1	The impact of purifying and background selection on the inference of population history:
2	problems and prospects
3	
4	Parul Johri ^{*, ‡} , Kellen Riall [*] , Hannes Becher [†] , Brian Charlesworth [†] , and Jeffrey D. Jensen ^{*, ‡}
5	
6	*School of Life Sciences, Arizona State University, Tempe, AZ 85287, USA
7	[†] Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, EH9
8	3FL, United Kingdom
9	[‡] Correspondence can be addressed to these authors
10	
11	Corresponding authors:
12	Parul Johri
13	pjohri1@asu.edu
14	
15	Jeffrey D. Jensen
16	Jeffrey.D.Jensen@asu.edu
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	

32 ABSTRACT

33 Current procedures for inferring population history are generally performed under the 34 assumption of complete neutrality - that is, by neglecting both direct selection and the effects of selection on linked sites. We here examine how the presence of direct purifying and background 35 36 selection may bias demographic inference by evaluating two commonly-used methods (MSMC 37 and *fastsimcoal2*), specifically studying how the underlying shape of the distribution of fitness 38 effects (DFE) and the fraction of directly selected sites interact with demographic parameter 39 estimation. The results show that, even after masking functional genomic regions, background 40 selection effects may result in the mis-inference of population growth under models of both 41 constant population size as well as decline. This effect is amplified as the strength of purifying 42 selection and the density of directly selected sites increases, as indicated by the distortion of the 43 site frequency spectrum and levels of nucleotide diversity at linked neutral sites. We also show 44 how simulated changes in background selection effects caused by population size changes can be 45 predicted analytically. We propose a potential method for correcting for the mis-inference of 46 population growth caused by selection. By treating the DFE as a nuisance parameter and 47 averaging across all potential realizations, we demonstrate that even directly selected sites may 48 be used to infer demographic histories with reasonable accuracy. 49 50 Keywords: demographic inference, background selection, distribution of fitness effects, MSMC, 51 *fastsimcoal2*, approximate Bayesian computation (ABC) 52 53 Running title: Demographic inference with selection 54 55 56 57 58 59 60 61 62

63 INTRODUCTION

64 The characterization of past population size change is a central goal of population genomic 65 analysis - with applications ranging from anthropological to agricultural to clinical (see review by Beichman *et al.* 2018). Furthermore, use of an appropriate demographic model provides a 66 67 necessary null model for assessing the impact of selection across the genome (e.g., Teshima et al. 2006; Thornton and Jensen 2007; Jensen et al. 2019). Multiple strategies have been proposed 68 69 for performing demographic inference, utilizing expectations related to levels of variation, the 70 site frequency spectrum, linkage disequilibrium, and within- and between-population relatedness 71 (e.g., Gutenkunst et al. 2009; Li and Durbin 2011; Lukic and Hey 2012; Harris and Nielsen 72 2013; Excoffier et al. 2013; Bhaskar et al. 2015; Sheehan and Song 2016; Ragsdale and 73 Gutenkunst 2017; Steinrücken et al. 2019; Kelleher et al. 2019; Speidel et al. 2019). 74 Although many methods perform well when evaluated under the standard assumption of 75 neutrality, it is difficult in practice to assure that the nucleotide sites used in empirical analyses experience neither direct selection nor the effects of selection at linked sites. For example, 76 77 inference is often performed using intergenic, 4-fold degenerate, or intronic sites. While there is 78 evidence for weak direct selection in all of these categories in multiple organisms (e.g., Haddrill 79 et al. 2005; Chamary and Hurst 2005; Andolfatto 2005; Lynch 2007; Zeng and Charlesworth 80 2010; Choi and Aquadro 2016; Jackson et al. 2017), it is also clear that such sites near or in 81 coding regions will also experience background selection (BGS; Charlesworth et al. 1993; 82 Charlesworth 2013), and may periodically be affected by selective sweeps as well (Messer and 83 Petrov 2013; Schrider et al. 2016). These effects are known to affect the local underlying 84 effective population size, and alter both the levels and patterns of variation and linkage 85 disequilibrium (Charlesworth et al. 1993; Kaiser and Charlesworth 2009; O'Fallon et al. 2010; 86 Charlesworth 2013; Nicolaisen and Desai 2013; Ewing and Jensen 2016; Johri et al. 2020). 87 However, commonly-used approaches for performing demographic inference that assume 88 complete neutrality, including fastsimcoal2 (Excoffier et al. 2013) and MSMC/PSMC (Li and 89 Durbin 2011; Schiffels and Durbin 2014), have yet to be thoroughly evaluated in the light of this 90 assumption, which is likely to be violated in practice. There are, however, some exceptions. 91 Rather than investigating existing software, Ewing and Jensen (2016) implemented an 92 approximate Bayesian (ABC) approach to quantify the impact of BGS effects, demonstrating 93 that weak purifying selection can generate a skew towards rare alleles that would be mis94 interpreted as population growth. Under certain scenarios, this resulted in a many-fold mis-

95 inference of population size change. However, the effects of the density of directly selected sites

96 and the shape of the distribution of fitness effects (DFE), which are probably of great

97 importance, have yet to be considered. Spanning the range of these potential parameter values is

98 important for understanding the implications for empirical application. For example, the

99 proportion of the genome experiencing direct purifying selection can vary greatly between

100 species, with estimates ranging from \sim 3-8% in humans, to \sim 12% in rice, to 37-53% in *D*.

101 melanogaster, to 47-68% in S. cerevisiae (Siepel et al. 2005; Liang et al. 2018). Further, many

102 organisms have highly compact genomes, with ~88% of the *E. coli* genome (Blattner et al. 1997)

and effectively all of many virus genomes, being functional (*e.g.*, >95% of the SARS-CoV-2

104 genome, Wu *et al.* 2020).

105 While such estimates allow us to approximate the effects of BGS in some model 106 organisms, in which recombination and mutation rates are well known, it is difficult to predict 107 these effects in the vast majority of study systems. Moreover, while the genome-wide mean of B, 108 a widely-used measure of BGS effects that measures the level of variability relative to neutral 109 expectation, can range from ~0.45 in *D. melanogaster* to ~0.94 in humans (Charlesworth 2013; 110 but see Pouyet et al. 2018), existing demographic inference approaches are usually applied 111 across organisms without considering this important source of differences in levels of bias. Here, we examine the effects of the DFE shape and functional density on two common demographic 112 113 inference approaches - the multiple sequentially Markovian coalescent (MSMC) and 114 *fastsimcoal2*. Finally, we propose an extension within the approximate Bayesian computation 115 (ABC) framework to address this issue, treating the DFE as a nuisance parameter and 116 demonstrating greatly improved demographic inference even when using directly selected sites 117 alone.

118

119

120 **RESULTS and DISCUSSION**

121

122 Effects of SNP numbers and genome size on inference under neutral equilibrium

123 The accuracy and performance of demographic inference was evaluated using two popular

methods, MSMC (Schiffels and Durbin 2014) and *fastsimcoal2* (Excoffier et al. 2013). In order

125 to assess performance, it was first necessary to determine how much genomic information is 126 required to make accurate inference when the assumptions of neutrality are met. Chromosomal 127 segments of varying sizes (10 Mb, 50 Mb, 200 Mb, and 1 Gb) were simulated under neutrality 128 and demographic equilibrium (*i.e.*, constant population size of 5000 diploid individuals) with 129 100 independent replicates each. For each replicate this amounted to the mean [SD] number of 130 segregating sites for each diploid individual being 1,944 [283], 9,996 [418], 40,046 [957] and 131 200,245 [1887]; for 50 diploid individuals, these values were 10,354 [225], 51,863 [567], 132 207,118 [1139] and 1,035,393 [2476] for 10 Mb, 50 Mb, 200 Mb and 1 Gb, respectively. Use of 133 MSMC resulted in incorrect inferences for all segments smaller than 1 Gb (Supp Figure 1). 134 Specifically, very strong recent growth was inferred instead of demographic equilibrium, 135 although ancestral population sizes were correctly estimated. In addition, when 2 or 4 diploid 136 genomes were used for inference, MSMC again inferred a recent many-fold growth for all 137 segment sizes even when the true model was equilibrium, but performed well when using 1 138 diploid genome for large segments (Supp Figure 1).

139 When using *fastsimcoal2* to perform demographic inference, parameters were accurately 140 estimated for all chromosomal segment sizes when the correct model (*i.e.*, equilibrium) was 141 specified (left panel of Supp Figure 2. However, when model selection was performed using a 142 choice of three models (equilibrium, instantaneous size change, and exponential size change), a 143 minority of replicates of the 10 Mb chromosomal segment incorrectly inferred a size change, 144 although the current population sizes were estimated accurately (middle panel of Supp Figure 2). 145 When the instantaneous bottleneck model was added to the set of competing models, a similarly 146 incorrect model choice was observed at segment sizes less than or equal to 50 Mb (right panel of 147 Supp Figure 2). These results suggest caution when performing inference on smaller regions or 148 genomes, specifically when the number of SNPs are less than ~200,000 per single diploid 149 individual, or less than 1 million per 50 diploid individuals. Extra caution should be used when 150 interpreting population size changes inferred by MSMC when using more than 1 diploid 151 individual.

Given this performance, all further analyses were restricted to characterizing
demographic inference on data that roughly matched the structure and size of the human genome
for every diploid individual, 22 chromosomes of size 150 Mb each were simulated, which
amounted to roughly 3 Gb of total sequence. Ten independent replicates of each parameter

combination were performed throughout, and inference utilized 1 and 50 diploid individuals forMSMC and *fastsimcoal2*, respectively.

158

159 Effect of the strength of purifying selection on demographic inference

160 In order to test demographic inference in the presence of BGS, all 22 chromosomes were

- simulated with exons of size 350 bp each, with varying sizes of introns and intergenic regions
- 162 (see Methods) being used in order to vary the fraction (5%, 10% and 20%) of the genome under
- 163 selection. Because the strength of selection acting on deleterious mutations affects how far the
- 164 effects of BGS extend, demographic inference was evaluated for various DFEs (Table 1). The
- 165 DFE was modelled as a discrete distribution with four fixed classes: $0 \le |2N_e s| < 1$, $1 \le |2N_e s| \le 1$

166 10, $10 < |2N_{es}| \le 100$ and $100 < |2N_{es}| \le 2N_{e}$. The fitness effects of mutations were uniformly

167 distributed within each bin and the DFE shape was altered by varying the proportion of

168 mutations belonging to each fixed class, given by f_0 , f_1 , f_2 , and f_3 , respectively (see Methods).

