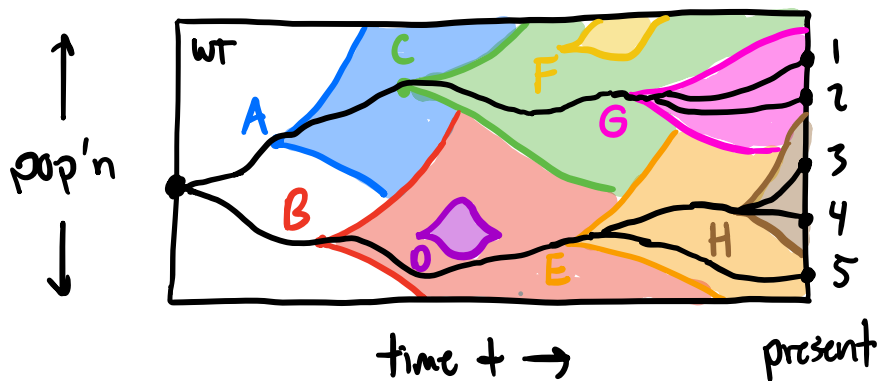# Chapter 15

# Linked selection and clonal interference

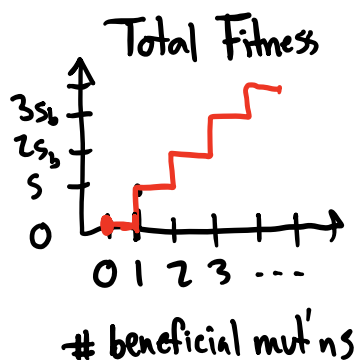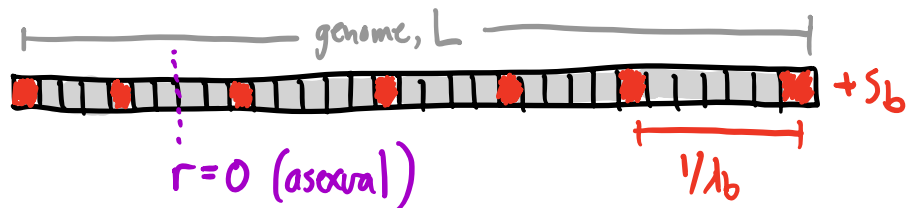# Linked selection & clonal interference (a.k.a. "Hill-Robertson Interference")



=) can't be reduced to $L=1$ or $L=2$ model (collective phase)

=) Most progress only recently, w/ big contribution from physicists

[e.g. Tsimring et al PRL '96, Rouzine et al '03, Desai + Fisher '07, ...]
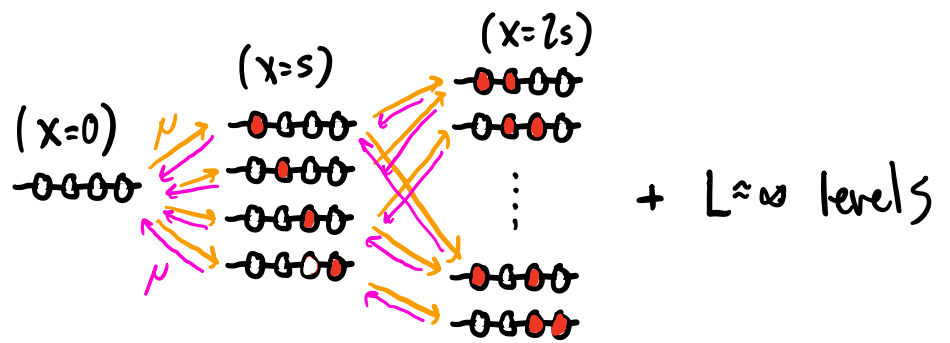
=) Analytical progress enabled by starting w/ very simple model:

## "Staircase" Model



① All mutations provide same benefit ($s_b$)

② Occur @ total rate $U_b \equiv L \lambda_b N$

③ Never run out (e.g. $L \lambda_b \to \infty$, $N \to 0$)

# Genotype network:



$(x=0)$  $(x=s)$  $(x=2s)$

+ $L \simeq \infty$ levels

# Key simplification:

"fitness class"  $f(k,t) \equiv \sum_{|\vec{g}|=k} f(\vec{g},t)$



$(x=0)$  $(x=s)$  $(x=2s)$  $(x=ks)$

$\Rightarrow$ coarse-grained SDE ( 1+1 dimensional  vs  $2^L+1$ dim.)

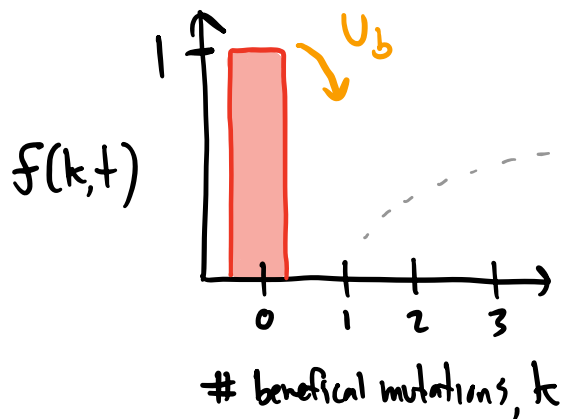$$\frac{\partial f(k)}{\partial t} = \underbrace{s_b(k - \bar{k}(t))f(k)}_{\text{selection (nonlinear)}} + \underbrace{U_b\left[f(k-1) - f(k)\right]}_{\text{mutation}}$$

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

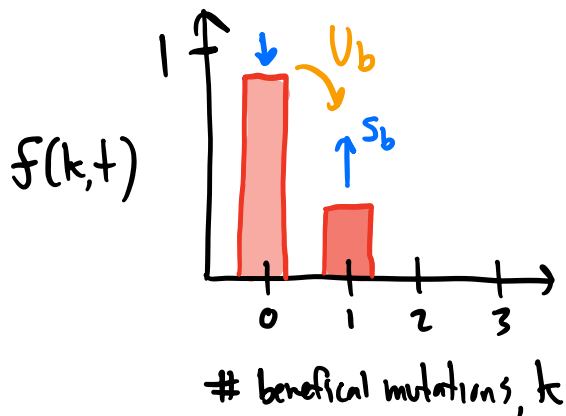$\Rightarrow$ let's consider behavior when $Ns_b \gg NU_b \gg 1$

(e.g. yeast barcode experiment in HW 4 Problem #1)

① Start w/ wildtype population @ $t=0$

$f(k,t)$



$U_b$

0   1   2   3

# beneficial mutations, $k$

$$\frac{df(1)}{dt} \approx sf(1) + U_b + \sqrt{\frac{f(1)}{N}} \eta_1(t)$$

② First-step mutations ($k=1$) establish & grow exponentially

$f(k,t)$



$U_b$

$s_b$

0   1   2   3

# beneficial mutations, $k$

$$\Rightarrow f(1,t) \approx \frac{U_b}{s_b}\left(e^{s_b t} - 1\right)$$

(deterministic approx good @ first, since $NU_b \gg 1$)

③ Double mutants establish before single mutants take over,

$\bar{k}(t)$

$f(k,t)$



$s_b - \bar{k}(t)$

$2s_b - \bar{k}(t)$

0   1   2   3

# beneficial mutations, $k$

$\Rightarrow$ clonal interference!

$$\left(\text{since} \int_0^{\tau_{1/2}} Nf(1,t)\cdot U_b \cdot s_b \, dt \sim NU_b \gg 1\right)$$

$\Rightarrow$ Is deterministic approx still useful?

$$\frac{\partial f(k)}{\partial t} = s_b\left(k - \bar{k}(t)\right)f(k) + U_b\left[f(k-1) - f(k)\right] + \sqrt{\frac{f(k)}{N}}\eta(k) - f(k)\sum_k \sqrt{\frac{f(k')}{N}}\eta(k')$$

selection (nonlinear)      mutation      genetic drift $\to 0$

$\Rightarrow$ can show: $\quad f_{det}(k,t) = \frac{1}{k!}\cdot\left[\frac{U_b}{s_b}\left(e^{s_b t} - 1\right)\right]^k \cdot e^{-\frac{U_b}{s_b}\left(e^{s_b t} - 1\right)}$ $\times$

$\Rightarrow$ Not self-consistent! $\quad\Rightarrow$ Predicts $s_b \bar{k}(t) \approx U_b e^{s_b t}$

(eventually $\underline{all}$ $f(k,t) \ll 1/N$ !)

