1 **From drift to draft: How much do beneficial mutations actually contribute to**

2 **predictions of Ohta's slightly deleterious model of molecular evolution?**

3

4

5 **Jun Chen**[*,†] **Sylvain Glémin**[*,§]**, Martin Lascoux**[*,‡]

6 [*]Program in Plant Ecology and Evolution, Department of Ecology and Genetics,

7 Evolutionary Biology Centre, Uppsala University, 75236 Uppsala, Sweden

8 § Université de Rennes, CNRS, ECOBIO [Ecosystèmes, Biodiversité, Evolution] -

9 UMR 6553, F-35000 Rennes, France

10

11 †Present address: College of Life Sciences, Zhejiang University, Hangzhou, Zhejiang

12 310058, China

13

14

15 ‡Author for correspondence: Martin.Lascoux@ebc.uu.se

**Abstract**

Since its inception in 1973 the slightly deleterious model of molecular evolution, aka the Nearly Neutral Theory of molecular evolution, remains a central model to explain the main patterns of DNA polymorphism in natural populations. This is not to say that the quantitative fit to data is perfect. In a recent study CASTELLANO *et al.* (2018) used polymorphism data from *D. melanogaster* to test whether, as predicted by the Nearly Neutral Theory, the proportion of effectively neutral mutations depends on the effective population size ($N_e$). They showed that a nearly neutral model simply scaling with $N_e$ variation across the genome could not explain alone the data but that consideration of linked positive selection improves the fit between observations and predictions. In the present article we extended their work in two main directions. First, we confirmed the observed pattern on a set of 59 species, including high quality genomic data from 11 animal and plant species with different mating systems and effective population sizes, hence *a priori* different levels of linked selection. Second, for the 11 species with high quality genomic data we also estimated the full Distribution of Fitness Effects (DFE) of mutations, and not solely the DFE of deleterious mutations. Both $N_e$ and beneficial mutations contributed to the relationship between the proportion of effectively neutral mutations and local $N_e$ across the genome. In conclusion, the predictions of the slightly deleterious model of molecular evolution hold well for species with small $N_e$. But for species with large $N_e$ the fit is improved by incorporating linked positive selection to the model.

**Keywords**: Nearly Neutral Theory, Distribution of Fitness Effects, beneficial mutations, linked selection

**Introduction**


The year 2018 saw the celebration of the 50[th] anniversary of the Neutral Theory of molecular evolution (called simply the Neutral Theory thereafter). At 50 years of age, the Neutral Theory is still shrouded in controversies, some pronouncing it dead and overwhelmingly rejected by facts (Kern and Hahn 2018) while others see it as very much alive and kicking (Nei *et al.* 2010; Jensen *et al.* 2019). As a quick glance at major textbooks in population genetics and at the literature would suggest, it seems fair to say that the Neutral Theory is certainly not totally dead. Even if it undoubtedly did lose some of its initial appeal it continues to play a central role in population genetics, a position well summarized by Kreitman (1996) in his spirited essay "The neutral theory is dead. Long live the Neutral Theory". Shortcomings of the Neutral Theory were already noted in the 1970s and the Neutral Theory has itself evolved. Indeed, its inadequacy to fully explain the data, in particular the constancy of the molecular clock, was already noted in 1973, leading Tomoko Ohta (1973) to propose the Nearly Neutral Theory of molecular evolution. In contrast to the Neutral Theory where most mutations are assumed to be neutral or strongly deleterious, the Nearly Neutral Theory assigns much more prominence to the contribution to standing polymorphism of mutations that are weakly selected and effectively neutral (Ohta 1992; Ohta and Gillespie 1996). Weakly selected mutations can be slightly deleterious or slightly beneficial, but as noted by Kreitman (1996) the best developed of the weak selection models primarily considers slightly deleterious mutations and was therefore christened by him "the slightly deleterious model". This is the model that we will be testing in most of the present paper.


Like the Neutral Theory, however, the Nearly Neutral Theory still assumes that "only a minute fraction of DNA changes in evolution are adaptive in nature" (Kimura 1983). Under this view, polymorphism is thought to be mostly unaffected by positive selection, except around the few recently selected beneficial alleles (selective sweeps). This was already at variance with the view put forward by Gillespie (e.g. Gillespie 2004) that assigned a greater role to linked positive selection in shaping polymorphism (see also CORBETT-DETIG *et al.*

3

75   2015) and is in even stronger contrast with the claim by Kern and Hahn (2018)

76   that "natural selection has played the predominant role in shaping within- and

77   between-species genetic variation" and that "the ubiquity of adaptive variation

78   both within and between species" leads to the rejection of the universality of the

79   Neutral Theory. In a far more nuanced assessment of the Neutral Theory and its

80   contribution, Jensen *et al.* (2018) argued that the effects of linked selection could

81   readily be incorporated in the Nearly Neutral framework. The heart of the

82   dispute, either today or in the early days of the Nearly Neutral Theory, is about

83   the degree to which each category of mutations contributed directly and

84   indirectly to genetic variation within- and between-species.

85

86   A core prediction of the Nearly Neutral Theory is that the fraction of mutations

87   affected by selection depends on $N_e$ (Ohta 1973). $N_e$ can vary among species but

88   also within a genome because of linked selection (reviewed in Ellegren and

89   Galtier 2016). The effect of selection against deleterious mutations on linked

90   neutral variants – background selection (Charlesworth *et al.* 1993) – is often

91   modeled by a simple re-scaling of $N_e$ but except in specific situations  effects of

92   linked selection are more complex and there is not a single re-scaling (Barton

93   1995; Zeng 2013; Comeron 2017; Cvijovic et al. 2018; Torres et al, 2019). In the

94   case of beneficial mutations, for instance, the interference depends both on the

95   beneficial effect of the sweeping mutation and on selection acting at linked sites

96   (Barton 1995; Weissman and Barton 2012).

97

98   Evidence that linked positive selection and not only direct selection on slightly

99   deleterious and beneficial mutations contributed to the relationship between the

100   fraction of mutations affected by selection and $N_e$ has recently been obtained by

101   Castellano *et al.* (2018). Using two *Drosophila melanogaster* genome re-

102   sequencing datasets, Castellano *et al.* (2018) tested a prediction of the slightly

103   deleterious model first obtained by Kimura (1979) and then extended by Welch *et*

104   *al.* (2008). Welch *et al.* (2008) showed that if one considers only deleterious

105   mutations, the logarithm of the ratio of nucleotide diversity at non-synonymous

106   and synonymous amino acid changes is linearly related to the logarithm of the

107   effective population size and that the slope of this log-log regression line is equal

108    to the shape parameter of the Distribution of Fitness Effects (DFE), $\beta$, if the DFE

109    of deleterious mutations is modeled by a Gamma distribution:

110

111    $ln(\pi_N/\pi_S) \approx -\beta \, ln(N_e) + constant$      [Eq. 1a]

112

113    where $\pi_N$ is the nucleotide diversity at non-synonymous sites and $\pi_S$ is the

114    nucleotide diversity at synonymous sites.

115

116    Or, rewriting this expectation by using $\pi_S$ as a proxy for $N_e$:

117

118    $ln(\pi_N/\pi_S) \approx -\beta \, ln(\pi_S) + constant'$      [Eq. 1b]

119

120    The second equation holds only if variation in $\pi_S$ solely depends on variation in

121    $N_e$, and that there is no correlation between the mutation rate and $N_e$. It should

122    also be pointed out that the DFE used here only considers deleterious mutations,

123    as estimated for instance by DFE-alpha (Eyre-Walker and Keightley 2009). A

124    direct test of this prediction using among-species comparison can be problematic

125    if mutation rates cannot be controlled for. To circumvent this problem,

126    Castellano *et al.* (2018) used within genome variation in $N_e$, under the reasonable

127    assumption that variation in mutation rates are negligible compared to variation

128    in $N_e$ across a genome. They found (see also James et al. 2017) that the slope was

129    significantly steeper than expected under a simple scaling of $N_e$ and simulations

130    indicated that linked positive selection, but not background selection, could

131    explain this discrepancy. The effect of linked selection on the relationship

132    between $\pi_N/\pi_S$ and $\pi_S$ is twofold. First it increases stochasticity in allele

133    frequencies, or, in other words, decreases the local effective population size.

134    Second, linked selection leads to non-equilibrium dynamics.  Genetic diversity

135    will recover faster for deleterious than neutral mutations, altering the

136    relationship between $\pi_N/\pi_S$ and $\pi_S$ (Brandvain and Wright, 2016; Do et al. 2015;

137    Gordo and Dionisio 2005; Vigué and Eyre-Walker 2019). More precisely, the

138    more a region is affected by selective sweeps, the lower $\pi_S$ is and the higher

139    $\pi_N/\pi_S$ is compared to the equilibrium expectation: this effect makes the slope

140    steeper compared to the equilibrium expectation.

141

142 In the present paper, we first confirmed the observed pattern on the set of 59

143 species used in Chen *et al.* (2017). We then used 11 high quality genomic

144 datasets for which an outgroup is available to test whether the results obtained

145 by Castellano *et al.* (2018) hold more generally and, in particular, in species with

146 much smaller effective sizes than *D. melanogaster*, and with different levels of

147 linkage disequilibrium. While we adopted the same general approach than

148 Castellano *et al.* (2018), our analysis differed from theirs in one important

149 respect. In their study, Castellano *et al.* (2018) only characterized the DFE of

150 deleterious mutations. We, instead, used a newly developed approach, *polyDFE*

151 (Tataru *et al.* 2017), that also considers positive mutations, which is expected to

152 improve the estimation of the shape of the DFE of deleterious mutations and to

153 disentangle the direct effects of both positive and negative selection.

