

# 1 Strong neutral sweeps occurring during a population contraction

2

3 Antoine Moinet<sup>1,2,3</sup>, Stephan Peischl<sup>\*1,2</sup> and Laurent Excoffier<sup>\*2,3</sup>

4

5 1. Interfaculty Bioinformatics Unit, University of Bern, Baltzerstrasse 6, 3012 Bern,  
6 Switzerland;

7 2. Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland;

8 3. Institute of Ecology and Evolution, University of Bern, Baltzerstrasse 6, 3012 Bern,  
9 Switzerland;

## 10 Abstract

11 A strong reduction in diversity around a specific locus is often interpreted as a recent rapid  
12 fixation of a positively selected allele, a phenomenon called a selective sweep. Rapid fixation of  
13 neutral variants can however lead to similar reduction in local diversity, especially when the  
14 population experiences changes in population size, e.g., bottlenecks or range expansions. The  
15 fact that demographic processes can lead to signals of nucleotide diversity very similar to  
16 signals of selective sweeps is at the core of an ongoing discussion about the roles of  
17 demography and natural selection in shaping patterns of neutral variation. Here we  
18 quantitatively investigate the shape of such neutral valleys of diversity under a simple model of  
19 a single population size change, and we compare it to signals of a selective sweep. We  
20 analytically describe the expected shape of such “neutral sweeps” and show that selective  
21 sweep valleys of diversity are, for the same fixation time, wider than neutral valleys. On the  
22 other hand, it is always possible to parametrize our model to find a neutral valley that has the  
23 same width as a given selected valley. We apply our framework to the case of a putative  
24 selective sweep signal around the gene *Quetzalcoat1* in *D. melanogaster* and show that the  
25 valley of diversity in the vicinity of this gene is compatible with a short bottleneck scenario  
26 without selection. Our findings provide further insight in how simple demographic models can  
27 create valleys of genetic diversity that may falsely be attributed to positive selection.

## 28 Introduction

29 Past demography and natural selection play a critical role in shaping extant genetic diversity. A  
30 central question in population genetics is to quantify their respective impact on observed  
31 genomic diversity. Because selection interferes with demographic estimates and vice versa,  
32 estimation of one of these two components is difficult without accounting for the other  
33 (Charlesworth *et al.* 1993, 1995; Kaiser and Charlesworth 2009; O'Fallon *et al.* 2010;  
34 Charlesworth 2013; Nicolaisen and Desai 2013; Johri *et al.* 2020, 2021b). Moreover, the relative  
35 importance of demography and selection as determinants of genome wide diversity is currently  
36 hotly debated, and may vary extensively among species (Corbett-Detig *et al.* 2015; Rousselle *et*  
37 *al.* 2018; Pouyet and Gilbert 2019; Galtier and Rousselle 2020). It has been shown that selection  
38 and demography can leave very similar footprints on the genetic diversity of a population  
39 (Andolfatto and Przeworski 2000; Teshima *et al.* 2006; Thornton and Jensen 2007; Johri *et al.*  
40 2021a). Disentangling the effects of demography and selection is therefore crucial to avoid  
41 erroneous inference of evolutionary scenarios from genomic data (Jensen *et al.* 2005; Wares  
42 2009; Mathew and Jensen 2015; Johri *et al.* 2020).

43 Hard selective sweeps lead to valleys of strongly reduced diversity around positively selected  
44 sites due to the hitchhiking of linked neutral loci (Maynard Smith and Haigh 1974), such  
45 observations of strong depletions of diversity in some genomic regions are often interpreted as  
46 due to past episode of positive selection, because the probability to observe a fast fixation of a  
47 neutral variant in a population of constant size is extremely low. However, during a range  
48 expansion for instance, some neutral or even mildly deleterious mutations can go quickly to  
49 fixation due to the low effective size of populations on the front of the range (Edmonds *et al.*  
50 2004; Klopstein *et al.* 2006; Hallatschek and Nelson 2008; Peischl *et al.* 2013), a phenomenon  
51 termed allele surfing (Klopstein *et al.* 2006). Theoretical studies have shown that the average  
52 neutral diversity on the wave front decays exponentially as the range expands (Hallatschek and  
53 Nelson 2008), similarly to what happens when a population experiences a sudden decay of the  
54 population size, i.e. a population contraction, due to a drastic change in the environment for  
55 example. In both cases, a mutation appearing when the population size is shrinking might go  
56 quickly to fixation, inducing a strong decrease of diversity in the surrounding genomic region,

57 whereas the average level of diversity might stay quite high depending on the strength and the  
58 duration of the contraction. As a result, the coalescent tree of alleles sampled in a population  
59 with strongly reduced effective population size will have short external branches, and long  
60 internal branches, depending on the parameters of the model (Excoffier *et al.* 2009). The site  
61 frequency spectrum associated to such a tree resembles a neutral SFS, but with a lack of rare  
62 alleles and an excess of high frequency sites, i.e. it becomes “flatter” (Sousa *et al.* 2014; Peischl  
63 and Excoffier 2015). The footprint left by the rapid fixation of a neutral allele on the  
64 surrounding genomic diversity, might thus be like that of a positively selected allele sweeping  
65 through a constant size population.

66 The expected shape of nucleotide diversity in genomic regions surrounding a site undergoing a  
67 rapid neutral fixation has been investigated analytically and numerically. Tajima (1990) studied  
68 the reduction of diversity during a neutral fixation at a given recombination distance from the  
69 fixing site. His results rely on rigorous mathematical arguments based on diffusion theory, but  
70 no closed form solution is provided for the shape of a neutral sweep. Johri *et al.* (2021a)  
71 described the valley of diversity occurring around a neutral fixation using an approach  
72 introduced for selective sweeps, assuming that the evolution of the allele frequency is that of a  
73 selected allele except in the initial stochastic phase. Here, we extend this work by inferring the  
74 dynamics of fixation of neutral alleles after a population contraction and we examine their  
75 effects on neighboring regions of the genome. We provide an analytical result for the expected  
76 coalescence time as a function of the recombination distance from the locus undergoing a fast  
77 fixation. Importantly, our results apply regardless of the process driving the allele going to  
78 fixation (neutrality, positive selection, background selection), as it only relies on the typical  
79 trajectory of an allele going to fixation in a given time, even though this trajectory differs  
80 depending on the underlying driver of this fixation (i.e., neutrality or selection). We compare  
81 our results against simulations and find that they hold for a wide range of realistic parameter  
82 combinations. We compare our results about the signature of neutral sweeps to patterns  
83 expected under selective sweeps and discuss potential differences between the signatures that  
84 could potentially allow us to discriminate between neutral and selective processes for a given  
85 demographic scenario. Finally, we investigate the similarity between the genomic signature of

86 an allele going to fixation either selectively or neutrally and observe that a selective sweep  
87 signal can in principle be replicated in a neutral model with an appropriate choice of  
88 demographic parameters. To illustrate this point, we examine a classical example of a selective  
89 sweep found in the genome of *D. melanogaster* around the *QtzI* gene (Rogers *et al.* 2010). We  
90 conclude that strong diversity depletions in the genome of a population, often attributed to the  
91 effect of positive selection, can be obtained with demographic effects only, and we call for  
92 caution when trying to detect signals of adaptation from genomic data, adding support to  
93 previous studies reaching similar conclusions (Thornton and Jensen 2007; Crisci *et al.* 2013;  
94 Jensen *et al.* 2019).

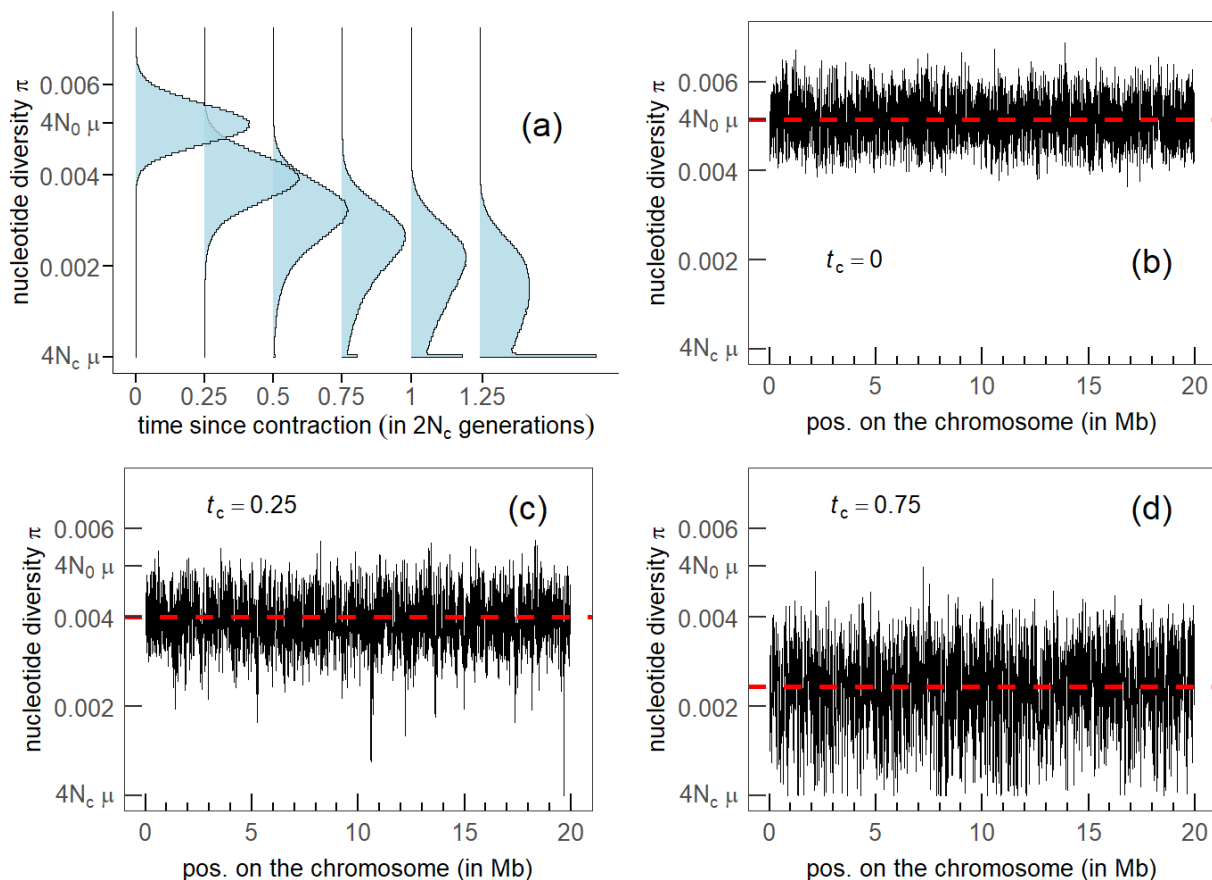
## 95 Model

96 We model here the effect of an instantaneous population contraction on genomic diversity.  
97 Throughout the whole manuscript, time is measured backwards. We assume that  $t_c$  generations  
98 before the present, the population size instantaneously dropped from  $N_0$  diploid individuals to  
99  $N_c$  individuals with  $N_c < N_0$ . We assume a standard coalescent model (Kingman 1982a; b) with  
100 discrete non-overlapping generations, random mating, monoecious individuals, and no  
101 selection. Two haplotypes sampled in the current population at time  $t = 0$  have, as we go  
102 backwards in time, a constant probability  $(2N_c)^{-1}$  of coalescing at each generation, for the first  $t_c$   
103 generations, and then this probability switches to  $(2N_0)^{-1}$  as we enter the ancestral  
104 uncontracted population. We can approximate the distribution of coalescence time  $T$  of these  
105 two haplotypes as a piecewise exponential distribution (see Appendix) with expected value:

$$106 \quad E[T] = 2(N_0 - N_c) e^{-t_c/2N_c} + 2N_c. \quad (1)$$

107 We see that the expected coalescence time decreases exponentially with the age of the  
108 contraction  $t_c$  and that it approaches  $2N_c$  for a very old contraction. Coalescence times cannot  
109 be measured directly from empirical data, but they are closely related to nucleotide diversity  $\pi$ .  
110 Under the infinitely many sites model, the number of nucleotide differences between two  
111 homologous DNA segments is proportional to their coalescence time  $T$  as  $\pi = 2\mu T$ , where  $\mu$  is  
112 the total mutation rate for the whole segment. Multiplying eq. (1) by  $2\mu$  shows that an

113 instantaneous population contraction leads to an exponential decrease of the expected  
114 nucleotide diversity along the genome with the age of the contraction  $t_c$ . However, it does not  
115 inform us on the distribution of nucleotide diversity  $\pi$  along the genome, or on spatially  
116 correlated patterns of diversity such as local depletion or excess of diversity relative to the  
117 expectation.



