1    # Strong neutral sweeps occurring during a population contraction

2

3    Antoine Moinet[1,2,3], Flavia Schlichta[2,3], Stephan Peischl*[1,2,4] and Laurent
4    Excoffier*[2,3]

5

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

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

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

11   4.   Corresponding author, *stephan.peischl@bioinformatics.unibe.ch*;

12   *. These authors contributed equally.

13

15

## Abstract

17   A strong reduction in diversity around a specific locus is often interpreted as a recent rapid

18   fixation of a positively selected allele, a phenomenon called a selective sweep. Rapid fixation of

19   neutral variants can however lead to similar reduction in local diversity, especially when the

20   population experiences changes in population size, e.g., bottlenecks or range expansions. The

21   fact that demographic processes can lead to signals of nucleotide diversity very similar to

22   signals of selective sweeps is at the core of an ongoing discussion about the roles of

23   demography and natural selection in shaping patterns of neutral variation. Here we

24   quantitatively investigate the shape of such neutral valleys of diversity under a simple model of

25   a single population size change, and we compare it to signals of a selective sweep. We

26   analytically describe the expected shape of such "neutral sweeps" and show that selective

27  sweep valleys of diversity are, for the same fixation time, wider than neutral valleys. On the

28  other hand, it is always possible to parametrize our model to find a neutral valley that has the

29  same width as a given selected valley. Our findings provide further insight in how simple

30  demographic models can create valleys of genetic diversity similar to those attributed to

31  positive selection.


## Introduction

33  Past demography and natural selection play a critical role in shaping extant genetic diversity. A

34  central question in population genetics is to quantify their respective impact on observed

35  genomic diversity. Because selection interferes with demographic estimates and vice versa,

36  estimation of one of these two components is difficult without accounting for the other

37  (Charlesworth *et al.* 1993, 1995; Kaiser and Charlesworth 2009; O'Fallon *et al.* 2010;

38  Charlesworth 2013; Nicolaisen and Desai 2013; Johri *et al.* 2020, 2021b). Moreover, the relative

39  importance of demography and selection as determinants of genome wide diversity is currently

40  hotly debated, and may vary extensively among species (Corbett-Detig *et al.* 2015; Rousselle *et*

41  *al.* 2018; Pouyet and Gilbert 2019; Galtier and Rousselle 2020). It has been shown that selection

42  and demography can leave very similar footprints on the genetic diversity of a population

43  (Andolfatto and Przeworski 2000; Teshima *et al.* 2006; Thornton and Jensen 2007; Johri *et al.*

44  2021a). Disentangling the effects of demography and selection is therefore crucial to avoid

45  erroneous inference of evolutionary scenarios from genomic data (Jensen *et al.* 2005; Wares

46  2009; Mathew and Jensen 2015; Johri *et al.* 2020).

47  Hard selective sweeps lead to valleys of strongly reduced diversity around positively selected

48  sites due to the hitchhiking of linked neutral loci (Maynard Smith and Haigh 1974), such

49  observations of strong depletions of diversity in some genomic regions are often interpreted as

50  due to past episode of positive selection, because the probability to observe a fast fixation of a

51  neutral variant in a population of constant size is extremely low. However, during a range

52  expansion for instance, some neutral or even mildly deleterious mutations can go quickly to

53  fixation due to the low effective size of populations on the front of the range (Edmonds *et al.*

54  2004; Klopfstein *et al.* 2006; Hallatschek and Nelson 2008; Peischl *et al.* 2013), a phenomenon

55  termed allele surfing (Klopfstein *et al.* 2006). Theoretical studies have shown that the average

56  neutral diversity on the wave front decays exponentially as the range expands (Hallatschek and

57  Nelson 2008), similarly to what happens when a population experiences a sudden decay of the

58  population size, i.e. a population contraction, due to a drastic change in the environment for

59  example. In both cases, a mutation appearing when the population size is shrinking might go

60  quickly to fixation, inducing a strong decrease of diversity in the surrounding genomic region,

61  whereas the average level of diversity might stay quite high depending on the strength and the

62  duration of the contraction. As a result, the coalescent tree of alleles sampled in a population

63  with strongly reduced effective population size will have short external branches, and long

64  internal branches, depending on the parameters of the model (Excoffier *et al.* 2009). The

65  average site frequency spectrum associated to such a tree resembles a neutral SFS, but with a

66  lack of rare alleles and an excess of high frequency sites, i.e. it becomes "flatter" (Sousa *et al.*

67  2014; Peischl and Excoffier 2015). The footprint left by the rapid fixation of a neutral allele on

68  the surrounding genomic diversity, might thus be like that of a positively selected allele

69  sweeping through a constant size population.

70  The expected shape of nucleotide diversity in genomic regions surrounding a site undergoing a

71  rapid neutral fixation has been investigated analytically and numerically. Tajima (1990) studied

72  the reduction of diversity during a neutral fixation at a given recombination distance from the

73  fixing site. His results rely on rigorous mathematical arguments based on diffusion theory, but

74  no closed form solution is provided for the shape of a neutral sweep. Johri *et al.* (2021a)

75  described the valley of diversity occurring around a neutral fixation using an approach

76  introduced for selective sweeps, assuming that the evolution of the allele frequency is that of a

77  selected allele except in the initial stochastic phase. Here, we extend this work by inferring the

78  dynamics of fixation of neutral alleles after a population contraction and we examine their

79  effects on neighboring regions of the genome. We provide an analytical result for the expected

80  coalescence time as a function of the recombination distance from the locus undergoing a fast

81  fixation. Importantly, our results apply regardless of the process driving the allele going to

82  fixation (neutrality, positive selection, background selection), as it only relies on the typical

83  trajectory of an allele going to fixation in a given time, even though this trajectory differs

84    depending on the underlying driver of this fixation (i.e., neutrality or selection). We compare

85    our results against simulations and find that they hold for a wide range of realistic parameter

86    combinations. We compare our results about the signature of neutral sweeps to patterns

87    expected under selective sweeps and discuss potential differences between the signatures that

88    could potentially allow us to discriminate between neutral and selective processes for a given

89    demographic scenario. Finally, we investigate the similarity between the genomic signature of

90    an allele going to fixation either selectively or neutrally and observe that a selective sweep

91    signal can in principle be replicated in a neutral model with an appropriate choice of

92    demographic parameters. We conclude that strong diversity depletions in the genome of a

93    population, often attributed to the effect of positive selection, can be obtained with

94    demographic effects only, and we call for caution when trying to detect signals of adaptation

95    from genomic data, adding support to previous studies reaching similar conclusions (Thornton

96    and Jensen 2007; Crisci *et al.* 2013; Jensen *et al.* 2019).