169 Three DFEs highly skewed towards a particular class were initially used to assess the impact of

- 170 the strength of selection on demographic inference (with the remaining mutations equally
- 171 distributed amongst the other three classes):

172 DFE1: a DFE in which 70% of mutations have weakly deleterious fitness effects (*i.e.*, $f_1 = 0.7$),

173 DFE2: a DFE in which 70% of mutations have moderately deleterious fitness effects (*i.e.*, $f_2 =$

174 0.7), and

175 DFE3: a DFE in which 70% of mutations have strongly deleterious fitness effects (*i.e.*, $f_3 = 0.7$).

176 In order to understand the effects of BGS, exonic sites were masked, and only linked

177 neutral intergenic sites were used for demographic inference. The three demographic models

- 178 examined were (1) demographic equilibrium, (2) a 30-fold exponential growth, mimicking the
- 179 recent growth experienced by European human populations, and (3) ~6-fold instantaneous
- 180 decline, mimicking the out-of-Africa bottleneck in human populations (Figure 1a). Although
- 181 these models were parameterized using previous estimates of human demographic history (Supp
- 182 Table 1; Gutenkunst *et al.* 2009), these basic demographic scenarios are applicable to many
- 183 organisms, although the magnitudes of population size changes in this case may represent an
- 184 extreme.
- 185 Under demographic equilibrium, when 20% of the genome experiences direct selection
 186 (with masking of the directly selected sites), we found the true population size to be

187 underestimated, as well recent population growth mis-inferred (Figure 1), with a larger bias in 188 MSMC than *fastsimcoal2*. Stronger growth is inferred with larger fitness effects of mutations. 189 Conversely, when the true demographic model is characterized by recent 30-fold growth, 190 demographic inference is accurate and performs equally well for both MSMC and *fastsimcoal2*, 191 with the exception of a slight underestimation of the ancestral population size for all DFE types. 192 When the true model is population decline, weakly deleterious mutations alone did not affect 193 inference drastically and it was possible to recover the true model. However, moderately and 194 strongly deleterious mutations resulted in an underestimation of population size and an inference 195 of strong growth, to the extent that population decline would be misinterpreted as growth. We 196 further tested the effect of BGS on demographic inference when changes in population size were 197 less severe, namely, when population growth and decline was only 2-fold, with qualitatively 198 similar results (Supp Figure 3).

199 Finally, given the strong evidence that most organisms have a bi-modal DFE with a 200 significant proportion of strongly deleterious or lethal mutations (Sanjuán 2010; Jacquier et al. 201 2013; Kousathanas and Keightley 2013; Bank et al. 2014; Charlesworth 2015), we investigated 202 the effect of this strongly deleterious class further. Thus, for comparison with the above, we 203 simulated a rather extreme case in which 30% or 50% of exonic mutations are strongly deleterious with fitness effects uniformly sampled between $100 \le 2N_{anc}s \le 2N_{anc}$, with the 204 205 remaining mutations being neutral (*i.e.*, DFE5 and DFE6; see Table 1). As with the above 206 results, both equilibrium and decline models were falsely inferred as growth, with an order of 207 magnitude underestimation of the true population size (Figure 2).

- In sum, neglecting BGS frequently results in the inference of population growth, almost
 regardless of the true underlying demographic model.
- 210

211 Effects of density and inclusion/exclusion of directly selected sites on inference

212 Although we have shown that the presence of purifying selection biases demographic inference,

213 the extent of mis-inference necessarily depends on the fraction of the genome experiencing direct

selection. We therefore compared models in which 5%, 10% or 20% of the genome was

215 functional. For this comparison, equal proportions of mutations in each DFE bin were assumed

- 216 (*i.e.*, $f_0 = f_1 = f_2 = f_3 = 0.25$) corresponding to DFE4 (Table 1). As before, when the true model
- 217 was exponential growth, inference was unbiased, with a slight underestimation of ancestral

population size when 20% of the genome experiences selection (Figure 3). Population decline was inferred reasonably well if less than 10% of the genome experiences direct selection, but could be mis-inferred as growth with greater functional density, as shown in Figure 3. Similarly, the extent to which population size is under-estimated at demographic equilibrium increases with the fraction of the genome under selection. Finally, it is noteworthy that many changes in population size that were falsely inferred were greater than 2-fold in size, suggesting the need for great caution when inferring such changes from real data.

- Importantly, the results presented do not significantly differ between inference performed while including directly selected sites (*i.e.*, no masking of functional regions; Supp Figure 4) versus inference performed using linked neutral sites (*i.e.*, masking functional regions; Figures 1-3). These results suggest that the exclusion of exonic sites, which is often assumed to provide a sufficiently neutral dataset to enable accurate demographic inference, is not necessarily a satisfactory solution unless gene density is low.
- 231

232 Effect of heterogeneity in recombination rates, mutation rates, and repeat masking

233 Variation in recombination and mutation rates, as well as the masking of repeat regions, may 234 also affect demographic inference procedures. We evaluated this issue by simulating 235 heterogeneity in both mutation and recombination rates (based on estimated human genome 236 maps, as described in the Methods section), and masking 10% of each simulated segment 237 drawing from the empirical distribution of repeat lengths in the human genome (Supp Figure 5). 238 In general, inferences under neutrality (Supp Figures 6-8) as well as under BGS (Supp Figures 9-239 11) were unaffected under all demographic models. These investigations suggest that such 240 demographic inference is robust to genome-wide variation in rates of recombination and 241 mutation, as well as the presence of repeat elements. Thus, serious mis-inference is more likely 242 to be caused by selection. These observations also suggest that simulations performed with mean 243 rates of recombination and mutation, as performed in this study, are sufficient to evaluate biases 244 caused by BGS.

245

246

Effects of BGS on diversity and the SFS under various demographic models: theoretical expectations versus simulation results

250 To better understand how BGS can lead to different biases in the inference of population history, 251 we investigated the extent of BGS effects under all three demographic models, with respect to 252 both the reduced diversity relative to neutrality, as well as the shape of the SFS at linked neutral 253 sites. First, we found that B, the diversity relative to the neutral expectation, differed among 254 demographic scenarios, disproportionately amplifying mis-inference under equilibrium and 255 decline (Figure 4). After a population decline, B was lower than that before the size change; 256 while after population expansion, B increased relative to that in the ancestral population, 257 sometimes approaching 1 (Figure 4). This may seem paradoxical, given that the magnitude of the 258 scaled selection coefficient (2 N_es) decreases with decreasing N_e (*i.e.*, the efficacy of purifying 259 selection decreases, and could thus be expected to result in larger values of B under population 260 decline). Conversely, with increasing N_e , B may be reduced.

261 However, these expectations apply only once a population has maintained a given N_e for 262 sufficient time such that equilibrium has been approached. During the initial stages of population 263 size change, and shortly afterwards, the dynamics of B tend to show a trend opposite to this long-264 term expectation (see also Figure 5 of Torres et al. 2020). This is because differences in Ne 265 caused by different initial levels of BGS cause differences in the rates of response to changes in population size – a small value of N_e (corresponding to low B) results in a faster response 266 267 compared with a high value (Fay and Wu 1999; Hey and Harris 1999; Pool and Nielsen 2007; 268 Pool and Nielsen 2009; Campos et al. 2014; Torres et al. 2020). The relative diversity values 269 observed with different initial equilibrium *B* values after a short period of population size change 270 may thus be very different from both the initial and final equilibrium values, so that the apparent 271 *B* values estimated by comparing diversities with and without BGS differ from the equilibrium 272 values. The overall effect is that there is an apparent increase in *B* immediately following a 273 population decline, and a decrease immediately following an expansion. Analytical models 274 describing these effects are presented in the Appendix. These models used the simulated values 275 of B at equilibrium before the population size changes to predict the apparent B values at the 276 ends of the periods of size change (see the Methods and Appendix). It can be seen from Figure 4 277 that there is good agreement between these predictions and the simulation results.

278 Because several demographic estimation methods are based on fitting a demographic 279 model to the SFS, it is also of interest to understand whether BGS can skew the SFS to different 280 extents under different demographic models. Although it is well known that BGS causes a skew 281 of the SFS towards rare variants under equilibrium models (Charlesworth et al. 1995; Nicolaisen 282 and Desai 2013), the effect of BGS on the SFS with population size change has not been much 283 explored (but see Johri et al. 2020 and Torres et al. 2020). As shown in Figure 5, with a 284 population size decline, the SFS of derived alleles is more skewed towards rare variants when 285 BGS is operating, especially when B is initially small, since the effects of BGS work in 286 opposition to the effects of the population size reduction. This difference in the left skew of the 287 SFS with and without BGS is much less noticeable in the case of population expansion, since 288 here the effects of BGS and the expansion act in a similar direction.

289 As with the estimates of the apparent B values discussed above, analytical predictions of 290 expected SFS after an instantaneous/ exponential change in population size can be made, using 291 the values of B at equilibrium before the population size change (see the Methods), using the 292 formulae of Polanski and Kimmel (2003) and Polanski et al. (2003) for the purely neutral case, 293 as described in the Methods section. Although the shape of the SFS is affected by both 294 demography and BGS, the impact of BGS is often comparatively small. Figure 5 shows that the 295 overall shape of the SFS is predicted reasonably well by the analytical results, although 296 deviations are to be expected for the rare allele classes, which are the most sensitive to 297 demographic change and selection. Overall, the results imply that BGS is more likely to bias 298 demographic inference post-decline compared with post-increase, consistent with the 299 performance of the methods described above. However, it should be noted that the exact patterns 300 observed will depend on the timing of population size changes relative to the time of sampling, 301 as well as the value of B before the size change. For example, with a past step increase in 302 population size, a genomic region with a sufficiently low *B* may not show signs of an expansion, 303 because all coalescent events will have been completed before the time of expansion. The 304 patterns described here thus represent only a small subset of the possibilities. In addition, all else 305 being equal, in populations with large long-term population sizes (and thus low values of *B*), 306 BGS would be expected to result in even larger biases than those observed here. 307

308

309 A potential solution: averaging across all possible DFEs

310 As shown above, demographic inference can be strongly affected by currently unaccounted for 311 BGS effects, as well as by direct purifying selection. A potential solution is thus to correct for 312 these effects when performing inference of population history. A widely-used approach to 313 estimating direct selection effects, DFE-alpha, takes a stepwise approach to inferring 314 demography, by using a presumed neutral class (synonymous sites); conditional on that 315 demography, it then estimates the parameters of the DFE (Keightley and Eyre-Walker 2007; 316 Eyre-Walker and Keightley 2009; Schneider et al. 2011; Keightley and Jackson 2018). However, 317 this approach does not include the possibility of effects of selection at linked sites, which can 318 result an over-estimate of population growth, as well as general mis-inference of the DFE (Johri

319 et al. 2020).

