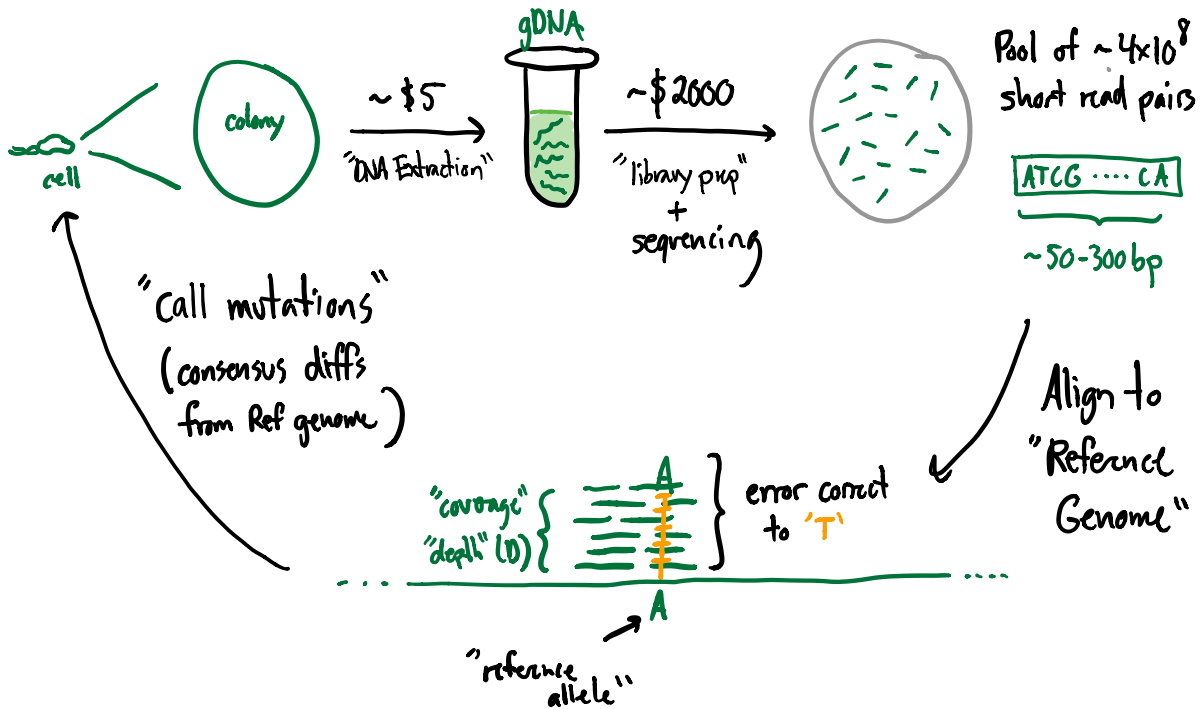


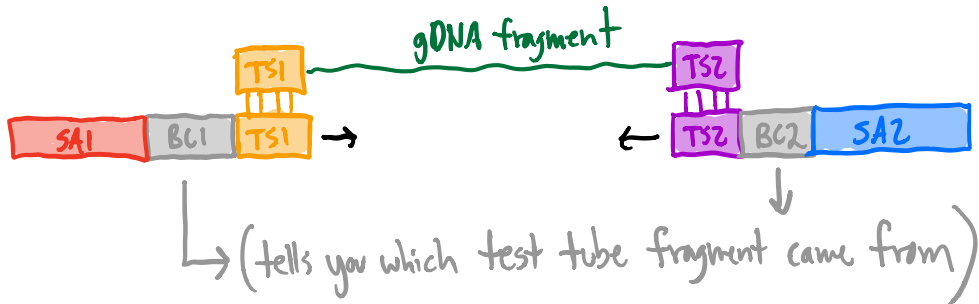
Announcements: PSET 3 DUE 2/23, Office Hours Today (12:30-2pm)
Anita Problem Session Fri (Slack)