154

155 **Material & Methods**

156

157 *Genomic data and regression of $\pi_N/\pi_S$ over $\pi_S$*

158

159 In a first step we re-analyzed the 59 species from Chen *et al.* (2017), which

160 included 34 animals and 28 plant species. We estimated the DFE using folded site

161 frequency spectra with the same method as in Chen *et al.* (2017) and calculated

162 the slope (regression coefficient of $\log(\pi_N/\pi_S)$ over $\log(\pi_S)$ as described in the

163 next paragraph. For DFE estimation using folded SFS the model assumes a

164 gamma distribution for deleterious mutations and takes demography (or

165 sampling or any departure from equilibrium) into account by introducing *n-1*

166 nuisance parameters for an SFS of size *n* (the corresponding code was provided

167 in Chen *et al.* (2017)). In later analyses that required unfolded site frequency

168 spectra, we retained 11 species with high quality genomic datasets and with an

169 available outgroup. These eleven species are given in Table 1. They include both

170 animal and plant species with contrasted levels of nucleotide polymorphism and

171 mating systems. For each of the eleven species, we aligned short reads to the

172 genome using BWA-mem (Li and Durbin 2010) and sorted the alignment using

173 SAMtools. PCR duplicates were removed and INDELs were realigned using GATK

174   toolkit (McKenna *et al.* 2010). HaplotypeCaller was used for individual genotype
175   identification and joint SNP calling was performed across all samples using
176   GenotypeGVCFs. Variant and invariant sites were kept only if genotypes of all
177   individuals were successfully identified (Carson *et al.* 2014). We collected Single
178   Nucleotide Polymorphism (SNPs) in all CDS regions and calculated genetic
179   diversity of 4-fold and 0-fold sites as proxies for polymorphism at synonymous
180   ($\pi_S$) and non-synonymous sites ($\pi_N$). Sites were all masked with 'N' and excluded
181   from further computation in the following four cases: heterozygous sites in
182   selfing species, sites with more than two variants, variants at sites within five
183   bases of a flanking INDEL, and missing individuals. We applied the same SNP
184   sampling strategy as in James et al. (2017) and Castellano *et al.* (2018) in order
185   to remove potential dependency between estimates of $\pi_N/\pi_S$ and $\pi_S$. In brief, we
186   first split all synonymous SNPs into three groups (S1, S2, and S3) using a
187   hypergeometric sampling based on the total number of synonymous sites. To bin
188   genes and reduce the difference in number of SNPs in each bin, we ranked genes
189   according to their Watterson's estimate of nucleotide diversity ($\theta_{S1}$) and grouped
190   these ranked genes into 20 bins each representing approximately 1/20 of the
191   total number of synonymous SNPs. We then used $\pi_{S2}$ to estimate the $\pi_N/\pi_S$
192   (mean $\pi_N$ divided by mean $\pi_s$ in each bin) ratio and $\pi_{S3}$ as an independent
193   estimate of the genetic diversity of each bin.

194

195   We calculated the slope of the linear regression (*l*) of the log-transformed value
196   of the $\pi_N/\pi_S$ ratio on the log-transformed value of $\pi_S$, using the "lm" function in R
197   (R Core Team 2018). In pilot runs on 59 species (population data of Chen *et al.*
198   (2017)), the estimates of *l* showed extensive variation depending on, among
199   other things, the qualities of genome sequencing, read depth, annotation and SNP
200   calling. Thus, we selected 11 species for which a high-quality genome sequence
201   and an outgroup were available. Individuals were selected from the same genetic
202   background, i.e. admixture or population structure were carefully removed. At
203   least 20 alleles (i.e. 10 individuals for outcrossing species or 20 for selfing
204   species) were retained from a single ancestral cluster defined in
205   Admixture/Structure analysis in the original publication. For the two *Capsella*
206   species, we performed Admixture analysis for both species separately.  A series

207    of quality controls for $l$ calculation were performed as described in the following.

208    The longest transcript for each gene model was kept only if it contained both

209    start and stop codons (putative full length) and no premature stop codons. SNPs

210    flanking five bases of INDEL were masked to avoid false positive calls. A grid of

211    filtering criteria (see details in Table S2) was also implemented on each species

212    based on sequence similarity against Swiss-Prot database (e-value, bit-score,

213    query coverage) and sequencing quality (sites with low read depth or ambiguous

214    variants). We selected the filtering criteria in order to maximize the adjusted $R^2$

215    in the log-log regression of $\pi_N/\pi_S$ on $\pi_S$. By doing so we aimed to reduce the error

216    introduced by annotation and quality difference between model and non-model

217    organisms. Also, to evaluate the variance introduced by random sampling and

218    grouping of SNPs, we performed 1,000-iteration bootstraps to get the bootstrap

219    bias-corrected mean and 95% confidence intervals for $l$ calculations.

220

221    *Estimates of the distributions of fitness effects*

222

223    The distribution of fitness effects (DFE) for all non-synonymous mutations

224    across the genome was first calculated by considering only deleterious

225    mutations. We first re-used the DFE parameters estimated in 59 animal and

226    plant species in Chen *et al.* (2017) that assumes that only neutral and slightly

227    deleterious mutations contribute to genetic diversity. In brief, in this previous

228    study the DFE was modeled using a gamma distribution with mean $S_d$ and shape

229    parameter $\beta$. Folded site frequency spectra (SFS) were compared between

230    synonymous and nonsynonymous sites and demography (or any departure from

231    equilibrium) was taken into account by introducing $n-1$ nuisance parameters for

232    an unfolded SFS of size $n$, following the method proposed by Eyre-Walker *et al.*

233    (2006). The possible issues and merits of this approach compared to those

234    based on an explicit (albeit very simplified) demographic model have been

235    discussed previously and the method introduced by Eyre-Walker *et al.* (2006)

236    has proved to be relatively efficient (Eyre-Walker and Keightley 2007; Tataru *et*

237    *al.* 2017). The calculations were carried out using an in-house Mathematica

238    script implementing the method of Eyre-Walker *et al.* (2006) provided in

239    supplementary S2 file of Chen *et al.* (2017).

240    However, for species with large effective population sizes, like *D. melanogaster*,

241    ignoring the effects of beneficial mutations could distort the DFE to a great

242    extent and lead to a wrong estimate of $\beta$. Therefore, we further estimated the

243    DFE under a full model that takes both deleterious and beneficial mutations into

244    account (Tataru *et al.* 2017) using unfolded SFS for 11 species. Briefly, the model

245    mixes the gamma distribution of deleterious mutations (shape=$\beta$, mean=$S_d$) with

246    an exponential distribution of beneficial mutations (mean=$S_b$), in proportions of

247    (1-$p_b$) and $p_b$, respectively. The unfolded SFS was calculated for the 11 retained

248    species, for which a closely related outgroup with similar sequencing quality was

249    available to polarize the SFS. Ancestral state was assigned as the state of the

250    outgroup if the outgroup was monomorphic for one of the two variants, and the

251    derived allele frequency was calculated from this polarization. Otherwise (in the

252    case of missing data, polymorphic site or third allele in the outgroup) the site

253    was masked. The percentage of SNPs that could not be polarized and were

254    masked varied between 0 and 29.3% with a mean of 4.6% and a median value of

255    0.5% (Table S2).

256    In addition, since polarization errors could remain, the error rate of the ancestral

257    state assignment ($\varepsilon_{an}$) was also taken into account in *PolyDFE*. The "gamma" DFE

258    (that only considers deleterious mutations) and the full DFE were estimated for

259    each species. In both cases a nuisance parameter was also fitted to account for

260    possible mis-assignment errors in SNP ancestral allele estimation (a step

261    required to obtain the unfolded SFS). Note that, although we used outgroups to

262    polarize SFSs, we did not use divergence but only polymorphism to estimate the

263    effect of beneficial mutations. This is at the cost of larger variance in estimates

264    but it avoids the (potentially strong) bias due to ancient variations in $N_e$ that

9

265    cannot be captured by modeling recent changes in population size (Rousselle et

266    al. 2018). When comparing the estimates of the DFE among several species, the

267    problem arises that the best model is not necessarily the same for all species (the

268    best model can include or not beneficial mutations and include or not

269    polarization errors). Comparisons cannot be fairly done if all species do not

270    share the same model. Alternatively, estimations under an over-parameterized

271    model can lead to large variance and extreme values. To circumvent this problem,

272    we used a model averaging procedure where each parameter of interest ($\beta$, $S_b$, $S_d$,

273    and $p_b$) is estimated as a weighted mean of estimates obtained under four models:

274    the Gamma DFE and the full DFE models, including polarization errors or not.

275    The weights given to the estimate from model $k$ is $w_k = e^{-1/2\Delta\text{AIC}_k}$ where

276    $\Delta\text{AIC}_m = \text{AIC}_m - \text{AIC}_{min}$ with AIC being the Akaike Information Criterion and

277    $\text{AIC}_{min}$ the minimum AIC among the four models (Posada and Buckley 2004). All

278    calculations were performed using the software *polyDFE* and the associated R

279    script (Tataru *et al.* 2017).

280

281    *Expectations under different selection models*

282

283    Independently of possible indirect effects of selective sweeps, [Eq. 1] only

284    considers deleterious mutations, in line with the initial view of the Nearly

285    Neutral Theory where beneficial mutations negligibly contribute to

286    polymorphism (Ohta 1973). Giving more weight to beneficial mutations slightly

287    modified the relationship between the slope of the linear regression, *l,* and the

288    shape parameter, $\beta$. For beneficial mutations only, the equivalent of [Eq. 1] is

289    simply (see Appendix):

290

291    $ln(\pi_N/\pi_S) \approx +\beta_b\,ln(N_e) + constant$    [Eq. 2]

292

10

293     where $\beta_b$ is the shape of the distribution of beneficial mutations, still assuming a
294     gamma distribution, so $\beta_b$ would be 1 in the statistical framework we used. Thus,
295     the $\pi_N/\pi_S$ ratio increases with $N_e$, so that considering beneficial mutations the
296     global $\pi_N/\pi_S$ decreases more slowly than when only deleterious mutations are
297     taken into account. Thus, with beneficial mutations the slope will always be
298     lower than without. For the majority of species beneficial mutations are rare
299     ($p_b \ll 1$) and thus $b$ (thereafter we define $b$ = -$l$) is approximately equal to $\beta$. For
300     those with a relatively high proportion of beneficial mutations, direct positive
301     selection should result in a flattened slope, i.e. a smaller value of $b$ than $\beta$. As we
302     mostly observed the reverse pattern, $b > \beta$, the observed discrepancy cannot be
303     explained by the direct effect of beneficial mutations.
304
305     *Trends across the genome and tests for selection*
306
307     For each of the 20 bins defined above and ranked according to their mean
308     synonymous nucleotide diversity we calculated $\beta$, $p_b$ and $S_b$ values and a
309     summary statistic of the site frequency spectrum, Tajima's D (Tajima 1989).
310     Tajima's D tests for an excess of rare over intermediate variants compared to the
311     frequencies expected under the standard coalescent and was calculated from
312     synonymous sites Demography does affect Tajima's D and can explain the
313     difference among species. However, a negative Tajima's D is also expected under
314     recurrent selective sweeps (Jensen *et al.* 2005; Pavlidis and Alachiotis 2017) and
315     should be more negative in genomic regions more strongly affected by linked
316     positive selection. Background selection can also affect Tajima's D in the same
317     direction but much more weakly (Charlesworth et al. 1995). Independently of
318     the species mean value, we thus expect a strong positive relationship between
319     recombination and Tajima's D in species where linked positive selection is
320     prominent.
321
322     *Forward simulations under selective sweep scenario*
323
324     The code developed by Castellano et al (2018) which is based on forward
325     simulations using the software SLiM, version 3.2.1 (Haller and Messer 2019) was

11

326    modified to assess the effect of parameters $p_b$, $S_b$, and N on $b$ and Tajima's D.

