

Selective sweeps under dominance and self-fertilisation

Matthew Hartfield^{1,2,*}, Thomas Bataillon²

1 Department of Ecology and Evolutionary Biology, University of Toronto,
Ontario, Canada.

2 Bioinformatics Research Centre, Aarhus University, 8000 Aarhus, Denmark.

* matthew.hartfield@birc.au.dk

Running Head: Sweeps under dominance and selfing

Key words: Adaptation; Dominance; Self-fertilisation; Selective Sweeps; *SLC24A5*

1 Abstract

2 A major research goal in evolutionary genetics is to uncover loci experiencing
3 adaptation from genomic sequence data. One approach relies on finding ‘selective
4 sweep’ patterns, where segregating adaptive alleles reduce diversity at linked neu-
5 tral loci. Recent years have seen an expansion in modelling cases of ‘soft’ sweeps,
6 where the common ancestor of derived variants predates the onset of selection. Yet
7 existing theory assumes that populations are entirely outcrossing, and dominance
8 does not affect sweeps. Here, we develop a model of selective sweeps that considers
9 arbitrary dominance and non-random mating via self-fertilisation. We investigate
10 how these factors, as well as the starting frequency of the derived allele, affect
11 average pairwise diversity, the number of segregating sites, and the site frequency
12 spectrum. With increased self-fertilisation, signatures of both hard and soft sweeps
13 are maintained over a longer map distance, due to a reduced effective recombina-
14 tion rate and faster fixation times of adaptive variants. We also demonstrate that
15 sweeps from standing variation can produce diversity patterns equivalent to hard
16 sweeps. Dominance can affect sweep patterns in outcrossing populations arising
17 from either a single novel mutation, or from recurrent mutation. It has little effect
18 where there is either increased selfing or the derived variant arises from stand-
19 ing variation, since dominance only weakly affects the underlying adaptive allele
20 trajectory. Different dominance values also alters the distribution of singletons
21 (derived alleles present in one sample). We apply models to a sweep signature at
22 the *SLC24A5* gene in European humans, demonstrating that it is most consistent
23 with an additive hard sweep. These analyses highlight similarities between certain
24 hard and soft sweep cases, and suggest ways of how to best differentiate between

25 related scenarios. In addition, self-fertilising species can provide clearer signals
26 of soft sweeps than outcrossers, as they are spread out over longer regions of the
27 genome.

28 **Author Summary**

29 Populations adapt by fixing beneficial mutations. As a mutation spreads, it drags
30 linked neutral variation to fixation, reducing diversity around adaptive genes. This
31 footprint is known as a ‘selective sweep’. Adaptive variants can appear either from
32 a new mutation onto a single genotype; from recurrent mutation onto different
33 genotypes; or from existing genetic variation. Each of these sources leaves subtly
34 different selective sweep patterns in genetic data, which have been explored under
35 simple biological cases. We present a general model of selective sweeps that in-
36 cludes self-fertilisation (where individuals produce both male and female gametes
37 to fertilise one another), and dominance (where fitness differences exist between
38 one and two gene copies within an individual). Soft sweep patterns are spread out
39 over longer genetic regions in self-fertilising individuals, while dominance mainly
40 affects sweeps in outcrossers from either a single or recurrent mutation. Applying
41 models to a sweep signal associated with human skin pigmentation shows that
42 this mutation was likely introduced into Eurasia from Africa in very few numbers.
43 These models demonstrate to what extent soft sweeps can be detected in genome
44 data, and how self-fertilising organisms can be good study systems for determining
45 the extent of different adaptive modes.

46 Introduction

47 Inferring adaptation from nucleotide sequence data is a major research goal in evo-
48 lutionary genetics. The earliest models focussed on the scenario where a beneficial
49 mutation appeared in the population in a single copy before rapidly spreading to
50 fixation. Linked neutral mutation would then ‘hitchhike’ to fixation with the adap-
51 tive variant, reducing diversity around the selected locus [1, 2]. Hitchhiking also
52 causes a rapid increase in linkage disequilibrium at flanking regions to the selected
53 site, although it is minimal when measured either side of the beneficial muta-
54 tion [3–5]. These theoretical expectations have spurred the creation of summary
55 statistics for detecting sweeps, based on finding regions of the genome exhibiting
56 extended runs of homozygosity [6–10].

57 Classic hitchhiking models consider ‘hard’ sweeps, where the common ancestor
58 of adaptive alleles occurs after its appearance [11]. Yet the last fifteen years have
59 seen a focus on quantifying ‘soft’ sweeps, where the most recent common ancestor
60 of the beneficial allele arose before the variant became selected for (reviewed in
61 [11–13]). Soft sweeps can originate from beneficial mutations being introduced
62 by recurrent mutation [14, 15], or from existing standing variation that was either
63 neutral or deleterious [16–22]. A key property of soft sweeps is that the beneficial
64 variant is present on multiple genetic backgrounds as it sweeps to fixation, so
65 different haplotypes are present around the derived allele. This property is often
66 used to detect soft sweeps in genetic data [23–28]. Soft sweeps have been inferred
67 in several organisms, including *Drosophila* [25, 26], humans [23, 29], maize [30]
68 and the malaria pathogen *Plasmodium falciparum* [31], although determining how
69 extensive soft sweeps are in nature remains a contentious issue [32].

70 Up to now, almost all models of selective sweeps have made the same simpli-
71 fying assumptions. In particular, there have been few analyses considering how
72 dominance affects sweep signatures. In a simulation study, Teshima and Prze-
73 worski [33] determined how recessive mutations spend a long period of time at low
74 frequencies, increasing the amount of recombination that acts on derived haplo-
75 types, weakening signatures of ‘hard’ sweeps. Fully recessive mutations may need
76 a long time to reach a high enough frequency so that they can be picked up by
77 genome scans for adaptive loci [34]. Ewing *et al.* [35] have carried out a general
78 mathematical analysis of dominance on ‘hard’ sweeps on genetic diversity. Yet the
79 impact that dominance has on ‘soft’ sweeps has yet to be explored in depth.

80 In addition, existing models have overwhelmingly assumed that populations are
81 sexual, with individuals haplotypes freely mixing between individuals. Different
82 reproductive modes alters how alleles are inherited over subsequent generations
83 and spread over time, therefore altering the hitchhiking effect. In particular, there
84 is a renewed interest in studying the mechanisms of adaptation in self-fertilising
85 species [36]. Self-fertilisation, where male and female gametes produced from the
86 same individual can fertilise each other, is prevalent amongst angiosperms [37],
87 some animals [38] and fungi [39]. Different levels of self-fertilisation is known
88 to affect overall adaptation rates. Dominant mutations are likelier to fix than
89 recessive ones in outcrossers, as they have a higher initial selection advantage
90 [40]. Yet recessive alleles can fix more easily in selfers than in outcrossers as they
91 rapidly create homozygote mutations [41, 42]. Hence the effects of dominance
92 and self-fertilisation become strongly intertwined, so it is important to consider
93 both together. Furthermore, a decrease in effective recombination rates in selfers
94 [43] can amplify the effects of linked selection, making it likelier that deleterious

95 mutations hitchhike to fixation with adaptive alleles [44], or nearby beneficial
96 alleles are lost if one is already spreading through the population [45].

97 Self-fertilisation is also known to affect the degree to which adaptation proceeds
98 from *de novo* mutation, or from standing variation. In a constant-sized population,
99 fixation of beneficial mutations from standing variation (either neutral or deleteri-
100 ous) is generally less likely in selfers as lower levels of diversity are maintained [46].
101 Yet if adaptation from standing variation does occur, then the beneficial variant
102 fixes more quickly in selfers than outcrossers, hence signatures of soft sweeps could
103 become more marked [42, 46].

104 Furthermore, adaptation from standing variation becomes likelier in selfers
105 under ‘evolutionary rescue’ scenarios, where swift adaptation needed to prevent
106 population extinction. This is because the population size is greatly reduced, so
107 the waiting time for the appearance of *de novo* rescue mutations becomes ex-
108 cessively long. Hence only adaptive mutations present in standing variation can
109 contribute to preventing population extinction [46]. High selfing rates can further
110 aid this process by creating beneficial homozygotes more rapidly than in outcross-
111 ing populations [47]. Therefore there is potential for soft sweeps to act in selfing
112 organisms.

113 However, little data currently exists on the extent of soft sweeps in self-fertilisers.
114 Many selfing organisms exhibit sweep-like regions, including *Arabidopsis thaliana*
115 [48–50]; *Caenorhabditis elegans* [51]; *Medicago truncatula* [52]; and *Microbotryum*
116 fungi [53]. Detailed analyses of these regions has been hampered by a lack of theory
117 on how hard and soft sweep signatures should manifest themselves under different
118 levels of self-fertilisation and dominance. Previous studies have only focussed on
119 special cases; Hedrick [54] analysed linkage disequilibrium caused by a hard sweep

120 under self-fertilisation, while Schoen *et al.* [55] modelled sweep patterns caused
121 by modifiers that altered the mating system in different ways. A knowledge of
122 expected diversity patterns following different types of sweeps can also be used to
123 create more realistic statistical models for finding and quantifying novel adaptive
124 candidate loci, while accounting for the mating system.

125 We present here a general model of selective sweeps. We determine the ge-
126 netic diversity present following a sweep from either a *de novo* mutation, or from
127 standing variation. The model assumes an arbitrary level of dominance and self-
128 fertilisation. We first present general results for the probability of how genetic
129 samples, carrying a recently-fixed beneficial mutation, are affected by recombi-
130 nation, dominance and selfing. We next determine how key summary statistics
131 (pairwise diversity; number of segregating sites; and the site frequency spectrum)
132 are affected by this general sweep model from standing variation. These results
133 are compared to an alternative soft-sweep case where adaptive alleles arise via
134 recurrent mutation. We also include a simulation study of how the distribution
135 of singletons are affected under different sweep scenarios, complementing a re-
136 cent study that used singleton densities to detect recent human adaptation [56].
137 We end by applying models to determine the history of a selective sweep at the
138 *SLC24A5* gene in humans, to evaluate the evolutionary history of this adaptation,
139 and determine if evidence exists for either non-additive dominance or a soft sweep
140 signature.

141 Results

142 Model Outline

143 We consider a diploid population of size N (carrying $2N$ haplotypes in total).
144 Individuals reproduce by self-fertilisation with probability σ , and outcross with
145 probability $1 - \sigma$. The level of self-fertilisation can also be captured by the in-
146 breeding coefficient $F = \sigma/(2 - \sigma)$ [57, 58]. There are two biallelic loci A , B with
147 a recombination rate r between them. Locus A represents a region where neutral
148 polymorphism accumulates under an infinite-sites model. Locus B determines fit-
149 ness differences, carrying an allele that initially segregates at low frequency for a
150 sizeable period of time. We are agnostic as to whether this allele is neutral or sub-
151 ject to weak selection, but note that an allele subject to strong purifying selection
152 would have only recently appeared in the population, which we do not consider.
153 Once the allele reaches a frequency f_0 it becomes advantageous, with selective
154 advantage $1 + hs$ in heterozygote form and $1 + s$ as a homozygote, with $0 \leq h \leq 1$
155 and $s > 0$. We further assume that selection is strong (i.e., $N_e hs \gg 1$) so that the
156 sweep trajectory can be modelled deterministically. Table 1 lists notation used in
157 the model analysis.

158 Our overall goal is to determine how the emergence of an adaptive allele from
159 standing variation at locus B affects genealogies underlying polymorphism at locus
160 A . We model the genetic histories at A while considering the genetic background
161 of neutral alleles (i.e., whether they are linked to the selected derived allele or
162 ancestral neutral allele at locus B). A schematic of the process is shown in Fig 1.
163 We follow the approach of Berg and Coop [21] and, looking backwards in time,
164 break down the allele history into two phases. The first phase (the ‘sweep phase’)

Symbol	Usage
N	Population size (with $2N$ haplotypes)
σ	Proportion of matings that are self-fertilising
F	Wright's inbreeding coefficient, $\sigma/(2 - \sigma)$ [57, 58]
N_e	Effective population size, equal to $N/(1 + F)$ with selfing [59]
A, B	Loci carrying neutral, selected alleles
r	Recombination rate between loci A, B
r_{eff}	'Effective' recombination rate, approximately equal to $r(1 - F)$ with selfing [43]
R	$2Nr$, the population-level recombination rate
f_0	Frequency at which the derived allele at B becomes advantageous
$f_{0,A}$	'Accelerated' effective starting frequency of B appearing as a single copy, conditional on fixation
s	Selective advantage of derived allele at B
h	Dominance coefficient of derived allele at B
t	Number of generations in the past from the present day
τ_{f_0}	Time in the past when derived locus became beneficial
$p(t)$	Frequency of beneficial allele at time t
P_{NR}	Probability that neutral marker does not recombine onto ancestral background during sweep phase
$P_{NR}(i n)$	Probability that i of n neutral markers do not recombine during sweep phase
H_l, H_h	'Effective' dominance coefficient for allele at low, high frequency
P_{coal}	Probability that two samples coalesce in the standing phase
$P_{coal,M}$	Probability that two samples coalesce instead of arising by different mutations
π	Pairwise diversity at site (π_0 is expected value without selection)
π_{SV}	Pairwise diversity following sweep from standing variation
π_M	Pairwise diversity following sweep from recurrent mutation
\tilde{s}	'Effective' selection coefficient to map hard sweep onto standing variation cases
$P_{ESF}(k i)$	Ewens' Sampling Formula for the probability of k ancestral backgrounds formed from i non-recombined lineages
$\mathbb{E}(T_{tot})$	Expected time covered by entire genealogy
$\mathbb{E}(S)$	Expected number of segregating sites
μ	Probability of neutral mutation occurring per site per generation
μ_b	Probability of beneficial mutation occurring at target locus per generation
$\theta = 4N_e\mu$	Population level neutral mutation rate
$\Theta_b = 2N_e\mu_b$	Population level beneficial mutation rate

Table 1. Glossary of Notation.

165 considers the derived allele at B being selectively favoured and spreading through
166 the population. The length of this phase is assumed to be sufficiently short ($t \sim$
167 $1/s$) so that no samples coalesce during this time, but they can recombine onto
168 the ancestral background. The second phase (the ‘standing phase’) assumes that
169 the derived allele is present at a fixed frequency f_0 . Here, the two samples can
170 either coalesce, or one of them recombines onto the ancestral background. Berg
171 and Coop [21] showed that this assumption allows traditional coalescent results to
172 be used to infer genetic properties of the sweep, after appropriate rescaling of the
173 coalescent rate by f_0 .

174 For tightly linked loci ($r \rightarrow 0$), the relatively rapid fixation time of the derived
175 variant makes it unlikely for unique polymorphisms to arise on different haplo-
176 types, reducing neutral diversity. Further from the target locus, recombination
177 can transfer allele copies at A away from the selected background to the ancestral
178 background, so diversity reaches neutral levels.

179 Self-fertilisation creates two key differences compared to traditional outcrossing
180 models. First, the effective population size and recombination rate are scaled by
181 factors $1/(1+F)$ and $1-F$ respectively [43,58]. Second, the trajectory of adaptive
182 alleles, which determines expected diversity patterns following adaptation, depends
183 on the levels of self-fertilisation (σ) and dominance (h). A goal of this analyses will
184 be to determine how these processes interact to affect neutral variation following
185 a sweep, and therefore the ability to detect different types of recent adaptation.

186 Throughout, analytical solutions are compared to results obtained from Wright-
187 Fisher forward-in-time stochastic simulations. The simulation procedure itself is
188 described in the ‘Methods’ section.

189 **Probability of no recombination during sweep phase**

190 Looking back in time following a sweep, sites linked to the beneficial allele can re-
191 combine onto the ancestral genetic background, so they exhibit the same diversity
192 as putatively neutral regions. Let $p(t)$ be the frequency of the adaptive mutation
193 at time t , defined as the number of generations prior to the present day. Further
194 define $p(0) = 1$ (i.e., the allele is fixed at the present day), and τ_{f_0} the time in the
195 past when the derived variant became beneficial (i.e., $p(\tau_{f_0}) = f_0$). If the neutral
196 locus lies at a recombination distance r from the derived variant, then the proba-
197 bility that it will not recombine onto a neutral background is $1 - r(1 - p(t))$ [21].
198 We also define $r = r_{eff} = r(1 - F)$, which is the ‘effective’ recombination rate af-
199 ter accounting for the increased homozygosity created due to self-fertilisation [43].
200 More exact r_{eff} terms exist [60,61], but they are approximately equal to $r(1 - F)$
201 over short map distances. Using these more exact terms do not improve the accu-
202 racy of the analytical model relative to simulations for the parameters used (data
203 not shown).

204 Over τ_{f_0} generations, the total probability that a single lineage does not re-
205 combine onto a neutral background, P_{NR} , equals:

$$\begin{aligned} P_{NR} &= \prod_{t=0}^{\tau_{f_0}} (1 - r_{eff}(1 - p(t))) \\ &\approx \exp\left(-r_{eff} \int_{t=0}^{\tau_{f_0}} (1 - p(t)) dt\right) && \text{since } r_{eff} \ll 1 \\ &\approx \exp\left(-r_{eff} \int_{p=1}^{f_0} \frac{(1 - p(t))}{dp/dt} dp\right) && \text{integrating over } p \end{aligned} \tag{1}$$

206 We can calculate P_{NR} for general levels of self-fertilisation if the selection co-
207 efficient is not too weak (i.e., $1/N_e \ll s \ll 1$). Here the rate of change of the allele
208 frequency is given by [42]:

$$\frac{dp}{dt} = -sp(1-p)(F+h-Fh+(1-F)(1-2h)p) \quad (2)$$

209 Note the negative factor in Eq 2 since we are looking back in time. By substituting
210 Eq 2 into Eq 1, we obtain the following analytical solution for P_{NR} :

$$\begin{aligned} P_{NR} &= \exp\left(-\frac{r_{eff}}{H_l s} \log\left(1 + \frac{H_l}{H_h} \left(\frac{1}{f_0} - 1\right)\right)\right) \\ &= \left(1 + \frac{H_l}{H_h} \left(\frac{1}{f_0} - 1\right)\right)^{-r_{eff}/(H_l s)} \end{aligned} \quad (3)$$

211 Here, $H_l = F+h-Fh$, $H_h = 1-h+Fh$ are the ‘effective’ dominance coefficients
212 when the beneficial variant is at a low or high frequency [42]. We can understand
213 Eq 3 as follows. The beneficial mutation takes $(1/H_l s) \log(1 + (H_l/H_h)(1/f_0 - 1))$
214 generations to go to fixation from initial frequency f_0 . The rate at which the allele
215 spreads depends on the ratio of the effective dominance coefficients H_l , H_h . These
216 terms mediate the relative amount of time a beneficial allele spends at low and
217 high frequencies, affecting the probability that a neutral marker recombines away
218 from the selected background. Looking back in time, a proportion r_{eff} of neutral
219 markers become unlinked from the beneficial allele each generation, so when the
220 allele reaches its starting frequency f_0 a proportion P_{NR} of neutral markers remain
221 linked to it [62].