⇒) Instead, if we <u>simulate</u> model, observe "travelling wave":

$\bar{k}(t) \rightarrow$

$f(k,t)$

fitness class $k$

$U_b$

⟩ after $\tau$ gens

⇓

one "click"

$f(k,t+\tau)$

$s_b$
$2s_b$
$5s_b$
$U_b$

"nose" ≡ $q$ classes above $\bar{k}(t)$

⇒) What determines $\tau(N, s_b, U_b)$ & $q(N, U_b, s_b)$?

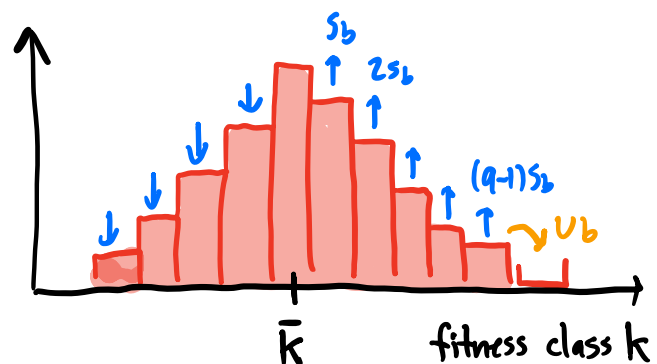⇒) Today: heuristic analysis [~ Desai & Fisher 2007]

applies when: $N s_b \gg N U_b \gg 1$ + $1 \ll s_b^{-1}$ & $q \gg 1$

Leads to simplifications:

① <mark>mutations</mark> only important for establishing new "nose"

(since $s_b \gg U_b$)

$s_b$
$2s_b$
$(q-1)s_b$
$U_b$

$\bar{k}$         fitness class $k$
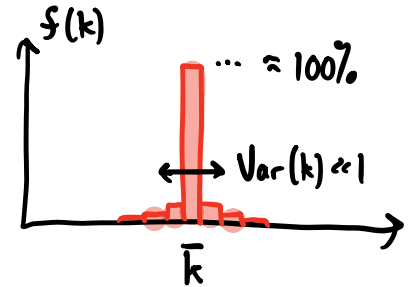
② Genetic drift only important for establishing new nose

(since $\tau \gg 1/_{qs_b}$, individual mutations establish before next click.)

③ most of pop'n is near $k \approx \bar{k}(t)$



Problem 3 of HW 4:

$$\frac{\partial \langle \bar{k} \rangle}{\partial t} = \left\langle \sum_k k \frac{\partial f(k)}{\partial t} \right\rangle = \left\langle \sum_k s_b (k - \bar{k})^2 f(k,t) \right\rangle$$

$\underbrace{\phantom{xxxxxxxx}}$
$\equiv 1/\tau$

$\underbrace{\phantom{xxxxxxxxx}}$
$s_b \, Var(k)$

$$\Rightarrow \quad Var(k) = \frac{1}{s_b \tau} \ll 1 \quad \text{(by assumption)}$$

④ Also implies that $\bar{k}(t)$ clicks suddenly:

$\Rightarrow$ i.e. for most $t \in 0, \tau$ $\Rightarrow$ $\bar{k}(t) = \bar{k}(0)$

$\Rightarrow$ everyone grows as $f(k,t) \sim f(k,0) e^{[k-\bar{k}(0)]st}$

$\Rightarrow$ Now we have all ingredients to understand wave :



$\bar{k}$     fitness class $k$

$\Rightarrow$ in one click $(\tau)$, must

establish new nose

$\Downarrow$

$$f(\bar{k}+q, t) = \frac{1}{Nqs} e^{qs(t-\tau)}$$

$\Rightarrow$ $\boxed{\tau \approx \text{establishment time of nose class !}}$

$f(\bar{k}+q-1, t) = \frac{1}{Nqs} e^{(q-1)st}$

$(q-1)s$

$qs$

$f_1(t) \sim \frac{1}{Nqs} e^{qs(t-\tau_1)}$

$f_2 \sim \frac{1}{Nqs} e^{qs(t-\tau_2)}$

$f_j \sim \frac{1}{Nqs} e^{qs(t-\tau_j)}$

multiple mutant lineages

$\log \left( \frac{1}{Nqs} \right)$

time $t$

$$\Rightarrow f(\bar{k}+q,t) = \sum_{j=0}^{J_{max}} f_j(t) \equiv \frac{1}{Nqs} e^{qs(t-\tau)}$$

establishment time
for whole class.

$\Rightarrow$ $j^{th}$ successful mutant establishes when:

$$\underbrace{\int_0^{\tau_j} NU_b \cdot f_{q-1}(t) \cdot qs_b \, dt \sim O(j)}$$

→ Note: extra little bit will be important below!

$$\int_0^{\tau_j} NU_b \cdot \frac{1}{Nqs} e^{(q-1)st} \cdot qs_b \, dt = \frac{U_b}{qs_b} e^{(q-1)s_b\tau_j} \sim O(k)$$

$$\Rightarrow \quad \boxed{\tau_j = \frac{1}{(q-1)s_b} \log\left( \frac{s_b}{U_b} \cdot q \cdot j \right)}$$

$$\Rightarrow \text{Note:} \quad \tau_j = \underbrace{\frac{1}{(q-1)s_b} \log\left( \frac{s_b}{U_b} \cdot q \right)}_{\tau_1} + \underbrace{\frac{1}{(q-1)s_b} \cdot \log(j)}_{\tau_j - \tau_1}$$

$$\tau_1 \gg \tau_j - \tau_1$$

(most time spent waiting for first mut'n)

$\Rightarrow$ <u>many</u> mutations establish in quick succession $\left( \delta t \sim \frac{1}{qs_b} \ll \tau \right)$

$\Rightarrow$ Typical size of $j^{th}$ lineage:

$$\Rightarrow \quad f_j(t) \sim \frac{1}{Nqs} e^{qs(t-\tau_j)} = \frac{e^{qst}}{Nqs} \left( \frac{s_b q_j}{U_b} \right)^{-1-\frac{1}{q-1}}$$

<span style="color:red">$\uparrow$<br>extra bit<br>will be<br>important!</span>

$\Rightarrow$ Size of entire nose class:

$$f(\bar{k}+q, t) = \sum_{j=1}^{J_{max}} f_j(t) = \frac{1}{Nqs} e^{qst} \left( \frac{s_b}{U_b} \right)^{-\frac{q}{q-1}} \sum_{j=1}^{J_{max}} \frac{1}{(q \cdot j)^{1+\frac{1}{q-1}}} \quad \nearrow 1$$

$$\equiv \frac{1}{Nqs} e^{qs(t-\tau)} \qquad \text{(set equal!)}$$

$\Rightarrow$ <span style="background:#7eb8e8">Time to establish <u>new</u> nose:</span> $\quad \tau = \frac{1}{(q-1)s} \log\left( \frac{s_b}{U_b} \right)$

$$\text{vs} \quad \tau_j \equiv \frac{1}{(q-1)s_b} \log\left| \frac{s_b}{U_b} \cdot q \cdot j \right)$$

$\left[ \begin{array}{l} \text{Note: } \tau < \tau_j \text{ b.c.} \\ \text{multiple mutations} \\ \text{contribute } \rho \text{ once} \end{array} \right]$

<u>One task remaining</u>... how to determine $q(N, s_b, U_b)$?

$\Rightarrow$ follow new nose over time:

$$f(\bar{k}+q, \tau) \approx \frac{1}{Nqs} \xrightarrow{\tau} \frac{1}{Nqs} e^{(q-1)s\tau} \xrightarrow{\tau} \frac{1}{Nqs} e^{(q-1)s\tau + (q-2)s\tau} \rightarrow \cdots$$

(right after est.)

