# Inferring parameters of the distribution of fitness effects of new mutations when beneficial mutations are strongly advantageous and rare

5  Tom R. Booker[1,2]

6

7  1. Department of Forest and Conservation Sciences, University of British Columbia, Vancouver,

8  Canada.

9  2. Biodiversity Research Centre, University of British Columbia, Vancouver, Canada

10

11  Correspondence: booker@zoology.ubc.ca

12

13

14

15

16

17

18

19

# Abstract

Characterising the distribution of fitness effects (DFE) for new mutations is central in evolutionary genetics. Analysis of molecular data under the McDonald-Kreitman test has suggested that adaptive substitutions make a substantial contribution to between-species divergence. Methods have been proposed to estimate the parameters of the distribution of fitness effects for positively selected mutations from the unfolded site frequency spectrum (uSFS). However, when beneficial mutations are strongly selected and rare, they may make little contribution to standing variation and will thus be difficult to detect from the uSFS. In this study, I analyse uSFS data from simulated populations subject to advantageous mutations with effects on fitness ranging from mildly to strongly beneficial. When advantageous mutations are strongly selected and rare, there are very few segregating in populations at any one time. Fitting the uSFS in such cases leads to underestimates of the strength of positive selection and may lead researchers to false conclusions regarding the relative contribution adaptive mutations make to molecular evolution. Fortunately, the parameters for the distribution of fitness effects for harmful mutations are estimated with high accuracy and precision. The results from this study suggest that the parameters of positively selected mutations obtained by analysis of the uSFS should be treated with caution and that variability at linked sites should be used in conjunction with standing variability to estimate parameters of the distribution of fitness effects in the future.

# Introduction

44

45    Characterising the distribution of fitness effects for beneficial mutations is central in evolutionary

46    biology. The rate and fitness effects of advantageous mutations may determine important

47    evolutionary processes such as how variation in quantitative traits is maintained (Hill, 2010), the

48    evolution of sex and recombination (Otto, 2009) and the dynamics of evolutionary rescue in

49    changing environments (Orr & Unckless, 2014). However, despite its central role in evolution,

50    relatively little is known about the distribution of fitness effects (DFE) for advantageous mutations

51    in natural populations. The DFE for advantageous mutations can be estimated from data obtained

52    via targeted mutation or from mutation accumulation experiments (e.g. Bank, Hietpas, Wong,

53    Bolon, & Jensen, 2014; Böndel et al., 2019; reviewed in Bailey & Bataillon, 2016), but such efforts

54    may be limited to laboratory systems. Alternatively, estimates of the DFE can be obtained for

55    natural systems using population genetic methods.

56

57    When natural selection is effective, beneficial alleles are promoted to eventual fixation while

58    deleterious variants are maintained at low frequencies. Migration, mutation, selection and genetic

59    drift interact to shape the distribution of allele frequencies in a population (Wright, 1937).

60    Parameters of the DFE for both advantageous and deleterious mutations can be estimated by

61    modelling population genomic data, specifically the site frequency spectrum (SFS). The SFS is the

62    distribution of allele frequencies present in a sample of individuals drawn from a population. By

63    contrasting the SFS for a class of sites expected to be subject to selection with that of a neutral

64    comparator, one can estimate the parameters of the DFE if selected mutations are segregating in

65    the population of interest (reviewed in Eyre-Walker & Keightley, 2007). Typically, the DFE for

66    nonsynonymous sites in protein coding genes is estimated using synonymous sites as the neutral

67    comparator. Several methods have been proposed that estimate the DFE for deleterious

3

68    mutations from the SFS under the assumption that beneficial mutations contribute little to

69    standing genetic variation (e.g. Barton & Zeng, 2018; Boyko et al., 2008; Keightley & Eyre-Walker,

70    2007; Tataru, Mollion, Glemin, & Bataillon, 2017).

71

72    The DFE for deleterious mutations can be used when estimating $\alpha$, the proportion of between-

73    species divergence attributable to adaptive evolution (Eyre-Walker & Keightley, 2009). $\alpha$ can be

74    estimated by rearranging the terms of the McDonald-Kreitman test (MK-test), which assesses the

75    extent of positive selection. Under strong purifying selection, the ratio of divergence at

76    nonsynonymous sites ($d_N$) to that of synonymous sites ($d_S$) should be exactly equal to the ratio of

77    nucleotide diversity at nonsynonymous ($\pi_N$) and synonymous sites ($\pi_S$)(McDonald & Kreitman,

78    1991). Adaptive evolution of protein sequences may contribute to $d_N$ such that $d_N/d_S > \pi_N/\pi_S$.

79    Charlesworth (1994) suggested rearranging the terms of the MK-test to estimate the excess $d_N$ due

80    to positive selection ($\alpha$) as

$$\alpha = 1 - d_S\pi_N/d_N\pi_S.$$

82    Slightly deleterious alleles may contribute to both standing genetic variation and between-species

83    divergence, estimates of $\alpha$ may therefore be refined by subtracting the contribution that

84    deleterious alleles make to both polymorphism and divergence and this can be calculated using

85    the DFE for harmful mutations (Eyre-Walker & Keightley, 2009). Application of such methods to

86    natural populations suggest that $\alpha$ is of the order of 0.5 in a large variety of animal taxa (Galtier,

87    2016). However, if adaptive evolution is as frequent as MK-test analyses suggest, the assumption

88    that advantageous alleles contribute little to standing variation may be violated and ignoring them

89    could lead to biased estimates of the DFE (Tataru et al., 2017).

90

91    When advantageous alleles contribute to standing variation, parameters of the DFE for both

92    deleterious and beneficial mutations can be estimated from the SFS (Schneider et al., 2011; Tataru

4

93   et al., 2017). When data from an outgroup species are available, variable sites within a focal

94   species can be polarised as either ancestral or derived and the *unfolded* SFS (uSFS) can be

95   obtained. Inference of ancestral/derived states is, however, potentially error-prone (Keightley &

96   Jackson, 2018). The uSFS is a vector of length $2n$, where $n$ is the number of haploid genome copies

97   sampled. The $i^{th}$ entry of the uSFS is the count of derived alleles observed at a frequency $i$ in the

98   sample. Note that when outgroup data are not available, alleles cannot be polarised and the

99   distribution of minor allele frequencies (known as the *folded* SFS) is analysed. There is limited

100  power to detect positive selection from the SFS, so the DFE for beneficial mutations is often

101  modelled as a discrete class of mutational effects, with one parameter specifying the fitness

102  effects of beneficial mutations, $\gamma_a = 2N_e s_a$ where $N_e$ is the effective population size and $s_a$ is the

103  positive selection coefficient in homozygotes, and another specifying the proportion of new

104  mutations that are advantageous, $p_a$. Estimates of $\gamma_a$ and $p_a$ for nonsynonymous sites have only

105  been obtained a handful of species, and these are summarised in Table 1. The positive selection

106  parameter estimates that have been obtained for mice and *Drosophila* are fairly similar (Table 1).

107  Note that the estimates for humans obtained by Castellano et al, (2019) did not provide a

108  significantly greater fit to the observed data than did a model with no positive selection.

109  Furthermore, Castellano et al, (2019) estimated the parameters for numerous great ape species,

110  the parameters shown for humans are representative of the estimates for all taxa they analysed.

111
112  **Table 1** Estimates of the parameters of positive selection obtained from the uSFS for
113  nonsynonymous sites.

| Common name | Scientific name | $\gamma_a$ | $p_a$ | Reference | Method used[¶] |
|---|---|---|---|---|---|
| House mouse | *Mus musculus castaneus* | 14.5 | 0.0030 | Booker & Keightley, (2018) | *DFE-alpha* |
| Fruit fly | *Drosophila melanogaster* | 23.0 | 0.0045 | Keightley et al, (2016) | *DFE-alpha* |
| Humans | *Homo sapiens* | 0.0064[†] | 0.000025 | Castellano et al, (2019) | *polyDFE* |

114  ¶ - DFE-alpha implements the analysis methods described by Schneider et al., (2011), *polyDFE* implements the
115  methods described by Tataru et al., (2017)

116    † - Castellano et al., (2019) estimated the mean fitness effect for an exponential distribution of advantageous
117    mutational effects.

118

119    Depending on the rate and fitness effects of beneficial mutations, different aspects of population

120    genomic data may be more or less informative for estimating the parameters of positive selection.

121    As beneficial mutations spread through populations, they may carry linked neutral variants to high

122    frequency, causing selective sweeps (Barton, 2000). On the other hand, if advantageous mutations

123    have mild fitness effects, they may take a long time to reach fixation and make a substantial

124    contribution to standing genetic variation. Because of this, uSFS data and polymorphism data at

125    linked sites may both be informative for understanding the parameters of positive selection. For

126    example, Campos et al., (2017) used a model of selective sweeps to analyse the negative

127    correlation observed between $d_N$ and $\pi_S$ in *Drosophila melanogaster* and estimated $\gamma_a$ = 250 and

128    $p_a$ = 2.2 x $10^{-4}$, but this method assumes a constant population size.  An analysis of the uSFS from

129    the same dataset that modelled of population size change yielded estimates of $\gamma_a$ = 23 and $p_a$ =

130    0.0045 for nonsynonymous sites (Keightley et al., 2016). The sharp contrast between the two

131    studies' estimates of the positive selection parameters may due to different assumptions but

132    could potentially be explained if the DFE for advantageous mutations in *D. melanogaster* is

133    bimodal. If this were so, the different methods (i.e. sweep models versus uSFS analysis) may be

134    capturing distinct aspects of the DFE for advantageous mutations, or it could be that both models

135    are highly unidentifiable. The handful of studies that have attempted to estimate $\gamma_a$ and $p_a$ from

136    the uSFS have yielded similar estimates of positive selection (Table 1), which may indicate

137    commonalities in the DFE for beneficial mutations across taxa. On the other hand, uSFS analyses

138    may have only found evidence for mildly beneficial mutations because the approach is only

139    powered to detect weakly beneficial mutations. Indeed, verbal arguments have suggested that

140    rare strongly selected advantageous mutations, which may contribute little to standing variation,

141    will be undetectable by analysis of the uSFS (Booker & Keightley, 2018; Campos et al., 2017).

