

# In What Direction is Evolution Going?

Stanley Sawyer, Department of Mathematics

Are most new fixed mutations deleterious  
or advantageous?

“Naïve pan-selectionism” would say that  
most are advantageous.

Competing theories suggest that most are  
selectively neutral or mildly deleterious.

Collaborators:

Dan Hartl, Harvard University

Carlos Bustamante, Cornell University

Rob Kulathinal, Harvard University

John Parsch, University of Munich

Haley Abel, Washington University

Amei Amei, University of Nevada Los Vegas

and many others

McDonald-Kreitman tables: For DNA sequences from two closely-related species at one locus:

	mono. at diff. bases	poly. in either spp.
Replacement	$K_r$	$S_r$
Silent	$K_s$	$S_s$

(Replacement means that it changes an amino acid.)

Excess/Deficit of replacement fixed differences suggest possible positive/negative selection.

The polymorphisms at most loci are too sparse to make any conclusion, but we can use the Mantel-Haenszel test to aggregate over loci:

	mono. at diff. bases	poly. in either spp.	
Replacement	$K_r$	$S_r$	$T_r$
Silent	$K_s$	$S_s$	$T_s$
(Totals)	$K_t$	$S_t$	$N$

Over many independent loci ( $1 \leq i \leq L$ ):

$$Z = \sum_{i=1}^L \left( K_{ri} - EK_{ri} \right) / \sqrt{\sum_i \text{Var}(K_{ri})}$$

has an  $N(0, 1)$  distribution if there are no deficits or excesses, where  $EK_r = T_r K_t / N$  at the  $i^{\text{th}}$  locus. ( $Z$  is necessary to avoid Simpson's Paradox.)

*Drosophila*:  $L = 91$   $Z = 3.30$  ( $P < 0.001$ )

*D. melanogaster*:  $7 \leq m_i \leq 12$  sequences

*D. simulans*:  $n_i = 1$  sequence ( $i = \text{locus}$ )

(Baines et al 2007, Sawyer et al 2007)

*Arabidopsis*:  $L = 12$   $Z = -4.47$  ( $P < 10^{-5}$ )

*A. thaliana*:  $14 \leq m_i \leq 21$  sequences

*A. lyrata*:  $n_i = 1$  sequence

(Bustamante et al 2002)

Can we estimate the amount of selection involved?  
How to model?

Many events tend to happen on a scale of  $N_e$  generations, where  $N_e$  is the effective population size.

It is useful to consider five different kinds of mutations, where  $s$  is the rate of selection per generation:

- |       |                      |                     |
|-------|----------------------|---------------------|
| (i)   | $s < 0,  sN  \gg 1$  | Evolutionary lethal |
| (ii)  | $s < 0,  sN  = O(1)$ | Weakly deleterious  |
| (iii) | $s = 0$              | Neutral             |
| (iv)  | $s > 0,  sN  = O(1)$ | Weakly advantageous |
| (v)   | $s > 0,  sN  \gg 1$  | Hopeful monsters(?) |

Evolutionary lethal mutations can be ignored since they rapidly disappear in this time scale, and hopeful monsters are essentially never polymorphic.

This will be a theory of (i,ii,iii,iv). This ignores the most interesting mutations (v), but they may be rare.

Looking ahead, we will find that the expected proportions of beneficial mutations among nonlethal replacement mutations in a 91-locus *Drosophila* dataset are

New (nonlethal) mutations	7%
Polymorphic in samples	30%
Fixed differences	94%

The model: We assume

- All new mutations occur at a new site.
- Sites are unlinked; that is, are statistically independent. (Seems OK by forwards simulation for applications with two related species. Also, many loci show evidence of strong short-segment gene conversion, which could randomize sites.)
- Directional selection for each new mutant site, with no epistasis or dominance over sites.
- Silent sites are neutral. For replacement mutations,  $\gamma = (N_e)s$  for each new mutation is drawn from a normal distribution with parameters

$$N(\gamma_i, \sigma_w^2)$$

where  $\gamma_i$  depends on the  $i^{\text{th}}$  locus. (Bustamante et al 2002 has the same model with  $\gamma \equiv \gamma_i$ .)