327    More specifically, a 20-kb genomic region was simulated with a mutation rate of

328    $1 \times 10^{-6}$ to study the behavior of $b$ and Tajima's D under selective sweep scenarios

329    with varying parameters of $p_b$, $S_b$, and N. First, we simulated equal amounts of

330    neutral and deleterious mutations whose fitness effects were drawn from a

331    gamma distribution with a shape parameter 0.4 and a mean $s_d$ of -10. Different

332    percentages of beneficial mutations ($p_b$= 1%, 0.8%, 0.5%, 0.4%, 0.3%, 0.2%,

333    0.01%, 0.005% and 0) were drawn randomly from a distribution with a fixed $s_b$

334    of 1 to simulate loci experiencing selective sweeps at different frequency and we

335    then calculated $b$ (Fig. 5 of Castellano et al (2018)) and Tajima's D.  We also

336    investigated the behavior of $b$ and Tajima's D by varying $s_b$ (1, 0.5, 0.1), N (100,

337    500, 1000) and the recombination rate (Nr=0, 1e-3, 1e-2). Simulated values were

338    averaged across 50 samples, which were taken every 5N generations after an

339    initial burn-in period of 10N generations.

340

341    **Results**

342

343    *b and β are generally similar but the variance is large*

344

345    One of the most important predictions of the Nearly Neutral Theory is that the

346    proportion of effectively neutral mutations is a function of the effective

347    population size (Kimura and Ohta 1971; Ohta 1972; Ohta 1973; Ohta 1992). In

348    species with large effective population size, selection is efficient and the

349    proportion of effectively neutral mutations is small. Here we used the ratio of

350    genetic diversity at 0-fold over 4-fold degenerate sites ($\pi_N/\pi_S$) in protein coding

351    regions as a measure of the proportion of effectively neutral mutations and

352    examined the linearity between $\log(\pi_N/\pi_S)$ and $\log(\hat{N}_e)$ across the genomes of

353    59 species used in Chen *et al.* (2017). The slope (linear regression coefficient

354    between $\log(\pi_N/\pi_S)$ and $\log(\hat{N}_e)$) was negative for 51 of the 59 species ($l<0$),

355    although it was significantly different from zero at p=0.05 in less than half of the

356    species (28/59). The value of $l$ varied from -0.424 (*D. melanogaster*) to 0.22

357    (*Callithrix jacchus*) (Table S1). Since balancing selection can lead to both high $\pi_S$

358  and $\pi_N/\pi_S$, it can generate an increase in $\pi_N/\pi_S$ for high-$\pi_S$ bins. We thus

359  removed the five bins with the highest diversity and recalculated *l* values for all

360  species. This reduced the *l* values of 36 species and led to negative *l* values in 55

361  species.

362

363  We further examined the DFE for mutations across the genome in the same

364  datasets. A gamma distribution with two parameters, mean ($S_d$) and shape ($\beta$),

365  was used to describe the distribution of deleterious mutations under purifying

366  selection. Importantly, the contribution of beneficial mutations, even those under

367  weak selection that are potentially behaving neutrally, is ignored in this case.

368  Estimates of the shape parameter, $\beta$, varied from 0.01 (*C. jacchus*) to 0.347 (*D.*

369  *melanogaster*) but were only weakly correlated with effective population size

370  (Table S1).

371

372  Considering only deleterious mutations and assuming a simple scaling of $N_e$

373  variation across the genome, the slightly deleterious model predicts that the

374  value of the slope of the linear regression between $\log(\pi_N/\pi_S)$ and $\log(\hat{N}_e)$, *b*

375  (recall that *b* = -*l*), is equal to $\beta$ (Welch *et al.* 2008). The discrepancy between the

376  two might indicate a departure from this model, and Castellano *et al.* (2018)

377  suggested that in *D. melanogaster,* where the observed slope was steeper than

378  expected, the departure was caused by linked positive selection across the

379  genome. We observed a general consistency between $\beta$ and *b* as estimators of

380  effective neutrality (linear coef. = 1.04, intercept=0.007, p-value<2e-16, adjusted

381  $R^2$=0.35, Fig. 1A). The difference ($\Delta = b - \beta$) was small in 40 species and varied

382  from -0.1 to 0.1 (Fig. 1B). In 36 species (61%) *b* values were larger than $\beta$ and in

383  23 species (39%) $\beta$ was larger than *b*. However, the variation in $\Delta$ was not

384  explained by $\pi_S$ or $N_e$ as the adjusted $R^2$ was only 0.06.  Removing the five bins

385  with the highest diversity, the correlation between $\beta$ and *b* was still significant

386  (coef. 0.89, p-value=2.14e-6). The median value of $\Delta$ increased from 0.0085 to

387  0.045 but there was still no correlation between $\Delta$ and $\hat{N}_e$.

388

389  *The effects of quality control and full DFE model*

390

391    The variation in Δ may come from two sources. First, it can be due to the

392    estimation quality of $b$ and $\beta$. Tests have shown that quality control on

393    sequencing and SNP-calling can have a dramatic influence on $b$ calculations and

394    ignoring beneficial mutations in DFE model could also distort the estimates of $\beta$

395    (Tataru *et al.* 2017). Second, the variation in Δ can be caused by departures from

396    the assumptions underlying the simple version of the Nearly Neutral Theory, for

397    instance a larger role of direct or linked positive selection than assumed by the

398    theory.

399

400    To assess the relative importance of these two sources we selected 11 species

401    with genomic data of high quality and performed a series of stringent quality

402    controls (see details in M&M) before re-estimating $b$. This improved the

403    goodness of fit for the log linear regression between $\pi_N/\pi_S$ and $\pi_S$ across the

404    genome and $b$ estimates were significantly different from zero for all 11 species

405    (Table 1 and Fig. 2, see also details in Table S2 and Fig. S1). For estimating $\beta$, we

406    used closely related species to polarize the SFS and applied both the gamma DFE

407    model and the full DFE model implemented in *polyDFE*, which considers both

408    deleterious and beneficial mutations. Instead of choosing the best DFE model, an

409    average value weighted by the different models' AIC scores was calculated for

410    each parameter (Tataru and Bataillon 2019).

411

412    In this case we observed a better correlation between $b$ and $\beta$ (rho = 0.727, p-

413    value=0.011) than when we considered the 59 species and used only a gamma

414    DFE. In addition, considering beneficial mutations slightly increases $\beta$ estimates,

415    making them closer to $b$. However, the linear coefficient between $b$ and $\beta$ (1.26)

416    is significantly higher than one and the variation of Δ remains large (-0.026 ~

417    0.289) suggesting that some additional factors may lie behind the remaining

418    variation.

419

420    *The roles of effective population size and positive selection*

421

14

422 We then tested if the variation in Δ, where Δ=$b-\beta$, could simply reflect

423 differences in effective population size ($N_e$) among species. Estimates of $N_e$ were

424 obtained by rescaling $\pi_S$ using estimates of the mutation rate ($\mu$) from the

425 literature (see Table S3 for the sources of the $\mu$ estimates). When Δ is regressed

426 against log($\hat{N}_e$), log($\hat{N}_e$) explained up to 49% of the variance in Δ (p-

427 value=0.014). Considering the uncertainty in $\mu$, we also regressed Δ on log($\pi_S$),

428 and obtained similar results ($R^2$=0.41, p-value=0.019, Fig. 3).

429

430 Furthermore, we tested whether species with potentially more selective sweeps

431 show higher Δ, as predicted by Castellano *et al.* (2018). An explicit model of

432 selective sweeps is difficult to fit given the uncertainty about beneficial

433 mutations parameters and would require additional information, especially on

434 the recombination map of the different species. Alternatively, we qualitatively

435 reasoned that, in addition to be more frequent when the effective population is

436 large, the number of selective sweeps should increase with both the proportion

437 ($p_b$) and the mean strength of beneficial mutations ($S_b$). Log($S_b$) had a significant

438 and positive effect on Δ (p-value=0.0018, Fig. 3) and explained 64.3% of the

439 variance in Δ but the effect of $p_b$ was not significant (p-value=0.29). When

440 considered together, the effects of both log($S_b$) and log($\pi_S$) (or $\hat{N}_e$) in the joint

441 model explained up to 78% of the variance in Δ (p-value=0.0068 and 0.059,

442 respectively, Table 2). However, no significant effect of $p_b$ could be detected

443 either in the single regression model (p-value=0.29) or joint model with other

444 variables (p-value=0.15). The rate of adaptive evolution relative to the neutral

445 mutation rate, $\omega_a$ (Galtier 2016) combines the proportion ($p_b$) and the mean

446 strength of beneficial mutations ($S_b$) according to $\omega_a = p_S \times S_b / (1 - \exp(-S_b))$.

447 However, as for $p_b$ the effect of $\omega_a$ on Δ was not significant (p-value=0.17)

448 although the relationship is positive as expected.

449

450 *Trends across the genome and tests for selection*

451

452 Variation of DFE parameters across bins could also explain the difference

453 between $\beta$ and $b$ since the underlying assumption is that $\beta$ is constant across bins.