6

142

143    The studies describing the two most recently proposed methods for estimating the DFE for

144    beneficial mutations from the uSFS (Schneider et al., 2011; Tataru et al., 2017) performed

145    extensive simulations, but did not test cases of rare advantageous mutations with strong effects

146    on fitness. Testing this case is important, as studies that have analysed patterns of putatively

147    neutral genetic diversity across the genome have indicated that the DFE for advantageous

148    mutations contains strongly beneficial mutations in a variety of taxa (Booker & Keightley, 2018;

149    Campos et al., 2017; Elyashiv et al., 2016; Nam et al., 2017; Uricchio et al., 2019). Note that Tataru

150    et al., (2017) did simulate a population subject to frequent strongly beneficial mutations ($\gamma_a$ = 800

151    and $p_a$ = 0.02), but the parameter combination they tested may not be biologically relevant as the

152    proportion of adaptive substitutions it yielded was far higher than is typically estimated from real

153    data ($\alpha$ = 0.99). The limited parameter ranges tested in the simulations performed by Schneider et

154    al., (2011) and Tataru et al., (2017) leave a critical gap in our knowledge as to how uSFS based

155    methods perform when advantageous mutations are strongly selected and infrequent.

156

157    In this study, I use simulated datasets to fill this gap and examine how uSFS-based analyses

158    perform when beneficial mutations are strongly selected and rare. I simulate populations subject

159    to a range of positive selection parameters, including cases similar to those modelled by Tataru et

160    al., (2017) and cases where beneficial mutations are strongly selected but infrequent. It has been

161    pointed out that estimating selection parameters by modelling within species polymorphism along

162    with between-species divergence makes the assumption that the DFE has remained invariant since

163    the ingroup and outgroup began to diverge (Tataru et al., 2017). By analysing only the

164    polymorphism data, one can potentially avoid that problematic assumption. Using the state-of-

165    the-art package *polyDFE* v2.0 (Tataru & Bataillon, 2019), I analyse the uSFS data and estimate

166    selection parameters for all simulated datasets with or without divergence.  The results from this

167  study suggest that, when beneficial mutations are strongly selected and rare, analysis of the uSFS

168  results in spurious parameter estimates and the proportion of adaptive substitutions may be

169  poorly estimated.

# Methods

## Population genomic simulations

172  I tested the hypothesis that the parameters of infrequent, strongly beneficial mutations are

173  difficult to estimate by analysis of the uSFS using simulated datasets. Wright-Fisher populations of

174  $N_e$ = 10,000 diploid individuals were simulated using the forward-in-time package *SLiM* (v3.2;

175  Haller & Messer, 2019). Simulated chromosomes consisted of seven gene models, each separated

176  by 8,100bp of neutrally evolving sequence. The gene models consisted of five 300bp exons

177  separated by 100bp neutrally evolving introns. The gene models were based on those used by

178  Campos & Charlesworth, (2019), but unlike that study, I did not model the untranslated regions of

179  genes. Nonsynonymous sites were modelled by drawing the fitness effects for 2/3rds of mutations

180  in exons from a distribution of fitness effects (DFE), while the remaining 1/3 were strictly neutral

181  and used to model synonymous sites. The fitness effects of nonsynonymous mutations were

182  beneficial with probability $p_a$ or deleterious with probability $1 - p_a$. Beneficial mutations had a

183  fixed selection coefficient of $\gamma_a = 2N_e s_a$. The fitness effects of deleterious mutations were drawn

184  from a gamma distribution with a mean of $\gamma_d = 2N_e s_d = -2,000$ and a shape parameter of $\beta = 0.3$ ($s_d$

185  being the negative selection coefficient in homozygotes). The gamma distribution of deleterious

186  mutational effects was used for all simulated datasets and was based on results for

187  nonsynonymous sites in *Drosophila melanogaster* (Loewe & Charlesworth, 2006). Uniform rates of

188  mutation ($\mu$) and recombination ($r$) were set to 2.5 x $10^{-7}$ (giving $4N_e r = 4N_e \mu = 0.01$). Note that $\mu$

189  and $r$ are far higher than is biologically realistic for most eukaryotes, I scaled up these rates to

190    model a population with a large $N_e$ using simulations of 10,000 individuals. Across simulations I

191    varied the $\gamma_a$ and $p_a$ parameters and performed 2,000 replicates for each combination of

192    parameters. Thus, I simulated a dataset of 21Mbp of coding sequence for each combination of $\gamma_a$

193    and $p_a$ tested.

194

195    In this study, I assumed a discrete class of beneficial mutational effects rather than a continuous

196    distribution, which is likely unrealistic for most organisms. Theoretical arguments have been

197    proposed that the DFE for beneficial mutations that go to fixation should be exponential (Orr,

198    2003). However, the studies that have estimated the DFE for beneficial mutations from population

199    genetic data have often modelled discrete classes of effects (Campos et al., 2017; Elyashiv et al.,

200    2016; Keightley et al., 2016; Uricchio et al., 2019). I chose to model discrete selection coefficients

201    in the simulated datasets in order to better understand the limitations of the methods rather than

202    to accurately model the DFE for beneficial mutations.

203

204    To model the accumulation of nucleotide substitutions after the split of a focal population with an

205    outgroup, I recorded all substitutions that occurred in the simulations. Campos & Charlesworth,

206    (2019) analysed simulations very similar to those that I performed in this study and showed that

207    populations subject to beneficial mutations with $\gamma_a$ = 250 and $p_a$ = 0.0002 took 14$N_e$ generations

208    to reach mutation-selection-drift equilibrium. In this study I modelled a range of positive selection

209    parameters, so to ensure that my simulations reached equilibrium I performed 85,000 (34$N_e$)

210    generations of burn-in before substitutions were scored. The expected number of neutral

211    nucleotide substitutions that accumulate per site in $T$ generations is $d_{Neutral} = T\mu$. The point

212    mutation rate in my simulations was set to $\mu$ = 2.5 x $10^{-7}$ per site per generation, so I ran the

213    simulations for 200,000 generations beyond the end of the burn-in phase to model a neutral

214 divergence of $d_{Neutral}$ = 0.05. All variants present in the population sampled at a frequency of 1.0

215 were also scored as substitutions.

216

217 Using the 2,000 simulated datasets, I constructed 100 bootstraps by sampling with replacement.

218 From each bootstrap sample, I collated variants and constructed the uSFS for synonymous and

219 nonsynonymous sites for 20 diploid individuals.

220

221 Analysis of simulation data

222 I calculated several summary statistics from the simulated datasets. Firstly, I calculated pairwise

223 nucleotide diversity at synonymous sites ($\pi_s$) and expressed it relative to the neutral expectation

224 of $\pi_0 = 4N_e\mu$ = 0.01. Secondly, divergence at nonsynonymous sites for both advantageous ($dN_a$)

225 and deleterious mutations ($dN_d$) was used to calculate the observed proportion of adaptive

226 substitutions, $\alpha_{Obs} = dN_a/(dN_a + dN_d)$. Finally, I recorded the total number of beneficial mutations

227 segregating in simulated populations, $S_{Adv}$.

228

229 I estimated DFEs from simulated data by analysis of the uSFS using *polyDFE* (v2.0; Tataru &

230 Bataillon, 2019). *polyDFE* fits an expression for the uSFS expected under a full DFE to data from

231 putatively neutral and selected classes of sites and estimates parameters by maximum likelihood.

232 For each set of positive selection parameters, simulated uSFS data were analysed under "Model B"

233 in *polyDFE* (a gamma distribution of deleterious mutational effects plus a discrete class of

234 advantageous mutations). Initial parameters for the maximisation were calculated from the data

235 using the '-e' option and the uSFS was analysed either with or without divergence using the "-w"

236 option in *polyDFE*. Analysing the uSFS without divergence causes the selection parameters to be

237 inferred from polymorphism data alone. For each replicate, I tested whether the inclusion of

238    beneficial mutations in the DFE improved model fit using likelihood ratio tests between the best-

239    fitting model and a model with $p_a$ set to 0.0. Setting $p_a$ = 0.0 means that positive selection does

240    not influence the likelihood, so two fewer parameters are being estimated. Twice the difference in

241    log-likelihood between the full DFE model and the model with $p_a$ = 0.0 was tested against a $\chi^2$

242    distribution with 2 degrees of freedom. Likelihood surfaces were estimated by running *polyDFE*

243    using a grid of fixed values for DFE parameters.

244

245
246
247    Data Availability

248    All code and *SLiM* configuration files needed to reproduce the results shown in this study are

249    available at https://github.com/TBooker/PositiveSelection_uSFS.