222 Note that for the special case $F = 0$ and $h = 1/2$, $H_l = H_h = 1/2$ and Eq 3

223 reduces to $(1/f_0)^{-(2r/s)}$. This is a standard result for the reduction of diversity
224 following a sweep in outcrossing models with additive dominance [1, 21, 62, 63].

225 **Probability of coalescence from standing variation**

226 When the variant becomes advantageous at frequency f_0 , we expect $\sim 2Nf_0$ haplo-
227 types will carry it. We assume that f_0 remains fixed in time, so that different events
228 occur with constant probabilities. Berg and Coop [21] have shown this assumption
229 provides a good approximation to coalescent rates during the standing phase. The
230 outcome during the standing phase can therefore be determined by considering two
231 competing Poisson processes. The two haplotypes could coalesce; the waiting time
232 for this event is exponentially distributed with rate $1/(2N_e f_0) = (1 + F)/(2N f_0)$,
233 assuming N_e is reduced by a factor $1 + F$ due to self-fertilisation [59]. Alterna-
234 tively, one of the two samples could recombine onto the ancestral background; the
235 exponential mean time for this event is $2r_{eff}(1 - f_0)$ (note the factor of two as
236 there are two samples under consideration). For two competing exponential dis-
237 tributions with rates λ_1 and λ_2 , the probability of the first event occurring *given*
238 *an event happens* equals $\lambda_1/(\lambda_1 + \lambda_2)$ [64]. Hence the probability that two samples
239 coalesce instead of recombine equals:

$$P_{coal} = \frac{\frac{1+F}{2Nf_0}}{\frac{1+F}{2Nf_0} + 2r_{eff}(1 - f_0)} = \frac{1}{1 + 2R(1 - F)f_0(1 - f_0)/(1 + F)} \quad (4)$$

240 where $R = 2Nr$ is the population-scaled recombination rate. Note the presence of
241 the $(1 - F)/(1 + F) = 1 - \sigma$ term, reflecting how selfing reduces the relative effect
242 of recombination by this factor (by both increasing homozygosity, and reducing
243 N_e so coalescence becomes more likely). Hence for a fixed recombination rate R ,

244 samples are more likely to coalesce with increased self-fertilisation, limiting the
245 creation of different background haplotypes. Yet the same coalescent probability
246 can be recovered by increasing the recombination distance by a factor $1/(1 - \sigma)$;
247 that is, if a longer genetic region is analysed.

248 **Effective starting frequency from a de novo mutation**

249 When a new beneficial mutation appears at a single copy, it is highly likely to
250 go extinct by chance [40]. Beneficial mutations that increase in frequency faster
251 than expected when rare are more able to overcome this stochastic loss and reach
252 fixation. These beneficial mutations will hence display an apparent ‘acceleration’
253 in their logistic growth, equivalent to having a starting frequency that is greater
254 than $1/(2N)$ [1, 65–67]. In Section A of the Supplementary *Mathematica* file (S1
255 File; S2 File for PDF copy), we outline how to determine the ‘effective’ starting
256 frequency of hard sweeps that go to fixation, by comparing the sojourn time for the
257 deterministic process to the stochastic diffusion process. We determine that ‘hard’
258 sweeps that go to fixation have the following elevated effective starting frequency:

$$f_{0,A} = \frac{1 + F}{4N_s H_l} \quad (5)$$

259 where $H_l = F + h - Fh$ is the effective dominance coefficient for mutations at
260 a low frequency. This result is consistent with those obtained by Martin and
261 Lambert [67], who obtained a distribution of effective starting frequencies using
262 stochastic differential equations.

263 This acceleration effect can create substantial increases in the apparent f_0 .
264 The effect is strongest for recessive mutations; for example, for $N = 5,000$ and

265 $s = 0.05$ (as used in simulations below), $f_{0,A} = 0.01$ for recessive mutations with
266 $h = 0.1$, an 100-fold increase above $f_0 = 1/(2N) = 0.0001$. $f_{0,A}$ is more modest for
267 additive and dominant mutations; Eq 5 reduces to $1/2Ns$ with $h = 1/2$ or $F = 1$.
268 Hence sweeps from standing variation whose actual f_0 lies below $f_{0,A}$ will produce
269 sweep signatures that may appear similar to hard sweeps. As a consequence, in
270 simulations we use a minimum $f_0 = 0.02$ for adaptation from standing variation
271 cases, which lies above the highest possible value of $f_{0,A}$ for this parameter set.

272 Expected Pairwise Diversity

273 We can use P_{NR} and P_{coal} to calculate the expected pairwise diversity (denoted
274 π) present on a genetic fragment flanking a beneficial allele following a sweep.
275 Looking back in time, one of two possible outcomes can arise. Either two neutral
276 sites linked to the adaptive mutant do not recombine during the sweep phase,
277 and proceed to coalesce during the standing phase. This outcome occurs with
278 probability $P_{NR} \cdot P_{coal}$, creating identical genotypes ($\pi = 0$) since this process
279 occurs rapidly compared to the rate of neutral coalescence. Alternatively, one of
280 the two samples will recombine onto the ancestral background with probability
281 $1 - (P_{NR} \cdot P_{coal})$, so the samples will exhibit background neutral levels of diversity
282 ($\pi = \pi_0$). Hence expected diversity following a sweep equals:

$$\begin{aligned} \mathbb{E}\left(\frac{\pi}{\pi_0}\right) &= 1 - (P_{NR} \cdot P_{coal}) \\ &= \left(\frac{1}{1 + 2R(1 - F)f_0(1 - f_0)/(1 + F)}\right) \cdot \left(1 + \frac{H_l}{H_h} \left(\frac{1}{f_0} - 1\right)\right)^{-r(1-F)/(H_l s)} \end{aligned} \tag{6}$$

283 Eq 6 reflects similar formulas for diversity following soft sweeps in haploid
284 outcrossing populations [15, 21]. Fig 2 plots Eq 6 with different dominance, self-
285 fertilisation, and standing frequency values. The analytical solution fits well com-
286 pared to simulations, although some inaccuracies appear when the mutation ap-
287 pears from a single initial copy. Under complete outcrossing, baseline levels of
288 diversity are restored (i.e., $\pi/\pi_0 \rightarrow 1$) closer to the sweep origin for recessive mu-
289 tations ($h = 0.1$), compared to co-dominant ($h = 0.5$) or dominant ($h = 0.9$) mu-
290 tations. Hence recessive mutations produce weaker signatures of selective sweeps.
291 Dominant and co-dominant mutations produce similar reductions in genetic di-
292 versity, so these cases may be hard to differentiate between from diversity data
293 alone.

294 These patterns can be understood in terms of the underlying allele trajectories
295 (Fig 3). For outcrossing populations, recessive mutations spend most of the sojourn
296 time at low frequencies, maximising the number of recombination events over the
297 sweep history, restoring neutral variation. These trajectories mimic those of sweeps
298 from standing variation, which spend extended periods of time at low frequencies
299 in the standing phase. Conversely, dominant mutations spend most of their time at
300 a high frequency, so there is less chance for neutral markers to recombine onto the
301 ancestral background. Similar results were found by Teshima and Przeworski [33].

302 As the degree of self-fertilisation increases, sweep signatures become similar to
303 the co-dominant case as the derived allele is more likely to spread as a homozy-
304 gote, reducing the influence that dominance exerts over beneficial allele trajectories
305 (Fig 3(b)). In addition, sweep signatures stretch over longer physical regions due
306 to the reduced effective recombination rate [43]. Increasing f_0 also causes sweeps
307 with different dominance coefficients to produce comparable signatures. Here,

308 beneficial mutation trajectories become alike after conditioning on starting at an
309 elevated frequency. In particular, recessive mutations no longer spend the major-
310 ity of their sojourn times at low frequencies, reducing the probability that neutral
311 markers can recombine onto ancestral backgrounds (Fig 3(d)–(f)).

312 Overall, it appears that dominance only strongly affects diversity levels for
313 hard sweeps in outcrossing populations. With increased levels of self-fertilisation,
314 or if the mutation arises from standing variation, allele trajectories (and expected
315 diversity patterns) become similar across different dominance values.

316 **Different Sweep Scenarios can Yield Virtually Identical Signatures**

317 Visual inspection of Fig 2 suggests that different sweep scenarios can produce
318 equivalent reductions in genetic diversity. For example, reductions in diversity
319 caused by a recessive mutation ($h < 0.1$) might be similar to those caused by a
320 mutation with additive dominance ($h = 0.5$) but with a weaker selection coefficient.
321 Similarly, a sweep from standing variation can be mistaken for a weaker hard
322 sweep. Determining how different scenarios cause similar reductions in genetic
323 diversity is useful when testing the most plausible sweep model underlying observed
324 diversity patterns. Berg and Coop [21] argued that it was not possible to find an
325 ‘effective selection coefficient’ \tilde{s} that maps $\mathbb{E}(\pi/\pi_0)$ for a hard sweep onto results
326 expected under a sweep from standing variation. However, we demonstrate in
327 Section A of S3 File (with mathematical analyses in Section C of S1 File) how
328 the argument of Berg and Coop [21] relies on an approximation that only holds
329 when the population-level recombination rate is extremely low (specifically, when
330 $4Nr f_0(1 - f_0) \ll 1$).

331 In fact, a sweep arising from standing variation with selective advantage s

332 can be mapped onto a hard sweep with intensity \tilde{s} , with general self-fertilisation
 333 and $h = 1/2$ (it does not appear possible to obtain a solution for any h). We
 334 equate Eq 6 for general f_0 to the special hard-sweep case $f_0 = 1/2N$ with selection
 335 coefficient \tilde{s} (we do not use $f_{0,A}$ for the hard sweep to calculate tractable analytic
 336 solutions). After equating the two cases and solving for \tilde{s} , we obtain:

$$\begin{aligned} \tilde{s} &= -2r(1-\sigma)\log(2N) \left(\log \left(\left(\frac{1}{f_0} \right)^{-\frac{2r(1-\sigma)}{s}} \frac{1+r(1-\sigma)(2N-1)/(N)}{1+4Nr(1-\sigma)f_0(1-f_0)} \right) \right)^{-1} \\ &\approx -2r(1-\sigma)\log(2N) \left(\log \left(\left(\frac{1}{f_0} \right)^{-\frac{2r(1-\sigma)}{s}} \frac{1}{1+4Nr(1-\sigma)f_0(1-f_0)} \right) \right)^{-1} \end{aligned} \quad (7)$$

337 The approximation in Eq 7 assumes $r_{eff}(2N-1)/(N) \ll 1$. To understand
 338 \tilde{s} , recall that the expected reduction in diversity following a a hard sweep with
 339 $f_0 = 1/2N$ is $(2N)^{-2r(1-\sigma)/\tilde{s}}$ (Eq 6, assuming $H_l = H_h = (1+F)/2$ due to
 340 additive dominance, and $P_{coal} \approx 1$). Inverting this term and solving for \tilde{s} gives
 341 $\tilde{s} = -2r(1-\sigma)\log(2N)/\log(\mathbb{E}(\pi/\pi_0))$. Eq 7 is hence equivalent to the selective
 342 coefficient causing a hard sweep, given that the underlying diversity was actually
 343 shaped by a mutation arising from standing variation.

344 Fig 4(a) plots Eq 7 as a function of R , demonstrating that \tilde{s} increases with
 345 the recombination rate. \tilde{s} can be either less than or greater than s depending on
 346 f_0 and R . Increasing f_0 causes diversity to be restored closer to the beneficial
 347 allele as it is likelier that recombination occurs during the standing phase. Hence
 348 the $f_0 = 0.1$ case is equivalent to a hard sweep caused by a more weakly selected
 349 beneficial allele (Fig 4(b)).

350 In Section A of S3 File (with mathematical analyses in Section C of S1 File)

351 we show that for an outcrossing population with any f_0 , it is possible to find an
352 effective selection coefficient \tilde{s}_h so that a beneficial allele with $h = 1/2$ produces
353 an equivalent sweep pattern to a mutation with arbitrary dominance. We also
354 demonstrate that it is possible to find \tilde{s}_F to map a co-dominant sweep in an
355 outcrossing population onto an equivalent sweep under partial selfing, but this
356 mapping only holds for hard sweeps.

357 Overall, these results caution that it will be necessary to compare a broad
358 range of models when inferring the likeliest cause of selective sweep patterns, and
359 that identifiability issues are to be expected when trying to determine which sweep
360 model best fits diversity data. An example of these issues in relation to investi-
361 gating sweep patterns in humans will be outlined in the section “Application to a
362 selective sweep at the human *SLC24A5* gene”.

363 **Number of Segregating Sites**

364 We can also calculate the total time underlying the genealogy, $\mathbb{E}(T_{tot})$, and there-
365 fore the expected number of segregating sites $\mathbb{E}(S)$. We consider n samples of
366 the derived allele; looking back in time, i of these samples fail to recombine off
367 the derived background during the sweep. The probability of this event can be
368 drawn from a binomial distribution with probability P_{NR} . We denote this value
369 $P_{NR}(i|n) \sim Bin(n, P_{NR})$. Out of these i samples, let k of them recombine dur-
370 ing the sweep phase to create different ancestral backgrounds of the derived allele.
371 Berg and Coop [21] demonstrated how the number of lineages that recombine away

372 from the derived background can be determined using Ewens' Sampling Formula:

$$P_{ESF}(k|i) = S(i, k) \frac{R_{f_0}^k}{\prod_{l=1}^{i-1} (R_{f_0} + l)} \quad (8)$$

373 where $R_{f_0} = 4Nr f_0(1 - f_0)$ is the scaled recombination rate acting on the ancestral
 374 background at frequency f_0 , and $S(i, k)$ are Stirling numbers of the first kind
 375 [15, 21, 68]. Here, we use the rescaled version of R_{f_0} accounting for the reduced
 376 effective recombination rate and effective population size caused by self-fertilisation
 377 (see Eq 4):

$$P_{ESF}(k|i) = S(i, k) \frac{(2R(1 - F)f_0(1 - f_0)/(1 + F))^k}{\prod_{l=1}^{i-1} ((2R(1 - F)f_0(1 - f_0)/(1 + F)) + l)} \quad (9)$$

378 Finally, for the k neutral lineages created in the standing phase, along with the
 379 $n - i$ neutral lineages created in the sweep phase, the expected total time for the
 380 genealogy for all of them, in units of $2N_e$ generations, equals $\sum_{j=1}^{k+n-i-1} 1/j$ [69].
 381 The total time covered by the genealogy is the product of these three terms,
 382 summed over all possible outcomes:

$$\mathbb{E}(T_{tot}) = \sum_{i=0}^n P_{NR}(i|n) \sum_{k=0}^i P_{ESF}(k|i) \sum_{j=1}^{k+n-i-1} 1/j \quad (10)$$

383 $\mathbb{E}(S)$ is $\theta \mathbb{E}(T_{tot})$ where $\theta = 4N_e \mu$ is the population level mutation rate [70].
 384 Equivalent results for outcrossing populations are given by Pennings and Hermis-
 385 son [15, Eq. 15] for adaptation from recurrent mutation, and Berg and Coop [21,
 386 Eq. 10] for adaptation from standing variation. Both these derivations assume
 387 $k > 1$ in the standing phase, as it was argued that $\mathbb{E}(T_{tot}) = 0$ so no segregating
 388 polymorphisms exist. Since simulation results show that this outcome is possible

389 under low recombination rates, we do not include this conditioning in Eq 10.

390 Fig 5 plots $\mathbb{E}(S)$ alongside simulation results. The analytical solution provides
391 a good fit but tends to overestimate simulations, as also observed by Berg and
392 Coop [21]. Also note that fewer segregating sites are present with partial selfing,
393 due to a reduction in the net mutation rate $\theta = 4N_e\mu$ caused by lower N_e .

394 Site Frequency Spectrum

395 The calculations for $\mathbb{E}(S)$ can be extended to determine the full site-frequency
396 spectrum (SFS) following a sweep; that is, the probability that out of n samples,
397 $l = 1, 2 \dots n - 1$ of them carry derived alleles. The full derivation is outlined in
398 Section B of S3 File, and is similar to that used by Berg and Coop [21, Eq 15].
399 However we use a different approach when considering special cases where either
400 all or none of the sampled lineages recombine away from the derived background
401 during the sweep phase. In particular, if all lineages recombine away during the
402 sweep phase, then the SFS reduces to the neutral case; if none do then only a
403 singleton class is included to account for new mutations.

404 Fig 6 plots the expected SFS (Eq B14 in S3 File) alongside simulation results.
405 Analytical results fit simulation data well, although there can be a tendency for
406 it to underestimate the proportion of low- and high-frequency classes ($l = 1$ and
407 9 in Fig 6), and overestimate proportion of intermediate-frequency classes. Hard
408 sweeps in either outcrossers or partial selfers are characterised by a large amount
409 of singletons or highly-derived variants (Fig 6(a)), which is a typical selective
410 sweep signature [71, 72]. As the initial selected frequency f_0 increases, so does
411 the number of intermediate-frequency variants (Fig 6(b)). This signature is often

412 seen as a characteristic of ‘soft’ sweeps [15, 21], reflecting the increased number of
413 genetic backgrounds that the beneficial allele appears on. Yet recessive hard sweeps
414 ($h = 0.1$ and $f_0 = 1/2N$) can produce SFS profiles that are similar to sweeps from
415 standing variation, due to the increased number of recombination events occurring
416 over the timespan of the sweep, especially at low frequencies for long periods of
417 time. As with π/π_0 , SFS patterns will not unambiguously discriminate between
418 sweep scenarios.

419 With increased levels of self-fertilisation, both hard and soft sweep signatures
420 are recovered if measuring the SFS further away from the beneficial allele (Fig 6(c),
421 (d)). For example, a heightened number of intermediate-frequency alleles are ob-
422 served in a sweep from standing variation (Fig 6(d)). Here too, one has to analyse
423 a recombination distance that is $1/(1 - \sigma)$ times longer than in outcrossers to
424 observe soft-sweep behaviour.

425 In the Supplementary *Mathematica* file (Section E of S1 File) we plot SFS
426 results for other recombination distances. In particular, these results demonstrate
427 that with higher f_0 , the SFS becomes similar to the neutral case over a shorter
428 recombination distance than for hard sweeps, as reflected with results for expected
429 pairwise diversity (Eq 6).