118

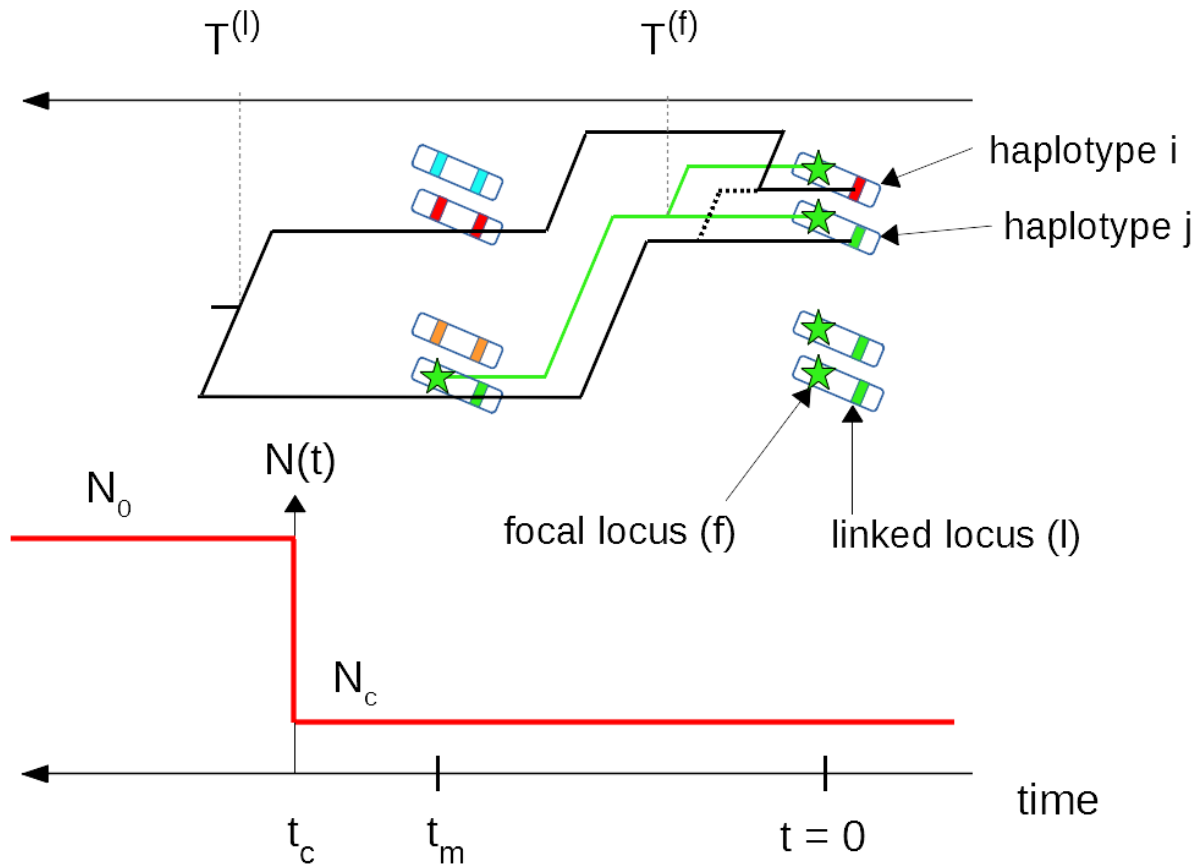
119 **Figure 1.** Nucleotide diversity of a population experiencing a contraction, as a function of the  
120 time  $t_c$  elapsed since the contraction, measured in units of  $2N_c$ . (a) distribution of nucleotide  
121 diversity as a function of time, nucleotide diversity along the chromosome at  $t_c = 0$  (panel b), at  
122  $t_c = 0.25$  (panel c) and at  $t_c = 0.75$  (panel d). Population size before contraction  $N_0 = 2.37 \times 10^6$   
123 and after contraction  $N_c = 4,400$ . Mutation rate  $\mu = 5.42 \times 10^{-10}$  per site per generation.  
124 Recombination rate  $r = 3.5 \times 10^{-8}$  per site per generation. Chromosome size  $L = 20$  Mb. Window  
125 size 10 Kb sliding at 1 Kb intervals. Sample size: 30 haplotypes. These parameters are taken from  
126 Rogers et al. (2010). Simulations were performed with fastsimcoal2 (Excoffier et al. 2021).

127

128 Fig. 1 shows the evolution of the distribution of  $\pi$  as a function of the time  $t_c$  elapsed since the  
129 contraction. For  $t_c = 0$ , there is no contraction, and the population size remains constant and  
130 equal to  $N_0$ . In this case we see (fig. 1a,1b,  $t_c = 0$ ) that the distribution of  $\pi$  is symmetric and  
131 centered at  $E[\pi] = 4N_0\mu$ . For an older contraction, we see that the distribution is not only  
132 shifted to lower values of diversity as expected from eq. (1), but that it also becomes strongly  
133 peaked around  $\pi = 4N_c\mu$ . This bimodality of the distribution can be understood intuitively in the  
134 following way. There are two possible types of coalescent trees for haplotypes sampled after  
135 the population contraction (note that the tree depends on the locus considered because of  
136 recombination). Indeed, the most recent common ancestor (MRCA) of the sample lived either  
137 before the contraction ( $T_{\text{MRCA}} > t_c$ ), or after the contraction ( $T_{\text{MRCA}} < t_c$ ). In the former case, the  
138 tree at this locus has long inner branches and short outer branches, whereas in the latter case,  
139 the tree is essentially a (short) neutral tree corresponding to a population of constant size  $N_c$   
140 (Excoffier *et al.* 2009). Both types of trees occur at different loci and correspond to the two  
141 observed modes in the distribution of the nucleotide diversity along the chromosome. The  
142 precise shape of the distribution of nucleotide diversity across sites depends on the relative  
143 frequency of both types of trees, which itself depends on the age of the contraction  $t_c$ . For a  
144 sample of size two, the probability that the MRCA lived after the contraction, that is,  $T_{\text{MRCA}} < t_c$   
145 is  $1 - e^{-t_c/2N_c}$ . For a larger sample of haplotypes, there is no closed form solution for this  
146 probability, but the trees rooted after the contraction are rare for  $t_c \ll 2N_c$  and very frequent  
147 when  $t_c \gg 2N_c$  (Tavaré 1984). Therefore, the evolution of the distribution of  $\pi$  for increasing  
148 contraction age  $t_c$  appears to be a transition from a unimodal distribution centered at  $4N_0\mu$  to  
149 another unimodal distribution centered at  $4N_c\mu$ , with both modes coexisting for intermediate  
150 ages (fig. 1). This bimodality has been pointed out previously in the context of population  
151 bottlenecks (Austerlitz *et al.* 1997); however, those studies mainly focused on long duration  
152 bottlenecks (the effect of a contraction or a bottleneck on nucleotide diversity is the same  
153 provided that the bottleneck is not yet finished, or that it finished very recently so that the  
154 effect of population recovery is negligible). In the present work, we investigate the effect of  
155 short contractions on the genetic diversity and make the claim that this short contraction  
156 regime is of particular interest as it can lead, such as in fig. 1c, to genomic signatures similar to

157 those generated by positive selection acting on a few sites in an otherwise neutral genome.  
158 More specifically, we want to quantitatively describe the reduction of diversity along the  
159 genome that is observed around a locus with a small  $T_{MRCA}$  (such as in fig. 1c in the regions  
160 around 10-11 and 19-20 Mb), where we observe a valley or trough of diversity. Akin to what is  
161 done for selective sweeps, we consider the (neutral) fast fixation of an allele and analyze the  
162 impact of hitchhiking on the genetic diversity of neighboring loci, and we refer to this process as  
163 a neutral sweep.

164 To investigate neutral sweeps in our model, we consider the following scenario:  $t_m$  generations  
165 ago a mutation occurred at a single site on the chromosome, which we call the focal site. We  
166 further assume that this mutation has just fixed in the population, i.e., that it was segregating at  
167 a frequency strictly lower than one in the last generation (at  $t = 1$ ), and has now (at  $t = 0$ ) a  
168 frequency equal to one. We assume that the population contraction occurred  $t_c$  generations  
169 ago, with  $t_c \geq t_m$ . As the mutant enters the population as a single allelic copy at the focal locus,  
170 defined as a non-recombining region surrounding the focal site, this copy is a common ancestor  
171 for all the copies ( $2N_c$ ) present at fixation. However, it is not necessarily the most recent  
172 common ancestor. Fig.2 shows a sketch of our model to help visualize how recombination can  
173 maintain diversity at linked loci around a locus where a new mutation quickly fixed in the  
174 population.



175

176 **Figure 2.** Instantaneous population contraction with a subsequent neutral fixation. A mutant  
 177 (green star) appeared  $t_m$  generations ago and has just fixed neutrally in a diploid population  
 178 that experienced a contraction  $t_c$  generations ago. We represent the population as a set of  $2N_c$   
 179 two-locus haplotypes that are painted so that the gene copies present at  $t = 0$  can be traced  
 180 back to  $t = t_m$ . Due to recombination, haplotype  $i$  carries a red gene copy at the linked locus at  
 181  $t = 0$ . Correspondingly, the coalescence time  $T^{(l)}$  of the haplotypes  $i$  and  $j$  at the linked locus  
 182 (black tree) is larger than  $t_m$ . On the other hand, the coalescence time  $T^{(f)}$  at the focal locus  
 183 (green tree) is smaller than  $t_m$  because at this locus all gene copies descend from the same  
 184 haplotype (due to the fixation of the focal mutation).