320 Building on this idea, Johri et al. (2020) recently proposed an approach that includes both 321 direct and background effects of purifying selection, simultaneously inferring the deleterious 322 DFE and demography. By utilizing the decay of BGS effects around functional regions, they 323 demonstrated high accuracy under the simple demographic models examined. Moreover, the 324 method makes no assumptions about the neutrality of synonymous sites, and can thus be used to 325 estimate selection acting on these sites, as well as in non-coding functional elements. However, 326 this computationally-intensive approach is specifically concerned with jointly inferring the DFE 327 and demographic parameters. As such, if an unbiased characterization of the population history 328 is the sole aim, this procedure may be needlessly involved. We thus here examine the possibility 329 of rather treating the DFE as an unknown nuisance parameter, averaging across all possible DFE 330 shapes, in order to assess whether demographic inference may be improved simply by correcting 331 for these selection effects without inferring their underlying parameter values. This approach 332 utilizes functional (*i.e.*, directly selected) regions, a potential advantage in populations for which 333 only coding data may be available (e.g., exome-capture data; see Jones and Good 2016), or more 334 generally in organisms with high gene densities.

In order to illustrate this approach, a functional genomic element was simulated under demographic equilibrium, 2-fold exponential population growth and 2-fold exponential population decline with four different DFE shapes (as described previously, and shown in Figure 6). A number of summary statistics were calculated (see Methods) for the entire functional region. Inference was first performed assuming strict neutrality, and inferring a one-epoch size

340 change (thus estimating the ancestral (N_{anc}) and current population sizes (N_{cur})). As was found

341 with the other inference approaches examined, population sizes were underestimated and a false

342 inference of population growth was observed in almost all cases when selective effects are

343 ignored (Figure 6).

344 Next, the assumption of neutrality was relaxed, and mutations were simulated with fitness 345 effects characterized by a discrete DFE, with the same fitness classes given above (f_0, f_1, f_2, f_3) . Values for f were drawn from a uniform prior between 0 and 1, such that $\sum f = 1$. Note that no 346 347 assumptions were made about which sites in the genomic region were functionally important, or 348 regarding the presence/absence of a neutral class. These directly selected sites were then used to 349 infer demographic parameters. We found that, by varying the shape of the DFE, averaging across 350 all realizations, and estimating only parameters related to population history, highly accurate 351 inference of modern and ancestral population sizes is possible (Figure 6). These results 352 demonstrate that, even if the true DFE of a population is unknown (as will always be the case in 353 reality), it is possible to infer demographic history with reasonable accuracy by approximately 354 correcting for these selective effects.

355

356

357 CONCLUSION

358 While commonly used approaches for inferring demography assume neutrality and independence 359 among segregating sites, it is very difficult to verify those assumptions empirically. In addition, 360 there is considerable evidence for wide-spread effects of selection on linked sites in many 361 commonly studied organisms (Hernandez et al. 2011; Cutter and Payseur 2013; Williamson et al. 362 2014; Elyashiv et al. 2016; Campos et al. 2017; Booker and Keightley 2018; Pouyet et al. 2018; 363 Torres et al. 2018; Castellano et al. 2020). As such, we explored how violations of the 364 assumption of neutrality may affect demographic inference, particularly with regard to the 365 underlying strength of purifying selection and the genomic density of directly selected sites. 366 Generally speaking, the neglect of these effects (*i.e.*, background selection) results in an 367 inference of population growth, with the severity of the growth model roughly scaling with 368 selection strength and density. Thus, when the true underlying model is in fact growth, 369 demographic mis-inference is not particularly severe; when the true underlying model is constant

370 size or decline, the mis-inference can be extreme, with a many-fold underestimation of

371 population size.

372 It is important to note that BGS effects extend to genomic distances in a way that is 373 positively related to the strength of purifying selection. For instance, strongly and moderately 374 deleterious mutations affect patterns of diversity at large genomic distances, whereas mildly 375 deleterious mutations primarily skew allele frequencies at adjacent sites. Thus, if intergenic 376 regions further away from exons are used to perform demographic inference, it is predominantly 377 strongly deleterious mutations that are likely to bias inferences. In contrast, if synonymous sites 378 are used to infer demographic history, mildly deleterious mutations will be most important. Thus, 379 as we here focused on relatively sparsely-coding genomes (resembling human-like gene 380 densities) and used intergenic sites for inference, moderately and strongly deleterious mutations 381 resulted in more severe mis-inference. In addition, the effect of decay of BGS due to mildly 382 deleterious mutations depends on multiple parameters. For instance, with an exon of length 500 383 bp and *Drosophila*-like parameters ($N_e = 10^6$; recombination rate = mutation rate = 10^{-8} / site / 384 generation), B will increase from 0.53 (at 10 bases from the end of the exon) to 0.94 at a distance 385 of 1000 bases. On the other hand, with human-like parameters ($N_e = 10^4$; recombination rate = 386 mutation rate = 10^{-8} / site / generation),) the corresponding change in B is only 0.981 to 0.982 387 (Supp Table 2). Thus, mildly deleterious mutations will have drastically different effects in this 388 regard depending on the underlying population parameters. Hence, a large impact of selection is 389 expected when such approaches are applied to synonymous sites in smaller and more compact 390 genomes with higher gene densities and a larger input of deleterious mutations at closely linked 391 sites.

392 Comparing the two inferences methods investigated here, it appears that *fastsimcoal2* is 393 less prone to inferring false fluctuations in population size caused by BGS. However, both 394 methods falsely infer growth when it is absent, with increasing severity as the density of coding 395 regions increases. The times of population growth inferred by both methods appear to be affected 396 in unpredictable ways. Further, we observed little difference in performance when using all 397 SNPs relative to when SNPs are thinned to be separated by 5kb (Supp Figure 12), presumably 398 because the models investigated do not generate strong LD. Overall, the degree of mis-inference 399 caused by a neglect of BGS is largely similar between the two methods.

400 However, it is noteworthy that even when all sites are strictly neutral, or only 5% of the 401 genome experiences direct selection, demographic equilibrium is mis-estimated by MSMC as a 402 series of size changes. The pattern of these erroneous size changes lend a characteristic shape to 403 the MSMC curve (*i.e.*, ancient decline and recent growth) which appears to resemble the 404 demographic history previously inferred for the Yoruban population (Schiffels and Durbin 405 2014), including the time at which changes in population size occurred (Supp Figure 13). 406 Previous work has demonstrated that the resulting demographic model does not in fact fit the 407 observed SFS in the Yoruban population (Beichman et al. 2017; Lapierre et al. 2017). A similar 408 shape has also been inferred in the vervet subspecies (Warren et al. 2015; Figure 4), in passenger 409 pigeons (Hung et al. 2014; Figure 2), in elephants (Palkopoulou et al. 2018; Figure 4), in 410 Arabidopsis (Fulgione et al. 2018; Figure 3), and in grapevines (Zhou et al. 2017; Figure 2A). 411 Although the inferred population size fluctuations under simulated neutrality are only 412 \sim 1.2-fold, in most empirical applications the fluctuations are of a somewhat larger magnitude (\sim 413 2-fold in pigeons, Arabidopsis, and grapevines). Nonetheless, this neutral performance of 414 MSMC under demographic equilibrium is concerning, and adds to the other previously published 415 cautions concerning the interpretation of MSMC results. For example, Mazet et al. (2016) and 416 Chikhi et al. (2018) demonstrated that under constant population size with hidden structure, 417 inference may suggest false size change (see also Orozco-terWengel 2016). In addition, MSMC 418 has been reported to falsely infer growth prior to instantaneous bottlenecks (Bunnefeld et al. 419 2015). In addition, we observed that if insufficient genomic data is utilized, or more than one 420 diploid genome is used to perform inference, MSMC falsely infers recent growth of varying 421 magnitudes (as previously observed by Beichman et al. 2017).

422 In sum, we find that the effects of purifying and background selection result in similar 423 demographic mis-inference across approaches, and that masking functional sites does not yield 424 accurate parameter estimates. In order to side-step many of these difficulties, our proposed 425 approach of inferring demography by averaging selection effects across all possible DFE shapes 426 within an ABC framework appears to be promising. Utilizing only functional regions, we found 427 a great improvement in accuracy, without making any assumptions regarding the true underlying 428 shape of the DFE or the neutrality of particular classes of sites. As such, this approach represents 429 a more computationally efficient avenue if only demographic parameters are of interest, and

430 ought to be particularly useful in the great majority of organisms in which unlinked neutral sites431 either do not exist, or are difficult to identify and verify.

- 432
- 433

434 METHODS

435 Simulations: Simulations were performed using SLiM 3.1 (Haller and Messer 2019) with 10 436 replicates per evolutionary scenario. For every replicate, 22 chromosomes of 150Mb each were 437 simulated, totaling ~ 3 Gb of information per individual genome (similar to the amount of 438 information in a human genome). Within each chromosome, 3 different types of regions were 439 simulated, representing non-coding intergenic, intronic, and exonic regions. Based on the NCBI 440 RefSeq human genome annotation, downloaded from the UCSC genome browser for hg19 (http://genome.ucsc.edu/; Kent et al. 2002), mean values of exon sizes and intron numbers per 441 442 gene were calculated. To represent mean values for the human genome (Lander et al. 2001), each 443 gene was comprised of 8 exons and 7 introns, and exon lengths were fixed at 350 bp. By varying 444 the lengths of the intergenic and intronic regions, three different genomic configurations with 445 varying densities of functional elements were simulated and compared - with 5%, 10% and 20% 446 of the genome being under direct selection - hereafter referred to as genome5, genome10, and 447 genome20, respectively. Genome5 was comprised of introns of 3000 bp and intergenic sequence 448 of 31000 bp, genome10 of introns of 1500 bp and intergenic sequence of 15750 bp, while 449 genome20 was comprised of introns of 600 bp and intergenic sequence of 6300 bp. The total 450 chromosome sizes of these genomes were approximately 150 Mb (150,018,599 bp, 150,029,949 451 bp, 150,003,699 bp in genome5, 10, and 20, respectively) with 2737, 5164, and 11278 genes per 452 chromosome in genome5, 10 and 20, respectively. In order to be conservative with respect to the 453 performance of existing demographic estimators, intronic and intergenic regions were assumed 454 to be neutral.