430 **Soft sweeps from recurrent mutation**

431 Until now, we have only focussed on a ‘soft’ sweep that arises from standing
432 variation. An alternative type of ‘soft’ sweep is one where recurrent mutation at the
433 selected locus introduces the beneficial allele onto different genetic backgrounds.
434 We can examine this case by modifying existing results. Pennings and Hermisson

435 [15] demonstrated that the expected reduction in pairwise diversity $\mathbb{E}(\pi/\pi_0) =$
436 $1 - [(P_{coal,M})(P_{NR})]$ where $P_{coal,M} = 1/(1+2\Theta_b)$ is the probability that two samples
437 are identical by descent instead of arising on different genetic backgrounds by
438 independent mutation events. Here, $\Theta_b = 2N_e\mu_b$ is the population level mutation
439 rate at the beneficial locus. We can compare the signatures of these two different
440 types of soft sweeps by using this solution, with P_{NR} as given by Eq 3 with $f_0 =$
441 $1/(2N)$, and $\Theta_b = 2N_e\mu_b = (2N\mu_b)/(1+F)$ in $P_{coal,M}$.

442 Fig 7(a), (b) compares $\mathbb{E}(\pi/\pi_0)$ in the standing variation case, and for the re-
443 current mutation case, under different levels of self-fertilisation. Several differences
444 are apparent. First, while dominance only weakly affects sweep signatures arising
445 from standing variation, it more strongly affects sweeps from recurrent mutation
446 in outcrossing populations, as the underlying allele trajectories are affected by
447 the level of dominance since each variant arises from an initial frequency $\sim 1/(2N)$
448 (Fig 3). Second, both models exhibit different behaviour close to the selected locus
449 ($R \rightarrow 0$). The recurrent mutation model has diversity levels that are greater than
450 zero, while the standing variation model exhibits no diversity. As R increases,
451 diversity reaches higher levels in the standing variation case than for the recurrent
452 mutation case. To determine the recombination rate when the recurrent mutation
453 model exhibits higher diversity than the standing variation model, we assume that
454 close to the adaptive mutant, it is very unlikely for samples to recombine during
455 the sweep phase (i.e., $P_{NR} \approx 1$). It remains to determine when $P_{coal,M}$ is higher
456 than P_{coal} under standing variation, which occurs when:

$$\begin{aligned} R \leq R_{Lim} &= \frac{\Theta_b}{f_0(1-f_0)(1-F)} \\ &\approx \frac{\Theta_b}{f_0(1-F)} \text{ for } f_0 \ll 1 \end{aligned} \quad (11)$$

457 Hence for a fixed Θ_b , the window where recurrent mutations creates higher
458 diversity near the selected locus increases for lower f_0 or higher F , since both
459 these factors reduces the potential for recombination to create new haplotypes
460 during the standing phase. Eq 11 accurately reflects when standing variation
461 sweeps exhibit higher diversity (Fig 7(a), (b)), but becomes inaccurate for $h = 0.1$
462 in outcrossing populations. Here, beneficial alleles have elevated fixation times, so
463 some recombination is likely to occur during the sweep phase. We also observe
464 that for higher selfing rates, the ratio of π_{SV} (diversity under sweep from standing
465 variation) to π_M (diversity under sweep from recurrent mutation) becomes higher
466 than in outcrossers (compare Fig 7(c) with 7(d)). This is because the effects of
467 sweeps arising from recurrent mutation on diversity becomes diluted over a longer
468 genetic distance under self-fertilisation, due to weakened effects of recombination.

469 We can also modify the expected SFS to account for recurrent mutation dur-
470 ing the standing phase (see Section B in S3 File for details). These calculations
471 verify that, close to the selected locus, sweeps from recurrent mutations show
472 more intermediate-frequency variants than sweeps from standing variation. The
473 situation is reversed once R exceeds R_{Lim} .

474 **Distance between singletons**

475 A selective sweep increases the mean distance between ‘singletons’, which are de-
476 rived alleles that are only observed on a single haplotype. This phenomena was
477 recently used to detect evidence for recent human adaptation [56]. We hence
478 ran computer simulations to investigate the distribution of distances between the
479 beneficial locus and the nearest singleton under different scenarios.

480 **Singleton distances in fixed sweeps**

481 We first measured the distance from the beneficial allele to the nearest singleton
482 across 50 samples taken from a fixed sweep. These distances are compared to those
483 obtained from the neutral background before a beneficial mutation was introduced
484 (see Fig 11(a) in the Methods for a schematic). Due to the computational limita-
485 tions of individual-based simulations, a large number of samples did not contain
486 singletons (Fig 8(a)). Focussing on samples containing singletons in the neutral
487 population, they are likelier to lie close to the target locus (Fig 8(b)). Sweeps
488 reduce the overall frequency of observed singletons, and also increases the distance
489 from the selected allele to the nearest singleton. These distributions are visibly
490 different for sweeps of different dominance effects; recessive mutations ($h = 0.1$)
491 cause a much stronger reduction in observed singleton densities than dominant
492 adaptations ($h = 0.9$). This behaviour likely arises as recessive mutations increase
493 in frequency closer to the present time, while dominant mutations reach a higher
494 frequency earlier on (Fig 3). The rapid increase in frequency of recessive mutations
495 in the recent past makes it even less likely for singletons to appear on selected back-
496 grounds. This result is reflected in the SFS, where hard sweeps caused by recessive

497 mutations also display a lower number of singletons (Fig 6(a)).

498 We showed that in outcrossing populations, a sweep arising from a recessive
499 or dominant mutation can cause the same reduction in diversity as that caused
500 by a co-dominant mutation, after rescaling the selection coefficient (Section A
501 in S3 File). Hence we next measured the distribution of singleton distances for
502 co-dominant sweeps but with different selection coefficients, to determine if sim-
503 ilar patterns are produced to cases with different dominance. Weakly-selected
504 mutations ($s = 0.01$) exhibit results that are similar to the neutral case, while
505 strongly-selected mutations ($s = 0.09$) show a clear reduction in singleton densi-
506 ties (Fig 8(c), (d)). These patterns are opposite to what is observed for recessive
507 and dominant mutations respectively, implying that singleton densities may pro-
508 vide clearer evidence regarding the dominance underlying a selective sweep.

509 **Singleton distances in partial sweeps**

510 We next investigated singleton distances from partial sweeps (i.e., those that have
511 not completely fixed in the population). Specifically, we look at sweeps that have
512 reached a frequency of 70% when they were sampled. The neutral expectation
513 was calculated by measuring singleton distances around SNP that lie between a
514 frequency of 65% – 75% (Fig 11(b) in the Methods). For the neutral case, there
515 are always more singletons observed on the derived background, since it is present
516 at a higher frequency (Fig 9(a)). Focussing on samples where a singleton was
517 observed, we then see that the distributions are similar between ancestral and
518 derived backgrounds (Fig 9(b)). On selected backgrounds, there are many more
519 samples not carrying singletons (Fig 9(a)). For samples carrying singletons, fewer
520 of them lie closer to the target locus on derived backgrounds, compared to ancestral

521 backgrounds. Furthermore, singleton distances are uniformly distributed along
522 the genetic tract on derived backgrounds, with visibly similar distances occurring
523 irrespective of the dominance level (Fig 9(b)). Hence while singleton distances can
524 provide evidence of ongoing adaptation, there appears to be very little power to
525 infer the dominance level of the mutation.

526 In Section C of S3 File, we show that increasing either f_0 or F weakens the
527 effect that h has on singleton distance distributions, in line with previous results
528 showing how an increase in either of these values weakens the effect that dominance
529 has on summary statistics. We also show that increasing the number of samples
530 under investigation (from 50 to 1000) weakens the ability of singleton distributions
531 to detect fixed sweeps as singleton distances will only be affected with very recent
532 (i.e., very strong) selection [56]. However, evidence of an ongoing sweep (i.e.,
533 one observed at a frequency of 0.7) can still be seen if taking a large number
534 of samples, as the distributions are markedly different between the ancestral and
535 derived backgrounds.

536 **Application to a selective sweep at the human *SLC24A5*** 537 **gene**

538 To demonstrate how these sweep models can be used to infer properties of genetic
539 adaptation, we reanalyse a selective sweep at the *SLC24A5* gene in European
540 human populations. The rs1426654 SNP harbours a G \rightarrow A substitution that is
541 strongly associated with skin pigmentation in Eurasian populations [73, 74]. It
542 was long assumed that the derived A mutation was only present at a negligible
543 frequency in Africa, yet recent data has shown it to be present at an elevated

544 frequency in East Africa [74]. These East African populations harbour the same
545 extended haplotype as in Eurasia, suggesting that the mutation was reintroduced
546 into Africa following the out-of-Africa human expansion. Nevertheless, the recent
547 discovery of these new haplotypes begs the question of whether the derived SNP
548 was introduced into Eurasia at an elevated frequency or not. Hence we performed
549 a maximum-likelihood fit of these analytical solutions to the sweep signature pro-
550 duced around the derived SNP in Europe, to determine whether it is consistent
551 with a hard sweep, or instead one from either standing variation or recurrent
552 mutation.

553 We downloaded diversity data from European populations in the 1000 Genomes
554 phase 3 release, and fitted models to diversity data around the derived SNP (see
555 Methods and Section G of S1 File for details). We implemented a nested model
556 comparison, to test for the presence of either a sweep from standing variation, or
557 from recurrent mutation. In both cases we also tested for the presence of non-
558 additive dominance ($h \neq 1/2$). Results are outlined in Table 2. For the standing
559 variation case, the best fitting model implicated that the sweep arose from a new
560 mutation (a ‘hard sweep’) with additive dominance, with a selection coefficient
561 $s = 0.065$ (see Fig 10(a) of a fit of this model to the sweep region). Models that
562 included an elevated initial frequency also estimated unrealistically high selection
563 coefficients, with s nearly equal to a thousand. These findings suggest that large
564 sweep signatures, such as those observed in the *SLC24A5* gene, are extremely
565 unlikely to be formed by adaptations arising from standing variation, in line with
566 theoretical work (see also [21]). It was also not possible to discern a sweep assuming
567 additive dominance from non-additive dominance; analysis of the likelihood surface
568 shows that a ridge of maximum likelihood exists for constant hs , reinforcing the

569 idea that it is not easy to discern non-additive dominance from diversity data alone
 570 (Section G of S1 File).

Model	Parameters	s	h (1/2)	x_0 (1/2 N_e)	Θ (0)	LL	ΔAIC
<i>HS, AD</i>	1	0.065	–	–	–	-4982.57	846
HS, NAD	2	0.15	0.18	–	–	-4982.57	848
SV, AD	2	815	–	0.017	–	-4207.29	NA
SV, NAD	3	933	0.82	0.017	–	-4207.29	NA
RM, AD	2	0.20	–	–	0.56	-4134.14	0
RM, NAD	3	0.26	0.37	–	0.56	-4134.14	2

Table 2. Results of maximum-likelihood model fitting of *SLC24A5* sweep signature. Results are presented for a hard sweep model (‘HS’); from standing variation (‘SV’), or from recurrent mutation (‘RM’). We also consider additive or non-additive dominance (denoted AD, NAD respectively). Numbers in brackets next to each parameter heading are the fixed values if they are not estimated for that particular model (as represented by a dash). ΔAIC is the difference in AIC between that model and the best fitting one (RM, AD, which is highlighted in bold). The italicised model HS, AD is the best fitting realistic model.

571 For the recurrent mutation model, the best-fitting model included a significant
 572 level of mutation at the target SNP ($\Theta = 0.56$). However, this high mutation rate
 573 leads to elevated diversity levels around the target SNP, which is not present in
 574 observed data (Fig 10(a)). We also tested for the presence of recurrent mutation
 575 by measuring H -statistics around the sweep region [25] (see Methods for formal
 576 definitions of these statistics), which measure the relative frequency of different
 577 haplotypes across samples. Specifically, high H_{12} , low H_2/H_1 values are consistent
 578 with a single haplotype fixing, in line with a hard sweep. Conversely, a reduced
 579 H_{12} and elevated H_2/H_1 values suggest multiple haplotypes fixing, which occurs
 580 following adaptation from standing variation or recurrent mutation. Fig 10(b)
 581 demonstrates that around the target SNP, H_1 is close to 1 while H_2/H_1 is near

582 zero. Both results indicate that a single haplotype has fixed around the target
583 SNP, which is not expected following a sweep from recurrent mutation [15]. It
584 seems that the recurrent mutation model had the highest likelihood due to spikes
585 of high diversity around the target SNP, which can be mistaken for a recurrent
586 mutation effect if not checked against other analyses.

587 These models assume a fixed population size, but it is known that humans
588 have a complex demographic history. European populations have likely undergone
589 a contraction following migration from Africa, followed by extensive population
590 growth [75]. To determine if this demography could have drastically affected our
591 inference of different sweep signatures, we ran simulations using MSMS with in-
592 ferred parameters to determine how sweep signatures are affected by this demo-
593 graphic history. Yet even under a growth-bottleneck model, a hard sweep model
594 fits the observed sweep pattern better than either of the soft sweep models, after
595 rescaling parameters by the different present-day N_e (Section D in S3 File, with
596 plots also available in Section G of S1 File).

597 Furthermore, the derived A allele is present in African populations but at a low
598 frequency (in 55 of 1063 African haplotypes in the 1000 Genomes dataset). This
599 begs the question of whether the derived allele was introduced into Eurasia, but
600 at too low a frequency to influence the maximum-likelihood model fit. Fig 10(c)
601 shows phylogenetic trees of 20Kb regions either surrounding the target SNP, or up-
602 stream, downstream of the SNP. We observe that most European samples carrying
603 the derived mutation cluster together, reflecting recent appearance and spread of
604 the derived allele. However, within these clades we also observe some African hap-
605 lotypes carrying the derived allele, suggesting that it was introduced into Eurasia
606 due to out-of-Africa migration.

607 Overall, our model analyses determined that the derived SNP at the *SLC24A5*
608 gene most likely followed ‘hard’ sweep dynamics. However, we also find evidence
609 for ancestral African haplotypes forming the basis of the sweep. Hence the likeliest
610 outcome is that the derived allele was introduced into Eurasia at a sufficiently low
611 frequency so that its sweep dynamics were indistinguishable from a hard sweep.
612 Given a selection coefficient of 0.065, co- dominance ($h = 0.5$) and $N_e = 10,000$,
613 Eq 5 predicts an $f_{0,A}$ of 0.7%. It is likely that the derived haplotype was introduced
614 at a lower frequency than this value.

615 Discussion

616 Summary of Theoretical Findings

617 While there has been many investigations into how different types of adaptation
618 can be detected from next-generation sequence data [11, 13, 76, 77], these models
619 assumed idealised sexually reproducing populations and beneficial mutations that
620 have additive dominance ($h = 0.5$). Here we have created a general model of
621 a selective sweep, with arbitrary levels of self-fertilisation and dominance. Our
622 principal focus is on comparing a ‘hard sweep’ arising from a single allele copy to a
623 ‘soft sweep’ arising from standing variation, but we have also considered the effect
624 of adaptation from recurrent mutation (Fig 7).

625 We find that the qualitative patterns of different selective sweeps under selfing
626 remain similar to expectations from classic outcrossing models. In particular, a
627 sweep from standing variation still creates an elevated number of intermediate-
628 frequency variants compared to a sweep from *de novo* mutation (Figs 6, 7). This

629 pattern is a known signature of a ‘soft sweep’ [11,13,15,21], meaning that common
630 statistical methods used for detecting them (e.g., observing an higher number of
631 haplotypes than expected [24, 25]) can, in principle, still be applied to selfing
632 organisms (but see the discussion below with regards to dominance). Under self-
633 fertilisation, these signatures are stretched over longer physical regions than in
634 outcrossers. These extensions arise as self-fertilisation affects gene genealogies
635 during both the sweep and standing phases, but in different ways. During the
636 sweep phase, beneficial alleles fix more rapidly under higher self-fertilisation as
637 homozygote mutations are created more quickly [41,42]. In addition, the effective
638 recombination rate is reduced by approximately $1 - F$ [43]. These two effects
639 mean that neutral variants linked to an adaptive allele are less likely to recombine
640 onto the neutral background during the sweep phase, as reflected in Eq 1 for
641 P_{NR} . During the standing phase, two samples are more likely to coalesce with
642 increased self-fertilisation since N_e is decreased by a factor $1/(1 + F)$ [59]. This
643 effect, combined with an reduced effective recombination rate, means that the
644 overall probability of recombination during the standing phase is reduced by a
645 factor $1 - \sigma$ (Eqs 4, 9, B14 in S3 File). Hence intermediate-frequency variants,
646 which could provide evidence of adaptation from standing variation, will be spread
647 out over longer genomic regions. The elongation of sweep signatures means soft
648 sweeps can be easier to detect in selfing organisms than in outcrossers, since the
649 differences in diversity caused by sweeps are spread out over longer regions.

650 We have also investigated how dominance affects soft sweep signatures, since
651 previous analyses have only focussed on how hard sweeps are affected with differ-
652 ent dominance effects [33–35]. In outcrossing organisms, recessive mutations leave
653 weaker sweep signatures than additive or dominant mutations as they spend more

654 time at low frequencies, increasing the amount of recombination that restores neu-
655 tral variation (Figs 2, 3). With increased self-fertilisation, dominance has less of an
656 impact on sweep signatures as most mutations are homozygotes (Fig 3). However,
657 dominance has different effects on separate types of ‘soft’ sweeps. Dominance only
658 weakly affects sweeps from standing variation, as trajectories of beneficial alleles
659 become similar once the variant’s initial frequency greatly exceeds $1/2N$ (Figs 2, 3).
660 Yet different dominance levels can affect sweep signatures if the beneficial allele is
661 reintroduced from recurrent mutation (Fig 7). Hence if one wishes to understand
662 how dominance affects selective sweep signatures, it is also important to consider
663 the type of selective sweep underlying observed genetic diversity. We also showed
664 how beneficial variants of different dominance values create distinct alterations
665 in the distances to the nearest singleton (Fig 8). These results suggest that the
666 distribution of low-frequency variants around a sweep can provide information on
667 the dominance value underlying it. Investigating the utility of singletons to de-
668 tect dominance effects seems a worthy future research direction, especially since in
669 our example of estimating properties of the *SLC24A5* sweep, it is tricky to infer
670 non-additive dominance from diversity data alone.

671 We also derived an ‘effective selection coefficient’ \tilde{s} so that sweeps from standing
672 variation will produce a pattern of diversity reduction equivalent to a hard sweep
673 (Eq 7; Fig 4), and an \tilde{s}_h so that a non-additive sweep in an outcrossing population
674 can be mapped onto a co-dominant sweep (Section A in S3 File). These derivations
675 imply that different types of sweep models can lead to similar outcomes, which may
676 prove problematic when making inferences from genomic data [78, Supplementary
677 Material]. Yet it may be apparent if some sweep signatures arise from standing
678 variation or not, if unrealistic parameters are needed to produce expected patterns

679 of diversity. In particular, for the *SLC24A5* sweep to appear from standing vari-
680 ation, the underlying selection coefficient must be unrealistically large (Table 2).
681 Hence adaptation from elevated standing variation (greater than 0.7%) is unlikely
682 for this case.