97    ## Model

98    We model here the effect of an instantaneous population contraction on genomic diversity.

99    Throughout the whole manuscript, time is measured backwards. We assume that $t_c$ generations

100   before the present, the population size instantaneously dropped from $N_0$ diploid individuals to

101   $N_c$ individuals with $N_c < N_0$. We assume a standard coalescent model (Kingman 1982a; b) with

102   discrete non-overlapping generations, random mating, monoecious individuals, and no

103   selection. Two haplotypes sampled in the current population at time $t = 0$ have, as we go

104   backwards in time, a constant probability $(2N_c)^{-1}$ of coalescing at each generation, for the first $t_c$

105   generations, and then this probability switches to $(2N_0)^{-1}$ as we enter the ancestral

106   uncontracted population. We can approximate the distribution of coalescence time $T$ of these

107   two haplotypes as a piecewise exponential distribution (see Appendix) with expected value:

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

108   We see that the expected coalescence time decreases exponentially with the age of the

109   contraction $t_c$ and that it approaches $2N_c$ for a very old contraction. Coalescence times cannot

110     be measured directly from empirical data, but they are closely related to nucleotide diversity $\pi$.

111     Under the infinitely many sites model, the number of nucleotide differences between two

112     homologous DNA segments is proportional to their coalescence time $T$ as $\pi = 2\mu T$, where $\mu$ is

113     the total mutation rate for the whole segment. Multiplying eq. (1) by $2\mu$ shows that an

114     instantaneous population contraction leads to an exponential decrease of the expected

115     nucleotide diversity along the genome with the age of the contraction $t_c$. However, it does not

116     inform us on the distribution of nucleotide diversity $\pi$ along the genome, or on spatially

117     correlated patterns of diversity such as local depletion or excess of diversity relative to the

118     expectation.

119

120

121     Fig. 1 shows the evolution of the distribution of $\pi$ as a function of the time $t_c$ elapsed since the

122     contraction. For $t_c = 0$, there is no contraction, and the population size remains constant and

123     equal to $N_0$. In this case we see (Fig. 1a,1b, $t_c = 0$) that the distribution of $\pi$ is symmetric and

124     centered at $E[\pi] = 4N_0\mu$. For an older contraction, we see that the distribution is not only

125     shifted to lower values of diversity as expected from eq. (1), but that it also becomes strongly

126     peaked around $\pi = 4N_c\mu$. This bimodality of the distribution can be understood intuitively in the

127     following way. There are two possible types of coalescent trees for haplotypes sampled after

128     the population contraction (note that the tree depends on the locus considered because of

129     recombination). Indeed, the most recent common ancestor (MRCA) of the sample lived either

130     before the contraction ($T_{MRCA} > t_c$), or after the contraction ($T_{MRCA} < t_c$). In the former case, the

131     tree at this locus has long inner branches and short outer branches, whereas in the latter case,

132     the tree is essentially a (short) neutral tree corresponding to a population of constant size $N_c$

133     (Excoffier *et al.* 2009). Both types of trees occur at different loci and correspond to the two

134     observed modes in the distribution of the nucleotide diversity along the chromosome. The

135     precise shape of the distribution of nucleotide diversity across sites depends on the relative

136     frequency of both types of trees, which itself depends on the age of the contraction $t_c$. For a

137     sample of size two, the probability that the MRCA lived after the contraction, that is, $T_{MRCA} < t_c$

138    is $1 - e^{-t_c/2N_c}$. For a larger sample of haplotypes, there is no closed form solution for this

139    probability, but the trees rooted after the contraction are rare for $t_c \ll 2N_c$ and very frequent

140    when $t_c \gg 2N_c$ (Tavaré 1984). Therefore, the evolution of the distribution of $\pi$ for increasing

141    contraction age $t_c$ appears to be a transition from a unimodal distribution centered at $4N_0\mu$ to

142    another unimodal distribution centered at $4N_c\mu$, with both modes coexisting for intermediate

143    ages (Fig. 1). This bimodality has been pointed out previously in the context of population

144    bottlenecks (Austerlitz *et al.* 1997); however, those studies mainly focused on long duration

145    bottlenecks (the effect of a contraction or a bottleneck on nucleotide diversity is the same

146    provided that the bottleneck is not yet finished, or that it finished very recently so that the

147    effect of population recovery is negligible). In the present work, we investigate the effect of

148    short contractions on the genetic diversity and make the claim that this short contraction

149    regime is of particular interest as it can lead, such as in Fig. 1c, to genomic signatures similar to

150    those generated by positive selection acting on a few sites in an otherwise neutral genome.

151    More specifically, we want to quantitatively describe the reduction of diversity along the

152    genome that is observed around a locus with a small $T_{MRCA}$ (such as in Fig. 1c in the regions

153    around 10-11 and 19-20 Mb), where we observe a valley or trough of diversity. Akin to what is

154    done for selective sweeps, we consider the (neutral) fast fixation of an allele and analyze the

155    impact of hitchhiking on the genetic diversity of neighboring loci, and we refer to this process as

156    a neutral sweep.

157    To investigate neutral sweeps in our model, we consider the following scenario: $t_m$ generations

158    ago a mutation occurred at a single site on the chromosome, which we call the focal site. We

159    further assume that this mutation has just fixed in the population, i.e., that it was segregating at

160    a frequency strictly lower than one in the last generation (at $t = 1$) and has now (at $t = 0$) a

161    frequency equal to one. We assume that the population contraction occurred $t_c$ generations

162    ago, with $t_c \geq t_m$. As the mutant enters the population as a single allelic copy at the focal locus,

163    defined as a non-recombining region surrounding the focal site, this copy is a common ancestor

164    for all the copies ($2N_c$) present at fixation. However, it is not necessarily the most recent

165    common ancestor. Fig.2 shows a sketch of our model to help visualize how recombination can

166    maintain diversity at linked loci around a locus where a new mutation quickly fixed in the

167    population.

168

# Results

## Average coalescence time at a linked locus

171    We can calculate the expected coalescence time $T^{(l)}$ of two randomly sampled haplotypes at a

172    linked locus as a function of the recombination rate $r$ from the focal locus. The idea is to

173    consider two haplotypes with a given coalescence time $T^{(f)}$ at the focal locus, and then follow

174    the genealogy of the gene copies carried by these two haplotypes at the linked locus backward

175    in time, while considering possible recombination events. The expected coalescent time at the

176    linked locus is then

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

177    where $\bar{x}_t$ is the average frequency of the mutant (derived) allele at the focal locus at time $t$

178    counting backward from present. A detailed derivation of this equation is given in Appendix A4.

179    The first term of the right-hand side of eq. (2) corresponds to cases where lineages escape the

180    neutral sweep due to recombination, and still have not coalesced after $t_m$ generations. In this

181    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

182    before the lineages coalesce, due to the contraction that happened $t_c$ - $t_m$ generations before

183    the focal mutation. The second term of the right-hand side of eq. (2) corresponds to cases

184    where the lineages cannot escape the sweep and are forced to coalesce at a time $T^{(l)} \le t_m$.

