



1 **ABSTRACT**

2 The question of the relative evolutionary roles of adaptive and non-adaptive processes has  
3 been a central debate in population genetics for nearly a century. While advances have been  
4 made in the theoretical development of the underlying models, and statistical methods for  
5 estimating their parameters from large-scale genomic data, a framework for an appropriate  
6 null model remains elusive. A model incorporating evolutionary processes known to be in  
7 constant operation - genetic drift (as modulated by the demographic history of the population)  
8 and purifying selection – is lacking. Without such a null model, the role of adaptive processes  
9 in shaping within- and between-population variation may not be accurately assessed. Here,  
10 we investigate how population size changes and the strength of purifying selection affect  
11 patterns of variation at neutral sites near functional genomic components. We propose a  
12 novel statistical framework for jointly inferring the contribution of the relevant selective and  
13 demographic parameters. By means of extensive performance analyses, we quantify the  
14 utility of the approach, identify the most important statistics for parameter estimation, and  
15 compare the results with existing methods. Finally, we re-analyze genome-wide population-  
16 level data from a *Zambian population of *Drosophila melanogaster**, and find that it has  
17 experienced a much slower rate of population growth than was inferred when the effects of  
18 purifying selection were neglected. Our approach represents an appropriate null model,  
19 against which the effects of positive selection can be assessed.

20

21

22 **Keywords:** Background selection, demographic inference, distribution of fitness effects,  
23 approximate Bayesian computation

24

25

26

27

28

29

30

31

32

33

34

## 1 INTRODUCTION

2 At the founding of population genetics in the early 20th century, Fisher, Haldane, and Wright  
3 developed much of the mathematical and conceptual framework underlying the study of  
4 population-level processes dictating variation observed within- and between-species.  
5 However, as evidenced by decades of published interactions, they held differing views  
6 regarding the relative importance of adaptive vs. non-adaptive processes in driving evolution.  
7 As pointed out by Crow (2008), these issues were not really resolved, but "rather they were  
8 abandoned in favor of more tractable studies." With the advent of the Neutral Theory  
9 (Kimura 1968, 1983; King and Jukes 1969; Ohta 1973), the evolutionary importance of  
10 stochastic effects due to finite population size, as earlier advocated by Wright, received  
11 renewed attention.

12 In the following decades, further theoretical developments as well as the availability  
13 of large-scale sequencing data have validated the important role of genetic drift (Kimura  
14 1983; Walsh and Lynch 2018). However, subsequent research on the indirect effects of  
15 selection on patterns of variability at linked neutral alleles has re-ignited previous debates  
16 (Kern and Hahn 2018; Jensen *et al.* 2019). In particular, it remains unclear whether the large  
17 class of strongly and weakly deleterious variants hypothesized under the Neutral Theory, and  
18 their effects on linked neutral sites (background selection, BGS), are sufficient to explain  
19 genome wide patterns of variation, or whether a substantial contribution from the effects of  
20 beneficial variants on linked neutral sites (*i.e.*, selective sweeps), is required.

21 The primary difficulty in answering this question stems from our lack of an  
22 appropriate neutral null model - that is, a model incorporating genetic drift as modulated by  
23 the demographic history of the population, as well as a realistic distribution of fitness effects  
24 summarizing the pervasive effects of both direct and indirect purifying selection. Without a  
25 model incorporating these evolutionary processes, which are certain to be occurring  
26 constantly in natural populations, it is not feasible to quantify the frequency with which  
27 adaptive processes may also be acting to shape patterns of polymorphism and divergence.

28 It can, however, be difficult to distinguish the individual contributions of positive and  
29 purifying selection from demographic factors such as changes in population size, as all of  
30 these evolutionary processes may leave similar imprints in population genetic data. For  
31 example, both purifying selection and population growth can distort gene genealogies of  
32 linked neutral sites in a similar fashion (Charlesworth *et al.* 1993; Kaiser and Charlesworth  
33 2009; O'Fallon *et al.* 2010; Charlesworth 2013; Nicolaisen and Desai 2013), and result in a  
34 skewing of the site frequency spectrum (SFS) towards rare variants. In fact, demographic

1 inference is often performed using either synonymous or intronic sites, which are close to  
2 sites in coding regions, but the contribution of the effects of selection at linked sites are  
3 generally ignored. Patterns of variation in these regions may be skewed by the effects of  
4 either negative selection (Zeng 2013; Ewing and Jensen 2016) or positive selection (Messer  
5 and Petrov 2013), and this could strongly affect the accuracy of the inferred demographic  
6 model (Ewing and Jensen 2016; Schrider *et al.* 2016). In other words, selection may cause  
7 demographic parameters to be mis-estimated in such a way that population size changes are  
8 over or under-estimated.

9 In addition, the extent of BGS can vary considerably across the genome. Although it  
10 is understood theoretically to be a function of the number and selective effects of directly  
11 selected sites, as well as the rate of recombination (Hudson and Kaplan 1995; Nordborg *et al.*  
12 1996; Charlesworth 1996, 2013), the interaction between these parameters and the underlying  
13 demographic history of the population remains poorly understood, even for simple models.  
14 Furthermore, existing analytical work (Zeng and Charlesworth 2010b; Zeng 2013; Nicolaisen  
15 and Desai 2013) has largely been done under the assumption of demographic equilibrium,  
16 and is often restricted to describing mutations of large effect. Thus, weak selection effects (on  
17 the order of  $|2N_e s| < 10$ ), which are thought to be common, may not be well captured by these  
18 predictions. Furthermore, in regions of low crossing over, interference between this class of  
19 mutations may result in even greater distortions of the underlying genealogies (Kaiser and  
20 Charlesworth 2009; O’Fallon *et al.* 2010; Good *et al.* 2014).

21 We first investigate the joint effects of demography, the shape of the distribution of  
22 fitness effects (DFE) of deleterious mutations, and the number of selected sites in shaping  
23 linked neutral variation. Next, we utilize the decay of background selection effects, by  
24 examining regions spanning coding / non-coding boundaries to jointly infer the DFE of the  
25 coding region and the demographic history of the population. By performing extensive  
26 performance analyses, quantifying both power and error associated with this approximate  
27 Bayesian (ABC) approach (Beaumont *et al.* 2002), the method is shown to perform well  
28 across arbitrary demographic histories and DFE shapes. Importantly, by utilizing patterns of  
29 variation and divergence across coding and non-coding boundaries, this approach avoids the  
30 assumption of synonymous site neutrality inherent to MK-style approaches - an assumption  
31 that has been shown to be strongly violated in many organisms of interest (Chamary and  
32 Hurst 2005; Lynch 2007; Zeng and Charlesworth 2010a; Lawrie *et al.* 2013; Choi and  
33 Aquadro 2016; Jackson *et al.* 2017) and which can result in serious mis-inference  
34 (Matsumoto *et al.* 2016). In applying this approach to genome-wide data from a Zambian

1 population of *Drosophila melanogaster*, results show that the Zambian population has  
2 experienced a very mild 1.2-fold growth, considerably less than previous estimates which did  
3 not account for the BGS-induced skew of the SFS. In addition, we estimate that ~25% of all  
4 mutations in exons are effectively neutral in this population, and we find little evidence for  
5 wide-spread selection on synonymous sites.

6  
7  
8

## 9 METHODS

10 **Simulations:** SLiM 3.1 (Haller and Messer 2019) was used to simulate a functional element  
11 of length  $L$ , which is flanked by neutral non-functional regions. The functional region  
12 experiencing purifying selection is given by a DFE that is modeled as a discrete distribution  
13 with four bins (Figure 1a) representing effectively neutral ( $|\gamma| < 1$ ), weakly deleterious ( $1 \leq$   
14  $|\gamma| < 10$ ), moderately deleterious ( $10 \leq |\gamma| < 100$ ), and strongly deleterious ( $100 \leq |\gamma| \leq$   
15  $10000$ ) classes of mutations, where  $\gamma = 2N_e s$ , and  $s$  is the selection coefficient for homozygous  
16 mutations. Semi-dominance is assumed, so that the fitness of mutant heterozygotes is exactly  
17 intermediate between the values for the two homozygotes (a dominance coefficient,  $h$ , of  
18 0.5). Fitness effects are assumed to follow a uniform distribution within each of the four bins.  
19 In order to infer the extent of purifying selection, we estimated the fraction of mutations in  
20 each bin, referred to as  $f_0, f_1, f_2$  and  $f_3$ , respectively (Figure 1a), such that  $0 \leq f_i \leq 1$ , and  $\sum_i f_i$   
21  $= 1$ , for  $i = 0, 1, 2$ , and  $3$ . In addition, in order to limit the computational complexity, we  
22 restricted values of  $f_i$  to multiples of 0.05 (*i.e.*,  $f_i \in \{0.0, 0.05, 0.10 \dots 0.95, 1.0\} \forall i$ ). These  
23 constraints allowed us to sample 1,771 different DFE realizations, and to work independently  
24 of any arbitrary assumption regarding DFE shape.

25

26 **Simulations under demographic equilibrium:** Simulations were performed for 4 different  
27 values of  $L$  – 0.5 kb, 1kb, 5kb, and 10kb. The intergenic regions were assumed to be 10kb  
28 and simulations were restricted to the intergenic region on one side of the functional region.  
29 For the purpose of power analyses and testing, we used population-genetic parameter values  
30 that approximately resemble *Drosophila* populations. Population size was assumed to be  $10^6$   
31 and the recombination rate ( $1 \times 10^{-8}$  per site per generation) and mutation rate ( $1 \times 10^{-8}$  per  
32 site per generation) were constant across the simulated region. Although we have not  
33 included gene conversion in this study, it will be an important addition in future studies. The

1 simulations were performed with 5000 ( $=N_{\text{sim}}$ ) diploid individuals and the recombination and  
2 mutation rates were scaled proportionally to maintain realistic values of  $N_e s$ .

3 We used a burn-in period of 80,000 generations, and an additional 20,000 ( $=4N_{\text{sim}}$ )  
4 generations were allowed for neutral evolution. For every set of parameter combination (*i.e.*,  
5  $f_0, f_1, f_2$  and  $f_3$ ) we performed 1000 replicate simulations, and summarized both the mean and  
6 variance of common summary statistics to perform the subsequent ABC analysis.

7

8 ***Simulations under non-equilibrium demography:*** Simulations with demographic changes  
9 were performed specifically to match the details of the *D. melanogaster* genome. A set of 94  
10 exons belonging to the *D. melanogaster* genome were chosen according to certain criteria  
11 (see Results). For each exon, simulations were performed using the length of the exon  
12 together with 4 kb of flanking intergenic sequence. The mutation rate was conservatively  
13 assumed to be  $3.0 \times 10^{-9}$  per site per generation (Keightley *et al.* 2014) although somewhat  
14 higher mutation rates have been estimated in other studies (Schrider *et al.* 2013; Assaf *et al.*  
15 2017). Ancestral and current population sizes were sampled from a uniform prior between  
16  $10^5$ - $10^7$  and  $f_i \in \{0.0, 0.05, 0.10 \dots 0.95, 1.0\}$  such that  $\sum_i f_i = 1$ , for  $i = 0, 1, 2$ , and 3.  
17 Nucleotide diversity at 4-fold degenerate sites was found to be 0.019 for the Zambian  
18 population of *D. melanogaster*, which would give an estimate of  $N_e$  of  $1.6 \times 10^6$ . A scaling  
19 factor of 320 corresponding to  $N_e/N_{\text{sim}} (=1.6 \times 10^6 / 5000)$  was used to perform all  
20 simulations with demographic changes. A total of 10 replicates were performed for each  
21 exon, resulting in 940 replicates for every parameter combination. These simulations were  
22 conducted using the computational resources of Open Science Grid (Pordes *et al.* 2007;  
23 Sfiligoi *et al.* 2009).

24

25 ***Calculation of summary statistics:*** First, we fit a logarithmic function to the recovery of  
26 nucleotide diversity ( $\pi$ ) around the functional region such that  $\pi = \text{slope} * \ln(x) + \text{intercept}$ ,  
27 where  $x$  is the distance of the site from the functional region in base pair. We used the *slope*  
28 and *intercept* of the fit to define the number of bases required for a 50%, 75%, and 90%  
29 recovery of nucleotide diversity, with 50% and below being defined as the “linked neutral”  
30 region and the 50% and above as the “neutral” region. This analysis provides for three non-  
31 overlapping regions: (1) functional (experiencing direct selection), (2) linked-neutral  
32 (experiencing observable levels of background selection), and (3) neutral (experiencing low /  
33 unobservable levels of background selection). The following statistics were calculated for

1 each of these three types of regions: nucleotide diversity ( $\pi$ ), Watterson's  $\theta$ , Tajima's  $D$ , Fay  
2 and Wu's  $H$  (both absolute and normalized), number of singletons, haplotype diversity, LD-  
3 based statistics ( $r^2$ ,  $D$ ,  $D'$ ), and divergence (*i.e.*, number of fixed mutations per site per  
4 generation). Simulations for any particular set of parameters were run with 1000 replicates  
5 and the mean and variance of the above statistics across replicates were used as summary  
6 statistics for ABC. In addition to these variables, six statistics summarizing the characteristics  
7 of the recovery of  $\pi$  in linked neutral regions were also included as summary statistics.  
8 Specifically,  $\pi = slope * \ln(distance) + intercept$  was fit and slope, intercept, maximum value of  
9  $\pi$ , and number of bases required for 50%, 75%, and 90% recovery were calculated and  
10 included as summary statistics. Together, these amount to 72 initial summary statistics. All  
11 statistics were calculated using the Python package pylibseq (Thornton 2003). The sample  
12 size was kept constant at 100 genomes (*i.e.*, 50 diploid individuals). It should be noted that  
13 some statistics are strongly dependent on the number of sites used in the calculations, and the  
14 size of linked and neutral regions varied for every set of parameter combination, though this  
15 effect is captured in the individual prior distributions.

16

17 **ABC:** We used an approximate Bayesian (ABC) approach, using the R package,  
18 "abc" (Csilléry *et al.* 2012), to co-estimate the DFE characterizing a functional region, as well  
19 as the population history characterizing the population in question. We used linear regression  
20 (aided by neural net to handle non-linearity) as well as ridge regression to infer posteriors,  
21 with a tolerance of 0.05 (*i.e.*, 5% of the total number of simulations are accepted by ABC to  
22 estimate the posterior probability of each parameter), a cross-validation set of 100  
23 simulations, and the weighted median as point estimates.

24

25 **Ranking of summary statistics:** Ranking of summary statistics was performed separately for  
26 both demographic equilibrium and non-equilibrium cases, using two different methods. The  
27 first approach consisted of performing Box-Cox transformations on all 72 summary statistics  
28 to correct for non-linear relations between statistics and parameters. The squared correlation  
29 coefficient,  $r^2$ , between the transformed statistics and parameters was then used to rank each  
30 statistic for every parameter separately and a statistic was considered to be significantly  
31 correlated with the parameter if the  $p$ -value was less than 0.05 (Bonferroni corrected for  
32 multiple testing). The second approach involved a modified version of the algorithm  
33 proposed by Joyce and Marjoram (2008) for ranking statistics. With this algorithm, we

1 started with the entire set of 72 statistics. Every statistic was removed from the set and cross-  
2 validation using 20 randomly sampled simulations was used to identify the statistic that  
3 corresponded to the least error (*i.e.*, the removal of which causes the least reduction in  
4 accuracy). The same algorithm was performed iteratively until only two statistics remained.  
5 This method was performed for each parameter separately, was replicated 10 times, and the  
6 average ranking across these replicates was used to obtain the final ranking. The second  
7 approach was extremely time consuming and was thus only used to rank the statistics to infer  
8 the DFE under demographic equilibrium.

9

10 **Comparison with DFE-alpha:** Simulations were performed under demographic-non-  
11 equilibrium models, with 100 replicates of 94 exons each, and ancestral population sizes of  
12 10,000 for all. Functional regions were simulated with 30% neutral sites, which were used to  
13 calculate the neutral SFS required by DFE-alpha. *Est\_dfe* (Schneider *et al.* 2011) was used on  
14 unfolded SFS to perform demographic inference and to infer the deleterious DFE. The  
15 proportion of adaptive mutations was fixed at 0.0. Final estimates of the DFE were obtained  
16 as  $N_w s$  where  $N_w$  is the weighted population size inferred by *est\_dfe*.

17

18 **Drosophila data application:** Release 5 of the *D. melanogaster* genome assembly (Hoskins  
19 *et al.* 2007) and annotation version 5.57 were used, downloaded from

20 [ftp://ftp.flybase.net/genomes/Drosophila\\_melanogaster/dmel\\_r5.57\\_FB2014\\_03/gff/](ftp://ftp.flybase.net/genomes/Drosophila_melanogaster/dmel_r5.57_FB2014_03/gff/).

21 Crossing over rates estimated by Comeron *et al.* (2012) for every exon and flanking  
22 intergenic region were obtained from the *D. melanogaster* Recombination Rate Calculator  
23 ([https://petrov.stanford.edu/cgi-bin/recombination-rates\\_updateR5.pl](https://petrov.stanford.edu/cgi-bin/recombination-rates_updateR5.pl)) (Fiston-Lavier *et al.*  
24 2010), and explicitly utilized for each specific region considered. These rates were halved to  
25 obtain sex-averaged rates of recombination (Campos *et al.* 2017) as all regions were  
26 restricted to autosomes. We excluded all genes that have a crossing over rate 5-fold larger or  
27 smaller than the average (*i.e.*, we used only genes with a crossing over rate of between 0.44  
28 and 11 cM/Mb). Consensus sequences of all Zambia lines were downloaded from  
29 <http://www.johnpool.net/genomes.html> (Lack *et al.* 2015). IBD tracks and admixture tracks  
30 were masked using scripts provided by the same site. Individuals with any known inversions  
31 were entirely excluded from the analysis (Kapopoulou *et al.* 2018).

32 The final set consisted of 76 haploid genomes. PhastCons scores calculated with  
33 respect to 15 insect taxa were downloaded from the UCSC genome browser  
34 (<https://genome.ucsc.edu/>). For each of the 94 exons, summary statistics were calculated



1 using pylibseq (Thornton 2003) for the coding region and for 2kb intergenic regions flanking  
2 both sides. In order to exclude sites in intergenic regions that might be under direct selection,  
3 a phastCons cutoff score of 0.8 was used to calculate all statistics. That is, sites that had a  
4 greater than or equal to 80% probability of being a conserved noncoding element identified  
5 by phastCons, were excluded when calculating statistics.

6 For the purpose of obtaining derived alleles and for calculating branch-specific rates  
7 of substitution, we used the ancestral sequence to the *D. melanogaster* genome provided to us  
8 by the authors of Kolaczkowski *et al.* (2011). The ancestral sequence reconstruction had been  
9 performed by maximum likelihood over 15 insect genomes available in the UCSC genome  
10 browser (Karolchik *et al.* 2004). Sites with missing ancestral sequence were excluded from  
11 analysis. Branch-specific rates of substitution (also referred to as divergence in this study)  
12 were calculated by identifying derived alleles that were fixed in the *D. melanogaster*  
13 Zambian population (*i.e.*, polymorphic sites were removed). After excluding sites with  
14 missing ancestral information, with IBD and admixture tracks, and which were likely to  
15 belong to a non-coding conserved element, we had on average 1062 sites per exon, 556 sites  
16 per linked region, and 666 sites per neutral regions.

17 It should be noted that for the purpose of performing inference using ABC,  
18 substitution rates in simulations were calculated per base pair for 25,000 generations. We  
19 thus normalized all rates obtained from simulations by the expected neutral substitution rate  
20 (*i.e.*,  $\mu_{sim}t_{sim} = 320 \times 3 \times 10^{-9} \times 25000 = 0.024$ , where  $\mu$  is the mutation rate and  $t$  is the number  
21 of generations). Divergence estimates from *D. melanogaster* were normalized by an expected  
22 neutral substitution rate of  $\mu t = 3 \times 10^{-9} \times 21333333$  (the estimated divergence time) = 0.064.  
23 In addition, inference was performed using divergence estimates only in the exonic regions.  
24 ABC inference for *Drosophila* was performed using the abc package in R, with linear  
25 regression aided by neural net with default parameters. Each inference was performed 50  
26 times, and the mean of point estimates obtained were reported as the final estimates of  
27 parameters.

28

### 29 **Data and code availability**

30 The following data will be made publicly available upon acceptance of the manuscript on  
31 [https://github.com/paruljohri/BGS\\_Demography\\_DFE](https://github.com/paruljohri/BGS_Demography_DFE). 1) Aligned sequences of the single-  
32 exon genes and their corresponding intergenic regions used in this study, including derived  
33 alleles and fixed substitutions. 2) Scripts to calculate statistics from simulations and from

1 empirical data as well as the code used to perform simulations. 3) Values of all calculated  
2 statistics obtained for all parameter combinations.

3

4

## 5 RESULTS AND DISCUSSION

6 **Recovery of nucleotide diversity as predicted under equilibrium:** The nucleotide site  
7 diversity ( $B$ ) at neutral sites with linkage to sites experiencing direct purifying selection can  
8 be obtained by modifying Equation 6 of Nordborg *et al.* (1996), which is of the form

9

$$10 \quad B = \frac{\pi}{\pi_0} \sim \exp \left[ -\int \int E(t, z) dz dt \right]$$

11

12 where  $\pi_0$  is the nucleotide diversity without selection and  $\pi$  is the nucleotide diversity with  
13 background selection effects. The exponent  $E$  is a function of the distribution of heterozygous  
14 selection coefficients ( $t = hs$ ) for deleterious mutations and the physical distance ( $z$ ) between  
15 the neutral and selected sites. Here,  $s$  is the selection coefficient, and  $h$  is the dominance  
16 coefficient.

17 For the purpose of the current study, Equation S1a of the SI of Campos and  
18 Charlesworth (2019) was modified to model a neutral site outside a gene, and which is a  
19 distance  $y$  basepairs from the end of the functional region. If the position of a selected site is a  
20 distance  $x$  basepairs from the end (in the opposite direction), the distance between the two  
21 sites is  $z = x + y$ . The basic equation for the exponent of the BGS function for a given  
22 selection coefficient,  $E(t)$ , was obtained as follows:

23

$$24 \quad E(t) = \frac{Ut}{l} \int_0^l \frac{dx}{[t + (g + r_c y)(1 - t) + r_c x(1 - t)]^2} \quad (1)$$

25

26 where  $U(t)$  is the total mutation rate to deleterious alleles over the entire gene,  $l$  is the length  
27 of the gene in basepairs,  $g$  is the rate of gene conversion, and  $r_c$  is the rate of crossing over per  
28 basepair. The crossover map is assumed to be linear, so that the net rate of recombination  
29 between the two sites is  $g + r_c z$ , and  $z$  is assumed to be sufficiently large that the effect of  
30 gene conversion is independent of  $z$ .

1 On integrating Equation 1 with respect to  $x$  between 0 and  $l$  (see the Appendix for  
 2 details), the exponent  $E$  as a function of the length of selected sites and the distance from the  
 3 functional element can be obtained:

$$\begin{aligned}
 4 \quad E(t) &= \frac{Ut}{r_c l(1-t)} \left\{ \frac{1}{[t + (g + r_c y)(1-t)]} - \frac{1}{[t + (g + r_c y)(1-t) + r_c l(1-t)]} \right\} \\
 5 \\
 6 \\
 7 \quad &= \frac{Ut}{[t + (g + r_c y)(1-t)][t + g(1-t) + r_c(y+l)(1-t)]} \quad (2)
 \end{aligned}$$

8  
 9 Note that the above equation shows that, if  $t$  is small compared with  $y$ , BGS effects  
 10 outside the coding region will be minimal.

11 We can integrate  $E(t)$  over the distribution of selection coefficients as described in the  
 12 Appendix. The expectation of  $E(t)$  is given by the sum of the following two terms:

$$13 \\
 14 \quad U[r_c l(1-a)]^{-1} \left\{ 1 + a[(1-a)(t_{i+1} - t_i)]^{-1} \ln \left[ \frac{a + (1-a)t_i}{a + (1-a)t_{i+1}} \right] \right\} \quad (3a)$$

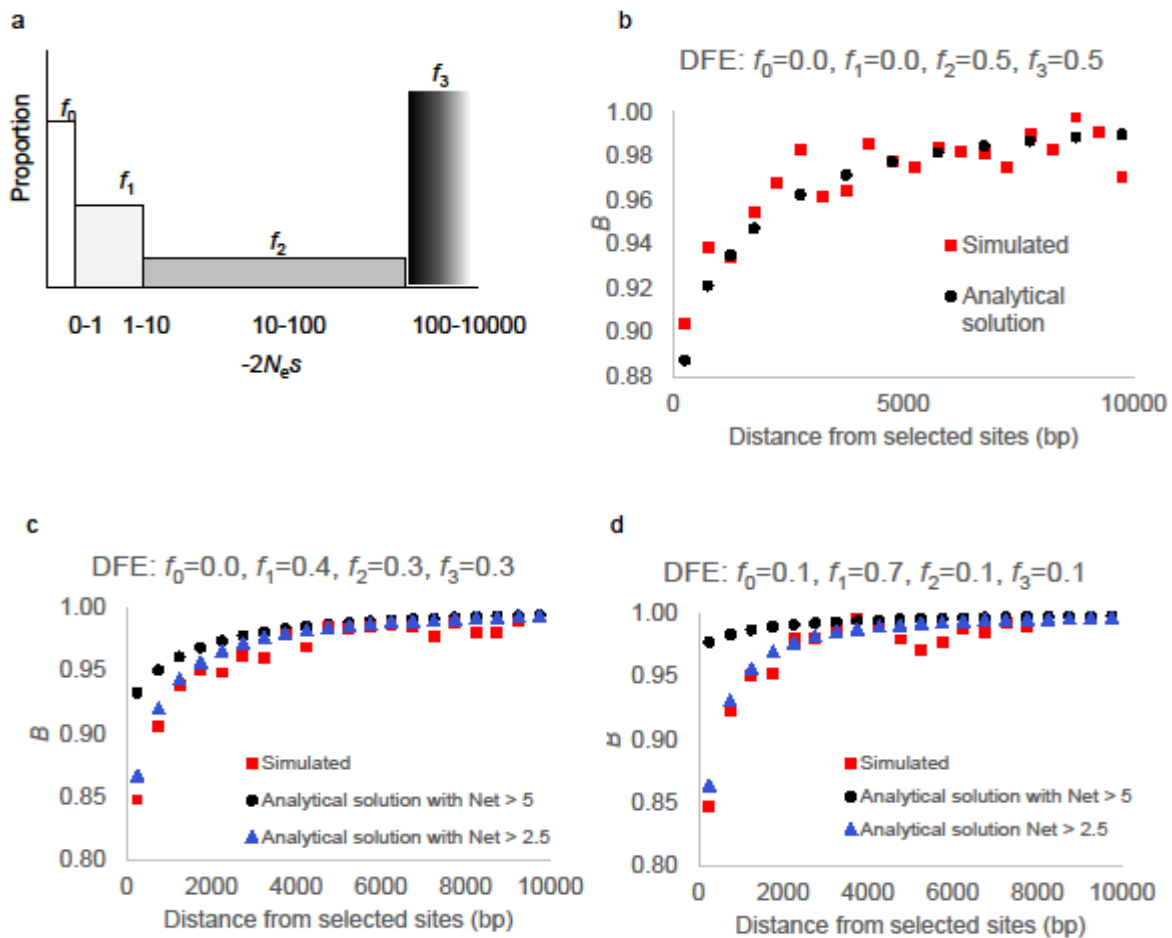
$$15 \\
 16 \quad - U[r_c l(1-b)]^{-1} \left\{ 1 + b[(1-b)(t_{i+1} - t_i)]^{-1} \ln \left[ \frac{b + (1-b)t_i}{b + (1-b)t_{i+1}} \right] \right\} \quad (3b)$$

17  
 18 where  $a = g + r_c y$  and  $b = g + r_c(y + l)$  and the  $t_i$ 's correspond to the boundary of the  
 19 discrete bins. For the case when  $b \ll 1$ , the sum of the two components is approximately  
 20 equal to:

$$21 \\
 22 \quad U(t_{i+1} - t_i)^{-1} \ln \left[ \frac{b + t_{i+1}}{b + t_i} \right] \quad (3c)$$

23  
 24 Figure 1b shows the theoretical results as well as those from simulations, for  $r = 10^{-6}$ ,  
 25  $l = 1000$ ,  $U = l\mu$ ,  $\mu = 10^{-6}$ ,  $g = 0$ ,  $t_0 = 0$ ,  $t_1 = 0.00005$ ,  $t_2 = 0.0005$ ,  $t_3 = 0.005$ , and  $t_4 = 0.5$ . It  
 26 should be noted that these derivations assume that  $N_e t \gg 1$ , which is violated by the presence  
 27 of the weakly deleterious DFE class ( $f_1$ ). Most studies deal with this assumption by ignoring  
 28 the contribution of mutations with  $N_e t < 5$  or 10 (Charlesworth 2013; Elyashiv *et al.* 2016;  
 29 Torres *et al.* 2019). As expected, we found a significant discordance between the simulated  
 30 and theoretically predicted values for the slope of the recovery of diversity as  $f_1$  increases

1 (Figure 2c and 2d, Supp Table 1). On including only mutations with  $N_{et} > 2.5$ , the diversity  
 2 patterns are mostly well explained, even when the DFE is highly skewed towards the weakly  
 3 deleterious class. In fact, it is interesting to note that a combination of high values of  $f_1$  and  $f_2$   
 4 can result in BGS effects that extend up to  $\sim 4$  kb, even for very short exons, although the  
 5 maximum reduction in diversity is around 10-15% (consistent with Charlesworth 2012,  
 6 Campos and Charlesworth 2019b) .  
 7



8  
 9  
 10  
 11  
 12  
 13  
 14  
 15  
 16  
 17  
 18  
 19  
 20

**Figure 1:** (a) An example of a discrete DFE with four classes of mutations. The proportion of each class of mutation,  $f_i$ , lies between 0 and 1. (b) Nucleotide site diversity relative to the neutral expectation ( $B = \pi/\pi_0$ ) as a function of the distance from the directly selected sites (length 1 kb), as predicted by the analytical solution (black points) and observed in simulations (red points). (c, d) Analytical predictions and simulated values for a DFE with larger contributions from the weakly deleterious class of mutations. Note that, for the analytical solutions, the two classes of results represent cases where mutations with  $2N_{et} < 5$  (black circles) and  $2N_{et} < 2.5$  (blue triangles) were ignored.