250

# Results

## Population genomic simulations

253    I performed simulations that modelled genes subject to mutation-selection-drift balance with

254    fitness effects drawn from a distribution that incorporated both deleterious and advantageous

255    mutations. The DFE for harmful mutations was constant, but I varied the fraction ($p_a$) and fitness

256    effects ($\gamma_a$) of beneficial mutations across simulated datasets (Table 2). For each set of

257    advantageous mutation parameters, 21Mbp of coding sequences was simulated, of which 14Mbp

258    were nonsynonymous and 7Mbp were synonymous sites. Variants present in the simulated

259    populations were used to construct the uSFS for a sample of 20 diploid individuals (Figure S1), a

260    sample size which is fairly typical of current population genomic datasets (e.g. Castellano et al.,

261    2019; Laenen et al., 2018; Williamson et al., 2014).

11

262

263

264

265

266

267

268

269

270

271 **Table 2** Parameters of positive selection assumed in simulations and the proportion of *polyDFE*
272 runs for which modelling positive selection gave a significantly better fit to the data.

| $\gamma_a$ | $p_a$ | $\gamma_a\, p_a$ | Proportion of likelihood ratio tests significant | | Proportion of analyses with gradient < 0.01 | |
|---|---|---|---|---|---|---|
| | | | *With divergence* | *Without divergence* | *With divergence* | *Without divergence* |
| 10 | | 0.001 | 0.02 | 0.07 | 0.11 | 0.71 |
| 50 | | 0.005 | 0.98 | 0.86 | 0.10 | 0.77 |
| 100 | 0.0001 | 0.01 | 0.98 | 0.02 | 0.03 | 0.58 |
| 500 | | 0.05 | 1.00 | 0.39 | 0.00 | 0.99 |
| 1,000 | | 0.10 | 1.00 | 1.00 | 0.00 | 0.71 |
| 10 | | 0.01 | 0.99 | 0.96 | 0.15 | 0.71 |
| 50 | | 0.05 | 1.00 | 1.00 | 0.06 | 0.98 |
| 100 | 0.001 | 0.10 | 1.00 | 1.00 | 0.00 | 0.97 |
| 500 | | 0.50 | 1.00 | 1.00 | 0.00 | 0.94 |
| 1,000 | | 1.00 | 1.00 | 1.00 | 0.00 | 0.71 |
| 10 | | 0.10 | 1.00 | 1.00 | 0.03 | 0.80 |
| 50 | | 0.50 | 1.00 | 1.00 | 0.02 | 0.99 |
| 100 | 0.01 | 1.00 | 1.00 | 1.00 | 0.02 | 0.95 |
| 500 | | 5.00 | 1.00 | 1.00 | 0.00 | 0.72 |
| 1,000 | | 10.0 | 1.00 | 1.00 | 0.00 | 0.41 |

273

274 Across simulations, the strength of selection acting on advantageous mutations ranged from $\gamma_a$ =

275 10 to $\gamma_a$ = 1,000. For a given $p_a$ parameter, increasing the strength of selection increased the

276 observed proportion of adaptive substitutions, $\alpha_{Obs}$ (Figure 1A). This is expected and is due to the

277 monotonic increasing relationship between fixation probability and the strength of positive

12

278    selection first described by Haldane (1927). Additionally, parameter combinations for which $\gamma_a p_a$

279    were equal had similar proportions of adaptive substitutions, for example compare $\gamma_a$ = 10 and $p_a$

280    = 0.01 to $\gamma_a$ = 1,000 and $p_a$ = 0.0001 (Figure 1A). This was also expected because the rate of

281    adaptive substitutions is proportional to $\gamma_a p_a$. In some datasets, particularly when $p_a$ = 0.01 and

282    advantageous mutations were very strongly selected (i.e. $\gamma_a \geq 500$), $\alpha_{Obs}$ exceeded 0.75, which is

283    higher than is typically estimated from empirical data (Galtier, 2016), so these parameter

284    combinations may not be biologically relevant.

285

286    The effects of selection at linked sites varied across simulated datasets. The DFE for deleterious

287    mutations was kept constant across simulations, so the extent of background selection should be

288    fairly similar across all parameter sets and thus variation in $\pi_S/\pi_0$ reflects the effects of selective

289    sweeps. Under neutrality $\pi_S/\pi_0$ had an expected value of 1.0 and I found that selection at linked

290    sites reduced nucleotide diversity below that expectation in all simulations (Figure 1B). Increasing

291    the fitness effects or frequency of advantageous mutations had a strong effect on genetic diversity

292    at synonymous sites, as shown by $\pi_S/\pi_0$ in Figure 1B. The highlighted points in Figure 1 indicate

293    parameter combinations for which $\gamma_a p_a$ = 0.01. As expected, $\alpha_{Obs}$ for these three parameter sets

294    was very similar (Figure 1A). Figure 1B shows that $\pi_S/\pi_0$ decreased across these three parameter

295    combinations as the strength of positive selection increased. Finally, differences in $p_a$ explained

296    most of the variation in the proportion of segregating advantageous mutations ($S_{Adv.}/S$) across

297    simulated datasets, but $S_{Adv.}/S$ also increased with the strength of positive selection (Figure 1C).

298    On the basis of these results, it is clear that there will be lower power to estimate positive

299    selection on the basis of standing variation when advantageous mutations are rare (i.e. $p_a$ =

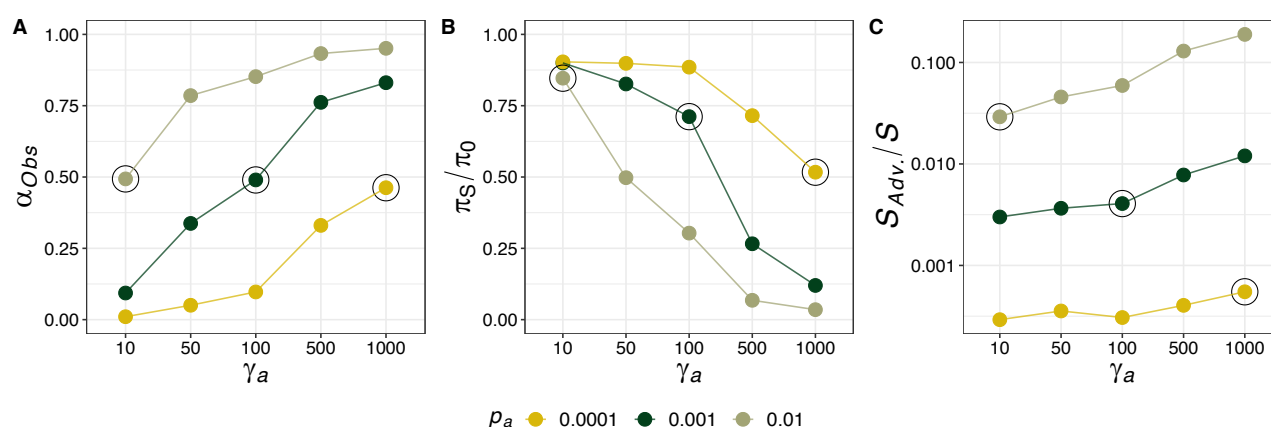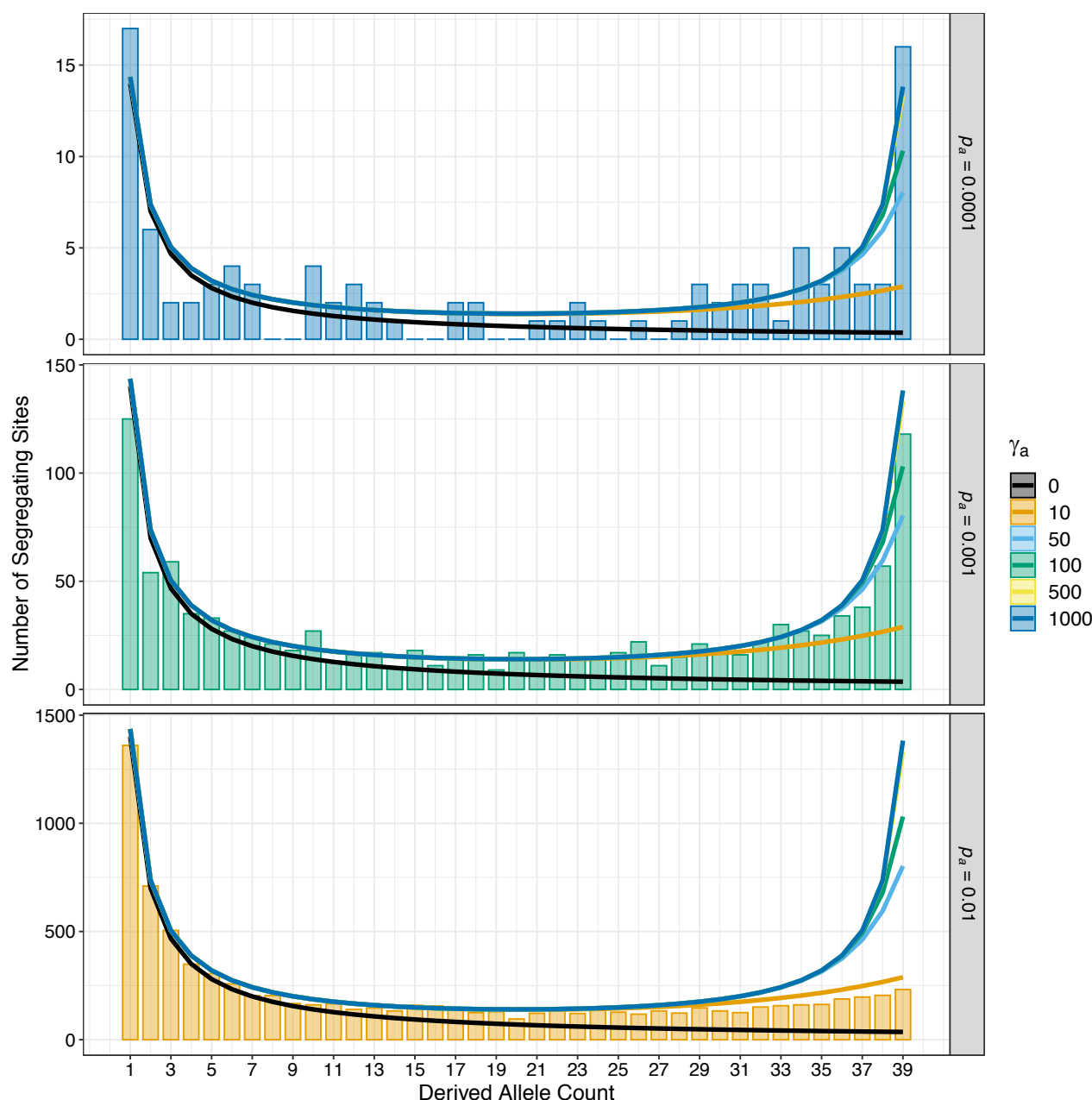300    0.0001) than when they are comparatively frequent (i.e. $p_a$ = 0.01).