15

454  We thus calculated $\beta$ for all 20 bins for the 11 species. Seven species had $\beta$ values

455  increasing weakly with genetic diversity (p-value<0.05, mean regression

456  coefficient 0.056) while *C. grandiflora* and *H. timareta* had a much faster increase

457  (regression coefficient =0.2 and 0.15, respectively, Table 3). In five species, the

458  slope was steeper than the maximum $\beta$ value, similar to what was obtained by

459  Castellano *et al.* (2018) in *Drosophila*. However, the slope was shallower than the

460  maximum $\beta$ value in the six remaining species and in five of them the maximum $\beta$

461  value was larger than 1 (Table 1). We also compared $p_b$ and $S_b$ values across bins.

462  In *A. thaliana* $p_b$ increased slowly with diversity whereas in *C. grandiflora*, *S.*

463  *huaylasense*, and *D. melanogaster* $p_b$ decreased significantly (p-value<0.05). In all

464  11 species, $S_b$ did not show any significant trend across bins. To more formally

465  test for the significance of these variations, we also divided the genomes into five

466  bins (to get enough power per bin) and tested the invariance of the DFE across

467  bins using likelihood ratio tests as implemented in *polyDFE*. For all species, a

468  model with independent DFE parameters for each bin is significantly better than

469  a model with shared parameters across bins (see Table S4).

470

471  For all 11 selected species we also calculated Tajima's D (Tajima 1989),

472  thereafter simply called D, in each bin to test for departure from neutrality

473  across the genome. Mean values of D were slightly negative across bins for most

474  species except *S. habrochaites*. For nine of the eleven species, D values increased

475  significantly with genetic diversity (Table 3).  Interestingly, we found a negative

476  and strong correlation of Tajima's D with $\log(S_b)$ for all 11 species (p-

477  value=0.0086, Pearson's correlation coef. =-0.74) but not with any other DFE

478  parameters. This is in agreement with the expectation that selective sweeps

479  decrease D. Background selection could also decrease D albeit to a lower extent.

480  We further tested the trends of positive and negative selection by calculating the

481  proportions of deleterious or beneficial mutations over all bins with selective

482  strength <-10 and >10, respectively. However, no significant trends were

483  identified for either type of direct selection.

484

485  We also tested whether alternative measures of the possible occurrence of

486  selective sweeps could explain a larger part of the variation in Δ. We used both

16

487     the mean Tajima's D and the among-genome regression coefficient of the

488     relationship between D and $\pi_S$ ($\rho_D$) as predictors. More negative D and stronger

489     positive regression coefficient between D and $\pi_S$ can be viewed as signature of

490     stronger hitchhiking effects. So we would expect to see a negative effect of D and

491     a positive effect of $\rho_D$ on the variation in $\Delta$. In combination with $\pi_S$ (or $\hat{N}_e$), both

492     D and $\rho_D$ indeed explained a significant part of the variation in $\Delta$ (adjusted

493     $R^2$=0.76, Table 2).

494

495     *Simulations*

496

497     Castellano et al. (2018) used forward simulations to assess the extent to which

498     selective sweeps made the slope the relationship between $\log(\pi_N/\pi_S)$ and $\log(\hat{N}_e)$

499     steeper and thereby could explain the discrepancy between the slope and the

500     shape parameter of the DFE, $\beta$. They tested varying proportions of adaptive

501     mutations (their Fig. 5). We extended their investigation to test the effect of

502     selective strength ($s_b$) on $b$ with a fixed $\beta$ (0.4) and how selective strength ($s_b$)

503     also affected estimates of Tajima's D. Without recombination (Nr=0), Fig. 4

504     shows that when $s_b$ increased from 0.1 to 1, $b$ increased from 0.46 to 0.72

505     ($\Delta$=0.06 to 0.32). As expected mean Tajima's D decreased from -0.36 to -0.77 as

506     $s_b$ increased and $\rho_D$ between D and $\pi_S$ increased (see also Table 4). We also

507     increased N from 100 to 500, and to 1000, and fixed the mean selective strength

508     at either $S_b$ = 10 or $S_d$ = -1000. With these parameters, the strength of selection

509     was not affected by N but the number of sweeps increased with N due to the

510     higher input of (beneficial) mutations. In this case $\Delta$ increased from 0.06 to 0.41

511     as N increased and Tajima's D again decreased (Table 4 and Fig. 5). With

512     recombination (Nr=1e-3 and Nr=1e-2), we noticed similar trends of $b$, D, and $\rho_D$

513     when $s_b$ or N are large enough to recover the significance of the linearity between

514     $\log(\pi_N/\pi_S)$ and $\log(\pi_S)$ (Fig. S2 and S3).

515

516     **Discussion**

517

518    The aim of the present study was to test quantitatively one of the predictions of
519    the Nearly Neutral Theory of molecular evolution or, more precisely, the slightly
520    deleterious model. More specifically, we used full genome datasets to test
521    whether the proportion of effectively neutral mutations varies with local
522    variation in $N_e$ across the genome and decreases linearly with increasing $N_e$ and
523    whether the slope is equal to the shape parameter of the DFE. The negative log
524    linear relationship between $\pi_N/\pi_S$ and $N_e$ observed in previous studies
525    (Gossmann et al. 2011; Murray et al. 2017; Castellano et al. 2018; Vigué and
526    Eyre-Walker 2019) was also observed in the present study, although the slope
527    was not always significantly negative and, when negative, could differ
528    significantly from the shape parameter of the DFE and be much steeper. The
529    latter was especially true in species with large effective population size and the
530    difference was correlated to the estimated mean strength of selection acting on
531    beneficial mutations. In the case of species with large effective population size
532    neglecting linked positive selection could therefore lead to a significant
533    quantitative discrepancy between predictions and observations. On the other
534    hand, the slightly deleterious model appears as a good approximation when the
535    effective population size is small.  Below we first consider possible caveats and
536    discuss the implications of the results for the relative importance of purifying
537    and adaptive selection in shaping the genetic diversity of species.

538

539    The discrepancy between the slope of the log linear relationship between $\pi_N/\pi_S$
540    and $N_e$ and $\beta$ could simply be due to difficulties in estimating them precisely. In
541    general, estimates of the DFE shape parameter, $\beta$, were rather stable compared
542    to estimates of the slope of the regression of $\log(\pi_N/\pi_S)$ over $\log(\pi_S)$, $b$, with the
543    variance of the former being half that of the latter independently of quality
544    control and whether the SFS was folded or unfolded. High variation in $b$
545    estimates may explain the fact that a significant correlation between $\pi_N/\pi_S$ and
546    $\pi_S$ could not be observed for all species, particularly those with low genetic
547    diversity (e.g. great apes). Therefore, a stringent quality control for read
548    alignment and SNP calling is necessary, even for *D. melanogaster*, where an
549    improvement of the fit in *l* calculation (linear regression adjusted $R^2$=0.79 to 0.95)
550    leads to a dramatic change in the estimate of $\Delta$ (from 0.077 to 0.29). Even if a

18

551    stringent quality control had been implemented, the goodness of fit for the log

552    linear regression leading to the estimation of $b$ would differ significantly from

553    species to species. The fit across the *D. melanogaster* and *A. thaliana* genomes

554    was almost perfect ($R^2 > 0.95$) while, at the other extreme, the fit was rather poor

555    in *S. habrochaites* ($R^2 = 0.38$). However, even among species for which the fit is

556    almost perfect ($R^2 > 0.95$) $b$ could vary rather dramatically: *D. melanogaster* had a

557    much larger $l$ (0.7) than *A. thaliana* (0.48), *C. rubella* (0.43), and *Z. mays* (teosinte,

558    0.29), whereas $\beta$ only changed marginally for these species. Not all species

559    though showed a significant negative linear relationship between $\pi_N/\pi_S$ and $\hat{N}_e$

560    and some even had positive slopes, especially for those of low diversity (e.g.

561    great apes, Fig 2). Therefore, besides purifying selection the slope is also likely to

562    be affected by additional factors. Factors that affect the likelihood to observe a

563    negative relationship between $\pi_N/\pi_S$ and $\hat{N}_e$ and its relationship with the DFE

564    parameters were thoroughly discussed by Castellano et al. (2018). Below we

565    highlight those that seem particularly relevant when considering a group of

566    species with contrasted levels of diversity as was done here. These factors are

567    the variation in $N_e$ estimates along the genome, which itself reflects the joint

568    distribution along the genome of recombination rate and density of selected sites,

569    the DFE, and the variation along the genome of the rate of adaptive evolution

570    (Castellano et al. 2018).

571

572    Lack of joint variation in recombination rate and selected sites seems to be an

573    unlikely cause for an absence of negative relationship between $\pi_N/\pi_S$ and $N_e$ as

574    such a relationship is observed in selfing species where this joint variation is

575    expected to the more limited than in outcrossing ones.  A possible source of

576    variance in $\beta$ could be that the single-sided gamma distribution does not describe

577    well the real DFE curves, at least not for all species, particularly when the DFE is

578    not unimodal (Tataru *et al.* 2017). For species like *D. melanogaster*, for instance,

579    there is mounting evidence of adaptive evolution (reviewed in Eyre-Walker 2006,

580    Sella et al. 2009). Therefore, it is necessary to consider the possible contribution

581    of beneficial mutations. The full DFE model provided a much better fit than the

582    gamma DFE that considers only deleterious mutations in *D. melanogaster* (log

583   likelihood= -187.3 versus -245.7, respectively). This was also true of some of the

584   outcrossing plants like *Capsella grandiflora*, and *Solanum huaylasense*. In all three

585   species $\beta$ estimates increased when estimated with the Full DFE instead of the

586   Gamma DFE, sometimes significantly (from 0.33 to 0.41 in *D. melanogaster*

587   (Rwanda) and 0.15 to 0.31 in *S. huaylasense*) and at other times only marginally

588   (0.27 to 0.30 in *C. grandiflora*). Taking beneficial mutations into account when

589   fitting the shape of the DFE can partly reduce the discrepancy between $\beta$

590   estimates and the slope of the regression. However, it is not sufficient as $\Delta$ was

591   positive in 10 over the 11 focal species we studied.