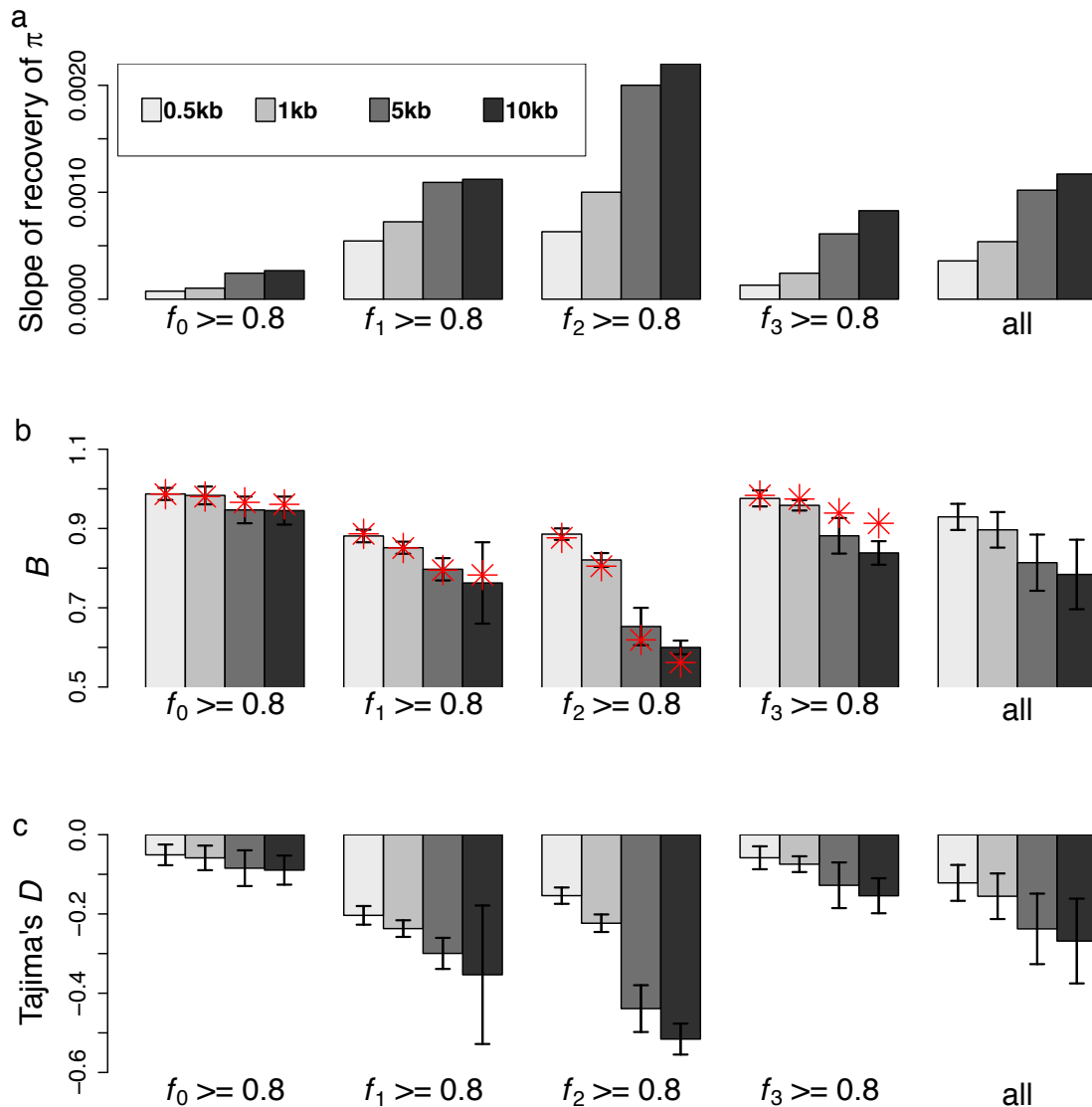
## 1 **Joint effects of demography, the DFE and the number of selected sites on linked neutral** 2 **variation**

3 While the above results show that the effect of BGS on linked neutral regions can be  
4 determined analytically, there are several reasons for investigating background selection  
5 effects using simulations. First, the analytical expressions neglect the contribution of very  
6 weakly deleterious mutations ( $|N_e t| < 2.5$ ), and do not predict the skew in the SFS. In addition,  
7 they assume demographic equilibrium, which is probably not true of natural populations.

8  
9 ***Effects of the shape of the DFE and number of selected sites:*** We first simulated 10kb  
10 neutral regions linked to functional regions of varying sizes, 0.5kb, 1kb, 5kb, and 10kb,  
11 assuming demographic equilibrium, as shown in Figure 2. By varying the contributions from  
12 each bin of selective effects, denoted by  $f_0, f_1, f_2$  and  $f_3$ , it was possible to sample all possible  
13 DFE shapes, as described in the Methods section. As expected from equation 3c, the  
14 reduction in diversity is non-linearly proportional to the number of selected sites for a given  
15 recombination rate. A larger number of selected sites increases both the total reduction in  
16 diversity and the slope of the recovery of diversity away from functional regions (Supp  
17 Figure 1). The maximum reduction in diversity in the linked neutral regions (immediately  
18 adjacent to the functional region), averaged across all DFE realizations, is approximately 8%,  
19 12%, 24%, and 29% for 0.5kb, 1kb, 5kb, and 10kb, respectively. Furthermore, for the chosen  
20 recombination rate, the median numbers of base pairs necessary to achieve a 50% recovery in  
21 diversity are 955, 1035, 1350, and 1650 bp, respectively, (Figure 2a).

22 The reduction of nucleotide diversity at closely linked neutral regions was maximized  
23 when the proportion of weakly deleterious mutation ( $f_1$ ) and moderately deleterious mutations  
24 ( $f_2$ ) was largest (Figure 2b, Supp Table 2). The effect is maximized when purifying selection  
25 is weak, allowing mutations to segregate in the population prior to being purged (Campos *et*  
26 *al.* 2017). Although weakly deleterious mutations ( $f_1$ ) only reduce variation slightly, they  
27 generate significant distortions in the SFS (Figure 2c), consistent with previous studies  
28 (Nordborg *et al.* 1996; Charlesworth 2012; Nicolaisen and Desai 2013). Moderately  
29 deleterious mutations ( $f_2$ ) result in the largest reduction in  $\pi$ , the highest rate of recovery of  $\pi$   
30 around functional regions, and the largest skew in the SFS towards rare variants. As  
31 expected, the proportion of strongly deleterious mutations ( $f_3$ ) does not greatly affect levels of  
32 linked neutral variation, and these mutations skew the SFS only slightly. Further, increasing  
33 the number of selected sites results in larger background selection effects for all DFE types,

1 as is to be expected. It should be noted that these generalizations of BGS effects are  
 2 dependent on how far from the selected sites the measurements are being considered. For  
 3 instance, the distance affected by deleterious mutations is expected to be an increasing  
 4 function of the size of their fitness effects. As we were interested in understanding BGS  
 5 effects caused by all classes of mutations, we focus our discussion to sites closer to the  
 6 functional boundary, where all classes of mutation are likely to have an impact.  
 7



8  
 9 **Figure 2:** Effects of BGS under demographic equilibrium. (a) The slope of the recovery of  
 10 nucleotide diversity in 10kb linked neutral regions flanking functional regions fitted such that  
 11  $\pi = slope * \ln(\text{distance from functional region}) + intercept$ , (b) nucleotide diversity in 500bp  
 12 linked neutral regions flanking functional regions relative to neutral expectation, and (c)  
 13 Tajima's  $D$  for 500bp linked neutral region flanking functional regions. All of the above are  
 14 shown for various sizes of functional elements (0.5-10 kb) and DFE shapes. The four DFE  
 15 shapes considered are  $f_i \geq 0.8$  for  $i=0,1,2,3$ , where more than 80% of mutations reside in  
 16 DFE class  $f_i$ , such that  $\sum f_j = 0.2$ , where  $j \neq i$ . The DFE category "all" represents an average over

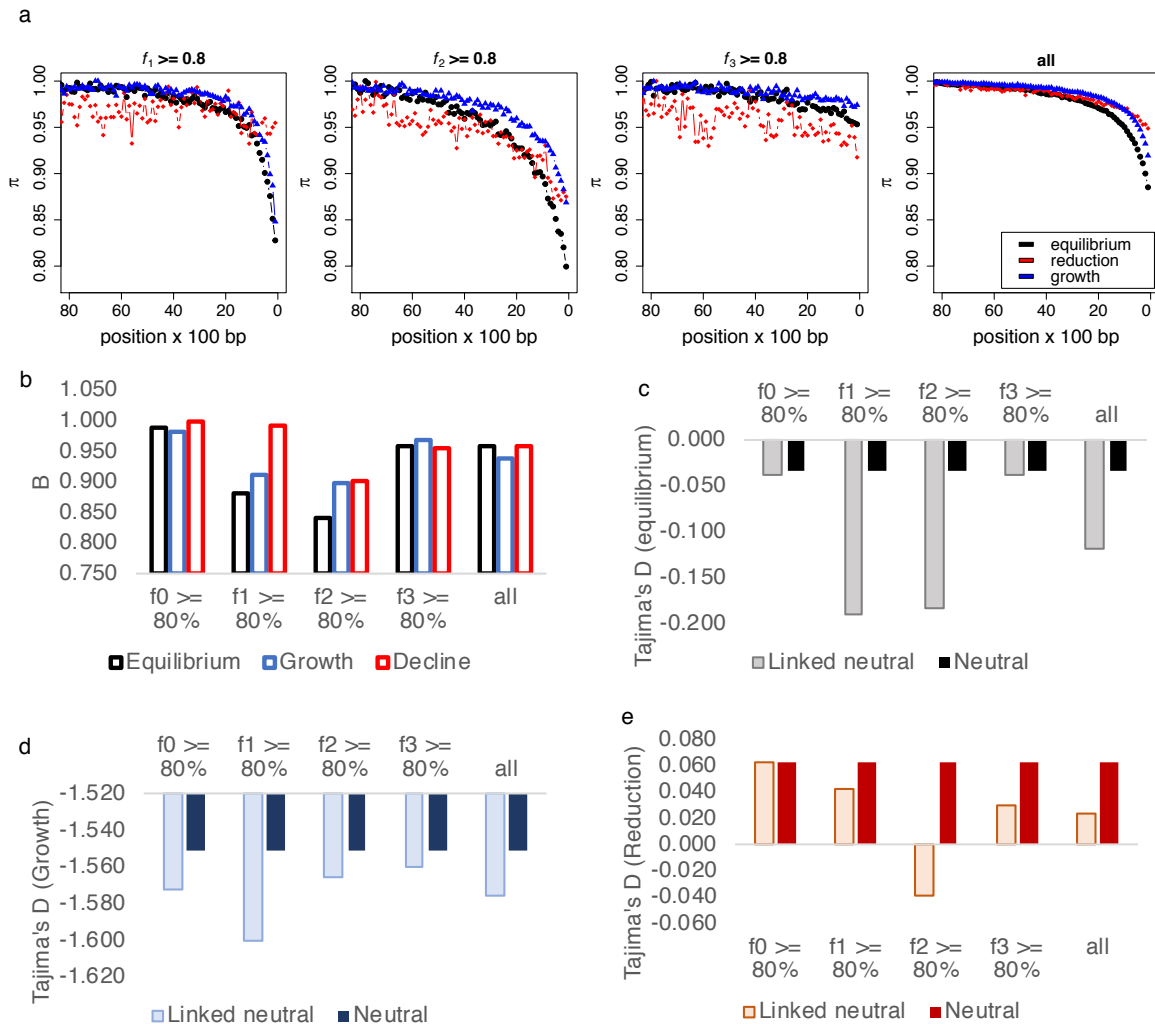
1 all possible DFE shapes. The error bars are  $2 \times$  standard deviation. Red points show the  
2 analytical prediction for (1)  $f_0=0.85, f_1=0.05, f_2=0.05, f_3=0.05$ , (2)  $f_0=0.05, f_1=0.85, f_2=0.05,$   
3  $f_3=0.05$ , (3)  $f_0=0.05, f_1=0.05, f_2=0.85, f_3=0.05$ , and (4)  $f_0=0.05, f_1=0.05, f_2=0.05, f_3=0.85$ .

4  
5  
6  
7 To summarize, at demographic equilibrium, neutral regions linked to functional  
8 regions under selection undergo a reduction in diversity and a skew in the site frequency  
9 spectrum, both of which depend on the underlying shape of the DFE and the number of  
10 directly selected sites (Charlesworth *et al.* 1993; Charlesworth 2013; Campos and  
11 Charlesworth 2019). Importantly for the sake of statistical inference, however, the three  
12 classes of deleterious mutation ( $f_1, f_2, f_3$ ) behave differently, suggesting the possibility of  
13 distinguishing their relative contributions (as discussed in the next section). Furthermore,  
14 these results again demonstrate the potentially important role of BGS in shaping patterns of  
15 neutral variation, highlighting the danger posed by ignoring these effects when performing  
16 demographic inference (see Ewing and Jensen 2016). Additionally, the dramatic difference in  
17 the extent of background selection effects as a function of the number of directly selected  
18 sites strongly implies the necessity of directly modeling exon sizes in any empirical  
19 application.

20  
21 ***Effects of demography and the shape of the DFE on background selection:*** We investigated  
22 the effects of BGS after recent changes in population size. Populations with the same  
23 ancestral population size ( $N_{anc}$ ) either experienced 10-fold exponential growth or contraction  
24 in the last  $4N_{anc}$  generations and BGS effects were compared to populations that remained in  
25 equilibrium throughout, for all possible DFE shapes. Both expansion and contraction result in  
26 reduced BGS effects (*i.e.*, there is an increase in  $B$  compared to equilibrium), irrespective of  
27 the shape of the DFE (Figure 3a, b). This observation suggests that the extent of BGS caused  
28 by functional elements may not only be determined by the strength of selection, but rather  
29 also by the demographic history of the population. Thus, demographic effects may in  
30 principle explain variable inferences among studies of the importance of purifying selection  
31 in shaping genome-wide patterns of variation (Cutter and Payseur 2013).

32  
33 Interestingly however, there is still a significant skew in the SFS at linked neutral sites caused  
34 by BGS after a population size change (Figure 3c-e). Thus, in more compact genomes, where  
35 background selection is pervasive, this suggests that methods which use the SFS to fit

- 1 demographic models may over-estimate growth and either under-estimate population
- 2 contraction or mis-classify contraction as expansion. It is also interesting to note that BGS
- 3 effects are largest under demographic equilibrium, such that constant population size is likely
- 4 to be inferred as population growth.



5 **Figure 3:** Effects of BGS under non-equilibrium demography. (a) The slope of recovery of  
6 nucleotide diversity in linked neutral regions for different DFE shapes under equilibrium  
7 demography (black), population expansion (blue), and contraction (red). (b) Nucleotide site  
8 diversity relative to neutral expectation over 500 bp of linked neutral regions flanking  
9 functional regions, for varying DFE shapes and three different demographic models -  
10 equilibrium (black), 10-fold exponential expansion (blue), and 10-fold exponential decline  
11 (red). (c) Tajima's  $D$  for the 500 bp linked neutral region flanking the functional region under  
12 equilibrium, (d) after a 10-fold expansion, and (e) after a 10-fold population size reduction.  
13 The four DFE shapes considered in all panels are  $f_i \geq 80\%$  for  $i=0,1,2,3$ , where more than  
14 80% of mutations reside in DFE class  $f_i$ . The DFE category "all" represents an average over  
15 all possible DFE shapes. For non-equilibrium demography,  $\gamma=2N_{anc}s$ , where  $N_{anc}$  is the  
16 ancestral population size.

17  
18  
19



## 1 **Inference of the DFE under demographic equilibrium**

2 The next question we investigated was whether the parameters of the DFE can be estimated  
3 using the set of summary statistics described in the Methods section. We first determined  
4 whether it is possible to distinguish the four different classes of the DFE under demographic  
5 equilibrium, using population genomic data and divergence from the closest outgroup  
6 species. The simulations involved functional regions of lengths  $L = 0.5\text{kb}$ ,  $1\text{kb}$ ,  $5\text{kb}$  and  
7  $10\text{kb}$ , with linked neutral regions of  $10\text{kb}$  and a discrete DFE as described previously. An  
8 approximate Bayesian (ABC) approach was implemented to quantify our ability to infer the  
9 four DFE parameters. The recovery of nucleotide diversity over linked neutral regions was  
10 used to calculate the number of bases ( $\pi_{50}$ ) required for diversity to recover to 50% of its  
11 maximum value observed (see Methods). The linked neutral region within  $\pi_{50}$  base pairs  
12 from the functional region was defined as “Linked”, and the remainder was defined as  
13 “Neutral” (Figure 4a). Statistics were calculated for three regions (Functional, Linked, and  
14 Neutral) separately and the means and variances across simulation replicates of each statistic  
15 were used to infer the four parameters. The simulation replicates signify independent loci in a  
16 genome. In the following sub-sections, we describe the performance of the method and its  
17 robustness to various model violations.

18

19 **Accuracy of inference:** All four DFE classes were estimated fairly accurately when using all  
20 statistics (Supp Figure 2a). However, under demographic equilibrium, the DFE is inferred  
21 much more accurately using statistics from the functional regions alone, thus side-stepping  
22 the need for the identification of linked neutral regions (Figure 4b, Supp Figure 2b). In both  
23 cases, the accuracy of inference is highest for the neutral class and lowest for  $f_2$  (*i.e.*, for  
24 moderately deleterious mutations), and improves significantly when the size of the functional  
25 region increases (Supp Figure 2). While using only functional regions to perform inference,  
26 the absolute difference between the true value and the estimated value of the neutral class is  
27 approximately 0.034, 0.030, 0.017, and 0.010 for functional sizes of 0.5kb, 1kb, 5kb and 10  
28 kb. That is, for 1kb regions the method cannot distinguish whether the neutral class of  
29 mutations comprise 30% or 33% of the DFE. For the moderately deleterious class this error is  
30 larger – 0.077, 0.060, 0.028, and 0.019, respectively. These absolute error values are not  
31 surprising, as the  $f_i$  in our simulations are multiples of 0.05 out of computational necessity.  
32 The accuracy of the estimates can thus be increased by sampling the parameter space  
33 more densely. The accuracy of estimation can also be evaluated using  $r^2$  between the true and

1 estimated values. For instance, for 1kb functional regions, the  $r^2$  values for  $f_0, f_1, f_2$  and  $f_3$  are  
2 0.93, 0.91, 0.89, and 0.87 respectively.

3 It should be noted that this approach does not distinguish between non-synonymous  
4 and synonymous mutations. Indeed, no assumption is made regarding which specific bases  
5 are neutral, nearly neutral, or deleterious in the coding region. Thus, this method can be used  
6 to estimate the DFE for any type of functional region, as well as to assess the non-neutrality  
7 of synonymous sites by comparing their frequency in a given coding region with the  
8 occupancy of the  $f_0$  class.

9

10 ***Effect of mis-specification of exon size and recombination rate:*** In view of these results, it  
11 is important to consider if accurate estimates depend on correctly specifying precise exon  
12 size, or whether it would be sufficient to generate priors assuming, for example, a mean exon  
13 length characterizing a genome. To quantify this effect, simulated data sampled from the  
14 priors was based on 1kb exons, while the test data were obtained from simulations based on  
15 alternative exon sizes. The error in inference of the DFE increases as the difference between  
16 exon sizes of the priors and that of the true sizes are increased (Supp Figure 3), with the  
17 highest in the moderately deleterious class ( $f_2$ ), although when exon sizes are sufficiently  
18 large, mis-specification of exon-size does not strongly impact performance. A similar  
19 approach was used to determine if the presence of another functional region (also 1kb in size)  
20 separated by an intron or intergenic region would skew inference. As expected, smaller intron  
21 sizes result in stronger mis-inference than larger ones, and intronic/ intergenic sizes larger  
22 than 4 kb performed essentially as well as those with no nearby functional exon (Supp Figure  
23 4). Moreover, a two-fold difference between assumed and actual recombination rates resulted  
24 in inflation of error dramatically (Supp Figure 5 and 6). Informatively, the direction of bias  
25 generated differs by DFE class (Supp Figure 6). For example, when true recombination rates  
26 are half of those assumed, the inferred weakly deleterious class is greatly inflated. As this  
27 class of mutations most strongly skews the linked neutral SFS, this mis-inference presumably  
28 arises from an attempt to fit stronger linked effects by inferring a higher proportion of  
29 mutations in this class, whereas in reality the increased BGS effects are being generated by  
30 fewer recombination events than are assumed.

31 These results highlight the importance of taking into account the specific exonic-  
32 intronic-intergenic structure of a particular genomic region of interest, nearby functional  
33 regions and the specific recombination rate. Although any configuration of these details may

1 be directly simulated, an alternative approach is simply to group exons of like size across a  
2 genome, and further reduce these to a group that is devoid of neighboring functional regions.

3

4

### 5 **Joint inference of purifying selection and demography, under non-equilibrium** 6 **conditions**

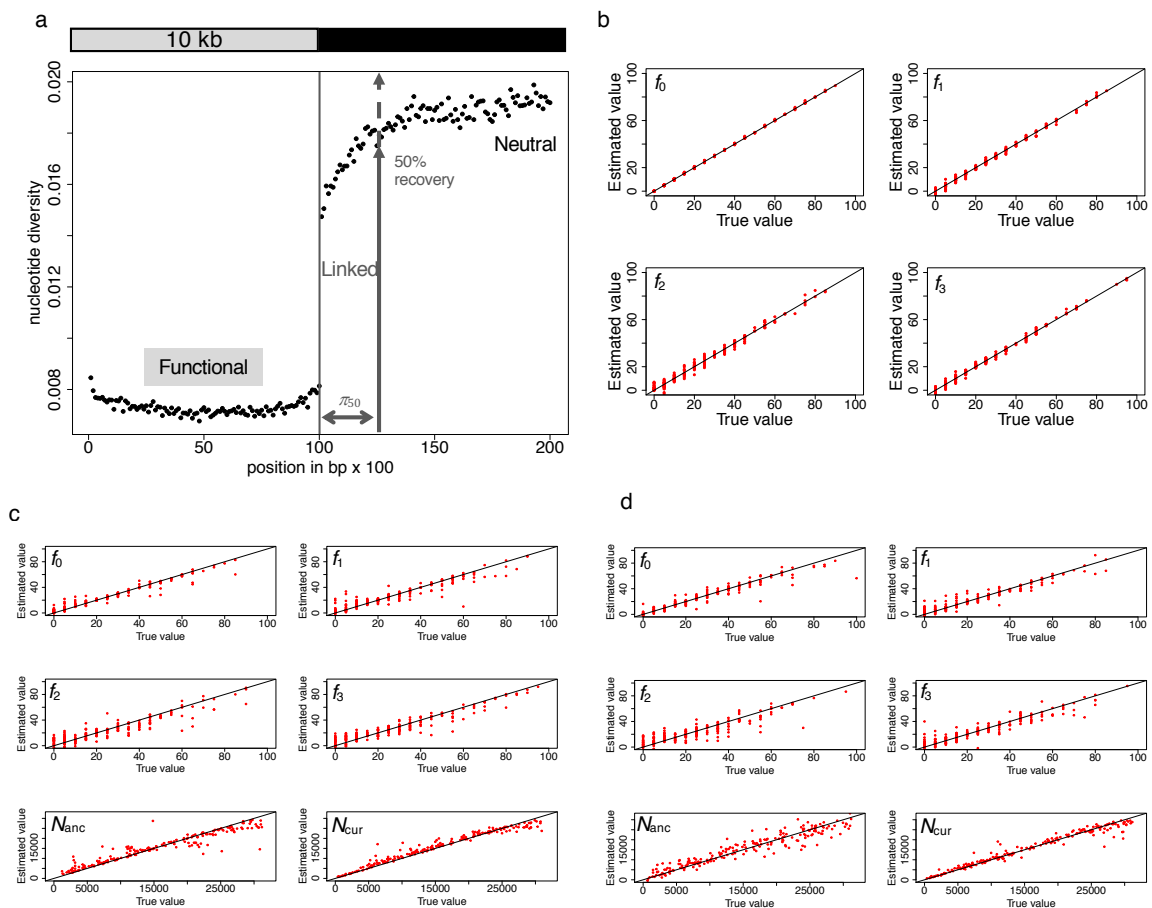
7 Based on the above results demonstrating that details of exon sizes and recombination rates  
8 are essential for accurate inference, we explicitly modeled both exon sizes and recombination  
9 rates when examining our ability to jointly infer demographic changes along with the DFE.  
10 As our example involved an African population of *D. melanogaster*, we chose single-exon  
11 genes that had more than 4kb non-coding regions flanking both sides and whose exon sizes  
12 were between 500-2000bp. For this specific set of 94 exons, we simulated functional regions  
13 with precise exon sizes linked to 4kb neutral regions and utilized the previously inferred local  
14 recombination rate for each exonic region in question. For every parameter combination, we  
15 performed 10 replicates of each of the 94 exon sizes (resulting in a total of 940 replicates per  
16 parameter combination), with their respective recombination rates and exon sizes, and  
17 summarized the resulting mean and variance of summary statistics.

18 Models of exponential population size expansion and contraction assumed various  
19 ancestral population size ( $N_{anc}$ ) and current population size ( $N_{cur}$ ), which were both sampled  
20 uniformly between  $10^5 - 10^7$ , as in previous studies (Duchen *et al.* 2013; Arguello *et al.*  
21 2019). As earlier work has inferred the duration of the expansion in Zambian populations to  
22 be of the order of  $\sim N_e$  generations, the time duration was scaled down and fixed to  $N_{sim}$   
23 (=5000 generations) in order to attempt to infer both historical and current population sizes.  
24 Thus, for this framework, we evaluated the estimates of six parameters:  $f_0, f_1, f_2, f_3, N_{anc}$ , and  
25  $N_{cur}$ .

26

27 **Accuracy of joint inference:** Encouragingly, the results demonstrated an ability to  
28 successfully co-estimate the DFE and both ancestral and current population sizes, using the  
29 set of coding and linked non-coding summary statistics described above (Figure 4c). Under  
30 non-equilibrium demography, the error in inference of the strongly deleterious class of  
31 mutations is larger. The absolute differences between true and estimated values were 0.019,  
32 0.027, 0.033, and 0.034 for the four DFE classes, respectively; the errors in ancestral and  
33 current sizes were 10.1% and 7.3% respectively. The  $r^2$  between the true and estimated  
34 values of  $f_0, f_1, f_2, f_3, N_{anc}$ , and  $N_{cur}$  were 0.97, 0.97, 0.95, 0.95, 0.99, and 0.99, respectively.

1            Nonetheless, the performance of the full 6-parameter estimation procedure is good,  
 2            without relying on the usual step-wise approach of first utilizing putatively neutral sites to  
 3            infer a demographic history, and then fixing that demographic history in order to estimate  
 4            DFE parameters. Interestingly, joint estimation is almost as accurate when using statistics  
 5            only from functional regions (Figure 4d), although it inflates the errors in the estimates of  $f_2$   
 6            and  $f_3$ . The absolute differences between the true and estimated values of  $f_0, f_1, f_2, f_3$  were  
 7            0.015, 0.025, 0.054, and 0.049, respectively, while the error in estimates of population sizes  
 8            increases to 23% and 8% for  $N_{anc}$  and  $N_{cur}$ , respectively. Although the error in ancestral  
 9            population size is quite large if only functional regions are used to co-estimate all six  
 10            parameters, the accuracy of inference can be improved significantly by adding more replicate  
 11            simulations of each parameter set to the ABC framework.  
 12



13  
 14  
 15 **Figure 4:** (a) Calculation of summary statistics across functional, linked and neutral regions.  
 16 (b) Accuracy of estimation (cross validation) of the four classes of the DFE using statistics  
 17 for functional regions only (size 1kb), under equilibrium demography. (c) Joint estimation of  
 18 population size changes and the DFE using all statistics. (d) Joint estimation of population  
 19 size changes and the DFE using statistics for functional regions only. The true proportions of  
 20 mutations in each DFE class and  $N_{anc}, N_{cur}$  are given on the X-axes, while the estimated

1 values are given on the Y-axes. Parameters are indicated on the upper left corner for each  
2 plot. Each dot represents one out of 200 different parameter combinations, sampled randomly  
3 from the entire set of simulations.

4

5

6

## 7 **Statistics that are important for distinguishing different classes of the DFE and** 8 **demography**

9 As it is important to understand which statistics may be necessary to distinguish between the  
10 effects of demography and the different classes of the DFE, two different approaches were  
11 used to rank statistics by their importance. First, statistics were simply ranked by their  
12 regression coefficient with respect to each parameter separately. Non-linear relationships  
13 were taken into account by using Box-Cox transformation, as suggested by Wegmann *et al.*  
14 (2009). With stationary population size, most of the top predictors of the fraction of neutral  
15 ( $f_0$ ) and strongly deleterious ( $f_3$ ) sites are statistics summarizing the functional region (Supp  
16 Table 3). The top four statistics for each parameter are displayed in Supp Figure 7. In  
17 addition, a modified method of Joyce and Marjoram (2008) was also employed to rank  
18 statistics (Supp Table 4) for equilibrium demography.

19 As expected, statistics that correlate most strongly with the fraction of neutral  
20 mutations are levels of divergence and the fraction of high frequency derived alleles, as  
21 summarized by  $\theta_H$  (Fu 1995; Fay and Wu 2000) in functional regions. As the weakly  
22 deleterious class of mutations generate BGS effects at closely linked sites, statistics in the  
23 functional and linked region are most strongly correlated with  $f_1$ . This also correlates most  
24 with  $H'$  in functional regions, a statistic that contrasts the proportion of high frequency  
25 derived variants with those of derived variants segregating at intermediate frequency (Fay  
26 and Wu 2000). Although this statistic was designed to identify selective sweeps, which may  
27 result in a larger proportion of high frequency derived alleles, it is highly predictive of the  
28 fraction of weakly deleterious class of mutations in the absence of positive selection. As  
29 shown previously, larger  $f_1$  generates a stronger skew in the linked neutral SFS towards rare  
30 variants and is thus also reflected in values of Tajima's  $D$  in the linked neutral region.  
31 Measures of linkage disequilibrium in the functional and linked neutral regions are also  
32 correlated with the weakly deleterious class of mutations.

33 Because the moderately deleterious class of mutations generates BGS effects that  
34 extend for larger distances than the more weakly selected class, the strongest correlates of

1 this class are generally statistics from the neutral region furthest from the directly selected  
2 sites. All the different summaries of the SFS -  $\theta_W$ ,  $\theta_\pi$ , and  $\theta_H$  - correlate with this parameter,  
3 as well as the total reduction in linked neutral diversity (given by the intercept of the  
4 regression fit of  $\pi = slope * \ln(distance) + intercept$ , where  $\pi$  is the diversity in linked neutral  
5 regions). The strongly deleterious class of mutations is correlated with the number of  
6 singletons and  $\theta_W$ , which is highly sensitive to singletons.

7 A similar analysis was performed on simulations under models of demographic non-  
8 equilibrium. Here, the DFE parameters are significantly correlated only with the statistics for  
9 functional regions (Supp Table 5 and 6). As expected intuitively, the statistics most highly  
10 correlated with the two demographic parameters are for the neutral linked regions. Ancestral  
11 population sizes correlate most with statistics that capture high frequency derived alleles in  
12 linked neutral as well as functional regions, as these represent older mutations; current  
13 population sizes correlate most with statistics that summarize LD. The same is true when  
14 ranked statistics are obtained only from functional regions. Because the class of moderately  
15 deleterious mutations and ancestral populations sizes are correlated with overlapping sets of  
16 statistics, the estimation of these two parameters is partially confounded. As such, LD-based  
17 statistics are essential in distinguishing between demography and purifying selection, and in  
18 distinguishing between ancestral and current population sizes. In addition, although the  
19 variances and means of the statistics are highly correlated, the variances play a more  
20 important role in estimating current population sizes.