683 **Soft sweeps from recurrent mutation or standing variation?**

684 Our theoretical results shed light onto how to distinguish between soft sweeps that
685 arise from either standing variation, or from recurrent mutation. Both models
686 are characterised by an elevated number of intermediate-frequency haplotypes, in
687 comparison to a hard sweep. Yet sweeps arising from recurrent mutation produces
688 intermediate-frequency haplotypes closer to the beneficial locus, while sweeps from
689 standing variation produce intermediate-frequency haplotypes further away from
690 the adaptive locus (Fig 7 and Section B in S3 File). Eq 11 provides a simple
691 condition for the recombination distance needed so a sweep from standing variation
692 exhibits higher diversity than one from recurrent mutation. The size of this region
693 increases under higher self-fertilisation.

694 This result has implications for inferring different types of sweeps. If multiple
695 swept haplotypes are present over long genetic distances, this observation im-
696 plies that the beneficial allele underlying the sweep likely originated from standing
697 variation as opposed to recurrent mutation. This phenomenon could explain the
698 elevated H_2/H_1 statistics, and reduced H_{12} values upstream of the *SLC24A5* SNP
699 (Fig 10(b)), especially given that we know the derived SNP to be present at a low
700 frequency in Africa. However, if this was truly a selective sweep arising from an
701 elevated starting frequency, we also expect elevated H_2/H_1 values downstream of
702 the SNP, which we do not observe. A simpler explanation for the elevated haplo-

703 type diversity is that the recombination rate is higher upstream of the SNP than
704 downstream, which has broken down the sweep signature to a greater extent in
705 this region (see Fig 12 in the Methods for the actual recombination map).

706 Different haplotype structure between sweeps from either standing variation or
707 recurrent mutation should be more pronounced in self-fertilising organisms, due
708 to the reduction in effective recombination rates. However, if investigating sweep
709 patterns over longer genetic regions, it becomes likelier that genetic diversity will
710 be affected by multiple beneficial mutations spreading throughout the genome.
711 Competing selective sweeps can lead to elevated diversity near a target locus for
712 two reasons. First, selection interference increases the fixation time of individual
713 mutations, allowing more recombination that can restore neutral diversity [79]. In
714 addition, competing selective sweeps can drag different sets of neutral variation to
715 fixation, creating asymmetric reductions in diversity [80]. Further investigations
716 of selective sweep patterns across long genetic distances will prove to be a rich area
717 of future research.

718 **Using models to determine properties of selective sweeps**

719 **Analysis of the *SLC24A5* sweep signature**

720 An emerging approach to quantifying properties of genetic adaptation involve fit-
721 ting sweep models to regions displaying high substitution rates compared to an
722 outgroup [78, 81, 82]. Inspired by these works, we demonstrated how the general
723 sweep models can be used to determine adaptation properties by applying them to
724 the *SLC24A5* gene in European humans. Overall, the sweep pattern best matches
725 a classic ‘hard’ sweep signature (Table 2; Fig 10). However, since the derived

726 allele is known to be present at a low frequency in Africa, it also appears that the
727 derived allele was introduced into Eurasia at a sufficiently low frequency so that
728 the resulting signature is equivalent to a ‘hard’ sweep, even if the mutation did not
729 appear after out-of-Africa migration (Fig 10(c)). This analyses demonstrates how
730 adaptive mutations arising from standing variation have to be present at a suffi-
731 ciently high frequency (above the ‘accelerated’ $f_{0,A}$ given by Eq 5) to be reliably
732 distinguished from classic hard sweeps. In addition, analysis of this specific sweep
733 region also demonstrates the utility of combining model fitting of genetic diversity
734 with other statistics (e.g., haplotype structure, phylogenetic relationships) to fully
735 work out the evolutionary history of individual selective sweeps.

736 One potential difficulty arising out of model analysis is that of estimating dom-
737 inance coefficients. Sweep models where h was non-additive did not explain the
738 data better than a co-dominant sweep. Nevertheless, there are several ad-hoc
739 reasons why the underlying mutation is likely to be approximately co-dominant.
740 Recessive hard sweeps appear similar to sweeps from standing variation (with a
741 weaker reduction in diversity at linked regions) and are heterozygous for long pe-
742 riods of time (Fig 3(a)). Hence the strong sweep signature, and high frequency
743 of the derived allele in European populations, makes it unlikely for this muta-
744 tion to be recessive. Similarly, strongly dominant mutations take a long period of
745 time to fully fix, in contrast to the observed near-fixation of the derived *SLC24A5*
746 SNP. It will be important to extent inference methods to more accurately quantify
747 dominance of adaptive mutations. One promising approach could be to analyse
748 singleton densities, which appear to differ under recessive and dominant sweeps
749 (Fig 8).

750 **Potential model applications to self-fertilising organisms**

751 Existing software for finding sweep signatures in nucleotide data are commonly
752 based on finding regions with a site-frequency spectrum matching what is ex-
753 pected under a selective sweep [83, 84]. The more general models developed here
754 can therefore be used to create more specific sweep-detection methods while ac-
755 counting for self-fertilisation. However, a recent analysis found that signatures of
756 soft sweeps can be incorrectly inferred if analysing genetic regions that flank hard
757 sweeps, which was named the ‘soft shoulder’ effect [85]. Due to the reduction in
758 recombination in selfers, these model results indicate that ‘soft-shoulder’ footprints
759 could be present over long genetic distances, and should be accounted for. One
760 remedy to this problem is to not just classify genetic regions as being subject to
761 either a ‘hard’ or ‘soft’ sweep, but also as being linked to a region subject to one
762 of these sweeps [27].

763 Further investigations of selective sweeps under self-fertilisation will also be
764 aided by the creation of new simulation methods that account for this mating
765 system. It is common to test sweep models by comparing results to coalescent
766 simulations of adaptation [86, 87], but existing simulations do not account for self-
767 fertilisation. Creating new simulation programs will prove important to further
768 explore other key properties of selective sweeps (e.g., haplotype structure, singleton
769 densities, power calculations) under selfing across larger sample and population
770 sizes. We therefore hope that these results will stimulate the creation of new
771 simulation and inference software to further explore how adaptation is affected by
772 different reproductive modes.

773 **Methods**

774 **Exact simulations, including dominance and self-fertilisation**

775 Simulations were coded in C and are based on Wright-Fisher population dynamics.
776 These are available in S4 File or online at [https://github.com/MattHartfield/](https://github.com/MattHartfield/DomSelfAdapt)
777 `DomSelfAdapt`. There exists N diploid individuals, each containing two haplo-
778 types consisting of a stretch of genetic material at which neutral mutations can
779 accumulate via an infinite-sites model. The far left hand side of the tract contains
780 the locus at which the beneficial allele can arise.

781 Each generation the entire population is replaced. First, the number of self-
782 fertilisation reproductions is drawn from a Binomial distribution with probability
783 σ . It is then decided which specific reproduction events will occur by selfing. To
784 create offspring, a first parent is chosen with probability proportional to its fitness,
785 then one of its two haplotypes is selected with equal probability. If selfing arises,
786 then the offspring's second haplotype is chosen from the same parent, which could
787 be the same as the first. Otherwise a second parent is selected, with probability
788 proportional to its fitness, then one of its haplotypes is chosen. The number of
789 recombination events per haplotype is drawn from a Poisson distribution with
790 mean r . Crossover locations are uniformly distributed over the fragment length.
791 Offspring haplotypes are subsequently created by initially copying over the first
792 sampled parental haplotype, then switching over to copying the second parental
793 haplotype after passing a recombination breakpoint. Selection and recombination
794 is repeated in this manner for all N individuals.

795 New neutral polymorphisms are then added. The number of mutations to be
796 added to the entire population is chosen from a Poisson distribution with mean

797 $2N\mu$. For each new mutation, it can appear in one of the $2N$ haplotypes with equal
798 probability, with its location selected from a uniform distribution. A ‘garbage-
799 collection’ routine is then executed to remove non-polymorphic loci. Fig 11(a)
800 outlines how polymorphisms are distributed in the simulation.

801 The simulation is split into two parts. A ‘burn-in’ phase is run first to generate
802 background neutral diversity, where the population evolves without any beneficial
803 alleles present for $20N$ generations. 100 different populations are created for each
804 neutral parameter set. In the second part, the adaptive mutation is introduced
805 into a single haplotype chosen at random; it is initially neutral until its frequency
806 matches or exceeds f_0 , at which point it has selective advantage s and dominance
807 coefficient h acting upon it. We can set $f_0 = 1/2N$ so that the mutation is
808 beneficial from the outset (a ‘hard’ sweep). The beneficial allele is then tracked
809 until it is either lost, or reaches the ‘census’ frequency at which the selective sweep
810 is analysed, after which we randomly sample haplotypes from the population to
811 create final outputs.

812 **Measuring mean pairwise diversity; number of segregating sites; site** 813 **frequency spectrum**

814 After the beneficial allele has gone to fixation, we sampled 10 haplotypes 10 times
815 from each burn-in population to create 1000 simulation estimates. For each of
816 these statistics, mutations are placed in one of 10 bins depending of the distance
817 from the sweep, with the relevant statistic calculated per bin. Mean values, along
818 with 95% confidence intervals, are calculated over all 1000 outputs.

819 **Measuring distances between singleton mutations**

820 We sampled 50 or 1000 haplotypes once from each base population, creating 100
821 total datasets. We also sample the same number of haplotypes from the burn-in
822 population to determine the neutral distribution of distances.

823 We investigated cases where the sweep has either gone to fixation, or where
824 the population is sampled after the beneficial allele exceeds a frequency of 0.7.
825 When the beneficial allele is sampled at fixation, the distance from the adaptive
826 locus to the nearest singleton is measured over all samples. The distance is nor-
827 malised to between 0 and 1, where 0 is the location of the selected locus and 1
828 the furthest right-hand edge. We also note how many samples did not contain sin-
829 gletons. When the sweep is sampled at a frequency of 0.7, we measure singleton
830 distances separately for samples carrying either the ancestral or derived allele. For
831 the neutral burn-in population, we first found derived alleles that were present at
832 a frequency between 0.65 and 0.75. For each of these, we measured the upstream
833 distance to the nearest singleton, if present. If not, we check if a singleton existed
834 downstream of the reference variant, and the singleton distance is calculated as the
835 distance of the nearest singleton from the left-hand edge of the genome, plus the
836 upstream distance from the reference variant to the right-hand edge (Fig 11(b)).
837 Summing distances in this manner is valid as we assume polymorphisms are uni-
838 formly distributed throughout the genome. Otherwise we noted if no singleton
839 existed on the haplotype.

840 **Human sweep data analyses**

841 **Data processing**

842 Data was retrieved from the 1000 Genomes phase 3 version 3 integrated call set
843 (<ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>) [88]. The
844 five European populations (CEU, FIN, GBR, IBN, TSI) were investigated; re-
845 lated individuals were removed ([ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/
846 release/20130502/20140625_related_individuals.txt](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/20140625_related_individuals.txt)) giving 503 total indi-
847 viduals. SNP data was obtained using *VCFtools* [89] over a 1Mb region, between
848 locations 47,930,001 and 48,930,000 on Chromosome 15 (the rs1426654 target SNP
849 is at location 48,426,484). Only biallelic SNPs in Hardy-Weinburg equilibrium
850 (with P -value greater than 10^{-6}) were retained; indels were removed. Pairwise
851 diversity was calculated in 20Kb bins over this region. Baseline diversity (i.e.,
852 that expected in the absence of a selective sweep) was determined by calculat-
853 ing mean diversity values at flanking regions both up- and downstream of the
854 sweep. Specifically, we measure the mean diversity between locations 47,930,001
855 and 48,220,000 upstream of the target SNP, and between locations 48,640,000 and
856 48,930,000 downstream of the target SNP (Fig 12(a)). Diversity estimates up- and
857 downstream were divided by the mean values between these regions (Fig 12(b)).
858 Sex-averaged recombination maps for each bin were obtained from Bh erer *et al.* [90]
859 (Fig 12(c)).

860 **Model fitting**

861 Sweep models were fitted to this diversity data using the maximum likelihood
862 procedure of Sattath *et al.* [81]. Two nested models were considered; one where a

863 sweep arose from standing variation (Equation 6), or where the sweep arose from
864 recurrent mutation (as described in the ‘Soft sweeps from recurrent mutation’
865 section). Since we are analysing human data assuming a fixed population size,
866 we set $F = 0$ and $N_e = 10,000$ [91]. Due to the large number of polymorphisms
867 per bin, we assume that observed pairwise diversity at recombination distance r is
868 normally distributed with mean values equal to the expected values given by the
869 models (denoted $m(r)$), and variance $v(r) = m(r)(1 - m(r))/n$ for n the number of
870 segregating sites in that bin. The log-likelihood for the data under these models,
871 as measured over all b bins, equals $-\sum_b(\log(2\pi v(r))/2 + (\hat{K}(r) - m(r))^2/(2v(r)))$,
872 where $\hat{K}(r)$ is the relative diversity in each bin.

873 Maximum likelihood for each model was found using the ‘FindMaximum’ func-
874 tion in *Mathematica* version 11.2 [92]. In all models we estimated the selection
875 coefficient s . We then used a nested model structure to determine if evidence
876 existed for non-additive dominance ($h \neq 1/2$); standing variation of the selective
877 sweep ($f_0 > 1/2N_e$); or recurrent mutation at the target SNP location ($\Theta \neq 0$).
878 We set options in ‘FindMaximum’ so that $s > 0$, and $0 < h < 1$, $f_0 > 1/2N_e$
879 and $\Theta > 0$ if these parameters were not fixed. We compared six models: (i) fixed
880 $h = 1/2$, $f_0 = 1/2N_e$ (hard sweep with additive dominance); (ii) variable h , fixed
881 $f_0 = 1/2N_e$ (hard sweep with non-additive dominance); (iii) fixed $h = 1/2$, vari-
882 able f_0 (standing variation sweep with additive dominance); (iv) variable h , f_0
883 (standing variation sweep with non-additive dominance); (v) fixed $h = 1/2$, vari-
884 able Θ (recurrent mutation with additive dominance) (iv) variable h , Θ (recurrent
885 mutation with non-additive dominance). Note that for hard sweep models, we do
886 not use $f_{0,A}$ to ensure a tractable model fit. Using $f_0 = 1/2N_e$ should not prove
887 problematic for inferring different types of sweeps, as long as estimated f_0 for the

888 standing variation cases lie above $f_{0,A}$, so the two cases can be differentiated. Since
889 estimated $f_0 \sim 1.7\%$ and $f_{0,A} \sim 0.7\%$, this condition is fulfilled.

890 To calculate the H statistics of Garud *et al.* [25], haplotype counts in each of the
891 20Kb bins were obtained using the ‘--hapcount’ function in *VCFtools*. From these
892 the relevant haplotype statistics were calculated per bin. Let there be K unique
893 haplotypes in a bin, ordered so that p_1 is the frequency of the most common
894 haplotype, p_2 the frequency of the second common haplotype, and so on. Then
895 $H_1 = \sum_i^K (p_i^2)$, $H_{12} = (p_1 + p_2)^2 + \sum_{i=3}^K (p_i)^2$, and $H_2 = H_1 - p_1^2$. We also calculated
896 the ratio H_2/H_1 .

897 Human Sweep Simulations

898 We ran simulations of the selective sweep using MSMS [87] to determine expected
899 diversity patterns under different demographic scenarios. To ensure tractable sim-
900 ulations, we simulated 100 haplotypes using a genetic region of length 200Kb, with
901 the selected site located in the middle of the region. The scaled neutral mutation
902 rate $4N_e\mu$ equalled 188.8 (assuming $N_e = 10,000$), reflecting a per-basepair rate
903 of 2.36×10^{-8} as recently used by Field *et al.* [56]; the scaled recombination rate
904 $2N_e r$ was set to 55.4 reflecting the sex-averaged recombination rate over the region
905 as determined by Bhérier *et al.* [90]. Three sweep scenarios were simulated: (i) a
906 hard sweep (ii) a sweep from standing variation with initial selected frequency
907 1.7% (iii) a sweep from recurrent mutation with $\Theta = 2N_e\mu_b = 0.56$. Input val-
908 ues reflect those obtained from the maximum likelihood model fitting. Simulations
909 were run assuming two demographic scenarios; either a constant population of size
910 $N_e = 10,000$, or a growth-bottleneck demographic mimicking human migration out
911 of Africa (parameters used are outlined in Fig 1(d) of Schrider *et al.* [93]). For

912 the latter model, other parameters were scaled by the present-day $N_e = 35,900$.
913 In both the growth-bottleneck models and constant-sized models assuming a 1.7%
914 starting frequency, MSMS requires the user to set a time in the past when se-
915 lection started acting on the beneficial mutation. In these cases, starting times
916 were set so that the sweep reached fixation in the present day. We also simulated
917 pairwise diversity from a neutral growth-bottleneck demographic scenario, to de-
918 termine expected baseline diversity in the absence of a selective sweep. All results
919 are averages over 1,000 simulation runs. A complete list of command lines and
920 parameters are outlined in S5 Table.

921 **Phylogenetic analyses**

922 Biallelic SNPs in Hardy-Weinburg equilibrium ($P > 10^{-6}$) were extracted from the
923 five European populations and the five African populations (ESN, GWD, LWK,
924 MSL, YRI) in the 1000 Genomes dataset, in bins of size 20Kb, from between
925 basepair locations 48,320,000–48,340,000, 48,420,000–48,440,000, and 48,500,000–
926 48,520,000 on chromosome 15. Distance matrices were then created for all pair-
927 wise comparisons of individuals, where the distance between two individuals is
928 defined as the sum of all differences over all segregating sites (e.g., a heterozygote-
929 homozygote difference at a SNP adds 1 to the distance; a derived homozygote-
930 ancestral homozygote difference adds 2). Phylogenetic trees were created from
931 these matrices by neighbour-joining, using the ‘nj’ function in the ‘ape’ package
932 for R [94,95].

933 **Supporting information**

934 **S1 File. Supplementary *Mathematica* File.** *Mathematica* notebook of al-
935 gebraic derivations and simulation comparisons (.nb format).

936 **S2 File. Supplementary *Mathematica* File (PDF).** *Mathematica* notebook
937 of algebraic derivations and simulation comparisons (.pdf format).