**Figure 1** Population genetic summary statistics collated across all simulated genes. $\alpha_{Obs}$ is the observed proportion of substitutions fixed by positive selection. $\pi_S/\pi_0$ is genetic diversity relative to neutral expectation ($\pi_0 = 0.01$). $S_{Adv.}/S$ is the proportion of segregating nonsynonymous sites that are advantageous in the simulated datasets.

## Analysis of the unfolded site frequency spectrum

Figure 2 shows the observed (bars) and expected (lines) distribution of derived allele frequencies for beneficial mutations segregating in simulated populations. The three panels of Figure 2 correspond to three parameter combinations for which $\gamma_a p_a = 0.01$ ($\gamma_a = 1,000$ and $p_a = 0.0001$, $\gamma_a = 100$ and $p_a = 0.001$ and $\gamma_a = 10$ and $p_a = 0.01$). The lines in each of the panels of Figure 2 show the analytical expectation for the uSFS of advantageous mutations calculated using Equation 2 from Tataru et al., (2017). The analytical expectation closely matches the observed data for all three combinations (Figure 2). However, for a given value of $p_a$, the analytical expectation for models with increasing fitness effects were very similar, which likely makes it difficult to distinguish them on the basis of polymorphism alone (Figure 2). For the three parameter sets shown in Figure 2, the overall contribution that advantageous alleles make to the uSFS for nonsynonymous sites is small relative to deleterious ones (Figure S1). Accurate estimation of positive selection parameters from the uSFS requires that the distribution of advantageous alleles

14

321    can be distinguished from deleterious variants, so when $p_a$ is small it seems likely that uSFS

322    analyses will be unable to easily distinguish competing models.



323

**Figure 2** The uSFS for advantageous mutations under different combinations of positive selection parameters. The three bar charts show observed uSFS from simulations that model positive selection parameters that yield similar $\alpha$. The lines in each panel show the expected frequency spectra for different strengths of beneficial mutations and were obtained using Equation 2 from Tataru et al., (2017).

329

330    When analysing a particular uSFS dataset in *polyDFE*, I either modelled the full DFE (i.e. a gamma

331    distribution of deleterious mutations and a discrete class of advantageous mutational effects), or

15

332   just a gamma DFE for harmful mutations (dDFE). I compared the two models using likelihood ratio

333   tests, which tested the null hypothesis that the fit of the full DFE model is similar to that of a

334   model containing only deleterious mutations. For each of the combinations of positive selection

335   parameters shown in Table 2, I ran *polyDFE* on uSFS data from 100 bootstrap replicates. When

336   modelling the full uSFS (i.e. with divergence), *polyDFE* identified models containing positive

337   selection consistently for all but one ($p_a$ = 0.0001 and $\gamma_a$ = 10) of the parameter combinations

338   tested (Table 2). When the DFE was inferred from polymorphism data alone (i.e. without

339   divergence), models containing positive selection were identified less often, particularly when

340   beneficial mutations were rare ($p_a$ = 0.0001; Table 2). Table 2 also shows the proportion of analysis

341   runs for which the gradient of the likelihood exceeded 0.1. The *polyDFE* manual (Tataru &

342   Bataillon, 2019) suggests that gradients >0 indicate that the program has hailed to identify a

343   unique likelihood maximum. When the full uSFS was modelled, the gradient of the likelihood was

344   frequently >0, indicating that the model did not converge on a unique optimum. When modelling

345   the uSFS without divergence, *polyDFE* reported gradients <0.01 for a large proportion of replicate
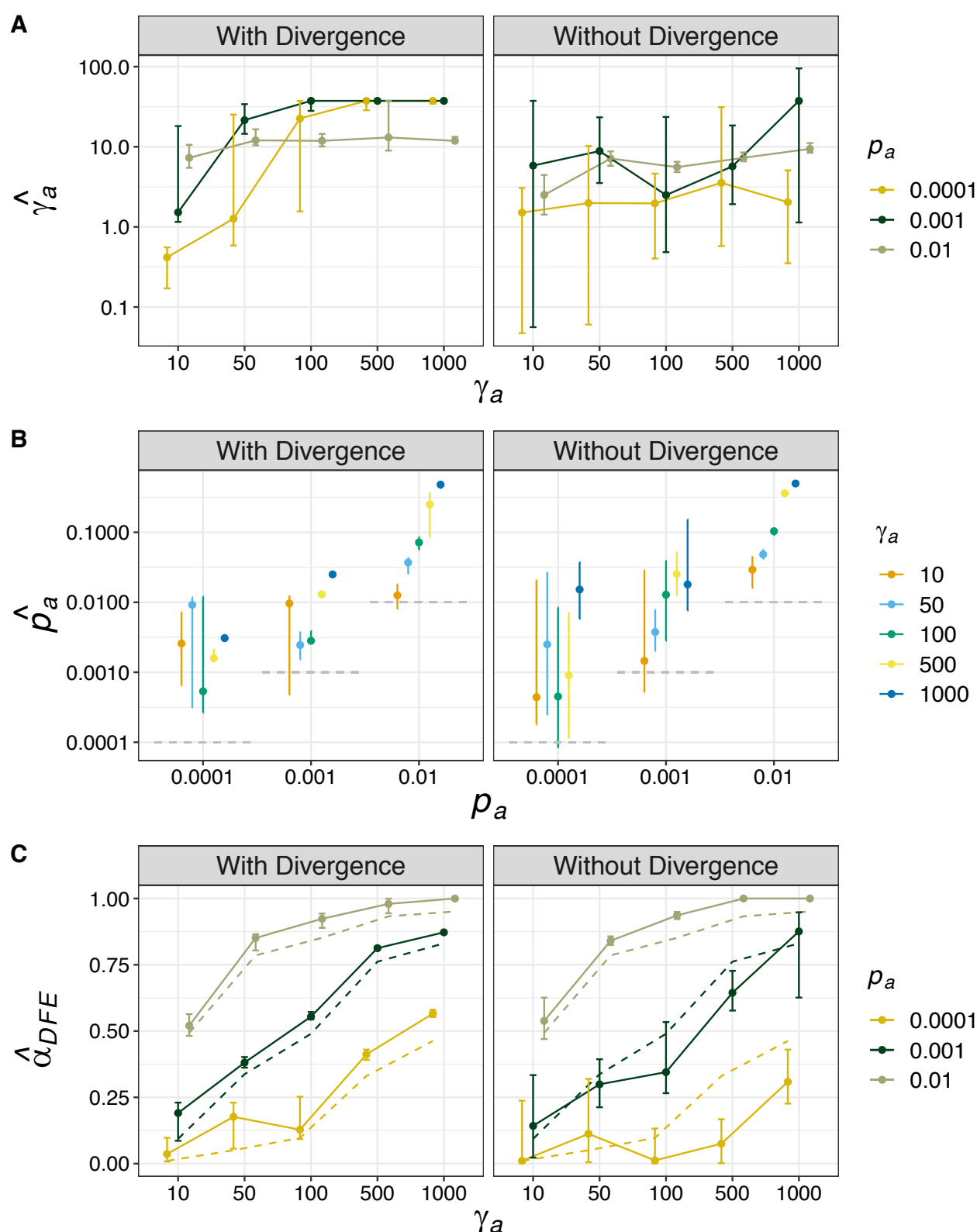
346   analyses (Table 2).

347

**Figure 3** Estimates of the parameters of advantageous mutations and the proportion of adaptive substitutions they imply from simulated datasets. A) $\gamma_a$ is the inferred selective effect of a new advantageous mutation; B) $p_a$ is the proportion of new mutations that are beneficial, the horizontal dashed grey lines indicate the simulated values in each case; C) $\alpha_{DFE}$ is the proportion of adaptive substitutions expected under the inferred DFE, the dashed lines indicate $\alpha_{Obs}$, the proportion of adaptive substitutions observed in the simulated datasets. Error bars indicate the 95% range of 100 bootstrap replicates.