21

## 22 **Comparison with DFE-alpha**

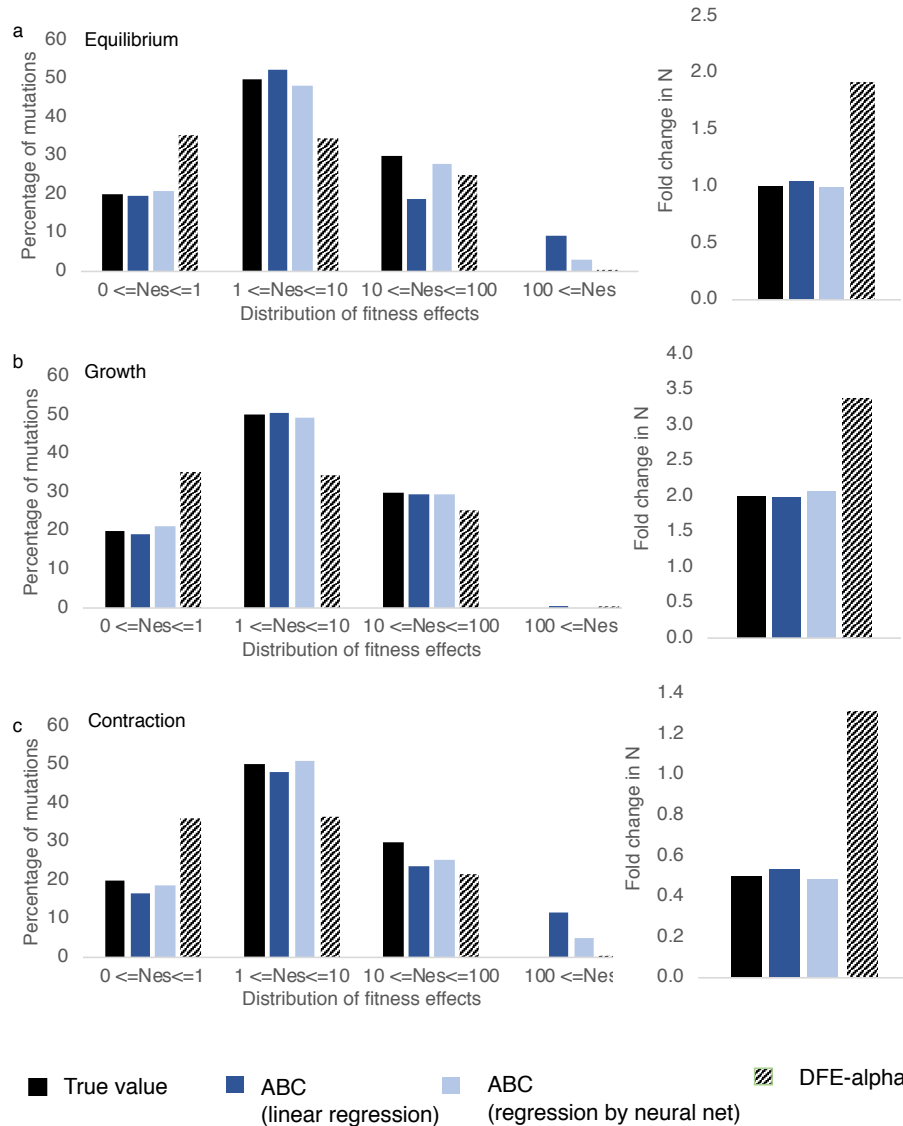
23 Although there are no other programs that simultaneously co-estimate both demographic and  
24 selection parameters, we compared the performance of our method to the step-wise approach  
25 of DFE-alpha (Keightley and Eyre-Walker 2007; Eyre-Walker and Keightley 2009;  
26 Schneider *et al.* 2011), a program used widely for the inference of the DFE. DFE-alpha  
27 assumes that synonymous sites are neutral and uses their site frequency spectrum to infer  
28 changes in demography. Conditional on the inferred demography and under the assumption  
29 that the deleterious selection coefficients follow a given distribution (generally gamma), the  
30 program infers the shape and rate parameter of the assumed distribution. We simulated  
31 demographic equilibrium, 2-fold population growth and 2-fold population contraction, and  
32 inferred the change in population size as well as the DFE using both ABC and DFE-alpha.  
33 Because DFE-alpha uses neutral sites to infer demography, in all cases we simulated a DFE

1 consisting of ~30% neutral mutations, which were used as a proxy for synonymous sites.  
2 These simulations were performed exactly as described previously for non-equilibrium  
3 conditions. Exons sizes between 500-2000 bp with flanking 4 kb linked neutral regions were  
4 simulated with recombination rates specific to the selected 94 exons (and a total of 940  
5 replicates for every parameter combination). DFE-alpha performs slightly better than ABC if  
6 the true DFE is indeed gamma distributed (Supp Figure 8) although our method is able to  
7 infer the DFE with very similar accuracy.

8 For a discrete DFE which is skewed towards highly deleterious mutations, DFE-alpha  
9 and ABC perform with similar accuracy. However, our method performs better if the DFE is  
10 skewed towards slightly deleterious mutations (*i.e.*, class  $f_1$ ) as shown in Supp Figure 9. It is  
11 important to note that, for the purpose of this comparison, simulations were run with numbers  
12 of directly selected sites between 500–2000 bp and ensuring that 30% of mutations were  
13 neutral, as the neutral mutations were required to estimate demography by DFE-alpha. Under  
14 these conditions, background selection results in a relatively small skew in the neutral SFS  
15 (see Campos and Charlesworth 2019).

16 As noted previously, a potential advantage of the methodology proposed here is that,  
17 by simultaneously estimating selection and demography, one is not required to make any  
18 assumptions about the neutrality of synonymous sites. We evaluated this feature by  
19 simulating a scenario where ~33% of the assumed neutral sites were actually experiencing  
20 weak direct selection. As weak purifying selection generates a larger fraction of rare variants  
21 than stronger selection, programs based on neutrality would be likely to falsely infer growth.  
22 As expected, DFE-alpha inferred 2-fold growth under demographic equilibrium, and in fact  
23 inferred slight growth even for a 2-fold contraction (Figure 5). The resulting DFE over-  
24 estimated the fraction of neutral mutations and under-estimated the fraction of weakly  
25 deleterious mutations. As noted previously, such mis-inference will increase with the density  
26 of selected sites. Our ABC approach, however, accurately estimated the proportion of neutral  
27 mutations present in the selected region (Figure 5), illustrating the importance of such joint  
28 inference.

29



1  
2 **Figure 5:** Comparison of the performance of the proposed ABC approach in the current study  
3 with DFE-alpha, under (a) demographic equilibrium, (b) exponential growth, and (c)  
4 exponential decline. In all cases, 30% of sites were assumed to be synonymous, out of which  
5 33% were weakly selected. Solid black bars are the true simulated values, dark blue bars give  
6 the ABC performance using ridge regression, and light blue bars give the ABC performance  
7 using linear regression aided by neural nets. Patterned bars show the performance of DFE-  
8 alpha. A total of 998,300 sites were analyzed in the functional region for each parameter  
9 combination, with approximately 332,767 representing synonymous and 665,533  
10 representing nonsynonymous sites.

11

12

13

#### 14 **Application to *Drosophila melanogaster***

15 The proposed method is well-suited to compact genomes, in which most sites may be  
16 experiencing purifying selection, but is computationally intensive for large genomic regions.



1 When simulating small genomic regions, the presence of nearby coding regions that are not  
2 included in the models can generate additional BGS effects and thus bias inference. We thus  
3 restricted our analyses to protein-coding exons in the *D. melanogaster* genome between 500  
4 to 2000 bp in length that are single exon genes, and are flanked on both sides by intergenic  
5 regions larger than 4 kb (with the latter two criteria chosen to avoid strong effects of linkage  
6 with other nearby functional elements). It should be noted that any genic structure could  
7 readily be chosen for inference by directly simulating the associated details when  
8 constructing the priors - we have simply chosen this realization in order to provide an  
9 illustrative application.

10 The recombination rates of both the 5' and 3' flanking intergenic regions are highly  
11 correlated (Supp Figure 10) and span a considerable magnitude (Supp Figure 11), with a  
12 mean rate of 2.21 cM/Mb (*i.e.*, the average recombination rate for these chosen single exon  
13 genes is very near the autosomal genome-wide average of 2.32 cM/Mb). We also verified  
14 that this set of genes was not unusual with regards to genome-wide coding sequence  
15 divergence (Supp Figure 12). Furthermore, because sites in intergenic regions in *D.*  
16 *melanogaster* may also experience direct selection (Halligan and Keightley 2006; Casillas *et*  
17 *al.* 2007), we used phastCons scores to exclude intergenic sites that may potentially be  
18 functionally important. All sites with a phastCons score larger than 0.8 were excluded (Siepel  
19 *et al.* 2005). Table 1 provides the observed summary statistics for each region class, where  
20 intergenic sites that had a greater than or equal to 80% probability of belonging to a  
21 conserved element (*i.e.*, with phastCons score  $\geq 0.8$ ) were excluded. It should be noted that,  
22 although there does not appear to be a large difference between divergence (*i.e.*, number of  
23 fixed substitutions specific to *D. melanogaster*) in exonic vs intergenic regions, this  
24 observation is consistent with previous studies (Table 1 in Andolfatto 2005). In addition,  
25 because we have restricted our analyses to sites where the ancestor of *D. melanogaster* could  
26 be predicted with high confidence, our analyses may be skewed towards more conserved  
27 sites, potentially resulting in lower divergence in intergenic regions. Previous estimates of  
28 divergence at 4-fold degenerate sites have been estimated to be roughly 0.05-0.06 (Halligan  
29 and Keightley 2006; Langley *et al.* 2012; Charlesworth *et al.* 2018), while that in coding  
30 regions to be 0.023 (Langley *et al.* 2012). Although our estimates are lower than previous  
31 estimates, this discrepancy is well explained by the larger number of individuals used to  
32 subtract polymorphic sites in this study (Supp Table 7). At a sample size of 1 individual  
33 (corresponding to pairwise divergence), our estimates of divergence at 4-fold degenerate sites  
34 is 0.05 and in coding regions is 0.023, consistent with previous studies. In addition, a very

1 similar reduction between pairwise divergence and polymorphism-adjusted divergence is  
2 observed in simulated data (Supp Table 8).

3 Interestingly, although previous studies have inferred ~2-fold growth in the Zambian  
4 population of *D. melanogaster* (Li and Stephan 2006; Laurent *et al.* 2011; Duchen *et al.*  
5 2013; Kapopoulou *et al.* 2018; Arguello *et al.* 2019), we infer only a 1.2-fold growth, with an  
6 ancestral  $N_e$  of 1,225,393 and current  $N_e$  of 1,357,760. In contrast to previous studies  
7 (Keightley and Eyre-Walker 2007; Huber *et al.* 2017), we infer a much larger proportion of  
8 mildly deleterious mutations and a smaller proportion of highly deleterious mutations (see  
9 Figure 6), with  $f_0 = 24.7\%$ ,  $f_1 = 49.4\%$ ,  $f_2 = 3.9\%$ , and  $f_3 = 21.9\%$  - but this reflects the fact  
10 that our procedure includes the possibility of selection on synonymous sites. As we have  
11 inferred the DFE for a select class of single exon genes, genes which have slightly higher  
12 divergence than average (Supp Figure 12), it is possible that these exons are experiencing  
13 weaker purifying selection compared to the genome-wide mean. Furthermore, because we  
14 have obtained the DFE of both coding sequences and UTR regions, 4-fold degenerate sites  
15 represent 12% of all sites, while UTR regions comprise 29% of all sites. Previous studies  
16 have estimated that roughly 6-10% of all mutations at non-synonymous sites may be  
17 effectively neutral. Thus, assuming that all 4-fold degenerate sites are neutral, ~40% of UTR  
18 regions are neutral (Andolfatto 2005; Campos *et al.* 2017), and ~6-10% of nonsynonymous  
19 mutations are neutral, we expect  $f_0$  to be ~27-30%. Encouragingly, we infer  $f_1 = 25\%$ . This  
20 observation implies that the majority of synonymous sites are not experiencing direct  
21 selection, consistent with previous results for *D. melanogaster* (Jackson *et al.* 2017).

22 In order to confirm whether our inferred parameters explain the observed *D.*  
23 *melanogaster* data, we simulated 10 replicates of each of the 94 exons using the parameter  
24 estimates, and evaluated whether the mean of the observed *D. melanogaster* values are in the  
25 5% tails of the distribution of statistics obtained via simulations. Our parameter estimates  
26 result in a very good fit to empirical *D. melanogaster* population data (Figure 6, Supp Figure  
27 13) for all three categories - functional (*i.e.*, exonic), linked (*i.e.*, non-coding region adjacent  
28 to exons) and neutral (*i.e.*, non-coding region adjacent to the linked region). Our parameter  
29 estimates fail, however, to explain the observed Tajima's  $D$  values (linked region  $p = 0.011$ ,  
30 neutral region  $p = 0.010$ ) and divergence (linked region  $p = 0.029$ , neutral region  $p=0.0$ ) in  
31 intergenic regions – though both are well fit in functional regions.

32 As both positive selection in exons and purifying selection in non-coding regions  
33 could partially drive these patterns, we investigated both of these model violations. Non-  
34 coding regions flanking 2kb of the selected exons (which were used to perform inference)

1 were found to have 777 sites that had phastCons scores greater than or equal to 0.8, with a  
2 mean and median length of 25 and 15 bp, respectively. We therefore simulated conserved  
3 elements in non-coding regions that were 20bp in length, uniformly distributed, and which  
4 made up 40% of the flanking neutral sites (*i.e.*, 800 sites in total). Conserved elements were  
5 simulated with weak ( $f_1 = 100\%$ ), moderate ( $f_2 = 100\%$ ) and strong ( $f_3 = 100\%$ ) purifying  
6 selection separately. Upon masking these sites, as was done in our *Drosophila* data analysis,  
7 there was no observed difference in the distribution of all statistics (Supp Figure 14),  
8 suggesting that background selection caused by small conserved elements does not  
9 significantly affect our inference, and in fact does not alter the fit of our inferred model to the  
10 data. Interestingly, without masking sites – that is, allowing sites that experience direct weak  
11 purifying selection to remain in the flanking sequence - our model is much better able to  
12 explain a lower Tajima's  $D$  and divergence in intergenic regions (Supp Figure 15). Thus, it  
13 appears likely that unaccounted-for weak purifying selection across multiple sites in  
14 intergenic regions could contribute to the discrepancy between statistics generated by our  
15 model and those observed in the data.

16 Next, we simulated positive selection under 4 different scenarios - representing rare  
17 and strong (1% of all mutations in exonic regions are beneficial with  $2N_e s = 1000$ ), common  
18 and strong (5% of mutations in exonic regions are beneficial with  $2N_e s = 1000$ ), common and  
19 weak (5% of mutations in exonic regions are beneficial with  $2N_e s = 10$ ) and rare and weak  
20 (1% of mutations in exonic regions are beneficial with  $2N_e s = 10$ ) selection. Interestingly, we  
21 find that, although strong positive selection, whether common or rare, better explains the  
22 lower Tajima's  $D$  values in intergenic regions, it also drastically alters the distribution of  
23 most other statistics, resulting in an overall much poorer fit (Supp Figure 16, 17). For  
24 instance, common and strong positive selection reduces  $\theta_H$  by an order of magnitude relative  
25 to our fitted model and drastically increases the variance while decreasing the mean of  
26 haplotype diversity. In contrast to strong positive selection, weakly positively selected  
27 mutations do not alter the distribution of Tajima's  $D$  in intergenic regions, but do slightly  
28 increase  $\theta_H$  in functional regions, which improves the fit to the observed data (Supp Figure  
29 18, 19). In addition, all cases of positive selection significantly increase divergence in  
30 functional regions. For comparison, we also simulated the two scenarios of positive selection  
31 used by Lange and Pool (2018) - 0.2% of all mutations are beneficial with  $2N_e s = 60$ , and  
32 0.00013% of all mutations are beneficial with  $2N_e s = 10000$ . As the frequency of positively  
33 selected alleles is lower in these scenarios, there was no observed difference between the  
34 distribution of statistics resulting from including or excluding positive selection (Supp Figure

1 20, 21). Thus, if the frequency of strongly positively selected mutations is much lower than  
 2 1%, our estimates of both demography and DFE shape should be unbiased, and the beneficial  
 3 fixations would be virtually undetectable. Future studies will further investigate the ability of  
 4 such an approach to quantify the occupancy of a beneficial mutational class.

5

6

7

8

9

10

11 **Table 1:** Statistics calculated for the 94 single-exon genes including their 3' flanking  
 12 intergenic sequences, for 76 individuals (devoid of any inversion) in the African Zambian  
 13 population. Sites with phastCons scores higher than 0.8 were excluded. Functional refers to  
 14 exons, linked refers to intergenic region (~ 1kb) adjacent to exons and neutral refers to  
 15 intergenic regions further away from exons that are adjacent to linked regions and ~1kb in  
 16 size (Figure 4a). Derived alleles were identified by polarizing alleles with respect to the  
 17 ancestral sequence of *D. melanogaster* obtained from ancestral reconstruction over 15 insect  
 18 species.

19

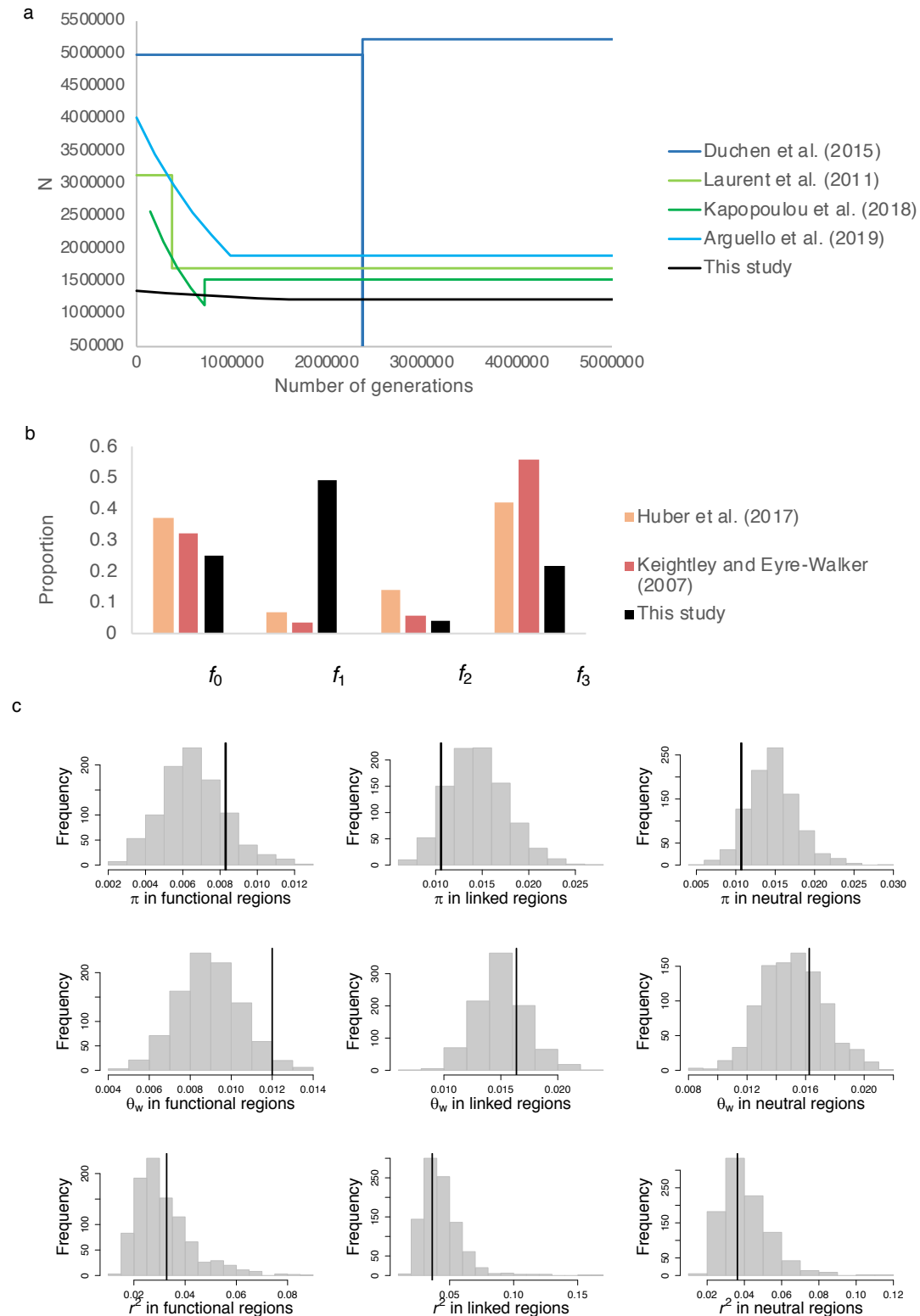
	mean			standard deviation		
	functional	linked	neutral	functional	linked	neutral
$\pi$	0.0083	0.0106	0.0107	0.0039	0.0042	0.0038
$\theta_w$	0.0120	0.0166	0.0162	0.0045	0.0053	0.0049
$\theta_H$	0.0088	0.0098	0.0097	0.0054	0.0053	0.0056
$H'$	-0.0633	0.0871	0.1169	0.5371	0.4118	0.3829
Tajima's $D$	-1.0537	-1.1469	-1.1103	0.5338	0.4874	0.4694
Singleton density	0.0215	0.0303	0.0307	0.0086	0.0116	0.0117
Haplotype diversity	0.9711	0.9680	0.9762	0.0452	0.0458	0.0444
$r^2$	0.0328	0.0364	0.0363	0.0109	0.0136	0.0128
$D$	0.0005	0.0005	0.0006	0.0009	0.0010	0.0012
Branch-specific divergence	0.01378	0.0156	0.0159	0.0075	0.0077	0.0071

20

21

22

23



1

2

3

4

5

6

7

**Figure 6:** Joint inference of demography and purifying selection in the Zambian population of *D. melanogaster*. (a) Demographic model inferred from previous studies (colored lines) and from the current study (black lines). (b) The distribution of fitness effects of deleterious mutations at coding regions (including synonymous and non-synonymous sites) as inferred by previous studies of other populations (colored bars) and at exonic sites of single-exons

1 genes as inferred by the current study (black bars). The X-axis is given by  $f_0$ :  $0 < |2N_e s| < 1$ ,  $f_1$ :  
2  $1 < |2N_e s| < 10$ ,  $f_2$ :  $10 < |2N_e s| < 100$ , and  $f_3$ :  $100 < |2N_e s| < 10000$ . For the previous studies, the DFE  
3 shown in this figure includes the fraction of synonymous sites in the neutral  $f_0$  class. (c)  
4 Distribution of key summary statistics ( $\pi$ ,  $\theta_w$ ,  $r^2$ ) in functional, linked and neutral regions  
5 upon simulating 100 replicates of 94 exons each under the inferred parameters. The vertical  
6 lines represent the values of the statistics obtained from 76 individuals of *D. melanogaster*  
7 from Zambia, after excluding non-coding sites with phastCons score  $\geq 0.8$ .

8

9

10

## 11 CONCLUSION

12 Independent of any dogmatic stance regarding the roles of adaptive vs. non-adaptive  
13 explanations for observed levels and patterns of DNA sequence variation and divergence, it  
14 has been widely accepted that natural populations are not at demographic equilibrium, but are  
15 often characterized by fluctuating population sizes and other demographic perturbations.  
16 Additionally, a rich empirical and experimental literature has clarified the pervasive  
17 importance of purifying selection in eliminating the constant input of deleterious variants. It  
18 has also been well-demonstrated that ignoring direct effects of purifying selection and its  
19 impact on linked sites can strongly bias demographic inference (Ewing and Jensen 2016), and  
20 that ignoring demographic effects biases estimates of parameters of selection (Jensen *et al.*  
21 2005; Thornton and Jensen 2007; Crisci *et al.* 2012, 2013). Yet, despite agreement that these  
22 processes are certain to be occurring constantly in populations and shaping patterns of  
23 variation and evolution, the construction of a statistical approach capable of simultaneously  
24 estimating parameters of the concerned processes has proven challenging. Here we provide  
25 one such approach, for which we demonstrate an ability to co-estimate the parameters of a  
26 generalized DFE along with those underlying the population history.

27

28 By fitting a four-parameter DFE model that includes weak, intermediate and strong purifying  
29 selection, as well as neutrally evolving sites, this approach avoids two common, and  
30 potentially perilous, assumptions: 1) synonymous sites are not assumed to be neutral,  
31 consistent with a growing body of literature (Chamary and Hurst 2005; Lynch 2007; Zeng  
32 and Charlesworth 2010a; Lawrie *et al.* 2013; Choi and Aquadro 2016; Jackson *et al.* 2017),  
33 and 2) the DFE is not assumed to follow a specific parameterized distribution, such as the  
34 widely-used gamma distribution.

35 Our results demonstrate that it is possible to jointly infer the deleterious DFE and past  
36 demographic changes using an ABC framework, by including various summary statistics

1 capturing aspects of the SFS, linkage disequilibrium and divergence, compared between  
2 coding and flanking non-coding sequence. Ancestral population sizes and the frequency of  
3 the most deleterious classes of the DFE are estimated with relatively low accuracy, whereas  
4 the current population sizes and the neutral mutation class are estimated with high accuracy.  
5 In addition, we demonstrated that, if synonymous sites are indeed experiencing substantial  
6 purifying selection, existing programs such as DFE-alpha will over-estimate recent growth  
7 and under-estimate the proportion of mildly deleterious mutations. Importantly, the approach  
8 proposed here performs equally well regardless of whether synonymous sites are neutral or  
9 selected. However, our approach continues to assume the neutrality of flanking non-coding  
10 regions, though putatively conserved sites were masked, and the impact of that masking on  
11 inference was thoroughly assessed via simulation.

12 Because we make no assumptions about which sites in the functional region of  
13 interest are neutral, it is in principle possible to estimate the DFE for any functional element  
14 using this methodology, including regulatory elements or functional regions with  
15 interdigitated sites experiencing direct selection. The results further suggest that the accurate  
16 co-estimation of these parameters is possible using only functional regions. Such an approach  
17 may be extremely useful in genomes for which it is difficult to characterize putatively neutral  
18 sites, as well as for compact genomes in which non-coding regions may be limited.

19 This approach can in principle be applied to any organism and functional class of  
20 interest, although power analyses suggest the utility of prior knowledge of the boundaries of  
21 functional regions and recombination rates. Here we have provided an illustrative example in  
22 *D. melanogaster*. The results suggest that the Zambian population has been largely stable in  
23 size, and that exonic regions have a large proportion of mildly deleterious mutations.

24 Although this result might seem surprising, the DFE inferred by the current method provides  
25 the distribution of selective effects over all sites, including synonymous sites and sites in  
26 UTRs. Hence, in comparing the DFE estimated in the current study with previous estimates  
27 of the neutral class of mutations, it appears unnecessary to invoke widespread selection on  
28 synonymous sites in *D. melanogaster*. This result is largely consistent with most previous  
29 studies (Akashi 1995; Jackson *et al.* 2017). For instance, our estimate of the strength of  
30 purifying selection acting on synonymous sites in the Zambian population is in line with  
31 earlier estimates for African populations (Zeng and Charlesworth 2010a; Jackson *et al.*  
32 2017).

33 In addition to the proposed inference framework, we have derived an analytical  
34 expression for the reduction in variation caused by background selection at neutral sites

1 outside functional regions for the case of a discrete DFE, making it feasible to obtain  
2 analytical predictions for any chosen DFE. Not only does a discrete DFE provide flexibility  
3 in inference, it may also be a more realistic representation of the true DFE (Kousathanas and  
4 Keightley 2013; Bank *et al.* 2014b). Although gamma distributions represent a reasonable fit  
5 to the DFE inferred from genome-wide studies (Eyre-Walker and Keightley 2007), the DFE  
6 will be mis-inferred if the true distribution is multimodal (Kousathanas and Keightley 2013),  
7 as has been observed widely (*e.g.*, in yeast (Bank *et al.* 2014a), viruses (Sanjuán 2010), and  
8 *E.coli* (Jacquier *et al.* 2013)). In addition, the best fitting parameterized continuous  
9 distribution appears to be extremely specific to the particular dataset being tested, and most  
10 alternative distributions fit the data nearly as well as the best fitting-distribution (Huber *et al.*  
11 2017; Kim *et al.* 2017). The discrete DFE proposed here thus reduces the number of  
12 necessary assumptions, and has been shown to perform well in the plausible scenario in  
13 which common assumptions are indeed violated (*e.g.*, if the true DFE is not gamma-  
14 distributed). Analytical results under demographic equilibrium and simulations under  
15 demographic non-equilibrium stress that the number of selected sites and the specific shape  
16 of the DFE (for instance the presence of mildly and moderately deleterious mutations) both  
17 decrease linked neutral variation around functional regions more than previously appreciated,  
18 and skew the SFS even when there is no reduction in diversity. Such variation in exon lengths  
19 and DFE shapes across a genome can increase variance of statistics in linked neutral regions,  
20 which could contribute to false positives when detecting positive selection using outlier  
21 approaches.

22         There are at least two important caveats worth considering, which will be the subject  
23 of future study. The first concerns the estimates of ancestral and current effective population  
24 sizes. As the effective population size varies across the genome in a fashion correlated with  
25 local recombination rates (Becher *et al.* 2020, in press), the estimates provided here ought to  
26 be viewed as a mean across the loci in question. While we have improved upon the common  
27 assumption of a singular genome-wide value by directly modeling each locus-specific  
28 recombination rate when performing inference, the general importance of this effect in  
29 demographic modeling remains in need of further study. The second concerns the inference  
30 of selection. This study represents a proof-of-concept in demonstrating that such  
31 simultaneous inference of demography and the DFE is feasible, thereby avoiding common  
32 assumptions underlying a step-wise inference approach. While this interplay of genetic drift  
33 and purifying selection is in fact sufficient alone to fit all features of the data (consistent with  
34 previous claims: Comeron 2014, 2017; Harris *et al.* 2018; Jensen *et al.* 2019), this is not the



1 same as claiming that positive selection is not also occurring. As our simulation results  
2 demonstrate, the addition of rare, weakly beneficial mutations is consistent with the data,  
3 though the inclusion of these parameters does not result in an improved fit. The question is  
4 less about presence/absence, than it is about statistical identifiability. Conversely, the  
5 addition of a strongly beneficial mutational class was found to be inconsistent with observed  
6 data. In order to investigate this further, future work will evaluate the ability to co-estimate a  
7 beneficial class of fitness effects within this framework. It should also be noted that the  
8 example chosen to highlight our approach focuses on only a subset of genes in the *D.*  
9 *melanogaster* genome, and there is no expectation that the observed DFE in this class will  
10 necessarily be universal across all coding regions in the population under consideration. In  
11 fact, means of scaled selection coefficients of deleterious mutations have been shown to be  
12 negatively correlated with divergence at nonsynonymous sites (Campos *et al.* 2017).  
13 Importantly however, this general inference approach accounting for these two dominant  
14 processes will be a valuable tool in future genomic scans, and this appropriate null is  
15 anticipated to greatly reduce the notoriously high false-positive rates associated with the  
16 identification of positively selected loci.