$\Rightarrow$ **After $q$ clicks, old nose is new mean!** (majority of pop'n)

$$f(q\tau) \sim \frac{1}{Nqs_b} e^{(q-1)s_b^2\tau + (q-2)s_b^2\tau + \cdots + s_b^2\tau} \sim \frac{1}{Nqs_b} e^{\frac{q^2 s_b^2 \tau}{2}} \sim \mathcal{O}(1)$$

$\Rightarrow$ system of 2 eqs for $\tau$ & $q$!

$$\frac{q^2 s_b \tau}{2} \approx \log(Ns_b) \quad + \quad \tau = \frac{1}{q s_b} \log\left(\frac{s_b}{U_b}\right)$$

$\Rightarrow$ solution: $\quad q = \dfrac{2\log(NS_b)}{\log\left(\frac{S_b}{U_b}\right)} \quad ; \quad \tau = \dfrac{1}{2S_b}\dfrac{\log^2\left(\frac{S_b}{U_b}\right)}{\log(NS_b)}$

$$\Rightarrow \quad \left\langle \dfrac{d\bar{X}}{dt} \right\rangle = \dfrac{S_b}{\tau} = \dfrac{2S_b^2\log(NS_b)}{\log^2(S_b/U_b)}$$

$\left(\text{compare to } \sim NU_b S_b^2 \text{ in } \textit{successive mutations} \text{ regime}\right)$
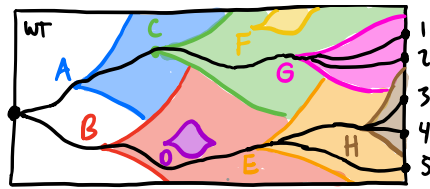
$\Rightarrow$ Self consistency: $\quad S\tau \gg 1 \quad + \quad q \gg 1$

$$\Rightarrow \quad \log\left(\frac{S_b}{U_b}\right) \ll \log(NS_b) \ll \log^2\left(\frac{S_b}{U_b}\right)$$

Note: used heuristic derivation here...
for formal analysis (using branching processes)
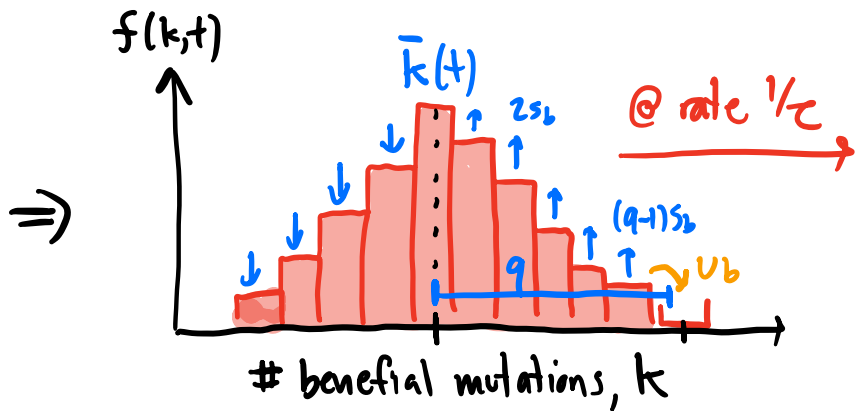see Appendix A and B below
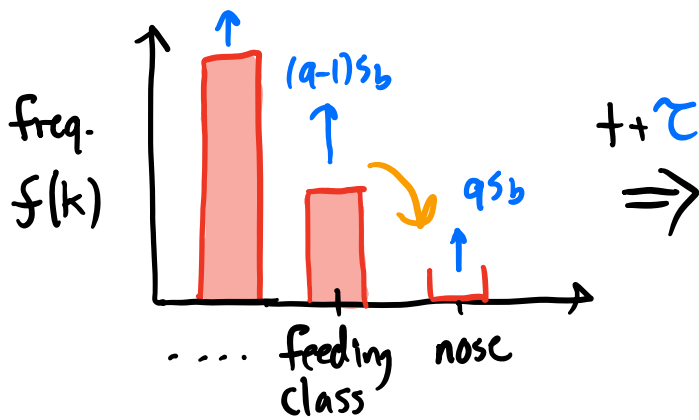
# Recap : clonal interference
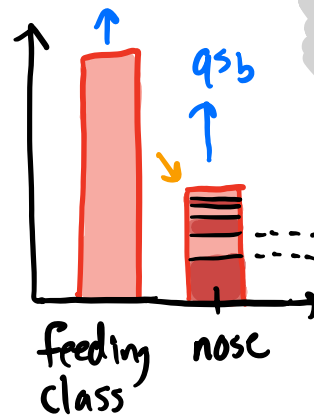


## "Staircase" model

← Genome, $L \gg 1$ →

↳ selected mut'ns, $+S_b$

total rate $U_b \equiv L\lambda_b N$

⇒

$f(k,t)$

$\bar{k}(t)$   $2s_b$   @ rate $1/\tau$

$(q-1)s_b$

$q$   $U_b$

\# benefial mutations, k

## key behavior occurs @ "nose":

freq. $f(k)$

$(q-1)s_b$   $qs_b$

.... feeding class   nose

$+\!+\tau$ ⇒

$qs_b$

feeding class   nose

**Multiple mutations contribute to nose!**
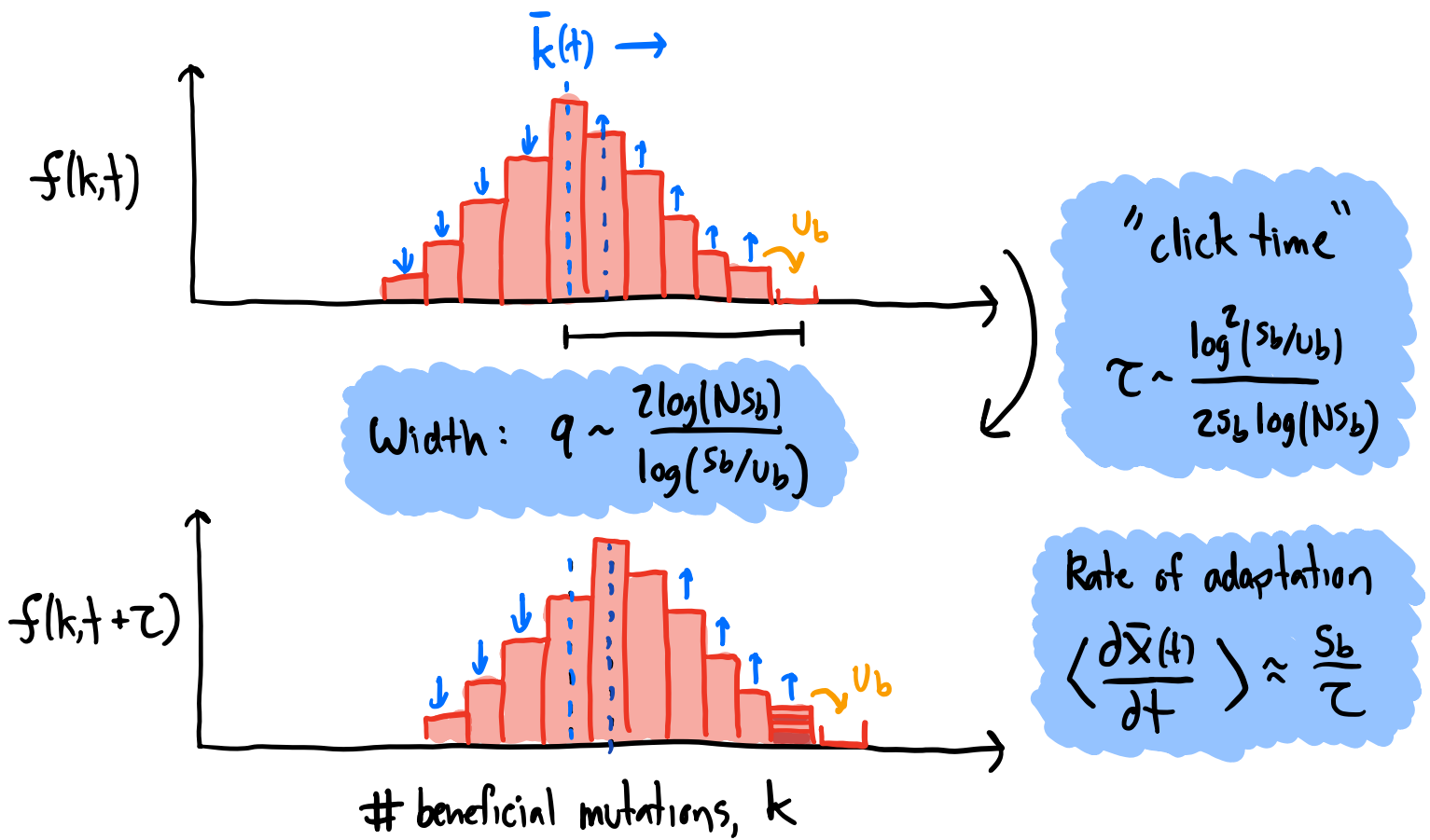
$$f_j(t) \sim \frac{1}{Nqs} e^{qs(t-\tau_j)}$$