357    Figures 3A and 3B show the parameters of positive selection estimated by analysis of uSFS from

358    simulated datasets. I found that when simulated beneficial mutations were mildly advantageous

359    ($\gamma_a$ = 10) but relatively frequent ($p_a$ = 0.01), both $\gamma_a$ and $p_a$ were estimated accurately regardless of

360    whether divergence was modelled or not (Figures 3A-B). This finding is consistent with both

361    Schneider et al., (2011) and Tataru et al., (2017). When $p_a$ = 0.01 and $\gamma_a$ > 10, the analysis of the

362    uSFS with or without divergence yielded very similar parameter estimates, but in both cases, the

363    strength of positive selection seemed to be positively correlated with the estimated $p_a$ (Figure 3).

364    In all cases, when beneficial mutations had $\gamma_a \geq 50$, neither $\gamma_a$ nor $p_a$ were accurately estimated

365    (Figure 3).

366

367    Tataru et al., (2017) pointed out that, if one had an estimate of the full DFE (i.e. with divergence),

368    the proportion of adaptive substitutions could be obtained by taking the ratio of the fixation

369    probability for a new beneficial mutation over the fixation probability for a random mutation

370    integrating over the full DFE (Equation 10; Tataru et al., 2017). The proportion of adaptive

371    substitutions obtained in this way is denoted $\alpha_{DFE}$. When modelling the full uSFS, $\alpha_{DFE}$ was

372    estimated with high accuracy, but with a slight upward bias (Figure 3C). When the DFE was

373    inferred without divergence $\alpha_{DFE}$ was underestimated when beneficial mutations were strongly

374    selected and rare (Figure 3).

375

376    In the presence of infrequent, strongly beneficial mutations the parameters of the DFE for

377    deleterious mutations estimated by *polyDFE* were very accurate (Figure S2). Estimates of the DFE

378    for harmful mutations were less accurate when beneficial mutations occurred with $p_a \geq 0.001$ and

379    $\gamma_a \geq 100$. This is presumably because in such cases recurrent selective sweeps eliminate a large

380    amount of neutral diversity and distort the distribution of standing genetic variation at

381    nonsynonymous sites. However, as stated above, the parameter range where the DFE for harmful

382    mutations was poorly estimated in this study may not be biologically relevant.
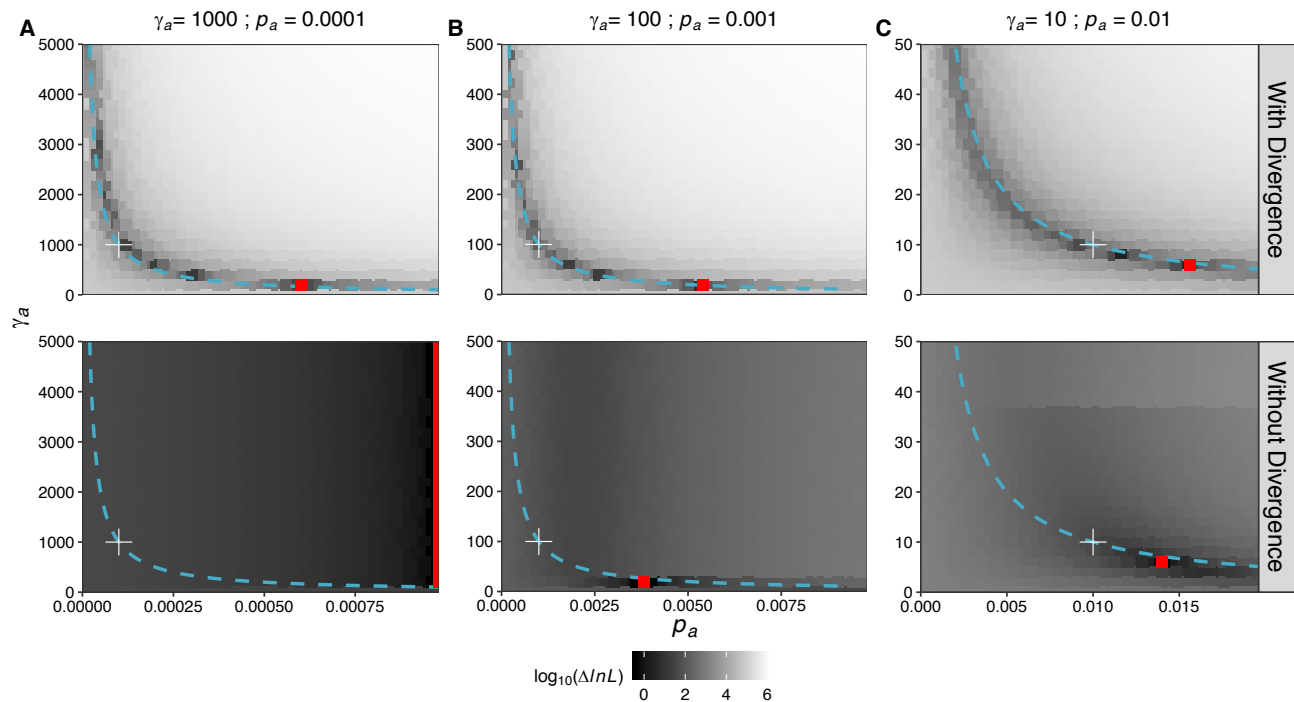
383

384    Model Identifiability

385



386
387    **Figure 4** The likelihood surface for the $\gamma_a$ and $p_a$ parameters for three simulated datasets. Hue
388    indicates differences in log likelihood between a particular parameter combination and the best-
389    fitting model. Best fitting models are indicated by red points and the true parameters are given
390    above the plots and indicated by the white plus signs on the likelihood surface. The relation $\gamma_a p_a$ =
391    0.1 is shown as a turquoise line and is constant across the three datasets shown.
392

393    It is very difficult to tease apart the parameters of positive selection from the uSFS by maximum

394    likelihood. Figure 4 shows the likelihood surface for the three sets of positive selection parameters

395    that satisfy the condition $\gamma_a p_a$ = 0.1. The proportion of adaptive substitutions is largely determined

396    by the product $\gamma_a p_a$ (Kimura & Ohta, 1971) and, as expected, the three parameter combinations

397    shown in Figure 4 all exhibit a similar $\alpha_{Obs}$ (Figure 1A). However, the extent by which neutral

398    genetic diversity is reduced and the number of segregating advantageous mutations differ

399    substantially across the three parameter combinations (Figure 1). The top row of panels in Figure 4

19

400 shows that when modelling the full uSFS, the likelihood surface closely tracks the relation $\gamma_a p_a$ =

401 0.1. Focussing on the top panel in Figure 4A, the maximum likelihood estimates (MLEs) of the

402 positive selection parameters (the red dot) are far from the true parameter values (indicated by

403 the plus sign), but the MLEs obtained satisfy $\gamma_a p_a$ = 0.1. The ridge in the likelihood surface

404 observed when modelling the full uSFS was described by both Schneider et al., (2011) and Tataru

405 et al., (2017). It comes about because between-species divergence carries information about $\alpha$,

406 and $\alpha$ is proportional to $\gamma_a p_a$.

407

408 Inferring the parameters of the DFE from polymorphism alone avoids the assumption of an

409 invariant DFE, but when doing so it may be difficult to distinguish competing models. Indeed,

410 across the three parameter combinations shown, values close to the truth were only obtained

411 from simulated data when $\gamma_a$ = 10 and $p_a$ = 0.01 (bottom panel Figure 4C). In the case of $\gamma_a$ = 1000

412 and $p_a$ = 0.0001, the likelihood surface about the true parameters was very flat (Figure 4A).

413 Increasing the $p_a$ parameter increased likelihood for all strengths of selection, so that the MLEs

414 shown in Figure 4A are simply the values with the highest $p_a$ in the range tested (the vertical red

415 line in Figure 4A). When $\gamma_a$ = 100 and $p_a$ = 0.001, the likelihood surface about the estimates was

416 steep, but the selection parameters identified by maximum likelihood were incorrect (Figure 4B).

417

# Discussion

419 In this study, I analysed simulated datasets modelling a range of positive selection parameter

420 combinations. I found that estimates of positive selection parameters obtained by analysis of the

421 uSFS were only accurate when beneficial mutations had $\gamma_a \leq 50$, under stronger selection the

422 individual parameters of positive selection were not accurately estimated (Figure 3). This is not

423 particularly surprising and is consistent with verbal arguments made in published studies (Booker

424 & Keightley, 2018; Campos et al., 2017). However, it is troubling that when beneficial mutations

425 are strongly selected and rare, the uSFS may often indicate a significant signal of positive

426 selection, but erroneous parameter estimates are obtained. If one were to analyse an empirical

427 dataset and estimate parameters of positive selection of the order $\gamma_a \sim 10$ and $p_a \sim 0.01$, it would

428 be difficult to know whether those were reflective of the true underlying parameters or an

429 artefact of strong selection.

430

431 On the basis of this study, it seems that researchers should treat parameters of positive selection

432 obtained by analysis of the uSFS with caution. The expected uSFS for advantageous mutations is

433 very similar when DFE models share the same $p_a$ parameter, and in such cases differing models

434 can only be distinguished by the density of high frequency derived variants (Figure 2). Polarization

435 error when estimating the uSFS can generate an excess in the number of high frequency variants

436 (Keightley & Jackson, 2018), so may generate a spurious signal of strong positive selection.