17

18

## 19 **ACKNOWLEDGEMENTS**

20 We would like to thank Rebecca Harris for discussions related to this project. This research  
21 was conducted using resources provided by Research Computing at Arizona State University  
22 (<http://www.researchcomputing.asu.edu>) and the Open Science Grid, which is supported by  
23 the National Science Foundation and the U.S. Department of Energy's Office of Science. We  
24 especially thank Lauren Michael and Christina Koch from the Open Science Grid, for all of  
25 their efforts to provide technical assistance. This work was funded by National Institutes of  
26 Health grant R01GM135899 to JDJ.

27

28

## 29 **DISCLOSURE DECLARATION**

30 The authors declare no conflicts of interests.

31

32

33

34

1 **APPENDIX**

2

3 **Derivation of the analytical expression for the reduction in diversity due to background**  
 4 **selection generated by a discrete distribution of fitness effects under demographic**  
 5 **equilibrium**

6 Because we model a discrete DFE with four fixed bins, with  $t$  being uniformly distributed  
 7 within each bin, the definite integral for each bin is the integral of  $E(t)$  with respect to  $t$ . We  
 8 thus have:

9 
$$\int_0^1 E(t)dt = \frac{f_0}{t_1-t_0} \int_{t_0}^{t_1} E(t)dt + \frac{f_1}{t_2-t_1} \int_{t_1}^{t_2} E(t)dt + \frac{f_2}{t_3-t_2} \int_{t_2}^{t_3} E(t)dt + \frac{f_3}{t_4-t_3} \int_{t_3}^{t_4} E(t)dt \quad (A1)$$

10

11 where the  $t_i$ 's correspond to the boundary of the discrete bins. The first integral is over  $t$  such  
 12 that  $0 \leq 2N_e s \leq 1$ , the second over  $t$  such that  $1 \leq 2N_e s \leq 10$ , the third such that  
 13  $10 \leq 2N_e s \leq 100$  and fourth as  $100 \leq 2N_e s \leq 10000$ . In our case,  $t_0 = 0$ ,  $t_1 = 0.00005$ ,  $t_2 =$   
 14  $0.0005$ ,  $t_3 = 0.005$ , and  $t_4 = 0.5$ . While this mirrors the DFE considered here, the same  
 15 procedure can be done for any set of bins for a given DFE.

16 Integrating  $E$  over a uniform distribution between  $t_0$  and  $t_1$ , with probability density ( $t_1$   
 17  $- t_0$ )<sup>-1</sup>, where  $a = g + r_c y$  and  $b = g + r_c(y + l)$ , we have:

18

19 
$$\int \frac{t dt}{(1-t)[a+t(1-a)]} = \int \left\{ \frac{t}{(1-t)} + \frac{t(1-a)}{[a+t(1-a)]} \right\} dt \quad (A2)$$

20

21 The second integral on the right-hand side of this equation can be evaluated by  
 22 substituting  $u = a + t(1 - a)$  for  $t$ , with  $t = (u - a)/(1 - a)$  and  $dt = du/(1 - a)$ . This gives:

23

24 
$$(1 - a) \int \frac{t(1-a)}{[a+t(1-a)]} dt = (1 - a)^{-1} \int u^{-1} (u - a) du = (1 - a)^{-1} [u - a \ln(u)] \quad (A3)$$

25

26 With the change in variable, the normalizing factor for the probability density  
 27 function is now  $(u_1 - u_0)^{-1} = (1 - a)^{-1}(t_1 - t_0)^{-1}$ . The contribution of this component to the  
 28 expectation of  $E(t)$  over the uniform distribution yields equation (3b) of the main text. A  
 29 similar expression can be written for the integral of  $-t/[(1 - t)[b + t(1 - b)]]$  in the first line of  
 30 equation (2). When adding this to the integral of  $t/[(1 - t)[a + t(1 - a)]]$ , the integrals involving  
 31  $1/(1 - t)$  cancel out, so this term simply contributes the following term to the expectation of  
 32  $E(t)$ , yielding equation (3b). The expectation of  $E(t)$  is the sum of these two terms.  $E(t)$  can

1 also be numerically integrated over a definite interval as specified above, with constant  
2 values as chosen in our simulations:  $r = 10^{-6}$ ,  $l = 500$  or  $1000$  or  $5000$  or  $10000$ ,  $U = l \times \mu$ ,  $\mu =$   
3  $10^{-6}$ , and  $g = 0$ .

4  
5  
6  
7  
8

## 7 REFERENCES

- 9 Akashi H., 1995 Inferring weak selection from patterns of polymorphism and divergence at  
10 “silent” sites in *Drosophila* DNA. *Genetics* 139: 1067–1076.
- 11 Andolfatto P., 2005 Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* 437:  
12 1149–1152. <https://doi.org/10.1038/nature04107>
- 13 Arguello J. R., S. Laurent, and A. G. Clark, 2019 Demographic history of the human  
14 commensal *Drosophila melanogaster*. *Genome Biol Evol* 11: 844–854.  
15 <https://doi.org/10.1093/gbe/evz022>
- 16 Assaf Z. J., S. Tilk, J. Park, M. L. Siegal, and D. A. Petrov, 2017 Deep sequencing of natural  
17 and experimental populations of *Drosophila melanogaster* reveals biases in the  
18 spectrum of new mutations. *Genome Res.* 27: 1988–2000.  
19 <https://doi.org/10.1101/gr.219956.116>
- 20 Bank C., R. T. Hietpas, A. Wong, D. N. Bolon, and J. D. Jensen, 2014a A Bayesian MCMC  
21 approach to assess the complete distribution of fitness effects of new mutations:  
22 uncovering the potential for adaptive walks in challenging environments. *Genetics*  
23 196: 841–852. <https://doi.org/10.1534/genetics.113.156190>
- 24 Bank C., G. B. Ewing, A. Ferrer-Admettla, M. Foll, and J. D. Jensen, 2014b Thinking too  
25 positive? Revisiting current methods of population genetic selection inference. *Trends*  
26 *in Genetics* 30: 540–546. <https://doi.org/10.1016/j.tig.2014.09.010>
- 27 Beaumont M. A., W. Zhang, and D. J. Balding, 2002 Approximate Bayesian Computation in  
28 Population Genetics. *Genetics* 162: 2025–2035.
- 29 Becher H., B. C. Jackson, and B. Charlesworth, 2020 Patterns of genetic variability in  
30 genomic regions with low rates of recombination. *bioRxiv* 739888 (in press).  
31 <https://doi.org/10.1101/739888>
- 32 Campos J. L., L. Zhao, and B. Charlesworth, 2017 Estimating the parameters of background  
33 selection and selective sweeps in *Drosophila* in the presence of gene conversion.  
34 *PNAS* 114: E4762–E4771. <https://doi.org/10.1073/pnas.1619434114>
- 35 Campos J. L., and B. Charlesworth, 2019 The effects on neutral variability of recurrent  
36 selective sweeps and background selection. *Genetics* 212: 287–303.  
37 <https://doi.org/10.1534/genetics.119.301951>

- 1 Casillas S., A. Barbadilla, and C. M. Bergman, 2007 Purifying selection maintains highly  
2 conserved noncoding sequences in *Drosophila*. *Mol. Biol. Evol.* 24: 2222–2234.  
3 <https://doi.org/10.1093/molbev/msm150>
- 4 Chamary J., and L. D. Hurst, 2005 Evidence for selection on synonymous mutations affecting  
5 stability of mRNA secondary structure in mammals. *Genome Biology* 6: R75.  
6 <https://doi.org/10.1186/gb-2005-6-9-r75>
- 7 Charlesworth B., M. T. Morgan, and D. Charlesworth, 1993 The effect of deleterious  
8 mutations on neutral molecular variation. *Genetics* 134: 1289–1303.
- 9 Charlesworth B., 1996 Background selection and patterns of genetic diversity in *Drosophila*  
10 *melanogaster*. *Genet. Res.* 68: 131–149.
- 11 Charlesworth B., 2012 The effects of deleterious mutations on evolution at linked sites.  
12 *Genetics* 190: 5–22. <https://doi.org/10.1534/genetics.111.134288>
- 13 Charlesworth B., 2013 Background Selection 20 Years onThe Wilhelmine E. Key 2012  
14 Invitational Lecture. *J Hered* 104: 161–171. <https://doi.org/10.1093/jhered/ess136>
- 15 Charlesworth B., J. L. Campos, and B. C. Jackson, 2018 Faster-X evolution: Theory and  
16 evidence from *Drosophila*. *Molecular Ecology* 27: 3753–3771.  
17 <https://doi.org/10.1111/mec.14534>
- 18 Choi J. Y., and C. F. Aquadro, 2016 Recent and long term selection across synonymous sites  
19 in *Drosophila ananassae*. *J Mol Evol* 83: 50–60. [https://doi.org/10.1007/s00239-016-](https://doi.org/10.1007/s00239-016-9753-9)  
20 [9753-9](https://doi.org/10.1007/s00239-016-9753-9)
- 21 Comeron J. M., R. Ratnappan, and S. Bailin, 2012 The many landscapes of recombination in  
22 *Drosophila melanogaster*. *PLOS Genetics* 8: e1002905.  
23 <https://doi.org/10.1371/journal.pgen.1002905>
- 24 Comeron J. M., 2014 Background selection as baseline for nucleotide variation across the  
25 *Drosophila* genome. *PLOS Genetics* 10: e1004434.  
26 <https://doi.org/10.1371/journal.pgen.1004434>
- 27 Comeron J. M., 2017 Background selection as null hypothesis in population genomics:  
28 insights and challenges from *Drosophila* studies. *Phil. Trans. R. Soc. B* 372:  
29 20160471. <https://doi.org/10.1098/rstb.2016.0471>
- 30 Crisci J. L., Y.-P. Poh, A. Bean, A. Simkin, and J. D. Jensen, 2012 Recent progress in  
31 polymorphism-based population genetic inference. *J. Hered.* 103: 287–296.  
32 <https://doi.org/10.1093/jhered/esr128>
- 33 Crisci J. L., Y.-P. Poh, S. Mahajan, and J. D. Jensen, 2013 The impact of equilibrium  
34 assumptions on tests of selection. *Front. Genet.* 4: 235.  
35 <https://doi.org/10.3389/fgene.2013.00235>
- 36 Crow J. F., 2008 Mid-century controversies in population genetics. *Annual Review of*  
37 *Genetics* 42: 1–16. <https://doi.org/10.1146/annurev.genet.42.110807.091612>

- 1 Csilléry K., O. François, and M. G. B. Blum, 2012 abc: an R package for approximate  
2 Bayesian computation (ABC). *Methods in Ecology and Evolution* 3: 475–479.  
3 <https://doi.org/10.1111/j.2041-210X.2011.00179.x>
- 4 Cutter A. D., and B. A. Payseur, 2013 Genomic signatures of selection at linked sites:  
5 unifying the disparity among species. *Nature Reviews Genetics* 14: 262–274.  
6 <https://doi.org/10.1038/nrg3425>
- 7 Duchen P., D. Živković, S. Hutter, W. Stephan, and S. Laurent, 2013 Demographic inference  
8 reveals African and European admixture in the North American *Drosophila*  
9 *melanogaster* population. *Genetics* 193: 291–301.  
10 <https://doi.org/10.1534/genetics.112.145912>
- 11 Elyashiv E., S. Sattath, T. T. Hu, A. Strutsovsky, G. McVicker, *et al.*, 2016 A genomic map  
12 of the effects of linked selection in *Drosophila*. *PLOS Genet* 12: e1006130.  
13 <https://doi.org/10.1371/journal.pgen.1006130>
- 14 Ewing G. B., and J. D. Jensen, 2016 The consequences of not accounting for background  
15 selection in demographic inference. *Molecular Ecology* 25: 135–141.  
16 <https://doi.org/10.1111/mec.13390>
- 17 Eyre-Walker A., and P. D. Keightley, 2007 The distribution of fitness effects of new  
18 mutations. *Nature Reviews Genetics* 8: 610–618. <https://doi.org/10.1038/nrg2146>
- 19 Eyre-Walker A., and P. D. Keightley, 2009 Estimating the rate of adaptive molecular  
20 evolution in the presence of slightly deleterious mutations and population size change.  
21 *Mol Biol Evol* 26: 2097–2108. <https://doi.org/10.1093/molbev/msp119>
- 22 Fay J. C., and C. I. Wu, 2000 Hitchhiking under positive Darwinian selection. *Genetics* 155:  
23 1405–1413.
- 24 Fiston-Lavier A.-S., N. D. Singh, M. Lipatov, and D. A. Petrov, 2010 *Drosophila*  
25 *melanogaster* recombination rate calculator. *Gene* 463: 18–20.  
26 <https://doi.org/10.1016/j.gene.2010.04.015>
- 27 Fu Y.-X., 1995 Statistical properties of segregating sites. *Theoretical Population Biology*  
28 172–197.
- 29 Good B. H., A. M. Walczak, R. A. Neher, and M. M. Desai, 2014 Genetic diversity in the  
30 interference selection limit. *PLOS Genetics* 10: e1004222.  
31 <https://doi.org/10.1371/journal.pgen.1004222>
- 32 Haller B. C., and P. W. Messer, 2019 SLiM 3: Forward genetic simulations beyond the  
33 Wright–Fisher model. *Mol Biol Evol* 36: 632–637.  
34 <https://doi.org/10.1093/molbev/msy228>
- 35 Halligan D. L., and P. D. Keightley, 2006 Ubiquitous selective constraints in the *Drosophila*  
36 genome revealed by a genome-wide interspecies comparison. *Genome Res* 16: 875–  
37 884. <https://doi.org/10.1101/gr.5022906>

- 1 Harris R. B., A. Sackman, and J. D. Jensen, 2018 On the unfounded enthusiasm for soft  
2 selective sweeps II: Examining recent evidence from humans, flies, and viruses.  
3 PLOS Genetics 14: e1007859. <https://doi.org/10.1371/journal.pgen.1007859>
- 4 Hoskins R. A., J. W. Carlson, C. Kennedy, D. Acevedo, M. Evans-Holm, *et al.*, 2007  
5 Sequence finishing and mapping of *Drosophila melanogaster* heterochromatin.  
6 Science 316: 1625–1628. <https://doi.org/10.1126/science.1139816>
- 7 Huber C. D., B. Y. Kim, C. D. Marsden, and K. E. Lohmueller, 2017 Determining the factors  
8 driving selective effects of new nonsynonymous mutations. PNAS 114: 4465–4470.  
9 <https://doi.org/10.1073/pnas.1619508114>
- 10 Hudson R. R., and N. L. Kaplan, 1995 Deleterious background selection with recombination.  
11 Genetics 141: 1605–1617.
- 12 Jackson B. C., J. L. Campos, P. R. Haddrill, B. Charlesworth, and K. Zeng, 2017 Variation in  
13 the intensity of selection on codon bias over time causes contrasting patterns of base  
14 composition evolution in *Drosophila*. Genome Biol Evol 9: 102–123.  
15 <https://doi.org/10.1093/gbe/evw291>
- 16 Jacquier H., A. Birgy, H. L. Nagard, Y. Mechulam, E. Schmitt, *et al.*, 2013 Capturing the  
17 mutational landscape of the beta-lactamase TEM-1. PNAS 110: 13067–13072.  
18 <https://doi.org/10.1073/pnas.1215206110>
- 19 Jensen J. D., Y. Kim, V. B. DuMont, C. F. Aquadro, and C. D. Bustamante, 2005  
20 Distinguishing between selective sweeps and demography using DNA polymorphism  
21 data. Genetics 170: 1401–1410. <https://doi.org/10.1534/genetics.104.038224>
- 22 Jensen J. D., B. A. Payseur, W. Stephan, C. F. Aquadro, M. Lynch, *et al.*, 2019 The  
23 importance of the Neutral Theory in 1968 and 50 years on: A response to Kern and  
24 Hahn 2018. Evolution 73: 111–114. <https://doi.org/10.1111/evo.13650>
- 25 Joyce P., and P. Marjoram, 2008 Approximately sufficient statistics and bayesian  
26 computation. Statistical Applications in Genetics and Molecular Biology 7.  
27 <https://doi.org/10.2202/1544-6115.1389>
- 28 Kaiser V. B., and B. Charlesworth, 2009 The effects of deleterious mutations on evolution in  
29 non-recombining genomes. Trends in Genetics 25: 9–12.  
30 <https://doi.org/10.1016/j.tig.2008.10.009>
- 31 Kapopoulou A., S. P. Pfeifer, J. D. Jensen, and S. Laurent, 2018 The demographic history of  
32 African *Drosophila melanogaster*. Genome Biol Evol 10: 2338–2342.  
33 <https://doi.org/10.1093/gbe/evy185>
- 34 Karolchik D., A. S. Hinrichs, T. S. Furey, K. M. Roskin, C. W. Sugnet, *et al.*, 2004 The  
35 UCSC Table Browser data retrieval tool. Nucleic Acids Res. 32: D493-496.  
36 <https://doi.org/10.1093/nar/gkh103>
- 37 Keightley P. D., and A. Eyre-Walker, 2007 Joint inference of the distribution of fitness  
38 effects of deleterious mutations and population demography based on nucleotide  
39 polymorphism frequencies. Genetics 177: 2251–2261.  
40 <https://doi.org/10.1534/genetics.107.080663>

- 1 Keightley P. D., R. W. Ness, D. L. Halligan, and P. R. Haddrill, 2014 Estimation of the  
2 spontaneous mutation rate per nucleotide site in a *Drosophila melanogaster* full-sib  
3 family. *Genetics* 196: 313–320. <https://doi.org/10.1534/genetics.113.158758>
- 4 Kern A. D., and M. W. Hahn, 2018 The neutral theory in light of natural selection. *Mol Biol*  
5 *Evol* 35: 1366–1371. <https://doi.org/10.1093/molbev/msy092>
- 6 Kim B. Y., C. D. Huber, and K. E. Lohmueller, 2017 Inference of the distribution of selection  
7 coefficients for new nonsynonymous mutations using large samples. *Genetics* 206:  
8 345–361. <https://doi.org/10.1534/genetics.116.197145>
- 9 Kimura M., 1968 Evolutionary rate at the molecular level. *Nature* 217: 624–626.  
10 <https://doi.org/10.1038/217624a0>
- 11 Kimura M., 1983 *The neutral theory of molecular evolution*. Cambridge University Press.
- 12 King J. L., and T. H. Jukes, 1969 Non-Darwinian Evolution. *Science* 164: 788–798.  
13 <https://doi.org/10.1126/science.164.3881.788>
- 14 Kolaczowski B., A. D. Kern, A. K. Holloway, and D. J. Begun, 2011 Genomic  
15 differentiation between temperate and tropical Australian populations of *Drosophila*  
16 *melanogaster*. *Genetics* 187: 245–260. <https://doi.org/10.1534/genetics.110.123059>
- 17 Kousathanas A., and P. D. Keightley, 2013 A comparison of models to infer the distribution  
18 of fitness effects of new mutations. *Genetics* 193: 1197–1208.  
19 <https://doi.org/10.1534/genetics.112.148023>
- 20 Lack J. B., C. M. Cardeno, M. W. Crepeau, W. Taylor, R. B. Corbett-Detig, *et al.*, 2015 The  
21 *Drosophila* Genome Nexus: A Population Genomic Resource of 623 *Drosophila*  
22 *melanogaster* Genomes, Including 197 from a Single Ancestral Range Population.  
23 *Genetics* 199: 1229–1241. <https://doi.org/10.1534/genetics.115.174664>
- 24 Lange J. D., and J. E. Pool, 2018 Impacts of recurrent hitchhiking on divergence and  
25 demographic inference in *Drosophila*. *Genome Biol Evol* 10: 1882–1891.  
26 <https://doi.org/10.1093/gbe/evy142>
- 27 Langley C. H., K. Stevens, C. Cardeno, Y. C. G. Lee, D. R. Schrider, *et al.*, 2012 Genomic  
28 variation in natural populations of *Drosophila melanogaster*. *Genetics* 192: 533–598.  
29 <https://doi.org/10.1534/genetics.112.142018>
- 30 Laurent S. J. Y., A. Werzner, L. Excoffier, and W. Stephan, 2011 Approximate bayesian  
31 analysis of *Drosophila melanogaster* polymorphism data reveals a recent colonization  
32 of Southeast Asia. *Molecular Biology and Evolution* 28: 2041–2051.  
33 <https://doi.org/10.1093/molbev/msr031>
- 34 Lawrie D. S., P. W. Messer, R. Hershberg, and D. A. Petrov, 2013 Strong purifying selection  
35 at synonymous sites in *D. melanogaster*. *PLOS Genetics* 9: e1003527.  
36 <https://doi.org/10.1371/journal.pgen.1003527>
- 37 Li H., and W. Stephan, 2006 Inferring the demographic history and rate of adaptive  
38 substitution in *Drosophila*. *PLOS Genetics* 2: e166.  
39 <https://doi.org/10.1371/journal.pgen.0020166>

- 1 Lynch M., 2007 *The Origins of Genome Architecture*. Sinauer Associates, Sunderland,  
2 Massachusetts.
- 3 Matsumoto T., A. John, P. Baeza-Centurion, B. Li, and H. Akashi, 2016 Codon usage  
4 selection can bias estimation of the fraction of adaptive amino acid fixations. *Mol*  
5 *Biol Evol* 33: 1580–1589. <https://doi.org/10.1093/molbev/msw027>
- 6 Messer P. W., and D. A. Petrov, 2013 Frequent adaptation and the McDonald–Kreitman test.  
7 *PNAS* 110: 8615–8620. <https://doi.org/10.1073/pnas.1220835110>
- 8 Nicolaisen L. E., and M. M. Desai, 2013 Distortions in genealogies due to purifying selection  
9 and recombination. *Genetics* 195: 221–230.  
10 <https://doi.org/10.1534/genetics.113.152983>
- 11 Nordborg M., B. Charlesworth, and D. Charlesworth, 1996 The effect of recombination on  
12 background selection. *Genet. Res.* 67: 159–174.
- 13 O’Fallon B. D., J. Seger, and F. R. Adler, 2010 A continuous-state coalescent and the impact  
14 of weak selection on the structure of gene genealogies. *Mol Biol Evol* 27: 1162–1172.  
15 <https://doi.org/10.1093/molbev/msq006>
- 16 Ohta T., 1973 Slightly deleterious mutant substitutions in evolution. *Nature* 246: 96–98.  
17 <https://doi.org/10.1038/246096a0>
- 18 Pordes R., D. Petravick, B. Kramer, D. Olson, M. Livny, *et al.*, 2007 The open science grid.  
19 *J. Phys.: Conf. Ser.* 78: 012057. <https://doi.org/10.1088/1742-6596/78/1/012057>
- 20 Sanjuán R., 2010 Mutational fitness effects in RNA and single-stranded DNA viruses:  
21 common patterns revealed by site-directed mutagenesis studies. *Philos Trans R Soc*  
22 *Lond B Biol Sci* 365: 1975–1982. <https://doi.org/10.1098/rstb.2010.0063>
- 23 Schneider A., B. Charlesworth, A. Eyre-Walker, and P. D. Keightley, 2011 A method for  
24 inferring the rate of occurrence and fitness effects of advantageous mutations.  
25 *Genetics* 189: 1427–1437. <https://doi.org/10.1534/genetics.111.131730>
- 26 Schrider D. R., D. Houle, M. Lynch, and M. W. Hahn, 2013 Rates and genomic  
27 consequences of spontaneous mutational events in *Drosophila melanogaster*.  
28 *Genetics* 194: 937–954. <https://doi.org/10.1534/genetics.113.151670>
- 29 Schrider D. R., A. G. Shanku, and A. D. Kern, 2016 Effects of linked selective sweeps on  
30 demographic inference and model selection. *Genetics* 204: 1207–1223.  
31 <https://doi.org/10.1534/genetics.116.190223>
- 32 Sfiligoi I., D. C. Bradley, B. Holzman, P. Mhashilkar, S. Padhi, *et al.*, 2009 The pilot way to  
33 grid resources using glideinWMS, pp. 428–432 in *2009 WRI World Congress on*  
34 *Computer Science and Information Engineering*.
- 35 Siepel A., G. Bejerano, J. S. Pedersen, A. S. Hinrichs, M. Hou, *et al.*, 2005 Evolutionarily  
36 conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 15:  
37 1034–1050. <https://doi.org/10.1101/gr.3715005>



- 1 Thornton K., 2003 Libsequence: a C++ class library for evolutionary genetic analysis.  
2 Bioinformatics 19: 2325–2327. <https://doi.org/10.1093/bioinformatics/btg316>
- 3 Thornton K. R., and J. D. Jensen, 2007 Controlling the false-positive rate in multilocus  
4 genome scans for selection. Genetics 175: 737–750.  
5 <https://doi.org/10.1534/genetics.106.064642>
- 6 Torres R., M. G. Stetter, R. Hernandez, and J. Ross-Ibarra, 2019 The temporal dynamics of  
7 background selection in non-equilibrium populations. bioRxiv.  
8 <https://doi.org/10.1101/618389>
- 9 Walsh B., and M. Lynch, 2018 *Evolution and Selection of Quantitative Traits*. Oxford  
10 University Press.
- 11 Wegmann D., C. Leuenberger, and L. Excoffier, 2009 Efficient approximate Bayesian  
12 computation coupled with Markov Chain Monte Carlo without likelihood. Genetics  
13 182: 1207–1218. <https://doi.org/10.1534/genetics.109.102509>
- 14 Zeng K., and B. Charlesworth, 2010a Studying patterns of recent evolution at synonymous  
15 sites and intronic sites in *Drosophila melanogaster*. J Mol Evol 70: 116–128.  
16 <https://doi.org/10.1007/s00239-009-9314-6>
- 17 Zeng K., and B. Charlesworth, 2010b The effects of demography and linkage on the  
18 estimation of selection and mutation parameters. Genetics 186: 1411–1424.  
19 <https://doi.org/10.1534/genetics.110.122150>
- 20 Zeng K., 2013 A coalescent model of background selection with recombination, demography  
21 and variation in selection coefficients. Heredity 110: 363–371.  
22 <https://doi.org/10.1038/hdy.2012.102>
- 23  
24  
25

**SUPPLEMENTARY FIGURES AND TABLES**

**Supp Table 1:** Prediction of diversity in linked neutral regions for two different DFE realizations, as predicted by equation 3a and 3b. Expected diversity is 0.02.

DFE	Length of coding:	Distance from selected site -->		
		100 bp	10 kb	100 kb
$f_0 = 0; f_1 = 0, f_2 = 50, f_3 = 50$	1 kb	0.0174857	0.019805	0.0199807
	5 kb	0.0150363	0.0192131	0.0199054
	10 kb	0.0141392	0.0187112	0.0198152
$f_0 = 0; f_1 = 0; f_2 = 80; f_3 = 20$	1 kb	0.0160713	0.0197472	0.0199892
	5 kb	0.0124679	0.0190189	0.0199474
	10 kb	0.0113456	0.0184565	0.0198979

**Supp Table 2:** Reduction in neutral and linked neutral diversity as a function of the DFE - as illustrated by considering DFE realizations in which one class is largely over-represented, for different exon lengths (0.5kb, 1kb, 5kb, and 10kb). The expected diversity under neutrality is 0.02.

	$f_0 \geq 80\%$	$f_1 \geq 80\%$	$f_2 \geq 80\%$	$f_3 \geq 80\%$
<b>Neutral Diversity (0.5 kb)</b>	0.01995	0.01982	0.01958	0.01975
<b>Linked neutral diversity (0.5kb)</b>	0.01976	0.01809	0.01795	0.01891
<b>Slope of recovery (0.5kb)</b>	0.00007	0.00054	0.00063	0.00013
<b>Neutral Diversity (1 kb)</b>	0.01992	0.01972	0.01926	0.01951
<b>Linked neutral diversity (1kb)</b>	0.01968	0.01756	0.01674	0.01906
<b>Slope of recovery (1kb)</b>	0.00010	0.00072	0.00100	0.00024
<b>Neutral Diversity (5 kb)</b>	0.01972	0.01934	0.01760	0.01814
<b>Linked neutral diversity (5 kb)</b>	0.01915	0.01634	0.01318	0.01714
<b>Slope of recovery (5kb)</b>	0.00024	0.00109	0.00200	0.00061
<b>Neutral Diversity (10kb)</b>	0.01960	0.01909	0.01673	0.01709
<b>Linked neutral diversity (10kb)</b>	0.01903	0.01615	0.01208	0.01579
<b>Slope of recovery (10kb)</b>	0.00027	0.00112	0.00220	0.00083

**Supp Table 3:** Statistics ranked by their importance in predicting the DFE classes under equilibrium using the correlation coefficients between the statistics and parameters.