## 185 Results

### 186 Average coalescence time at a linked locus

187 We can calculate the expected coalescence time  $T^{(l)}$  of two randomly sampled haplotypes at a  
 188 linked locus as a function of the recombination rate  $r$  from the focal locus. The idea is to  
 189 consider two haplotypes with a given coalescence time  $T^{(f)}$  at the focal locus, and then follow



190 the genealogy of the gene copies carried by these two haplotypes at the linked locus backward  
191 in time, while considering possible recombination events. The expected coalescent time at the  
192 linked locus is then

$$193 \quad E[T^{(l)}] = \left(1 - E\left[e^{-2r \sum_{t=1}^{T^{(f)}} (1 - \bar{x}_t)}\right]\right) (t_m + T_m) + E\left[T^{(f)} e^{-2r \sum_{t=1}^{T^{(f)}} (1 - \bar{x}_t)}\right] \quad (2)$$

194 where  $\bar{x}_t$  is the average frequency of the mutant (derived) allele at the focal locus at time  $t$   
195 counting backward from present. A detailed derivation of this equation is given in Appendix A4.  
196 The first term of the right-hand side of eq. (2) corresponds to cases where lineages escape the  
197 neutral sweep due to recombination, and still have not coalesced after  $t_m$  generations. In this  
198 case we need to wait on average  $T_m = 2(N_0 - N_c) e^{-(t_c - t_m)/2N_c} + 2N_c$  extra generations  
199 before the lineages coalesce, due to the contraction that happened  $t_c - t_m$  generations before  
200 the focal mutation. The second term of the right-hand side of eq. (2) corresponds to cases  
201 where the lineages cannot escape the sweep and are forced to coalesce at a time  $T^{(l)} \leq t_m$ .

## 202 Distribution of coalescence times at the focal locus

203 To evaluate eq. (2), we need to determine the probability distribution of the pairwise  
204 coalescence times  $T^{(f)}$  at the focal locus, as well as the expected frequency trajectory of the  
205 derived allele. Even though this allele fixes neutrally in a population of constant size (the  
206 contraction occurs prior to the mutation), the distribution of coalescent times at the focal locus  
207  $T^{(f)}$  departs from the usual exponential distribution for a neutral coalescent process because the  
208 allele fixes in exactly  $t_m$  generations, and hence the coalescence time for a randomly chosen  
209 pair of haplotypes is at most  $t_m$ . Slatkin (1996) investigated the coalescent process within a  
210 “mutant allelic class” that originated from a single mutation at a given time in the past. He  
211 derived exact analytical results for the average pairwise coalescence time, but the coalescence  
212 distribution itself can only be expressed with multidimensional integrals and obtaining a closed  
213 form expression does not appear feasible. We therefore use a different approach: given a  
214 particular fixation trajectory of the mutant allele, i.e. given the number of mutant copies  $N_\mu$  at  
215 each generation between  $t = 0$  and  $t = t_m$ , we can express the coalescence time distribution  
216 within the mutant allelic class, using the result of a coalescent in a population with a time-

217 dependent (but deterministic) size  $N_\mu(t)$  (Griffiths and Tavaré 1994). Averaging over all  
 218 possible trajectories of the mutation, we obtain:

$$219 \quad P(T^{(f)}) = \sum_{\{x_t\}} \left[ \frac{1}{2N_c x_{T^{(f)}}} \prod_{t=1}^{T^{(f)}-1} \left( 1 - \frac{1}{2N_c x_t} \right) \right] P(\{x_t\}) \quad (3a)$$

220 where  $x_t = N_\mu(t)/(2N_c)$  is the frequency of the mutant  $t$  generations from fixation, and  
 221  $P(\{x_t\})$  is the probability of a given trajectory.  $P(\{x_t\})$  can be evaluated (see Appendix A2) and  
 222 the sum in eq. (3a) can in principle be computed numerically; however, the number of  
 223 trajectories to consider is prohibitive. As a first approximation, we can replace  $x_t$  by its  
 224 expectation  $\bar{x}_t$ , i.e., we neglect the fluctuations of the trajectory around the mean to obtain

$$225 \quad P(T^{(f)}) \simeq \frac{1}{2N_c \bar{x}_{T^{(f)}}} \prod_{t=1}^{T^{(f)}-1} \left( 1 - \frac{1}{2N_c \bar{x}_t} \right). \quad (3b)$$

226 The last step is to determine the average trajectory of an allele fixing in exactly  $t_m$  generations.  
 227 Zhao *et al.* (2013) as well as Maruyama and Kimura (Maruyama and Kimura 1975) have  
 228 investigated the characteristic trajectory of an allele fixing in a given time but they do not  
 229 provide a closed form solution. Here, we use a different approach (also based on diffusion  
 230 theory to obtain an approximation for the average trajectory of an allele fixing in exactly  $t_m$   
 231 generations, starting from a frequency  $p_0$ . As detailed in the Appendix A2, we obtain

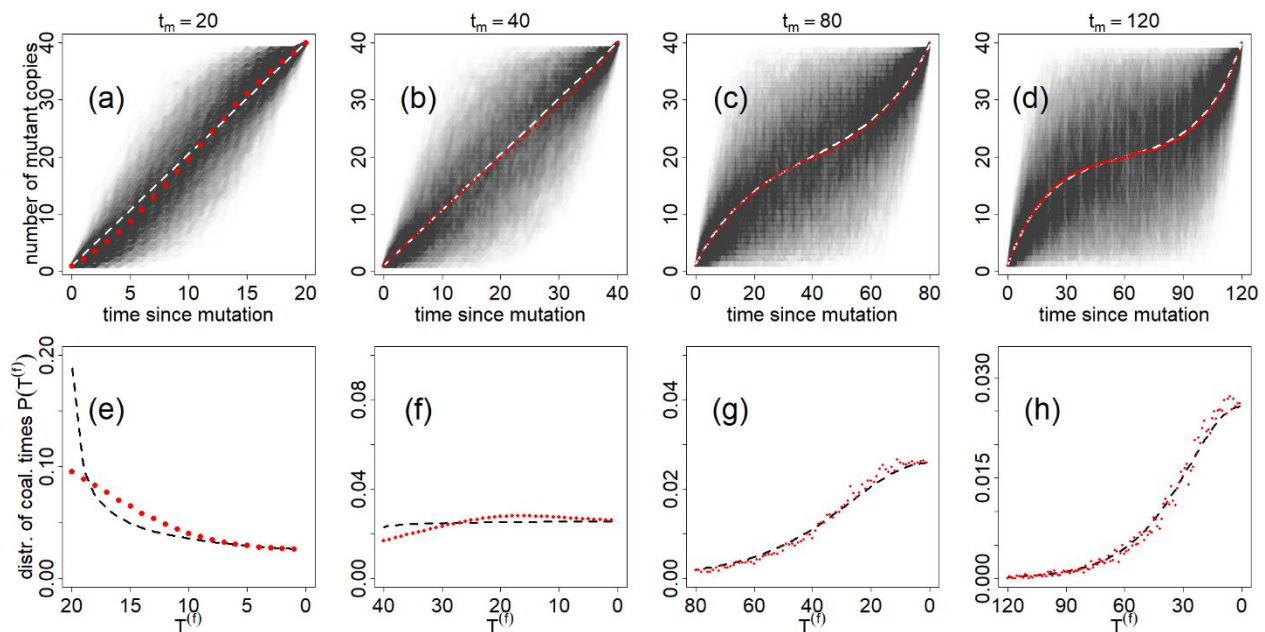
$$232 \quad \bar{x}_t = 1/2(1 - (1 - 2p_0)e^{-(t_m-t)/N_c} + e^{-t/N_c}), \quad (4a)$$

233 which is valid for  $t_m \gg 2N_c$ . For very fast fixations, i.e., when  $t_m \ll 2N_c$ , the frequency of the  
 234 allele increases approximately linearly as

$$235 \quad \bar{x}_t = 1 - (1 - p_0) \frac{t}{t_m}. \quad (4b)$$

236 We remind the reader that  $t$  is counted backwards from fixation. Fig. 3 compares equations (4a)  
 237 and (4b) to trajectories obtained from simulations of a Wright-Fisher diploid population. We  
 238 find good agreement between the simulations and the analytical results. Importantly, the  
 239 typical neutral trajectory for large values of the fixation time has an “inverse-sigmoid shape”  
 240 (fig. 3c), contrary to the typical sigmoid trajectory of a positively selected allele going to fixation  
 241 in a constant size population (see fig. 5a). This neutral trajectory occurs because, conditional on

242 non-loss, neutral alleles need to quickly escape loss at the beginning and remain at  
 243 intermediate frequencies to stay away from both fixation and loss until they eventually fix in  
 244 the population at  $t = 0$  (i.e. in exactly  $t_m$  generations). Fig. 3d-3f also shows the coalescence  
 245 time distribution for several values of the fixation time  $t_m$ . The comparison of the distribution of  
 246 pairwise coalescence time with numerical simulations of a Wright-Fisher model shows that our  
 247 approximation eq. (3b) is quite accurate but overestimates the probability of coalescence for  
 248 large coalescence times when  $t_m$  is small (fig. 3d). Notably, coalescence (simulated or  
 249 theoretical) is more probable at large times (i.e. when the mutant appeared) for short fixation  
 250 times (fig. 3d), whereas it is more probable at small times (i.e. close to fixation) for large  
 251 fixation times (fig. 3e). The coalescence rate within the mutant allelic class is given by the  
 252 inverse of the number of mutant copies and is for all values of the fixation time slightly more  
 253 than  $1/2N_c$  at the first generation. However, when the fixation time is short (fig. 3d), there is a  
 254 fast increase of the coalescence rate backwards in time, and many lineages are forced to  
 255 coalesce at  $t = t_m$ . When the fixation time is large (fig. 3f), the coalescence rate also increases  
 256 backwards in time, but the increase is much slower. In that case, most coalescence events  
 257 happen in much less than  $t_m$  generations, so that the early increase in frequency of the mutant  
 258 has almost no influence on the coalescence distribution.



259  
 260 **Figure 3.** Average frequency (a-d) and coalescence time distribution (e-h) of an allele fixing in a  
 261 diploid population of constant size  $N_c = 20$  in exactly  $t_m$  generations, starting as a single copy

262 (i.e.  $p_0 = (2N_c)^{-1}$ ). The red dots are the results of Wright-Fisher simulations, and the black and  
263 white dashed lines are calculated with eqs. (4b) (first and second columns) (4a) (third and fourth  
264 columns) and (3b). In panes (a-d) we show the variability of the fixation process by overlapping  
265 1780 fixing trajectories. The (numerically estimated) probability, for a mutant that appears at  
266 the onset of the contraction, to fix in less than  $t_m$  generations is 0.006, 0.16, 0.64 and 0.86 for  
267  $t_m = 20, 40, 80$  and  $120$  respectively (for this particular value of  $N_c$ ).