938 **S3 File. Supplementary Text File.** Additional results and figures pertain-
939 ing to effective reduction in diversity under different scenarios; deriving the site-
940 frequency spectrum; further results on singleton distributions; and simulation re-
941 sults of *SLC24A5* sweep region.

942 **S4 File. Simulation Code.** Forward-in-time simulation code written in C.
943 Also available from <https://github.com/MattHartfield/DomSelfAdapt>.

944 **S5 Table. Simulation Command Lines.** List of MSMS command lines used
945 to simulate a sweep at the *SLC24A5* region under different scenarios.

946 **Acknowledgments**

947 We would like to thank Sally Otto for providing information on the elevated effec-
948 tive starting frequency of beneficial mutations; Dan Schrider for providing informa-
949 tion on how to simulate demographic scenarios; and Mikkel Schierup for feedback
950 on the human sweep analyses.

References

1. Maynard Smith J, Haigh J. The hitch-hiking effect of a favourable gene. *Genet Res.* 1974;23:23–35.
2. Kaplan NL, Hudson RR, Langley CH. The “Hitchhiking Effect” Revisited. *Genetics.* 1989;123(4):887–899.
3. Thomson G. The effect of a selected locus on linked neutral loci. *Genetics.* 1977;85(4):753–788.
4. Innan H, Nordborg M. The Extent of Linkage Disequilibrium and Haplotype Sharing Around a Polymorphic Site. *Genetics.* 2003;165(1):437.
5. McVean GAT. The Structure of Linkage Disequilibrium Around a Selective Sweep. *Genetics.* 2007;175(3):1395–1406.
6. Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, et al. Detecting recent positive selection in the human genome from haplotype structure. *Nature.* 2002;419(6909):832–837.
7. Kim Y, Nielsen R. Linkage Disequilibrium as a Signature of Selective Sweeps. *Genetics.* 2004;167(3):1513–1524.
8. Voight BF, Kudaravalli S, Wen X, Pritchard JK. A Map of Recent Positive Selection in the Human Genome. *PLoS Biol.* 2006;4(3):e72.
9. Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. On Detecting Incomplete Soft or Hard Selective Sweeps Using Haplotype Structure. *Mol Biol Evol.* 2014;31(5):1275–1291.

10. Vatsiou AI, Bazin E, Gaggiotti OE. Detection of selective sweeps in structured populations: a comparison of recent methods. *Mol Ecol.* 2016;25(1):89–103.
11. Hermisson J, Pennings PS. Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation. *Methods Ecol Evol.* 2017;8(6):700–716.
12. Barrett RDH, Schluter D. Adaptation from standing genetic variation. *Trends Ecol Evol.* 2008;23(1):38–44.
13. Messer PW, Petrov DA. Population genomics of rapid adaptation by soft selective sweeps. *Trends Ecol Evol.* 2013;28(11):659–669.
14. Pennings PS, Hermisson J. Soft Sweeps II – Molecular Population Genetics of Adaptation from Recurrent Mutation or Migration. *Mol Biol Evol.* 2006;23(5):1076–1084.
15. Pennings PS, Hermisson J. Soft Sweeps III: The Signature of Positive Selection from Recurrent Mutation. *PLoS Genet.* 2006;2(12):e186.
16. Orr HA, Betancourt AJ. Haldane’s Sieve and Adaptation From the Standing Genetic Variation. *Genetics.* 2001;157(2):875–884.
17. Innan H, Kim Y. Pattern of polymorphism after strong artificial selection in a domestication event. *Proc Natl Acad Sci USA.* 2004;101(29):10667–10672.
18. Przeworski M, Coop G, Wall JD. The Signature of Positive Selection on Standing Genetic Variation. *Evolution.* 2005;59(11):2312–2323.

19. Hermisson J, Pennings PS. Soft Sweeps: Molecular Population Genetics of Adaptation From Standing Genetic Variation. *Genetics*. 2005;169(4):2335–2352.
20. Wilson BA, Petrov DA, Messer PW. Soft Selective Sweeps in Complex Demographic Scenarios. *Genetics*. 2014;198(2):669–684.
21. Berg JJ, Coop G. A Coalescent Model for a Sweep of a Unique Standing Variant. *Genetics*. 2015;201(2):707–725.
22. Wilson BA, Pennings PS, Petrov DA. Soft Selective Sweeps in Evolutionary Rescue. *Genetics*. 2017;205(4):1573–1586.
23. Peter BM, Huerta-Sanchez E, Nielsen R. Distinguishing between Selective Sweeps from Standing Variation and from a *De Novo* Mutation. *PLoS Genet*. 2012;8(10):e1003011.
24. Vitti JJ, Grossman SR, Sabeti PC. Detecting Natural Selection in Genomic Data. *Annu Rev Genet*. 2013;47(1):97–120.
25. Garud NR, Messer PW, Buzbas EO, Petrov DA. Recent Selective Sweeps in North American *Drosophila melanogaster* Show Signatures of Soft Sweeps. *PLoS Genet*. 2015;11(2):e1005004.
26. Garud NR, Petrov DA. Elevated Linkage Disequilibrium and Signatures of Soft Sweeps Are Common in *Drosophila melanogaster*. *Genetics*. 2016;203(2):863–880.
27. Schrider DR, Kern AD. S/HIC: Robust Identification of Soft and Hard Sweeps Using Machine Learning. *PLoS Genet*. 2016;12(3):e1005928.

28. Sheehan S, Song YS. Deep Learning for Population Genetic Inference. *PLoS Comput Biol.* 2016;12(3):e1004845.
29. Schrider DR, Kern AD. Soft Sweeps Are the Dominant Mode of Adaptation in the Human Genome. *Mol Biol Evol.* 2017;34(8):1863–1877.
30. Fustier MA, Brandenburg JT, Boitard S, Lapeyronnie J, Eguiarte LE, Vigouroux Y, et al. Signatures of local adaptation in lowland and highland teosintes from whole-genome sequencing of pooled samples. *Mol Ecol.* 2017;26(10):2738–2756.
31. Anderson TJC, Nair S, McDew-White M, Cheeseman IH, Nkhoma S, Bilgic F, et al. Population Parameters Underlying an Ongoing Soft Sweep in Southeast Asian Malaria Parasites. *Mol Biol Evol.* 2016;34(1):131–144.
32. Jensen JD. On the unfounded enthusiasm for soft selective sweeps. *Nat Commun.* 2014;5.
33. Teshima KM, Przeworski M. Directional Positive Selection on an Allele of Arbitrary Dominance. *Genetics.* 2006;172(1):713–718.
34. Teshima KM, Coop G, Przeworski M. How reliable are empirical genomic scans for selective sweeps? *Genome Res.* 2006;16(6):702–712.
35. Ewing G, Hermisson J, Pfaffelhuber P, Rudolf J. Selective sweeps for recessive alleles and for other modes of dominance. *J Math Bio.* 2011;63(3):399–431.
36. Hartfield M, Bataillon T, Glémin S. The Evolutionary Interplay between Adaptation and Self-Fertilization. *Trends Genet.* 2017;33(6):420–431.

37. Igic B, Kohn JR. The distribution of plant mating systems: study bias against obligately outcrossing species. *Evolution*. 2006;60(5):1098–1103.
38. Jarne P, Auld JR. Animals mix it up too: the distribution of self-fertilization among hermaphroditic animals. *Evolution*. 2006;60(9):1816–1824.
39. Billiard S, López-Villavicencio M, Devier B, Hood ME, Fairhead C, Giraud T. Having sex, yes, but with whom? Inferences from fungi on the evolution of anisogamy and mating types. *Biol Rev Camb Philos Soc*. 2011;86(2):421–442.
40. Haldane JBS. A Mathematical Theory of Natural and Artificial Selection, Part V: Selection and Mutation. *Math Proc Cambridge Philos Soc*. 1927;23(7):838–844.
41. Charlesworth B. Evolutionary Rates in Partially Self-Fertilizing Species. *Am Nat*. 1992;140(1):126–148.
42. Glémin S. Extinction and fixation times with dominance and inbreeding. *Theor Popul Biol*. 2012;81(4):310–316.
43. Nordborg M. Linkage Disequilibrium, Gene Trees and Selfing: An Ancestral Recombination Graph With Partial Self-Fertilization. *Genetics*. 2000;154(2):923–929.
44. Hartfield M, Glémin S. Hitchhiking of Deleterious Alleles and the Cost of Adaptation in Partially Selfing Species. *Genetics*. 2014;196(1):281–293.
45. Hartfield M, Glémin S. Limits to Adaptation in Partially Selfing Species. *Genetics*. 2016;203(2):959–974.

46. Glémin S, Ronfort J. Adaptation and Maladaptation in Selfing and Outcrossing Species: New Mutations Versus Standing Variation. *Evolution*. 2013;67(1):225–240.
47. Uecker H. Evolutionary rescue in randomly mating, selfing, and clonal populations. *Evolution*. 2017;71(4):845–858.
48. Long Q, Rabanal FA, Meng D, Huber CD, Farlow A, Platzer A, et al. Massive genomic variation and strong selection in *Arabidopsis thaliana* lines from Sweden. *Nat Genet*. 2013;45(8):884–890.
49. Huber CD, Nordborg M, Hermisson J, Hellmann I. Keeping It Local: Evidence for Positive Selection in Swedish *Arabidopsis thaliana*. *Mol Biol Evol*. 2014;31(11):3026–3039.
50. Fulgione A, Koornneef M, Roux F, Hermisson J, Hancock AM. Madeiran *Arabidopsis thaliana* Reveals Ancient Long-Range Colonization and Clarifies Demography in Eurasia. *Mol Biol Evol*. 2018;35(3):564–574.
51. Andersen EC, Gerke JP, Shapiro JA, Crissman JR, Ghosh R, Bloom JS, et al. Chromosome-scale selective sweeps shape *Caenorhabditis elegans* genomic diversity. *Nat Genet*. 2012;44(3):285–290.
52. Bonhomme M, Boitard S, San Clemente H, Dumas B, Young N, Jacquet C. Genomic Signature of Selective Sweeps Illuminates Adaptation of *Medicago truncatula* to Root-Associated Microorganisms. *Mol Biol Evol*. 2015;32(8):2097–2110.

53. Badouin H, Gladieux P, Gouzy J, Siguenza S, Aguilera G, Snirc A, et al. Widespread selective sweeps throughout the genome of model plant pathogenic fungi and identification of effector candidates. *Mol Ecol*. 2017;26(7):2041–2062.
54. Hedrick PW. Hitchhiking: A Comparison of Linkage and Partial Selection. *Genetics*. 1980;94(3):791–808.
55. Schoen DJ, Morgan MT, Bataillon T. How Does Self-Pollination Evolve? Inferences from Floral Ecology and Molecular Genetic Variation. *Philos Trans R Soc Lond B Biol Sci*. 1996;351(1345):1281–1290.
56. Field Y, Boyle EA, Telis N, Gao Z, Gaulton KJ, Golan D, et al. Detection of human adaptation during the past 2000 years. *Science*. 2016;354(6313):760–764.
57. Wright S. The genetical structure of populations. *Ann Eugen*. 1951;15:323–354.
58. Caballero A, Hill WG. Effects of Partial Inbreeding on Fixation Rates and Variation of Mutant Genes. *Genetics*. 1992;131(2):493–507.
59. Nordborg M, Donnelly P. The Coalescent Process With Selfing. *Genetics*. 1997;146(3):1185–1195.
60. Roze D. Diploidy, Population Structure, and the Evolution of Recombination. *Am Nat*. 2009;174(S1):S79–S94.
61. Roze D. Background Selection in Partially Selfing Populations. *Genetics*. 2016;203(2):937–957.

62. Barton NH. Genetic Hitchhiking. *Philos Trans R Soc Lond B Biol Sci.* 2000;355:1553–1562.
63. Stephan W, Wiehe THE, Lenz MW. The effect of strongly selected substitutions on neutral polymorphism: analytical results based on diffusion theory. *Theor Popul Biol.* 1992;41:237–254.
64. Wakeley J. *Coalescent theory: an introduction.* vol. 1. Greenwood Village, Colorado: Roberts & Company Publishers; 2009.
65. Barton NH. The effect of hitch-hiking on neutral genealogies. *Genet Res.* 1998;72:123–133.
66. Desai MM, Fisher DS. Beneficial Mutation-Selection Balance and the Effect of Linkage on Positive Selection. *Genetics.* 2007;176(3):1759–1798.
67. Martin G, Lambert A. A simple, semi-deterministic approximation to the distribution of selective sweeps in large populations. *Theor Popul Biol.* 2015;101:40–46.
68. Abramowitz M, Stegun IA. *Handbook of Mathematical Functions.* New York: Dover Publications, Inc.; 1970.
69. Watterson GA. On the number of segregating sites in genetical models without recombination. *Theor Popul Biol.* 1975;7(2):256–276.
70. Hudson RR. Gene Genealogies and the Coalescent Process. In: Futuyma DJ, Antonovics J, editors. *Oxford Surveys in Evolutionary Biology.* vol. 7. Oxford Univ. Press, Oxford; 1990. p. 1–42.

71. Braverman JM, Hudson RR, Kaplan NL, Langley CH, Stephan W. The hitchhiking effect on the site frequency spectrum of DNA polymorphisms. *Genetics*. 1995;140(2):783–796.
72. Kim Y, Stephan W. Detecting a Local Signature of Genetic Hitchhiking Along a Recombining Chromosome. *Genetics*. 2002;160(2):765–777.
73. Lamason RL, Mohideen MAPK, Mest JR, Wong AC, Norton HL, Aros MC, et al. SLC24A5, a Putative Cation Exchanger, Affects Pigmentation in Zebrafish and Humans. *Science*. 2005;310(5755):1782–1786.
74. Crawford NG, Kelly DE, Hansen MEB, Beltrame MH, Fan S, Bowman SL, et al. Loci associated with skin pigmentation identified in African populations. *Science*. 2017;358(6365).
75. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. Inferring the Joint Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. *PLoS Genet*. 2009;5(10):e1000695.
76. Pritchard JK, Di Rienzo A. Adaptation - not by sweeps alone. *Nat Rev Genet*. 2010;11(10):665–667.
77. Stephan W. Signatures of positive selection: from selective sweeps at individual loci to subtle allele frequency changes in polygenic adaptation. *Mol Ecol*. 2016;25(1):79–88.
78. Elyashiv E, Sattath S, Hu TT, Strutsosky A, McVicker G, Andolfatto P, et al. A Genomic Map of the Effects of Linked Selection in *Drosophila*. *PLoS Genet*. 2016;12(8):e1006130.

79. Kim Y, Stephan W. Selective Sweeps in the Presence of Interference Among Partially Linked Loci. *Genetics*. 2003;164(1):389–398.
80. Chevin LM, Billiard S, Hospital F. Hitchhiking Both Ways: Effect of Two Interfering Selective Sweeps on Linked Neutral Variation. *Genetics*. 2008;180(1):301–316.
81. Sattath S, Elyashiv E, Kolodny O, Rinott Y, Sella G. Pervasive Adaptive Protein Evolution Apparent in Diversity Patterns around Amino Acid Substitutions in *Drosophila simulans*. *PLoS Genet*. 2011;7(2):e1001302.
82. Halligan DL, Kousathanas A, Ness RW, Harr B, Eöry L, Keane TM, et al. Contributions of Protein-Coding and Regulatory Change to Adaptive Molecular Evolution in Murid Rodents. *PLoS Genet*. 2013;9(12):e1003995.
83. Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, Bustamante C. Genomic scans for selective sweeps using SNP data. *Genome Res*. 2005;15(11):1566–1575.
84. Boitard S, Schlötterer C, Futschik A. Detecting Selective Sweeps: A New Approach Based on Hidden Markov Models. *Genetics*. 2009;181(4):1567–1578.
85. Schrider DR, Mendes FK, Hahn MW, Kern AD. Soft Shoulders Ahead: Spurious Signatures of Soft and Partial Selective Sweeps Result from Linked Hard Sweeps. *Genetics*. 2015;200(1):267–284.

86. Spencer CCA, Coop G. SelSim: a program to simulate population genetic data with natural selection and recombination. *Bioinformatics*. 2004;20(18):3673–3675.
87. Ewing G, Hermisson J. MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics*. 2010;26(16):2064–2065.
88. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*. 2015;526(7571):68–74.
89. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. *Bioinformatics*. 2011;27(15):2156–2158.
90. Bhérer C, Campbell CL, Auton A. Refined genetic maps reveal sexual dimorphism in human meiotic recombination at multiple scales. *Nat Commun*. 2017;8:14994.
91. Jorde LB, Bamshad M, Rogers AR. Using mitochondrial and nuclear DNA markers to reconstruct human evolution. *BioEssays*. 1998;20(2):126–136.
92. Wolfram Research, Inc . *Mathematica Edition: Version 11.2*. Champaign, Illinois: Wolfram Research, Inc.; 2017.
93. Schrider DR, Shanku AG, Kern AD. Effects of Linked Selective Sweeps on Demographic Inference and Model Selection. *Genetics*. 2016;204(3):1207–1223.

94. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics*. 2004;20(2):289–290.
95. R Development Core Team. R: A Language and Environment for Statistical Computing; 2014. Available from: <http://www.R-project.org>.

Figures

Fig 1. A schematic of the model. The history of the derived variant is separated into two phases. The ‘standing phase’ (shown in light gray), is when the derived variant is segregating at a frequency f_0 for a long period of time. The ‘sweep phase’ (shown in dark gray) is when the variant becomes selected for and starts increasing in frequency. Dots on the right-hand side represent samples of the derived haplotype taken in the present day, with lines representing their genetic histories. Samples can recombine onto the ancestral background either during the sweep phase or the standing phase. Solid lines represent coalescent histories on the derived genetic background; dotted lines represent coalescent histories on the ancestral background.

Fig 2. Expected pairwise diversity following a selective sweep. Plots of $\mathbb{E}(\pi/\pi_0)$ as a function of the recombination rate scaled to population size $2Nr$. Lines are analytical solutions (Eq 6), points are simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$ (note μ is scaled by N in simulations, not N_e), and dominance coefficient $h = 0.1$ (red lines, points), 0.5 (black lines, points), or 0.9 (blue lines, points). Rate of self-fertilisation equals 0 ; 0.5 ; or 0.95 (note the x -axis range changes with the self-fertilisation rate). The sweep arose from either a single *de novo* mutation (actual $f_0 = 1/2N$; note we use $f_{0,A}$ in our model, as given by Eq 5), standing variation with $f_0 = 0.02$; or $f_0 = 0.05$. Further results are plotted in Section B of S1 File.

Fig 3. Beneficial allele trajectories. These were obtained by numerically evaluating the negative of Eq 2 forward in time. $N = 5,000$, $s = 0.05$, and h equals either 0.1 (red lines), 0.5 (black lines), or 0.9 (blue lines). Values of f_0 and self-fertilisation rates used are shown at the end of the relevant row and column. Note the different x -axis scales used in each panel. Further results are plotted in Section B of S1 File.