- The  $\gamma_i$  for loci are themselves drawn from a normal distribution

$$N(\mu_\gamma, \sigma_b^2)$$

This means that the distribution of  $\gamma$ s for new mutations is the same as a random-effects model in statistics.

*PRF model:* The probability of survival of a new mutant is approximately  $p(\gamma) = (1/N)(2\gamma/(1-\exp(-2\gamma)))$ , so that most new mutants are immediately lost. However, a proportion of these will eventually be fixed.

The sites in the general population that are polymorphic will vary from time to time, but there will always be a random set of sites that are polymorphic. These will have a random set of population site frequencies  $p$  for these random sites, all moving independently (since sites are unlinked).

In the limit as  $N \rightarrow \infty$ , with  $O(1)$  new mutations per generation, the result is a Poisson random field of population site frequencies.

If  $\gamma$  is fixed, the polymorphic population frequencies form a Poisson random field on  $0 < p < 1$  as  $N \rightarrow \infty$  with mean density

$$\theta_r \frac{1 - e^{-2\gamma(1-p)}}{1 - e^{-2\gamma}} \frac{dp}{p(1-p)} \quad (\text{Replacement})$$

$$\theta_s \frac{dp}{p} \quad (\text{Silent})$$

(Sawyer and Hartl 1992) Here  $\theta_r$  and  $\theta_s$  are the replacement and silent-site mutation rates per generation, and mutant replacement bases have a relative selective advantage of  $\gamma/N_e$ .

Fixations occur at the relative rates

$$\theta_r \frac{2\gamma}{1 - \exp(-2\gamma)} \quad \text{and} \quad \theta_s$$

These are *population* fixation rates and polymorphism frequencies. For *samples* of  $m$  and  $n$  sequences from two closely-related species, the counts  $K_r, S_r, K_s, S_s$  are independent Poisson with means

$$E(K_r) = \theta_r \left( \frac{2\gamma}{1 - e^{-2\gamma}} \right) (t + G(m) + G(n))$$

$$E(S_r) = \theta_r \left( \frac{2\gamma}{1 - e^{-2\gamma}} \right) (F(m) + F(n))$$

$$E(K_s) = \theta_s \left( t + \frac{1}{m} + \frac{1}{n} \right)$$

$$E(S_s) = \theta_s (L(m) + L(n))$$

Here  $t$  is the scaled divergence time of the two species and

$$G(n) = \int_0^1 (1-p)^{n-1} \frac{1 - e^{-2\gamma p}}{2\gamma p} dp$$

$$F(n) = \int_0^1 \frac{1 - p^n - (1-p)^n}{1-p} \frac{1 - e^{-2\gamma p}}{2\gamma p} dp$$

$$L(n) = \sum_{i=1}^{n-1} \frac{1}{i}$$

If the  $\gamma$ s are chosen independently from  $N(\gamma_i, \sigma_w^2)$  within each locus, the formulas for  $E(K_r)$  and  $E(S_r)$  are replaced by double integrals, with a Gaussian integral on the outside. The sampling formulas are valid if e.g.  $n = 1$ , as long as the other sample size  $m > 1$ . In that case, all of the polymorphism information comes from the species with  $m > 1$ .

The model allows  $\theta_{ri} \neq \theta_{si}$ , so that  $\theta_{ri}/(2\theta_{si}) = q_i$  gives an estimate of the average number of possible nonlethal amino-acid replacements at the  $i^{\text{th}}$  locus.

A *Drosophila* dataset (Pröschel et al 2006) with  $T = 91$  autosomal loci has  $m_i > 1$  sequences (at the  $i^{\text{th}}$  locus) from *D. melanogaster* in Zimbabwe and one sequence from *D. simulans* (from North Carolina). The model has three “local” parameters for each locus

$$\theta_{ri}, \theta_{si}, \gamma_i \quad 1 \leq i \leq 91$$

and four “global” parameters (shared by all loci)

$$\mu_\gamma, \sigma_b^2, \sigma_w^2, t$$

This makes  $3T + 4 = 277$  parameters for  $4T = 364$  observations.