455 Recombination and mutation rates were assumed to be equal at $1 \ge 10^{-8}$ /site / generation. 456 Neither crossover interference (see the discussion in Campos and Charlesworth 2019) nor gene 457 conversion were modeled. Exonic regions in the genomes experienced direct purifying selection 458 given by a discrete DFE comprised of 4 fixed classes (Johri et al. 2020), whose frequencies are 459 denoted by *f*: *f*₀ with $0 \le 2N_e s < 1$ (*i.e.*, effectively neutral mutations), *f*₁, with $1 \le 2N_e s < 10$ 460 (*i.e.*, weakly deleterious mutations), *f*₂, with $10 \le 2N_e s < 100$ (*i.e.*, moderately deleterious

461 mutations), and f_3 , with $100 \le 2N_{\rm e}s < 2N_{\rm e}$ (*i.e.*, strongly deleterious mutations), where $N_{\rm e}$ is the 462 effective population size and s is the reduction in fitness of the mutant homozygote relative to 463 wild-type. Within each bin, the distribution of s was assumed to be uniform. All mutations were 464 assumed to be semi-dominant. In all cases, the Ne corresponding to the DFE refers to the 465 ancestral effective population size. Six different types of DFE were simulated: DFE1: a DFE 466 skewed largely towards mildly deleterious mutations, given by $f_0=0.1$, $f_1=0.7$, $f_2=0.1$, $f_3=0.1$; 467 DFE2: a DFE skewed towards moderately deleterious mutations, $f_0=0.1$, $f_1=0.1$, $f_2=0.7$, $f_3=0.1$; 468 DFE3: a DFE skewed towards strongly deleterious mutations, $\hbar = 0.1$, $\hbar = 0.1$, $\hbar = 0.1$, $\hbar = 0.7$; 469 DFE4: a DFE with equal proportions of all mutations, $f_0=0.25$, $f_1=0.25$, $f_2=0.25$, $f_3=0.25$; 470 DFE5: a DFE with equal proportions of neutral and strongly deleterious mutations, $f_0 = 0.5$, 471 $f_1=0.0$, $f_2=0.0$, $f_3=0.5'$ and DFE6: a DFE with a majority of neutral mutations and a minority of 472 strongly deleterious mutations, $f_0=0.7$, $f_1=0.0$, $f_2=0.0$, $f_3=0.3$. 473 Three different demographic models were tested for each of these DFEs (Supp Table 1):

474 1) demographic equilibrium, 2) recent exponential 30-fold growth, resembling that estimated for
475 the human CEU population (Gutenkunst et al. 2009), and 3) ~6-fold instantaneous decline,

476 resembling the out-of-Africa bottleneck in humans (Gutenkunst et al. 2009).

477

478 *Running MSMC*: In order to quantify the effect of purifying selection on demographic

479 inference, we used entire chromosomes generated by SLiM to generate input files for MSMC.

480 For comparison, and to quantify the effect of BGS alone on demographic inference, we masked

the exonic regions to generate input files. For all parameters, MSMC was performed on a single

482 diploid genome, as the results for this case were the most accurate (Supp Figure 1). Input files

483 were made using the script ms2multihetsep.py provided in the msmc-tools-Repository

484 downloaded from https://github.com/stschiff/msmc-tools. MSMC1 and 2 were run as follows:

485 *msmc_1.1.0_linux64bit -t 5 -r 1.0 -o output_genomeID input_chr1.tab input_chr2.tab ...*

486 *input_chr22.tab.* Population sizes obtained from MSMC were plotted up to the maximum

487 number of generations obtained from MSMC, and the final value of the ancestral population size

488 was extended indefinitely as a dashed line.

489

490 *Running Fastsimcoal2*: In order to minimize the effects of linkage disequilibrium (LD), only
491 SNPs separated by 5 kb were used for inference, following Excoffier et al. (2013). Inference was

492 also performed by including all SNPs in order to assess the impact of violating the assumption of

- 493 no LD on inference. Site frequency spectra (SFS) were obtained for both sets of SNPs for all 10
- 494 replicates of every combination of demographic history and DFE. SNPs from all 22
- 495 chromosomes were pooled together to calculate the SFS. *Fastsimcoal2* was used to fit each SFS
- 496 to 4 distinct models: (a) equilibrium, which estimates only a single population size parameter
- 497 (*N*); (b) instantaneous size change (decline/growth), which fits 3 parameters ancestral
- 498 population size (N_{anc}), current population size (N_{cur}), and time of change (T); (c) exponential size
- 499 change (decline/growth), which also estimates 3 parameters Nanc, Ncur and T; and (d) an
- 500 instantaneous bottleneck model with 3 parameters $-N_{anc}$, intensity, and time of bottleneck. The
- 501 parameter search ranges for both ancestral and current population sizes in all cases were
- 502 specified to be uniformly distributed between 100-500000 individuals, while the parameter range
- 503 for time of change was specified to be uniform between 100-10000 generations in all models.
- The intensity of the bottleneck, as specified by the population size during the reduction, was also sampled uniformly from the range 100-500000 individuals. The following command line was used to run *fastsimcoal2*:
- 507 fsc26 -t demographic model.tpl -n 150000 -d -e demographic model.est -M -L 50 -q,
- 508 Model selection was performed as recommended by Excoffier et al. (2013). For each 509 demographic model, the maximum of maximum likelihoods from all replicates was used to 510 calculate the Akaike Information Criterion (AIC) = $2 \times$ number of parameters – $2 \times$
- 511 ln(likelihood). The relative likelihoods (Akaike's weight of evidence) in favor of the *t*^h model
 512 was then calculated by:
- 513 $w(i) = \frac{e^{-0.5\Delta_i}}{\sum_{i=1}^4 e^{-0.5\Delta_j}}$

514 where $\Delta_i = AIC_i - AIC_{min}$. The model with the highest relative likelihood was selected as the 515 best model, and the parameters estimated using that model were used to plot the final inferred 516 demography.

517

518 Simulations of variable recombination and mutation rates, and repeat masking: In order to 519 simulate variation in recombination and mutation rates, all 22 chromosomes were simulated by 520 mimicking chromosome 6 (~171Mb) of the human genome. Recombination rates obtained from 521 Yoruban populations were obtained from the UCSC genome browser, while the mutation rate 522 map (https://molgenis26.target.rug.nl/downloads/gonl public/mutation rate map/release2/) was 523 assumed to correspond to estimates obtained from *de novo* mutations (Francioli et al. 2015), as in 524 Castellano et al. (2020). Absolute values of mutation rates were normalized in order to maintain 525 the mean mutation rate across the genome at $\sim 1.0 \times 10^{-8}$ per site per generation. Recombination 526 and mutation rate estimates were taken from positions of approximately 10 Mb to 160 Mb, with 527 the recombination map starting at 10010063 bp and the mutation map starting at 10010001 bp. 528 Regions with missing data for either of the two estimates were simulated with rates 529 corresponding to the previous window, except for the case of centromeres in which no 530 recombination was assumed. In order to understand the effect of excluding centromeric regions 531 in empirical studies, the 4Mb region corresponding to the centromere was masked, 532 corresponding to 48.5 to 52.5 Mb of the simulated 150Mb chromosomes. In order to evaluate the 533 effect of masking repeat regions, random segments comprising 10% of each chromosome were 534 masked. The lengths of these segments were drawn from the lengths of repeat regions found in 535 the human genome (Supp Figure 5), as obtained from the repeat regions in the hg19 assembly of 536 the human genome from the UCSC genome browser.

537

538 **Performing inference by approximate Bayesian computation (ABC)**: ABC was performed 539 using the R package "abc" (Csilléry et al. 2010), and non-linear regression aided by a neural net 540 (used with default parameters as provided by the package) was used to correct for the 541 relationship between parameters and statistics (Johri et al. 2020). To infer posterior estimates, a 542 tolerance of 0.1 was applied (*i.e.*, 10% of the total number of simulations were accepted by ABC 543 in order to estimate the posterior probability of each parameter). The weighted medians of the 544 posterior estimates for each parameter were used as point estimates. ABC inference was 545 performed under two conditions: (1) complete neutrality, or (2) the presence of direct purifying 546 selection. In both cases only 2 parameters were inferred - ancestral (N_{anc}) and current (N_{cur}) 547 population sizes. However, in scenario 2, the shape of the DFE was also varied. Specifically, the 548 parameters f_0 , f_1 , f_2 , and f_3 were treated as nuisance parameters and were sampled such that $0 \le f_1$ 549 ≤ 1 , and $\Sigma_i f_i = 1$, for i = 0 to 3. In addition, in order to limit the computational complexity 550 involved in the ABC framework, values of f_i were restricted to multiples of 0.05 (*i.e.*, $f_i \in \{0.0, ..., ..., c_i\}$ 551 $(0.05, 0.10, \dots, 0.95, 1.0) \forall i$, which allowed us to sample 1,771 different DFE realizations. 552 Simulations were performed with functional genomic regions, and the demographic model was

553 characterized by 1-epoch changes in which the population either grows or declines exponentially 554 from ancestral to current size, beginning at a fixed time in the past.