592

593   Based on the prediction of the Nearly Neutral Theory with direct positive

594   selection (Equation 2), the proportion of beneficial mutations is the only factor

595   that could alter the relationship between $b$ and $\beta$ and should always result in a

596   larger $\beta$ compared to $b$. However, this is usually not the case as, on the contrary,

597   values of $b$ larger than $\beta$ have generally been reported (Chen *et al.* 2017; James *et*

598   *al.* 2017; Castellano *et al.* 2018). In this paper we systematically investigated this

599   relationship across the genomes of multiple species. Two thirds of the 59 species

600   and 10 out of the subset of eleven species that were selected for the high quality

601   of their genome, had larger $b$ than $\beta$ values. Hence direct positive selection is not

602   the main cause of the discrepancy.

603

604   Investigation of DFE parameter changes across bins may help to identify changes

605   in natural selection. Increasing $\beta$ values over bins could be a signal for stronger

606   positive selection in low diversity regions. Although the maximum $\beta$ value of

607   some species can be larger than $b$, $\beta$ grows slowly for most species and shows

608   hardly any pattern between species. Neither did $p_b$ or $S_b$. This lack of significant

609   trend in these parameters could simply be due to an increase in variance of their

610   estimates as only one twentieth of the total number of polymorphic sites were

611   used for DFE calculations in each bin. It could also again suggest that direct

612   selection is not the main cause of the discrepancy.

613

614   One of the main findings of the present study is that a large proportion of

615   variance in the discrepancy can be explained by the estimated strength of

616  positive selection, which can be regarded as an indication for linked selection,
617  such as selective sweeps or more generally hitchhiking effects.  To test for that,
618  we compared changes in Tajima's D and its among-genome correlation
619  coefficients over bins. As expected we observed a negative effect of D and a
620  positive effect of $\rho_D$ on $\Delta$, both suggesting the presence of linked selection, with
621  lower diversity at nearby sites and thus increased discrepancy between $b$ and $\beta$.
622  This is also in agreement with our simulations and those of Castellano et al.
623  (2018) that illustrate that hitchhiking effects can lower the genetic diversity at
624  nearby neutral or nearly neutral positions. These results can be understood
625  because selective sweep effects cannot simply be captured by a rescaling of $N_e$.
626  Selective sweeps not only reduce genetic diversity at linked sites but also distort
627  the coalescent genealogy (Fay and Wu 2000; Walsh and Lynch 2018; Campos
628  Parada and Charlesworth 2019), so that we cannot define a single $N_e$ in this
629  context (Weissman and Barton 2012). In particular, the scaling is not expected to
630  be the same for neutral or weakly selected polymorphisms. However, as far as
631  we know, there is no quantitative model predicting the value of the slope as a
632  function of DFE, rates of sweep and recombination rates, and such models still
633  need to be developed.
634
635  **Conclusions**
636
637   There are three major conclusions to the present study. First, the Nearly Neutral
638  Theory in its initial form may not explain all aspects of polymorphisms but,
639  almost 50 years after it was first proposed by Tomoko Ohta (Ohta 1973), it still
640  constitutes an excellent starting point for further theoretical developments
641  (Galtier 2016; Walsh and Lynch 2018). Second, considering linked beneficial
642  selection indeed helps to explain more fully polymorphism data, and this is
643  especially true for species with high genetic diversity. This can explain both
644  patterns of synonymous polymorphism (Corbett-Detig et al. 2015) and how
645  selection reduces non-synonymous polymorphism (Castellano et al. 2018, this
646  study). One could have a progressive increase of the effect of selective sweeps as
647  suggested by Walsh and Lynch (2018, chapter 8) with a shift from genetic drift to
648  genetic draft (Gillespie 1999; 2000; 2001). If so, we could have three domains.

649    For small population sizes, drift would dominate and the nearly neutral theory in
650    its initial form would apply. For intermediate population sizes beneficial
651    mutations would start to play a more important part, and finally for large
652    population sizes, the effect of selective sweeps would dominate and draft would
653    be the main explanation of the observed pattern of diversity. Third, our study
654    once more emphasizes the central importance of the DFE in evolutionary
655    genomics and we will likely see further developments in this area.

656

661

662    **Data availability:** The vcf files used in the present study are available on request.

663

22

664     **Table 1** Species and datasets used in the present study
665

| Species | Ref. | Outgroup | Ref. | Mating type | AIC | $b$ | $\beta_{full}$ | $\beta_{gamma}$ | $\beta_{max}$ |
|---|---|---|---|---|---|---|---|---|---|
| *A. thaliana* | Alonso-Blanco *et al.* (2016) | *A. lyrata* | (Novikova *et al.* 2016) | selfing | 231.3, 227.3 | 0.48 | 0.32 | 0.32 | 0.45 |
| *A. lyrata* | (Novikova *et al.* 2016) | *A. thaliana* | Alonso-Blanco *et al.* (2016) | outcrossing | 247.4, 243.4 | 0.50 | 0.35 | 0.34 | 0.36 |
| *C. rubella* | (Koenig *et al.* 2018) | *C. grandiflora* | (Agren *et al.* 2014) | selfing | 201.4, 200.3 | 0.43 | 0.39 | 0.26 | 2.86 |
| *C. grandiflora* | (Agren *et al.* 2014) | *C. rubella* | (Koenig *et al.* 2018) | outcrossing | **321.9,** 327.8 | 0.52 | 0.30 | 0.27 | 0.36 |
| *S. habrochaites* | Aflitos *et al.* (2014) | *S. lycopersicon* | Aflitos *et al.* (2014) | selfing | **141.5,** 148.1 | 0.21 | 0.23 | 0.13 | 3.61 |
| *S. huaylasense* | Aflitos *et al.* (2014) | *S. lycopersicon* | Aflitos *et al.* (2014) | outcrossing | **87.1,** 121.5 | 0.54 | 0.31 | 0.15 | 3.89 |
| *S. propinquum* | Mace *et al.* (2013) | *S. bicolor* | Mace *et al.* (2013) | selfing | 163.8, 159.8 | 0.37 | 0.26 | 0.26 | 0.34 |
| *Z. mays* (teosinte) | Chia *et al.* (2012) | *T. dactyloides* | Chia *et al.* (2012) | outcrossing | 208.1, 204.1 | 0.29 | 0.19 | 0.18 | 0.45 |
| *P. trichocarpa* | Evans *et al.* (2014) | *P. nigra* | (Faivre-Rampant *et al.* 2016) | outcrossing | 318.9, 319.6 | 0.42 | 0.22 | 0.16 | 2.21 |
| *D. melanogaster* | Huang *et al.* (2014) | *D. simulans* | Stanley and Kulathinal (2016) | outcrossing | **422.7,** 535.5 | 0.70 | 0.41 | 0.33 | 0.51 |
| *H. timareta* | Martin *et al.* (2013) | *H. melpomene* | Martin *et al.* (2013) | outcrossing | 208.2, 204.2 | 0.44 | 0.21 | 0.21 | 2.78 |

666     Note: AIC values were estimated by *polyDFE* for models with and without the effects of beneficial mutations, respectively (bold numbers showed significance <

667     0.05). The same applies to $\beta_{full}$ and $\beta_{gamma}$ as well. $\beta_{max}$ corresponds to the maximum value of those estimated by *polyDFE* for each ranked gene bin.

668 **Table 2** Summary table of multiple regression analyses of the effects of $\pi_S$ $S_b$,

669 Tajima's D, and $\rho_D$ on $\Delta$, the difference between $b$ and $\beta$.

670

| $\Delta \sim \pi_S + \log_{10}(S_b)$ | *Coef.* | *SE* | *t value* | *p-value* |
|---|---|---|---|---|
| Intercept | 0.14 | 0.031 | 4.69 | 0.0016** |
| $\pi_S$ | 7.93 | 2.96 | 2.68 | 0.028* |
| $\log_{10}(S_b)$ | 0.015 | 3.6e-3 | 4.24 | 0.0029** |
| p-value: 0.0008144 | Adjusted R²: 0.7888 | | | |
| $\Delta \sim \pi_S + D + \rho_D$ | | | | |
| Intercept | -0.031 | 0.035 | -0.87 | 0.41 |
| Tajima's D | -0.10 | 0.042 | -2.39 | 0.048* |
| $\rho_D$ | 0.0015 | 6.05e-4 | 2.56 | 0.038* |
| $\pi_S$ | 15.80 | 3.39 | 4.65 | 0.0040** |
| p-value: 0.002978 | Adjusted R²: 0.708 | | | |

671
672 ***: p<0.001, **: 0.001<p<0.01, *: 0.01<p<0.05, ˙: 0.05<p<0.1

673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699

700  **Table 3** Changes of summary statistics and DFE parameters across 20 rank gene

701  groups.

702

| | Tajima's D | | $\rho_\beta{}^a$ | $\rho p_b{}^a$ |
|---|---|---|---|---|
| | **median** | $\rho_D{}^a$ | | |
| *A. thaliana* | -0.38 | 20.10*** | 0.033*** | 9.65e-4** |
| *A. lyrata* | -0.60 | 30.13*** | 0.057* | 7.75e-5 |
| *C. rubella* | -0.28 | 15.75* | 0.039* | 8.26e-4 |
| *C. grandiflora* | -1.06 | 23.02** | 0.20*** | -3.53e-3˙ |
| *S. habrochaites* | 0.22 | -5.36 | 0.11 | -7.48e-3 |
| *S. huaylasense* | -0.17 | -8.59** | -0.32 | -5.54e-2*** |
| *S. propinquum* | -0.10 | 60.04*** | 0.075*** | 1.82e-3 |
| *Z. mays* | -0.52 | -0.39 | 0.055*** | 2.39e-3 |
| *P. trichocarpa* | -0.43 | 79.20*** | 0.079 | -2.80e-3 |
| *D. melanogaster* | -0.73 | 7.41** | 0.078*** | -3.81e-3*** |
| *H. timareta* | -0.10 | 6.58** | 0.15*** | 9.87e-4 |

703

704  a: $\rho$ is the slope of the regression of D ($\beta$, and $p_b$, respectively) over genetic

705  diversity across ranked groups of genes.