We used MCMC (Markov Chain Monte Carlo), which essentially estimates parameters by asking where most of the mass of the likelihood

$$L(\theta_{ri}, \theta_{si}, \gamma_i, \mu_\gamma, \sigma_b, \sigma_w, t, K_{ri}, S_{ri}, K_{si}, S_{si}) \quad (*)$$

is, viewed as a function of the parameters, with the observed data  $K_{ri}, S_{ri}, K_{si}, S_{si}$  held constant.

MCMC works by defining a Markov chain with (\*) as a stationary distribution and computing averages and quantiles over long runs of this Markov chain.

Technically speaking, we used  $n = 1,000,000$  “burnin” Markov chain steps to stabilize the parameters and then  $n = 20,000,000$  further steps, sampling only every  $10^{\text{th}}$  step to save memory. The last 10 “subchains”, each with 2,000,000 steps (200,000 samples), gave very similar results.

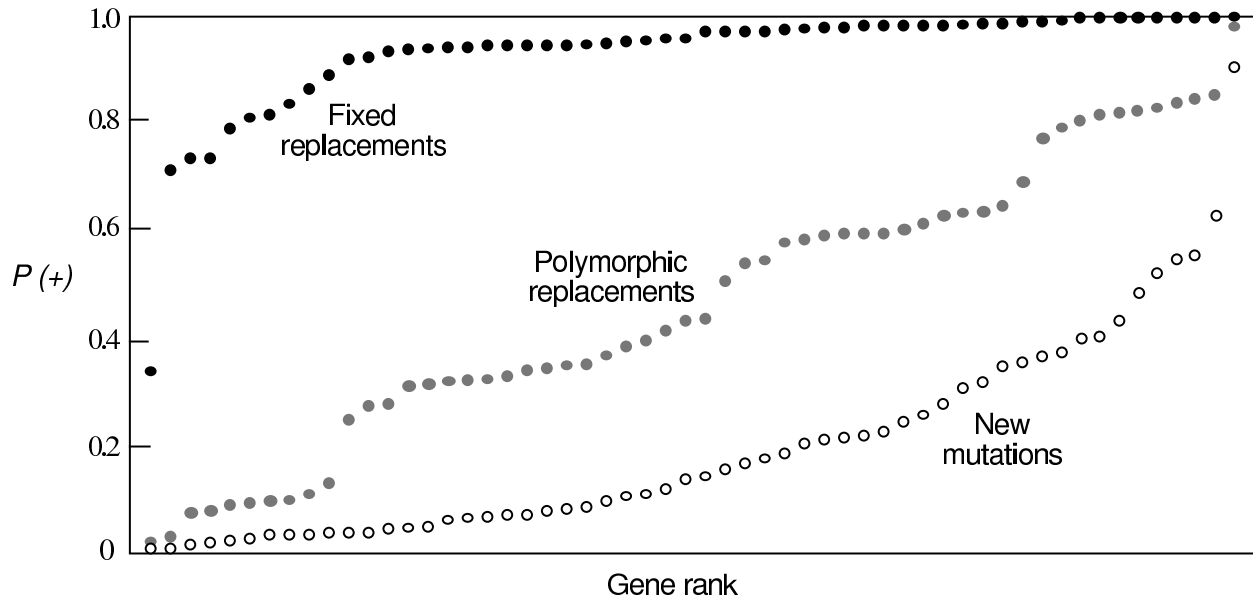
Results for global parameters were  
(2,000,000 samples, 10 subchains)

Var	Mean $\pm 1.96 \times$ SD	GR
$\mu_\gamma > 0$	$0.16 \pm 0.72$	1.009
$\mu_\gamma$	$-11.3 \pm 30.9$	1.020
$\sigma_w$	$6.91 \pm 11.5$	1.021
$\sigma_b$	$4.28 \pm 4.38$	1.018
$\sigma_w / (\sigma_b + \sigma_w)$	$0.55 \pm 0.28$	1.015
$t$	$4.48 \pm 0.44$	1.000