For the purpose of illustration, and for a contrast with the human-like parameter set above, parameters for ABC testing were selected to resemble those of *D. melanogaster* African populations. Priors on ancestral and current population sizes were drawn from a uniform distribution between 10^{5} - 10^{7} diploid individuals, while the time of change was fixed at 10^{6} (~*N*_e) generations. In order to simulate functional regions, 94 single-exon genes, as described in Johri *et al.* (2020) and provided in

561 <u>https://github.com/paruljohri/BGS_Demography_DFE/blob/master/DPGP3_data.zip</u>, were

562 simulated with recombination rates specific to those exons (<u>https://petrov.stanford.edu/cgi-</u>

563 <u>bin/recombination-rates_updateR5.pl</u>) (Fiston-Lavier et al. 2010; Comeron et al. 2012). Mutation

rates were assumed to be fixed at 3×10^{-9} per site per generation (Keightley et al. 2009;

565 Keightley et al. 2014).

566 All parameters were scaled by the factor 320 in order to decrease computational time, 567 using the principle first described by Hill and Robertson (1966), and subsequently by others 568 (Comeron and Kreitman 2002; Hoggart et al. 2007; Kaiser and Charlesworth 2009; Kim and 569 Wiehe 2009; Uricchio and Hernandez 2014; Campos and Charlesworth 2019). The scaled 570 population sizes thus ranged between $\sim 300-30000$ and were reported as scaled values in the 571 main text. One thousand replicate simulations were performed for every parameter combination 572 $(N_{\text{anc}}, N_{\text{cur}}, f_0, f_1, f_2, f_3)$; for performing ABC inference, 50 diploid genomes were randomly 573 sampled without replacement, and summary statistics were calculated using pylibseq 0.2.3 574 (Thornton 2003). Means and variances (between replicates) of the following summary statistics 575 were calculated across the entire exonic region: nucleotide site diversity (π), Watterson's θ , 576 Tajima's D, Fay and Wu's H (both absolute and normalized), number of singletons, haplotype 577 diversity, LD-based statistics (r^2, D, D') , and divergence (*i.e.*, number of fixed mutations per site 578 per generation after the burn-in period). As opposed to the above examples, in this inference 579 scheme only exonic data (i.e., directly selected sites) were utilized. Test datasets were generated 580 in exactly the same fashion as described above.

581

582 Analytical expectations for the relative site frequencies: To compute the expected relative
 583 frequencies of site frequency classes, the approach of Polanski and Kimmel (2003) was

584 followed. They describe a method for computing the "probability that a SNP has b mutant 585 bases", which is equivalent to the site frequency spectrum (SFS) of derived variants. This 586 method (their equations 3-10) allows for the specification of arbitrary population size histories 587 and sample sizes. For reasons of computational precision, a sample size of 10 diploid genomes 588 was chosen. The demographic scenarios detailed in Supplementary Table 1 were implemented as 589 piecewise functions of the effective population size (counting haploid genomes), and the effect 590 of BGS was included by scaling these functions by values of B before population size change as 591 obtained from the forwards-in-time simulations described above. A Mathematica notebook 592 detailing these results, will be made available upon publication (see data availability statement). 593 In addition, analytical expressions can be obtained for pairwise diversity values when there are 594 step changes or exponential growth in population size, as described in the Appendix and in an 595 example program that calculates diversity values after exponential growth. 596 597 **Data availability:** The following data will be made publicly available upon acceptance of the

598 manuscript: (1) Scripts used to perform simulations; (2) Input files used to run *fastsimcoal2*; (3)

599 Scripts used for plotting; (4) Plotted results of MSMC and *fastsimcoal2* for all models and

600 scenarios tested in this work; (5) A Mathematica (version 12.1) notebook detailing calculations

601 of analytical expectations for the relative SFS; (6) An example program (Fortran script)

602 demonstrating how to obtain analytical expressions for values of *B* after exponential growth. All

603 supplemental files will be made publicly available at

604 <u>https://github.com/paruljohri/demographic_inference_with_selection</u>.

605

606

607 Acknowledgements

We would like to thank Susanne Pfeifer for helpful discussions related to this project, and for
 feedback on the manuscript. This research was conducted using resources provided by Research

610 Computing at Arizona State University (http://www.researchcomputing.asu.edu) and the Open

611 Science Grid, which is supported by the National Science Foundation and the U.S. Department

612 of Energy's Office of Science. This work was funded by National Institutes of Health grant

613 R01GM135899 to JDJ.

614 **APPENDIX**

615 There are two simple scenarios for population size change for which simple explicit expressions for the

- 616 expected pairwise coalescent time or diversity can be obtained, without using the methodology of
- 617 Polanski and Kimmel (2003) and Polanski *et al.* (2003) – a step change in N or an exponential growth
- 618 in N. First consider the coalescent process for a step change, where the current and initial effective
- 619 population sizes are denoted by N_{e1} and N_{e0} , respectively. Let B be the background selection parameter
- 620 at the start of the process of change, corresponding to effective size N_{e0} . For convenience, time is
- 621 scaled in units of $2N_{e1}$ generations, and the time of the change in population size on this scale is
- 622 denoted by T_0 , counting back from the present time, T = 0. T_0 is assumed to be sufficiently small that B
- 623 remains approximately constant during the period since the change in size. Denote the ratio N_{e0}/N_{e1} by
- 624 *R*. The derivation for the case of a step change in population size is similar to that given by Pool and
- 625 Nielsen (2009) for the purpose of comparing X chromosomes and autosomes.

626 Between times T and T_0 , coalescence occurs at a rate B^{-1} on the chosen timescale, so that the 627 contribution from this period to the net coalescent time for a pair of alleles sampled at T = 0 is: 628

629
$$B^{-1} \int_0^{T_0} T \exp(-B^{-1}T) \, \mathrm{d}T = B - B \exp(-B^{-1}T_0) - T_0 \exp(-B^{-1}T_0)$$

630

631 There is a probability of $\exp(-B^{-1}T_0)$ that there is no coalescence when T lies between 0 and T_0 , after which coalescence occurs at a rate 1/BR, giving a net contribution to the coalescence time of: 632 633

- $(BR + T_0) \exp(-B^{-1}T_0)$ 634
- 635

636 The net coalescence time for the stepwise change with BGS is given by the sum of these 637 two expressions:

- 638
- 639

$$B[1 + (R - 1)\exp(-B^{-1}T_0)]$$
(1a)

- 640
- 641 If this expression is compared to the corresponding equation with B = 1, the apparent value of B
- 642 at the time of sampling of the pair of alleles is given by:

bioRxiv preprint doi: https://doi.org/10.1101/2020.04.28.066365. this version posted July 6, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. It is made available under a CC-BY-ND 4.0 International license.

643

644

$$B_{s} = \frac{B[1+(R-1)\exp(-B^{-1}T_{0})]}{[1+(R-1)\exp(-T_{0})]}$$
(1b)

645

Next, consider a process of exponential change in population size, starting at an initial effective size of N_{e0} at t_0 generations in the past and ending at size N_{e1} , such that the instantaneous growth rate rper generation is r=equal to $\ln(N_{e1}/N_{e0})/t_0$. The effective population size at time t in the past is $N_e(t) =$ $N_{e1}\exp(-rt)$; with BGS, the rate of coalescence at time t is $1/BN_e(t)$. As before, the BGS parameter is assumed to remain constant over the period of population size change. It follows that the probability of no coalescence by generation t in the past is:

652

653
$$P_{nc}(t) = \exp\left[-\int_0^t (2BN_{e1})^{-1} \exp(rt) \, dt\right] = \exp[c^{-1}(1-e^{rt})]$$
(2)

- 654
- 655 where $c = 2BNe_1r$.

656 The pre-growth period with $t > t_0$ contributes an expected coalescent time of 657 $(2BN_{e0} + t_0)P_{nc}(t_0)$, on the scale of generations.

Following Slatkin and Hudson (1991), to obtain the contribution from the period with $t > t_0$, it is convenient to measure time as $\tau = rt$. The probability of coalescence between τ and $\tau + d\tau$ is then given by:

- 661
- 662

$$P_c(\tau) = c^{-1} e^{\tau} \exp[c^{-1}(1 - e^{\tau})] d\tau$$
(3)

663

664 The contribution from this period to the expected coalescent time is given by the integral 665 of $\tau P_c(\tau)$ between 0 and τ_0 . Following Slatkin and Hudson (1991), by transforming to u =666 exp(τ), this contribution can be expressed as the following integral:

- 667
- 668

669

$$\bar{\tau}_1 = c^{-1} e^{c^{-1}} \int_1^{u_0} \ln(u) e^{-uc^{-1}} du$$
(4)

670 This integral can easily be evaluated numerically. The corresponding mean coalescent 671 time on the scale of generations is obtained by division by *r*, and the result can be added to 672 $(2BN_{e0} + t_0)P_{nc}(t_0)$, yielding the net expected coalescent time. By dividing the resulting

- 673 expression by the corresponding expression with B = 1, the apparent BGS effect at the time of
- 674 sampling can be obtained, in the same way as for the step change model.
- 675
- 676

677 FIGURES AND TABLES

- 678
- 679

680 **Table 1**: Proportion (f_i) of mutations in each class of the discrete distribution of fitness effects

681 (DFE) simulated in this study.

DFEI 0.1 0.7 0.1 0.1	
DFE2 0.1 0.1 0.7 0.1	
DFE3 0.1 0.1 0.1 0.7	
DFE4 0.25 0.25 0.25 0.25	
DFE5 0.5 0.0 0.0 0.5	
DFE6 0.7 0.0 0.0 0.3	





684 Figure 1: Inference of demography by MSMC (red lines; 10 replicates) and *fastsimcoal2* (blue lines; 10 replicates) with and without BGS, under demographic equilibrium (left column), 30-685 fold exponential growth (middle column), and ~6-fold instantaneous decline (right column). The 686 687 true demographic models are depicted as black lines, with the x-axis origin representing the 688 present day. (a) All genomic sites are strictly neutral. (b) Exonic sites experience purifying 689 selection specified by a DFE comprised largely of weakly deleterious mutations (DFE1: $\hbar = 0.1$, 690 $f_1 = 0.7$, $f_2 = 0.1$, $f_3 = 0.1$). (c) Exonic sites experience purifying selection specified by a DFE 691 comprised largely of moderately deleterious mutations (DFE2: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.7$, $f_3 = 0.7$, $f_4 = 0.7$, $f_5 = 0.7$, $f_7 = 0.7$, $f_8 = 0$ 692 0.1). (d) Exonic sites experience purifying selection specified by a DFE comprised largely of 693 strongly deleterious mutations (DFE3: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.7$). Exons represent 20% 694 of the genome, and exonic sites were masked/excluded when performing demographic inference, 695 quantifying the effects of BGS alone. The dashed lines represent indefinite extensions of the 696 ancestral population sizes. 697











711 Figure 3: Inference of demography by MSMC (red lines; 10 replicates) and *fastsimcoal2* (blue

12 lines; 10 replicates) in the presence of BGS with varying proportions of the genome under

selection, for demographic equilibrium (left column), exponential growth (middle column), and instantaneous decline (right column). Exonic sites were simulated with purifying selection with

all *f* values equal to 0.25 (DFE4), and were masked when performing inference. Directly

selected sites comprise (a) 20% of the simulated genome, (b) 10% of the simulated genome, and

(c) 5% of the simulated genome. The true demographic models are given by the black lines, with

the x-axis origin representing the present day. The dashed lines represent indefinite extensions of

- 719 the ancestral population sizes.
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727



728

Figure 4: *B* is the nucleotide site diversity (π) with BGS relative to its purely neutral expectation

730 (π_0) . The results are shown for (a) demographic equilibrium, (b) population growth, and (c)

population decline. For all cases, the ancestral *B* (*i.e.*, *B* pre-change in population size) is shown

in white bars, *B* post-change in population size is shown in solid gray bars, and the analytical

radius for the post-size change *B* is shown as red bars. Various DFEs were used: **DFE1**: *f*₀

734 = 0.1, $f_1 = 0.7$, $f_2 = 0.1$, $f_3 = 0.1$. **DFE2**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.7$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.7$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$. **DFE3**: $f_0 = 0.1$, $f_1 = 0.1$.