437 Analysis methods have been proposed which attempt to estimate the rate of polarisation error

438 when modelling the uSFS (Barton & Zeng, 2018; Tataru et al., 2017), but further study is required

439 to determine whether such methods reduce the signal of positive selection in uSFS-based

440 analyses. However, accounting for positive selection when analysing the uSFS yielded robust

441 estimates of the DFE for harmful mutations across the simulated datasets (Figure S2), although I

442 only examined a single DFE for harmful mutations in this study. Tataru et al., (2017) showed that

443 *polyDFE* accurately recovered the parameters of a range of DFE models if positive selection is

444 accounted for.

445

446 Estimates of $\alpha$ based on analysis of the uSFS may be biased when beneficial mutations are strongly

447 selected and infrequent. Calculating $\alpha$ using the rearranged MK-test makes the problematic

448 assumption that the DFE has remained invariant in the time since the focal species began to

21

449     diverge from the outgroup (Tataru et al., 2017). However, Tataru et al., (2017) pointed out that

450     one can avoid that assumption if $\alpha_{DFE}$ is calculated from a DFE estimated without divergence data.

451     In this study, estimates of $\alpha_{DFE}$ obtained when the full uSFS was analysed were very precise, but

452     with a slight upward bias (Figure 3). When simulated beneficial mutations were strongly selected

453     and rare, the parameters inferred using polymorphism data alone (i.e. without divergence) yielded

454     spurious estimates of $\alpha_{DFE}$ (Figure 3). When analysing datasets from real populations, $\alpha_{DFE}$ may not

455     capture the contribution that strongly beneficial mutations make to molecular evolution. This may

456     make it difficult to contrast $\alpha_{DFE}$ between species with large differences in $N_e$, because the number

457     of segregating advantageous mutations and thus ability to accurately estimate selection

458     parameters will depend on the population size.

459

460     The nature of the distribution of fitness effects for natural populations is largely unknown. In this

461     study, I analysed the uSFS data under the exact DFE model that had been simulated (i.e. a gamma

462     distribution of deleterious mutational effects plus a discrete class of beneficial effects). However,

463     when analysing empirical data, researchers have to make assumptions about the probability

464     distribution that best describes the DFE of the focal population. A gamma distribution is often

465     assumed for deleterious mutations as it is flexible and is described by only two parameters (Eyre-

466     Walker & Keightley, 2007). However, when analysing real data, one may bias their analyses by

467     strictly adhering to one particular family of probability distributions (Kousathanas & Keightley,

468     2013). In practice, model averaging provides a way to estimate key features of the DFE while

469     remaining agnostic to the exact shape that the distribution should take (Tataru & Bataillon, 2020).

470     However, if there is bias in the parameter estimates that are obtained across the models that one

471     tests, as is the case for strongly beneficial mutations, a biased average would result.

472

473   The simulations I performed in this study generated the ideal dataset for estimating parameters of

474   selection from the uSFS. I simulated 21Mbp of coding sites in which genotypes and whether sites

475   were selected or not was unambiguously known. When analysing real data this is not the case and

476   researchers often have to filter a large proportion of sites out of their analyses or choose to

477   analyse a subset of genes that have orthology with outgroups or other biological properties of

478   interest. Even with perfect knowledge, strongly beneficial mutations only represented a small

479   proportion of the standing genetic variation at nonsynonymous sites (Figure 1, S1). In addition, the

480   populations I simulated were randomly mating and had constant sizes over time. The results I

481   present in this study suggest that even with perfect knowledge of a population that adheres to the

482   assumptions of a Wright-Fisher model, it is inherently difficult to infer the parameters of strongly

483   beneficial mutations from the uSFS, particularly so when beneficial mutations occur infrequently.

484

485   ## Estimating parameters of positive selection from the uSFS versus

486   ## estimates from patterns of diversity

487   As discussed above, studies based on analysis of the uSFS and those based on selective sweep

488   models have yielded vastly different estimates of the parameters of positive selection. Patterns of

489   neutral genetic diversity in both humans and wild mice cannot be explained by the effects of

490   background selection alone, and in both species it has been suggested that strongly beneficial

491   mutations are required to explain the observed patterns (Booker & Keightley, 2018; Nam et al.,

492   2017). In the case of wild house mice, positive selection parameters obtained by analysis of the

493   uSFS do not explain dips in nucleotide diversity around functional elements (Booker & Keightley,

494   2018). Recently, Castellano et al. (2019) analysed the uSFS for nonsynonymous sites in great ape

495   species but did not find significant evidence for positive selection. In their dataset, Castellano et al.

496   (2019) had at least 8 haploid genome sequences for each of great ape species they analysed, and

497    they argued that they were underpowered to detect positive selection on the basis of the uSFS. In

498    this study, I analysed datasets of 20 diploid individuals and found that it was very difficult to

499    accurately capture positive selection parameters. Increasing the number of sampled individuals

500    even further may increase the power to estimate the strength of positive selection, but this study

501    suggests that the increase in power will depend on the underlying DFE. When $p_a$ is small, the

502    expected number of advantageous mutations present in the uSFS for 200 diploids is less than 10

503    for most frequency classes when14Mbp of nonsynonymous sites have been used to construct the

504    uSFS (Figure S3). Indeed, Figure S3 shows that even with very large sample sizes, the expected

505    uSFS for beneficial mutations are very similar and may only be distinguished on the basis of a small

506    number of high frequency derived alleles. Thus, it may be that the uSFS is inherently limited in the

507    information it carries on the DFE for beneficial mutations so other sources of information may

508    have to be used to accurately recover parameters.

509

510    In this study, I modelled beneficial mutations using a discrete class of selection coefficients when,

511    in reality, there is likely a continuous distribution of fitness effects. Indeed, studies in both humans

512    and *D. melanogaster* have found evidence for a bimodal distribution containing both strongly and

513    weakly beneficial mutations contributing to adaptive evolution using methods which incorporate

514    linkage information but do not explicitly estimate selection parameters (Elyashiv et al., 2016;

515    Uricchio et al., 2019). There are currently no methods that estimate the DFE using an analytical

516    expression for the uSFS expected under the combined effects of BGS and sweeps. Rather,

517    nuisance parameters or demographic models are used to correct for the contribution that

518    selection at linked sites may make to the shape of the SFS (Eyre-Walker, Woolfit, & Phelps, 2006;

519    Galtier, 2016; Tataru et al., 2017). However, as this study shows, the parameters of positive

520    selection are not reliably estimated when analysing the uSFS alone. A way forward may be in using

521    computational approaches to make use of all of the available data, while not necessitating an

24

522  expression for the uSFS expected under the combined effects of BGS, sweeps, population size

523  change and direct selection. An advance in this direction has recently been made by Uricchio et al.,

524  (2019) who developed an ABC method for estimating $\alpha$ which makes use of the distortions to the

525  uSFS generated by BGS and sweeps. By applying their method to data from humans, Uricchio et

526  al., (2019) found that $\alpha$ = 0.13 for nonsynonymous sites, 72% of which was generated by mildly

527  beneficial mutations and 28% by strongly beneficial mutations. However, the computational

528  approach developed by Uricchio et al., (2019) could readily be extended to model an arbitrarily

529  complex DFE for beneficial mutations. Their methods could be implemented in a machine-learning

530  context, with training data generated by forward-simulations that capture confounding factors

531  such as population structure and population size change as well as the effects of selection at

532  linked sites.

# Acknowledgements

534  I wish to extend gratitude to Peter Keightley, Michael Whitlock and Sam Yeaman for valuable

535  advice and mentorship. Thanks to Thomas Bataillon, Sam Yeaman and Peter Keightley for

536  comments on previous versions of the manuscript. Thanks to Paula Tataru and Thomas Bataillon

537  for help with *polyDFE*. Thanks to two anonymous reviewers for constructive feedback.

# References

539  Bailey, S. F., & Bataillon, T. (2016). Can the experimental evolution programme help us elucidate

540  the genetic basis of adaptation in nature? *Mol Ecol*, *25*(1), 203–218.

541  Bank, C., Hietpas, R. T., Wong, A., Bolon, D. N., & Jensen, J. D. (2014). A bayesian MCMC approach

542  to assess the complete distribution of fitness effects of new mutations: uncovering the

543  potential for adaptive walks in challenging environments. *Genetics*, *196*(3), 841–852.

544  Barton, H. J., & Zeng, K. (2018). New methods for inferring the distribution of fitness effects for

545      INDELs and SNPs. *Molecular Biology and Evolution*, *35*(6), 1536–1546.

546      Barton, N. H. (2000). Genetic hitchhiking. *Philos Trans R Soc Lond B Biol Sci*, *355*(1403), 1553–1562.

547      Böndel, K. B., Kraemer, S. A., Samuels, T., McClean, D., Lachapelle, J., Ness, R. W., … Keightley, P. D.

548      (2019). Inferring the distribution of fitness effects of spontaneous mutations in

549      *Chlamydomonas reinhardtii*. *PLOS Biology*, *17*(6), e3000192.

550      Booker, T. R., & Keightley, P. D. (2018). Understanding the factors that shape patterns of

551      nucleotide diversity in the house mouse genome. *Mol. Biol. Evol.*, *35*(12), 2971–2988.