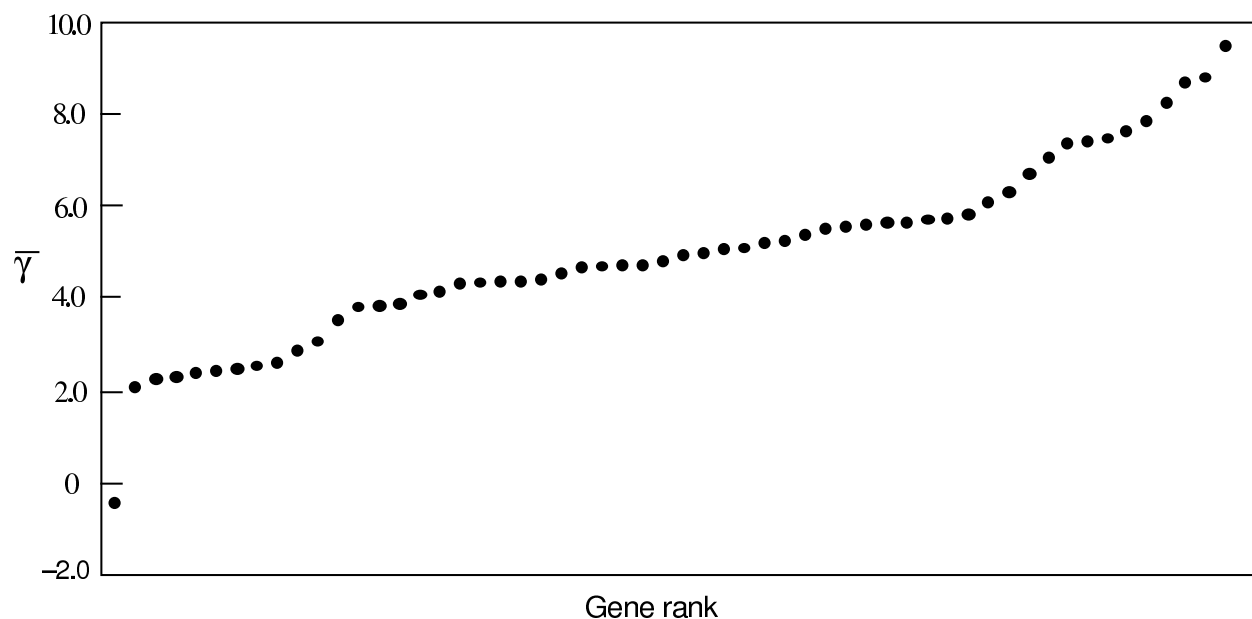
GR is a diagnostic for MCMC convergence.  $GR < 1.03$  is considered good.