Last time: Next-gen (Illumina) sequencing

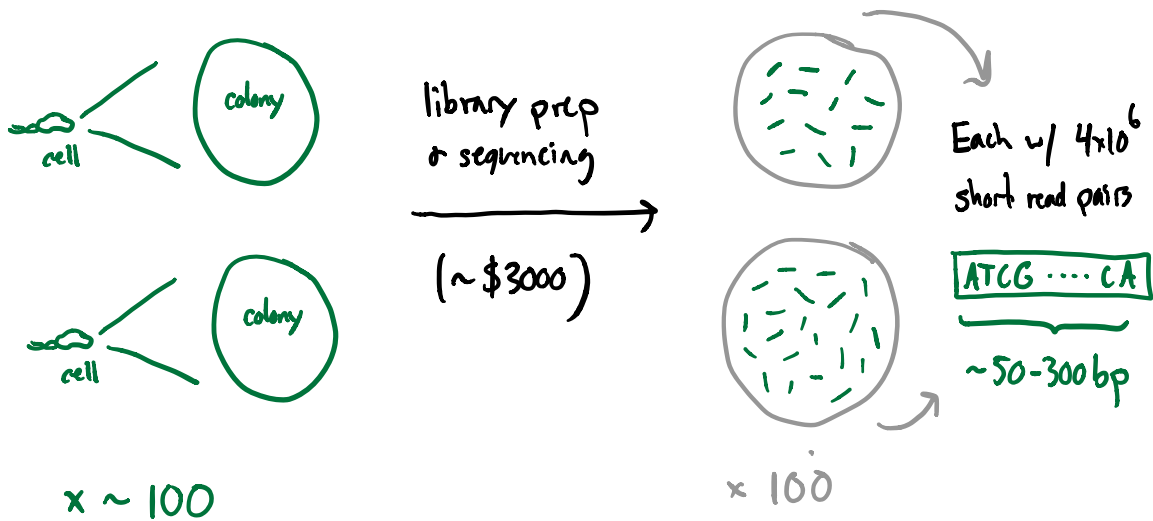


Overkill? $\sim 4 \times 10^8$ reads $\Rightarrow 10^{11}$ bp $\Rightarrow 10^5$ -fold coverage of E. coli genome

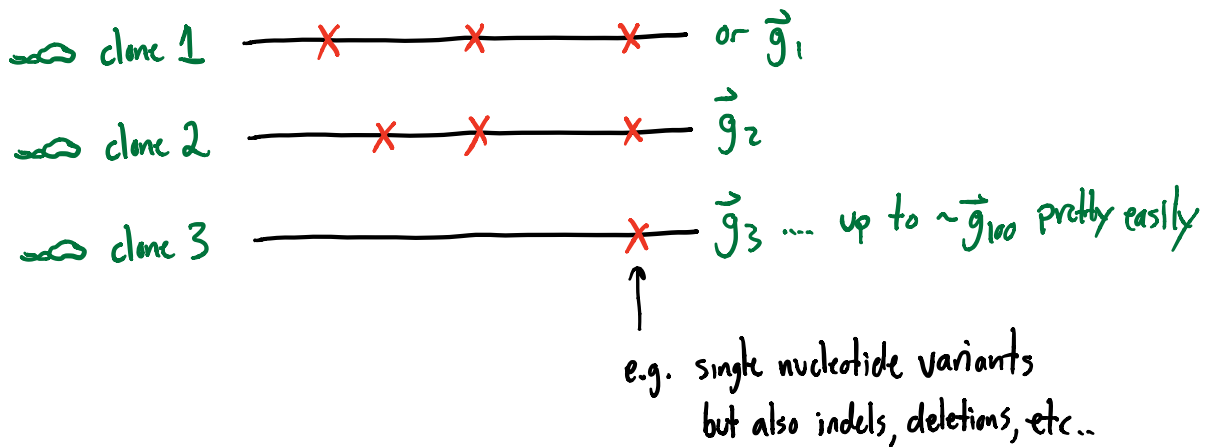
Solution: "multiplex": Add sample specific "barcode" sequence during library prep step



Upshot: can sequence ~100 E. coli libraries on one flow cell & get 300-fold coverage of E. coli genome



⇒ After aligning reads & detecting "true" mutations,
get sequences of genomes:



How are sampld genomes related to dist'n of genomes in pop'n?

⇒ let $n_{\vec{g}}$ = # sampled clones w/ genome \vec{g} (random from sampling)

$f(\vec{g})$ = frequency of genome \vec{g} in population
 (random from evolution)

⇒ then $\Pr[\{n_{\vec{g}}\} | n, \{f(\vec{g})\}] =$ multinomial distribution

↑
total # clones sampled

$$\propto \prod_{\vec{g}} \frac{f(\vec{g})^{n_{\vec{g}}}}{n_{\vec{g}}!}$$

⇒ Genotype space is huge!