$$\tau_j \sim \frac{1}{(q-1)s_b} \log\left(\frac{s_b \, qj}{U_b}\right)$$

## Total contribution:

$$f_{nose}(t) \equiv \sum_{j=1}^{\infty} f_j(t) \equiv \frac{1}{Nqs} e^{qs(t-\tau)} \quad\Rightarrow\quad \tau \sim \frac{1}{(q-1)s_b} \log\left(\frac{s_b}{U_b}\right)$$

$\Rightarrow$ Complete picture of dynamics of <u>fitness dist'n</u> :

$\overline{k}(t) \rightarrow$

$f(k,t)$



$u_b$

Width: $q \sim \dfrac{2\log(Ns_b)}{\log(s_b/u_b)}$

"click time"

$$\tau \sim \frac{\log^2(s_b/u_b)}{2s_b \log(Ns_b)}$$

$f(k,t+\tau)$



$u_b$

# beneficial mutations, k

Rate of adaptation

$$\left\langle \frac{\partial \overline{x}(t)}{\partial t} \right\rangle \approx \frac{s_b}{\tau}$$
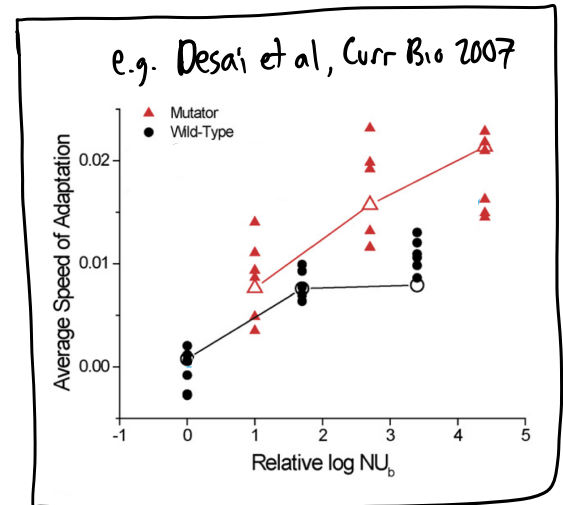
$\Rightarrow$ early tests for clonal interference in lab evolution experiments :
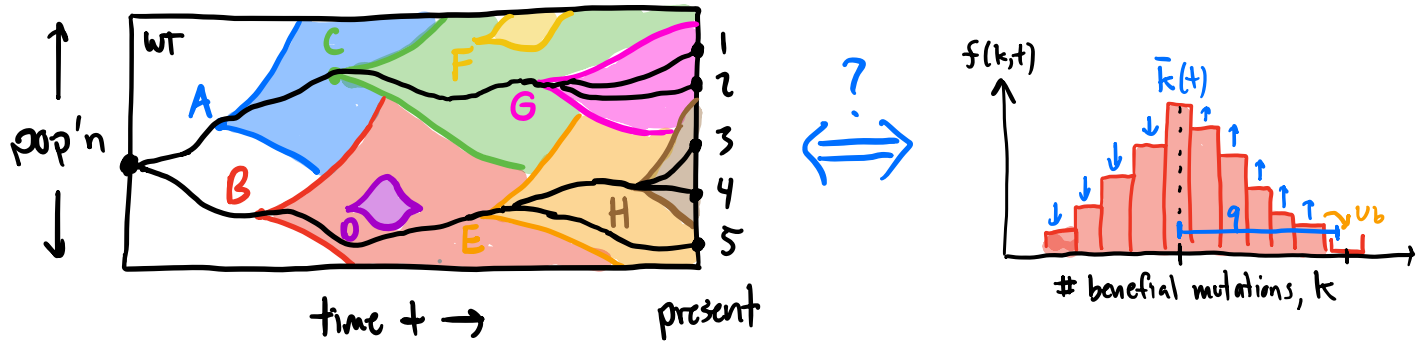
Successive mutations:

$$\left\langle \frac{d\overline{x}}{dt} \right\rangle \sim s_b^2 \cdot Nu_b$$

clonal interference :

$$\left\langle \frac{d\overline{x}}{dt} \right\rangle \sim s_b^2 \cdot \frac{\log(Ns_b)}{\log^2(s_b/u_b)}$$



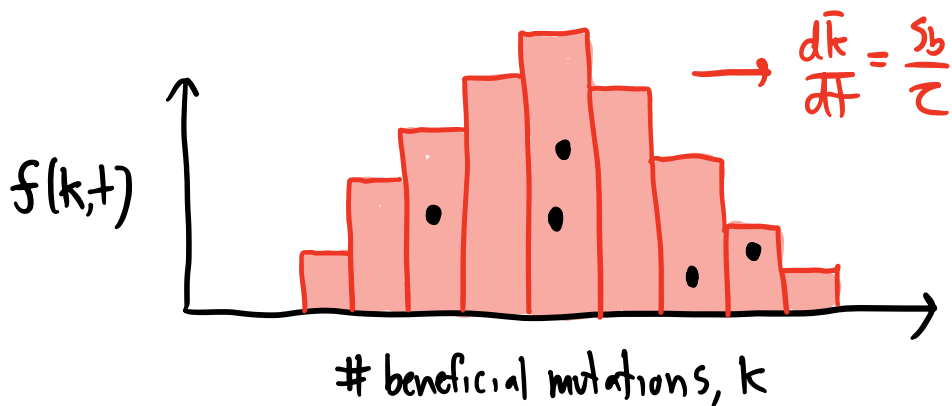e.g. Desai et al, Curr Bio 2007

Next: Can we use this picture to understand _genetic diversity_ backwards in time?



Answer: Yes we can! Let's start w/ some cartoons...

Step 1: draw sample of individuals from pop'n (present day)



$$\frac{d\bar{k}}{dt} = \frac{s_b}{\tau}$$

# Step 2: where was everyone <u>one</u> <u>click</u> ago?

$f(k,t)$

→ forward time

one click ago

reverse time

$f(k,t-\tau)$

# beneficial mutations, k

① can only coalesce if in same fitness class

② <u>But</u> little chance of coalescing in "bulk" of dist'n

(since $\tau \ll Nf_{q-1}(\tau)$, $Nf_{q-2}(\tau)$, etc.)