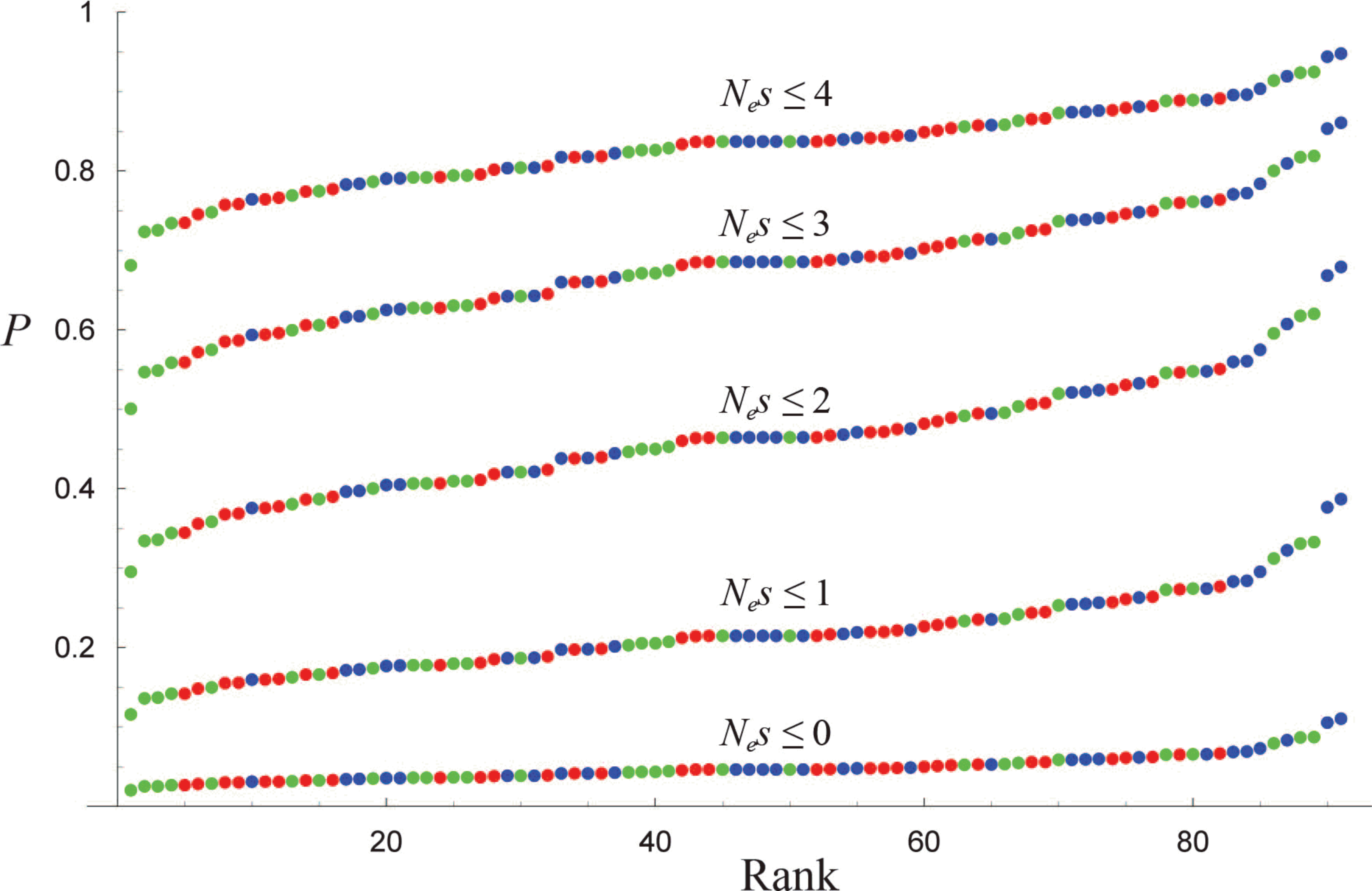
In the last subchain,  $\sigma_w / (\sigma_b + \sigma_w)$  had median 0.59 and varied in the range (0.24, 0.80) (middle 95% quantiles), so that about half of the  $\gamma$ -variability was within locus and about half was between locus.

*Fig 1:* For a 56-locus *Drosophila* dataset (Sawyer et al 2003), the proportions of mutations that are beneficial (based on averages of functions of  $(\gamma_i, \sigma_w)$  over the MCMC run):



*Fig 2:* Among fixed replacement mutations, the average scaled selection coefficients:





Robustness of model to various assumptions:  
(Checked by forwards simulations and other methods)

If linkage is tight, most parameter estimates are accurate but estimated confidence intervals are too small. (This means that conclusions that estimated parameters are nonzero can be wrong.)

If scaled recombination  $R = N_e \rho$  is of moderate size ( $R \geq 50$ ) and the breakpoints are uniformly distributed, conclusions are robust.

If the within-locus distribution of scaled selection coefficients of new replacement mutations is replaced by a double exponential (and so presumably also a gamma distribution) or by an even more heavy-tailed Student- $t$  distribution, then conclusions are qualitatively the same.

Thank you for coming.