Statistics ranked for $f_0$	$r^2$	Statistics ranked for $f_1$	$r^2$	Statistics ranked for $f_2$	$r^2$	Statistics ranked for $f_3$	$r^2$
func_div_sd	0.962	func_hprime_m	0.632	neu_thetaw_m	0.477	func_numSing_m	0.850
func_thetah_sd	0.962	func_rsqs_sd	0.511	neu_thetapi_m	0.463	func_numSing_sd	0.748
func_thetah_m	0.960	link_tajimasd_m	0.456	pi_intercept	0.425	func_thetaw_m	0.460
func_div_m	0.958	link_Dprime_m	0.441	neu_thetah_m	0.347	func_Dprime_sd	0.432
func_thetapi_sd	0.896	func_Dprime_sd	0.384	pi_slope	0.340	func_thetaw_sd	0.397
func_rsqs_m	0.873	pi_max	0.357	link_thetaw_m	0.320	func_hapdiv_m	0.334
func_Dprime_m	0.866	func_D_sd	0.297	func_tajimasd_m	0.309	func_rsqs_sd	0.317
func_thetapi_m	0.855	func_tajimasd_sd	0.280	func_rsqs_m	0.278	pi_max	0.297
func_D_m	0.847	func_hprime_sd	0.271	link_thetapi_m	0.252	func_hapdiv_sd	0.295
func_tajimasd_m	0.828	link_thetah_m	0.236	neu_div_m	0.240	link_tajimasd_m	0.267
func_hapdiv_sd	0.731	pi_slope	0.226	func_Dprime_m	0.235	neu_rsqs_m	0.264
func_thetaw_sd	0.729	link_thetapi_m	0.219	func_hapdiv_m	0.230	func_thetapi_m	0.258
func_hapdiv_m	0.683	func_numSing_m	0.200	link_div_m	0.229	link_hapdiv_sd	0.252
func_thetaw_m	0.663	func_numSing_sd	0.185	link_thetah_m	0.222	link_rsqs_sd	0.234
func_hprime_sd	0.578	link_div_m	0.183	func_hapdiv_sd	0.220	link_D_sd	0.230
func_hprime_m	0.553	pi_intercept	0.165	func_thetapi_m	0.215	link_hapdiv_m	0.218
func_D_sd	0.426	link_hapdiv_sd	0.159	func_thetapi_sd	0.210	link_Dprime_sd	0.217
pi_intercept	0.407	neu_rsqs_m	0.143	func_D_m	0.195	pi_slope	0.217
neu_thetaw_m	0.404	func_D_m	0.141	func_thetah_m	0.181	func_thetapi_sd	0.208
func_tajimasd_sd	0.381	link_rsqs_sd	0.126	func_D_sd	0.175	link_D_m	0.205
neu_thetapi_m	0.357	link_hapdiv_m	0.125	func_tajimasd_sd	0.171	neu_D_m	0.200
pi_slope	0.345	neu_D_m	0.120	func_thetah_sd	0.168	link_Dprime_m	0.191
func_numSing_sd	0.339	link_Dprime_sd	0.119	func_div_m	0.165	func_div_m	0.187
link_thetapi_m	0.326	link_D_sd	0.117	func_div_sd	0.159	link_hprime_sd	0.183
link_thetaw_m	0.320	func_Dprime_m	0.114	func_thetaw_sd	0.158	link_tajimasd_sd	0.182
link_thetah_m	0.315	link_hprime_sd	0.108	func_thetaw_m	0.148	pi_intercept	0.173
func_numSing_m	0.279	func_tajimasd_m	0.104	func_hprime_sd	0.136	func_div_sd	0.173
link_div_m	0.259	link_thetaw_m	0.103	func_rsqs_sd	0.128	link_rsqs_m	0.172
func_rsqs_sd	0.255	link_div_sd	0.099	neu_tajimasd_m	0.091	link_div_sd	0.164
neu_thetah_m	0.251	link_tajimasd_sd	0.098	link_tajimasd_m	0.081	func_thetah_m	0.163
link_Dprime_m	0.220	neu_D_sd	0.096	link_Dprime_m	0.059	link_thetapi_m	0.156
link_tajimasd_m	0.196	link_D_m	0.090	func_Dprime_sd	0.044	link_div_m	0.155
neu_div_m	0.171	neu_rsqs_sd	0.086	neu_Dprime_m	0.043	link_thetah_m	0.154
neu_numSing_m	0.076	link_numSing_sd	0.081	link_numSing_m	0.024	neu_rsqs_sd	0.153

neu_Dprime_m	0.053	neu_thetaw_m	0.076	neu_numSing_m	0.023	func_thetah_sd	0.153
neu_tajimasd_m	0.049	link_rsq_m	0.076	func_numSing_s d	0.021	link_thetah_sd	0.144
func_Dprime_sd	0.028	link_thetaw_sd	0.072	neu_hapdiv_m	0.011	neu_D_sd	0.142
link_numSing_m	0.027	link_thetah_sd	0.072	link_Dprime_sd	0.010	link_numSing_sd	0.140
neu_rsq_m	0.009	func_thetaw_m	0.062	neu_thetaw_sd	0.010	link_thetaw_sd	0.134
link_thetapi_sd	0.008	link_thetapi_sd	0.061	neu_thetapi_sd	0.009	link_thetapi_sd	0.127
neu_tajimasd_sd	0.008	func_rsq_m	0.058	link_D_m	0.008	func_tajimasd_sd	0.104
link_thetah_sd	0.007	func_hapdiv_m	0.053	link_D_sd	0.008	link_thetaw_m	0.101
neu_Dprime_sd	0.007	neu_hprime_sd	0.047	neu_numSing_sd	0.007	neu_numSing_m	0.098
neu_D_m	0.007	neu_div_sd	0.043	link_rsq_sd	0.007	func_D_sd	0.094
link_thetaw_sd	0.006	neu_Dprime_sd	0.039	link_tajimasd_sd	0.006	neu_div_sd	0.091
link_rsq_m	0.005	neu_thetapi_m	0.037	link_hapdiv_sd	0.005	neu_hprime_sd	0.078
neu_hprime_sd	0.005	neu_numSing_m	0.035	link_div_sd	0.005	neu_tajimasd_sd	0.075
neu_rsq_sd	0.005	neu_tajimasd_sd	0.034	link_rsq_m	0.005	neu_Dprime_sd	0.073
link_hapdiv_m	0.005	neu_thetah_sd	0.033	func_hprime_m	0.005	neu_thetah_sd	0.065
link_D_m	0.004	func_thetah_sd	0.033	link_hprime_sd	0.005	neu_thetaw_sd	0.053
link_numSing_sd	0.003	func_thetaw_sd	0.030	neu_thetah_sd	0.004	neu_thetapi_sd	0.053
neu_D_sd	0.003	link_hprime_m	0.029	func_numSing_ m	0.003	neu_thetaw_m	0.050
link_hprime_m	0.003	func_div_sd	0.028	link_hapdiv_m	0.002	neu_numSing_sd	0.037
neu_div_sd	0.003	neu_thetah_m	0.025	neu_div_sd	0.002	neu_hapdiv_m	0.033
link_D_sd	0.002	func_hapdiv_sd	0.024	neu_rsq_m	0.001	func_rsq_m	0.028
link_rsq_sd	0.002	neu_thetapi_sd	0.024	link_numSing_sd	0.001	neu_tajimasd_m	0.019
neu_hapdiv_m	0.002	func_thetah_m	0.023	pi_numbp50	0.001	func_hprime_sd	0.016
link_tajimasd_sd	0.001	neu_thetaw_sd	0.020	pi_numbp75	0.001	func_hprime_m	0.014
pi_max	0.001	func_div_m	0.020	pi_numbp90	0.001	neu_Dprime_m	0.014
link_hapdiv_sd	0.001	neu_div_m	0.018	link_thetah_sd	0.001	func_Dprime_m	0.012
link_hprime_sd	0.001	neu_hapdiv_m	0.015	neu_rsq_sd	0.001	neu_thetapi_m	0.012
neu_thetapi_sd	0.001	neu_numSing_sd	0.015	link_thetaw_sd	0.000	func_D_m	0.011
link_Dprime_sd	0.001	neu_Dprime_m	0.008	link_thetapi_sd	0.000	link_hprime_m	0.009
neu_hapdiv_sd	0.000	link_numSing_m	0.003	link_hprime_m	0.000	neu_thetah_m	0.005
link_div_sd	0.000	neu_tajimasd_m	0.003	neu_hprime_m	0.000	link_numSing_m	0.004
neu_hprime_m	0.000	func_thetapi_m	0.002	neu_D_m	0.000	neu_div_m	0.004
neu_thetah_sd	0.000	neu_hapdiv_sd	0.002	pi_max	0.000	neu_hapdiv_sd	0.004
neu_numSing_sd	0.000	func_thetapi_sd	0.001	neu_Dprime_sd	0.000	pi_numbp50	0.003
neu_thetaw_sd	0.000	pi_numbp50	0.001	neu_hprime_sd	0.000	pi_numbp75	0.003
pi_numbp50	0.000	pi_numbp75	0.001	neu_D_sd	0.000	pi_numbp90	0.003
pi_numbp75	0.000	pi_numbp90	0.001	neu_hapdiv_sd	0.000	neu_hprime_m	0.001
pi_numbp90	0.000	neu_hprime_m	0.000	neu_tajimasd_sd	0.000	func_tajimasd_m	0.001

**Supp Table 4:** Statistics ranked by their importance in predicting the DFE classes under equilibrium using a modified algorithm of Joyce and Marjoram (2008) and by averaging the ranking across 10 replicates for each parameter separately.

Statistics ranked for $f_0$	Avg rank	Statistics ranked for $f_1$	Avg rank	Statistics ranked for $f_2$	Avg rank	Statistics ranked for $f_3$	Avg rank
func div m	1.2	func thetah m	4.8	func thetaw m	2.8	func thetaw m	4.3
func thetah m	4.7	func thetaw m	6.1	func thetah m	4.7	func thetah m	5.2
func thetaw m	8.5	link thetaw m	7.9	link thetaw m	6.2	neu thetapi m	8.2
neu thetapi m	13.8	neu thetapi m	9.6	neu thetapi m	8	link thetaw m	10.6
func hprime m	14.7	link thetapi m	11	func_numSing_m	9.1	func_numSing_m	12.8
func tajimasd m	17.6	neu thetaw m	14.2	link thetapi m	11.4	neu thetaw m	12.8
link hapdiv m	18.7	func_numSing_m	17.4	link thetah m	12.7	link thetapi m	14.3
link thetaw m	18.7	link hprime m	18.8	neu thetah m	18.9	func hprime m	15.8
func rsq m	18.9	pi max	19.9	func Dprime sd	21.5	link hprime m	18.3
link thetapi m	19.8	pi numbp90	20.6	func thetapi m	21.5	pi intercept	18.5
link tajimasd m	20	func thetapi m	20.7	link div sd	22.7	func Dprime sd	19.5
func_numSing_m	20.3	func div m	22	func div sd	23.2	pi max	19.8
neu rsq m	20.7	pi numbp50	23	pi intercept	23.8	link thetah m	20
link rsq m	21	link thetah m	23.2	pi numbp90	24.1	pi numbp90	20.1
func hapdiv m	23	func hprime m	23.9	pi slope	24.4	neu Dprime sd	21
neu numSing m	23	pi intercept	24.3	neu Dprime sd	24.6	pi numbp75	22.1
func D m	25.3	func Dprime sd	24.7	link Dprime sd	24.8	func thetapi m	24.2
link thetah m	25.3	func div sd	24.8	neu thetaw m	24.8	neu hprime m	24.5
func thetapi m	26.2	pi slope	26.5	link hprime m	25.6	link hapdiv sd	24.7
link hprime m	26.7	func tajimasd m	26.9	func_numSing_s d	26.9	pi numbp50	24.9
func Dprime sd	26.8	neu hprime m	27.5	link D sd	27.1	neu D sd	25
link D m	27.3	func D sd	27.7	pi numbp75	27.4	neu thetah m	26.3
neu hprime m	28	link div sd	27.8	func hprime m	28.3	func D sd	26.4
pi numbp75	28.4	neu thetah m	29.7	pi max	28.9	func tajimasd m	26.5
pi intercept	28.8	link D sd	31.8	link numSing sd	29	pi slope	26.5
func div sd	29.1	neu rsq m	32	pi numbp50	29	link div sd	28.6
link_numSing_m	29.6	link tajimasd m	32.9	neu D sd	30.1	neu div sd	29.2
neu D m	29.7	neu div sd	33.6	func rsq sd	30.4	func div sd	30.1
neu thetah m	30	pi numbp75	34.1	neu rsq sd	31.7	func rsq sd	31.1
pi max	30.7	link rsq sd	34.5	neu div sd	32.2	neu rsq sd	32.7
neu rsq sd	31.5	neu D sd	34.6	func D sd	32.7	neu tajimasd sd	33.6
neu thetaw m	32.5	func D m	34.9	func tajimasd m	33.6	link rsq sd	33.7
func Dprime m	36	link Dprime sd	35.6	neu hprime m	34.6	neu hapdiv sd	33.7
func D sd	36.5	func rsq m	35.9	func hapdiv sd	35.2	link Dprime sd	34.3

pi numbp90	38.2	link hapdiv m	36	func thetah sd	36.1	link D sd	37.2
link Dprime sd	38.4	func_numSing_s d	36.4	link hapdiv sd	36.4	func_numSing_s d	37.6
link hapdiv sd	38.6	neu Dprime sd	37.5	link thetapi sd	36.4	neu_numSing_sd	38.8
func hapdiv sd	38.9	link rsq m	39.4	neu hapdiv sd	36.5	link tajimasd m	40.8
pi numbp50	39.1	neu hapdiv sd	40.7	func hprime sd	38	neu thetah sd	41.2
neu D sd	39.2	neu_numSing_sd	41.3	func tajimasd sd	38.1	link thetaw sd	41.3
link rsq sd	40.3	neu_numSing_m	42.1	link hprime sd	40.7	func hapdiv sd	41.7
link div sd	40.6	link D m	42.7	func div m	41.1	link_numSing_sd	42
func rsq sd	41.9	func Dprime m	43	neu tajimasd sd	41.1	func thetapi sd	42.8
pi slope	43.6	link thetaw sd	44.3	link rsq sd	41.7	link hprime sd	42.8
link_numSing_sd	43.8	link hapdiv sd	44.4	neu_numSing_sd	42.1	func thetaw sd	44.8
neu Dprime sd	43.8	neu hprime sd	44.7	link tajimasd sd	42.7	func rsq m	44.9
neu hapdiv sd	44.2	func tajimasd sd	44.9	neu thetaw sd	43.2	neu hprime sd	45
link tajimasd sd	44.5	link hprime sd	45.1	neu thetapi sd	43.6	func Dprime m	45.1
func_numSing_s d	44.7	func hapdiv sd	45.4	func D m	44.6	func hprime sd	45.3
link Dprime m	44.7	func rsq sd	45.5	neu hprime sd	44.8	func tajimasd sd	46.5
neu_numSing_sd	45	func thetaw sd	46.7	link div m	46.5	link tajimasd sd	46.8
neu hapdiv m	45.9	link tajimasd sd	46.7	link Dprime m	47.7	func div m	47.3
neu tajimasd m	46.9	neu rsq sd	46.7	link hapdiv m	48.5	link thetah sd	48.2
func tajimasd sd	47.2	neu thetah sd	47.4	func Dprime m	49.7	func thetah sd	48.4
link D sd	47.9	link Dprime m	47.5	neu div m	50	func D m	48.8
neu div sd	48.4	func thetah sd	47.8	link rsq m	50.1	link D m	48.9
func thetaw sd	49.2	neu tajimasd sd	47.9	neu Dprime m	50.1	neu thetaw sd	49.6
link thetaw sd	49.6	link_numSing_sd	48.2	func thetaw sd	50.5	link div m	50.1
neu tajimasd sd	49.8	func hprime sd	49.2	func thetapi sd	51.1	neu D m	50.4
func thetapi sd	51.3	neu thetaw sd	49.4	link thetaw sd	51.3	neu div m	51.1
func hprime sd	52.7	neu D m	49.5	link thetah sd	51.4	link rsq m	51.3
neu Dprime m	52.8	func hapdiv m	49.9	func hapdiv m	51.8	neu rsq m	54
neu thetaw sd	55	link thetapi sd	50.1	neu D m	52.9	link thetapi sd	54.2
neu thetah sd	56.3	link div m	51.2	neu thetah sd	53	neu thetapi sd	54.8
neu hprime sd	56.7	func thetapi sd	51.7	func rsq m	53.3	neu Dprime m	56.5
neu div m	57.1	link thetah sd	52	link tajimasd m	54.1	link Dprime m	57.4
link hprime sd	57.5	neu thetapi sd	53.8	link D m	56.5	neu tajimasd m	58.4
link thetah sd	58.6	neu div m	56.4	neu rsq m	58.3	func hapdiv m	59.4
neu thetapi sd	59.5	neu tajimasd m	57.5	neu_numSing_m	63.4	link_numSing_m	61.2
link div m	60	link_numSing_m	57.9	neu tajimasd m	65.3	neu_numSing_m	62.5
func thetah sd	60.6	neu Dprime m	59.9	link_numSing_m	65.6	link hapdiv m	62.6
link thetapi sd	62.5	neu hapdiv m	64.3	neu hapdiv m	67.9	neu hapdiv m	69





func hapdiv m	0.21	func thetah sd	0.05	func rsq m	0.04	func tajimasd sd	0.06
func_numSing_sd	0.14	func D sd	0.03	func numSing sd	0.03	func thetah m	0.06
func rsq m	0.14	func hapdiv m	0.02	func numSing m	0.02	func rsq sd	0.02
func hprime sd	0.12	func thetapi sd	0.01	func D m	0.02	func hprime sd	0.01
func_numSing_m	0.10	func D m	0.01	func hprime sd	0.01	func D sd	0.01
func D m	0.08	func hapdiv sd	0.01	func rsq sd	0.01	func rsq m	0.01
func rsq sd	0.06	func thetapi m	0.01	func D sd	0.01	func Dprime m	0.00
func D sd	0.03	func rsq m	0.01				
func_tajimasd_sd	0.02	func thetaw m	0.01				
		func thetaw sd	0.01				

**Supp Table 6:** Ranking of statistics when distinguishing between demography and purifying selection. Statistics significantly correlated with parameters of demography when statistics from all regions are used, and when only functional statistics are used for ranking. Significance was evaluated with  $p < 0.05$  with Bonferonni correction.

Ranking using all statistics				Ranking using statistics calculated from functional regions.			
Statistics ranked for Nanc	$r^2$	Statistics ranked for Neur	$r^2$	Statistics ranked for Nanc	$r^2$	Statistics ranked for Neur	$r^2$
neu thetah m	0.99	neu hapdiv m	0.92	func thetapi m	0.230	func rsq sd	0.659
link thetah m	0.99	link hapdiv m	0.91	func thetah m	0.229	func rsq m	0.618
neu thetapi m	0.93	link numSing m	0.84	func thetapi sd	0.196	func numSing m	0.577
link thetapi m	0.93	neu numSing m	0.84	func thetaw sd	0.192	func D sd	0.557
link thetapi sd	0.93	neu rsq m	0.83	func thetah sd	0.179	func tajimasd sd	0.487
link thetah sd	0.93	link rsq m	0.82	func thetaw m	0.159	func numSing sd	0.482
neu thetah sd	0.91	neu rsq sd	0.79	func Dprime m	0.136	func Dprime sd	0.435
neu thetapi sd	0.91	link rsq sd	0.79	func tajimasd m	0.135	func D m	0.374
neu thetaw sd	0.90	neu hprime sd	0.78	func hapdiv sd	0.121	func hprime sd	0.368
link thetaw sd	0.90	link hprime sd	0.78	func hapdiv m	0.078	func tajimasd m	0.275
neu thetaw m	0.72	link hapdiv sd	0.76	func hprime m	0.069	func Dprime m	0.216
link thetaw m	0.71	neu hapdiv sd	0.74	func numSing sd	0.042	func hapdiv m	0.165
link div sd	0.49	link D sd	0.68	func rsq m	0.041	func thetaw m	0.157
neu div sd	0.45	func rsq sd	0.66	func Dprime sd	0.032	func hapdiv sd	0.129
neu Dprime m	0.45	link numSing sd	0.64	func numSing m	0.025	func thetaw sd	0.070
link Dprime m	0.44	neu numSing sd	0.64	func rsq sd	0.012	func hprime m	0.061
link tajimasd m	0.43	link tajimasd sd	0.64	func tajimasd sd	0.009	func div m	0.022
neu tajimasd m	0.43	neu D sd	0.64	func D m	0.008	func div sd	0.017
neu Dprime sd	0.41	func rsq m	0.62			func thetapi m	0.008
link Dprime sd	0.38	neu tajimasd sd	0.62				
link hprime m	0.35	neu div m	0.59				
neu hprime m	0.34	link D m	0.58				
link div m	0.31	func numSing m	0.58				
neu div m	0.30	link div m	0.57				
neu numSing sd	0.25	func D sd	0.56				
link numSing sd	0.25	neu D m	0.56				
func thetapi m	0.23	func tajimasd sd	0.49				
func thetah m	0.23	func numSing sd	0.48				
func thetapi sd	0.20	func Dprime sd	0.44				

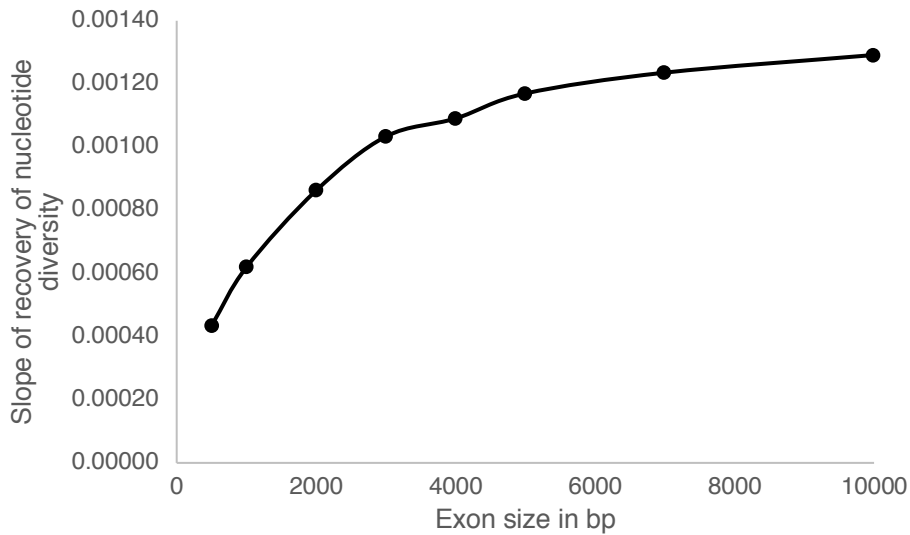
neu tajimasd sd	0.19	link tajimasd m	0.43				
func thetaw sd	0.19	neu tajimasd m	0.43				
func thetah sd	0.18	link Dprime sd	0.39				
link tajimasd sd	0.16	neu Dprime sd	0.38				
func thetaw m	0.16	func D m	0.37				
func Dprime m	0.14	func hprime sd	0.37				
func tajimasd m	0.14	link Dprime m	0.35				
func hapdiv sd	0.12	neu Dprime m	0.35				
neu numSing m	0.11	neu hprime m	0.34				
link numSing m	0.11	link hprime m	0.34				
link rsq sd	0.08	func tajimasd m	0.27				
neu rsq sd	0.08	link thetaw m	0.22				
func hapdiv m	0.08	func Dprime m	0.22				
link rsq m	0.08	neu thetaw m	0.21				
func hprime m	0.07	func hapdiv m	0.16				
neu rsq m	0.07	func thetaw m	0.16				
func numSing sd	0.04	neu div sd	0.15				
func rsq m	0.04	link div sd	0.14				
func Dprime sd	0.03	func hapdiv sd	0.13				
func numSing m	0.03	func thetaw sd	0.07				
func rsq sd	0.01	func hprime m	0.06				
link D m	0.01	link thetaw sd	0.04				
link hapdiv sd	0.01	neu thetaw sd	0.03				
link hapdiv m	0.01	link thetapi m	0.03				
func tajimasd sd	0.01	neu thetapi m	0.03				
func D m	0.01	func div m	0.02				
neu D m	0.01	func div sd	0.02				
		neu thetapi sd	0.01				
		func thetapi m	0.01				
		neu thetah sd	0.01				
		link thetapi sd	0.01				

**Supp Table 7:** The rate of fixed differences (*i.e.*, polymorphism-adjusted divergence) for different site types in *D. melanogaster*, where different numbers of individuals from the Zambia population were used to identify the set of polymorphic sites.

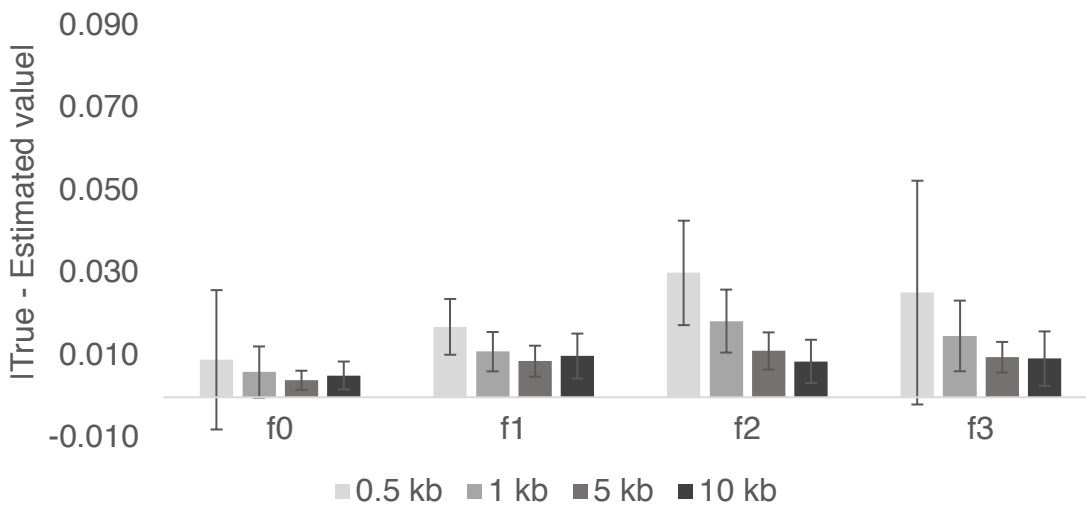
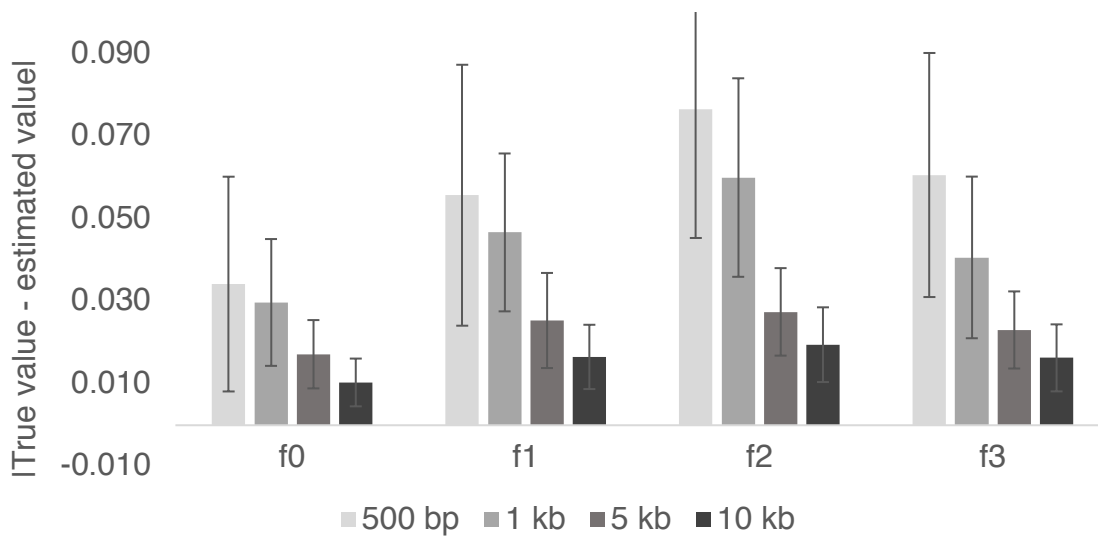
	Sample size:					
	1	2	5	15	30	76
<b>exon</b>	0.0238	0.0198	0.0170	0.0160	0.0159	0.0153
<b>coding</b>	0.0228	0.0182	0.0157	0.0146	0.0141	0.0135
<b>4-fold degenerate</b>	0.0497	0.0423	0.0349	0.0316	0.0311	0.0300
<b>0-fold degenerate</b>	0.0182	0.0123	0.0108	0.0102	0.0098	0.0094

**Supp Table 8:** The increase in divergence values obtained when calculating pairwise divergence (corresponding to sample size of 1) relative to when a much larger number of individuals are used to calculate the rate of fixed differences with exclusion of polymorphic sites.

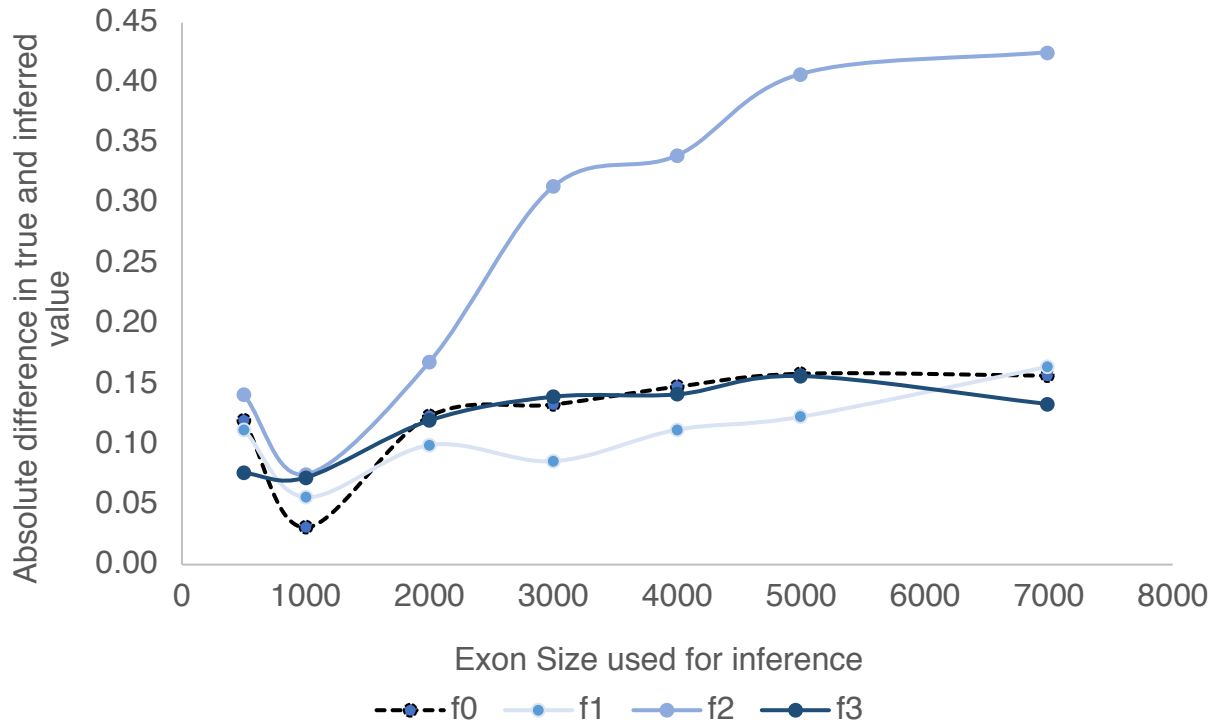
	Sample size:		
	1	76	100
<i>D. melanogaster</i> exon	1.551	1.000	
<i>D. melanogaster</i> 4-fold degenerate	1.658	1.000	
Simulated exon	1.736		1.000
Simulated neutral	1.634		1.000



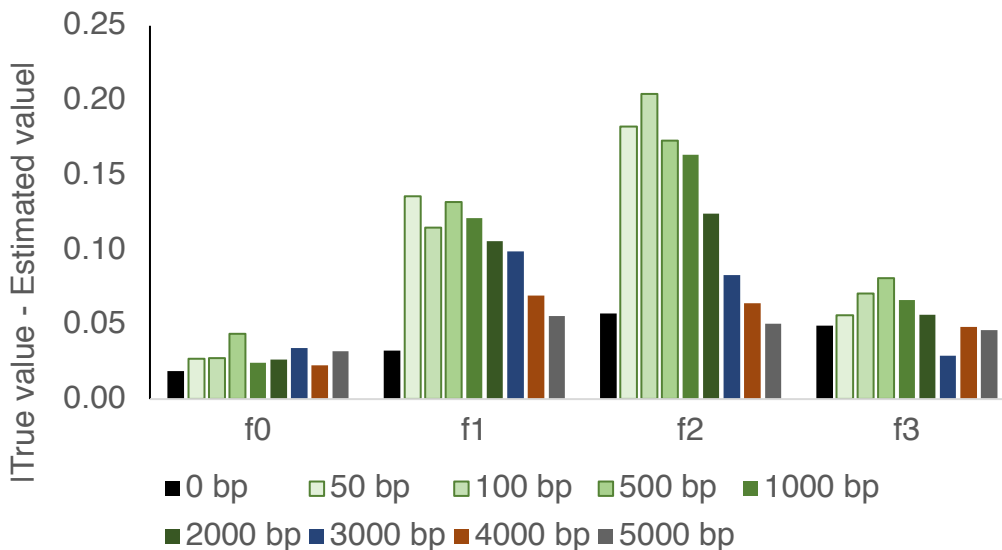
**Supp Figure 1:** Increase in the slope of recovery of diversity near functional regions of varying sizes. Larger values of slope represent a steeper recovery, concordant with larger reduction in diversity observed in the non-coding region.