⇒ often coarse-grain to summary statistics. $\begin{cases} 1 & \text{if mut @ site } \ell \\ 0 & \text{else} \end{cases}$

e.g. $n_\ell = \# \text{ individuals w/ mutation @ site } \ell = \sum_{\vec{g}} g_\ell n_{\vec{g}}$

⇒ can show $\Pr[n_\ell | n, \{f(\vec{g})\}] = \text{Binomial}(n, f_\ell)$

$$f_\ell = \sum_{\vec{g}} g_\ell f(\vec{g})$$

e.g. total # of mutations separating 2 genomes

⇒ since depends on genome length, often normalized by L :

$$\frac{\# \text{ mutations between 2 random clones}}{L} = \begin{cases} \text{"heterozygosity" } (\pi) & \text{if from same pop'n} \\ \text{"divergence" } (d) & \text{if from diff "species" (or isolated sub-pops)} \end{cases}$$

e.g. heterozygosity (π) in humans is $\sim 10^{-3}$

divergence (d) between humans + chimps is $\sim 10^{-2}$

heterozygosity (π) between E. coli in different humans is $\sim 10^{-2}$

Can we relate π to genotype distribution, $f(\vec{j})$?

Note that:

$$\pi = \frac{1}{L} \sum_{e=1}^L \left[g_{1e}(1-g_{2e}) + (1-g_{1e})g_{2e} \right]$$

$\nearrow i \in \{0,1\}$

$$\langle \pi | \{f(\vec{j})\} \rangle_{\text{sampling}} = \frac{1}{L} \sum_{e=1}^L \left[\underbrace{\langle g_{1e}(1-g_{2e}) \rangle}_{f_e(1-f_e)} + \underbrace{\langle (1-g_{1e})g_{2e} \rangle}_{(1-f_e)f_e} \right]$$

$$= \frac{1}{L} \sum_{e=1}^L 2f_e(1-f_e)$$

Remember:

$$\rightarrow f_e \equiv \sum_{\vec{j}} g_e f(\vec{j})$$

\Rightarrow averaging over f_e :

$$\langle \pi \rangle = \frac{1}{L} \sum_{e=1}^L \langle 2f_e(1-f_e) \rangle = \frac{1}{L} \sum_{e=1}^L \int 2f_e(1-f_e) P(f_e) df_e$$

e.g. if genome is collection of neutral sites

$$\Rightarrow p(f_e) \approx \frac{2N_e N}{f_e} \quad (\text{quasi-stationary dist'n})$$

$$\langle \pi \rangle = \int 2f(1-f) \frac{2N_e N}{f} df = \boxed{2N_e N}^*$$

$$\hookrightarrow \text{fit } N_e \equiv \frac{\langle \pi \rangle}{2N} = \frac{10^{-3}}{2 \times 10^{-8}}$$

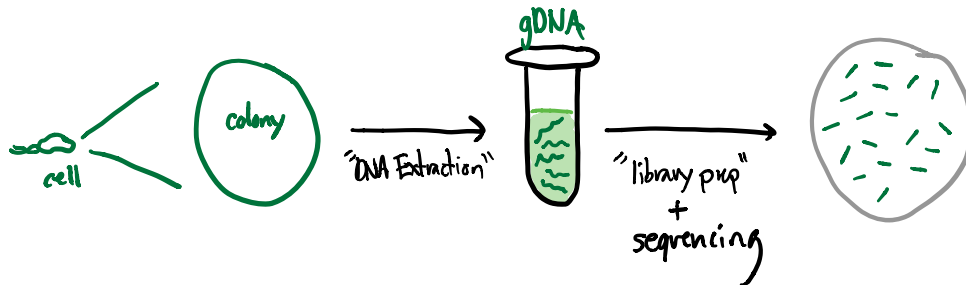
→ population can only spread out so far in sequence space!

⇒ Variance of π ⇒ much more complicated!