Fig 4. Effective reductions in diversity under different scenarios. (a) \tilde{s} (Eq 7) as scaled to s , as a function of $R = 2Nr$. $f_0 = 1/2N$ (black line), 0.02 (red line) or 0.1 (blue line). (b) Plot of π/π_0 (Eq 6) using \tilde{s} . Solid lines represent $f_0 = 1/2N$ (black line), 0.02 (red line) or 0.1 (blue line). Points are Eq 6 assuming $f_0 = 1/2N$, but using \tilde{s} (Eq 7) evaluated for $f_0 = 0.02$ (circles) or 0.1 (squares). The population is outcrossing; similar results exist for partial selfing ($\sigma = 0.5$) if measuring over a longer recombination distance (Section A of S2 File). Other parameters are $N = 5,000$, $s = 0.05$, $h = 0.5$.

Fig 5. Expected number of segregating sites following a selective sweep. A plot of $\mathbb{E}(S)$, as a function of the recombination rate scaled to population size $2Nr$. Lines are analytical solutions (Eq 10 multiplied by θ), points are simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$ (so $\theta = 4N_e\mu$ per bin is 4 for $\sigma = 0$, 3 for $\sigma = 0.5$, and 2.1 for $\sigma = 0.95$), and dominance coefficient $h = 0.1$ (red lines, points), 0.5 (black lines, points), or 0.9 (blue lines, points). Rate of self-fertilisation σ equals 0, 0.5, or 0.95 as denoted on the right-hand side; note the different x -axes ranges. The sweep arose from either a single *de novo* mutation or standing variation with $f_0 = 0.05$, as denoted at the top of the figure. Further results are plotted in Section D of S1 File.

Fig 6. Expected site frequency spectrum, in flanking regions to the adaptive mutation, following a selective sweep. Lines are analytical solutions (Eq B14 in the Supplementary Material), points are simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$ (so the effective mutation rate per bin is 4 for $\sigma = 0$ and 3 for $\sigma = 0.5$), and dominance coefficient $h = 0.1$ (red lines, points), 0.5 (black lines, points), or 0.9 (blue lines, points). The neutral SFS is also included for comparisons (grey dashed line). Rate of self-fertilisation $\sigma = 0$ or 1/2, as denoted on the right-hand side. The SFS is measured at a recombination distance of $R = 6$ for $\sigma = 0$, or $R = 11$ for $\sigma = 0.5$. The sweep arose from either a single *de novo* mutation or standing variation with $f_0 = 0.05$, as denoted above the panels. Results for other recombination distances are in Section E of S1 File.

Fig 7. Comparing sweeps from recurrent mutation to those from standing variation. (a), (b): Comparing the reduction in diversity following a ‘soft’ sweep (Eq 6), from either standing variation ($f_0 = 0.05$, solid lines) or recurrent mutation (using $P_{coal,M}$ with $\Theta_b = 0.2$, dashed lines). $N = 5,000$, $s = 0.05$, and dominance coefficient $h = 0.1$ (red lines), 0.5 (black lines), or 0.9 (blue lines). Populations are either outcrossing (a) or highly selfing ($\sigma = 0.95$; (b)). (c), (d): Plotting the ratio of the diversity following a sweep from standing variation (π_{SV}) to one from recurrent mutation (π_M). Parameters for each panel are as in (a) and (b) respectively. Vertical dashed black line indicates R_{Lim} (Eq 11), the predicted recombination rate where $\pi_{SV}/\pi_M = 1$ (horizontal dashed line in (c), (d)). Note the different x -axis lengths between panels (a), (c) and (b), (d). Results are also plotted in Section F of S1 File.

Fig 8. Histograms of distances from the selected locus to the nearest singleton. The distance is scaled to the maximum length of the sampled genome (e.g., a distance of 0.5 means that a singleton lies halfway along the sampled haplotype). A distance “>1” indicates that no singleton was observed, and therefore lies beyond the sampled haplotype. x -axis annotations denote the mid-point of each bin (e.g. ‘0.05’ indicates distance of 0 to less than 0.1). Distances are measured from either the neutral burn-in population, or one where a ‘hard’ sweep ($f_0 = 1/2N$) has fixed. $N = 5,000$, $F = 0$, $4N\mu = 40$, $R = 2Nr = 4$ across the whole genetic sample. (a), (c) are log-counts of the distances for all samples over all 100 simulations; (b), (d) are the frequency of distances over samples where a singleton was observed. In (a), (b) $s = 0.05$ and the dominance coefficient h varies, with values as given in the plot legend. For (c), (d), $h = 0.5$ and s varies, with values as given in the plot legend.

Fig 9. Plots of distances from the selected locus to the nearest singleton, for a partial sweep. Distances are measured from either the neutral burn-in population (grey dashed lines), or one where a ‘hard’ sweep ($f_0 = 1/2N$) has reached a frequency of 70% (coloured lines). (a) Ratio of the log-counts of the distances for derived and ancestral alleles. (b) Ratio of the frequency of singleton distances for derived and ancestral alleles for each bin (e.g. ‘0.05’ indicates distance of 0 to less than 0.1). Measurements are taken over all samples in all simulations. All plots are log-counts of the distances for all samples over all simulations. $N = 5,000$, $F = 0$, $4N\mu = 40$, $R = 2Nr = 4$ across the whole genome. In sweep cases, $s = 0.05$ with the dominance coefficient $h = 0.1$ (red lines), 0.5 (black lines) or 0.9 (blue lines). Black dashed line indicates the 1-to-1 ratio, where the derived and ancestral classes have the same frequency.

Fig 10. Analysis of the *SLC24A5* sweep signature in humans. (a) Plot of diversity around the derived SNP in the *SLC24A5* gene, scaled to baseline values (see Methods for details), as a function of the distance from the target SNP as measured in basepairs. Negative values denote distance upstream of the target SNP; positive values denote downstream distances. Red dashed line denotes the ‘hard sweep’ model; blue dashed line is the recurrent mutation model. (b) Plot of two haplotype statistics, H_{12} (black line) and H_2/H_1 (red line) over the sweep region. (c) Unrooted phylogenetic trees of European and African samples from the 1000 Genomes dataset at different distances from the target SNP; covered distances are denoted in the headings. Arrows indicate instances where African haplotypes carrying the derived SNP (blue triangles) are present in the clade of European samples that cluster due to the sweep.

Fig 11. Schematic of how neutral polymorphisms accumulate in simulations. (a) The selected locus is located at the far left-hand side, with a neutral tract stretching out to its right. Polymorphisms accumulate along this tract, with locations standardised to be between 0 and 1. The recombination rate per reproduction is drawn from a Poisson distribution with mean r . ‘Singletons’ are polymorphisms where the derived allele is present in only one sample, with one present at location 0.65. (b) Measuring singleton distances using segregating target SNPs at a reference point. In the top sample the nearest singleton is located upstream of the target SNP, with distance 0.3 between them. In the bottom sample the singleton is located downstream of the SNP. Hence the total distance is given as the upstream distance to the right-hand edge (0.5), plus the distance of the singleton from the left-hand edge (0.1), giving a total distance of 0.6.

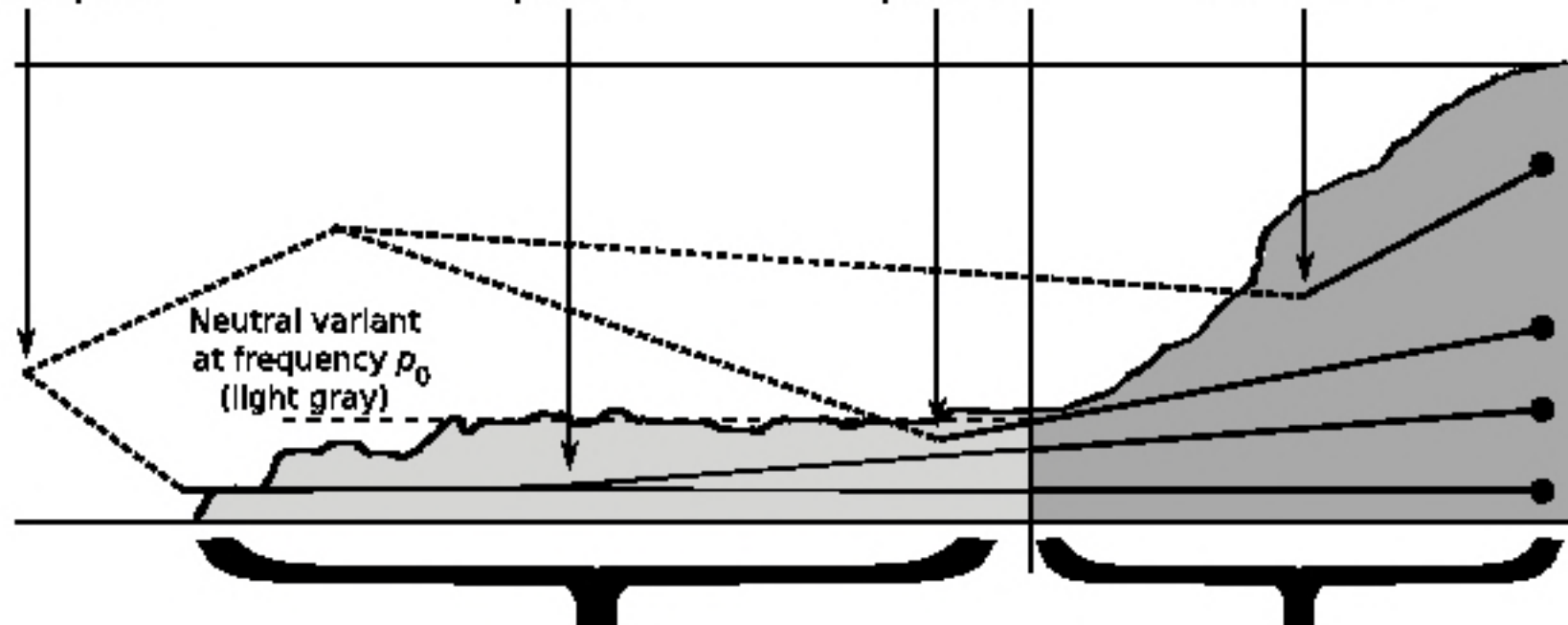
Fig 12. Diversity and recombination data around the *SLC24A5* sweep region. (a) Plot of raw pairwise diversity in 20Kb bins, as a function of distance from the target SNP. Dashed grey lines show mean diversity values when measured either upstream or downstream of the target SNP. (b) Relative diversity measurements, after dividing raw diversity measurements by the mean values from either up- or downstream of the target SNP. (c) Cumulative recombination distance from the target SNP, as obtained from Bhérier *et al.* [90], scaled by $2N = 20,000$.

Most recent
common ancestor
of all samples

Coalescence
during standing
phase

Recombination
during standing
phase

Recombination
during sweep



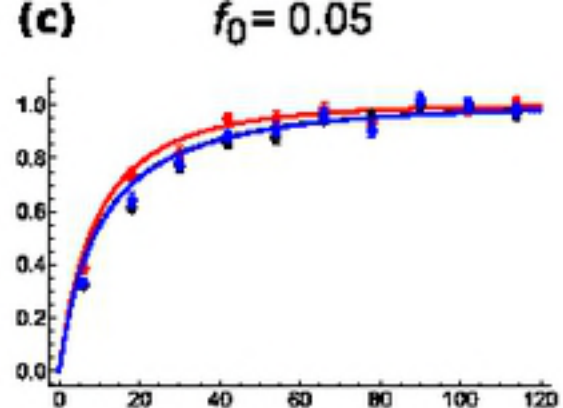
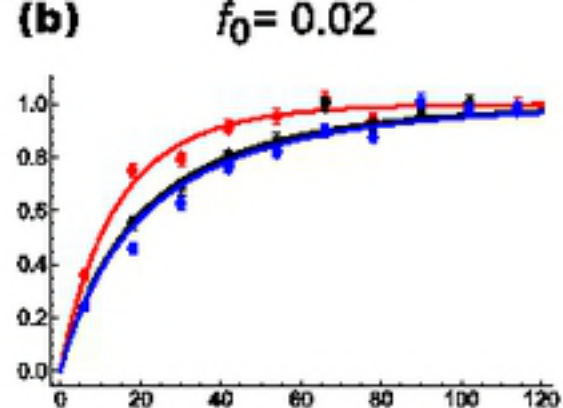
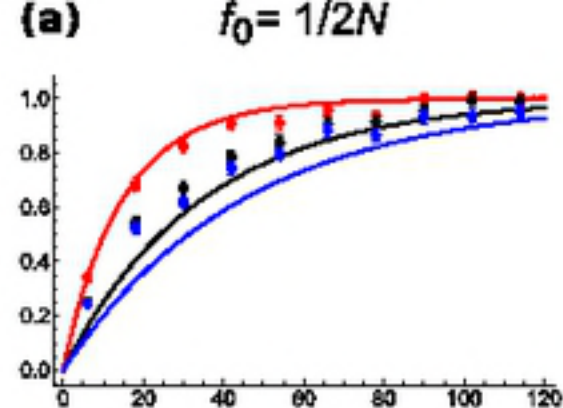
Neutral variant
at frequency p_0
(light gray)

Beneficial variant
at frequency p
(dark gray)

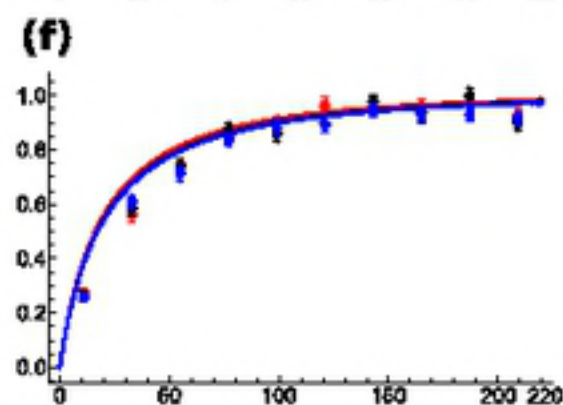
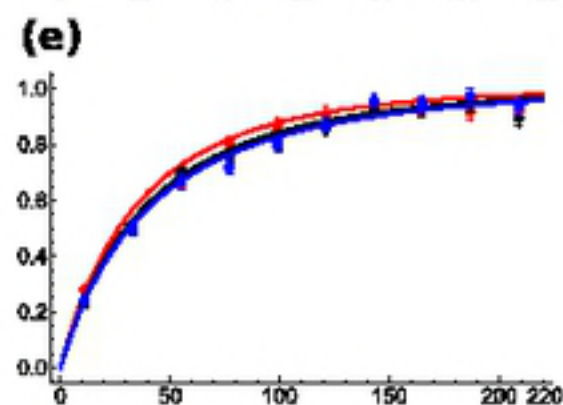
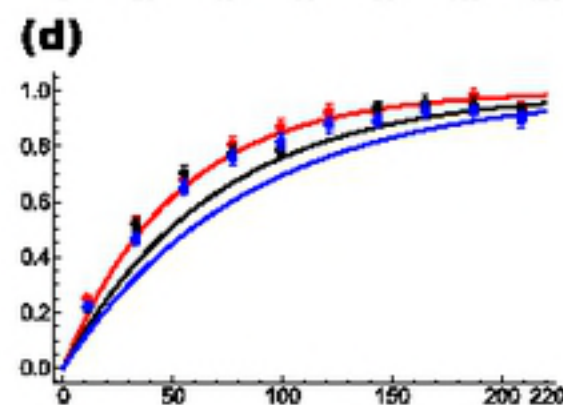
*Standing Phase:
Mutation is Neutral*

*Sweep Phase:
Mutation Selected For*

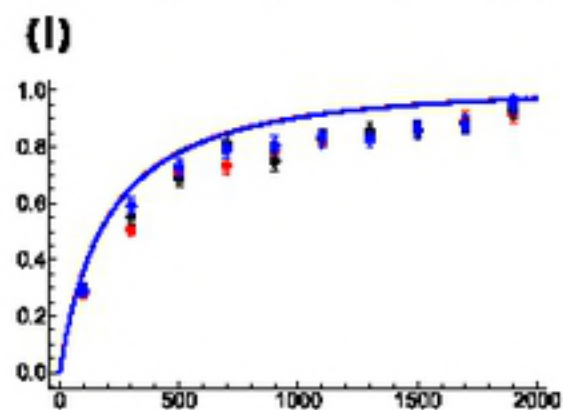
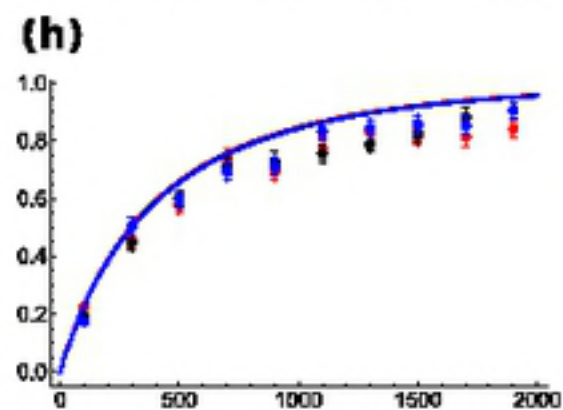
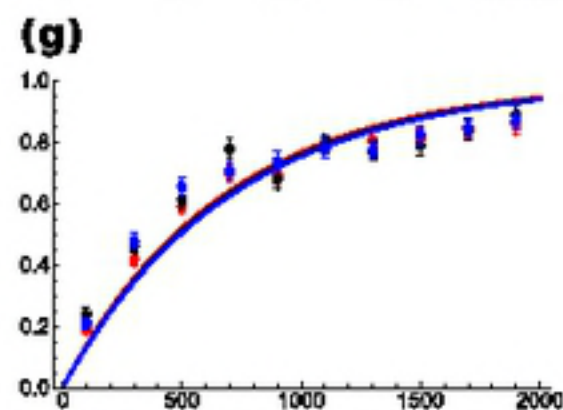
Expected
reduction
in diversity,
 $E(\pi/\pi_0)$



$\sigma = 0$
($F = 0$)