$f(k,t)$

one click ago

reverse time

$f(k,t-\tau)$

two clicks ago

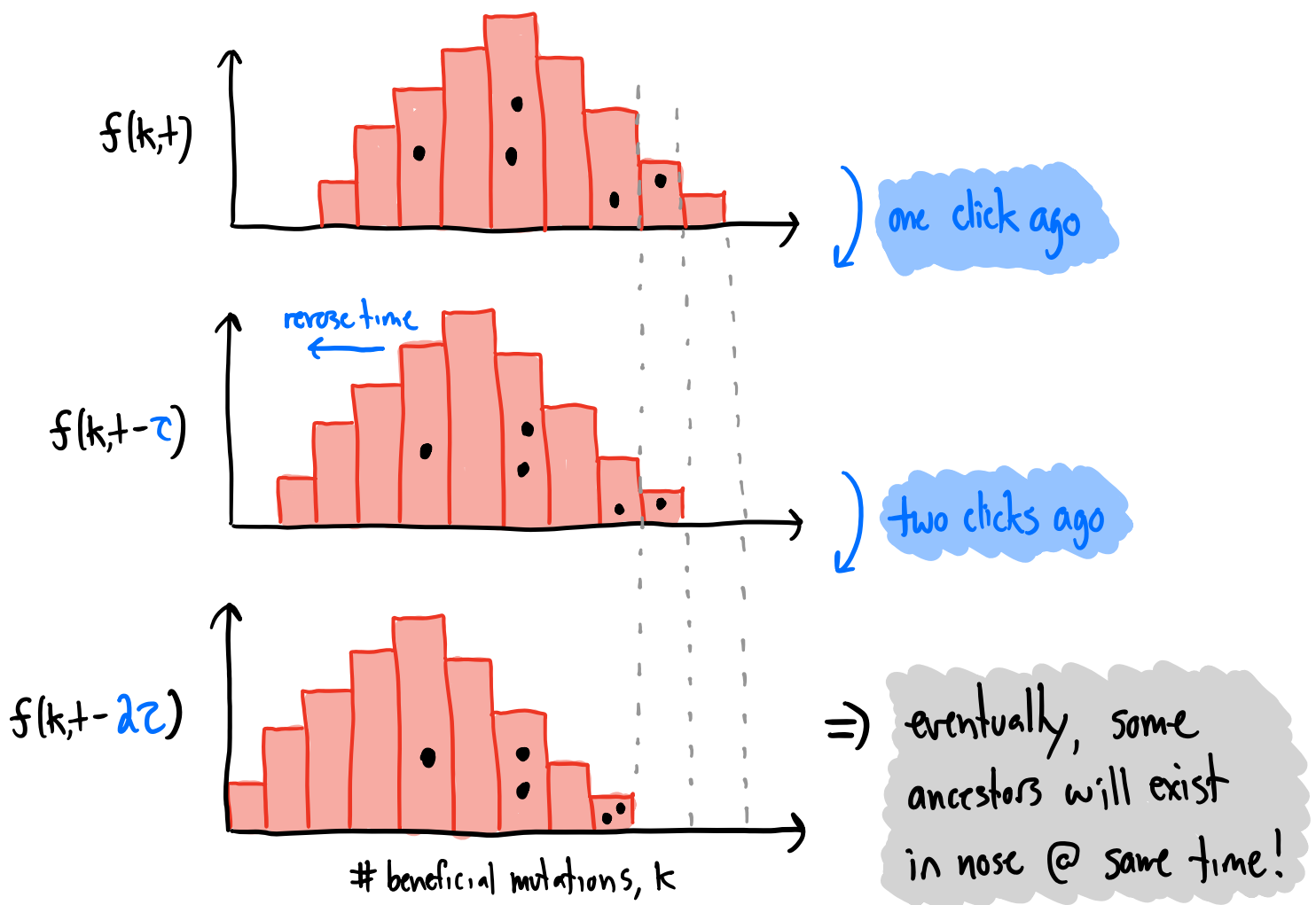$f(k,t-2\tau)$

# beneficial mutations, k

$\Rightarrow$ eventually, some ancestors will exist in nose @ same time!

Two possible scenarios:

① Individuals are from separate lineages in the nose

$U_b$

$f_{nose}(t)$

$f_j(t)$

feeding class    nose

$\Downarrow$

separate ancestors in feeding class

(distinct mut'n events)

② Individuals from <u>same</u> lineage in nose

$\Downarrow$

common ancestor
in feeding class

(coalescence w/in $\tau$ gens)

$$\Rightarrow \text{Probability: } p_c(2) = \sum_{j=1}^{\infty} \left( \frac{f_j(t)}{f_{nose}(t)} \right)^2 = \sum_{=1}^{\infty} \left[ \frac{\frac{1}{Nqs_b} e^{qs_b(t-\tau_j)}}{\frac{1}{Nqs_b} e^{qs_b(t-\tau)}} \right]^2$$

$$= \sum_{j=1}^{\infty} e^{-2qs_b(\tau_j - \tau)}$$

$\Uparrow$

only depends on establishment times $\tau_j$ !

⇒ if we plug-in $\boxed{\text{typical}}$ values of $\tau_j$ & $\tau$ from heuristics:

$$\tau_j \sim \frac{1}{(q-1)s_b} \log\left(\frac{s_b \, q \, j}{U_b}\right) ; \qquad \tau \sim \frac{1}{(q-1)s_b} \log\left(\frac{s_b}{U_b}\right)$$

$$\Rightarrow \quad P_c(2) = \sum_{j=1}^{\infty} e^{-2qs_b(\tau_j - \tau)} = \sum_{j=1}^{\infty} (qj)^{-\frac{2q}{q-1}} \approx \frac{1}{q^2}$$

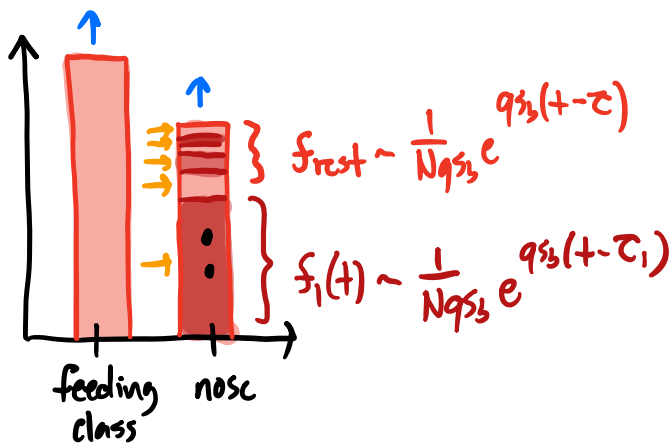⇒ suggests coalescence after $\sim q^2$ clicks $\left(T_{MRCA} \sim q^2 \tau\right)$

⇒ missing key part of puzzle : fluctuations

⇒ coalescence rare for $\boxed{\text{typical}}$ lineage sizes,
but small chance of having anomalously early mutant
where coalescence is much more likely!

e.g. if **first** successful mutation occurs when $\tau_1 \lesssim$ typical $\tau$ ...



$$f_{rest} \sim \frac{1}{Nqs_b} e^{qs_b(t-\tau)}$$

$$f_1(t) \sim \frac{1}{Nqs_b} e^{qs_b(t-\tau_1)}$$

feeding class    nosc

$$\frac{f_1(t)}{f_{nose}(t)} \gtrsim \mathcal{O}(1)$$

$$P_c(2) \sim \mathcal{O}(1) !$$

$\Rightarrow$ not a huge shift in **time**: typically, $\tau_1 - \tau \sim \frac{\log(q)}{qs_b}$

so $\Delta\tau_1$ is $\ll \frac{1}{s_b} \ll \tau$ (i.e. $\ll$ click time)

$\Rightarrow$ occurs w/ total probability:

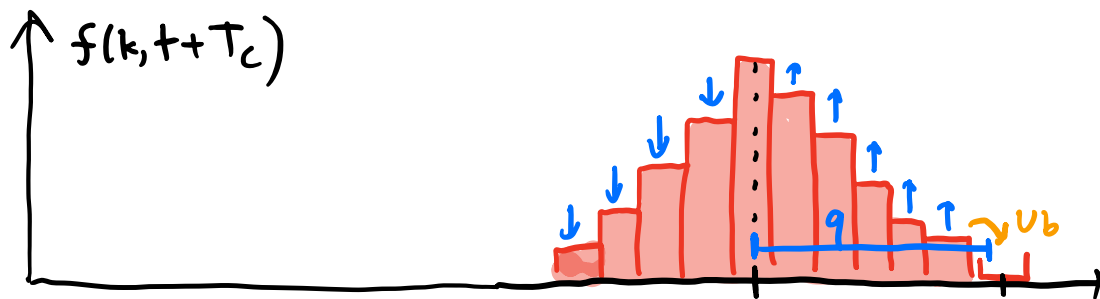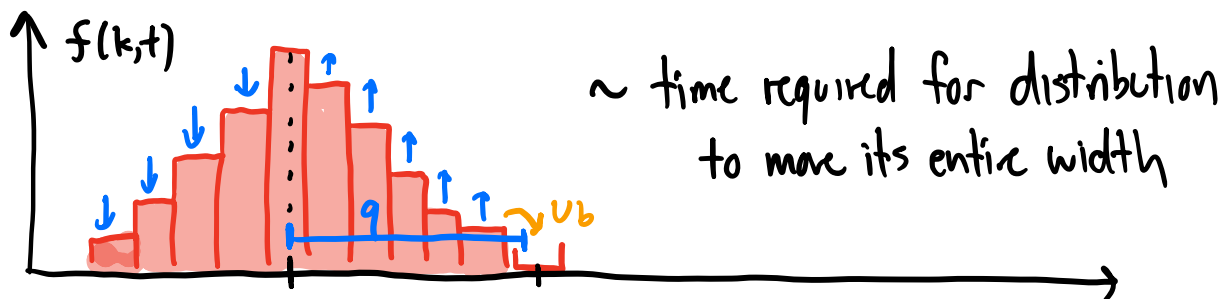$$P_{jackpot} \sim \int_0^\tau d\tau_1 \, NU_b f_{q-1}(t) \cdot qs_b$$

$$\sim \int_0^\tau d\tau_1 \, \cancel{N}U_b \cdot \frac{e^{(q-1)s_b t}}{\cancel{N}qs_b} \cdot \cancel{qs_b} \sim \frac{U_b}{(q-1)s_b} e^{(q-1)s_b\tau}$$

$$\sim \frac{1}{q}$$

$\Rightarrow$ $P_{jackpot} \sim 1/q$ $\left( \gg 1/q^2 \Rightarrow \begin{array}{l} \text{more likely to coalesce} \\ \text{via } \underline{\text{rare}} \text{ jackpot than} \\ \text{normal establishment process} \end{array} \right)$

$\Rightarrow$ typical coalescence after $1/P_{jackpot} \sim q$ clicks

$\Rightarrow$ coalescent timescale $T_c \equiv q\tau \sim \frac{1}{s_b} \log\left(\frac{s_b}{U_b}\right)$



$\sim$ time required for distribution to move its entire width

# benefial mutations, k

$\Rightarrow$ fluctuations were crucial for determining $T_c$ !

$\Rightarrow$ coalescence is "bursty":

e.g. in larger sample size $n$:

$$P_c(n \to 1) = \left( \frac{f_1(t)}{f_1(t) + f_{rest}(t)} \right)^n$$



$$\approx \begin{cases} \sim 1 & \text{if } f_1(t) \gtrsim n \cdot f_{rest}(t) \\ \ll 1 & \text{else} \end{cases}$$

$\Rightarrow$ $$P_{jackpot}(n) = \int_0^{\tau - \log(n)/q s_b} d\tau_1 \; N U_b f_{q-1}(t) \cdot q s_b \sim \frac{1}{qn}$$

$\Rightarrow$ i.e. multiple mergers likely! 

\* For "formal" treatment, see Appendix C ...

Another interesting feature of genealogies + travelling wave:

⇒ consider same
example:



# beneficial mutations

⇒ which individual's <u>descendents</u> are more likely
to take over pop'n in <u>future</u>?

⇒ e.g. $5 \to 4 \to 2,3 \to 1$

⇒ now let's try to "simulate" genealogy...

-3τ

-2τ

-τ

present

relative fitness, k-k̄(t)

backward in time

MRCA

1 2 3 4 5

3τ
2τ
τ

⇒ time (+burstiness) of coalescence in past

⇒ info about fitness in present

⇒ forecasts about who takes over in future!

# Predicting evolution from the shape of genealogical trees

Richard A Neher[1]*, Colin A Russell[2], Boris I Shraiman[3]*

[1]Evolutionary Dynamics and Biophysics, Max Planck Institute for Developmental Biology, Tübingen, Germany; [2]Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom; [3]Kavli Institute for Theoretical Physics, University of California, Santa Barbara, Santa Barbara, United States

⟹ implemented this idea for HA gene in influenza
(data from Problem #1 in HW1)



e.g. flu strains from 2006/07

("local branching index, LBI")

worst — best

highest inferred "fitness" (LBI)

e.g. flu strains from 2006-2007

+ 2007-2008

← new strains

best guess
from last season

0.005

Evaluating performance over multiple flu seasons:

Predictions

Random guess

$\Delta$(prediction) to next season

1.0

0.5

1995    1997    1999    2001    2003    2005    2007    2009    2011    2013

year

Best guess
(in hindsight)

# What about recombination?

recombination
rate $r$

↳ selected mut'ns, $+s_b$

total rate $U_b \equiv L\lambda_b N$

$r = 0$

$\Rightarrow$

$\bar{k}(t)$

$q$

$\sim U_b$

\# beneficial mutations, $k$

$r \gg N \cdot \mu\lambda_b \cdot s$

⊠ ··· ⊠

$\Rightarrow$

WT

A

B

C

time →

$r = 0$

$L$

$K(t)$

$q$ $\sim U_b$

# benefial mutations, $k$

In between?

$\leftarrow \ell^* \rightarrow$

$\otimes \quad \otimes \cdots \otimes$

$U_{b,eff} = \ell^* \lambda_b \mu$ ?

Blocks are ~ independent

... but multiple mut'ns / block!

$r \gg N \cdot \mu \lambda_b \cdot s$

$\otimes \cdots \otimes$

WT

A    B    C

time $\rightarrow$

If true, need:

① w/in blocks, recombination should be rare! ($r \hat{=} 0$)

$$\Rightarrow \quad r\ell^* \cdot T_c(N, s_b, U_{eff}(\ell^*)) \ll 1$$

② between blocks, recombination should be frequent!

$$\Rightarrow \quad r\ell^* \cdot T_c \gg 1$$

$\Rightarrow$ can we (almost) satisfy **both** w/ $r\ell^* \cdot T_c \sim O(1)$ ?

# Linkage block ansatz



$$\ell^* \sim {}^1/_{rT_c}$$

$$U_{b,eff} = \ell^* \lambda_b \mu$$

$$T_c \sim \frac{1}{s_b} \log\left(\frac{s_b}{U_{b,eff}}\right)$$

# benefial mutations, k

$\Rightarrow$ Self consistency: $\quad T_c \sim \frac{1}{s_b} \log\left(\frac{s_b}{\mu\lambda_b} \cdot rT_c\right)$

$\Rightarrow$ solution: $\quad T_c \sim \frac{1}{s_b} \log\left(\frac{r}{\mu\lambda_b}\right)$

$\Rightarrow \quad \ell^* \sim \frac{s_b}{r} \log^{-1}\left(\frac{r}{\mu\lambda_b}\right)$

$\Rightarrow$ self consistent if $\quad NU_{b,eff} \log(Ns_b) \gg 1 \quad \& \quad U_{eff,b} \ll s_b$

$\Rightarrow \quad N \cdot N\lambda_b \cdot s_b \gg r \gg \mu\lambda_b$

# Appendix A:   Formal analysis of the nose class

$\Rightarrow$ we can understand the establishment of the nose class more formally using the branching process framework that we studied in the 1st half of the course



$\Rightarrow$ Under our assumptions, nose can be described by LBP model:

$$\frac{df_q}{dt} = X_c(t) f_q + U_b f_{q-1}(t) + \sqrt{\frac{f_q}{N}} \eta(t)$$

w/ $f_q(0) = 0$ & time-varying:

selection:   $X_c(t) = \left[ q - \bar{k}(t) \right] s_b$

$+$

mutation:   $U_b f_{q-1}(t) = \frac{U_b}{2Nqs_b} e^{\int_0^t (X_c(t') - s_b) dt}$

$\Rightarrow$ In their analysis, Desai & Fisher (2007) assumed that $\bar{k}(t) \approx 0$ throughout the establishment period, so that $X_c(t) \equiv q s_b$ & $f_{q-1}(t) = \frac{1}{Nqs} e^{(q-1)s_b t}$

$\Rightarrow$ Let's see how far we can get by <u>relaxing</u> this approx & explicitly modeling the "click" of $\bar{k}(t)$...