## Distribution of coalescence times at the focal locus

186    To evaluate eq. (2), we need to determine the probability distribution of the pairwise

187    coalescence times $T^{(f)}$ at the focal locus, as well as the expected frequency trajectory of the

188    derived allele. Even though this allele fixes neutrally in a population of constant size (the

189    contraction occurs prior to the mutation), the distribution of coalescent times at the focal locus

190    $T^{(f)}$ departs from the usual exponential distribution for a neutral coalescent process because the

191    allele fixes in exactly $t_m$ generations, and hence the coalescence time for a randomly chosen

192    pair of haplotypes is at most $t_m$. Slatkin (1996) investigated the coalescent process within a

193    "mutant allelic class" that originated from a single mutation at a given time in the past. He

194    derived exact analytical results for the average pairwise coalescence time, but the coalescence

195    distribution itself can only be expressed with multidimensional integrals and obtaining a closed

196    form expression does not appear feasible. We therefore use a different approach: given a

197    particular fixation trajectory of the mutant allele, i.e. given the number of mutant copies $N_\mu$ at

198    each generation between $t = 0$ and $t = t_m$, we can express the coalescence time distribution

199    within the mutant allelic class, using the result of a coalescent in a population with a time-

200    dependent (but deterministic) size $N_\mu(t)$ (Griffiths and Tavaré 1994). Averaging over all

201    possible trajectories of the mutation, we obtain:

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

202    where $x_t = N_\mu(t)/(2N_c)$ is the frequency of the mutant $t$ generations from fixation, and

203    $P(\{x_t\})$ is the probability of a given trajectory. $P(\{x_t\})$ can be evaluated (see Appendix A2) and

204    the sum in eq. (3a) can in principle be computed numerically; however, the number of

205    trajectories to consider is prohibitive. As a first approximation, we can replace $x_t$ by its

206    expectation $\overline{x}_t$, i.e., we neglect the fluctuations of the trajectory around the mean to obtain

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

207    The last step is to determine the average trajectory of an allele fixing in exactly $t_m$ generations.

208    Zhao *et al.* (2013) as well as Maruyama and Kimura (Maruyama and Kimura 1975) have

209    investigated the characteristic trajectory of an allele fixing in a given time but they do not

210    provide a closed form solution. Here, we use a different approach (also based on diffusion

211    theory to obtain an approximation for the average trajectory of an allele fixing in exactly $t_m$

212    generations, starting from a frequency $p_0$. As detailed in the Appendix A2, we obtain

$$\overline{x}_t = 1/2\big(1 - (1 - 2p_0)e^{-(t_m-t)/N_c} + e^{-t/N_c}\big), \qquad (4a)$$

213      which is valid for $t_\mathrm{m} \gg 2N_c$. For very fast fixations, i.e., when $t_\mathrm{m} \ll 2N_c$, the frequency of the

214      allele increases approximately linearly as

$$\overline{x}_t = 1 - (1 - p_0)\frac{t}{t_\mathrm{m}}. \qquad\qquad (4b)$$

215      We remind the reader that $t$ is counted backwards from fixation. Fig. 3 compares equations (4a)

216      and (4b) to trajectories obtained from simulations of a Wright-Fisher diploid population. We

217      find good agreement between the simulations and the analytical results. Importantly, the

218      typical neutral trajectory for large values of the fixation time has an "inverse-sigmoid shape"

219      (Fig. 3c), contrary to the typical sigmoid trajectory of a positively selected allele going to fixation

220      in a constant size population (see Fig. 5a). This neutral trajectory occurs because, conditional on

221      non-loss, neutral alleles need to quickly escape loss at the beginning and remain at

222      intermediate frequencies to stay away from both fixation and loss until they eventually fix in

223      the population at $t = 0$ (*i.e.* in exactly $t_\mathrm{m}$ generations). Fig. 3e-3h also shows the coalescence

224      time distribution for several values of the fixation time $t_\mathrm{m}$. The comparison of the distribution of

225      pairwise coalescence time with numerical simulations of a Wright-Fisher model shows that our

226      approximation eq. (3b) is quite accurate but overestimates the probability of coalescence for

227      large coalescence times when $t_\mathrm{m}$ is small (Fig. 3d). Notably, coalescence (simulated or

228      theoretical) is more probable at large times *(i.e.* when the mutant appeared) for short fixation

229      times (Fig. 3d), whereas it is more probable at small times (i.e. close to fixation) for large

230      fixation times (Fig. 3e). The coalescence rate within the mutant allelic class is given by the

231      inverse of the number of mutant copies and is for all values of the fixation time slightly more

232      than $1/2N_c$ at the first generation. However, when the fixation time is short (Fig. 3e), there is a

233      fast increase of the coalescence rate backwards in time, and many lineages are forced to

234      coalesce at $t = t_\mathrm{m}$. When the fixation time is large (Fig. 3h), the coalescence rate also increases

235      backwards in time, but the increase is much slower. In that case, most coalescence events

236      happen in much less than $t_\mathrm{m}$ generations, so that the early increase in frequency of the mutant

237      has almost no influence on the coalescence distribution.

238

239 ## Effect of a neutral sweep on linked diversity

240 Combining equations (3b), (4a) with eq. (2) allows us to get an approximation for the average

241 coalescence time at linked loci. Since the derivation of eq. (2) assumes that there is at most one

242 recombination event in the genealogy of a randomly chosen pair of gene copies, we expect it to

243 be only accurate for small values of the recombination rate $r$. For large values of $r$ we use a

244 heuristic approach combining the result of eq. (2), which is accurate for small $r$, and the

245 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

246 in eq. (1). We fit the trough of diversity with an exponential function of the form:

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

247 where the coefficients $c = 1 - E\big[\mathrm{T}^{(f)}\big]/T_0$ and $a = 2E[(t_m + T_m - \mathrm{T}^{(f)})\sum_{t=1}^{\mathrm{T}^{(f)}}(1 - \overline{x}_t)]/(T_0 -$

248 $E\big[\mathrm{T}^{(f)}\big])$ are obtained by imposing that eqs. (2) and (5) coincide for small values of $r$ (using a

249 linear expansion in $r$). On Fig. 4 we compare the result of eq. (5) to Wright-Fisher simulations

250 with two recombining loci. We see in Fig. 4a that the exponential function fits the data

251 accurately at large values of the recombination distance, but that the fit is biased for intermediate

252 values of $r$. In Fig. 4b we see that the approximation is very good for low values of the

253 recombination distance, although there still is a slight bias. This discrepancy at small $r$ can be

254 corrected (solid lines in Fig. 4) if we use numerical estimations of $\overline{x}_t$ and $\mathrm{P}\big(\mathrm{T}^{(f)}\big)$, instead of eqs.

255 (4) and (3b), to evaluate eq. (5).

256

257 We observe, as expected, on Fig. 4 that the troughs of diversity induced by neutral sweeps are

258 wider and deeper for short fixation times. Similarly to what happens after a selective sweep,

259 there is less opportunity for linked loci to escape the sweep by recombination and maintain

260 diversity when the fixation is fast. In addition, the diversity level at the center of the valley is

261 given by the average coalescence time at the focal locus, which quickly decreases for small

262 fixation times $t_m$.