735 = 0.1, $f_2 = 0.1$, $f_3 = 0.7$. **DFE4**: $f_0 = f_1 = f_2 = f_3 = 0.25$. **DFE5**: $f_0 = 0.5$, $f_1 = 0.0$, $f_2 = 0.0$, $f_3 = 0.0$, $f_4 = 0.0$, $f_5 = 0.0$

- 736 0.5. **DFE6:** $f_0 = 0.7$, $f_1 = 0.0$, $f_2 = 0.0$, $f_3 = 0.3$. Exonic sites comprised ~10% of the genome,
- roughly mimicking the density of the human genome.

bioRxiv preprint doi: https://doi.org/10.1101/2020.04.28.066365. this version posted July 6, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. It is made available under a CC-BY-ND 4.0 International license.



738

Figure 5: The site frequency spectrum (SFS) of derived allele frequencies at neutral sites from 10 diploid genomes under (a) demographic equilibrium, (b) population growth, and (c) population decline, under the same DFEs as shown in Figure 4. The x-axis indicates the number of sample alleles (out of 20) carrying the derived variant. Exonic sites comprised ~10% of the genome, roughly mimicking the density of the human genome. The red solid circles give the values predicted analytically with a purely neutral model, but taking the simulation values of *B* into account in order to quantify the effective population size.



749 Figure 6: Comparison of estimates of ancestral (N_{anc}) and current (N_{cur}) population sizes when 750 assuming neutrality vs when varying the DFE shape as a nuisance parameter, using an ABC framework. Inference is shown for demographic equilibrium (left column), 2-fold exponential 751 752 growth (middle column), and 2-fold population decline (right column), for five separate DFE shapes that define the extent of direct purifying selection acting on the genomic segment for 753 754 which demographic inference is performed: (a) neutrality ($f_0 = 1$, $f_1 = 0$, $f_2 = 0$, $f_3 = 0$), (b) weak purifying selection (DFE1: $f_0 = 0.1$, $f_1 = 0.7$, $f_2 = 0.1$, $f_3 = 0.1$), (c) moderately strong purifying 755 756 selection (DFE2: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.7$, $f_3 = 0.1$), (d) strong purifying selection (DFE3: $f_0 = 0.1$) 757 0.1, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.7$), and (e) a DFE in which all classes of mutations are equally frequent (DFE4: $f_0 = f_1 = f_2 = f_3 = 0.25$). In each, the horizontal lines give the true values (black 758 759 for $N_{\rm anc}$; and gray for $N_{\rm cur}$) and the box-plots give the estimated values. Black and gray boxes 760 represent estimates when assuming neutrality, while red boxes represent estimates when the DFE is treated as a nuisance parameter. 761

762 **REFERENCES**

763 764 Andolfatto P. 2005. Adaptive evolution of non-coding DNA in Drosophila. Nature 437:1149-765 1152. 766 Bank C, Ewing GB, Ferrer-Admettla A, Foll M, Jensen JD. 2014. Thinking too positive? 767 Revisiting current methods of population genetic selection inference. Trends Genet. 768 30:540-546. 769 Beichman AC, Huerta-Sanchez E, Lohmueller KE. 2018. Using genomic data to infer historic 770 population dynamics of nonmodel organisms. Annu. Rev. Ecol. Evol. Syst. 49:433-456. 771 Beichman AC, Phung TN, Lohmueller KE. 2017. Comparison of single genome and allele 772 frequency data reveals discordant demographic histories. G3 7:3605–3620. Bhaskar A, Wang YXR, Song YS. 2015. Efficient inference of population size histories and 773 774 locus-specific mutation rates from large-sample genomic variation data. Genome Res. 775 25:268-279. 776 Blattner FR, Plunkett G, Bloch CA, Perna NT, Burland V, Riley M, Collado-Vides J, Glasner 777 JD, Rode CK, Mayhew GF, et al. 1997. The complete genome sequence of Escherichia 778 coli K-12. Science 277:1453-1462. 779 Booker TR, Keightley PD. 2018. Understanding the factors that shape patterns of nucleotide 780 diversity in the house mouse genome. Mol. Biol. Evol. 35:2971-2988. 781 Bunnefeld L, Frantz LAF, Lohse K. 2015. Inferring bottlenecks from genome-wide samples of 782 short sequence blocks. Genetics 201:1157-1169. 783 Campos JL, Charlesworth B. 2019. The effects on neutral variability of recurrent selective 784 sweeps and background selection. Genetics 212:287-303. 785 Campos JL, Halligan DL, Haddrill PR, Charlesworth B. 2014. The relation between 786 recombination rate and patterns of molecular evolution and variation in Drosophila 787 melanogaster. Mol. Biol. Evol. 31:1010-1028. 788 Campos JL, Zhao L, Charlesworth B. 2017. Estimating the parameters of background selection 789 and selective sweeps in Drosophila in the presence of gene conversion. Proc. Natl. Acad. 790 Sci. U.S.A. 114:E4762-E4771. 791 Castellano D, Eyre-Walker A, Munch K. 2020. Impact of mutation rate and selection at linked 792 sites on DNA variation across the genomes of humans and other Homininae. Genome 793 Biol. Evol. 12:3550-3561. 794 Chamary J, Hurst LD. 2005. Evidence for selection on synonymous mutations affecting stability 795 of mRNA secondary structure in mammals. Genome Biol. 6:R75.

796 797	Charlesworth B. 2013. Background selection 20 years on. The Wilhelmine E. Key 2012 invitational lecture. <i>J. Hered.</i> 104:161–171.
798 799	Charlesworth B. 2015. Causes of natural variation in fitness: evidence from studies of Drosophila populations. Proc. Natl. Acad. Sci. U.S.A. 112:1662–1669.
800 801	Charlesworth B, Morgan MT, Charlesworth D. 1993. The effect of deleterious mutations on neutral molecular variation. <i>Genetics</i> 134:1289–1303.
802 803	Charlesworth D, Charlesworth B, Morgan MT. 1995. The pattern of neutral molecular variation under the background selection model. <i>Genetics</i> 141:1619–1632.
804 805 806	Chikhi L, Rodríguez W, Grusea S, Santos P, Boitard S, Mazet O. 2018. The IICR (inverse instantaneous coalescence rate) as a summary of genomic diversity: insights into demographic inference and model choice. <i>Heredity</i> 120:13–24.
807 808	Choi JY, Aquadro CF. 2016. Recent and long term selection across synonymous sites in Drosophila ananassae. J. Mol. Evol. 83:50–60.
809 810	Comeron JM, Kreitman M. 2002. Population, evolutionary and genomic consequences of interference selection. <i>Genetics</i> 161:389–410.
811 812	Comeron JM, Ratnappan R, Bailin S. 2012. The many landscapes of recombination in Drosophila melanogaster. PLoS Genet. 8:e1002905.
813 814	Csilléry K, Blum MGB, Gaggiotti OE, François O. 2010. Approximate Bayesian Computation (ABC) in practice. <i>Trends Ecol. Evol.</i> 25:410–418.
815 816	Cutter AD, Payseur BA. 2013. Genomic signatures of selection at linked sites: unifying the disparity among species. <i>Nat. Rev. Genet.</i> 14:262–274.
817 818 819	Elyashiv E, Sattath S, Hu TT, Strutsovsky A, McVicker G, Andolfatto P, Coop G, Sella G. 2016. A genomic map of the effects of linked selection in <i>Drosophila</i> . <i>PLoS Genet</i> . 12:e1006130.
820 821	Ewing GB, Jensen JD. 2016. The consequences of not accounting for background selection in demographic inference. <i>Mol. Ecol.</i> 25:135–141.
822 823	Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust demographic inference from genomic and SNP data. <i>PLoS Genet</i> . 9:e1003905.
824 825 826	Eyre-Walker A, Keightley PD. 2009. Estimating the rate of adaptive molecular evolution in the presence of slightly deleterious mutations and population size change. <i>Mol. Biol. Evol.</i> 26:2097–2108.
827 828	Fay JC, Wu CI. 1999. A human population bottleneck can account for the discordance between patterns of mitochondrial versus nuclear DNA variation. <i>Mol. Biol. Evol.</i> 16:1003–1005.

829 830	Fiston-Lavier A-S, Singh ND, Lipatov M, Petrov DA. 2010. <i>Drosophila melanogaster</i> recombination rate calculator. <i>Gene</i> 463:18–20.
831 832 833	Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, Renkens I, van Duijn CM, Swertz M, Wijmenga C, van Ommen G, et al. 2015. Genome-wide patterns and properties of de novo mutations in humans. <i>Nat Genet</i> 47:822–826.
834 835 836	Fulgione A, Koornneef M, Roux F, Hermisson J, Hancock AM. 2018. Madeiran Arabidopsis thaliana reveals ancient long-range colonization and clarifies demography in Eurasia. <i>Mol. Biol. Evol.</i> 35:564–574.
837 838 839	Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. <i>PLoS Genet.</i> 5:e1000695.
840 841 842	Haddrill PR, Thornton KR, Charlesworth B, Andolfatto P. 2005. Multilocus patterns of nucleotide variability and the demographic and selection history of <i>Drosophila melanogaster</i> populations. <i>Genome Res.</i> 15:790–799.
843 844	Haller BC, Messer PW. 2019. SLiM 3: Forward genetic simulations beyond the Wright–Fisher model. <i>Mol. Biol. Evol.</i> 36:632–637.
845 846	Harris K, Nielsen R. 2013. Inferring demographic history from a spectrum of shared haplotype lengths. <i>PLoS Genet</i> . 9:e1003521.
847 848 849	Hernandez RD, Kelley JL, Elyashiv E, Melton SC, Auton A, McVean G, Project 1000 Genomes, Sella G, Przeworski M. 2011. Classic selective sweeps were rare in recent human evolution. <i>Science</i> 331:920–924.
850 851	Hey J, Harris E. 1999. Population bottlenecks and patterns of human polymorphism. <i>Mol Biol Evol</i> 16:1423–1426.
852 853	Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. <i>Genet. Res.</i> 8:269–294.
854 855 856	Hoggart CJ, Chadeau-Hyam M, Clark TG, Lampariello R, Whittaker JC, Iorio MD, Balding DJ. 2007. Sequence-level population simulations over large genomic regions. <i>Genetics</i> 177:1725–1731.
857 858 859	Hung C-M, Shaner P-JL, Zink RM, Liu W-C, Chu T-C, Huang W-S, Li S-H. 2014. Drastic population fluctuations explain the rapid extinction of the passenger pigeon. <i>Proc. Nat.</i> <i>Acad. Sci. U.S.A.</i> 111:10636–10641.
860 861 862	Jackson BC, Campos JL, Haddrill PR, Charlesworth B, Zeng K. 2017. Variation in the intensity of selection on codon bias over time causes contrasting patterns of base composition evolution in <i>Drosophila</i> . <i>Genome Biol. Evol.</i> 9:102–123.