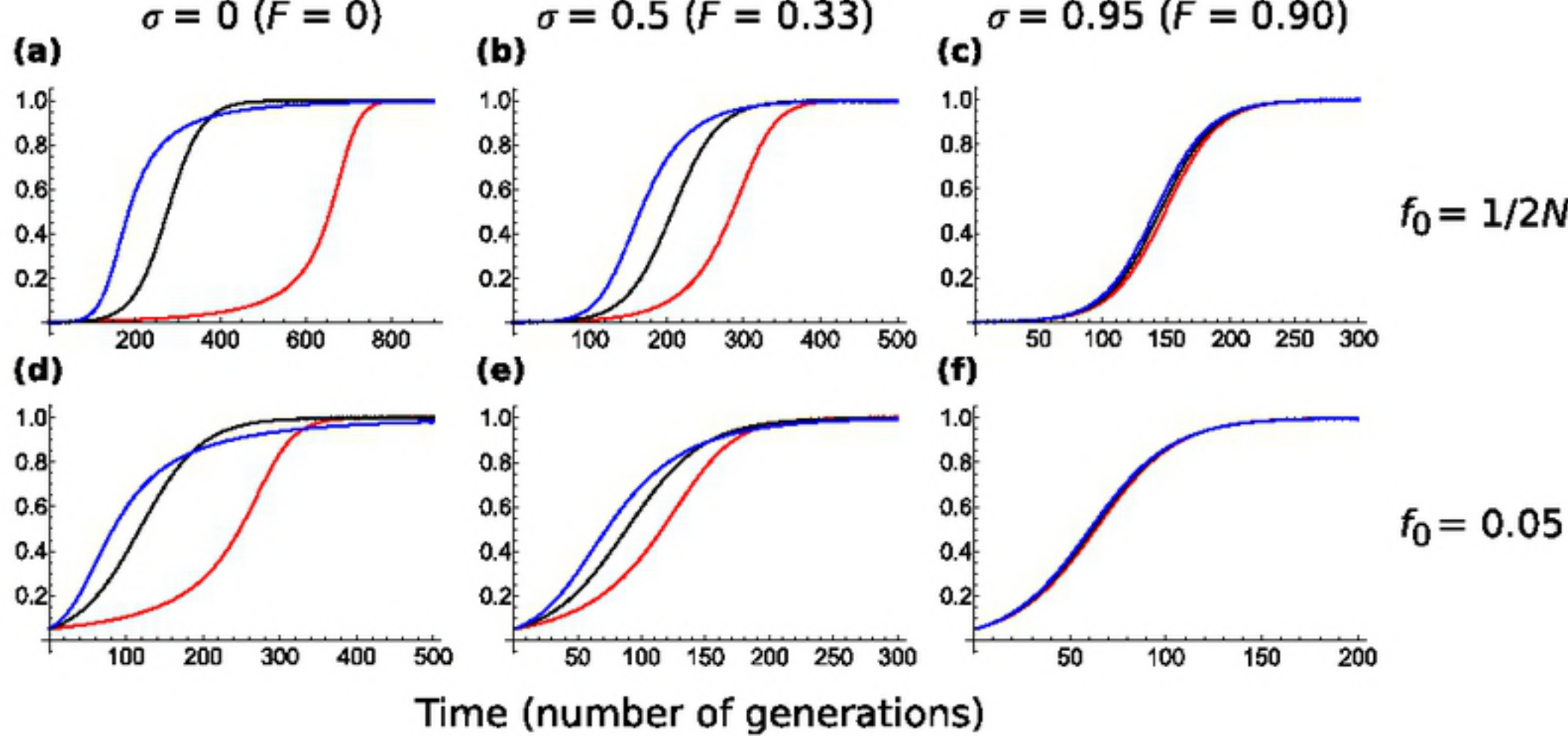
$\sigma = 0.5$
($F = 0.33$)

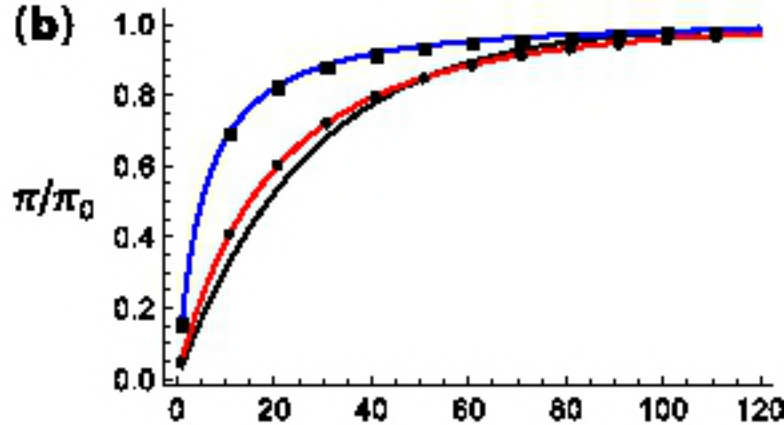
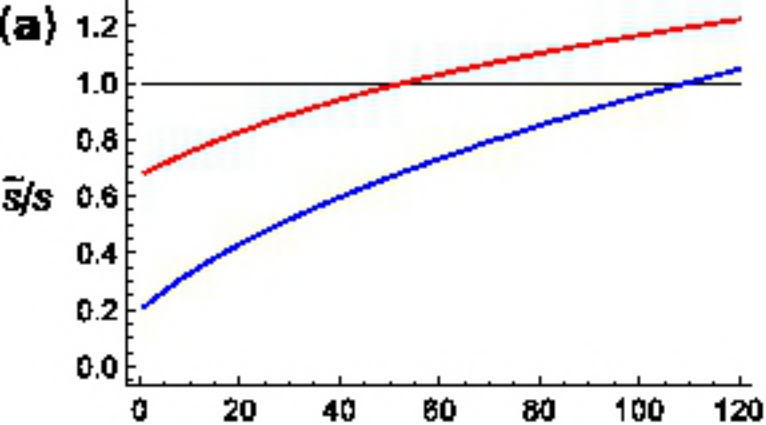


$\sigma = 0.95$
($F = 0.90$)

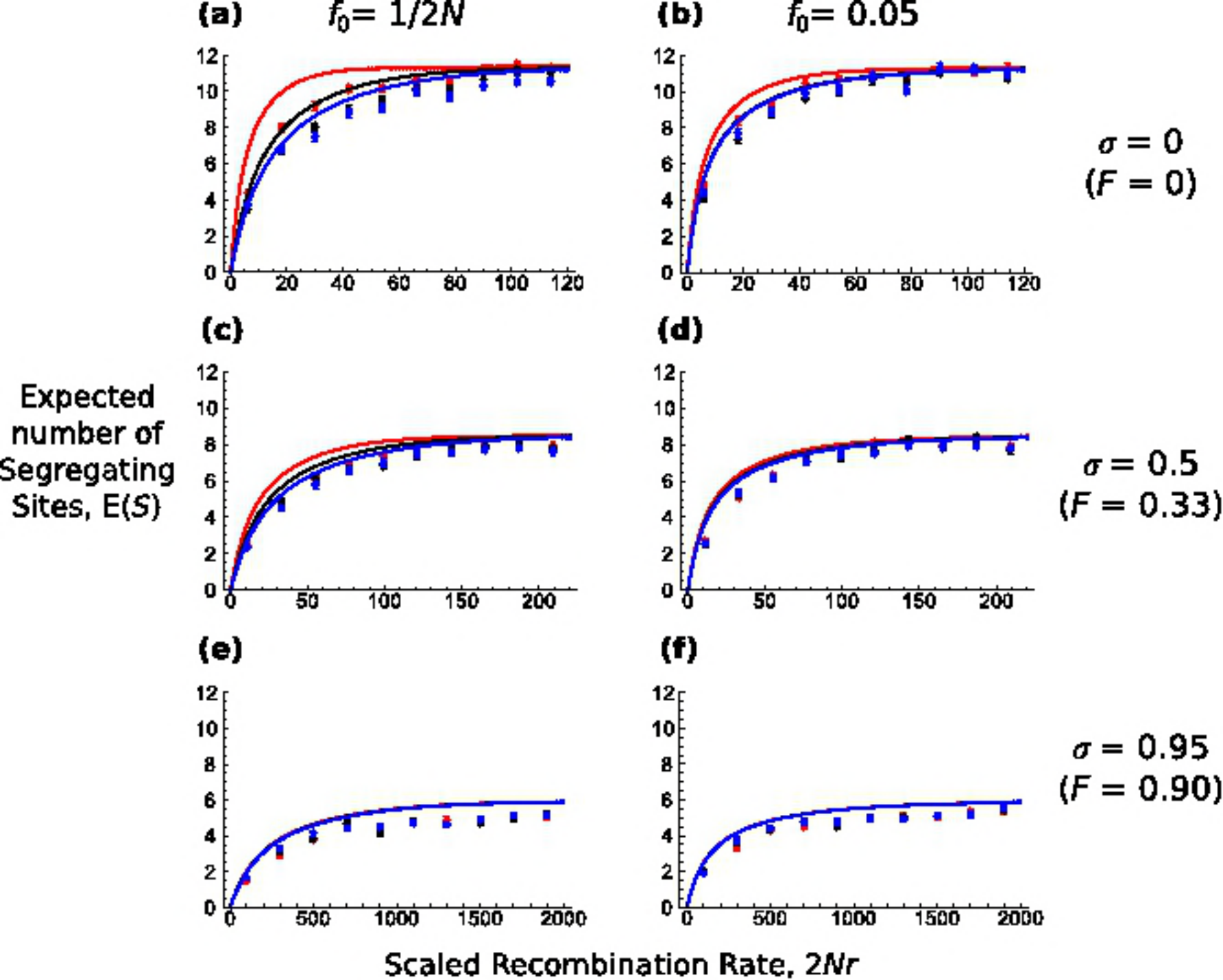
Scaled Recombination Rate, $2Nr$

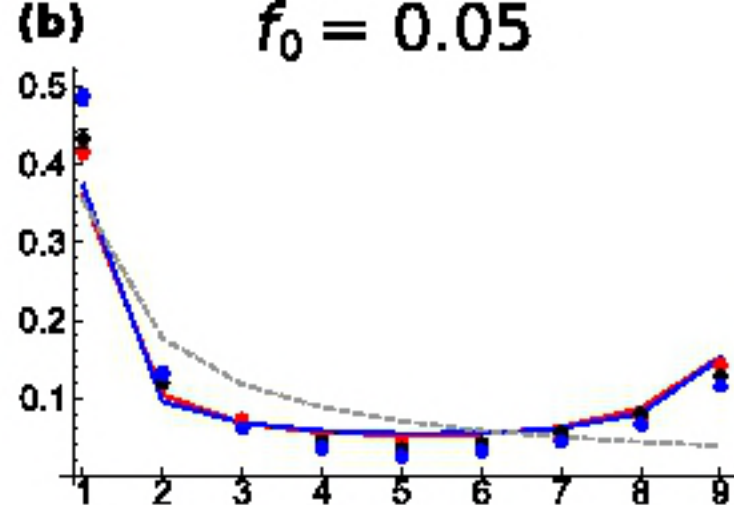
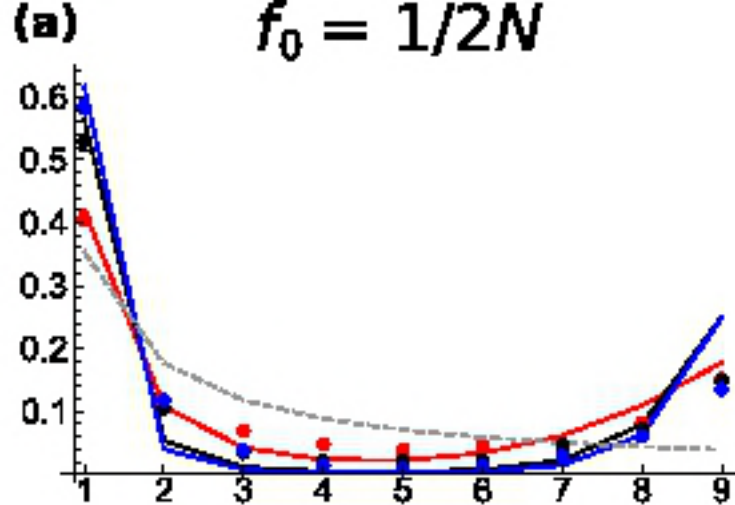
Beneficial
allele
frequency



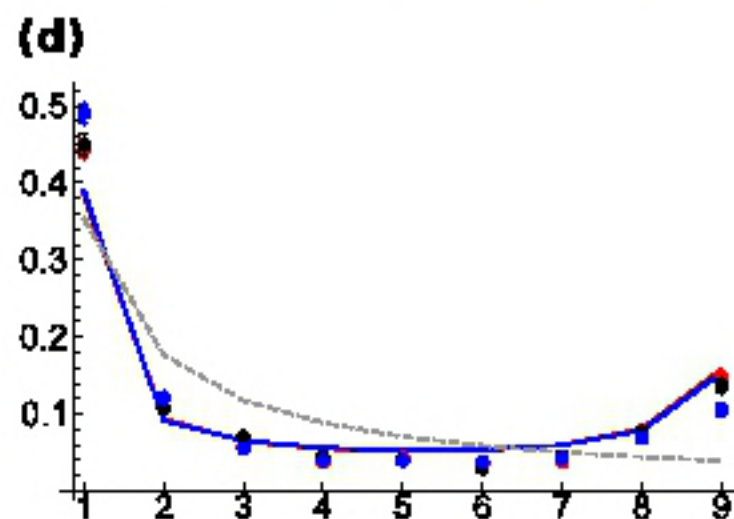
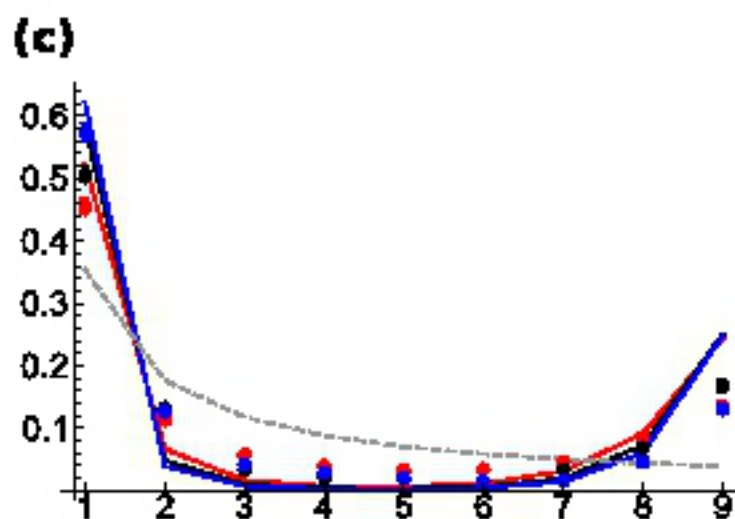


Scaled Recombination Rate, $2Nr$



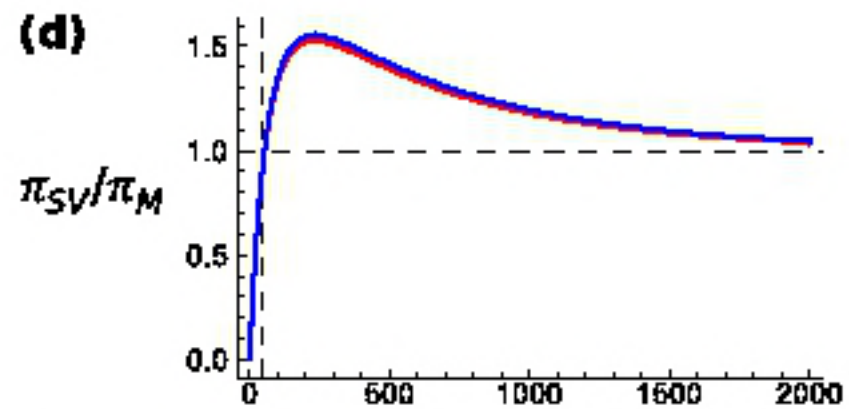
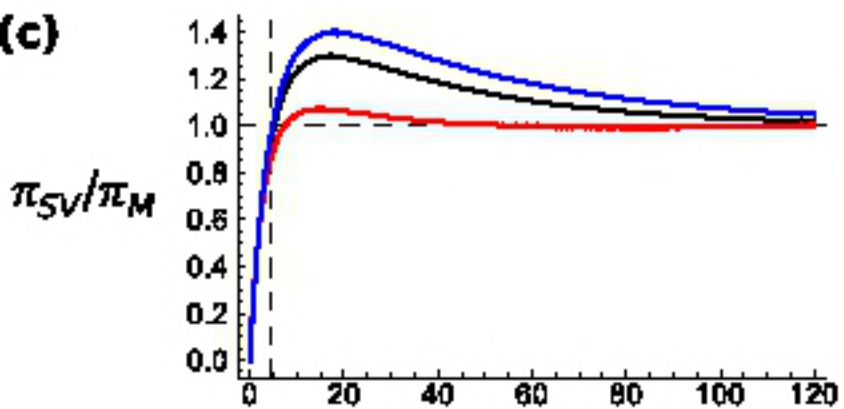
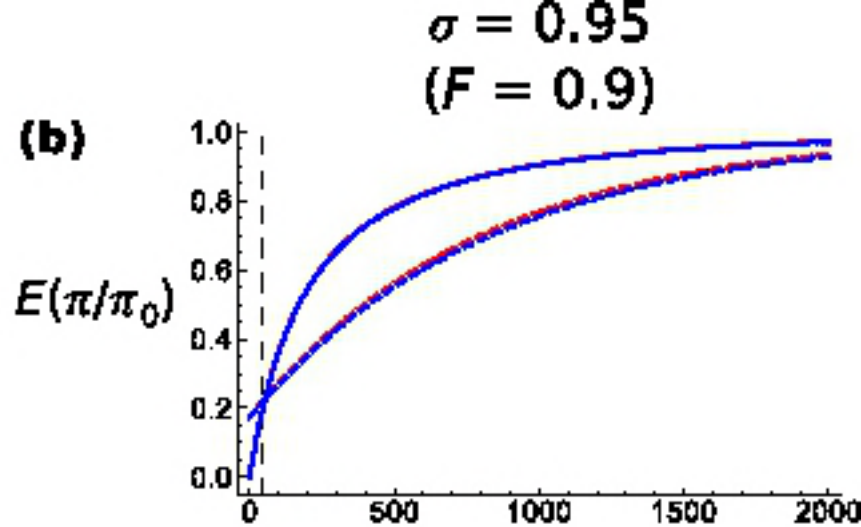
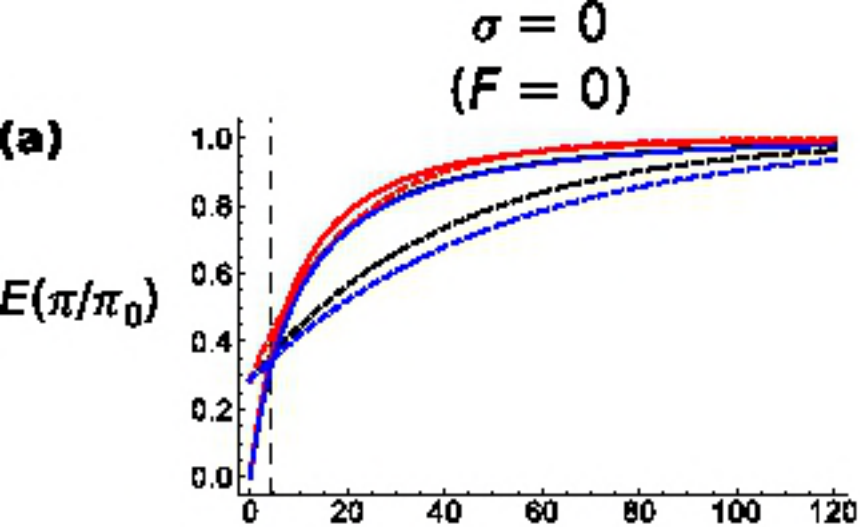