10

263 Comparison of neutral sweeps and selective sweeps

264 Since we did not make any assumption regarding the process driving the mutant allele to

265 fixation when deriving the average coalescence time at linked loci (eq. (2)) and the coalescence

266 time distribution at the focal locus (eq. (3b)), our framework allows us to directly compare the

267 signatures of different processes that can drive mutations to fixation in a given number of

268 generations. We illustrate this by comparing the effect of neutral and hard selective sweeps on

269 linked diversity. Later we will discuss how neutral sweeps compare to a larger variety of

270 scenarios (e.g. background selection, small selection coefficients, or dominant alleles). Here we

271 assume that the neutral and selected fixations occurred over the same time interval, that is in

272 both cases in exactly $t_m$ generations. The selected fixation is assumed to be codominant ($h$=0.5)

273 and occurs on an autosomal locus in a randomly mating diploid population of constant size $N_1$,

274 and we consider a strong selection strength ($2N_1s >> 1$) so that the allele frequency follows the

275 deterministic trajectory

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

276

277 where the fixation time is given by $t_m(s) = 2\log(4N_1 s)/s$ (Barton 1995). Then combining eqs. (5),

278 (3b) and (6), we can compute the average coalescence time at linked loci as a function of the

279 recombination distance $r$ to the focal locus, after replacing $T_m$, the average coalescence time at $t$

280 $= t_m$, by $2N_1$ in eq. (5) and $N_c$ by $N_1$ in eq. (3b). This approach yields results similar to

281 Charlesworth (2020), where the author investigated signals of selective sweeps correcting for

282 coalescent events that happen during the sweep, thus going beyond the common assumption of a

283 star tree structure at the focal locus. For sake of simplicity in the neutral case, we consider that

284 the mutant appeared at the time of the contraction, i.e. $t_m = t_c$. Furthermore, we will assume that

285 the average coalescence times (and consequently the genetic diversity) are equal in both

286 scenarios, i.e. that $T_0 = 2N_1$ which implies that

287
$$N_0(\mathrm{t_m}) = (N_1 - N_c)\,e^{t_m/2N_c} + N_c\,.\qquad\qquad (7)$$

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

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

290

291

292

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

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

295 allele in the sense that when one is increasing, the other is "resting" and vice versa. These

296 different trajectories translate into different coalescence distributions at the focal locus (Fig.

297 5b). If selection drives the fixation of the mutation, the distribution of coalescence time is

298 peaked at large coalescence times. In contrast, in the neutral case the distribution is skewed

299 towards small coalescence times. Correspondingly, the coalescence tree for the selected case

300 has a star-like structure (Hermisson and Pennings 2017), whereas the tree for the neutral case

301 has shorter outer branches. Therefore, for a given recombination distance, there will be fewer

302 recombinations on the neutral tree because it has a much smaller total length. As

303 recombination helps maintain diversity at linked loci, we would expect neutral troughs of

304 diversity to be wider than in the selected case. However, this is at odds with the valleys of

305 diversity observed in Fig. 5c, where the selective trough is wider than the neutral trough. Even

306 though recombinations occur less frequently on the neutral tree as compared to a selected

307 tree, a recombination on the neutral tree is more likely to lead to a change of genomic

308 background from derived to ancestral allele due to the inverse sigmoid neutral trajectory of the

309 derived allele. Recombination on the neutral tree will thus more often lead to a lineage

310 escaping the sweep, resulting in more efficient recovery of diversity in the neutral case for a

311 given genomic distance from the focal locus. Furthermore, we see that the trough is deeper in

312 the neutral case (Fig. 5c), since the average coalescence time is smaller at the focal site due to

313 the smaller total length of the coalescence tree.

314

315    To determine if these differences between selective and neutral troughs hold for other fixation

316    times and population sizes, we define two quantities that characterize the shape of a trough, as

317    well as its propensity to be detected in real data: i) the trough relative depth and ii) the width of

318    the trough. The relative depth is defined as the difference between the background level of

319    diversity and the diversity at the focal locus, divided by the background diversity, and the width

320    is measured at half depth, *i.e.* halfway between the background diversity and the diversity at

321    the focal locus. On Fig. 6 we plot the relative depth of neutral and selective troughs as a

322    function of their width for different fixation times $t_m$, calculated with our analytical expressions.

323    We see that the neutral troughs are not only always narrower than the selective troughs for the

324    same value of $t_m$, but also deeper. This is due to differences in the focal tree structure between

325    the selective case and the neutral case as well as difference in the ancestral background level in

326    both cases, as explained above. For very short fixation times (corresponding to selection

327    coefficients larger than 0.1), there is almost no difference between troughs generated by

328    selective and neutral sweeps. Indeed, for such values of $t_m$, in both cases the focal coalescence

329    tree is essentially a star tree because the increase in frequency is very fast, and the ancestral

330    backgrounds of diversity, $2N_0$ and $2N_1$, are also practically equal. Note however that at small $t_m$

331    the corresponding value of the selection coefficient $s$ (see legend of Fig. 6) may be

332    unrealistically high. For realistic values of the selection coefficient/fixation time, the neutral

333    troughs tend to be quite deep but narrow, whereas selective troughs are wider and their depth

334    decreases quickly for low selection coefficients. From Fig. 6, we see that the shape of a neutral

335    trough is generally different from a selective sweep signal, but in practice those differences

336    might be hidden due to the noise inherent present in real genomic data, and it might be

337    difficult to decide whether a genomic signal is a due to a neutral sweep or a selective sweep.

338

## Discussion

340    It has repeatedly been suggested that strong depletions of diversity in the genome are not

341    necessarily due to the presence of positive selection (Johri *et al.* 2020), and can also be the

342    result of demographic effects only, such as the allele surfing phenomenon occurring at the front

343    of a range expansion (Klopfstein *et al.* 2006). In this work, we considered a model of population

344    contraction to analyze quantitatively the genomic signature of the rapid fixation of a mutation

345    during a population contraction, but it should also apply in case of range expansions or

346    recurrent founder events by considering the harmonic mean of population sizes. Taking a step

347    further from previous work that focused on the impact of range expansion on mere allele

348    frequencies, we have studied here the impact of a neutral allele fixation on neighboring

349    genomic diversity. We show that the diversity profile around a recently fixed locus crucially

350    depends on the frequency trajectories of the allele going to fixation, and we outline the fact

351    that neutrally fixing alleles have an inverse-sigmoid trajectory (Fig. 3d), as compared to the

352    standard sigmoid frequencies observed for positively selected alleles. For the same fixation

353    time, this difference translates into different genomic signatures (see figs. 5c and 6). Our results

354    demonstrate that there is a short period after a demographic contraction (or during a range

355    expansion) where observed profiles of genomic diversity would look like those usually

356    attributed to selection (Fig. 1c), and that selective sweep signals can be mimicked by neutrally

357    fixing mutations without the need to invoke complex histories of population size changes.

358    Our results allow for a systematic comparison of selective and neutral troughs of diversity, and

359    we used our results to investigate trough shapes for range of neutral and selected scenarios

360    (see Fig. 6), which in principle can be used to decide whether a given empirical trough is due to

361    selection or demography, and to infer the corresponding parameters. However, we did not

362    consider the whole spectrum of possible selection scenarios. It would be indeed interesting to

363    use our results to study cases of background selection, small selection coefficients, and a

364    variety of dominance coefficients. All these cases should have their own characteristic

365    trajectories of fixation, and hence potentially different genomic signatures. In addition, in our

366    model we do not consider mutations that fixed in the past (we always assume that the allele

367    has just reached fixation), nor do we consider mutations appearing before the population

368    contraction, i.e., with $t_m > t_c$. The average coalescence time in the former case can be expressed

369    as a function of the coalescence time at fixation using conditional probabilities, and we can

370    show that a sweep signal vanishes exponentially with the time elapsed since fixation (see

371    Appendix A4). In the latter case, we can solve the problem by considering the number of gene

372    copies at $t_c$ that descend from the original copy that appeared at $t_m$. One could extend our

373    results by considering an allele starting from an arbitrary number of copies at $t_c$, akin to soft

374    selective sweeps; however, the analytic calculations are complex, and we leave this study for

375    future research. In any case, those additional scenarios must be considered when trying to infer

376    models from the study of troughs found in empirical data. Another phenomenon that renders

377    the inference of parameters cumbersome is a possible interference between troughs. Indeed,

378    when two loci fix neutrally in the population, the genetic diversity in the region between those

379    loci will be influenced by both fixations and will differ from the diversity expected in the vicinity

380    of a single fixing locus. As in the case of interference between the fixation of selected alleles

381    (Weissman and Barton 2012), this should limit the number of independent neutral fixations.