$\quad \Rightarrow$ will be harder because time-varying fitness

$$X_c(t) = q s_b - s_b \bar{k}(t)$$

$\Rightarrow$ From our discussion in class, can take

$$\bar{k}(t) \equiv \frac{e^{s_b(t-t_c)}}{1 + e^{s_b(t-t_c)}}$$

where $t_c$ is the time that $\bar{k}(t)$ clicks.

( later we will imagine that $t_c$ is close to $\tau$ ...)

$\Rightarrow$ From SDE, the generating function $H_f(z,t) = \left\langle e^{-z \cdot f_q(t)} \right\rangle$
satisfies the PDE:

$$\frac{\partial H_f}{\partial t} = \left[ X_c(t)\, z - \frac{z^2}{2N} \right] \frac{\partial H_f}{\partial z} - z\, v_b f_{q-1}(t)\, H_f$$

$$\text{w/ initial condition } H_f(z,0) = 1$$

$\Rightarrow$ can solve w/ method of characteristics:

$$\text{define}: \quad \mathcal{Y}(t_R) = \log\left[ H_f\left(z(t_R),\, t-t_R\right) \right]$$

$$\text{w/ } \mathcal{Y}(t) = 0, \quad z(0) = z,$$

$$\mathcal{Y}(0) \equiv \log H_f(z,t)$$

$\Rightarrow$ $\mathcal{Y}$ satisfies: $\quad \dfrac{d\mathcal{Y}}{dt_R} = -\dfrac{\frac{\partial H_f}{\partial t}}{H_f} + \dfrac{\frac{\partial H_f}{\partial z}}{H_f}\left(\dfrac{\partial z}{\partial \tau_R}\right)$

$$\Rightarrow \quad \frac{\partial y}{dt_R} = \left\{ \frac{dz}{dt_R} - \left[ x_c(t-t_R)z - \frac{z^2}{2N} \right] \right\} \frac{\partial \log H_f}{\partial z} + z(t_R) U_b f_{q-1}(t-t_R)$$

$$\Rightarrow \quad \text{if} \quad \frac{dz}{dt_R} = x_c(t-t_R)z - \frac{z^2}{2N} \quad \& \quad z(0) = z$$

$$\Rightarrow \quad y(t_R) = y(0) + \int_0^{t_R} z(t_R') U_b f_{q-1}(t-t_R) \, dt_R'$$

$$\Rightarrow \quad \log H_f(z,t) = - \int_0^t z(\tau) U_b f_{q-1}(t-\tau) \, d\tau$$

where

$$\frac{dz}{d\tau} = x_c(t-\tau)z - \frac{z^2}{2N}, \quad z(0) = z$$

$\Rightarrow$ solution for characteristic curve is given by:

$$z(\tau) = \frac{z \, e^{\int_0^\tau x_c(t-\tau') \, d\tau'}}{1 + \frac{z}{2N} \int_0^\tau e^{\int_0^{\tau'} x_c(t-\tau'') \, d\tau''} \, d\tau'}$$

(can plug in & check...)

So
$$H_f(z, t) = \exp\left[ -\int_0^t \frac{z \, U_b f_{q-1}(t-\tau) \, e^{\int_0^\tau x_c(t-\tau')d\tau'}}{1 + \frac{z}{2N}\int_0^\tau e^{\int_0^{\tau'} x_c(t-\tau'')d\tau''} d\tau'} \, d\tau \right]$$

$$= \exp\left[ -\int_0^t \frac{z \cdot U_b f_{q-1}(u) \, e^{\int_u^t x_c(u')du'}}{1 + \frac{z}{2N}\int_u^t e^{\int_{u'}^t x_c(u'')du''} du'} \, du \right]$$

$\Rightarrow$ again, helpful to define $\nu(t)$ s.t. $f_q(t) \equiv \frac{\nu(t)}{2Nqs_b} e^{\int_0^t x_c(t')dt'}$

$\Rightarrow$ $H_\nu(z, t) \equiv \left\langle e^{-z \cdot \nu(t)} \right\rangle \equiv H_f\left( 2Nqs_b \, e^{-\int_0^t x_c(t')dt'} z, t \right)$

$\Rightarrow$ $H_\nu(z, t) = \exp\left[ -\int_0^t \frac{z \, U_b f_{q-1}(u) \, 2Nqs_b \, e^{-\int_0^u x_c(u')du'}}{1 + qsz \cdot \int_u^t du' \, e^{-\int_0^{u'} x_c(u'')du''}} \, du \right]$

$\Rightarrow$ Similar to single-locus case, we expect $\nu(t)$
to approach constant value $\nu$ @ long times

$$\Rightarrow \quad H_v(z) \equiv \lim_{t \to \infty} H_v(z,t)$$

$$\Rightarrow \quad \log H_v(z) = -\int_0^\infty \frac{z \cdot U_b e^{-s_b t} \, dt}{1 + z \cdot q s_b \int_t^\infty dt' \, e^{-\int_0^{t'} x_c(t'') dt''}}$$

$\Rightarrow$ Now we have to plug in our expression for $X_c(t)$:

$$X_c(t) = q s_b - \frac{s_b e^{s_b(t-t_c)}}{1 + e^{s_b(t-t_c)}}$$

$$\Rightarrow \quad \int_0^t x_c(t') dt' = q s_b t - \log\left[\frac{1 + e^{s_b(t-t_c)}}{1 + e^{-s_b t_c}}\right]$$

$$\Rightarrow \quad e^{-\int_0^t x_c(t') dt'} = e^{-q s_b t}\left[\frac{1 + e^{s_b(t-t_c)}}{1 - e^{-s_b t_c}}\right]$$

$$\Rightarrow \quad q s_b \int_t^\infty dt' \, e^{-\int_0^{t'} x_c(t'') dt''} = \frac{e^{-q s_b t}}{1 + e^{-s_b t_c}} + \left(\frac{q}{q-1}\right) \frac{e^{-q s_b t} \, e^{s_b(t-t_c)}}{1 + e^{-s_b t_c}}$$

and hence:

$$\log H_v(z) = - \int_0^\infty \frac{z \cdot U_b e^{-S_b t} \, dt}{1 + z \cdot e^{-q S_b t} \left[ 1 + e^{S_b(t - t_c)} \left( \frac{q}{q-1} \right) \right]}$$

$\left( \text{where we have assumed that the click time } t_c \text{ is } \gg \frac{1}{S_b} \right)$

$\Rightarrow$ for large $q$ & <u>relevant</u> values of $z$, this integral
will be dominated by times w/in $O\left(\frac{1}{S_b}\right)$ of $\tau$.

$\Rightarrow$ can extend lower limit of integral to $t = -\infty$
w/o much error...

$\Rightarrow$ if $t_c$ is also w/in $O\left(\frac{1}{S_b}\right)$ of $\tau$, we can expand
$e^{S_b(t - t_c)}$ term in denominator, so that

$$\log H_v(z) \approx - \int_{-\infty}^\infty \frac{z \cdot U_b e^{-S_b t} \, dt}{1 + 2 \cdot z \cdot e^{-q S_b t}}$$

changing variables to $\xi = (2z)^{1/q} e^{-s_b t}$,

$$\log H_v(z) = \exp\left[ -\frac{U_b}{s_b} \cdot z^{1-\frac{1}{q}} \cdot \left( 2^{-\frac{1}{q}} \int_0^\infty \frac{\xi \, d\xi}{1+\xi^q} \right) \right]$$

$$\nearrow 1 + O\left(\frac{1}{q}\right)$$

$$\Rightarrow \quad H_v(z) = e^{-\frac{U_b}{s_b} z^{1-\frac{1}{q}}}$$

$$\Rightarrow \text{ typical value of } v \text{ occurs when } H_v\left(z = \frac{1}{v^*}\right) = e^{-1}$$

$$\Rightarrow \quad v^* = \left( \frac{s_b}{U_b} \right)^{\frac{q}{q-1}}$$

$$\Rightarrow \text{ substituting into } f_q(t) = \frac{v}{2Nqs} e^{qst} \equiv \frac{e^{qs(t-\tau)}}{2Nqs}$$

$$\Rightarrow \text{ typical value of } f_q^*(t) = \frac{e^{s_b t}}{Nqs_b} \left( \frac{s_b}{U_b} \right)^{\frac{q}{q-1}}$$

⇒ typical value of establishment time:

$$\tau^* = \frac{1}{(q-1)s_b} \log\left(\frac{s_b}{U_b}\right)$$

⇒ consistent w/ results from

simpler heuristic argument!