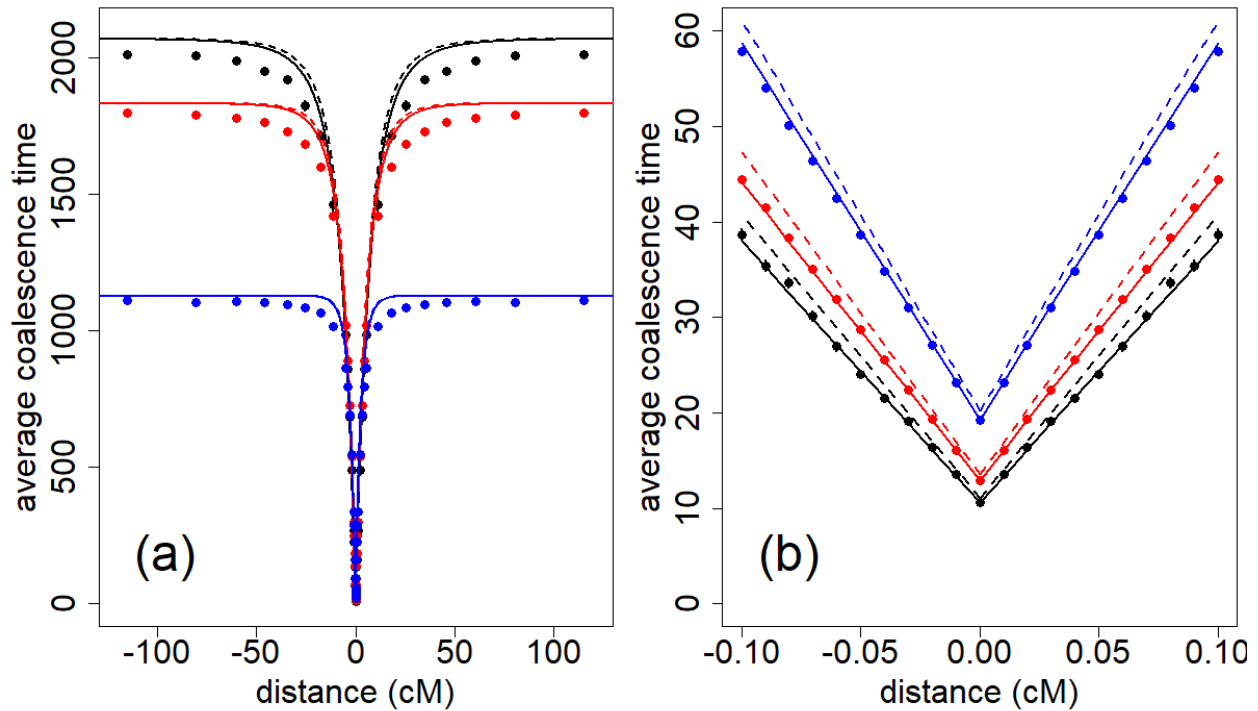
## 268 Effect of a neutral sweep on linked diversity

269 Combining equations (3b), (4a) with eq. (2) allows us to get an approximation for the average  
270 coalescence time at linked loci. Since the derivation of eq. (2) assumes that there is at most one  
271 recombination event in the genealogy of a randomly chosen pair of gene copies, we expect it to  
272 be only accurate for small values of the recombination rate  $r$ . For large values of  $r$  we use a  
273 heuristic approach combining the result of eq. (2), which is accurate for small  $r$ , and the  
274 expected diversity at unlinked loci, which is equal to  $T_0 = 2(N_0 - N_c) e^{-t_c/2N_c} + 2N_c$  as stated  
275 in eq. (1). We fit the trough of diversity with an exponential function of the form:

$$276 \quad E[T^{(l)}](r) = T_0(1 - ce^{-ar}), \quad (5)$$

277 where the coefficients  $c = 1 - E[T^{(f)}]/T_0$  and  $a = 2E[(t_m + T_m - T^{(f)}) \sum_{t=1}^{T^{(f)}} (1 - \bar{x}_t)] / (T_0 -$   
278  $E[T^{(f)}])$  are obtained by imposing that eqs. (2) and (5) coincide for small values of  $r$  (using a  
279 linear expansion in  $r$ ). On fig. 4 we compare the result of eq. (5) to Wright-Fisher simulations  
280 with two recombining loci. We see in fig. 4a that the exponential function fits the data accurately  
281 at large values of the recombination distance, but that the fit is biased for intermediate values of  
282  $r$ . In fig. 4b we see that the approximation is very good for low values of the recombination  
283 distance, although there still is a slight bias. This discrepancy at small  $r$  can be corrected (solid  
284 lines in fig. 4) if we use numerical estimations of  $\bar{x}_t$  and  $P(T^{(f)})$ , instead of eqs. (4) and (3b), to

285 evaluate eq. (5).



286

287 **Figure 4.** Average coalescence time at a linked locus, as a function of the recombination  
 288 distance from the focal locus where a mutant fixed in exactly  $t_m$  generations, starting from a  
 289 single copy  $t_m$  generations ago.  $t_m = 15$  in black,  $t_m = 20$  in red and  $t_m = 40$  in blue. The dots are  
 290 calculated with two-locus WF simulations, and compared to eq. (5) with either a numerical  
 291 estimation (solid lines) or a theoretical estimation (dashed lines) of  $\bar{x}_t$  and  $P(T^{(t)})$ .  $N_c = 20$ .  $N_0 =$   
 292 1500. The population experienced a contraction  $t_c = t_m$  generations ago.

293

294 We observe, as expected, on fig. 4 that the troughs of diversity induced by neutral sweeps are  
 295 wider and deeper for short fixation times. Similarly to what happens after a selective sweep,  
 296 there is less opportunity for linked loci to escape the sweep by recombination and maintain  
 297 diversity when the fixation is fast. In addition, the diversity level at the center of the valley is  
 298 given by the average coalescence time at the focal locus, which quickly decreases for small  
 299 fixation times  $t_m$ .

### 300 Comparison of neutral sweeps and selective sweeps

301 Since we did not make any assumption regarding the process driving the mutant allele to  
 302 fixation when deriving the average coalescence time at linked loci (eq. (2)) and the coalescence

303 time distribution at the focal locus (eq. (3b)), our framework allows us to directly compare the  
304 signatures of different processes that can drive mutations to fixation in a given number of  
305 generations. We illustrate this by comparing the effect of neutral and hard selective sweeps on  
306 linked diversity. Later we will discuss how neutral sweeps compare to a larger variety of  
307 scenarios (e.g. background selection, small selection coefficients, or dominant alleles). Here we  
308 assume that the neutral and selected fixations occurred over the same time interval, that is in  
309 both cases in exactly  $t_m$  generations. The selected fixation is assumed to be codominant ( $h=0.5$ )  
310 and occurs on an autosomal locus in a randomly mating diploid population of constant size  $N_1$ ,  
311 and we consider a strong selection strength ( $2N_1s \gg 1$ ) so that the allele frequency follows the  
312 deterministic trajectory

$$313 \quad \bar{x}_t = \frac{1}{1 + (2N_1 - 1) e^{-2(1-t/t_m) \log(2N_1)}}, \quad (6)$$

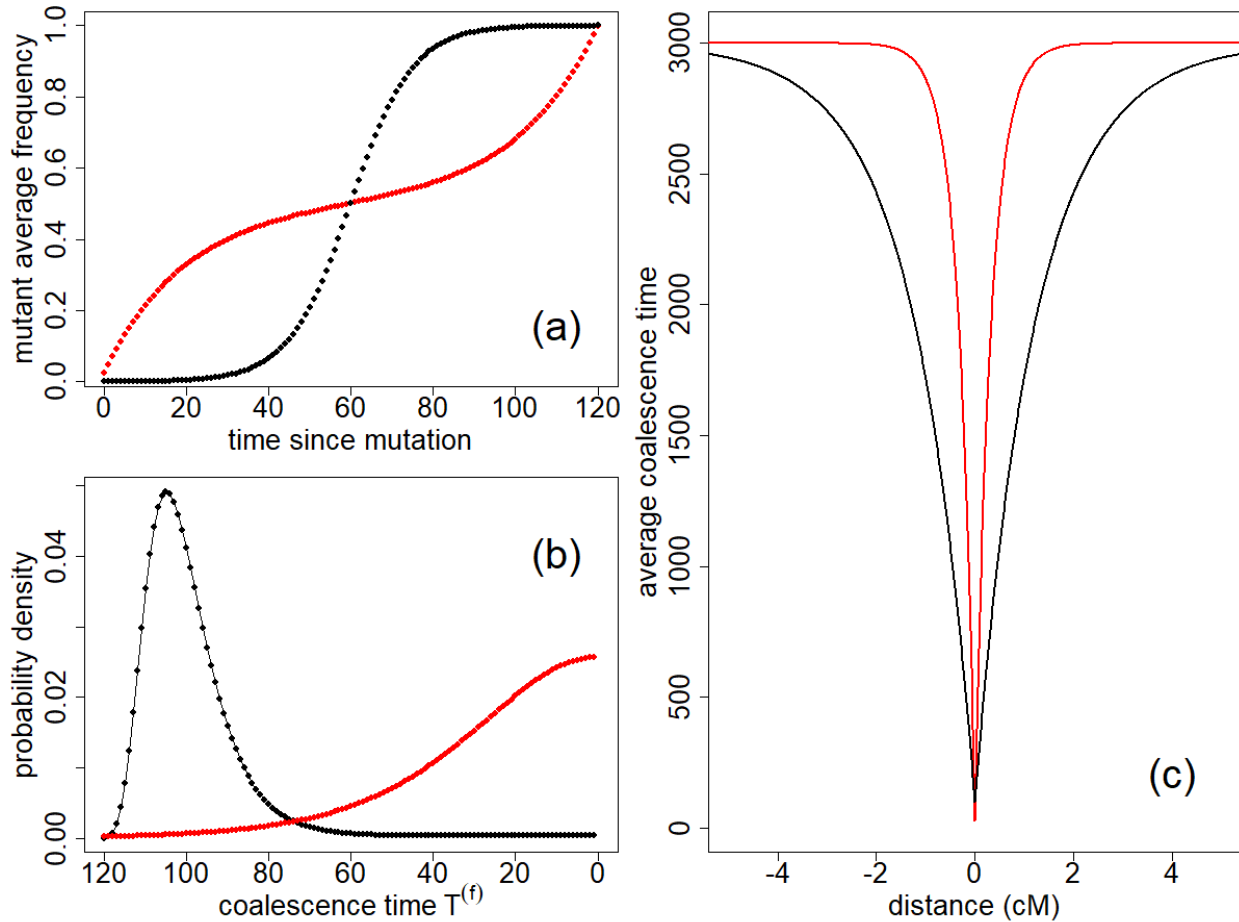
314  
315 where the fixation time is given by  $t_m(s) = 2\log(4N_1 s)/s$  (Barton 1995). Then combining eqs. (5),  
316 (3b) and (6), we can compute the average coalescence time at linked loci as a function of the  
317 recombination distance  $r$  to the focal locus, after replacing  $T_m$ , the average coalescence time at  $t$   
318  $= t_m$ , by  $2N_1$  in eq. (5) and  $N_c$  by  $N_1$  in eq. (3b). This approach yields results similar to  
319 Charlesworth (2020), where the author investigated signals of selective sweeps correcting for  
320 coalescent events that happen during the sweep, thus going beyond the common assumption of a  
321 star tree structure at the focal locus. For sake of simplicity in the neutral case, we consider that  
322 the mutant appeared at the time of the contraction, i.e.  $t_m = t_c$ . Furthermore, we will assume that  
323 the average coalescence times (and consequently the genetic diversity) are equal in both  
324 scenarios, i.e. that  $T_0 = 2N_1$  which implies that

325 
$$N_0(t_m) = (N_1 - N_c) e^{t_m/2N_c} + N_c . \quad (7)$$

326 In the neutral case we want the diversity to remain as high as  $4N_1\mu$  after the contraction, which

327 is possible only if the ancestral diversity was even higher, i.e. we have in general  $N_0 > N_1 > N_c$ .