382    The effect of trough interference is stronger for neighboring troughs, and the probability to

383    observe close troughs depends on the relative frequency of troughs along the genome, which

384    itself depends on the distribution of the $T_{MRCA}$. In Fig. 1d for example, the distribution of $T_{MRCA}$

385    has a mode centered around $4N_c$ (not shown) and correspondingly the nucleotide diversity is

386    peaked around $4N_c\,\mu$. As a result, we see many regions of the chromosome with a low diversity.

387    It is likely that those troughs interfere with each other and that they do not correspond to the

388    profile of an isolated trough. On the other hand, in Fig. 1c, the first mode of the $T_{MRCA}$

389    distribution is truncated because $t_c$ is much smaller than $4N_c$, and only $T_{MRCA}$s equal or close to

390    $t_c$ are observed (plus all the $T_{MRCA}$s corresponding to the second mode centered at $4N_0$). In this

391    case there is no interference and the (rare) troughs, such as the one in Fig. 7, are correctly

392    fitted by their theoretical expectation. Those considerations imply that, even though we know

393    the forward in time probability that an allele will fix in $t_m$ generations, it is difficult to infer the

394    parameters of a fixation scenario from a single observed neutral valley of diversity. It appears

395    therefore difficult to perform model selection from a single trough signal, i.e., to decide

396    whether a particular trough is due to selection or demographic effects, because alternative

397    demographic scenarios that we did not consider here could also lead to similar signals.

398    We performed simulations to investigate the signature of a neutral rapid fixation on the Site

399    Frequency Spectrum (SFS) (see Appendix . We chose demographic parameters such that

400    troughs are not numerous along the genome, and leave a strong footprint on genomic diversity.

15

401    Out of 10,000 simulations of 20Mb chromosomes, only 432 exhibit a (single) region of highly

402    reduced diversity (here arbitrarily set to less than 7% of the background diversity). By averaging

403    over all these valleys of diversity, we calculated the average SFS observed in a 15Kb window at

404    the center of the valley, and obtained a U-shape SFS, which is also expected around a selective

405    sweep (Huber *et al.* 2016). However, contrary to a fixation driven by selection (Suppl. Figure

406    S1), the SFS around a neutral fixation shows a slight excess of variants at intermediate

407    frequencies. This is probably due to the fact that some neutral haplotypes have spent more time

408    at intermediate frequencies before going to fixation than selected haplotypes that rapidly

409    "jump" from very low to very high frequencies (see Fig. 5a). Note also that the background

410    (genome wide) SFS away from neutral sweeps has a global excess of intermediate and high

411    frequency variants compared to a constant size population. This excess of high frequency

412    variants is typical of populations having gone through a recent population size reduction or a

413    bottleneck (Marth *et al.* 2004) due to the higher coalescence rate during the population

414    contraction. These differences in expected SFS around neutral and selected sweeps could help

415    decide whether regions of low diversity observed in empirical data are due to selection or to

416    demographic processes. However, since very few variants are usually observed in the vicinity of

417    single troughs, the empirical SFS in such a region might be too noisy to confidently identify the

418    cause of the diversity reduction. In principle, if several troughs of diversity were observed in a

419    genome, one could use the distribution of trough shapes and pooled SFS expected under a

420    given simple demographic model and a distribution of fitness effect to compare neutral and

421    selection models under a likelihood framework, but such an exploration is beyond the scope of

422    the present paper.

423    In conclusion, our results suggest that any empirical valley of diversity found in empirical data

424    can be reproduced neutrally with a population contraction using appropriate parameters. One

425    could argue that this identifiability problem disappears once the true evolutionary history is

426    correctly inferred. However, inferring the true demographic history requires precise knowledge

427    about how selection has shaped genomic diversity (Johri *et al.* 2020). In humans, for instance, it

428    has been estimated that roughly 95 % of genomic diversity is affected by some form of non-

429    neutral forces such as background selection or biased gene conversion (Pouyet *et al.* 2018)

430    potentially biasing demographic inference (Ewing and Jensen 2016). These considerations

431    indicate than genome scans in search for signals of adaptation might be more affected by past

432    demography than previously thought. We thus believe that despite current advances using

433    supervised machine learning or similar approaches (Schrider and Kern 2018), it remains

434    important to further study the effect of neutral fixations in various demographic scenarios using

435    localized genomic approaches such as the present analytical work (Johri *et al.* 2021b), as well as

436    with controlled experiments on real living organisms where both the selected locus and the

437    population history are known (Orozco-terWengel *et al.* 2012). Such work will be critical in order

438    to develop more appropriate evolutionary null models for statistical inference (Hahn 2008;

439    Johri *et al.* 2020).

## 440    Data availability

441    The authors affirm that all data necessary for confirming the conclusions of the article are

442    present within the article, figures, and tables.

## 448    Competing interest

449    None to declare

450

451 ## Bibliography

452

453 Andolfatto P., and M. Przeworski, 2000 A genome-wide departure from the standard neutral

454         model in natural populations of Drosophila. Genetics 156: 257–268.

455 Austerlitz F., B. Jung-Muller, B. Godelle, and P.-H. Gouyon, 1997 Evolution of coalescence

456         times, genetic diversity and structure during colonization. Theor. Popul. Biol. 51: 148–

457         164.

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

459 Charlesworth B., M. T. Morgan, and D. Charlesworth, 1993 The effect of deleterious mutations

460         on neutral molecular variation. Genetics 134: 1289–1303.

461 Charlesworth D., B. Charlesworth, and M. T. Morgan, 1995 The pattern of neutral molecular

462         variation under the background selection model. Genetics 141: 1619–1632.

463 Charlesworth B., 2013 Background selection 20 years on: the Wilhelmine E. Key 2012

464         invitational lecture. J. Hered. 104: 161–171.

465 Charlesworth B., 2020 How Good Are Predictions of the Effects of Selective Sweeps on Levels

466         of Neutral Diversity? Genetics 216: 1217–1238.

467 Corbett-Detig R. B., D. L. Hartl, and T. B. Sackton, 2015 Natural selection constrains neutral

468         diversity across a wide range of species. PLoS Biol. 13: e1002112.

469 Crisci J. L., Y.-P. Poh, S. Mahajan, and J. D. Jensen, 2013 The impact of equilibrium

470         assumptions on tests of selection. Front. Genet. 4: 235.

471    Edmonds C. A., A. S. Lillie, and L. Luca Cavalli-Sforza, 2004 Mutations arising in the wave

472         front of an expanding population. Proc. Natl. Acad. Sci. U. S. A. 101: 975–979.