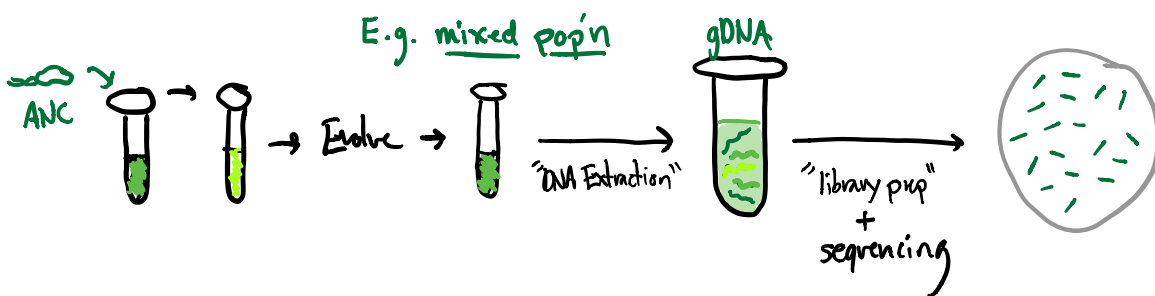
⇒ correlations between $g_e \sim g_{e'}$

⇒ will see more later!

So far, have focused on clones



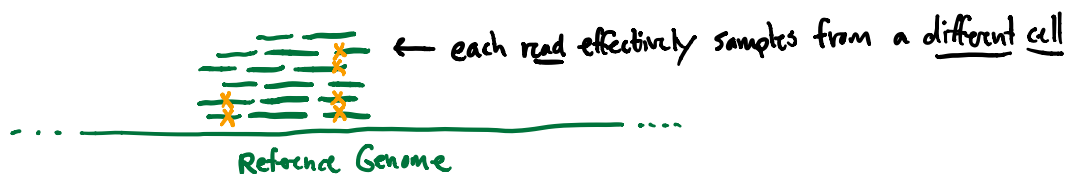
⇒ lots of other things we could put in our library prep...



⇒ known as "pooled sequencing" / "metagenomic sequencing"

⇒ in this case, assembly very hard!

⇒ Payoff comes from reference mapping:



e.g. if $A_e \equiv$ # reads w/ mutations @ site e

$$\Rightarrow \Pr[A_e | D_e, \{f(j)\}] = \text{Binomial}(D_e, f_e) + \text{sequencing/PCR errors}$$

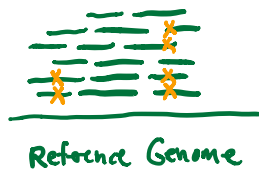
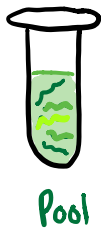
↑
total coverage
@ site e

\Rightarrow since we can sequence ~ 100 E. coli genomes
@ $>100x$ coverage in 1 run of Illumina sequencing

\Rightarrow can effectively sample ~ 100 clones
 $\sim 100x$ more cheaply by sequencing pools!

(much cheaper way to track freqs of individual mut'n's)

Downsides:

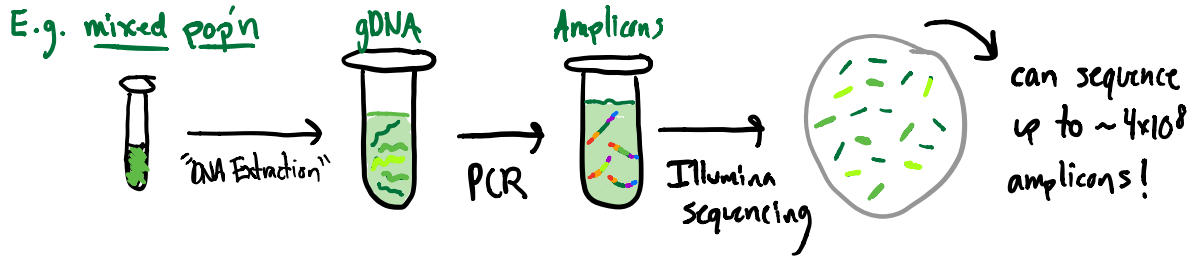


① sequencing errors!