706  ***: $p<0.001$, **: $0.001<p<0.01$, *: $0.01<p<0.05$, ˙: $0.05<p<0.1$

707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730

731

732 **Table 4** Results of forward simulations showing the effect of linked positive selection on $b$, $\Delta$ and summary statistics of the site frequency
733 spectrum for different values of the mean selective value of beneficial mutations, $S_b$ and the population size, N. $\rho_D$ is the correlation between
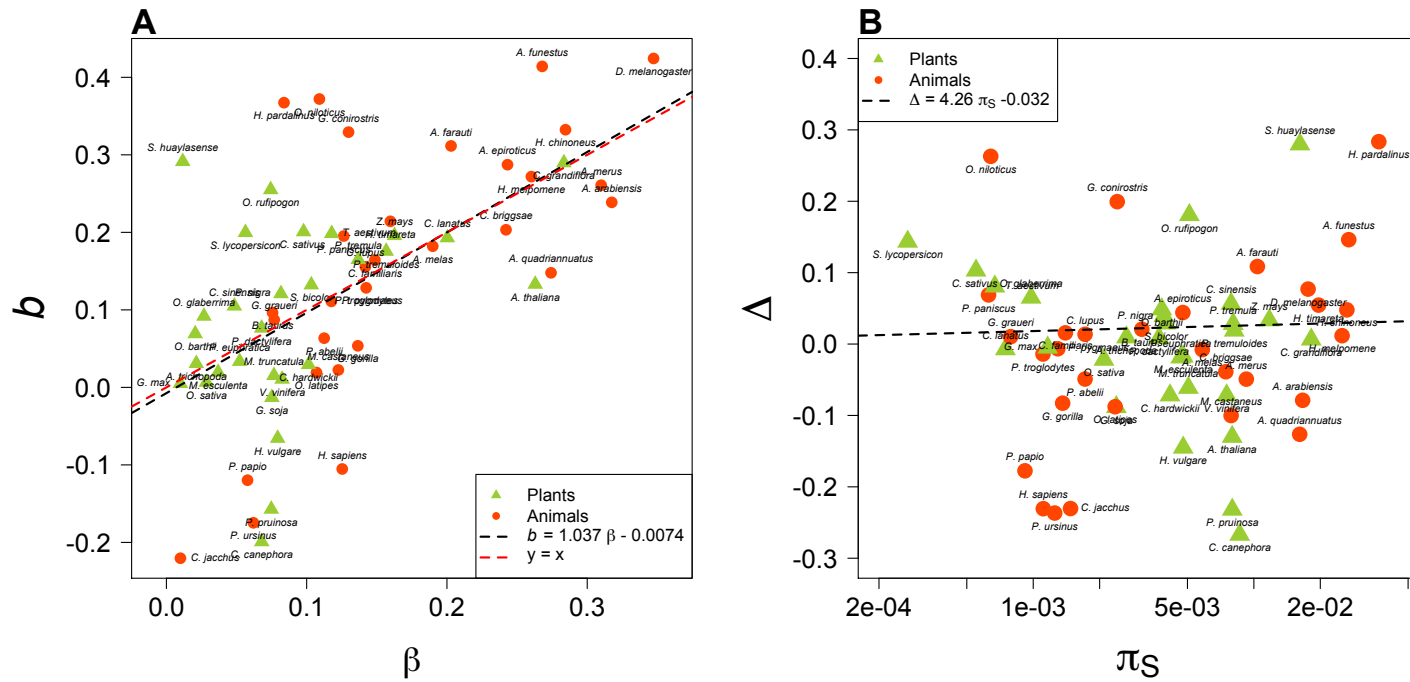734 $\pi_S$ and Tajima's D.

|  | N | $S_b$ | $S_d$ | $\beta$ | $b$ | $\Delta$ | $\pi_S$ | $\pi_N/\pi_S$ | $\rho_D$ | TD |
|---|---|---|---|---|---|---|---|---|---|---|
| Nr=0 | 100 | 20 | 1000 | 0.4 | 0.49 | 0.09 | 1.39 | 0.091 | 874.6 | -0.46 |
|  | 100 | 50 | 1000 | 0.4 | 0.61 | 0.21 | 1.18 | 0.094 | 909.9 | -0.70 |
|  | 100 | 100 | 1000 | 0.4 | 0.72 | 0.32 | 1.06 | 0.111 | 994.2 | -0.77 |
|  | 100 | 10 | 1000 | 0.4 | 0.46 | 0.06 | 1.52 | 0.082 | 739.9 | -0.36 |
|  | 500 | 10 | 1000 | 0.4 | 0.65 | 0.25 | 5.72 | 0.09 | 228.6 | -0.77 |
|  | 1000 | 10 | 1000 | 0.4 | 0.81 | 0.41 | 10.35 | 0.094 | 132.4 | -0.92 |
| Nr=1e-3 | 100 | 20 | 1000 | 0.4 | 0.06 | -0.34 | 1.64 | 0.076 | 662.5 | -0.18 |
|  | 100 | 50 | 1000 | 0.4 | 0.63 | 0.23 | 1.48 | 0.087 | 738.1 | -0.28 |
|  | 100 | 100 | 1000 | 0.4 | 0.72 | 0.32 | 1.17 | 0.097 | 966.8 | -0.58 |
|  | 100 | 10 | 1000 | 0.4 | 0.09 | 0.031 | 1.70 | 0.075 | 1011.1 | -0.12 |
|  | 500 | 10 | 1000 | 0.4 | 0.61 | 0.21 | 7.54 | 0.084 | 163.9 | -0.26 |
|  | 1000 | 10 | 1000 | 0.4 | 0.68 | 0.28 | 13.67 | 0.083 | 99.7 | -0.37 |
| Nr=1e-2 | 100 | 20 | 1000 | 0.4 | 0.43 | 0.03 | 1.74 | 0.077 | 739.3 | -0.048 |
|  | 100 | 50 | 1000 | 0.4 | 0.63 | 0.23 | 1.67 | 0.081 | 917.6 | -0.12 |
|  | 100 | 100 | 1000 | 0.4 | 0.78 | 0.38 | 1.61 | 0.084 | 898.4 | -0.15 |
|  | 100 | 10 | 1000 | 0.4 | 0.33 | -0.07 | 1.76 | 0.080 | 325.7 | -0.011 |
|  | 500 | 10 | 1000 | 0.4 | 0.69 | 0.29 | 8.55 | 0.073 | 165.4 | -0.06 |

| 1000 | 10 | 1000 | 0.4 | 0.99 | 0.59 | 16.7 | 0.072 | 86.3 | -0.12 |

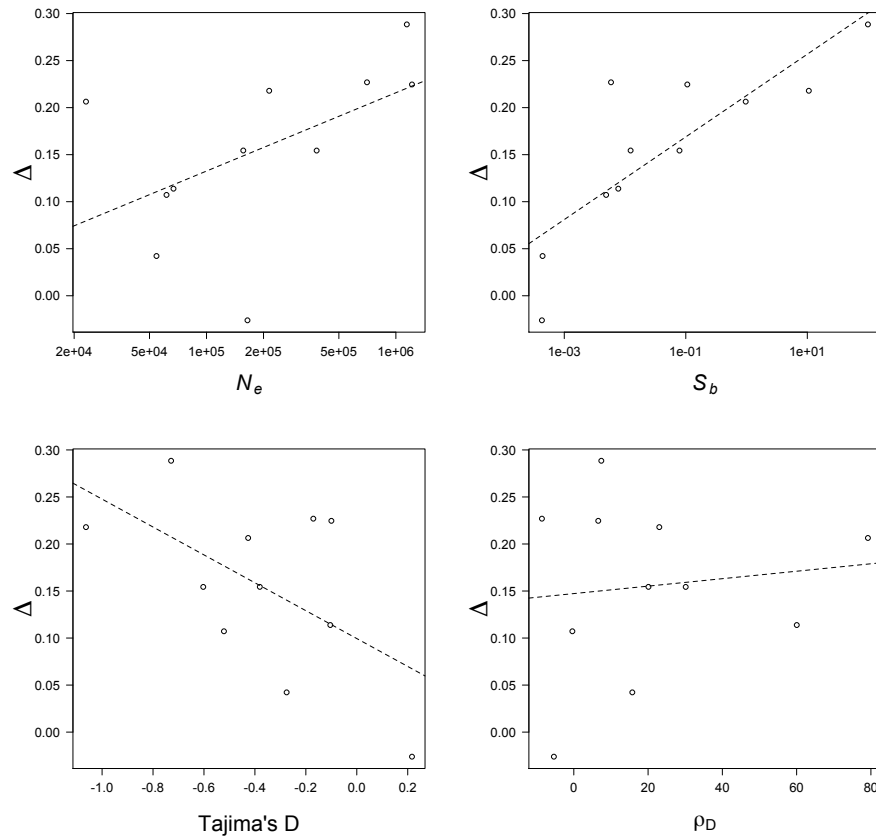**Figures**

**Fig. 1** (A) The correlation between *b* and the shape parameter of the DFE, *β,* from the 59 species in Chen *et al.* (2017). The observed slope of the regression of log($\pi_N$/ $\pi_S$) over log($\pi_S$), *l=-b*. (B) The distribution of Δ (=*b-β*) against genetic diversity at synonymous sites. *β* values were estimated from DFE models with only deleterious mutations considered (the gamma distribution).