473    Ewens W. J., 2004 *Mathematical Population Genetics: I. Theoretical Introduction*. Springer,

474         New York, NY.

475    Ewing G. B., and J. D. Jensen, 2016 The consequences of not accounting for background

476         selection in demographic inference. Mol. Ecol. 25: 135–141.

477    Excofffier L., N. Marchi, D. A. Marques, R. Matthey-Doret, A. Gouy, *et al.*, 2021 fastsimcoal2:

478         demographic inference under complex evolutionary scenarios. Bioinformatics.

479         https://doi.org/10.1093/bioinformatics/btab468

480    Excoffier L., M. Foll, and R. J. Petit, 2009 Genetic Consequences of Range Expansions. Annu.

481         Rev. Ecol. Evol. Syst. 40: 481–501.

482    Galtier N., and M. Rousselle, 2020 How Much Does Ne Vary Among Species? Genetics 216:

483         559–572.

484    Griffiths R. C., and S. Tavaré, 1994 Ancestral inference in population genetics. Stat. Sci. 9: 307–

485         319.

486    Hahn M. W., 2008 Toward a selection theory of molecular evolution. Evolution 62: 255–265.

487    Hallatschek O., and D. R. Nelson, 2008 Gene surfing in expanding populations. Theor. Popul.

488         Biol. 73: 158–170.

489    Haller B. C., and P. W. Messer, 2019 SLiM 3: Forward genetic simulations beyond the Wright-

490         Fisher model. Mol. Biol. Evol. 36: 632–637.

491    Hermisson J., and P. S. Pennings, 2017 Soft sweeps and beyond: understanding the patterns and

492        probabilities of selection footprints under rapid adaptation. Methods Ecol. Evol. 8: 700–

493        716.

494    Huber C. D., M. DeGiorgio, I. Hellmann, and R. Nielsen, 2016 Detecting recent selective sweeps

495        while controlling for mutation rate and background selection. Mol. Ecol. 25: 142–156.

496    Jensen J. D., Y. Kim, V. B. DuMont, C. F. Aquadro, and C. D. Bustamante, 2005 Distinguishing

497        between selective sweeps and demography using DNA polymorphism data. Genetics

498        170: 1401–1410.

499    Jensen J. D., B. A. Payseur, W. Stephan, C. F. Aquadro, M. Lynch, *et al.*, 2019 The importance

500        of the Neutral Theory in 1968 and 50 years on: A response to Kern and Hahn 2018.

501        Evolution 73: 111–114.

502    Johri P., B. Charlesworth, and J. D. Jensen, 2020 Toward an Evolutionarily Appropriate Null

503        Model: Jointly Inferring Demography and Purifying Selection. Genetics 215: 173–192.

504    Johri P., B. Charlesworth, E. K. Howell, M. Lynch, and J. D. Jensen, 2021a Revisiting the

505        Notion of Deleterious Sweeps. Genetics. https://doi.org/10.1093/genetics/iyab094

506    Johri P., K. Riall, H. Becher, L. Excoffier, B. Charlesworth, *et al.*, 2021b The Impact of

507        Purifying and Background Selection on the Inference of Population History: Problems

508        and Prospects. Mol. Biol. Evol. 38: 2986–3003.

509    Kaiser V. B., and B. Charlesworth, 2009 The effects of deleterious mutations on evolution in

510        non-recombining genomes. Trends Genet. 25: 9–12.

511    Kingman J. F. C., 1982a The coalescent. Stochastic Process. Appl. 13: 235–248.

512    Kingman J. F. C., 1982b On the genealogy of large populations. J. Appl. Probab. 19A: 27–43.

513    Klopfstein S., M. Currat, and L. Excoffier, 2006 The fate of mutations surfing on the wave of a

514         range expansion. Mol. Biol. Evol. 23: 482–490.

515    Marth G. T., E. Czabarka, J. Murvai, and S. T. Sherry, 2004 The Allele Frequency Spectrum in

516         Genome-Wide Human Variation Data Reveals Signals of Differential Demographic

517         History in Three Large World Populations. Genetics 166: 351–372.

518    Maruyama T., and M. Kimura, 1975 Moments for sum of an arbitrary function of gene frequency

519         along a stochastic path of gene frequency change. Proc. Natl. Acad. Sci. U. S. A. 72:

520         1602–1604.

521    Mathew L. A., and J. D. Jensen, 2015 Evaluating the ability of the pairwise joint site frequency

522         spectrum to co-estimate selection and demography. Front. Genet. 6: 268.

523    Maynard Smith J., and J. Haigh, 1974 The hitch-hiking effect of a favorable gene. Genet. Res.

524         23: 23–35.

525    Nicolaisen L. E., and M. M. Desai, 2013 Distortions in genealogies due to purifying selection

526         and recombination. Genetics 195: 221–230.

527    O'Fallon B. D., J. Seger, and F. R. Adler, 2010 A continuous-state coalescent and the impact of

528         weak selection on the structure of gene genealogies. Mol. Biol. Evol. 27: 1162–1172.

529   Orozco-terWengel P., M. Kapun, V. Nolte, R. Kofler, T. Flatt, *et al.*, 2012 Adaptation of

530       Drosophila to a novel laboratory environment reveals temporally heterogeneous

531       trajectories of selected alleles. Mol. Ecol. 21: 4931–4941.

532   Peischl S., I. Dupanloup, M. Kirkpatrick, and L. Excoffier, 2013 On the accumulation of

533       deleterious mutations during range expansions. Mol. Ecol. 22: 5972–5982.

534   Peischl S., and L. Excoffier, 2015 Expansion load: recessive mutations and the role of standing

535       genetic variation. Mol. Ecol. 24: 2084–2094.

536   Pouyet F., S. Aeschbacher, A. Thiéry, and L. Excoffier, 2018 Background selection and biased

537       gene conversion affect more than 95% of the human genome and bias demographic

538       inferences. Elife 7: e36317.

539   Pouyet F., and K. J. Gilbert, 2019 Towards an improved understanding of molecular evolution:

540       the relative roles of selection, drift, and everything in between. arXiv [q-bio.PE].

541   Rogers R. L., T. Bedford, A. M. Lyons, and D. L. Hartl, 2010 Adaptive impact of the chimeric

542       gene Quetzalcoatl in Drosophila melanogaster. Proc. Natl. Acad. Sci. U. S. A. 107:

543       10943–10948.

544   Rousselle M., M. Mollion, B. Nabholz, T. Bataillon, and N. Galtier, 2018 Overestimation of the

545       adaptive substitution rate in fluctuating populations. Biol. Lett. 14.

546       https://doi.org/10.1098/rsbl.2018.0055

547   Schrider D. R., and A. D. Kern, 2018 Supervised Machine Learning for Population Genetics: A

548       New Paradigm. Trends Genet. 34: 301–312.

549    Slatkin M., 1996 Gene genealogies within mutant allelic classes. Genetics 143: 579–587.

550    Sousa V., S. Peischl, and L. Excoffier, 2014 Impact of range expansions on current human

551        genomic diversity. Curr. Opin. Genet. Dev. 29: 22–30.