$\sigma = 0$
($F = 0$; $R = 6$)

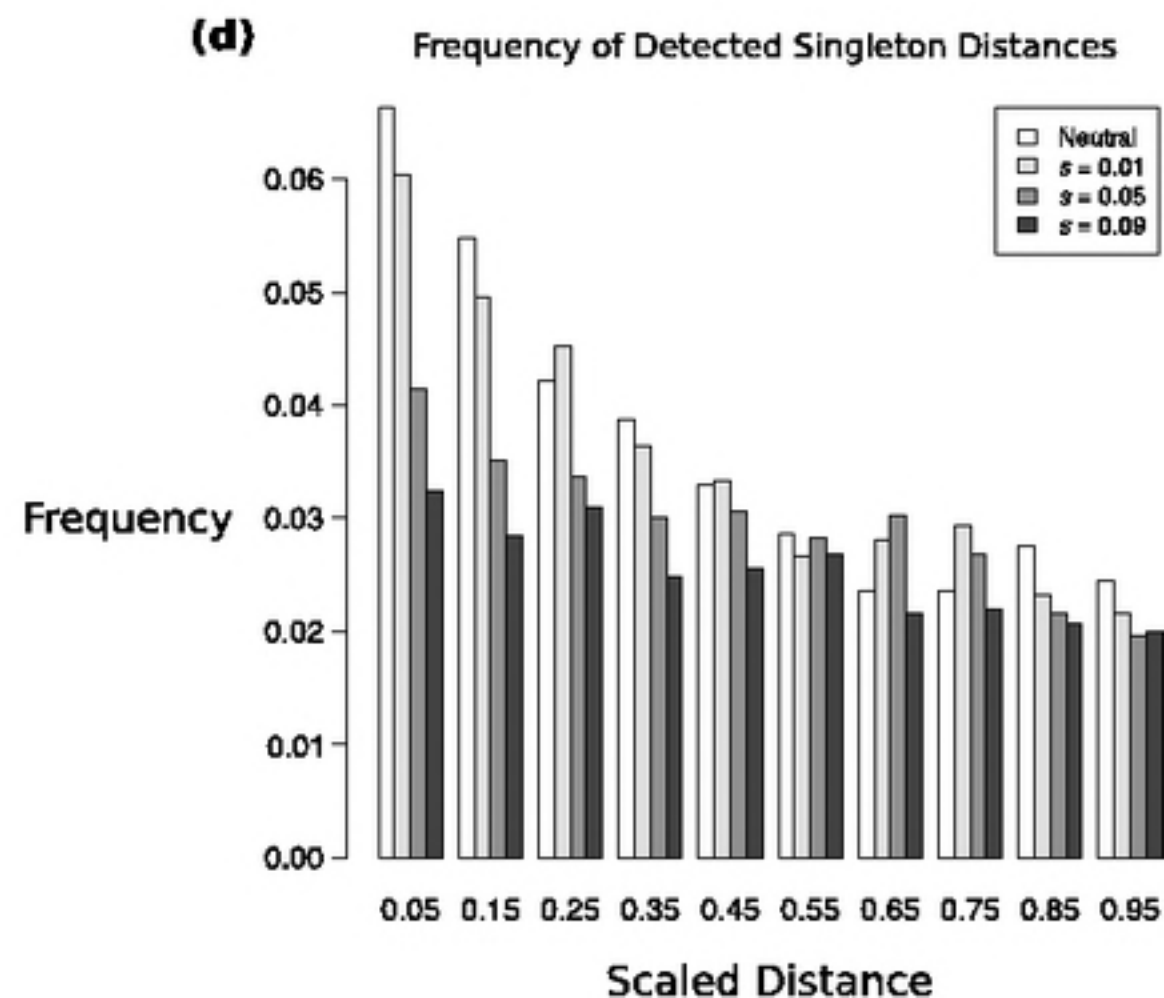
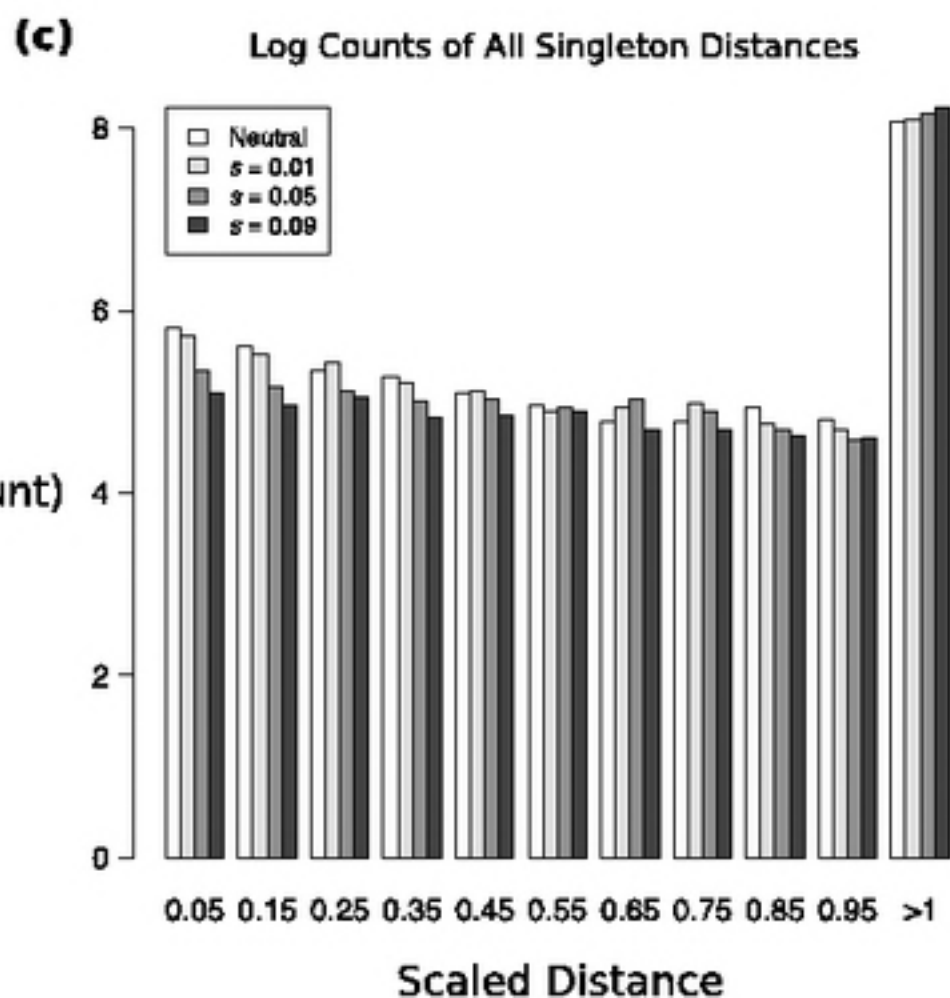
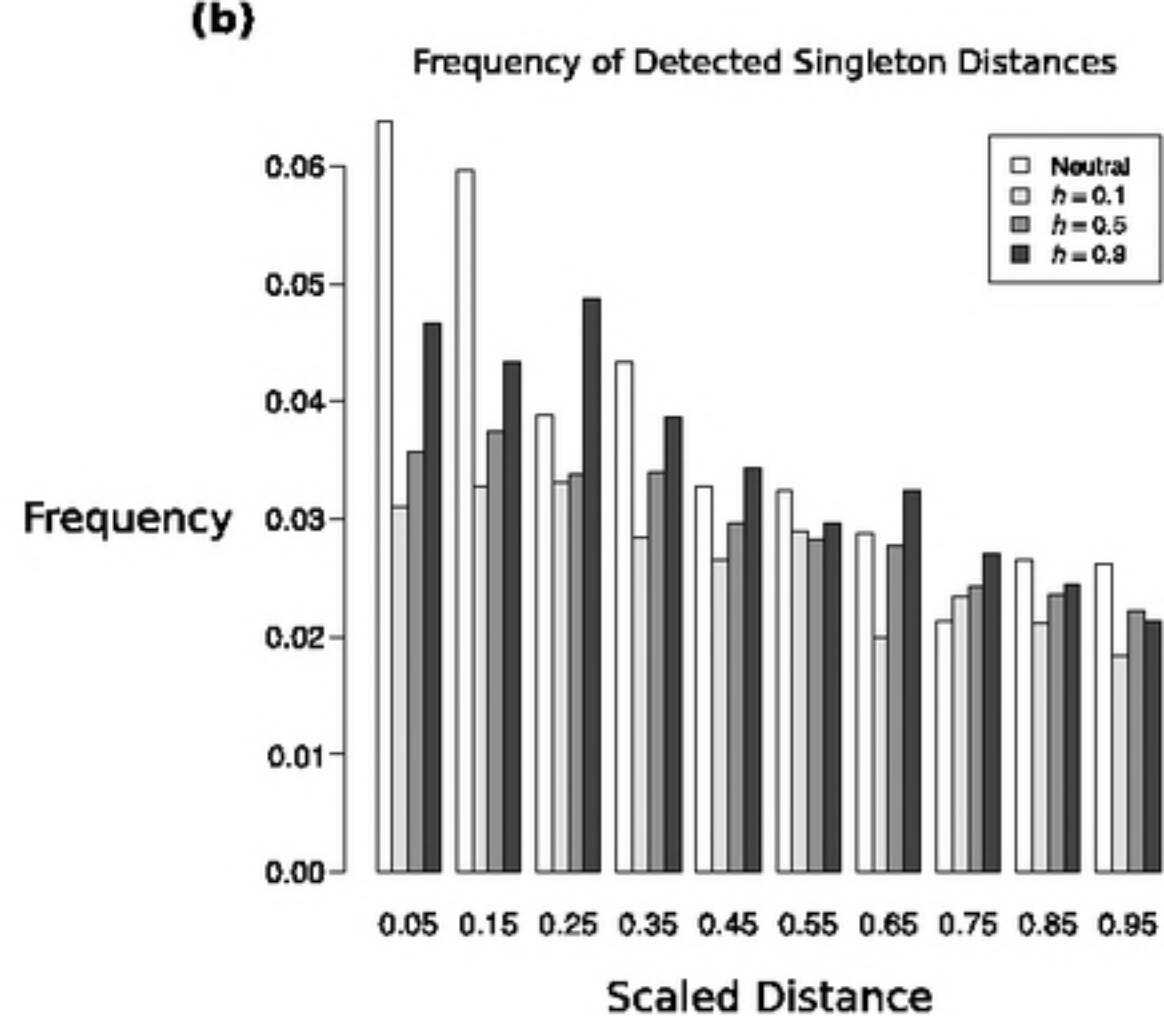
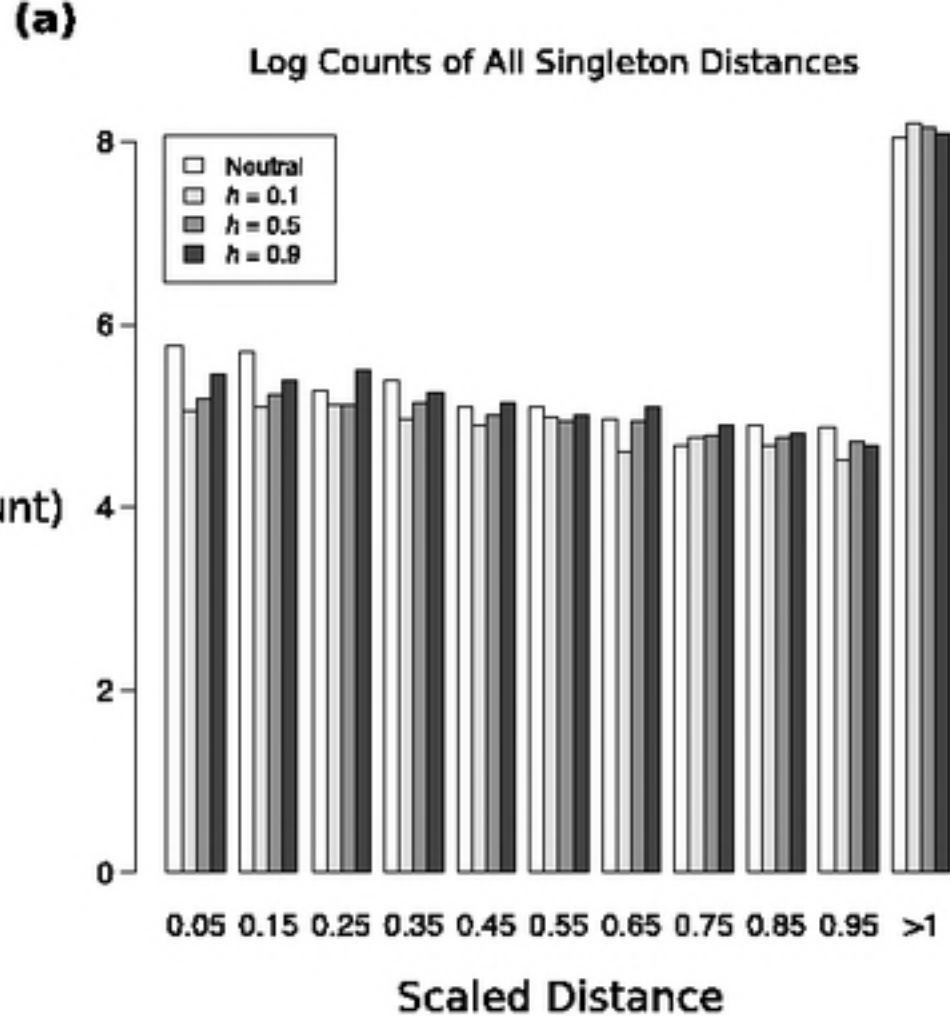


$\sigma = 0.5$
($F = 0.33$; $R = 11$)

Derived Allele Count



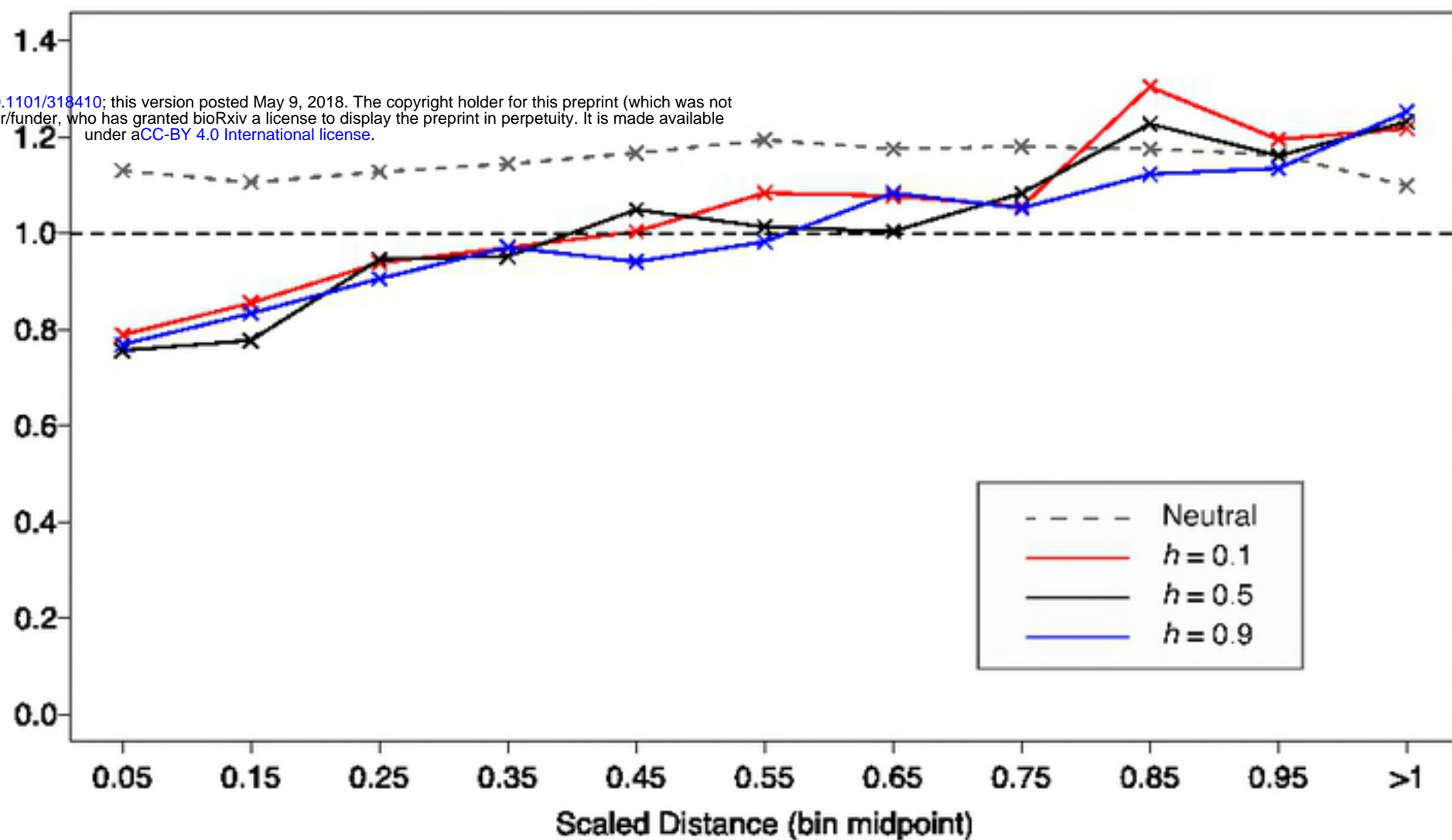
Scaled Recombination Rate, $2Nr$



(a)**All Samples**

bioRxiv preprint doi: <https://doi.org/10.1101/318410>; this version posted May 9, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

$\frac{\text{Log(Derived Counts)}}{\text{Log(Ancestral Counts)}}$

**(b)****Observed Samples**

Ratio of class frequencies
(Derived/Ancestral)