328



329

330 **Figure 5.** Comparison between troughs of diversity resulting from a selective sweep (black) and

331 a neutral sweep (red), for the same fixation time  $t_m = 120$  (corresponding to  $s \approx 0.1$  in the

332 selective case). Frequency of the fixing allele as a function of time (a), coalescence time

333 distribution (b) and diversity around the fixing site along the genome using eq. (5) (c).  $N_1 = 1500$ ,

334  $N_c = 20$  and  $N_0 = 2.97 \times 10^4$ .

335

336 In fig. 5a, we compare the mutant average frequency as a function of time for a selected and a

337 neutral fixation. The dynamics of the neutral fixation is the opposite of that of the selected

338 allele in the sense that when one is increasing, the other is “resting” and vice versa. These

339 different trajectories translate into different coalescence distributions at the focal locus (fig.

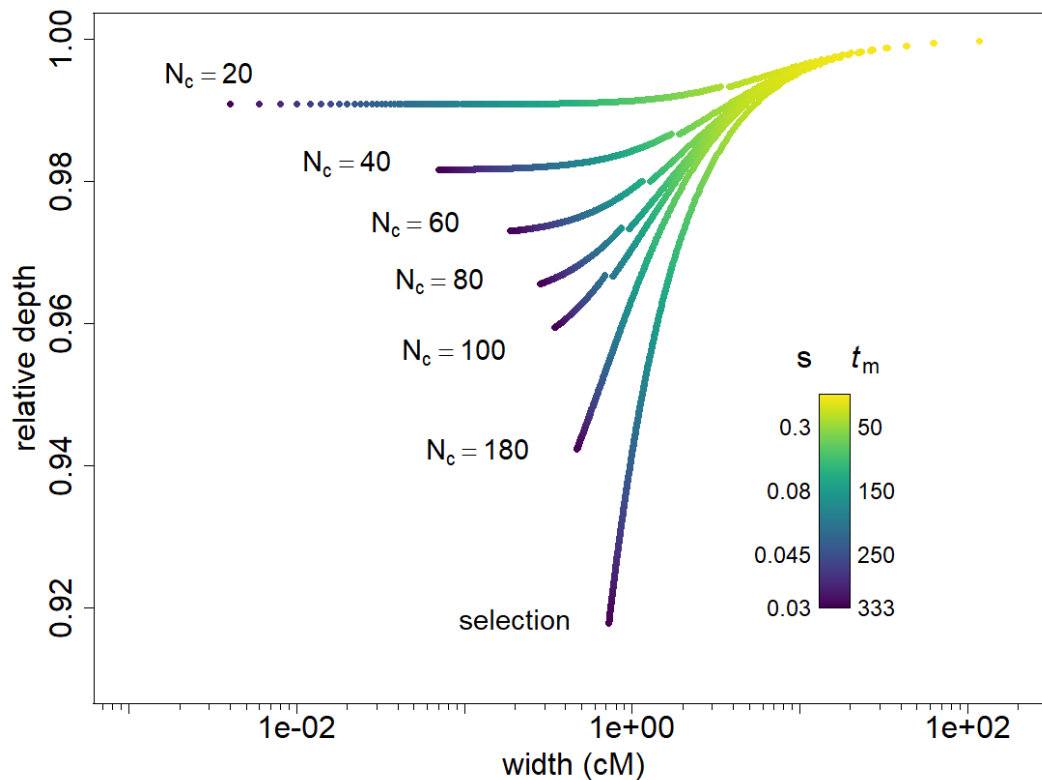


340 5b). If selection drives the fixation of the mutation, the distribution of coalescence time is  
341 peaked at large coalescence times. In contrast, in the neutral case the distribution is skewed  
342 towards small coalescence times. Correspondingly, the coalescence tree for the selected case  
343 has a star-like structure (not shown), whereas the tree for the neutral case has shorter outer  
344 branches. Therefore, for a given recombination distance, there will be fewer recombinations on  
345 the neutral tree because it has a much smaller total length. As recombination helps maintain  
346 diversity at linked loci, we would expect neutral troughs of diversity to be wider than in the  
347 selected case. However, this is at odds with the valleys of diversity observed in fig. 5c, where  
348 the selective trough is wider than the neutral trough. In fact, even though recombination is less  
349 abundant on the neutral tree, it is more efficient at recovering diversity. Indeed, if at a linked  
350 locus a pair of lineages escapes the sweep due to recombination, it takes on average an extra  
351  $2N_1$  generations, counted backwards from generation  $t = t_m$  when the mutant appeared, for  
352 them to coalesce in the selective case, and an extra  $2N_0$  generations in the neutral case. As  $N_0 >$   
353  $N_1$  two lineages escaping the sweep due to recombination have a larger coalescence time in the  
354 neutral case, and correspondingly a larger diversity, which explains why the neutral valley of  
355 diversity is narrower. Furthermore, we see that the trough is deeper in the neutral case (fig. 5c),  
356 since the average coalescence time is smaller at the focal site due to the smaller total length of  
357 the coalescence tree.

358  
359 To determine if these differences between selective and neutral troughs hold for other fixation  
360 times and population sizes, we define two quantities that characterize the shape of a trough, as  
361 well as its propensity to be detected in real data: i) the trough relative depth and ii) the width of  
362 the trough. The relative depth is defined as the difference between the background level of  
363 diversity and the diversity at the focal locus, divided by the background diversity, and the width  
364 is measured at half depth, i.e. halfway between the background diversity and the diversity at  
365 the focal locus. On fig. 6 we plot the relative depth of neutral and selective troughs as a  
366 function of their width for different fixation times  $t_m$ , calculated with our analytical expressions.  
367 We see that the neutral troughs are not only always narrower than the selective troughs for the  
368 same value of  $t_m$ , but also deeper. This is due to differences in the focal tree structure between



369 the selective case and the neutral case as well as difference in the ancestral background level in  
370 both cases, as explained above. For very short fixation times (corresponding to selection  
371 coefficients larger than 0.1), there is almost no difference between troughs generated by  
372 selective and neutral sweeps. Indeed, for such values of  $t_m$ , in both cases the focal coalescence  
373 tree is essentially a star tree because the increase in frequency is very fast, and the ancestral  
374 backgrounds of diversity,  $2N_0$  and  $2N_1$ , are also practically equal. Note however that at small  $t_m$   
375 the corresponding value of the selection coefficient  $s$  (see legend of fig. 6) may be  
376 unrealistically high. For realistic values of the selection coefficient/fixation time, the neutral  
377 troughs tend to be quite deep but narrow, whereas selective troughs are wider and their depth  
378 decreases quickly for low selection coefficients. From fig. 6, we see that the shape of a neutral  
379 trough is generally different from a selective sweep signal, but in practice those differences  
380 might be hidden due to the noise inherent present in real genomic data, and it might be  
381 difficult to decide whether a genomic signal is a due to a neutral sweep or a selective sweep.



382

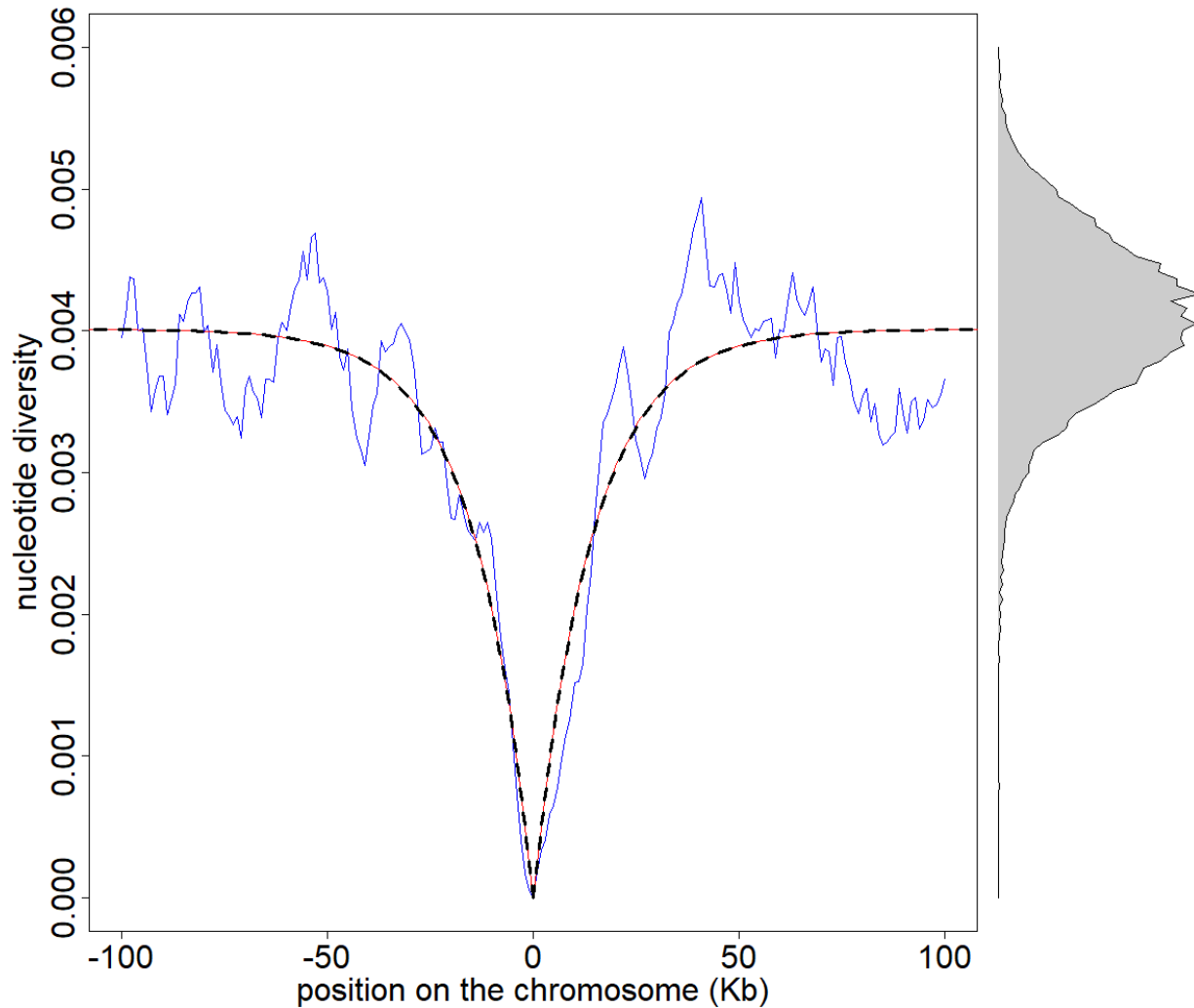
383 **Figure 6.** Relative depth as a function of the width of the diversity troughs, for different values  
384 of  $t_m$  and  $N_c$  in the neutral case and for selective scenarios with identical fixation times.  $t_m$  goes  
385 from 1 to 333 by increments of 1, the corresponding values of the selection coefficient  $s$  are

386 indicated on the left of the legend bar (for all of them we have  $N_1 s \gg 1$ ).  $N_1 = 1500$ .  $N_0$  is given  
387 by eq. (7) and depends on  $N_c$  and  $t_m$ . The jumps in the neutral curves for  $N_c = 20, 40, 60, 80$  and  
388 100 are due to the use of two different approximations for the frequency of the mutant, eqs.  
389 (4a) and (4b) and are located at  $t_m = 2N_c$ .