552    Tajima F., 1990 Relationship between DNA polymorphism and fixation time. Genetics 125:

553        447–454.

554    Tavaré S., 1984 Line-of-descent and genealogical processes, and their applications in population

555        genetics models. Theor. Popul. Biol. 26: 119–164.

556    Teshima K. M., G. Coop, and M. Przeworski, 2006 How reliable are empirical genomic scans for

557        selective sweeps? Genome Res. 16: 702–712.

558    Thornton K. R., and J. D. Jensen, 2007 Controlling the false-positive rate in multilocus genome

559        scans for selection. Genetics 175: 737–750.

560    Wares J. P., 2009 Evolutionary dynamics of transferrin in Notropis. J. Fish Biol. 74: 1056–1069.

561    Weissman D. B., and N. H. Barton, 2012 Limits to the rate of adaptive substitution in sexual

562        populations. PLoS Genet. 8: e1002740.

563    Zhao L., M. Lascoux, A. D. J. Overall, and D. Waxman, 2013 The characteristic trajectory of a

564        fixing allele: a consequence of fictitious selection that arises from conditioning. Genetics

565        195: 993–1006.

566

567    *Figure 1*. Nucleotide diversity of a population experiencing a contraction, as a function of the
568    time $t_c$ elapsed since the contraction, measured in units of $2N_c$. (a) distribution of nucleotide

569    *diversity as a function of time, nucleotide diversity along the chromosome at $t_c$ = 0 (panel b), at*
570    *$t_c$ = 0.25 (panel c) and at $t_c$ = 0.75 (panel d). Population size before contraction $N_0$ = 2.37×10⁶*
571    *and after contraction $N_c$ = 4,400. Mutation rate μ = 5.42×10⁻¹⁰ per site per generation.*
572    *Recombination rate r = 3.5×10⁻⁸ per site per generation. Chromosome size L = 20 Mb. Window*
573    *size 10 Kb sliding at 1 Kb intervals. Sample size: 30 haplotypes. These parameters are taken from*
574    Rogers *et al.* (2010). *Simulations were performed with fastsimcoal2 (Excofffier et al. 2021).*

575

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

585

586    ***Figure 3***. *Average frequency (a-d) and coalescence time distribution (e-h) of an allele fixing in a*
587    *diploid population of constant size $N_c$ = 20 in exactly $t_m$ generations, starting as a single copy*
588    *(i.e. $p_0 = (2N_c)^{-1}$). The red dots are the results of Wright-Fisher simulations, and the black and*
589    *white dashed lines are calculated with eqs. (4b) (first and second columns) (4a) (third and fourth*
590    *columns) and (3b). In panes (a-d) we show the variability of the fixation process by overlapping*
591    *1780 fixing trajectories. The (numerically estimated) probability, for a mutant that appears at*
592    *the onset of the contraction, to fix in less than $t_m$ generations is 0.006, 0.16, 0.64 and 0.86 for*
593    *$t_m$ = 20, 40, 80 and 120 respectively (for this particular value of $N_c$).*

594

595    ***Figure 4***. *Average coalescence time at a linked locus, as a function of the recombination*
596    *distance from the focal locus where a mutant fixed in exactly $t_m$ generations, starting from a*
597    *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*
598    *calculated with two-locus WF simulations, and compared to eq. (5) with either a numerical*
599    *estimation (solid lines) or a theoretical estimation (dashed lines) of $\overline{x}_t$ and $P\left(T^{(f)}\right)$. $N_c$ = 20. $N_0$ =*
600    *1500. The population experienced a contraction $t_c$ = $t_m$ generations ago.*

601

602    ***Figure 5***. *Comparison between troughs of diversity resulting from a selective sweep (black) and*
603    *a neutral sweep (red), for the same fixation time $t_m = 120$ (corresponding to s ≈ 0.1 in the*
604    *selective case). Frequency of the fixing allele as a function of time (a), coalescence time*
605    *distribution (b) and diversity around the fixing site along the genome using eq. (5) (c). $N_1$ = 1500,*
606    *$N_c$ = 20 and $N_0$ = 2.97×10⁴.*

607

608 *Figure 6. Relative depth as a function of the width of the diversity troughs, for different values*
609 *of $t_m$ and $N_c$ in the neutral case and for selective scenarios with identical fixation times. $t_m$ goes*
610 *from 1 to 333 by increments of 1, the corresponding values of the selection coefficient s are*
611 *indicated on the left of the legend bar (for all of them we have $N_1 s >> 1$). $N_1 = 1500$. $N_0$ is given*
612 *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*
613 *100 are due to the use of two different approximations for the frequency of the mutant, eqs.*
614 *(4a) and (4b) and are located at $t_m = 2N_c$.*

615

## Appendix

616

617 **A1. Coalescence distribution after a contraction**

618 We want to determine the coalescence time of two lineages in a population that experienced a

619 contraction $t_m$ generations ago, from a diploid size $N_0$ to $N_c$. As we go backward in time, the

620 coalescence rate switches from $(2N_c)^{-1}$ to $(2N_0)^{-1}$ at $T = t_c$. The probability distribution might

621 still be approximated by a piecewise exponential density:

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

622 The corresponding expectation for this distribution is

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

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

624 In this section time is counted forward from the mutation, which appears after the contraction,

625 so that during the fixation the diploid population size is constant and equal to $N_c$. We condition

626 on the fixation time $t_m$ of the mutant. We define the trajectory of a mutant as the list of

627  frequencies at all generations: $\{x_t\} = \left(x_0, x_1, \ldots, x_{t_m-1}, x_{t_m}\right)$. We assume that the mutant fixes

628  in exactly $t_m$ generations, starting from a frequency $p_0$, i.e. $x_0 = p_0$, $0 < x_{t_m-1} < 1$ and

629  $x_{t_m} = 1$. The probability that the mutant follows a given trajectory might be expressed as the

630  product of the transition probabilities

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

631  For an unconditional Wright Fisher model, $P(i, t \to j, t+1)$ is the probability to have j copies of

632  the new allele at $t+1$ given that there were i copies at t. We note $P_t(i \to j)$ for brevity. If we

633  only consider trajectories fixing in exactly $t_m$ generations and starting from a number $2N_c\,p_0$ of

634  copies at $t = 0$, then the transition probabilities are not equal to the transitions of the

635  unconditional Wright-Fisher model. However, thanks to Bayes theorem, we can write

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

636  From the first to the second line, we use the Markov property. The three terms involved in the

637  right-hand side of this equation can be approximated thanks to diffusion theory. In this

638  framework, the probability for an allele to fix in $t_m$ generations, given that there were i copies

639  at time t is approximately (Ewens 2004)

$$P(\text{fix in } t_m \mid i_t) = \frac{3}{2N_c}\left(1 - \frac{i}{2N_c}\right)\frac{i}{2N_c}\; e^{-(t_m - t)/2N_c} \qquad (S2)$$

640  The term $P_t(i \to j)$ is the unconditional binomial transition probability of the Wright Fisher

641  model (which does not depend on t). In principle, eq. (S1) can be used to compute the exact

642  distribution of coalescence times at the focal locus, using eq. (3a). However, the huge number

643  of possible trajectories fixing in $t_m$ generations ($(2N_c - 1)^{t_m-1}$) makes the average over

644  trajectories impossible to evaluate numerically. For this reason, we use the approximation in

645  eq. (3b).