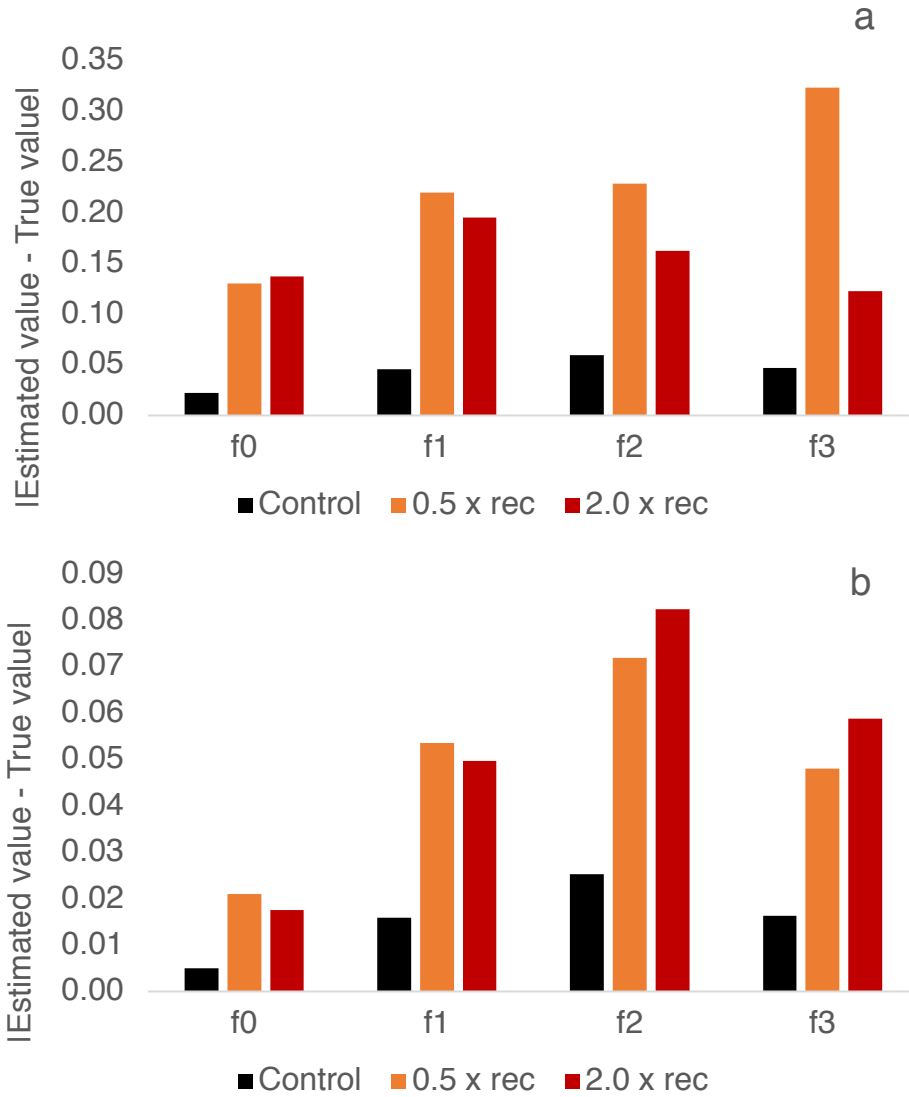
**Supp Figure 2:** Absolute difference between true and inferred value of parameters characterizing the DFE for 0.5 kb, 1 kb, 5 kb, and 10 kb functional regions. The upper panel displays the error in inference when using all statistics, while the lower uses only functional regions.



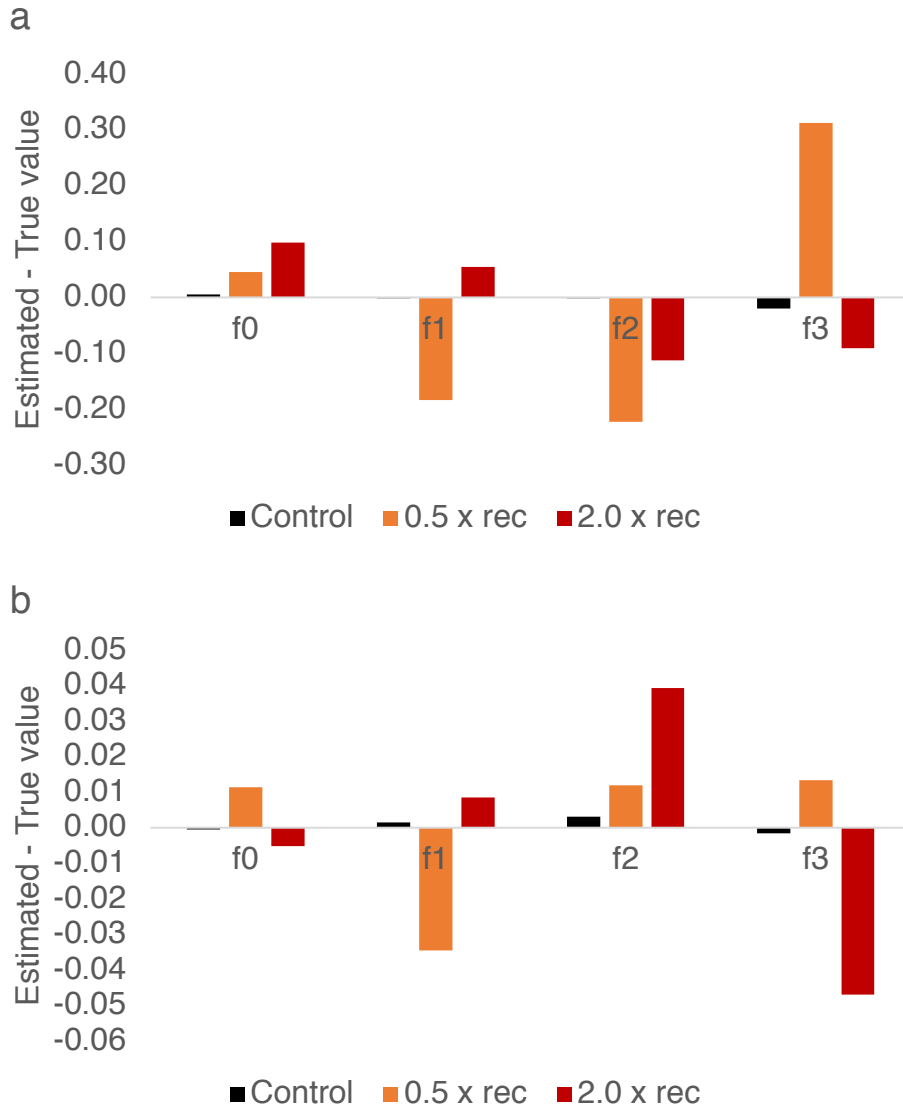
**Supp Figure 3:** Decrease in accuracy of inference for different DFE classes as the exon size assumed for inference is mis-specified. In this figure, the assumed exon size was 1kb, and the X-axis gives the true exon size.



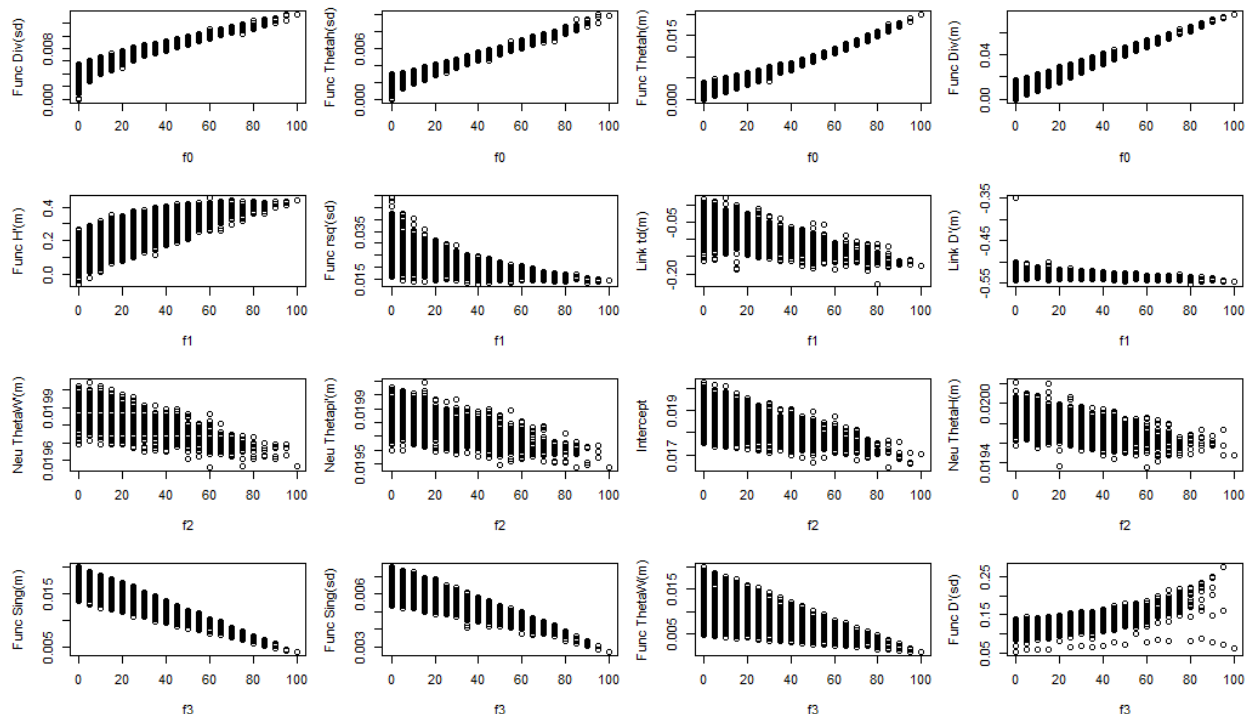
**Supp Figure 4:** Mis-inference of DFE in the presence of an additional unaccounted for 1 kb functional region near the target 1 kb exon used for inference. The intron/ intergenic distance between the two exons varies from 50-5000 bp, as shown by different colored bars. “0 bp” represents the negative control where there is no additional 1 kb exon present.



**Supp Figure 5:** Absolute difference between the true and estimated value of the DFE class, when the true recombination rate is half of that assumed for inference (orange) and when the true value is twice that assumed for inference (red), using a) all statistics and b) statistics only pertaining to the functional region.

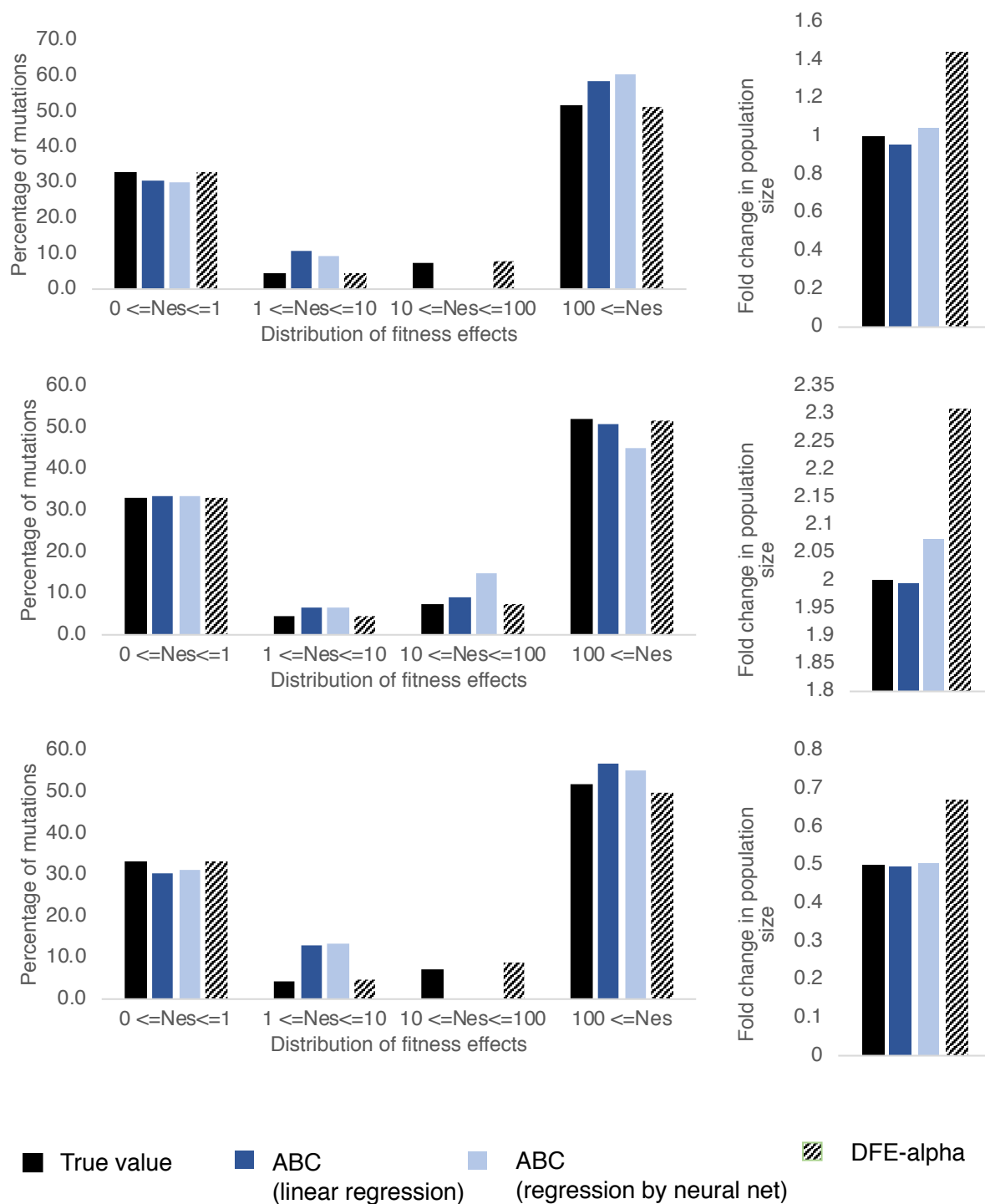


**Supp Figure 6:** Following Supp Figure 5, the direction of bias in inference of the DFE classes upon mis-specification of the recombination rate, using a) all statistics and b) statistics only pertaining to the functional region.

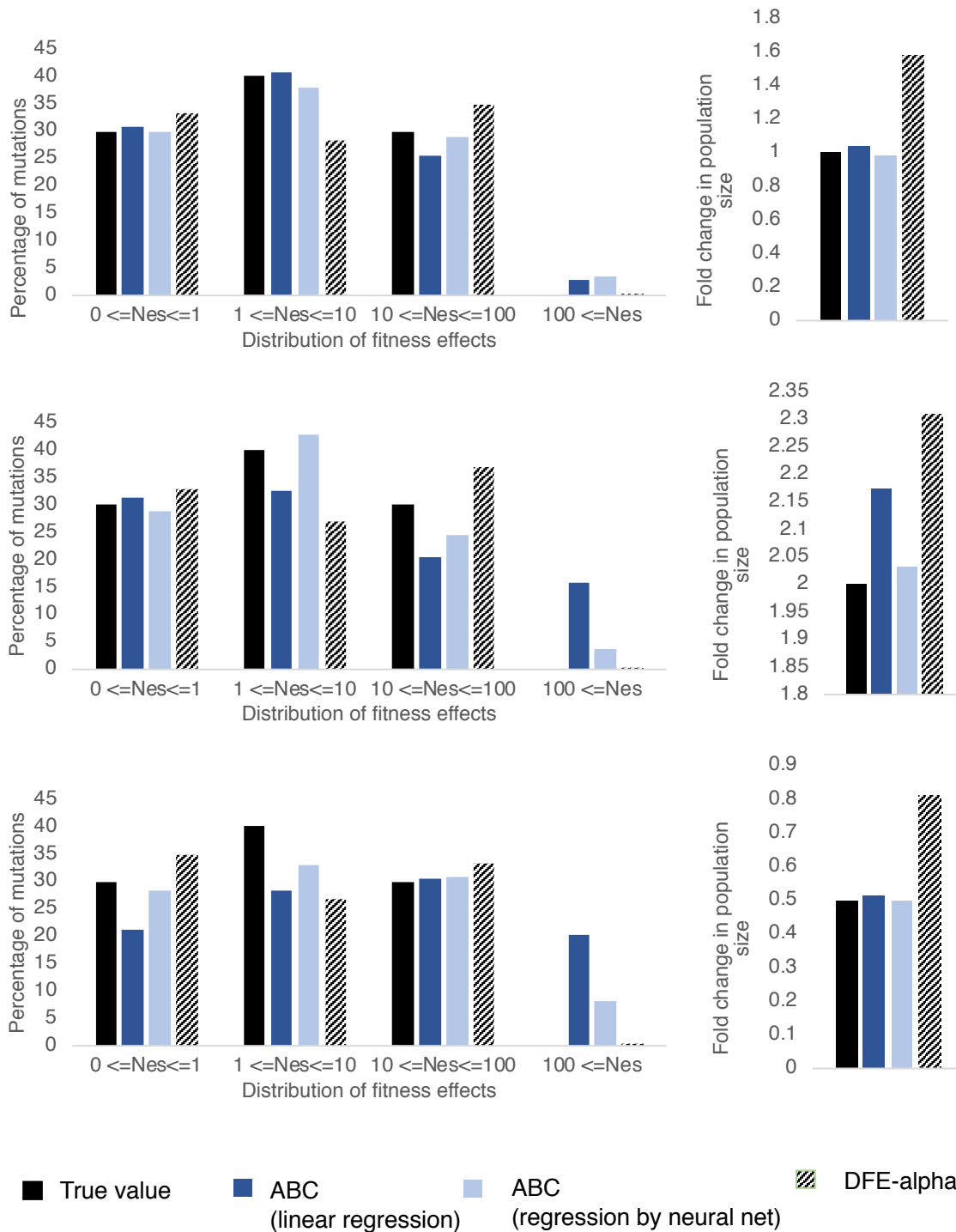


**Supp Figure 7:** Correlation of the top 4 statistics with parameters characterizing the DFE under demographic equilibrium. “Func” corresponds to the functional region, “Link” to the immediately linked region and “Neu” to the less linked region.

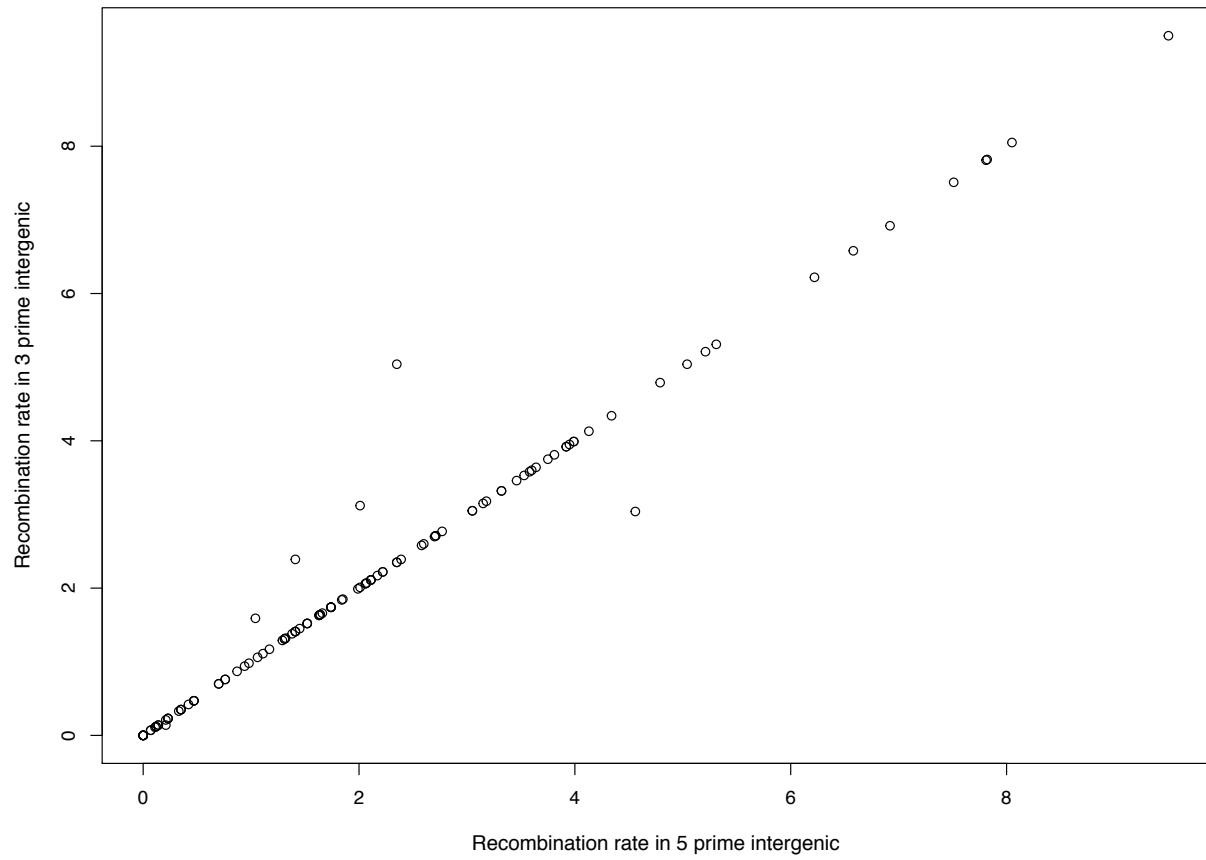




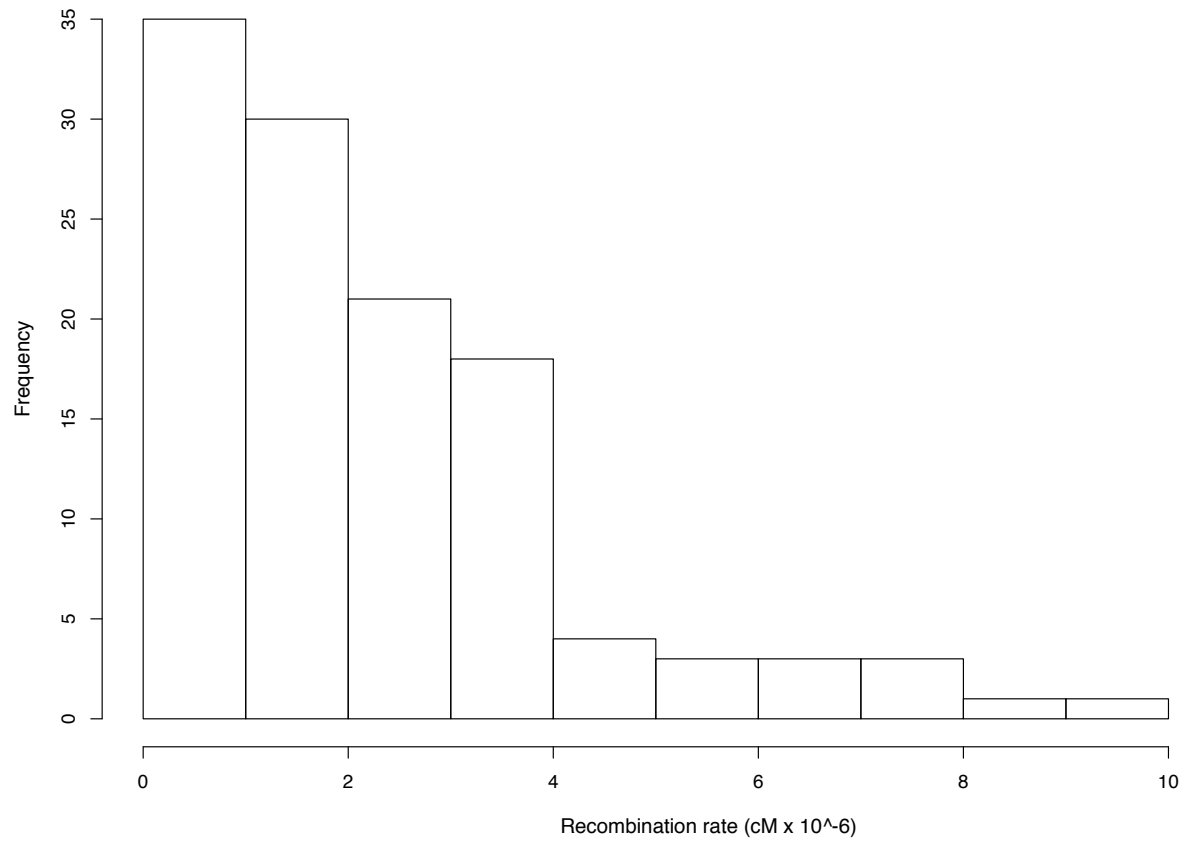
**Supp Figure 8:** Inference of demography and the DFE by the approach proposed here and DFE-alpha, when the true shape of the DFE is gamma distributed, for equilibrium (top panel), growth (middle panel), and decline (bottom panel). Solid black bars show the true value simulated, dark blue bars show our ABC performance using ridge regression, light blue bars show the ABC performance using linear regression aided by neural net. Patterned bars show the performance of DFE-alpha.



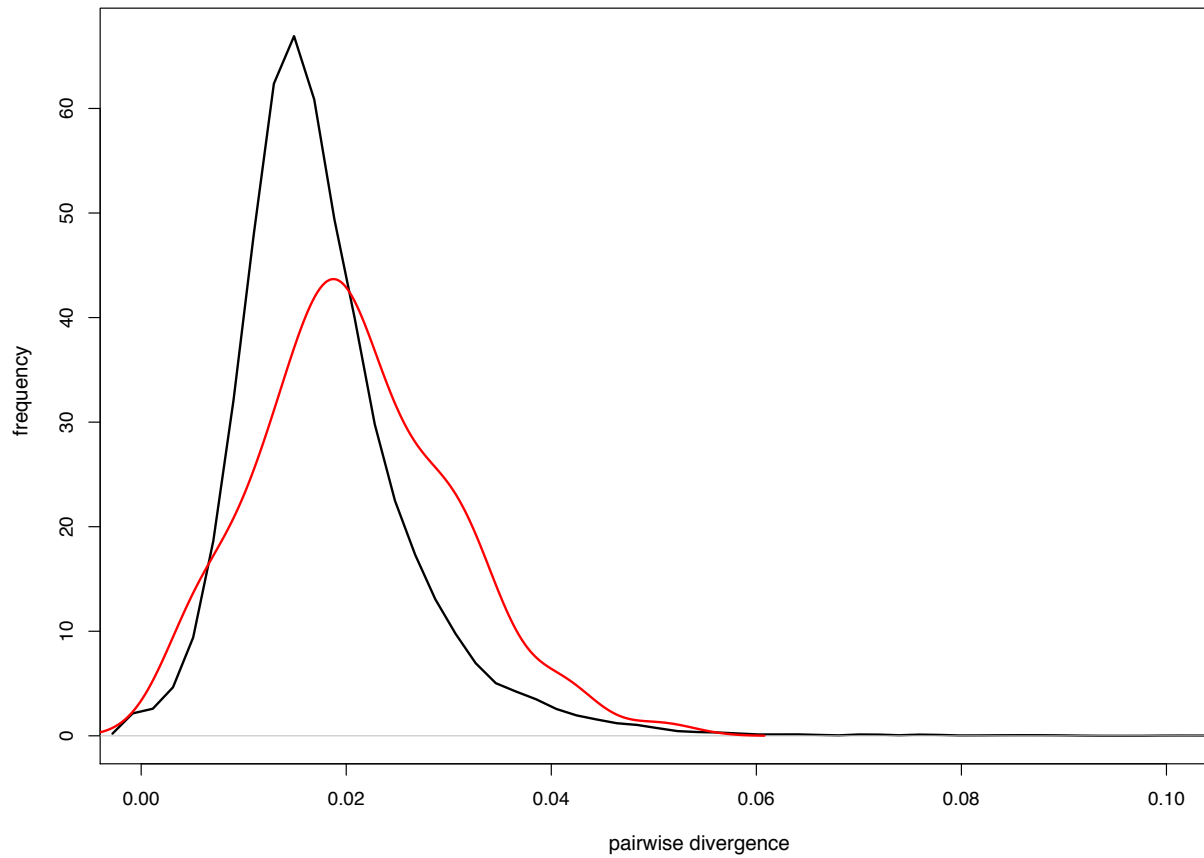
**Supp Figure 9:** Inference of demography and the DFE when the true shape of the DFE is discrete and skewed towards slightly deleterious class of mutations, for equilibrium (top panel), growth (middle panel), and decline (bottom panel). Solid black bars show the true value simulated, dark blue bars show the ABC performance using ridge regression, light blue bars show the ABC performance using linear regression aided by neural net. Patterned bars show the performance of DFE-alpha.