### 390 Is the Qtzl trough in *D. melanogaster* a neutral trough?

391 A region with reduced nucleotide diversity around the *Quetzalcoatl* gene identified in  
392 *Drosophila melanogaster* was judged compatible with a selective sweep (Rogers *et al.* 2010). A  
393 hard sweep model (Kaplan *et al.* 1989) was fitted assuming a constant population size of  $N_1 =$   
394  $1.85 \times 10^6$  diploid individuals and it was inferred that a positively selected allele fixed in the  
395 population  $1.5 \times 10^5$  generations ago ( $1.5 \times 10^4$  years) due to a selective advantage of  $s = 0.0098$   
396 (corresponding to a fixation time of more than 300 years). Using our theory, we fitted the data  
397 under a neutral demographic scenario of recent population size change that can generate  
398 neutral troughs with the same width and almost the same relative depth (less than 0.1%  
399 difference) as the *Quetzalcoatl* trough. To infer the demographic parameters, we measure the  
400 width of the selective sweep curve used to fit the data in (Rogers *et al.* 2010) and find a set of  
401 values of  $(N_c, t_m)$  that define a neutral trough with the same width. We then impose that  $t_m/2N_c$   
402  $= 0.25$  so that the troughs are rare yet observable along the chromosome as explained on fig. 1,  
403 and we obtain  $t_m = 2200$  and  $N_c = 4400$ . In fig. 7 we show a trough generated during a  
404 population contraction corresponding to these inferred values, using the software fastsimcoal2  
405 (Excoffier *et al.* 2021). We see that the neutral sweep fit is almost indistinguishable from the  
406 selective sweep fit because they not only have the same width, but also practically the same  
407 depth. Note that this simulated trough can be also seen in fig. 1c in the region 19-20 Mb. The  
408 same approach can be used to generate neutral troughs with a broad range of width and depth,  
409 which implies that in most cases, an alternative demographic neutral scenario can be  
410 compatible with a trough that is putatively due to selection. In practice, model inference does  
411 not rely solely on the fitting of a single trough, and genome wide information must be used.  
412 Therefore, we do not exclude here the possibility of the presence of adaptation in the  
413 *Quetzalcoatl* gene, but rather make the general warning that valleys of diversity do not  
414 necessarily indicate the presence of positive selection.

415 The authors affirm that all data necessary for confirming the conclusions of the article are  
416 present within the article, figures, and tables.



417

418 **Figure 7.** Trough of nucleotide diversity observed on a 20 Mb chromosome simulated with  
419 *fastsimcoal2*. The population experienced a contraction 2200 generations ago and the (diploid)  
420 population size was reduced from  $N_0 = 2.37 \times 10^6$  to  $N_c = 4400$ . The nucleotide diversity (blue line)  
421 is calculated on a sample of 30 haplotypes from our simulation. The black dashed line is the  
422 expected diversity (eq. (5)) for an allele that just fixed neutrally in the population, starting as a  
423 single copy 2200 generations ago. The red line is the expectation of a hard selective sweep with  
424 selection coefficient  $s = 0.0098$ . On the right we plot the distribution of nucleotide diversity for  
425 the whole chromosome. The mutation rate  $\mu = 5.42 \times 10^{-10}$  per site per generation, and  
426 recombination rate  $r = 3.5 \times 10^{-8}$  per site per generation were taken from (Rogers et al. 2010).  
427 The nucleotide diversity (in blue) is calculated for sliding windows of 10 Kb at 1 Kb intervals.

## 428 Discussion

429 It has repeatedly been suggested that strong depletions of diversity in the genome are not  
430 necessarily due to the presence of positive selection (Johri *et al.* 2020), and can also be the  
431 result of demographic effects only, such as the allele surfing phenomenon occurring at the front  
432 of a range expansion (Klopfstein *et al.* 2006). In this work, we considered a model of population  
433 contraction to analyze quantitatively the genomic signature of the rapid fixation of a mutation  
434 during a population contraction. Taking a step further from previous work that focused on the  
435 impact of range expansion on mere allele frequencies, we have studied here the impact of a  
436 neutral allele fixation on neighboring genomic diversity. We show that the diversity profile  
437 around a recently fixed locus crucially depends on the frequency trajectories of the allele going  
438 to fixation, and we outline the fact that neutrally fixing alleles have an inverse-sigmoid  
439 trajectory (fig. 3c), as compared to the standard sigmoid frequencies observed for positively  
440 selected alleles. For the same fixation time, this difference translates into different genomic  
441 signatures (see figs. 5c and 6). Our results demonstrate that there is a short period after a  
442 demographic contraction (or during a range expansion) where observed profiles of genomic  
443 diversity would look like those usually attributed to selection (fig. 1c), and that selective sweep  
444 signals can be mimicked by neutrally fixing mutations without the need to invoke complex  
445 histories of population size changes.

446 Our results allow for a systematic comparison of selective and neutral troughs of diversity, and  
447 we used our results to investigate trough shapes for range of neutral and selected scenarios  
448 (see fig. 6), which in principle can be used to decide whether a given empirical trough is due to  
449 selection or demography, and to infer the corresponding parameters. However, we did not  
450 consider the whole spectrum of possible selection scenarios. It would be indeed interesting to  
451 use our results to study cases of background selection, small selection coefficients, and a  
452 variety of dominance coefficients. All these cases should have their own characteristic  
453 trajectories of fixation, and hence potentially different genomic signatures. In addition, in our  
454 model we do not consider mutations that fixed in the past (we always assume that the allele  
455 has just reached fixation), nor do we consider mutations appearing before the population  
456 contraction, i.e., with  $t_m > t_c$ . The average coalescence time in the former case can be expressed

457 as a function of the coalescence time at fixation using conditional probabilities, and we can  
458 show that a sweep signal vanishes exponentially with the time elapsed since fixation (see  
459 Appendix A4). In the latter case, we can solve the problem by considering the number of gene  
460 copies at  $t_c$  that descend from the original copy that appeared at  $t_m$ . One could extend our  
461 results by considering an allele starting from an arbitrary number of copies at  $t_c$ , akin to soft  
462 selective sweeps; however, the analytic calculations are complex, and we leave this study for  
463 future research. In any case, those additional scenarios must be considered when trying to infer  
464 models from the study of troughs found in empirical data. Another phenomenon that renders  
465 the inference of parameters cumbersome is a possible interference between troughs. Indeed,  
466 when two loci fix neutrally in the population, the genetic diversity in the region between those  
467 loci will be influenced by both fixations and will differ from the diversity expected in the vicinity  
468 of a single fixing locus. As in the case of interference between the fixation of selected alleles  
469 (Weissman and Barton 2012), this should limit the number of independent neutral fixations.  
470 The effect of trough interference is stronger for neighboring troughs, and the probability to  
471 observe close troughs depends on the relative frequency of troughs along the genome, which  
472 itself depends on the distribution of the  $T_{MRCA}$ . In fig. 1d for example, the distribution of  $T_{MRCA}$   
473 has a mode centered around  $4N_c$  (not shown) and correspondingly the nucleotide diversity is  
474 peaked around  $4N_c \mu$ . As a result, we see many regions of the chromosome with a low diversity.  
475 It is likely that those troughs interfere with each other and that they do not correspond to the  
476 profile of an isolated trough. On the other hand, in fig. 1c, the first mode of the  $T_{MRCA}$   
477 distribution is truncated because  $t_c$  is much smaller than  $4N_c$ , and only  $T_{MRCA}$ s equal or close to  
478  $t_c$  are observed (plus all the  $T_{MRCA}$ s corresponding to the second mode centered at  $4N_0$ ). In this  
479 case there is no interference and the (rare) troughs, such as the one in fig. 7, are correctly fitted  
480 by their theoretical expectation. Those considerations imply that, even though we know the  
481 forward in time probability that an allele will fix in  $t_m$  generations, it is difficult to infer the  
482 parameters of a fixation scenario from a single observed neutral valley of diversity. It appears  
483 therefore difficult to perform model selection from a single trough signal, i.e., to decide  
484 whether a particular trough is due to selection or demographic effects, because alternative  
485 demographic scenarios that we did not consider here could also lead to similar signals. In

486 principle, if several troughs of diversity were observed in a genome, one could use the  
487 distribution of trough shapes expected under a given simple demographic model and a  
488 distribution of fitness effect to compare neutral and selection models under a likelihood  
489 framework.

490 In conclusion, our results suggest that any empirical valley of diversity found in empirical data  
491 can be reproduced neutrally with a population contraction using appropriate parameters. One  
492 could argue that this identifiability problem disappears once the true evolutionary history is  
493 correctly inferred. However, inferring the true demographic history requires precise knowledge  
494 about how selection has shaped genomic diversity (Johri *et al.* 2020). In humans, for instance, it  
495 has been estimated that roughly 95 % of genomic diversity is affected by some form of non-  
496 neutral forces such as background selection or biased gene conversion (Pouyet *et al.* 2018)  
497 potentially biasing demographic inference (Ewing and Jensen 2016). These considerations  
498 indicate that genome scans in search for signals of adaptation might be subject to stronger  
499 false positive rates than previously thought. We thus believe that despite current advances  
500 using supervised machine learning or similar approaches (Schridder and Kern 2018), it remains  
501 important to further study the effect of neutral fixations in various demographic scenarios using  
502 localized genomic approaches such as the present analytical work (Johri *et al.* 2021b); as well as  
503 with controlled experiments on real living organisms where both the selected locus and the  
504 population history are known (Orozco-terWengel *et al.* 2012). Such work will be critical in order  
505 to develop more appropriate evolutionary null models for statistical inference (Hahn 2008;  
506 Johri *et al.* 2020).