646  We consider here the probability that the allele has frequency x at time t, given that it started

647  at frequency $p_0$ at t = 0. Again if we only consider trajectories that fix in exactly $t_m$

648  generations, this probability is not equal to the neutral diffusive result. However, similarly to

649  the previous section, we can use Bayes theorem:

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

650  From diffusion theory (Ewens 2004), we also have

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

651  which is a second order expansion of an infinite series involving vanishing exponential terms

652  ($e^{-k(k+1)t/4N_c}$ for all k $\geq$ 1). This expansion is thus valid in the limit of large times t $\gg$ $2N_c$. We

653  deduce that the probability that an allele fixing in $t_m$ generations has frequency x at time t is

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

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

655  This expression is valid for $t_m \gg t \gg 2N_c$, and does not allow one to estimate the frequency

656  close to fixation. If we evaluate this expression for a given value of *t*, we must assume that $t_m$ is

657  much larger than *t* (otherwise (S2) is not accurate). It implies that we cannot evaluate the

658  frequency close to fixation, because wherever we "look", the fixation is always much later in

659  time. Consequently, we see that $E[x_t]$ tends to 1/2 when *t* is very large, which is the only

660  possible value for an average frequency infinitely far away from both fixation (at *t* = $t_m$) and loss

661  (at *t* = 0) . However, we know that the frequency should be symmetric, *i.e.* the allele should on

662  average approach fixation in the same way it escapes loss, because the neutral fixation of a

663  derived allele is the same as the loss of the ancestral allele. We thus write

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

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

665  observations)

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

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

667    We now return to the scenario of Fig. 2, with a backward in time approach. Using Bayes

668    theorem, we express the coalescence time of two haplotypes at the linked locus $T^{(l)}$,

669    conditioning on the coalescence time at the focal locus $T^{(f)}$

$$\mathrm{P}\big(T^{(l)}\big) = \int_0^{t_m} \mathrm{P}\big(T^{(l)}\,\big|\,T^{(f)}\big)\, P\big(T^{(f)}\big)\mathrm{d}T^{(f)} = \mathrm{E}\big[\mathrm{P}\big(T^{(l)}\,\big|\,T^{(f)}\big)\big]$$

670    We assume that the linked locus is close to the focal locus on the chromosome, more precisely

671    that the recombination rate r is very small $r \ll 1$, so that we consider at most one

672    recombination, occurring on one of the two focal lineages. We distinguish cases where there is

673    no recombination between $t = 0$ and $t = T^{(f)}$, cases where the allele at the linked locus

674    recombines (somewhere between $t = 0$ and $t = T^{(f)}$) onto a haplotype carrying the ancestral

675    allele at the focal locus, and cases where the allele at the linked locus recombines onto a

676    haplotype carrying the derived allele at the focal locus. We call the second and third case

677    heterozygous and homozygous recombination, respectively, referring to the zygosity at the

678    focal locus of the recombining pair of haplotypes (note that are three haplotypes, the two first

679    ones have a coalescence time $T^{(f)}$, and the third one recombines with one of these two). If there

680    is no recombination, then the coalescence time is the same for both loci, $T^{(l)} = T^{(f)}$. To treat the

681    case with a homozygous recombination, it is convenient to name the haplotypes: *i* and *j*

682    coalesce at $T_{ij}^{(f)} = T^{(f)}$ at the focal locus, and *k* is a third haplotype, onto which the linked allele

683    recombines (coming from *i*). The linked allele carried by *j* stays on the same haplotype (no more

684    than one recombination), and after recombining onto *k*, the linked allele initially carried by *i*

685    also stays on *k* (again, at most one recombination). This implies that those two linked alleles

686    coalesce at $T_{ij}^{(l)} = T_{jk}^{(f)}$. This time is in general different than $T_{ij}^{(f)}$, however on average $T_{jk}^{(f)}$ tand

687    $T_{ij}^{(f)}$ are equal (averaging over all possible coalescence trees at the focal locus). This implies that

688    we can treat the case with homozygous recombination as if there was no recombination. If

689    there is a heterozygous recombination between *i* and *k*, at some generation between $t = 0$ and $t$

28

690     $= T^{(f)}$, then the linked alleles still have not coalesced at $t = t_m$ because after the recombination

691     one of them is linked to a derived focal allele and the other one to an ancestral focal allele (and

692     they stay linked because there is at most one recombination). In that case, $T_{ij}^{(l)}$ is equal to $t_m$

693     plus a random time given by (on average) $T_m$, and is independent of $T_{ij}^{(f)}$. Using again Bayes

694     theorem and the previous results to write

$$
\begin{aligned}
\mathrm{P}\left(T^{(l)}\,\middle|\,T^{(f)}\right) = {}& \mathrm{P}\left(T^{(l)}\,\middle|\,T^{(f)}, \quad one\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right)P\left(one\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right) \\
& + \mathrm{P}\left(T^{(l)}\,\middle|\,T^{(f)}, \quad no\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right)P\left(no\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right) \\
= {}& f_m\left(T^{(l)} - t_m\right)\left[1 - P\left(no\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right)\right] \\
& + \delta\left(T^{(l)} - T^{(f)}\right)P\left(no\ het.\,\mathrm{rec.\ in}\left[0,\,T^{(f)}\right]\right)
\end{aligned}
$$

695     Where $\delta(\cdot)$ is the Dirac delta function, and $f_m$ is the unconditional coalescence distribution of a

696     pair of lineages sampled at $t = t_m$, i.e. it is equal to the function $f_0$ introduced above but

697     replacing $t_c$ by $t_c - t_m$ (note also that $f_m(t) = 0$ if $t < 0$). We then have to evaluate the

698     probability that there is no heterozygous recombination. At generation $t$ (counted backward)

699     the probability that a linked allele recombines onto a haplotype carrying the ancestral allele at

700     the focal locus is $r(1 - x_t)$, where $x_t$ is the frequency of the derived allele at the focal locus,

701     we deduce that the probability that there is no heterozygous recombination on either lineage is

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

702     This probability depends explicitly on the allele trajectory, which means that rigorously, all the

703     calculations should be conditioned on a given trajectory, and then averaged over all

704     trajectories. To allow for mathematical tractability, and to avoid heavy expressions, we consider

705     that as a good approximation $x_t = \overline{x}_t$. Finally we obtain

$$P\big(T^{(l)}\big) = E\left[\delta\big(T^{(l)} - T^{(f)}\big) \exp\left(-2r \sum_{t=1}^{T^{(f)}} (1 - x_t)\right)\right]$$

$$+ f_m\big(T^{(l)} - t_m\big) E\left[1 - \exp\left(-2r \sum_{t=1}^{T^{(f)}} (1 - x_t)\right)\right]$$

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

707

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

709    **generations ago**

710    Thanks to Bayes theorem we can write

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

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

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

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

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

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

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

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

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

718

719    **A5. Site frequency spectrum around a neutral trough compared to a selective trough**

720

$T^{(l)}$  $T^{(f)}$

haplotype i

haplotype j

$N_0$

N(t)

$N_c$

focal locus (f)  linked locus (l)

time

$t_c$  $t_m$  $t = 0$