\vec{g})$$

incoming recombinants

outgoing recombinants.

nonlinear!

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

Putting everything together, general multilocus model looks like:

$$\frac{\partial f(\vec{g})}{\partial t} = \left[ X(\vec{g}) - \bar{X}(t) \right] f(\vec{g}) + \sum_{\vec{g}'} M(\vec{g}' \to \vec{g}) f(\vec{g}') - M(\vec{g} \to \vec{g}') f(\vec{g})$$

selection (nonlinear)          mutation (linear, "local")

$$+ \varrho \sum_{\vec{g}_F, \vec{g}_D} T(\vec{g}_F, \vec{g}_D \to \vec{g}) f(\vec{g}_F) f(\vec{g}_D) - \varrho f(\vec{g})$$

recombination (nonlinear, non-local)

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

genetic drift (stochastic)

<u>Problem</u>: No exact solution for stationary dist'n, $P_{fix}$, etc.
— even for $L = 2$!

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

<u>Question</u>: Given parameters ("knobs") $L, N, X(\vec{g}), M, \varrho, T$

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

$$\frac{\partial f(\vec{g})}{\partial t} = \sim (x - \bar{x}) + \sim L * \mu$$
$$+ \sim \varrho + \sim \frac{z}{\sqrt{N}}$$

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

② in physics, might be primed to take $N \to \infty$ 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 $t \to \infty$:



$f(t)$

$\boxed{N\mu \ll 1}$

$\langle f \rangle$ $\Rightarrow$ mean field approx is **bad**!

$f(t)$

$\boxed{N_N \gg 1}$

$\langle f \rangle$ $\Rightarrow$ mean field approx is **good**!

time $t$            $p(f)$

**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. $L \sim 1000\,bp \Rightarrow 2^L \sim 10^{300}$ !

$\Rightarrow$ large L will always look like

(@ least in some dimensions)

$p(f)$

"case 2"

$f$

$\Rightarrow$ noise always relevant!

Need to look for other
approximations of SDE...

$$\frac{\partial S(\vec{g})}{\partial f} = \sim(x - \bar{x}) \; + \; \sim L * \mu$$
$$+ \; \sim \rho \; + \; \sim \frac{Z}{\sqrt{N}}$$

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

$\Rightarrow$ even if L>>1, if behavior "looks like" L=1 case,

$\Rightarrow$ can use what we already know...

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

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