## 507 Appendix

### 508 A1. Coalescence distribution after a contraction

509 We want to determine the coalescence time of two lineages in a population that experienced a  
510 contraction  $t_m$  generations ago, from a diploid size  $N_0$  to  $N_c$ . As we go backward in time, the  
511 coalescence rate switches from  $(2N_c)^{-1}$  to  $(2N_0)^{-1}$  at  $T = t_c$ . The probability distribution might  
512 still be approximated by a piecewise exponential density:

$$\begin{aligned}
 f_0(T) &= \frac{1}{2N_c} \exp\left(-\frac{T}{2N_c}\right) \text{ for } 0 < T < t_c \\
 &= \frac{1}{2N_0} \exp\left(-\frac{t_c}{2N_c}\right) \exp\left(-\frac{T-t_c}{2N_0}\right) \text{ for } T \geq t_c
 \end{aligned}$$

The corresponding expectation for this distribution is

$$\begin{aligned}
 E[T] = T_0 &= \int_0^{\infty} T f_0(T) dT \\
 &= 2N_0 e^{-t_c/2N_c} + 2N_c(1 - e^{-t_c/2N_c})
 \end{aligned}$$

## A2. Average frequency of an allele fixing in exactly $t_m$ generations

In this section time is counted forward from the mutation, which appears after the contraction, so that during the fixation the diploid population size is constant and equal to  $N_c$ . We condition on the fixation time  $t_m$  of the mutant. We define the trajectory of a mutant as the list of frequencies at all generations:  $\{x_t\} = (x_0, x_1, \dots, x_{t_m-1}, x_{t_m})$ . We assume that the mutant fixes in exactly  $t_m$  generations, starting from a frequency  $p_0$ , i.e.  $x_0 = p_0$ ,  $0 < x_{t_m-1} < 1$  and  $x_{t_m} = 1$ . The probability that the mutant follows a given trajectory might be expressed as the product of the transition probabilities

$$P(\{x_t\}) = \prod_{t=0}^{t_m-1} P(i, t \rightarrow j, t+1 \mid \text{fix in } t_m, p_0)$$

For an unconditional Wright Fisher model,  $P(i, t \rightarrow j, t+1)$  is the probability to have  $j$  copies of the new allele at  $t+1$  given that there were  $i$  copies at  $t$ . We note  $P_t(i \rightarrow j)$  for brevity. If we only consider trajectories fixing in exactly  $t_m$  generations and starting from a number  $2N_c p_0$  of copies at  $t=0$ , then the transition probabilities are not equal to the transitions of the unconditional Wright-Fisher model. However, thanks to Bayes theorem, we can write

$$\begin{aligned}
 P_t(i \rightarrow j \mid \text{fix in } t_m, p_0) &= \frac{P_t(\text{fix in } t_m \mid i \rightarrow j, p_0) P_t(i \rightarrow j \mid p_0)}{P(\text{fix in } t_m \mid p_0)} \\
 &= \frac{P(\text{fix in } t_m \mid j_{t+1}) P_t(i \rightarrow j)}{P(\text{fix in } t_m \mid p_0)} \quad (S1)
 \end{aligned}$$

531 From the first to the second line, we use the Markov property. The three terms involved in the  
 532 right-hand side of this equation can be approximated thanks to diffusion theory. In this  
 533 framework, the probability for an allele to fix in  $t_m$  generations, given that there were  $i$  copies  
 534 at time  $t$  is approximately (Ewens 2004, taking the time derivative of eq. 5.39)

$$535 \quad P(\text{fix in } t_m | i_t) = \frac{3}{2N_c} \left(1 - \frac{i}{2N_c}\right) \frac{i}{2N_c} e^{-(t_m-t)/2N_c}$$

536 The term  $P_t(i \rightarrow j)$  is the unconditional binomial transition probability of the Wright Fisher  
 537 model (which does not depend on  $t$ ). In principle, eq. (S1) can be used to compute the exact  
 538 distribution of coalescence times at the focal locus, using eq. (3a). However, the huge number  
 539 of possible trajectories fixing in  $t_m$  generations ( $(2N_c - 1)^{t_m-1}$ ) makes the average over  
 540 trajectories impossible to evaluate numerically. For this reason, we use the approximation in  
 541 eq. (3b).

542 We consider here the probability that the allele has frequency  $x$  at time  $t$ , given that it started  
 543 at frequency  $p_0$  at  $t = 0$ . Again if we only consider trajectories that fix in exactly  $t_m$   
 544 generations, this probability is not equal to the neutral diffusive result. However, similarly to  
 545 the previous section, we can use Bayes theorem:

$$546 \quad P(x_t | \text{fix in } t_m, p_0) = \frac{P(\text{fix in } t_m | x_t)P(x_t | p_0)}{P(\text{fix in } t_m | p_0)}$$

547 From diffusion theory (Ewens 2004, eq. 5.11), we also have

$$548 \quad P(x_t | p_0) = 6p_0(1 - p_0) e^{-t/2N_c} (1 + 5(1 - 2p_0)(1 - 2x)e^{-t/N_c})$$

549 which is a second order expansion of an infinite series involving vanishing exponential terms  
 550 ( $e^{-k(k+1)t/4N_c}$  for all  $k \geq 1$ ). This expansion is thus valid in the limit of large times  $t \gg 2N_c$ . We  
 551 deduce that the probability that an allele fixing in  $t_m$  generations has frequency  $x$  at time  $t$  is

$$552 \quad P(x_t | \text{fix in } t_m, p_0) = 6x(1 - x) (1 + 5(1 - 2p_0)(1 - 2x)e^{-t/N_c})$$

$$553 \quad \text{which yields } E[x_t | \text{fix in } t_m, p_0] = 1/2(1 - (1 - 2p_0)e^{-t/N_c})$$



554 This expression is valid for  $t_m \gg t \gg 2N_c$ , and does not allow to estimate the frequency close  
555 to fixation (we see that  $E[x_t]$  tends to  $1/2$  as time grows). However, invoking a symmetry  
556 argument we may write

$$557 \quad E[x_t | \text{fix in } t_m, p_0] = 1/2(1 - (1 - 2p_0)e^{-t/N_c} + e^{-(t_m-t)/N_c})$$

558 When  $t_m \ll 2N_c$ , we can use a linear approximation for the trajectory (based on the numerical  
559 observations)

$$560 \quad E[x_t | \text{fix in } t_m, p_0] = p_0 + (1 - p_0) \frac{t}{t_m}$$

### 561 **A3. Coalescence distribution at linked loci around a neutral fixation**

562 We now return to the scenario of fig. 2, with a backward in time approach. Using Bayes  
563 theorem, we express the coalescence time of two haplotypes at the linked locus  $T^{(l)}$ ,  
564 conditioning on the coalescence time at the focal locus  $T^{(f)}$

$$565 \quad P(T^{(l)}) = \int_0^{t_m} P(T^{(l)} | T^{(f)}) P(T^{(f)}) dT^{(f)} = E[P(T^{(l)} | T^{(f)})]$$

566 We assume that the linked locus is close to the focal locus on the chromosome, more precisely  
567 that the recombination rate  $r$  is very small  $r \ll 1$ , so that we consider at most one  
568 recombination, occurring on one of the two focal lineages. We distinguish cases where there is  
569 no recombination between  $t = 0$  and  $t = T^{(f)}$ , cases where the allele at the linked locus  
570 recombines (somewhere between  $t = 0$  and  $t = T^{(f)}$ ) onto a haplotype carrying the ancestral  
571 allele at the focal locus, and cases where the allele at the linked locus recombines onto a  
572 haplotype carrying the derived allele at the focal locus. We call the second and third case  
573 homozygous and heterozygous recombination respectively, referring to the zygosity at the focal  
574 locus of the recombining pair of haplotypes (note that are three haplotypes, the two first ones  
575 have a coalescence time  $T^{(f)}$ , and the third one recombines with one of these two). If there is no  
576 recombination, then the coalescence time is the same for both loci,  $T^{(l)} = T^{(f)}$ . To treat the case  
577 with a homozygous recombination, it is convenient to name the haplotypes:  $i$  and  $j$  coalesce at  
578  $T_{ij}^{(f)} = T^{(f)}$  at the focal locus, and  $k$  is a third haplotype, onto which the linked allele recombines

579 (coming from  $i$ ). The linked allele carried by  $j$  stays on the same haplotype (no more than one  
 580 recombination), and after recombining onto  $k$ , the linked allele initially carried by  $i$  also stays on  
 581  $k$  (again, at most one recombination). This implies that those two linked alleles coalesce at  $T_{ij}^{(l)} =$   
 582  $T_{jk}^{(f)}$ . This time is in general different than  $T_{ij}^{(f)}$ , however on average  $T_{jk}^{(f)}$  and  $T_{ij}^{(f)}$  are equal  
 583 (averaging over all possible coalescence trees at the focal locus). This implies that we can treat  
 584 the case with homozygous recombination as if there was no recombination. If there is a  
 585 heterozygous recombination between  $i$  and  $k$ , at some generation between  $t = 0$  and  $t = T^{(f)}$ ,  
 586 then the linked alleles still have not coalesced at  $t = t_m$  because after the recombination one of  
 587 them is linked to a derived focal allele and the other one to an ancestral focal allele (and they  
 588 stay linked because there is at most one recombination). In that case,  $T_{ij}^{(l)}$  is equal to  $t_m$  plus a  
 589 random time given by (on average)  $T_m$ , and is independent of  $T_{ij}^{(f)}$ . Using again Bayes theorem  
 590 and the previous results to write

$$\begin{aligned}
 591 \quad P(T^{(l)} | T^{(f)}) &= P(T^{(l)} | T^{(f)}, \text{ one het. rec. in } [0, T^{(f)}])P(\text{one het. rec. in } [0, T^{(f)}]) \\
 592 &\quad + P(T^{(l)} | T^{(f)}, \text{ no het. rec. in } [0, T^{(f)}])P(\text{no het. rec. in } [0, T^{(f)}]) \\
 593 &= f_m(T^{(l)} - t_m)[1 - P(\text{no het. rec. in } [0, T^{(f)}])] \\
 594 &\quad + \delta(T^{(l)} - T^{(f)})P(\text{no het. rec. in } [0, T^{(f)}])
 \end{aligned}$$

595 Where  $\delta(\cdot)$  is the Dirac delta function, and  $f_m$  is the unconditional coalescence distribution of a  
 596 pair of lineages sampled at  $t = t_m$ , i.e. it is equal to the function  $f_0$  introduced above but  
 597 replacing  $t_c$  by  $t_c - t_m$  (note also that  $f_m(t) = 0$  if  $t < 0$ ). We then have to evaluate the  
 598 probability that there is no heterozygous recombination. At generation  $t$  (counted backward)  
 599 the probability that a linked allele recombines onto a haplotype carrying the ancestral allele at  
 600 the focal locus is  $r(1 - x_t)$ , where  $x_t$  is the frequency of the derived allele at the focal locus,  
 601 we deduce that the probability that there is no heterozygous recombination on either lineage is

$$\begin{aligned}
 602 \quad P(\text{no het. rec. in } [0, T^{(f)}]) &= \prod_{t=1}^{T^{(f)}} (1 - r[1 - x_t])^2 \\
 &\simeq \exp\left(-2r \sum_{t=1}^{T^{(f)}} (1 - x_t)\right)
 \end{aligned}$$

603 This probability depends explicitly on the allele trajectory, which means that rigorously, all the  
 604 calculations should be conditioned on a given trajectory, and then averaged over all  
 605 trajectories. To allow for mathematical tractability, and to avoid heavy expressions, we consider  
 606 that as a good approximation  $x_t = \bar{x}_t$ . Finally we obtain

$$\begin{aligned}
 607 \quad P(T^{(l)}) = & E \left[ \delta(T^{(l)} - T^{(f)}) \exp \left( -2r \sum_{t=1}^{T^{(f)}} (1 - x_t) \right) \right] \\
 608 \quad & + f_m(T^{(l)} - t_m) E \left[ 1 - \exp \left( -2r \sum_{t=1}^{T^{(f)}} (1 - x_t) \right) \right]
 \end{aligned}$$

609 The expectation corresponding to this distribution yields eq. (2).

610

611 **A4. Average coalescence time at a linked locus around a mutation that completed fixation  $t_{\text{fix}}$**   
 612 **generations ago**

613 Thanks to Bayes theorem we can write

$$614 \quad E[T^{(l)}] = E[T^{(l)} | T^{(l)} < t_{\text{fix}}] P(T^{(l)} < t_{\text{fix}}) + E[T^{(l)} | T^{(l)} > t_{\text{fix}}] P(T^{(l)} > t_{\text{fix}})$$

615 i.e. we distinguish coalescence events happening in less than  $t_{\text{fix}}$  generations or more than  $t_{\text{fix}}$

616 generations. In the former case, the coalescence is neutral, unconditional (the fixation is

617 completed) and happens in a population of constant size  $N_c$  which means that

618  $E[T^{(l)} | T^{(l)} < t_{\text{fix}}]$  and  $P(T^{(l)} < t_{\text{fix}})$  can be worked out from the neutral exponential

619 distribution. On the other hand,  $E[T^{(l)} | T^{(l)} > t_{\text{fix}}]$  is equal to  $t_{\text{fix}}$  plus the expectation from eq.