- Jacquier H, Birgy A, Nagard HL, Mechulam Y, Schmitt E, Glodt J, Bercot B, Petit E, Poulain J,
 Barnaud G, et al. 2013. Capturing the mutational landscape of the beta-lactamase TEM-1. *Proc. Natl. Acad. Sci. U.S.A.* 110:13067–13072.
- Jensen JD, Payseur BA, Stephan W, Aquadro CF, Lynch M, Charlesworth D, Charlesworth B.
 2019. The importance of the Neutral Theory in 1968 and 50 years on: A response to Kern
 and Hahn 2018. *Evolution* 73:111–114.
- Johri P, Charlesworth B, Jensen JD. 2020. Toward an evolutionarily appropriate null model:
 jointly inferring demography and purifying selection. *Genetics* 215:173–192.
- Jones MR, Good JM. 2016. Targeted capture in evolutionary and ecological genomics. *Mol. Ecol.* 25:185–202.
- Kaiser VB, Charlesworth B. 2009. The effects of deleterious mutations on evolution in non recombining genomes. *Trends Genet*. 25:9–12.
- Keightley PD, Eyre-Walker A. 2007. Joint inference of the distribution of fitness effects of
 deleterious mutations and population demography based on nucleotide polymorphism
 frequencies. *Genetics* 177:2251–2261.
- Keightley PD, Jackson BC. 2018. Inferring the probability of the derived vs. the ancestral allelic
 state at a polymorphic site. *Genetics* 209:897–906.
- Keightley PD, Ness RW, Halligan DL, Haddrill PR. 2014. Estimation of the spontaneous
 mutation rate per nucleotide site in a *Drosophila melanogaster* full-sib family. *Genetics* 196:313–320.
- Keightley PD, Trivedi U, Thomson M, Oliver F, Kumar S, Blaxter ML. 2009. Analysis of the
 genome sequences of three *Drosophila melanogaster* spontaneous mutation accumulation
 lines. *Genome Res.* 19:1195–1201.
- Kelleher J, Wong Y, Wohns AW, Fadil C, Albers PK, McVean G. 2019. Inferring wholegenome histories in large population datasets. *Nat. Genet.* 51:1330–1338.
- Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler and D. 2002.
 The Human Genome Browser at UCSC. *Genome Res.* 12:996–1006.
- Kim Y, Wiehe T. 2009. Simulation of DNA sequence evolution under models of recent
 directional selection. *Brief. Bioinformatics* 10:84–96.
- Kousathanas A, Keightley PD. 2013. A comparison of models to infer the distribution of fitness
 effects of new mutations. *Genetics* 193:1197–1208.
- Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle
 M, FitzHugh W, et al. 2001. Initial sequencing and analysis of the human genome.
 Nature 409:860–921.

897 898	Lapierre M, Lambert A, Achaz G. 2017. Accuracy of demographic inferences from the site frequency spectrum: the case of the Yoruba population. <i>Genetics</i> 206:439–449.
899 900	Li H, Durbin R. 2011. Inference of human population history from individual whole-genome sequences. <i>Nature</i> 475:493–496.
901 902 903	Liang P, Saqib HSA, Zhang X, Zhang L, Tang H. 2018. Single-base resolution map of evolutionary constraints and annotation of conserved elements across major grass genomes. <i>Genome Biol. Evol.</i> 10:473–488.
904 905	Lukic S, Hey J. 2012. Demographic inference using spectral methods on SNP data, with an analysis of the human out-of-Africa expansion. <i>Genetics</i> 192:619–639.
906 907	Lynch M. 2007. The Origins of Genome Architecture. illustrated ed. Sunderland, Massachusetts: Sinauer Associates
908 909 910	Mazet O, Rodríguez W, Grusea S, Boitard S, Chikhi L. 2016. On the importance of being structured: instantaneous coalescence rates and human evolutionlessons for ancestral population size inference? <i>Heredity</i> 116:362–371.
911 912	Messer PW, Petrov DA. 2013. Frequent adaptation and the McDonald–Kreitman test. <i>Proc. Natl. Acad. Sci. U.S.A.</i> 110:8615–8620.
913 914	Nicolaisen LE, Desai MM. 2013. Distortions in genealogies due to purifying selection and recombination. <i>Genetics</i> 195:221–230.
915 916	O'Fallon BD, Seger J, Adler FR. 2010. A continuous-state coalescent and the impact of weak selection on the structure of gene genealogies. <i>Mol. Biol. Evol.</i> 27:1162–1172.
917 918	Orozco-terWengel P. 2016. The devil is in the details: the effect of population structure on demographic inference. <i>Heredity</i> 116:349–350.
919 920 921	Palkopoulou E, Lipson M, Mallick S, Nielsen S, Rohland N, Baleka S, Karpinski E, Ivancevic AM, To T-H, Kortschak RD, et al. 2018. A comprehensive genomic history of extinct and living elephants. <i>Proc. Nat. Acad. Sci. U.S.A.</i> 115:E2566–E2574.
922 923	Polanski A, Bobrowski A, Kimmel M. 2003. A note on distributions of times to coalescence, under time-dependent population size. <i>Theor Popul Biol</i> 63:33–40.
924 925 926	Polanski A, Kimmel M. 2003. New explicit expressions for relative frequencies of single- nucleotide polymorphisms with application to statistical inference on population growth. <i>Genetics</i> 165:427–436.
927 928	Pool JE, Nielsen R. 2007. Population size changes reshape genomic patterns of diversity. <i>Evolution</i> 61:3001–3006.
929	Pool JE, Nielsen R. 2009. Correction for Pool and Nielsen (2007). Evolution 63:1671.

930 Pouyet F, Aeschbacher S, Thiéry A, Excoffier L. 2018. Background selection and biased gene 931 conversion affect more than 95% of the human genome and bias demographic 932 inferences. Veeramah K, Wittkopp PJ, Gronau I, editors. eLife 7:e36317. 933 Ragsdale AP, Gutenkunst RN. 2017. Inferring demographic history using two-locus statistics. 934 Genetics 206:1037-1048. 935 Sanjuán R. 2010. Mutational fitness effects in RNA and single-stranded DNA viruses: common 936 patterns revealed by site-directed mutagenesis studies. Philos. Trans. R. Soc. Lond., B, 937 Biol. Sci. 365:1975–1982. 938 Schiffels S, Durbin R. 2014. Inferring human population size and separation history from 939 multiple genome sequences. Nat. Genet. 46:919-925. 940 Schneider A, Charlesworth B, Eyre-Walker A, Keightley PD. 2011. A method for inferring the 941 rate of occurrence and fitness effects of advantageous mutations. Genetics 189:1427-942 1437. 943 Schrider DR, Shanku AG, Kern AD. 2016. Effects of linked selective sweeps on demographic 944 inference and model selection. Genetics 204:1207-1223. 945 Sheehan S, Song YS. 2016. Deep learning for population genetic inference. PLoS Comput. Biol. 946 12:e1004845. 947 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, 948 Hillier LW, Richards S, et al. 2005. Evolutionarily conserved elements in vertebrate, 949 insect, worm, and yeast genomes. Genome Res. 15:1034-1050. 950 Slatkin M, Hudson RR. 1991. Pairwise comparisons of mitochondrial DNA sequences in stable 951 and exponentially growing populations. Genetics 129:555-562. 952 Speidel L, Forest M, Shi S, Myers SR. 2019. A method for genome-wide genealogy estimation 953 for thousands of samples. Nat. Genet. 51:1321-1329. 954 Steinrücken M, Kamm J, Spence JP, Song YS. 2019. Inference of complex population histories 955 using whole-genome sequences from multiple populations. Proc. Natl. Acad. Sci. U.S.A. 956 116:17115-17120. 957 Teshima KM, Coop G, Przeworski M. 2006. How reliable are empirical genomic scans for 958 selective sweeps? Genome Res 16:702-712. 959 Thornton K. 2003. Libsequence: a C++ class library for evolutionary genetic analysis. 960 *Bioinformatics* 19:2325–2327. 961 Thornton KR, Jensen JD. 2007. Controlling the false-positive rate in multilocus genome scans for selection. Genetics 175:737-750. 962

963 964	Torres R, Stetter MG, Hernandez RD, Ross-Ibarra J. 2020. The temporal dynamics of background selection in nonequilibrium populations. <i>Genetics</i> 214:1019–1030.
965 966	Torres R, Szpiech ZA, Hernandez RD. 2018. Human demographic history has amplified the effects of background selection across the genome. <i>PLoS Genet</i> . 14:e1007387.
967 968	Uricchio LH, Hernandez RD. 2014. Robust forward simulations of recurrent hitchhiking. <i>Genetics</i> 197:221–236.
969 970 971	Warren WC, Jasinska AJ, García-Pérez R, Svardal H, Tomlinson C, Rocchi M, Archidiacono N, Capozzi O, Minx P, Montague MJ, et al. 2015. The genome of the vervet (<i>Chlorocebus aethiops sabaeus</i>). <i>Genome Res.</i> 25:1921–1933.
972 973 974	Williamson RJ, Josephs EB, Platts AE, Hazzouri KM, Haudry A, Blanchette M, Wright SI. 2014. Evidence for widespread positive and negative selection in coding and conserved noncoding regions of <i>Capsella grandiflora</i> . <i>PLoS Genet</i> . 10:e1004622.
975 976 977	Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian J-H, Pei Y-Y, et al. 2020. A new coronavirus associated with human respiratory disease in China. <i>Nature</i> 579:265–269.
978 979	Zeng K, Charlesworth B. 2010. Studying patterns of recent evolution at synonymous sites and intronic sites in <i>Drosophila melanogaster</i> . J. Mol. Evol. 70:116–128.
980 981 982	Zhou Y, Massonnet M, Sanjak JS, Cantu D, Gaut BS. 2017. Evolutionary genomics of grape (<i>Vitis vinifera ssp. vinifera</i>) domestication. <i>Proc. Nat. Acad. Sci. U.S.A.</i> 114:11715–11720.
983 984 985	

SUPPLEMENTARY



Supp Figure 1: Performance of MSMC under neutrality and demographic equilibrium when using 1, 2 and 4 diploid individuals for inference and with varying chromosome sizes: (a) 10 Mb, (b) 50 Mb, (c) 200 Mb, (d) 1 Gb. The numbers of replicates for each panel varied slightly around 100. MSMC runs with 4 diploid individuals could not be obtained due to the long computational times required. As shown, MSMC has a tendency to take a common shape, often falsely indicating recent population growth.