# Appendix B: How many lineages contribute to new nose?

Recall in heuristic argument, we had:

$$f_{nose}(t) \equiv \sum_{j=1}^{J_{max}} f(t) = \frac{1}{Nqs} e^{qs_b(t-\tilde{\tau})} \cdot \sum_{j=1}^{J_{max}} \frac{1}{qj}^{1+\frac{1}{q}}$$

& argued that sum over $k$ converged to $\approx 1$.

$\Rightarrow$ Let's look @ this more carefully...

$\Rightarrow$ if $J_{max} \gg 1$ (will revisit below)

$$\Rightarrow \sum_{j=1}^{J_{max}} \frac{1}{qj}\frac{1}{j^{1+\frac{1}{q}}} \simeq \int_1^{max} \frac{dk}{qj^{1+\frac{1}{q}}} = 1 - e^{-\frac{1}{q}\log J_{max}}$$

Thus, sum converges to 1 provided that
$\log J_{max}$ is large compared to $q$

$\Rightarrow$ how does this translate to establishment times $\tilde{\tau}_j$?

recall that $\tau_j - \tau_1 \sim \frac{1}{q s_b} \log\left(j\right)$, so condition becomes:

$$\Rightarrow \quad \tau_{jmax} - \tau_1 \sim \frac{1}{q s_b} \log\left(J_{max}\right) \gg \frac{1}{s_b}$$

Thus, mutations that establish $\gg \frac{1}{s_b}$ after $\tau_1$ have <u>negligible</u> contribution to $f_{nose}(t)$, $\tau$, etc

$\Rightarrow$ since $\frac{1}{s_b} \ll \tau$, this happens <u>long</u> before next click.

$\Rightarrow$ can take $J_{max} \simeq \infty$ w/o losing any accuracy

i.e., can pretend that <u>infinite</u> # of muts contribute to establishment of new nose.

# Appendix C: formal analysis of coalescence in the nose

Recall: main result for stochastic size of nose:

$$f_{nose}(t) \equiv \frac{\nu}{2Nqs_b} e^{qs_b t} \implies H_{\nu}(z) \approx e^{-\frac{U_b}{s_b} z^{1 - 1/q}}$$

(supplement of lecture 19)

Let's fine-grain this further:

$$\implies \text{let } f_\ell(t) \equiv \text{freq of lineage in nose founded by beneficial mutation @ site } \ell$$

$$\implies \text{then } H_{\nu_\ell}(z) \approx e^{-\frac{\nu}{s_b} z^{1 - 1/q}}$$

$\implies$ Probability that 2 individuals coalesce = probability that they came from same lineage:

$$\Rightarrow \quad P_c(2) = \left\langle \sum_{\ell=1}^{L_{Hb}} \left( \frac{f_\ell(t)}{\sum_{\ell'} f_\ell(t)} \right)^2 \right\rangle = \left\langle \sum_{\ell=1}^{L_{Hb}} \left( \frac{v_\ell}{\sum_{\ell'} v_{\ell'}} \right)^2 \right\rangle$$

$\Rightarrow$ **Trick:** using $\displaystyle\int_0^\infty \frac{\lambda^\alpha}{\Gamma(\alpha)} z^{\alpha-1} e^{-\lambda z}\, dz = 1$, can write as

$$P_c(2) = \left\langle \sum_\ell \left( \frac{v_\ell}{\sum_{\ell'} v_{\ell'}} \right)^2 \right\rangle = \left\langle \sum_\ell v_\ell^2 \int_0^\infty dz\, z\, e^{-\left(\sum_{\ell'} v_{\ell'}\right) z} \right\rangle$$

$$= \sum_\ell \left\langle \int_0^\infty dz \cdot z \cdot \left( v_\ell^2 e^{-z v_\ell} \right) \cdot \prod_{\ell' \neq \ell} e^{-v_{\ell'} z} \right\rangle$$

$$= \sum_\ell \int_0^\infty dz \cdot z \cdot \frac{\partial^2 H_{v_\ell}(z)}{\partial z^2} \cdot \prod_{\ell' \neq \ell} H_{v_{\ell'}}(z)$$

$\Rightarrow$ using results above for $H_{v_\ell}(z)$ and $H_v(z)$,

$$\Rightarrow \quad \frac{\partial H_{v_\ell}}{\partial z} = -\frac{\mu_\ell}{S_b}\left(1-\tfrac{1}{q}\right) z^{-1/q} H_{v_\ell}(z)$$

$$\Rightarrow \quad \frac{\partial^2 H_{\nu_\ell}}{\partial z^2} = \frac{1}{q} \frac{N_\ell}{S_b} \left(1 - \frac{1}{q}\right) z^{-1-\frac{1}{q}} H_{\nu_\ell}(z) + O(N_\ell^2)$$

so that:

$$P_c(2) = \sum_\ell \int_0^\infty dz \cdot z \cdot \frac{\partial^2 H_{\nu_\ell}(z)}{\partial z^2} \cdot \prod_{\ell' \neq \ell} H_{\nu_{\ell'}}(z)$$

$$= \frac{1}{q} \sum_\ell \frac{N_\ell}{U_b} \int_0^\infty dz \; \frac{U_b}{S_b} \left(1 - \frac{1}{q}\right) z^{-\frac{1}{q}} H_\nu(z) \overset{\simeq \prod_\ell H_{\nu_\ell}(z)}{\prod_\ell H_{\nu_\ell}(z)}$$

$$= \frac{1}{q} \int_0^\infty - \frac{\partial H_\nu(z)}{\partial z} \; = \; \frac{1}{q} \left[ H\overset{1}{(0)} - H\overset{0}{(\infty)} \right]$$

$$= \frac{1}{q} \qquad \text{as desired}$$

Can do same thing for larger samples:

$$P_c(n) = \left\langle \sum_\ell \left(\frac{\nu_\ell}{\sum_{\ell'} \nu_{\ell'}}\right)^n \right\rangle = \sum_\ell \int_0^\infty dz \cdot \frac{(-1)^n z^n}{\Gamma(n)} \frac{\partial^n H_{\nu_\ell}(z)}{\partial z^n} \prod_{\ell' \neq \ell} H_{\nu_{\ell'}}(z)$$

$$\Rightarrow \quad \frac{\partial^n H_{\nu_\ell}(z)}{\partial z^n} = \frac{N_\ell}{S_b} \frac{(-1)^n}{q} \left(1 - \frac{1}{q}\right) \prod_{k=1}^{n-2} \left(k + \frac{1}{q}\right) z^{-n+1-\frac{1}{q}} H_{\nu_\ell}(z) + \mathcal{O}(N_\ell^2)$$

$$\approx \frac{N_\ell}{S_b} \frac{(-1)^n (n-2)!}{q} z^{-n+1-\frac{1}{q}} H_{\nu_\ell}(z)$$

$$\Rightarrow \quad P_c(n) = \frac{1}{q} \frac{\Gamma(n-1)}{\Gamma(n)} \cdot \sum_\ell \frac{N_\ell}{U_b} \cdot \int_0^\infty dz \cdot \frac{-\partial H_\nu(z)}{\partial z}$$

$$\approx \frac{1}{q(n-1)}$$

$$\Rightarrow \quad P_c(n) = \frac{P_c(2)}{n-1} \quad \leftarrow \quad \text{also known as Bolthausen-Sznitman coalescent (BSC)}$$