② which mutations are in same cells? ("linkage information")



Can also sequence pools of amplicons :



Two common targets :



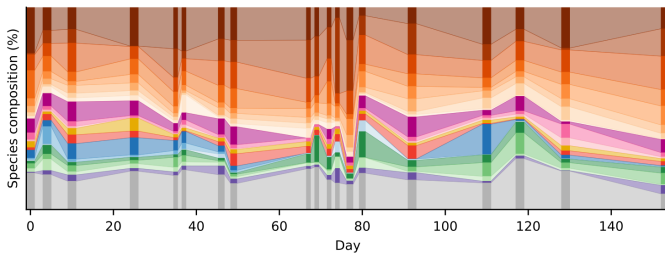
all bacteria have it
 ↳ track abundances of different species



can be added to lab strains
 ↳ track closely related strains



E.g. 1 person's microbiome over time

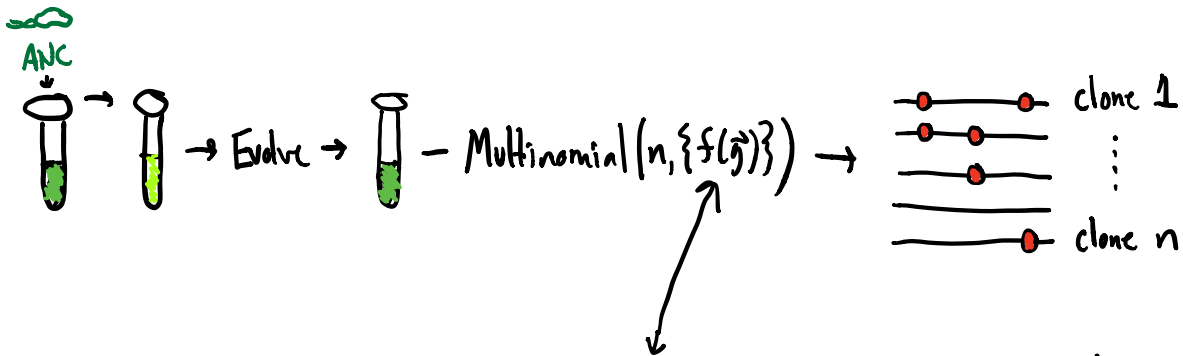


- *Bacteroides vulgatus*
- *Bacteroides coprocola*
- *Bacteroides uniformis*
- *Bacteroides cellulosilyticus*
- *Bacteroides eggerthii*
- *Bacteroides faecis*
- *Bacteroides massiliensis*
- *Bacteroides caccae*
- *Alistipes sp*
- *Alistipes onderdonkii*
- *Alistipes finegoldii*
- *Parabacteroides distasonis*
- *Paraprevotella clara*
- *Butyrivibrio crossotus*
- *Coprococcus sp*
- *Coprococcus comes*
- *Eubacterium rectale*
- *Eubacterium siraeum*
- *Eubacterium eligens*
- *Phascolarctobacterium sp*
- Other



E.g. Homework Problem

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



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

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

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

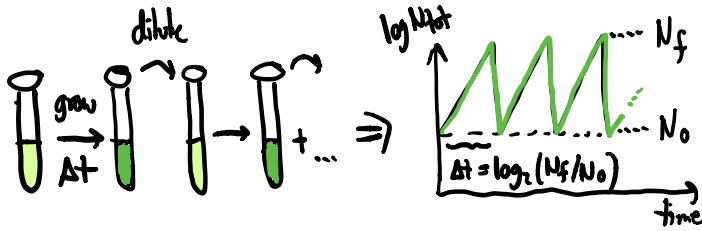
$L=2$: $\vec{g} = \underbrace{(0,0)}_{\text{WT}}, \underbrace{(1,0), (0,1)}_{\text{single mutants}}, \underbrace{(1,1)}_{\text{double mutant}}$

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

⋮

etc.

Can we generalize our serial dilution (& diffusion) models?



$$\frac{df(\vec{g})}{dt} = ???$$

① Genetic drift: first assume no growth rate differences...
(& no mutations)

\Rightarrow After 1 day of growth:

$$f(\vec{g}) \Rightarrow \frac{f(\vec{g})e^{r\Delta t}}{\sum_{\vec{g}'} f(\vec{g}')e^{r\Delta t}} = \frac{f(\vec{g})}{\sum_{\vec{g}'} f(\vec{g}')} = f(\vec{g})$$

\Rightarrow After dilution step:

$$n(\vec{g}, t+\Delta t) \sim \text{Poisson}(N_0 f(\vec{g})) \Rightarrow f(\vec{g}, t+\Delta t) = \frac{n(\vec{g}, t+\Delta t)}{\sum_{\vec{g}'} n(\vec{g}', t+\Delta t)}$$