Supp Figure 2: Performance of *fastsimcoal2* under neutrality and demographic equilibrium with varying chromosome sizes: (a) 10 Mb, (b) 50 Mb, (c) 200 Mb, (d) 1 Gb and when different model choices are provided: *Left panel*: the correct model of constant population size is specified; *Middle panel*: model selection was performed between 3 models – equilibrium, instantaneous size change, and exponential size change; *Right panel*: model selection was performed between 4 models – equilibrium, instantaneous size change, exponential size change, and instantaneous bottleneck. Inference was performed using 50 diploid individuals and the inferred population size estimates of the best model are plotted (blue lines). The numbers of replicates for each panel varied slightly around 100. The true model is shown in black.



Supp Figure 3: Inferred demography from MSMC (red lines) and *fastsimcoal2* (blue lines) in the presence of background selection, with the true DFE shown to the left of the panel, for 2-fold instantaneous decline (right column) and 2-fold exponential growth (left column). In this case, 20% of the genome experiences direct selection. The true demographic models are shown as black lines.



Supp Figure 4: Inference of demography by MSMC (red lines; 10 replicates) and *fastsimcoal2* (blue lines; 10 replicates) under demographic equilibrium (left column), 30-fold exponential growth (middle column), and ~6-fold instantaneous decline (right column) in the presence of direct purifying selection (*i.e.*, directly selected sites are not masked). The true demographic model is depicted in black lines. Exonic sites experience purifying selection specified by the following DFEs: (a) DFE1: $f_0 = 0.1$, $f_1 = 0.7$, $f_2 = 0.1$, $f_3 = 0.1$, (b) DFE2: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.7$, $f_3 = 0.1$, (c) DFE3: $f_0 = 0.1$, $f_1 = 0.1$, $f_2 = 0.1$, $f_3 = 0.7$, (d) DFE4: $f_0 = 0.25$, $f_1 = 0.25$, $f_2 = 0.25$, $f_3 = 0.25$, (e) DFE5: $f_0 = 0.7$, $f_1 = 0.0$, $f_2 = 0.0$, $f_3 = 0.3$, (f) DFE6: $f_0 = 0.5$, $f_1 = 0.0$, $f_2 = 0.0$, $f_3 = 0.5$. In this case, 20% of the genome is under selection.



Supp Figure 5: Distribution of lengths of repeat regions in the human genomes (hg19). Shown above is the distribution of lengths up to 1000 bp, although lengths of repeat regions range between 6 to 160602 bp.



Supplementary Figure 6: Performance of demographic inference by MSMC (red lines) and *fastsimcoal2* (blue lines) under different scenarios of neutrality when the true model is equilibrium: (a) there is variation in recombination and mutation rates and the centromeric region is masked, (c) there is variation in recombination and mutation rates, and short regions resembling repeats (comprising 10% of each chromosome) are randomly masked across the genome, and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome. The maximum and minimum fold change detected in every scenario is indicated on the upper right corner.



Supplementary Figure 7: Performance of demographic inference by MSMC (red lines) and *fastsimcoal2* (blue lines) under different scenarios of neutrality when the true model is 30-fold exponential growth: (a) there is variation in recombination and mutation rates across the genome, (b) there is variation in recombination and mutation rates, and the centromeric region is masked, (c) there is variation in recombination and mutation rates, and short regions resembling repeats are randomly masked across the genome (comprising of 10% of each chromosome), and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome.



Supplementary Figure 8: Performance of demographic inference by MSMC under different scenarios of neutrality when the true model is 6-fold instantaneous decline: (a) there is variation in recombination and mutation rates across the genome, (b) there is variation in recombination and mutation rates, and the centromeric region is masked, (c) there is variation in recombination sand mutation rates, and short regions resembling repeats are randomly masked across the genome (comprising of 10% of each chromosome), and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome.



Supplementary Figure 9: Performance of demographic inference by MSMC (red lines) and *fastsimcoal2* (blue lines) in the presence of background selection, under different scenarios when the true model is equilibrium: (a) there is variation in recombination and mutation rates, (b) there is variation in recombination and mutation rates and the centromeric region is masked, (c) there is variation in recombination and mutation rates, and short regions resembling repeats (comprising 10% of each chromosome) are randomly masked across the genome, and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome. Exons comprise of 20% of the genome, experience purifying selection given by DFE4 ($f_0 = f_1 = f_2 = f_3 = 0.25$) and are masked when performing inference.



Supplementary Figure 10: Performance of demographic inference by MSMC (red lines) and *fastsimcoal2* (blue lines) in the presence of background selection, under different scenarios when the true model is 30-fold exponential growth: (a) there is variation in recombination and mutation rates, (b) there is variation in recombination and mutation rates and the centromeric region is masked, (c) there is variation in recombination and mutation rates, and short regions resembling repeats (comprising 10% of each chromosome) are randomly masked across the genome, and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome. Exons comprise of 20% of the genome, experience purifying selection given by DFE4 ($f_0 = f_1 = f_2 = f_3 = 0.25$), and are masked when performing inference.



Supplementary Figure 11: Performance of demographic inference by MSMC (red lines) and *fastsimcoal2* (blue lines) in the presence of background selection, under different scenarios when the true model is a 6-fold instantaneous decline: (a) there is variation in recombination and mutation rates, (b) there is variation in recombination and mutation rates and the centromeric region is masked, (c) there is variation in recombination and mutation rates, and short regions resembling repeats (comprising 10% of each chromosome) are randomly masked across the genome, and (d) there is variation in recombination and mutation rates, and the centromere as well as small-sized repeats are randomly masked across the genome. Exons comprise of 20% of the genome, experience purifying selection given by DFE4 ($f_0 = f_1 = f_2 = f_3 = 0.25$), and are masked when performing inference.





Supp Figure 12: Scenarios where demographic inference by *fastsimcoal2* resulted in different best models when performed using all SNPs (right column) and when SNPs were thinned to be separated by 5 kb (left column). The best model was defined as the model corresponding to the lowest AIC over all ten replicates. The DFEs are indicated on the left. (a) Direct purifying selection under DFE4 (exonic sites are not masked) in 5% of the genome; (b) Direct purifying selection under DFE5 (exonic sites are not masked) in 5% of the genome; (c) Direct purifying selection under DFE1 (exonic sites are not masked) in 10% of the genome; (d) Direct purifying selection under DFE6 (exonic sites are not masked) in 10% of the genome; (e) Background selection under DFE5 (*i.e.*, exonic sites are masked) in which 5% of the genome is exonic; (f) Background selection under DFE1 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic; (h) Background selection under DFE2 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic; (h) Background selection under DFE2 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic; (h) Background selection under DFE2 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic; (h) Background selection under DFE4 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic; (h) Background selection under DFE4 (*i.e.*, exonic sites are masked) in which 20% of the genome is exonic.



Supp Figure 13: Inference of demographic history by MSMC. Top panel / red line: simulations in which the true model is constant population size, and 50% of new mutations in exons are strongly deleterious with the remainder being neutral, where exons comprise 5% of the genome. Bottom panel / green line: the empirical estimate of population history of the YRI population inferred with MSMC by Schiffels and Durbin (2014). The x-axis is in years (assuming a generation time of 30 years). Note that the y-axes are on different scales, and the magnitude of change observed in the empirical data is considerably larger in the simulated data. Thus, this comparison is only meant to illustrate this common shape taken in MSMC plots (and see similar shapes in, for example, vervets (Warren *et al.* 2015; Figure 4) and passenger pigeons (Hung *et al.* 2014; Figure 2)).

	Demographic	Ancestral	Current	Time of change
	mouers	population size	population size	in generations
1	Equilibrium	10,000	10,000	NA
2	Exponential growth	1000	30,000	850
3	Instantaneous decline	12,300	2,100	4,750

Supp Table 1: Parameters underlying the human-like demographic models considered.

Supp Table 2: Nucleotide diversity in the presence of BGS relative to that under neutrality (*B*) calculated for a neutral site distance *y* bases from the end of a gene/exon of length 500 bp. The exon experiences purifying selection with strength $2N_{es}=10$, where N_{e} is the effective population size and *s* is the reduction in fitness. Shown below is *t=hs* where *h* is the dominance coefficient (assumed to be 0.5 here), *r* is the recombination rate per site per generation and *u* is the mutation rate per site per generation. *B* was calculated by using Equation 2 of Johri *et al.* (2020).

	$t (2N_{\rm e}s)$					
Ne	=10)	r	и	<i>B</i> (y=1)	<i>B</i> (y=10)	<i>B</i> (y=1000)
10^{4}	0.00025	1.00×10^{-8}	1.00×10^{-8}	0.9805	0.9806	0.9820
		1.00×10^{-6}	1.00×10^{-6}	0.5152	0.5312	0.9444
10 ⁶	2.5×10^{-6}	1.00×10^{-8}	1.00×10^{-8}	0.5152	0.5312	0.9445
		1.00×10^{-6}	1.00×10^{-6}	0.4920	0.8227	0.9992