620 (5) which we note here  $E[T^{(l)}](t = t_{\text{fix}})$ . We obtain

$$621 \quad E[T^{(l)}] = 2N_c(1 - e^{-t_{\text{fix}}/2N_c}) + E[T^{(l)}](t = t_{\text{fix}}) e^{-t_{\text{fix}}/2N_c}$$

622 We see that the sweep signal vanishes exponentially with the time elapsed since fixation.

## 623 Acknowledgment

624 This work was partially supported by a Swiss NSF grant No 310030\_188883 to LE. We are  
625 grateful to Montgomery Slatkin, Brian Charlesworth and Jeff Jensen for their helpful comments.

## 626 Bibliography

627

628 Andolfatto P., and M. Przeworski, 2000 A genome-wide departure from the standard neutral  
629 model in natural populations of *Drosophila*. *Genetics* 156: 257–268.

630 Austerlitz F., B. Jung-Muller, B. Godelle, and P.-H. Gouyon, 1997 Evolution of coalescence  
631 times, genetic diversity and structure during colonization. *Theor. Popul. Biol.* 51: 148–  
632 164.

633 Barton N. H., 1995 Linkage and the limits to natural selection. *Genetics* 140: 821–841.

634 Charlesworth B., M. T. Morgan, and D. Charlesworth, 1993 The effect of deleterious mutations  
635 on neutral molecular variation. *Genetics* 134: 1289–1303.

636 Charlesworth D., B. Charlesworth, and M. T. Morgan, 1995 The pattern of neutral molecular  
637 variation under the background selection model. *Genetics* 141: 1619–1632.

638 Charlesworth B., 2013 Background selection 20 years on: the Wilhelmine E. Key 2012  
639 invitational lecture. *J. Hered.* 104: 161–171.

640 Charlesworth B., 2020 How Good Are Predictions of the Effects of Selective Sweeps on Levels  
641 of Neutral Diversity? *Genetics* 216: 1217–1238.

- 642 Corbett-Detig R. B., D. L. Hartl, and T. B. Sackton, 2015 Natural selection constrains neutral  
643 diversity across a wide range of species. *PLoS Biol.* 13: e1002112.
- 644 Crisci J. L., Y.-P. Poh, S. Mahajan, and J. D. Jensen, 2013 The impact of equilibrium  
645 assumptions on tests of selection. *Front. Genet.* 4: 235.
- 646 Edmonds C. A., A. S. Lillie, and L. Luca Cavalli-Sforza, 2004 Mutations arising in the wave  
647 front of an expanding population. *Proc. Natl. Acad. Sci. U. S. A.* 101: 975–979.
- 648 Ewens W. J., 2004 *Mathematical Population Genetics: I. Theoretical Introduction*. Springer,  
649 New York, NY.
- 650 Ewing G. B., and J. D. Jensen, 2016 The consequences of not accounting for background  
651 selection in demographic inference. *Mol. Ecol.* 25: 135–141.
- 652 Excoffier L., N. Marchi, D. A. Marques, R. Matthey-Doret, A. Gouy, *et al.*, 2021 fastsimcoal2:  
653 demographic inference under complex evolutionary scenarios. *Bioinformatics*.  
654 <https://doi.org/10.1093/bioinformatics/btab468>
- 655 Excoffier L., M. Foll, and R. J. Petit, 2009 Genetic Consequences of Range Expansions. *Annu.*  
656 *Rev. Ecol. Evol. Syst.* 40: 481–501.
- 657 Galtier N., and M. Rousselle, 2020 How Much Does Ne Vary Among Species? *Genetics* 216:  
658 559–572.
- 659 Griffiths R. C., and S. Tavaré, 1994 Ancestral inference in population genetics. *Stat. Sci.* 9: 307–  
660 319.
- 661 Hahn M. W., 2008 Toward a selection theory of molecular evolution. *Evolution* 62: 255–265.

- 662 Hallatschek O., and D. R. Nelson, 2008 Gene surfing in expanding populations. *Theor. Popul.*  
663 *Biol.* 73: 158–170.
- 664 Jensen J. D., Y. Kim, V. B. DuMont, C. F. Aquadro, and C. D. Bustamante, 2005 Distinguishing  
665 between selective sweeps and demography using DNA polymorphism data. *Genetics*  
666 170: 1401–1410.
- 667 Jensen J. D., B. A. Payseur, W. Stephan, C. F. Aquadro, M. Lynch, *et al.*, 2019 The importance  
668 of the Neutral Theory in 1968 and 50 years on: A response to Kern and Hahn 2018.  
669 *Evolution* 73: 111–114.
- 670 Johri P., B. Charlesworth, and J. D. Jensen, 2020 Toward an Evolutionarily Appropriate Null  
671 Model: Jointly Inferring Demography and Purifying Selection. *Genetics* 215: 173–192.
- 672 Johri P., B. Charlesworth, E. K. Howell, M. Lynch, and J. D. Jensen, 2021a Revisiting the  
673 Notion of Deleterious Sweeps. *Genetics*. <https://doi.org/10.1093/genetics/iyab094>
- 674 Johri P., K. Riall, H. Becher, L. Excoffier, B. Charlesworth, *et al.*, 2021b The Impact of  
675 Purifying and Background Selection on the Inference of Population History: Problems  
676 and Prospects. *Mol. Biol. Evol.* 38: 2986–3003.
- 677 Kaiser V. B., and B. Charlesworth, 2009 The effects of deleterious mutations on evolution in  
678 non-recombining genomes. *Trends Genet.* 25: 9–12.
- 679 Kaplan N. L., R. R. Hudson, and C. H. Langley, 1989 The “hitchhiking effect” revisited.  
680 *Genetics* 123: 887–899.
- 681 Kingman J. F. C., 1982a The coalescent. *Stochastic Process. Appl.* 13: 235–248.

- 682 Kingman J. F. C., 1982b On the genealogy of large populations. *J. Appl. Probab.* 19A: 27–43.
- 683 Klopstein S., M. Currat, and L. Excoffier, 2006 The fate of mutations surfing on the wave of a  
684 range expansion. *Mol. Biol. Evol.* 23: 482–490.
- 685 Maruyama T., and M. Kimura, 1975 Moments for sum of an arbitrary function of gene frequency  
686 along a stochastic path of gene frequency change. *Proc. Natl. Acad. Sci. U. S. A.* 72:  
687 1602–1604.
- 688 Mathew L. A., and J. D. Jensen, 2015 Evaluating the ability of the pairwise joint site frequency  
689 spectrum to co-estimate selection and demography. *Front. Genet.* 6: 268.
- 690 Maynard Smith J., and J. Haigh, 1974 The hitch-hiking effect of a favorable gene. *Genet. Res.*  
691 23: 23–35.
- 692 Nicolaisen L. E., and M. M. Desai, 2013 Distortions in genealogies due to purifying selection  
693 and recombination. *Genetics* 195: 221–230.
- 694 O’Fallon B. D., J. Seger, and F. R. Adler, 2010 A continuous-state coalescent and the impact of  
695 weak selection on the structure of gene genealogies. *Mol. Biol. Evol.* 27: 1162–1172.
- 696 Orozco-terWengel P., M. Kapun, V. Nolte, R. Kofler, T. Flatt, *et al.*, 2012 Adaptation of  
697 *Drosophila* to a novel laboratory environment reveals temporally heterogeneous  
698 trajectories of selected alleles. *Mol. Ecol.* 21: 4931–4941.
- 699 Peischl S., I. Dupanloup, M. Kirkpatrick, and L. Excoffier, 2013 On the accumulation of  
700 deleterious mutations during range expansions. *Mol. Ecol.* 22: 5972–5982.

- 701 Peischl S., and L. Excoffier, 2015 Expansion load: recessive mutations and the role of standing  
702 genetic variation. *Mol. Ecol.* 24: 2084–2094.
- 703 Pouyet F., S. Aeschbacher, A. Thiéry, and L. Excoffier, 2018 Background selection and biased  
704 gene conversion affect more than 95% of the human genome and bias demographic  
705 inferences. *Elife* 7: e36317.
- 706 Pouyet F., and K. J. Gilbert, 2019 Towards an improved understanding of molecular evolution:  
707 the relative roles of selection, drift, and everything in between. *arXiv [q-bio.PE]*.
- 708 Rogers R. L., T. Bedford, A. M. Lyons, and D. L. Hartl, 2010 Adaptive impact of the chimeric  
709 gene *Quetzalcoatl* in *Drosophila melanogaster*. *Proc. Natl. Acad. Sci. U. S. A.* 107:  
710 10943–10948.
- 711 Rousselle M., M. Mollion, B. Nabholz, T. Bataillon, and N. Galtier, 2018 Overestimation of the  
712 adaptive substitution rate in fluctuating populations. *Biol. Lett.* 14.  
713 <https://doi.org/10.1098/rsbl.2018.0055>
- 714 Schrider D. R., and A. D. Kern, 2018 Supervised Machine Learning for Population Genetics: A  
715 New Paradigm. *Trends Genet.* 34: 301–312.
- 716 Slatkin M., 1996 Gene genealogies within mutant allelic classes. *Genetics* 143: 579–587.
- 717 Sousa V., S. Peischl, and L. Excoffier, 2014 Impact of range expansions on current human  
718 genomic diversity. *Curr. Opin. Genet. Dev.* 29: 22–30.
- 719 Tajima F., 1990 Relationship between DNA polymorphism and fixation time. *Genetics* 125:  
720 447–454.



- 721 Tavaré S., 1984 Line-of-descent and genealogical processes, and their applications in population  
722 genetics models. *Theor. Popul. Biol.* 26: 119–164.
- 723 Teshima K. M., G. Coop, and M. Przeworski, 2006 How reliable are empirical genomic scans for  
724 selective sweeps? *Genome Res.* 16: 702–712.
- 725 Thornton K. R., and J. D. Jensen, 2007 Controlling the false-positive rate in multilocus genome  
726 scans for selection. *Genetics* 175: 737–750.
- 727 Wares J. P., 2009 Evolutionary dynamics of transferrin in *Notropis*. *J. Fish Biol.* 74: 1056–1069.
- 728 Weissman D. B., and N. H. Barton, 2012 Limits to the rate of adaptive substitution in sexual  
729 populations. *PLoS Genet.* 8: e1002740.
- 730 Zhao L., M. Lascoux, A. D. J. Overall, and D. Waxman, 2013 The characteristic trajectory of a  
731 fixing allele: a consequence of fictitious selection that arises from conditioning. *Genetics*  
732 195: 993–1006.