743

**Fig. 2** The regression of $\log(\pi_N/\pi_S)$ over $\log(\pi_S)$ for self-fertilizing *Arabidopsis thaliana* (dots) and its outcrossing relative *A. lyrata* (triangles).
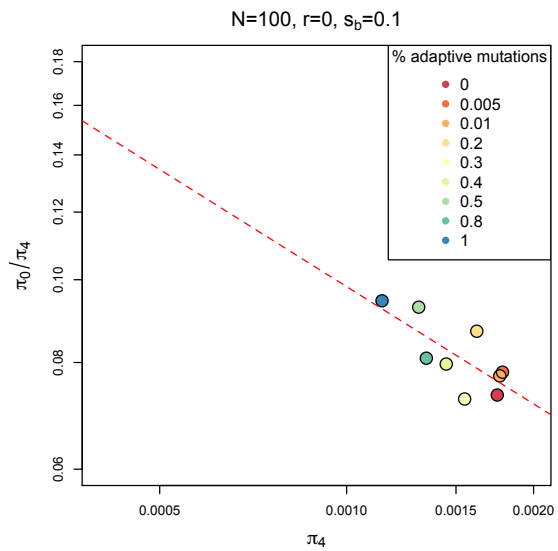
**Fig. 3** The relationship between Δ (=$b$-$\beta$) and effective population size, $N_e$, selective strength, $S_b$, Tajima's D and the trend of D across bins $\rho_D$ for 11 selected species. Dotted lines showed the linear regression line. $\beta$ and $S_b$ values were estimated from full DFE models with both deleterious and beneficial mutations considered (full DFE model with both gamma and exponential distributions).
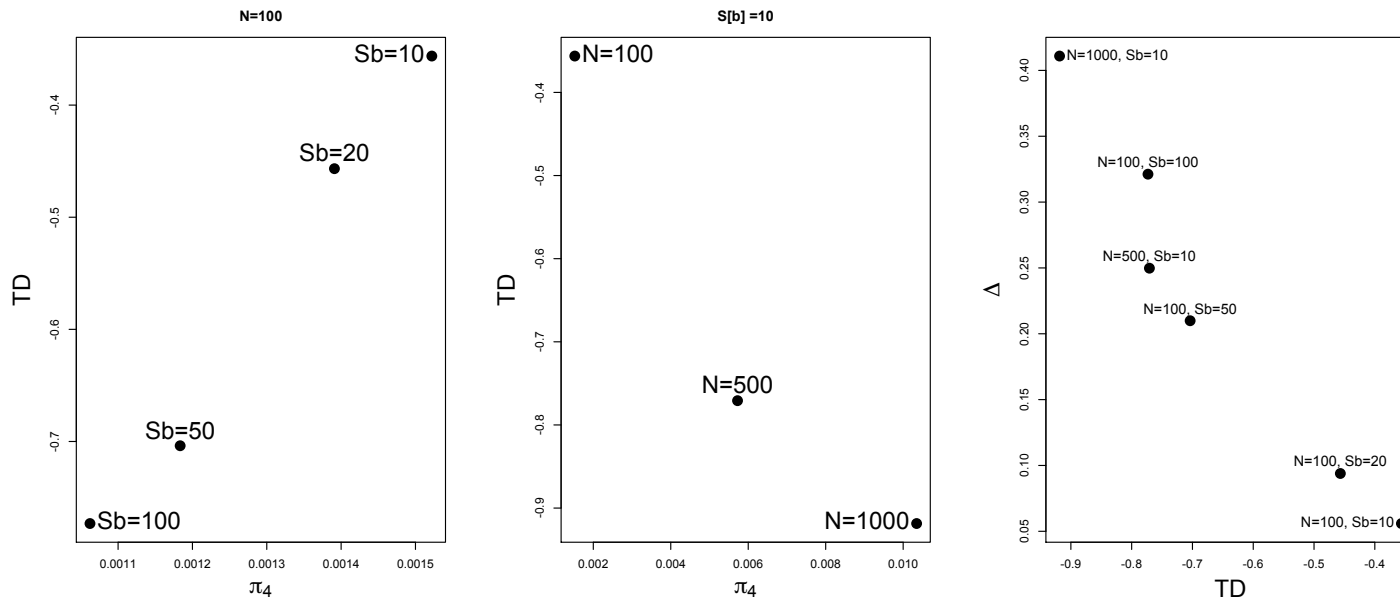
752    **Fig. 4** Effect of linked positive selection on the relationship between $\log(\pi_N/\pi_S)$ and $\log(N_e)$ and Tajima's D. Upper row: The linear

753    regression coefficient ($b$) between $\log(\pi_N/\pi_S)$ and $\log(N_e)$ increases with increasing positive selective strength (from left to right). The

754    red lines are the regression lines for each case. To facilitate comparisons among figures, and illustrate how the slope gets steeper as $s_b$

755    increases the regression lines corresponding to $s_b=0.1$ and/or $s_b=0.5$ values are reported with gray lines. Lower row: The red lines for

756    Tajima's D panels indicate the mean values.

757

758



759

**Fig. 5** The correlation between Tajima's D and $\pi_S$ depending on $S_b$ (left panel) and N (middle panel); Correlation between $\Delta$ and Tajima's

D (right panel). In all three cases the results were obtained with forward simulations in Slim assuming no recombination.

762

763

**Supplementary Information**

**Supplementary table**


**Supplementary table legends**


**Table S1**.  The 59 species used to compare the difference between -*l* and *β* assuming a gamma model for DFE. See Chen et al. (2017) for further details.


**Table S2**. Details of the 11 species used in the current study to compare the difference between -*l* and *β* assuming a full model (gamma + exponential) for the DFE.


**Table S3**. Mutation rates used for 11 species used in the current study for estimation of $N_e$.


**Table S4**. Test for the invariance of DFE parameter estimates across bins by comparing the log-likelihoods of independent estimates for each bin against those of shared estimates.

## References

Aflitos, S., E. Schijlen, H. de Jong, D. de Ridder, S. Smit *et al.*, 2014 Exploring genetic variation in the tomato (Solanum section Lycopersicon) clade by whole-genome sequencing. Plant Journal 80**:** 136-148.

Agren, J. A., W. Wang, D. Koenig, B. Neuffer, D. Weigel *et al.*, 2014 Mating system shifts and transposable element evolution in the plant genus Capsella. Bmc Genomics 15.

Alonso-Blanco, C., J. Andrade, C. Becker, F. Bemm, J. Bergelson *et al.*, 2016 1,135 Genomes Reveal the Global Pattern of Polymorphism in Arabidopsis thaliana. Cell 166**:** 481-491.

Barton, N. H., 1995 Linkage and the Limits to Natural-Selection. Genetics 140**:** 821-841.

Brandvain, Y., & Wright, S. I. 2016 The limits of natural selection in a nonequilibrium world. Trends in Genetics 32:201–210.

Campos Parada, J. L., and B. Charlesworth, 2019 The effects on neutral variability of recurrent selective sweeps and background selection. bioRxiv**:** 358309.

Carson, A. R., E. N. Smith, H. Matsui, S. K. Braekkan, K. Jepsen *et al.*, 2014 Effective filtering strategies to improve data quality from population-based whole exome sequencing studies. Bmc Bioinformatics 15.

Castellano, D., J. James and A. Eyre-Walker, 2018 Nearly Neutral Evolution Across the Drosophila melanogaster Genome. Molecular Biology and Evolution**:** msy164-msy164.

Charlesworth, B., M. T. Morgan and D. Charlesworth, 1993 The Effect of Deleterious Mutations on Neutral Molecular Variation. Genetics 134**:** 1289-1303.

Chen, J., S. Glemin and M. Lascoux, 2017 Genetic Diversity and the Efficacy of Purifying Selection across Plant and Animal Species. Molecular Biology and Evolution 34**:** 1417-1428.

Chia, J. M., C. Song, P. J. Bradbury, D. Costich, N. de Leon *et al.*, 2012 Maize HapMap2 identifies extant variation from a genome in flux. Nature Genetics 44**:** 803-U238.

Comeron, J.M. 2017 Background selection as null hypothesis in population genomics: insights and challenges from Drosophila studies. Phil. Trans. R. Soc. B 372: 20160471.

Coop, G. 2016 Does linked selection explain the narrow range of genetic diversity across species? bioRxiv doi: https://doi.org/10.1101/042598.

Corbett-Detig, R. B., D. L. Hartl and T. B. Sackton, 2015 Natural Selection Constrains Neutral Diversity across A Wide Range of Species. Plos Biology 13.

Cvijovic, I., B.H. Good and M.M. Desai, 2018 The effect of strong purifying selection on genetic diversity. Genetics 209: 1235-1278.

Do, R., D. Balick, H. Li, I. Adzhubei, S. Sunyaev and D. Reich 2015 No evidence that selection has been less effective at removing deleterious mutations in Europeans than in Africans. Nature genetics 47: 126

Ellegren, H., and N. Galtier, 2016 Determinants of genetic diversity. Nature Reviews Genetics 17**:** 422-433.

Evans, L. M., G. T. Slavov, E. Rodgers-Melnick, J. Martin, P. Ranjan *et al.*, 2014 Population genomics of Populus trichocarpa identifies signatures of selection and adaptive trait associations. Nature Genetics 46**:** 1089-1096.

Eyre-Walker, A., 2006 The genomic rate of adaptive evolution. Trends in Ecology & Evolution 21**:** 569-575.

Eyre-Walker, A., and P. D. Keightley, 2007 The distribution of fitness effects of new mutations. Nature Reviews Genetics 8**:** 610-618.

Eyre-Walker, A., and P. D. Keightley, 2009 Estimating the Rate of Adaptive Molecular Evolution in the Presence of Slightly Deleterious Mutations and Population Size Change. Molecular Biology and Evolution 26**:** 2097-2108.

Faivre-Rampant, P., G. Zaina, V. Jorge, S. Giacomello, V. Segura *et al.*, 2016 New resources for genetic studies in Populus nigra: genome-wide SNP discovery and development of a 12k Infinium array. Molecular Ecology Resources 16**:** 1023-1036.

Fay, J. C., and C. I. Wu, 2000 Hitchhiking under positive Darwinian selection. Genetics 155**:** 1405-1413.

Galtier, N., 2016 Adaptive Protein Evolution in Animals and the Effective Population Size Hypothesis. Plos Genetics 12.

Gillespie, J. H., 1999 The role of population size in molecular evolution. Theoretical Population Biology 55**:** 145-156.

Gillespie, J. H., 2000 Genetic drift in an infinite population: The pseudohitchhiking model. Genetics 155**:** 909-919.

Gillespie, J. H., 2001 Is the population size of a species relevant to its evolution? Evolution 55**:** 2161-2169.

Gillespie, J. H., 2004 *Population genetics : a concise guide.* Johns Hopkins University Press, Baltimore, Md.

Gordo, I., Dionisio, F. 2005 Nonequilibrium model for estimating parameters of deleterious mutations. *Physical Review. E, Statistical, Nonlinear, and Soft Matter Physics* 71: 031907.

Gossmann, T. I., Woolfit, M., Eyre-Walker, A. 2011 Quantifying the variation in the effective population size within a genome. *Genetics* 189:1389-1402.