552      Boyko, A. R., Williamson, S. H., Indap, A. R., Degenhardt, J. D., Hernandez, R. D., Lohmueller, K. E.,

553      … Bustamante, C. D. (2008). Assessing the evolutionary impact of amino acid mutations in the

554      human genome. *PLoS Genet*, *4*(5), e1000083.

555      Campos, J L, Zhao, L., & Charlesworth, B. (2017). Estimating the parameters of background

556      selection and selective sweeps in *Drosophila* in the presence of gene conversion. *Proc Natl*

557      *Acad Sci*, *114*(24), E4762–E4771.

558      Campos, José Luis, & Charlesworth, B. (2019). The effects on neutral variability of recurrent

559      selective sweeps and background selection. *Genetics*, *212*(1), 287–303.

560      Castellano, D., Macià, M. C., Tataru, P., Bataillon, T., & Munch, K. (2019). Comparison of the full

561      distribution of fitness effects of new amino acid mutations across great apes. *Genetics*,

562      genetics. 302494.2019.

563      Charlesworth, B. (1994). The effect of background selection against deleterious mutations on

564      weakly selected, linked variants. *Genetical Research*, *63*(03), 213.

565      Elyashiv, E., Sattath, S., Hu, T. T., Strutsovsky, A., McVicker, G., Andolfatto, P., … Sella, G. (2016). A

566      genomic map of the effects of linked selection in *Drosophila*. *PLoS Genet*, *12*(8), e1006130.

567      Eyre-Walker, A., & Keightley, P. D. (2007). The distribution of fitness effects of new mutations. *Nat*

568      *Rev Genet*, *8*(8), 610–618.

569      Eyre-Walker, A., & Keightley, P. D. (2009). Estimating the rate of adaptive molecular evolution in

570 the presence of slightly deleterious mutations and population size change. *Mol Biol Evol*,

571 *26*(9), 2097–2108.

572 Eyre-Walker, A., Woolfit, M., & Phelps, T. (2006). The distribution of fitness effects of new

573 deleterious amino acid mutations in humans. *Genetics*, *173*(2), 891–900.

574 Galtier, N. (2016). Adaptive protein evolution in animals and the effective population size

575 Hypothesis. *PLoS Genet.*, *12*(1), e1005774.

576 Haldane, J. B. S. (1927). A Mathematical Theory of Natural and Artificial Selection, Part V: Selection

577 and Mutation. *Mathematical Proceedings of the Cambridge Philosophical Society*, *23*(7), 838–

578 844.

579 Haller, B. C., & Messer, P. W. (2019). SLiM 3: Forward genetic simulations beyond the Wright-

580 Fisher model. *Mol. Biol. Evol.*, *36*(3), 632–637.

581 Hill, W. G. (2010). Understanding and using quantitative genetic variation. *Philosophical*

582 *Transactions of the Royal Society B: Biological Sciences*, Vol. 365, pp. 73–85.

583 Keightley, P D, Campos, J. L., Booker, T. R., & Charlesworth, B. (2016). Inferring the frequency

584 spectrum of derived variants to quantify adaptive molecular evolution in protein-coding

585 genes of *Drosophila melanogaster*. *Genetics*, *203*(2), 975–984.

586 Keightley, P D, & Eyre-Walker, A. (2007). Joint inference of the distribution of fitness effects of

587 deleterious mutations and population demography based on nucleotide polymorphism

588 frequencies. *Genetics*, *177*(4), 2251–2261.

589 Keightley, Peter D, & Jackson, B. C. (2018). Inferring the probability of the derived vs. the ancestral

590 allelic state at a polymorphic site. *Genetics*, *209*(3), 897–906.

591 Kimura, M., & Ohta, T. (1971). *Theoretical aspects of population genetics*. Princeton Univ. Press.

592 Kousathanas, A., & Keightley, P. D. (2013). A comparison of models to infer the distribution of

593 fitness effects of new mutations. *Genetics*, *193*(4), 1197–1208.

594 Laenen, B., Tedder, A., Nowak, M. D., Toräng, P., Wunder, J., Wötzel, S., … Slotte, T. (2018).

27

595   Demography and mating system shape the genome-wide impact of purifying selection in

596   Arabis alpina. *Proceedings of the National Academy of Sciences of the United States of*

597   *America*, *115*(4), 816–821.

598   Loewe, L., & Charlesworth, B. (2006). Inferring the distribution of mutational effects on fitness in

599   *Drosophila*. *Biol Lett*, *2*(3), 426–430.

600   McDonald, J. M., & Kreitman, M. (1991). Adaptive protein evolution at the *Adh* locus in *Drosophila*.

601   *Nature*, *351*.

602   Nam, K., Munch, K., Mailund, T., Nater, A., Greminger, M. P., Krützen, M., … Schierup, M. H.

603   (2017). Evidence that the rate of strong selective sweeps increases with population size in the

604   great apes. *Proceedings of the National Academy of Sciences of the United States of America*,

605   *114*(7), 1613–1618.

606   Orr, H. A. (2003). The distribution of fitness effects among beneficial mutations. *Genetics*, *163*(4),

607   1519–1526.

608   Orr, H. A., & Unckless, R. L. (2014). The population genetics of evolutionary rescue. *PLoS Genetics*,

609   *10*(8).

610   Otto, S. P. (2009). The evolutionary enigma of sex. *American Naturalist*, *174*(SUPPL. 1).

611   Schneider, A., Charlesworth, B., Eyre-Walker, A., & Keightley, P. D. (2011). A method for inferring

612   the rate of occurrence and fitness effects of advantageous mutations. *Genetics*, *189*(4), 1427–

613   1437.

614   Tataru, P., & Bataillon, T. (2019). polyDFEv2.0: testing for invariance of the distribution of fitness

615   effects within and across species. *Bioinformatics*, *35*(16), 2868–2869.

616   Tataru, P., & Bataillon, T. (2020). polyDFE: Inferring the distribution of fitness effects and

617   properties of beneficial mutations from polymorphism data. In *Methods in Molecular Biology*

618   (Vol. 2090, pp. 125–146).

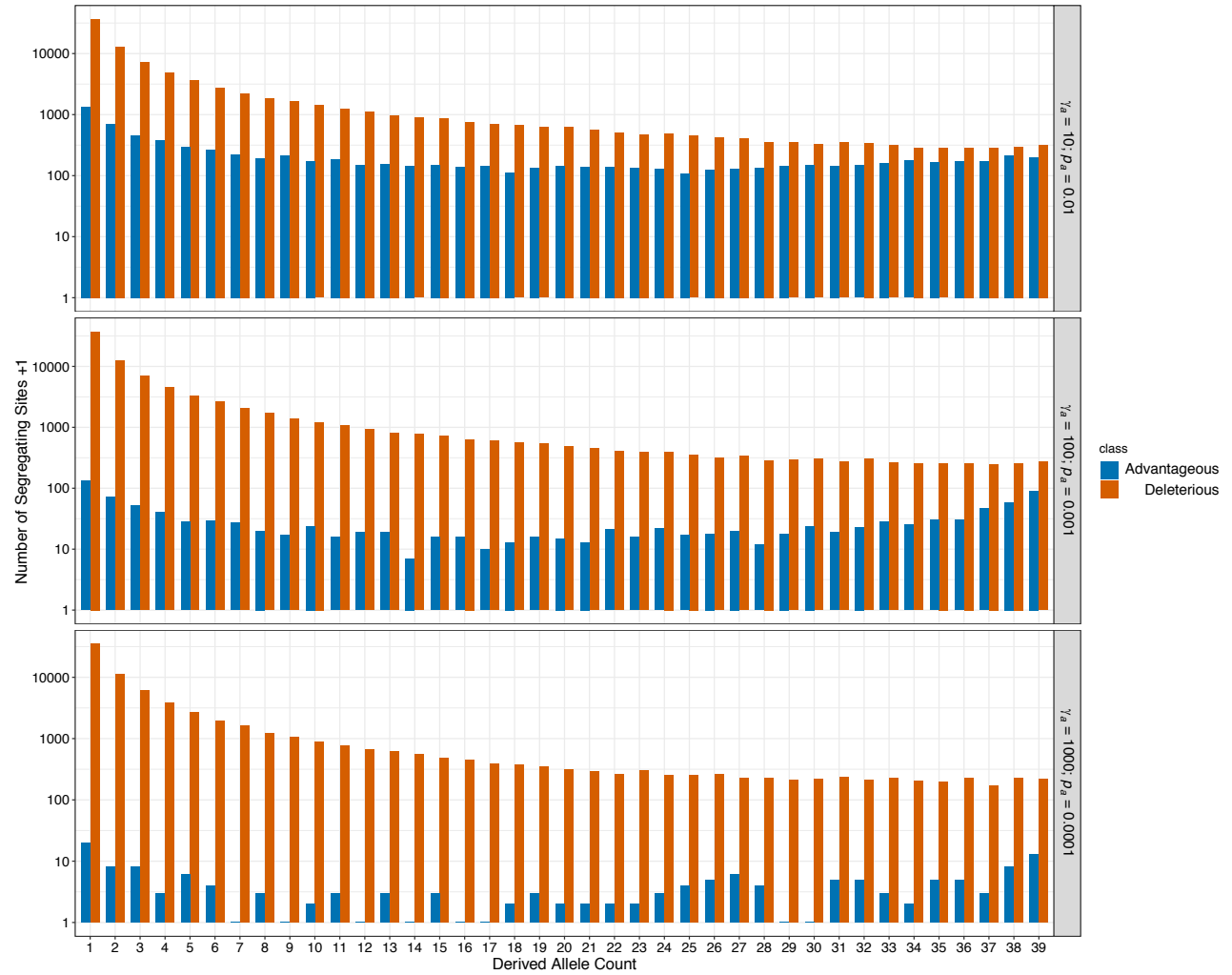619   Tataru, P., Mollion, M., Glemin, S., & Bataillon, T. (2017). Inference of distribution of fitness effects

620    and proportion of adaptive substitutions from polymorphism data. *Genetics*, *207*(3), 1103–

621    1119.