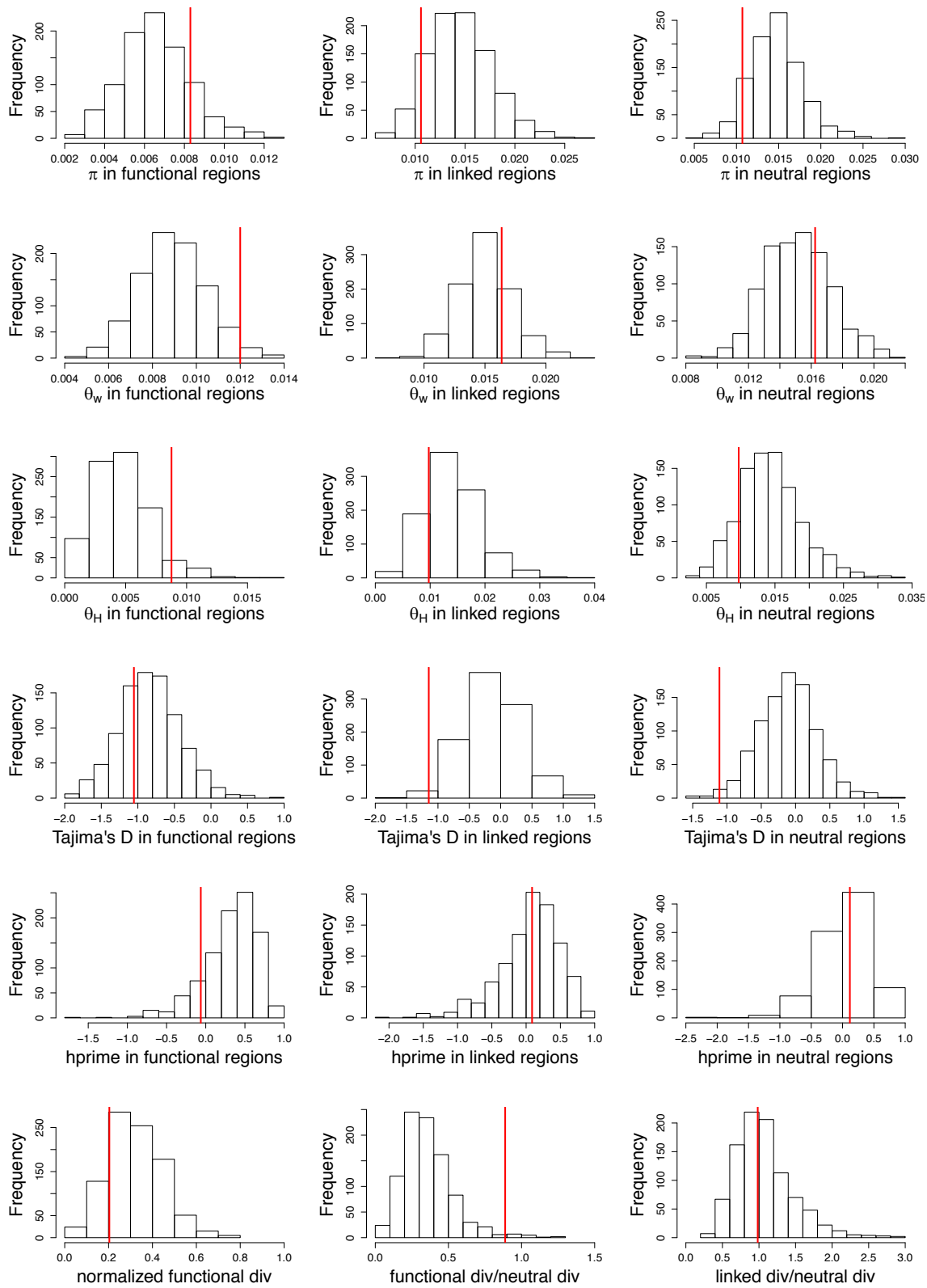
**Supp Figure 10:** Correlation of recombination rates at 5 prime flanking intergenic versus that in 3 prime flanking intergenic of all 94 exons chosen for analysis in *D. melanogaster*.

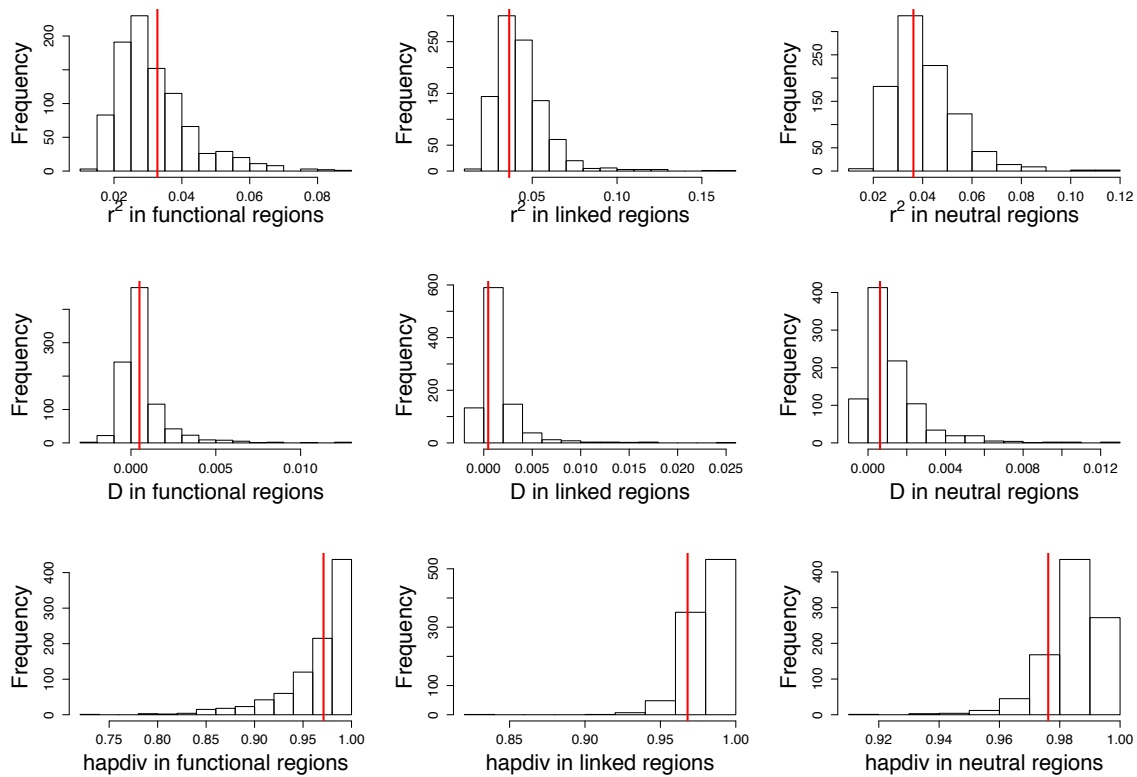


**Supp Figure 11:** Distribution of the rate of recombination in cM/Mb for all 94 exons selected for analysis in *D. melanogaster*.

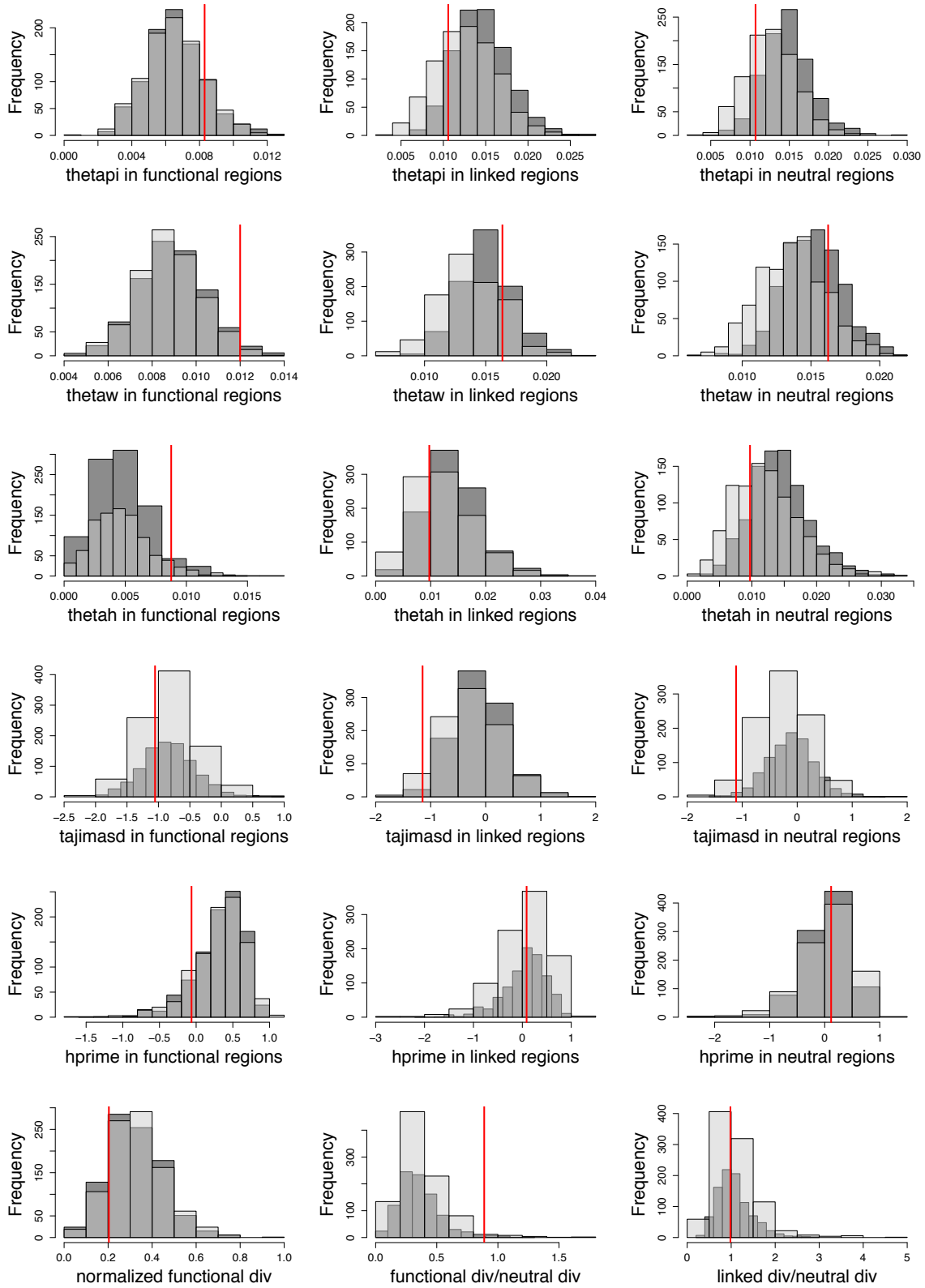


**Supp Figure 12:** Distribution of divergence per site of single-exon genes that have flanking intergenic regions larger than 4 kb (in red), and for all genes (in black), from *D. melanogaster*.

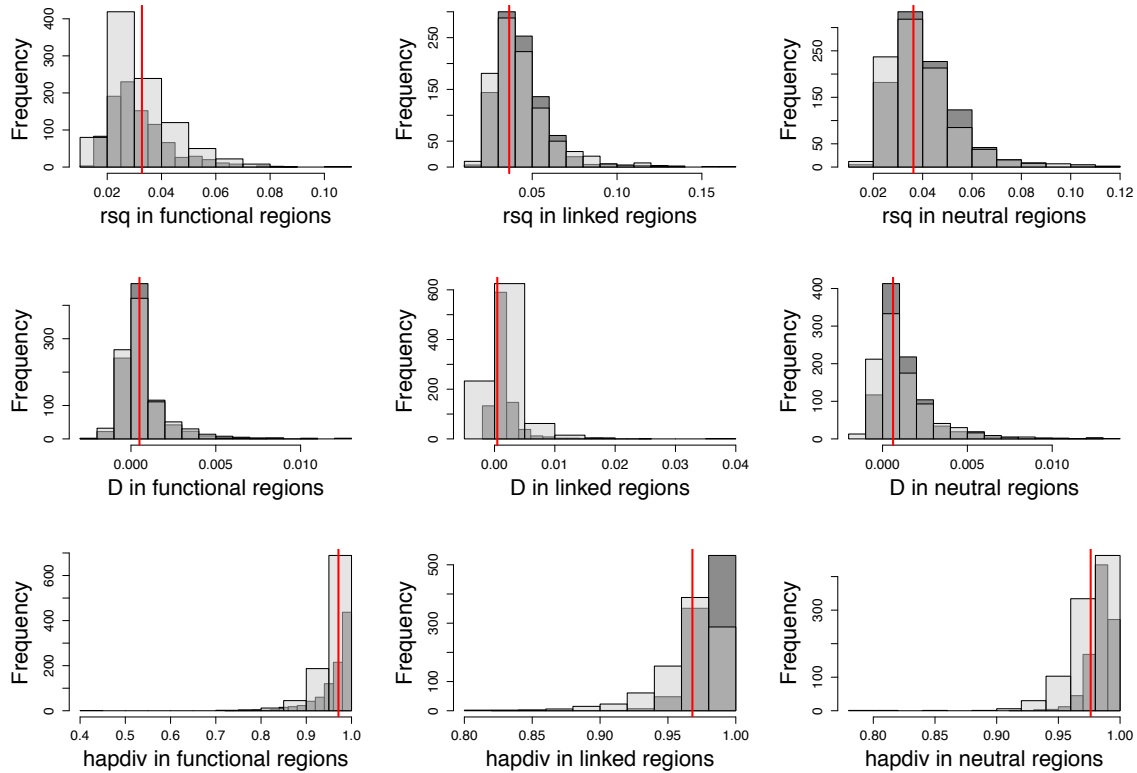




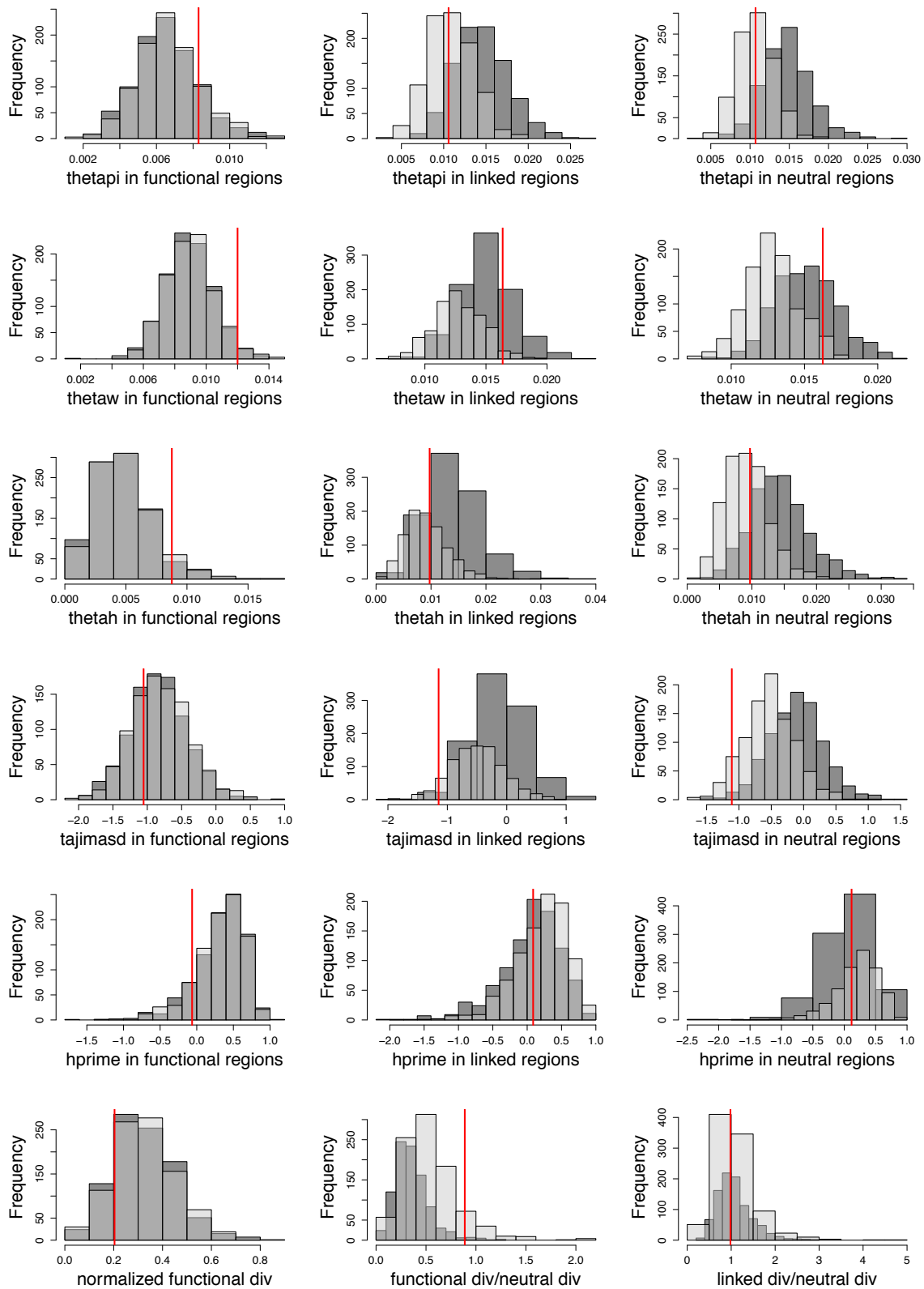
**Supp Figure 13:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{anc} = 1,225,393$ ,  $N_{cur} = 1,357,760$ ). Red lines indicate the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ .

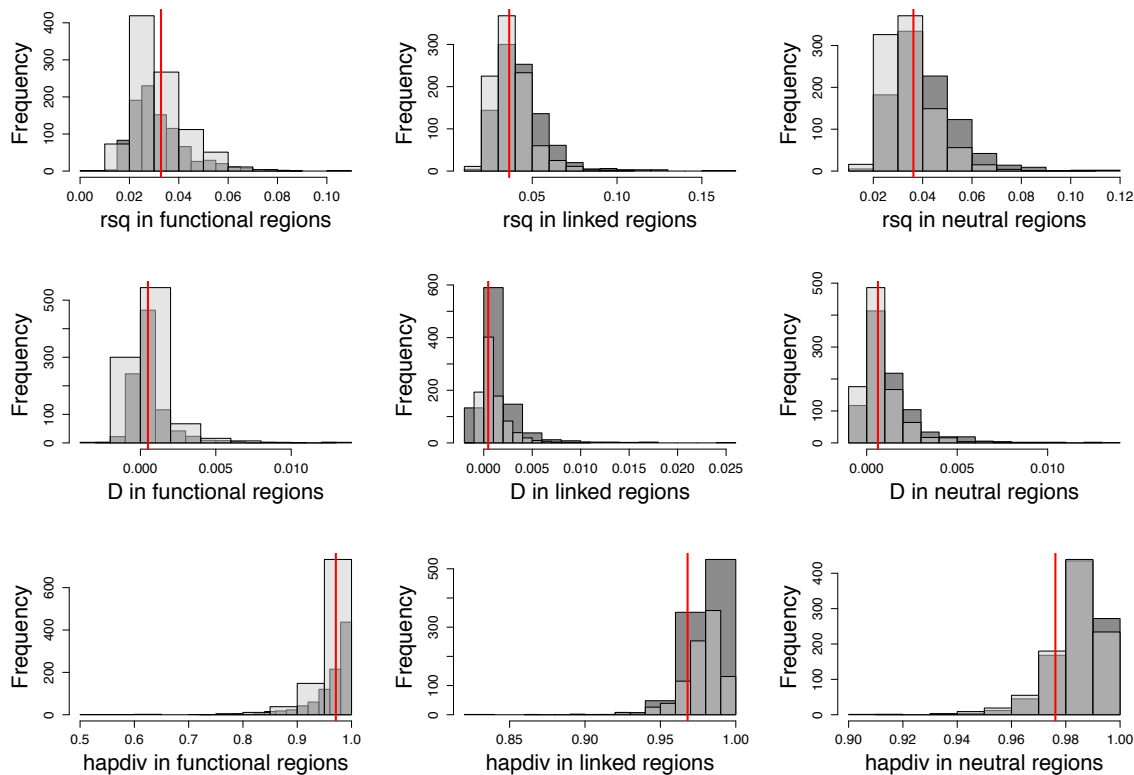




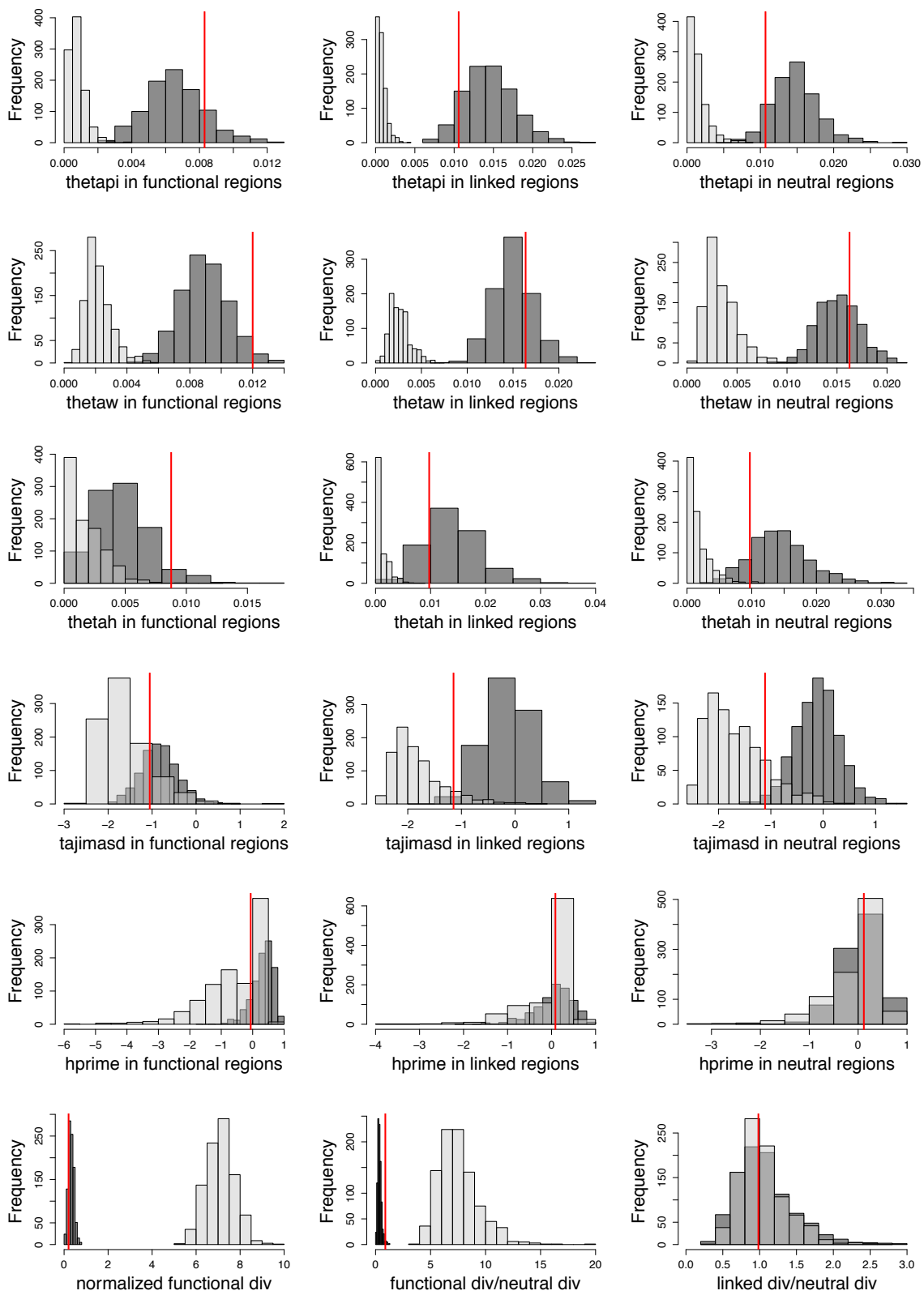


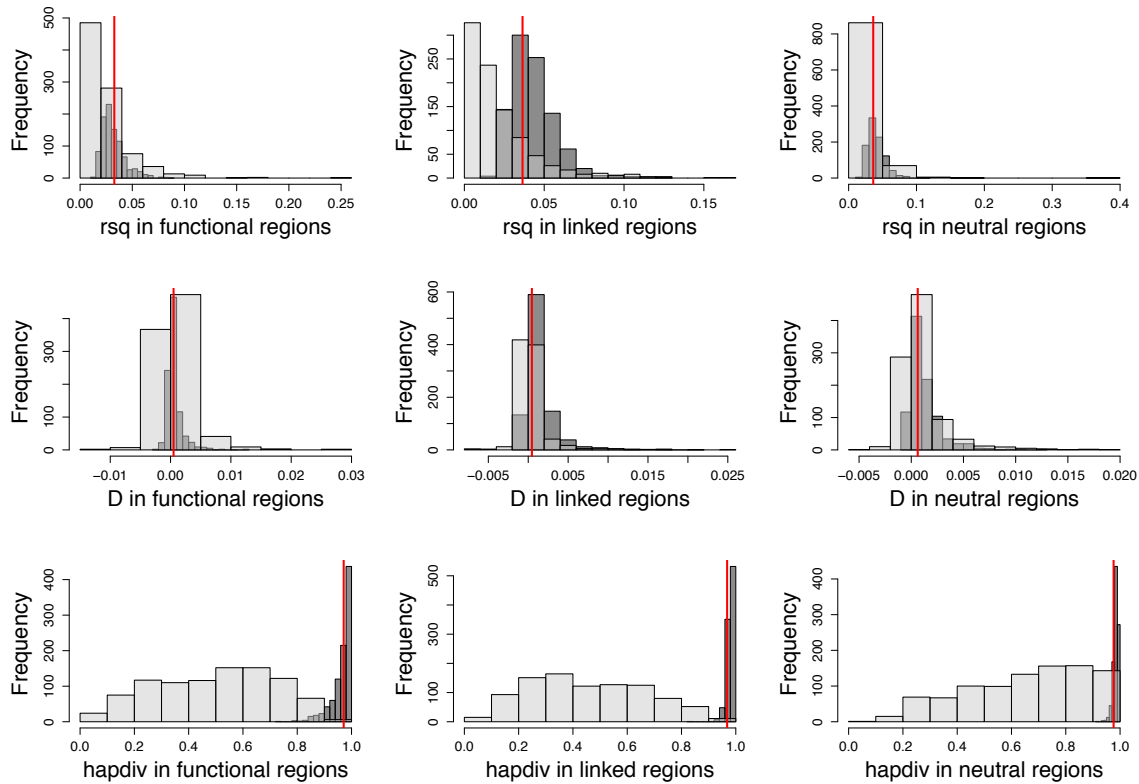
**Supp Figure 14:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{anc} = 1,225,393$ ,  $N_{cur} = 1,357,760$ ). In this case, conserved elements that represent 40% of non-coding regions were simulated to experience purifying selection with the class of mutations that result in strongest BGS effects ( $-100 < 2N_e s < -10$ ) and these sites were masked while calculating statistics. Red line indicates the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no selection on non-coding regions and light grey bars represent simulations with selection on non-coding regions.



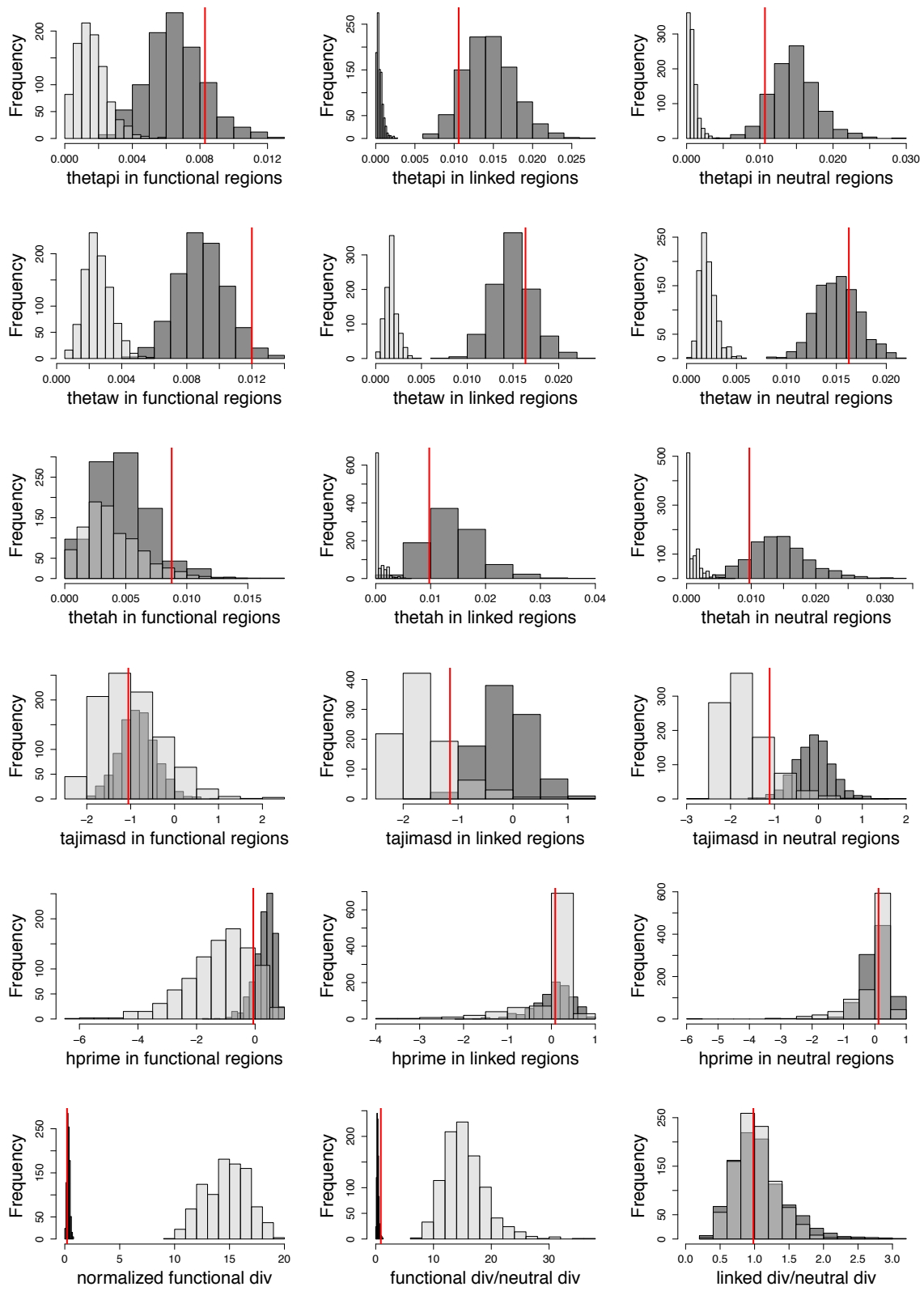


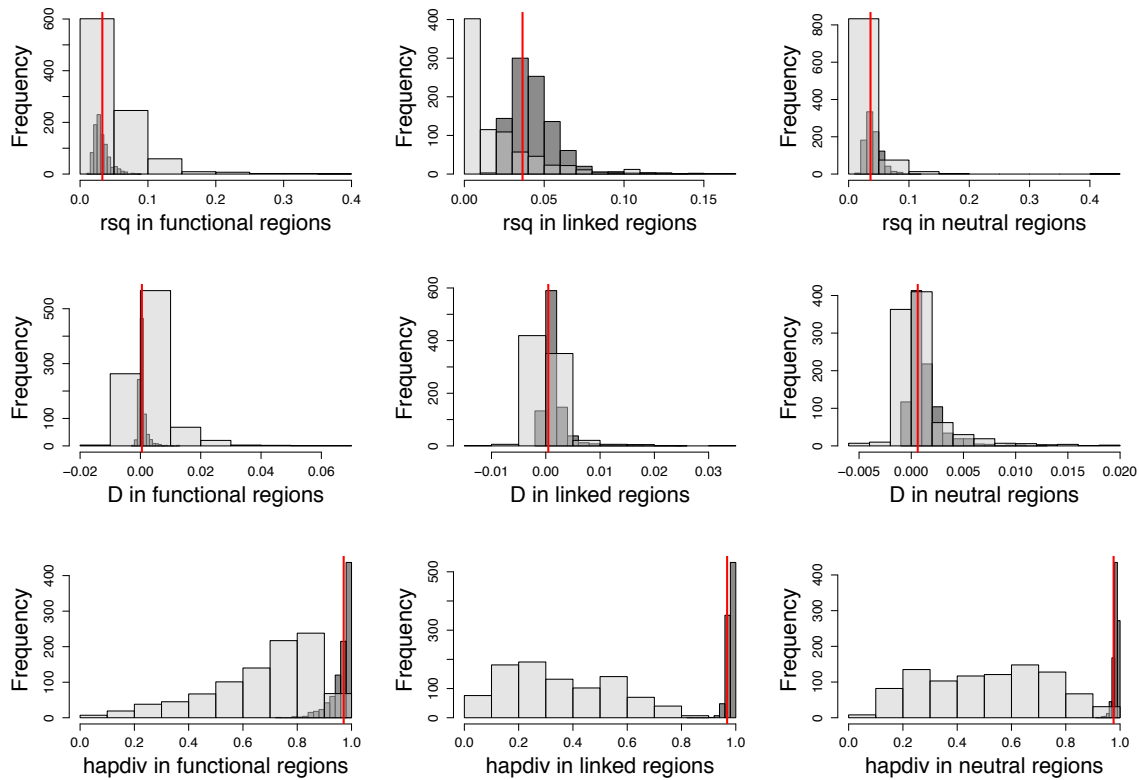
**Supp Figure 15:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{anc} = 1,225,393$ ,  $N_{cur} = 1,357,760$ ). In this case, conserved elements that represent 40% of non-coding regions were simulated to experience weak purifying selection ( $-10 < 2N_e s < -1$ ) and these sites were included while calculating statistics. Red line indicates the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no selection on non-coding regions and light grey bars represent simulations with selection on non-coding regions.