Haller, B. C., and P. W. Messer, 2019 SLiM 3: Forward Genetic Simulations Beyond the Wright-Fisher Model. Molecular Biology and Evolution 36**:** 632-637.

Huang, W., A. Massouras, Y. Inoue, J. Peiffer, M. Ramia *et al.*, 2014 Natural variation in genome architecture among 205 Drosophila melanogaster Genetic Reference Panel lines. Genome Research 24**:** 1193-1208.

James, J., D. Castellano and A. Eyre-Walker, 2017 DNA sequence diversity and the efficiency of natural selection in animal mitochondrial DNA. Heredity 118**:** 88-95.

Jensen, J. D., B.A. Payseur, W. Stephan, C.F. Aquadro, M. Lynch *et al.*, 2018 The Importance of the Neutral Theory in 1968 and 50 years on. [submitted].

Jensen, J. D., Y. Kim, V. B. DuMont, C. F. Aquadro and C. D. Bustamante, 2005 Distinguishing between selective sweeps and demography using DNA polymorphism data. Genetics 170**:** 1401-1410.

Jensen, J. D., B. A. Payseur, W. Stephan, C. F. Aquadro, M. Lynch *et al.*, 2019 The importance of the Neutral Theory in 1968 and 50 years on: A response to Kern and Hahn 2018. Evolution 73**:** 111-114.

Kern, A. D., and M. W. Hahn, 2018 The Neutral Theory in Light of Natural Selection. Molecular Biology and Evolution 35**:** 1366-1371.

881  Kimura, M., 1979 Model of Effectively Neutral Mutations in Which Selective
882       Constraint Is Incorporated. Proceedings of the National Academy of
883       Sciences of the United States of America 76: 3440-3444.
884  Kimura, M., 1983 *The Neutral Theory of Molecular Evolution*. Cabridge, UK:
885       Cambridge Univ. Press.
886  Kimura, M., and T. Ohta, 1971 Protein Polymorphism as a Phase of Molecular
887       Evolution. Nature 229: 467-&.
888  Koenig, D., J. Hagmann, R. Li, F. Bemm, T. Slotte *et al.*, 2018 Long-term balancing
889       selection drives evolution of immunity genes in Capsella. bioRxiv.
890  Kreitman, M., 1996 The neutral theory is dead. Long live the neutral theory.
891       Bioessays 18: 678-683.
892  Li, H., and R. Durbin, 2010 Fast and accurate long-read alignment with Burrows-
893       Wheeler transform. Bioinformatics 26: 589-595.
894  Mace, E. S., S. S. Tai, E. K. Gilding, Y. H. Li, P. J. Prentis *et al.*, 2013 Whole-genome
895       sequencing reveals untapped genetic potential in Africa's indigenous
896       cereal crop sorghum. Nature Communications 4.
897  Martin, S. H., K. K. Dasmahapatra, N. J. Nadeau, C. Salazar, J. R. Walters *et al.*, 2013
898       Genome-wide evidence for speciation with gene flow in Heliconius
899       butterflies. Genome Research 23: 1817-1828.
900  McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis *et al.*, 2010 The
901       Genome Analysis Toolkit: A MapReduce framework for analyzing next-
902       generation DNA sequencing data. Genome Research 20: 1297-1303.
903  Murray, G.G.R., A. E. R. Soares, B. J. Novak, N. K. Schaefer, J. A. Cahill, A. J. Baker, J.
904       R. Demboski, A. Doll, R. R. Da Fonseca, T. L. Fulton, M. T. P. Gilbert, P. D.
905       Heintzman, B. Letts, G. McIntosh, B.L. O'Connell, M. Peck, M.-L. Pipes, E. S.
906       Rice, K. M. Santos, A. G. Sohrweide, S. H. Vohr, R. B. Corbett-Detig, R. E.
907       Green and B. Shapiro 2017 Natural selection shaped the rise and fall of
908       passenger pigeon genomic diversity. Science *358*(6365), 951–954.
909  Nei, M., Y. Suzuki and M. Nozawa, 2010 The Neutral Theory of Molecular
910       Evolution in the Genomic Era. Annual Review of Genomics and Human
911       Genetics, Vol 11 11: 265-289.
912  Novikova, P. Y., N. Hohmann, V. Nizhynska, T. Tsuchimatsu, J. Ali *et al.*, 2016
913       Sequencing of the genus Arabidopsis identifies a complex history of
914       nonbifurcating speciation and abundant trans-specific polymorphism.
915       Nature Genetics 48: 1077-+.
916  Ohta, T., 1972 Population Size and Rate of Evolution. Journal of Molecular
917       Evolution 1: 305-314.
918  Ohta, T., 1973 Slightly Deleterious Mutant Substitutions in Evolution. Nature
919       246: 96-98.
920  Ohta,T., 1977 Extension to the neutral mutation random drift hypothesis, pp.
921       148-167 in *Molecular Evolution and Polymorphism,* edited by M. Kimura
922       National Institute of Genetics, Mishima, Japan.
923  Ohta, T., 1992 The Nearly Neutral Theory of Molecular Evolution. Annual Review
924       of Ecology and Systematics 23: 263-286.
925  Ohta, T., and J. H. Gillespie, 1996 Development of neutral and nearly neutral
926       theories. Theoretical Population Biology 49: 128-142.
927  Pavlidis, P., and N. Alachiotis, 2017 A survey of methods and tools to detect
928       recent and strong positive selection. Journal of Biological Research-
929       Thessaloniki 24.

930 Posada, D., and T. R. Buckley, 2004 Model selection and model averaging in phylogenetics: Advantages of akaike information criterion and Bayesian approaches over likelihood ratio tests. Systematic Biology 53**:** 793-808.

933 R Core Team, 2018 R: A language and environment for statistical computing. R Foundation for Statistical Computing, pp. R Foundation for Statistical Computing, Vienna, Austria.

936 Rousselle, M., M. Mollion, B. Nabholz, T. Bataillon, and N. Galtier, 2018 Overestimation of the adaptive substitution rate in fluctuating populations. *Biology Letters*, *14*(5).

939 Sella, G., D.A. Petrov, M. Przeworski and P. Andolfatto 2009 Pervasive natural selection in the Drosophila genome? PLoS genetics 5: e1000495

941 Stanley, C. E., and R. J. Kulathinal, 2016 Genomic signatures of domestication on neurogenetic genes in Drosophila melanogaster. Bmc Evolutionary Biology 16.

944 Tajima, F., 1989 Statistical-Method for Testing the Neutral Mutation Hypothesis by DNA Polymorphism. Genetics 123**:** 585-595.

946 Tataru, P., M. Mollion, S. Glemin and T. Bataillon, 2017 Inference of Distribution of Fitness Effects and Proportion of Adaptive Substitutions from Polymorphism Data. Genetics 207**:** 1103-1119.

949 Torres, R., M.G. Stetter, R.D. Hernandez and J. Ross-Ibarra, 2019 The temporal dynamics of background selection in non-equilibrium populations. BioRxiv **doi:** https://doi.org/10.1101/618389

952 Vigué L, Eyre-Walker A. 2019 The comparative population genetics of *Neisseria meningitidis* and *Neisseria gonorrhoeae.* PeerJ 7:e7216

954 Walsh, B., and M. Lynch, 2018 *Evolution and Selection of Quantitative Traits*. Oxford University Press     .

956 Weissman, D. B., and N. H. Barton, 2012 Limits to the Rate of Adaptive Substitution in Sexual Populations. Plos Genetics 8.

958 Welch, J. J., A. Eyre-Walker and D. Waxman, 2008 Divergence and Polymorphism Under the Nearly Neutral Theory of Molecular Evolution. Journal of Molecular Evolution 67**:** 418-426.

961 Zeng, K. 2013. A coalescent model of background selection with recombination, demography and variation in selection coefficients. Heredity 110: 363-371

964

965

966

967

968    **APPENDIX**

969

970    In a constant population with population size $N_e$, $\pi_S = 4N_e\mu$ and $\pi_N$ is given by

971    (Sawyer and Hartl 1992):

972    $$\pi_N = 2N_e\mu \int_0^1 2x(1-x)H(S,x)dx \quad (A1)$$

973    where

974    $$H(S,x) = \frac{1-e^{-S(1-x)}}{x(1-x)(1-e^{-S})} \qquad (A2)$$

975    is the mean time a new semidominant mutation of scaled selection coefficient $S =$

976    $4N_e s$ spends between $x$ and $x + dx$ (Wright 1938). For constant selection $S$, by

977    integrating (A1) and dividing by $4N_e\mu$, we have:

978    $$\frac{\pi_N}{\pi_S} = f(S) = \frac{2}{1-e^{-S}} - \frac{2}{S} \qquad (A3)$$

979    (A3) is valid for both positive and negative fitness effect. If we consider only

980    beneficial mutations with a gamma distribution of effects, with mean $S_b$ and

981    shape $\beta_b$: $\phi(S_b, \beta, S) = e^{-\frac{S\beta_b}{S_b}} S^{\beta-1} \left(\frac{\beta_b}{S_b}\right)^{\beta_b} /\Gamma(\beta_b)$, we can use the same approach

982    as Welch et al. (2008) to show that:

$$\frac{\pi_N}{\pi_S} = \int_0^\infty f(S)\phi(S_b, \beta_b, S)\, dS$$

983    $$= \frac{1}{\beta_b-1}\left(\frac{\beta_b}{S_b}\right)^{\beta_b}\left(\xi\left(\beta_b-1,\frac{\beta_b}{S_b}+1\right) + (\beta_b-1)\xi\left(\beta_b,\frac{\beta_b}{S_b}\right) - \xi\left(\beta_b-1,\frac{\beta_b}{S_b}\right)\right) \quad (A4)$$

984    where $\xi(x,y)$ is the Hurwith Zeta function. (A4) can be approximated under the

985    realistic assumption that $\frac{\beta_b}{S_b} \ll 1$ and taking Taylor expansion of (A4) in $\frac{\beta_b}{S_b}$

986    around 0. We thus obtain:

987    $$\frac{\pi_N}{\pi_S} \approx (2\pi)^{\beta_b}\left(\frac{S_b}{\beta_b}\right)^{\beta_b} \quad (A5)$$

988    which leads to equation [eq. 2] in the main text.