622  Uricchio, L. H., Petrov, D. A., & Enard, D. (2019). Exploiting selection at linked sites to infer the rate

623    and strength of adaptation. *Nat Ecol Evol*, *3*(6), 977–984.

624  Williamson, R. J., Josephs, E. B., Platts, A. E., Hazzouri, K. M., Haudry, A., Blanchette, M., & Wright,

625    S. I. (2014). Evidence for widespread positive and negative selection in coding and conserved

626    noncoding regions of *Capsella grandiflora*. *PLoS Genetics*, *10*(9), e1004622.

627  Wright, S. (1937). The distribution of gene frequencies in populations. *Proceedings of the National*

628    *Academy of Sciences*, *23*(6), 307–320.

629

630

631

632

633

634

635

636

637

638

639

640

# Supplementary Material



**Figure S1** The observed uSFS for nonsynonymous sites for three sets of positive selection parameters. The distribution of deleterious mutations is shown in orange and the distribution of advantageous mutations is shown in blue. For the purposes of visualising the data on a log scale, the number of segregating sites is shown +1.
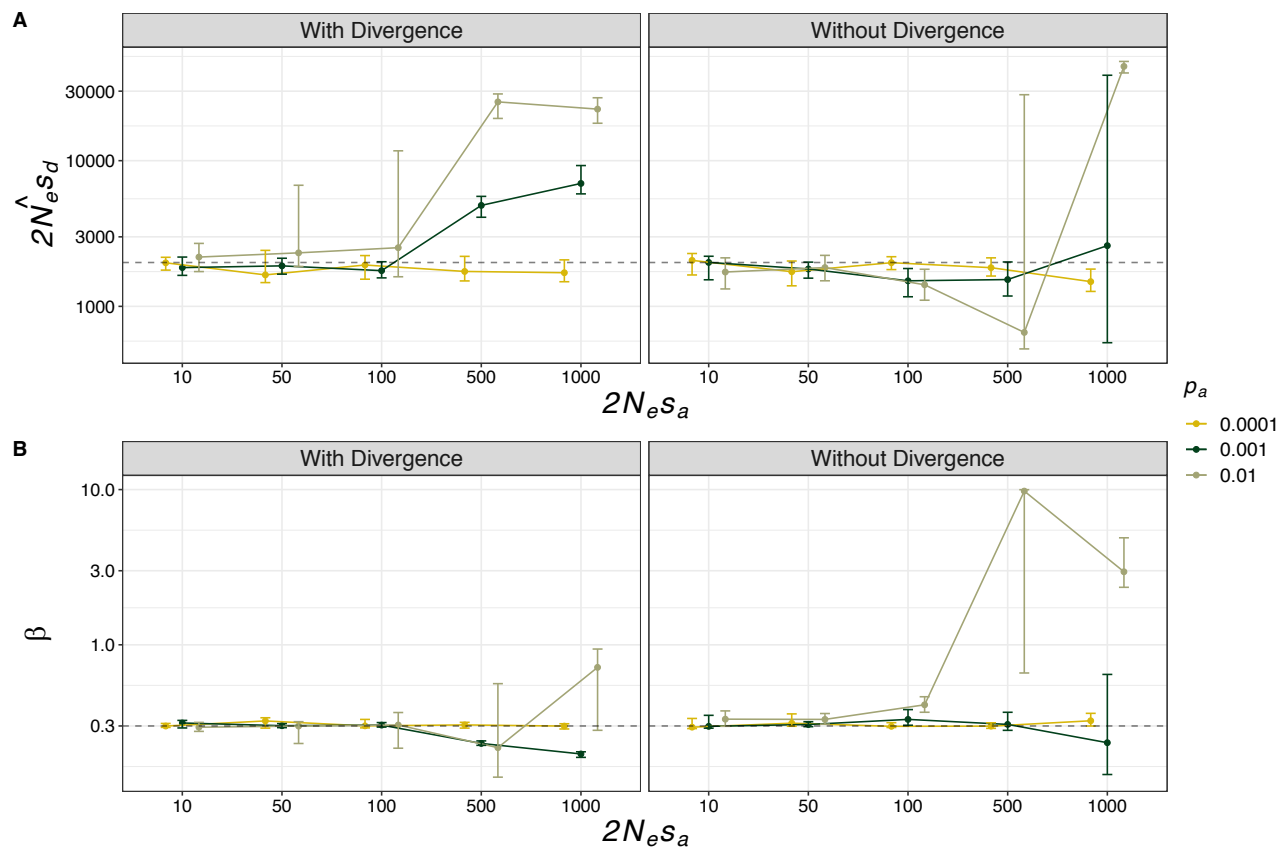
647
648
649 **Figure S2** Parameter estimates for the DFE for deleterious mutations obtained from simulated
650 datasets. A) the mean effect of a deleterious mutation and b) the shape parameter of the gamma
651 distribution. Error bars indicate the 95% range of 100 bootstrap replicates.
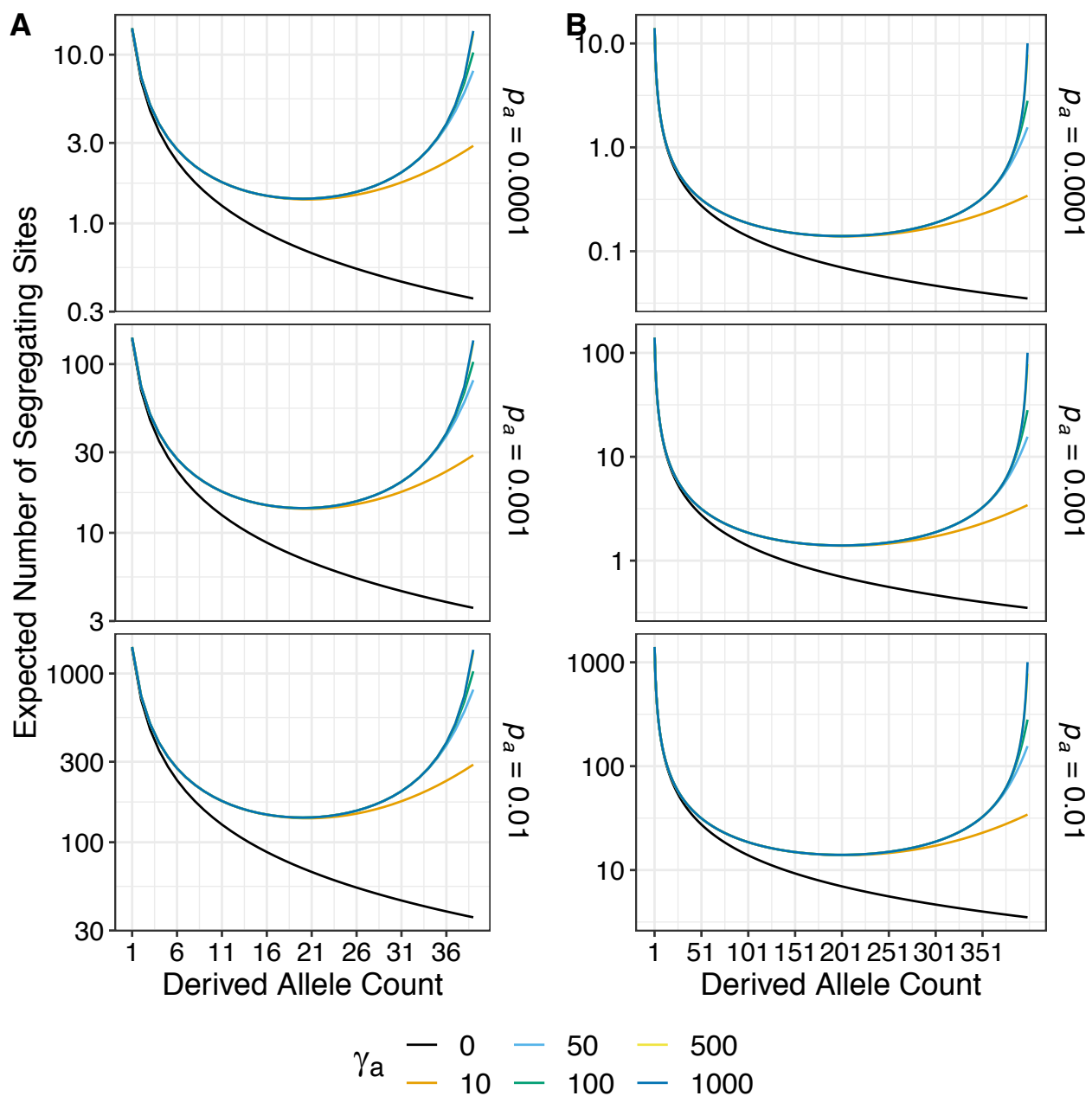652

**Figure S3** The expected uSFS for beneficial alleles. Panel A shows the expected uSFS for a sample of 20 diploid individuals, and panel B shows the uSFS for 200 diploid individuals.