\Rightarrow if repeat our Taylor expansions from before ($N_0 \gg \text{large}$)

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

where $z_{\vec{g}}$ are $N(0,1)$ random variables w/

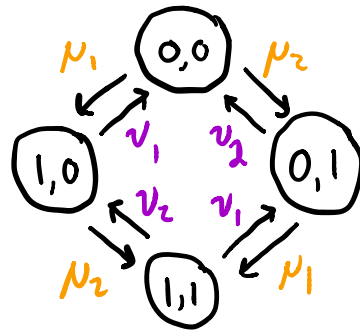
$$\langle z_{\vec{g}} \rangle = 0, \quad \langle z_{\vec{g}} z_{\vec{g}'} \rangle = \begin{cases} 0 & \text{if } \vec{g} \neq \vec{g}' \\ 1 & \text{if } \vec{g}' = \vec{g} \end{cases}$$

\Rightarrow correlations ensure that $f(\vec{g}, t)$ is normalized:

$$\begin{aligned} \sum_{\vec{g}} f(\vec{g}, t + \delta t) &= \sum_{\vec{g}} f(\vec{g}, t) + \sum_{\vec{g}} \sqrt{\frac{f(\vec{g}) \delta t}{N_c}} z_{\vec{g}} - \sum_{\vec{g}} f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}') \delta t}{N_c}} z_{\vec{g}'} \\ &= 1 \end{aligned}$$

② Mutations:

⇒ easiest to start w/ $L=2$:



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

⇒ generalizing results from $L=1$ case, after 1 dilution:

$$n(1,0,t+\Delta t) \sim \text{Poisson} \left(N_0 f(1,0,t) + N_0 \Delta t \left[\overbrace{\mu_1 f(0,0,t) + \nu_2 f(1,1,t)}^{\text{mutations into genotype}} \right] - N_0 \Delta t \left[\underbrace{\mu_2 f(1,0,t) + \nu_1 f(1,0,t)}_{\text{mutations out of genotype}} \right] \right)$$

⇒ continuum limit:

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

+ noise from drift.

↗ linear in * genotype freqs.

⇒ larger L's are similar, but more work to write out...

⇒ one way is:

$$\left[\frac{\delta f(\vec{g})}{\delta t} \right]_{\text{mut}} \equiv \sum_{\substack{\text{nearest} \\ \text{neighbors} \\ \vec{g}'}} \sum_{e=1}^L \left[\overbrace{\mu_e f(\vec{g}') g_e (1-g_e) + \nu_e f(\vec{g}') (1-g_e) g'_e}^{\text{mutations into genotype}} \right] - \sum_{e=1}^L \left[\overbrace{\mu_e f(\vec{g}) (1-g_e) + \nu_e f(\vec{g}) g_e}^{\text{mutations out of genotype}} \right]$$

$$\left(\frac{\delta f(\vec{g})}{\delta t} \right)_{\text{mut}} \equiv \sum_{\vec{g}'} \left[\underbrace{M(\vec{g}' \rightarrow \vec{g})}_{2^L \times 2^L \text{ matrix of mut'n rates}} f(\vec{g}') - \underbrace{M(\vec{g} \rightarrow \vec{g}')}_{2^L \times 2^L \text{ matrix of mut'n rates}} f(\vec{g}) \right]$$

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

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

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:

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

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

$$\Rightarrow f(\vec{g}, t + \Delta t) \approx f(\vec{g}, t) + [X(\vec{g}) - \bar{X}(t)] f(\vec{g}, t) \Delta t$$

\downarrow

population mean fitness

note: not an ensemble avg

$$\langle \bar{X}(t) f(\vec{g}, t) \rangle$$

$$\neq \langle \bar{X}(t) \rangle \langle f(\vec{g}, t) \rangle$$

$$\bar{X}(t) \equiv \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t)$$

$$\text{Normalized: } \sum_{\vec{g}} f(\vec{g}, t + \Delta t) = \sum_{\vec{g}} f(\vec{g}, t) + \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t) - \sum_{\vec{g}} f(\vec{g}) \sum_{\vec{g}'} X(\vec{g}') f(\vec{g}') = 1$$