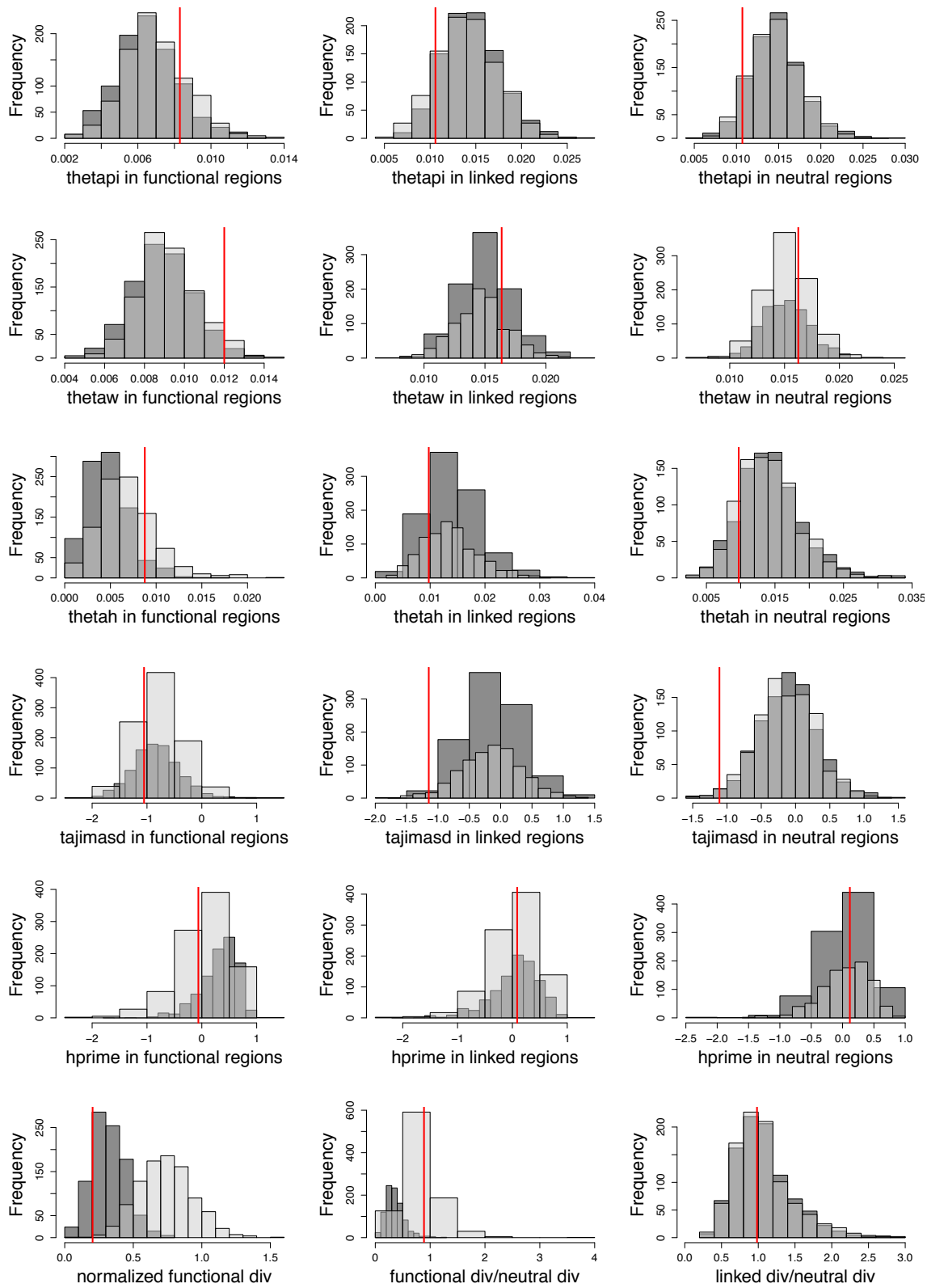


**Supp Figure 16:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{\text{anc}} = 1,225,393$ ,  $N_{\text{cur}} = 1,357,760$ ). Functional regions were simulated to experience rare (1%) and strong positive selection ( $2N_{\text{anc}s} = 1000$ ). Red lines indicate the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.

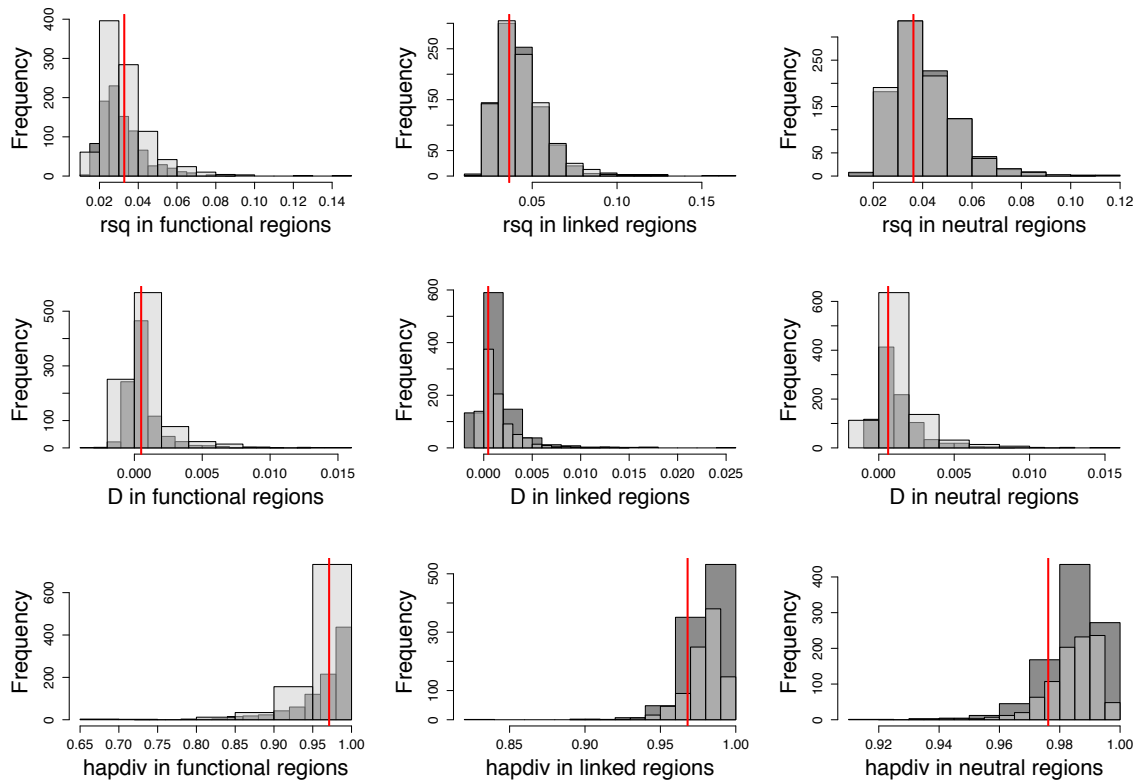




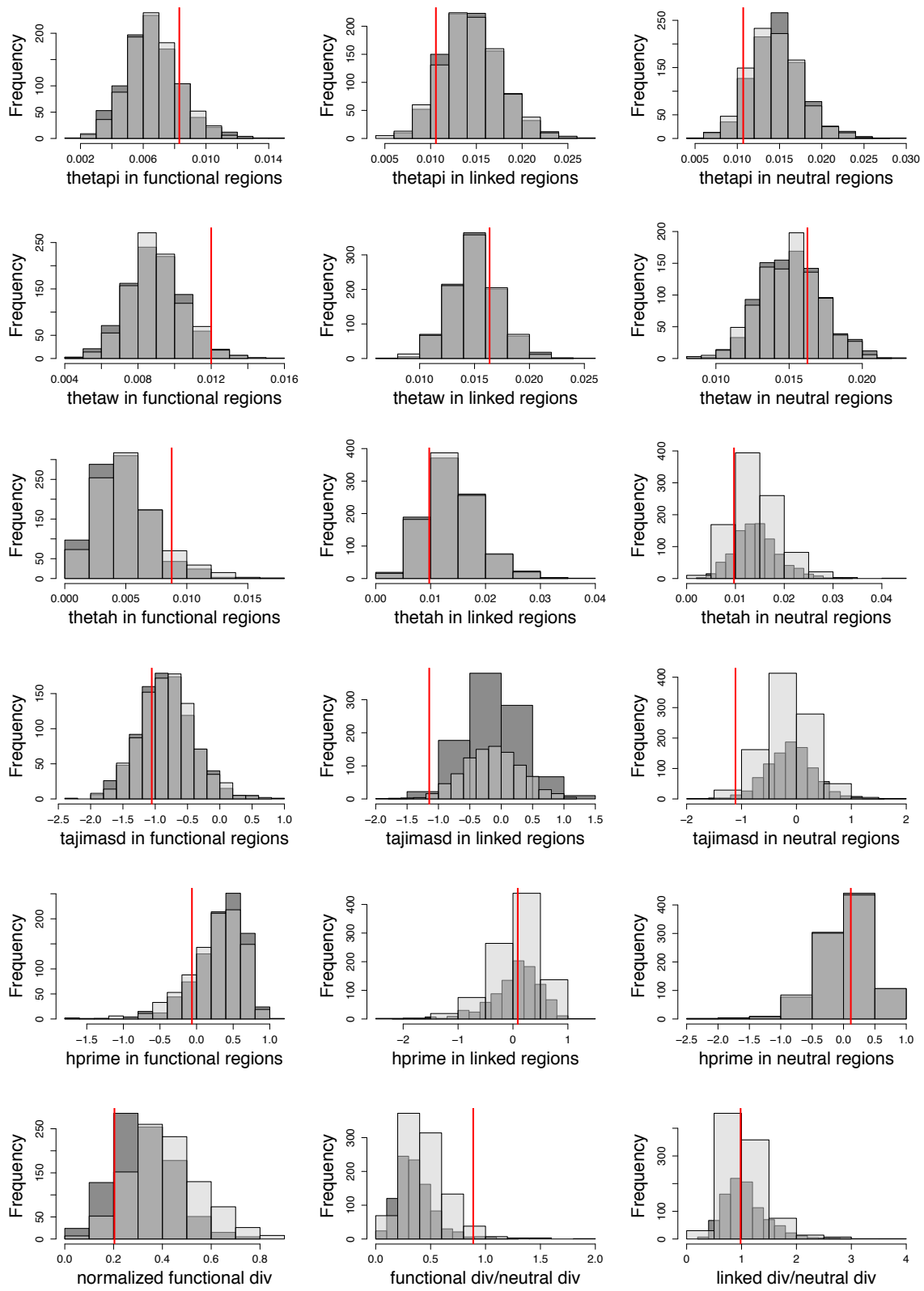
**Supp Figure 17:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1=0.49$ ,  $f_2=0.04$ ,  $f_3=0.22$ ,  $N_{anc}=1,225,393$ ,  $N_{cur} = 1,357,760$ ). Functional regions were simulated to experience common (5%) and strong positive selection ( $2N_{anc}s = 1000$ ). Red line indicates the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.

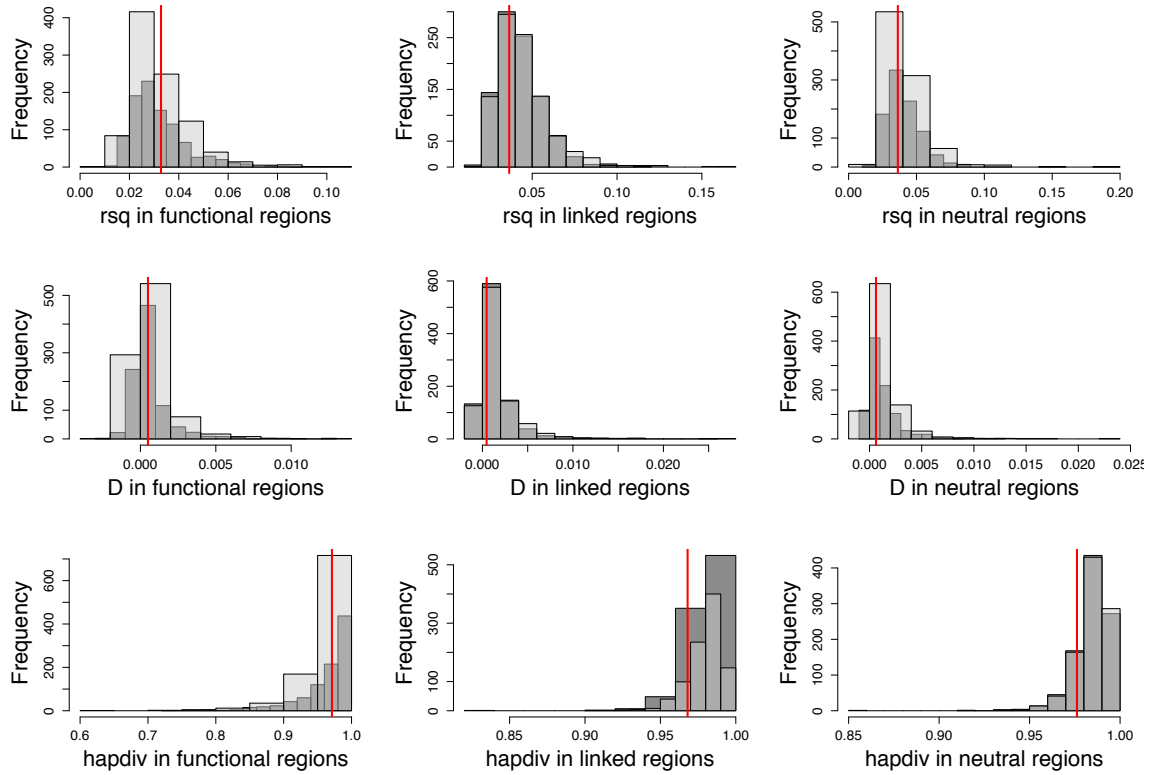




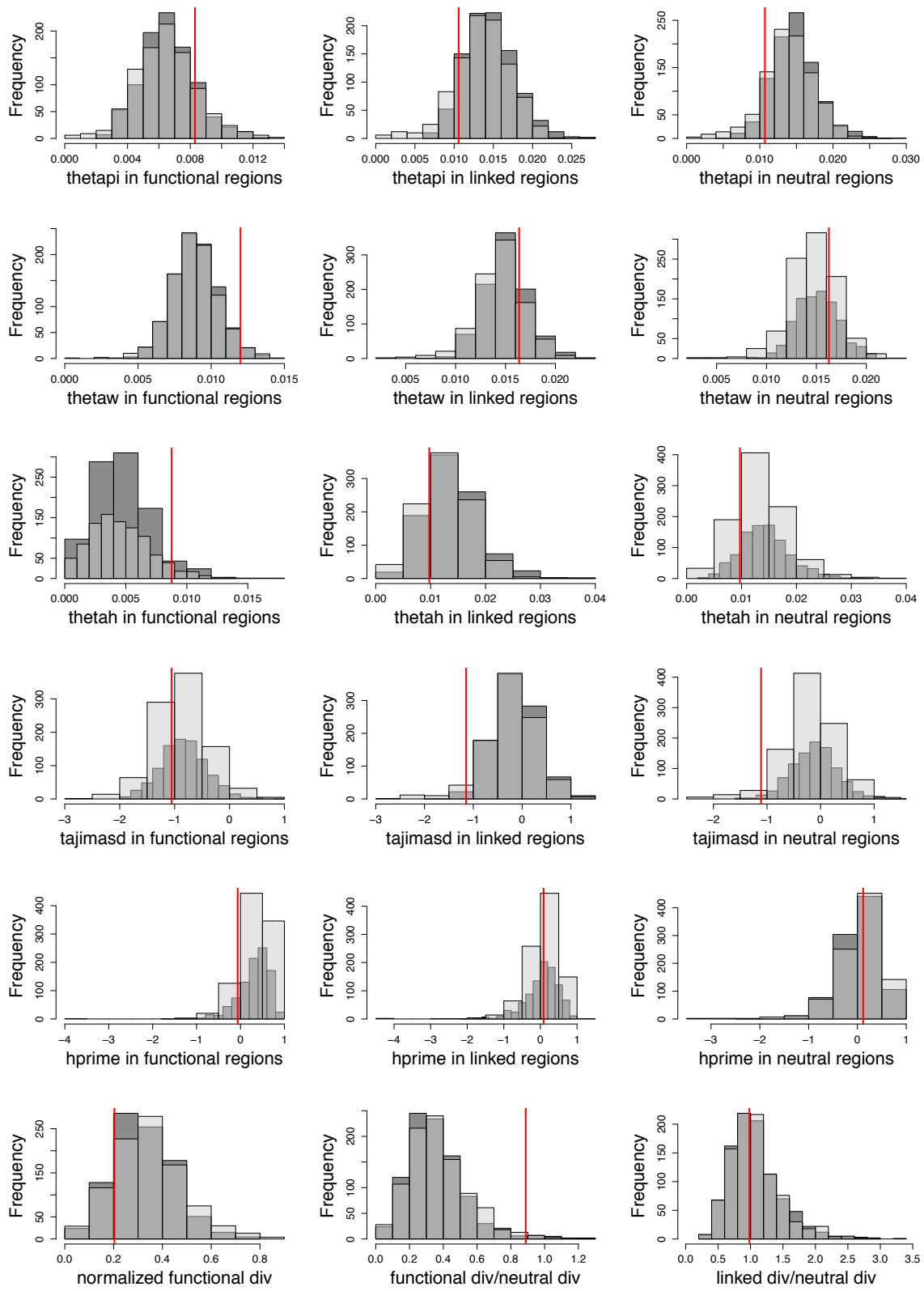


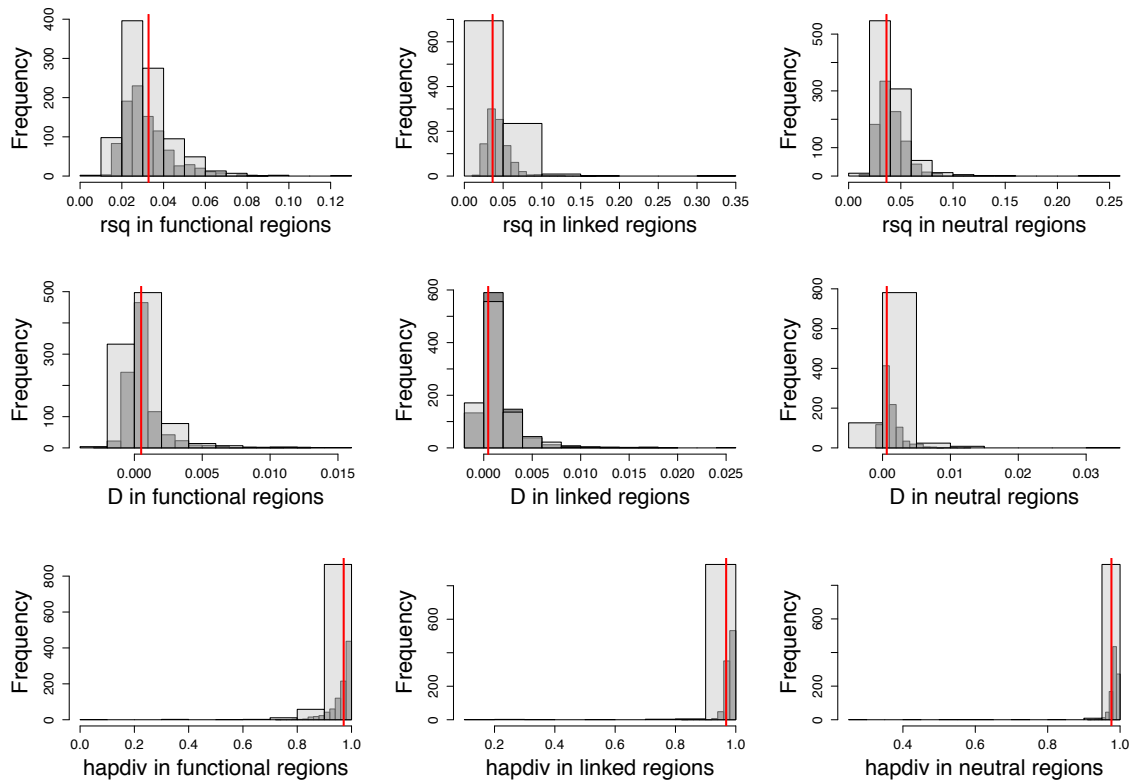
**Supp Figure 18:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{anc} = 1,225,393$ ,  $N_{cur} = 1,357,760$ ). Functional regions were simulated to experience common (5%) and weak positive selection ( $2N_{anc}s = 10$ ). Red lines indicate the value observed in 76 individuals of *Drosophila melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.



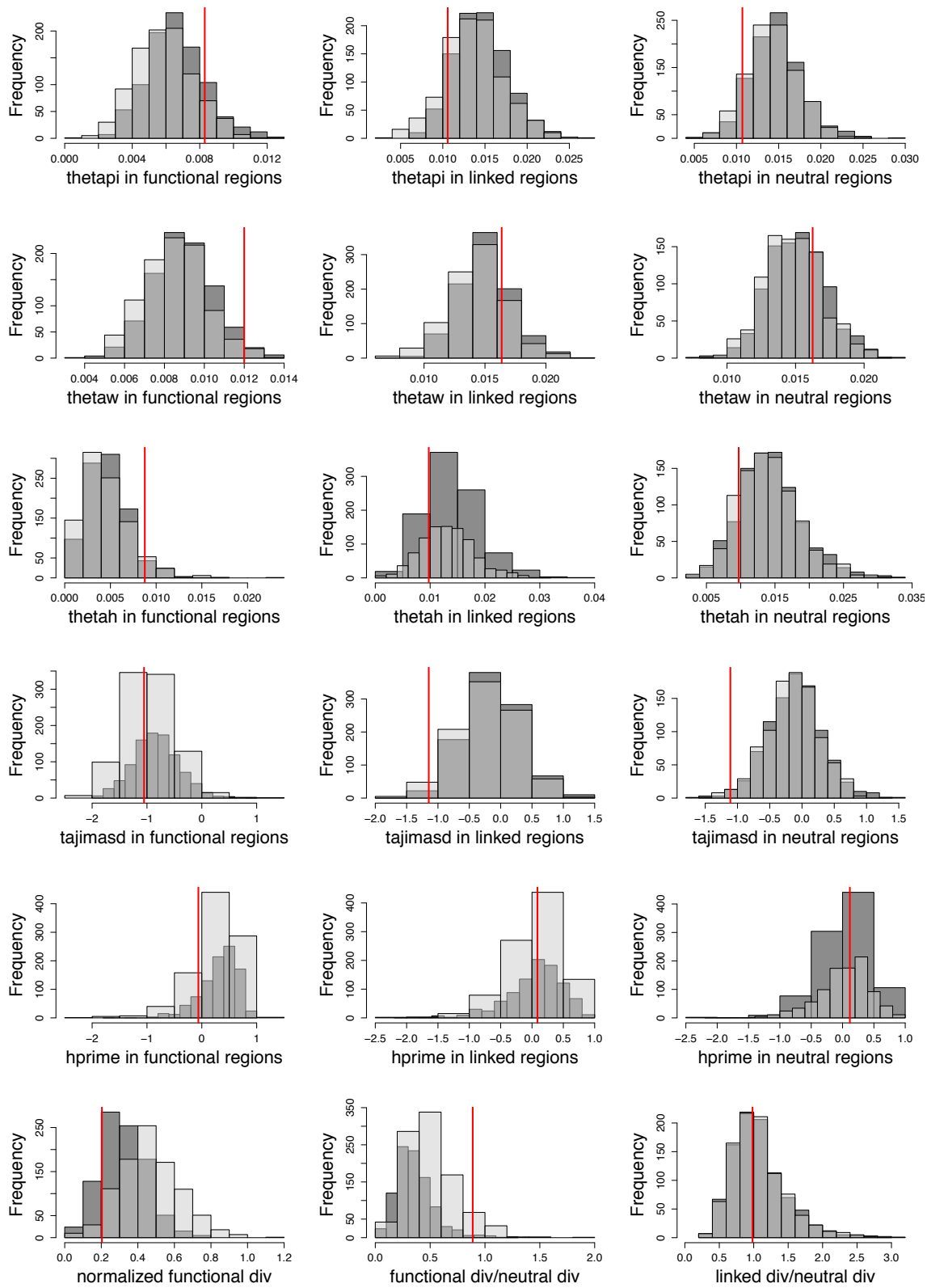


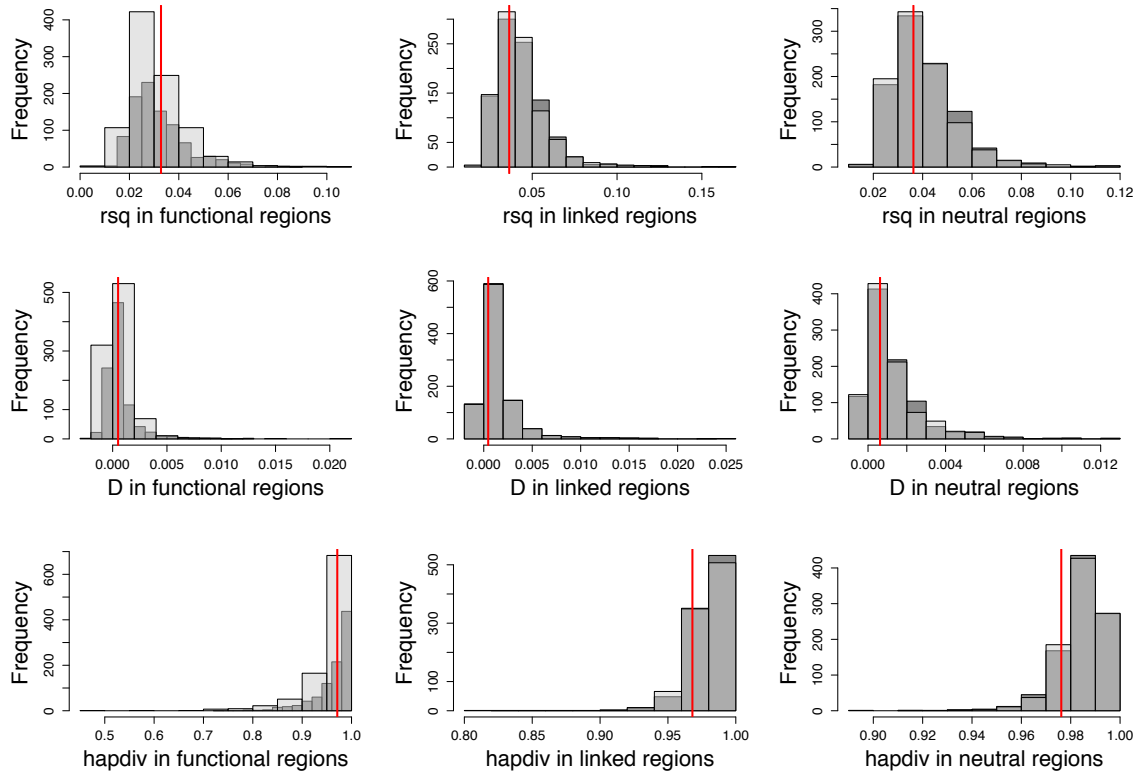
**Supp Figure 19:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1 = 0.49$ ,  $f_2 = 0.04$ ,  $f_3 = 0.22$ ,  $N_{anc} = 1,225,393$ ,  $N_{cur} = 1,357,760$ ). Functional regions were simulated to experience rare (1%) and weak positive selection ( $2N_{anc}s = 10$ ). Red lines indicate the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.





**Supp Figure 20:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1=0.49$ ,  $f_2=0.04$ ,  $f_3=0.22$ ,  $N_{anc}=1,225,393$ ,  $N_{cur} = 1,357,760$ ). Functional regions were simulated to experience rare ( $1.28 \times 10^{-4}$  %) and strong positive selection ( $2N_{anc}s = 10000$ ) as in Lange and Pool (2018). Red lines indicate the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.





**Supp Figure 21:** Distribution of summary statistics calculated from 94 exons simulated with 100 replicates each using our inferred model (*i.e.*,  $f_0 = 0.25$ ,  $f_1=0.49$ ,  $f_2=0.04$ ,  $f_3=0.22$ ,  $N_{anc}=1,225,393$ ,  $N_{cur} = 1,357,760$ ). Functional regions were simulated to experience rare (0.2%) and weak positive selection ( $2N_{anc}s = 60$ ) as in Lange and Pool (2018). Red lines indicate the value observed in 76 individuals of *D. melanogaster* from Zambia, after excluding sites with phastCons score  $\geq 0.8$ . Dark grey bars represent no positive selection and light grey bars represent simulations with positive selection in functional regions.