

Selective sweeps under dominance and inbreeding

Matthew Hartfield^{1,2,3,*}, Thomas Bataillon²

1 Department of Ecology and Evolutionary Biology, University of Toronto, Ontario M5S 3B2, Canada.

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

3 Institute of Evolutionary Biology, The University of Edinburgh, Edinburgh EH9 3FL, United Kingdom.

* m.hartfield@ed.ac.uk

Running Head: Sweeps under dominance and inbreeding

Key words: Adaptation; Dominance; Self-fertilisation; Selective Sweeps; Population Genetics.

Abstract

1
2 A major research goal in evolutionary genetics is to uncover loci experi-
3 encing positive selection. One approach involves finding ‘selective sweeps’
4 patterns, which can either be ‘hard sweeps’ formed by *de novo* mutation, or
5 ‘soft sweeps’ arising from recurrent mutation or existing standing variation.
6 Existing theory generally assumes outcrossing populations, and it is unclear
7 how dominance affects soft sweeps. We consider how arbitrary dominance
8 and inbreeding via self-fertilisation affect hard and soft sweep signatures.
9 With increased self-fertilisation, they are maintained over longer map dis-
10 tances due to reduced effective recombination and faster beneficial allele
11 fixation times. Dominance can affect sweep patterns in outcrossers if the
12 derived variant originates from either a single novel allele, or from recurrent
13 mutation. These models highlight the challenges in distinguishing hard and
14 soft sweeps, and propose methods to differentiate between scenarios.

15 Introduction

16 Inferring adaptive mutations from nucleotide polymorphism data is a major re-
17 search goal in evolutionary genetics, and has been subject to extensive modelling
18 work to determine the footprints they leave in genome data (Stephan 2019). The
19 earliest models focussed on a scenario where a beneficial mutation arose as a
20 single copy before rapidly fixing. Linked neutral mutations then ‘hitchhike’ to
21 fixation with the adaptive variant, reducing diversity around the selected locus
22 (Maynard Smith and Haigh 1974; Kaplan *et al.* 1989). Hitchhiking also increases
23 linkage disequilibrium at regions flanking the selected site, by raising the haplo-
24 type carrying the selected allele to high frequency. It is minimal when measured
25 at sites either side of the selected mutation (Thomson 1977; Innan and Nordborg
26 2003; McVean 2007). These theoretical expectations have spurred the creation of
27 summary statistics for detecting sweeps, usually based on finding genetic regions
28 exhibiting extended haplotype homozygosity (Sabeti *et al.* 2002; Kim and Nielsen
29 2004; Voight *et al.* 2006; Ferrer-Admetlla *et al.* 2014; Vatsiou *et al.* 2016), or an
30 increase in high frequency derived variants (Fay and Wu 2000; Kim and Stephan
31 2002; Nielsen 2005; Boitard *et al.* 2009; Yang *et al.* 2018; Fujito *et al.* 2018).

32 Classic hitchhiking models consider ‘hard’ sweeps, where the common ancestor
33 of an adaptive allele occurs after the onset of selection (Hermisson and Pennings
34 2017). Recent years have seen a focus on ‘soft’ sweeps, where the most recent com-
35 mon ancestor of a beneficial allele appeared before it became selected for (reviewed
36 by Barrett and Schluter (2008); Messer and Petrov (2013); Hermisson and Pennings
37 (2017)). Soft sweeps can originate from beneficial mutations being introduced by
38 recurrent mutation at the target locus (Pennings and Hermisson 2006a,b), or orig-

39 inating from existing standing variation that was either neutral or deleterious (Orr
40 and Betancourt 2001; Innan and Kim 2004; Przeworski *et al.* 2005; Hermisson and
41 Pennings 2005; Wilson *et al.* 2014; Berg and Coop 2015; Wilson *et al.* 2017). A
42 key property of soft sweeps is that the beneficial variant is present on multiple
43 genetic backgrounds as it sweeps to fixation, so different haplotypes may carry the
44 derived allele. This property is often used to detect soft sweeps in genetic data
45 (Peter *et al.* 2012; Vitti *et al.* 2013; Garud *et al.* 2015; Garud and Petrov 2016;
46 Schrider and Kern 2016; Sheehan and Song 2016; Harris *et al.* 2018a; Kern and
47 Schrider 2018; Harris and DeGiorgio 2018, 2019). Soft sweeps have been reported
48 in *Drosophila* (Karasov *et al.* 2010; Garud *et al.* 2015; Garud and Petrov 2016; Vy
49 *et al.* 2017), humans (Peter *et al.* 2012; Schrider and Kern 2017), maize (Fustier
50 *et al.* 2017), *Anopheles* mosquitoes (Xue *et al.* 2019), and pathogens including
51 *Plasmodium falciparum* (Anderson *et al.* 2016) and HIV (Pennings *et al.* 2014;
52 Williams and Pennings 2019). Yet determining how extensive soft sweeps are in
53 nature remains a contentious issue (Jensen 2014; Harris *et al.* 2018b).

54 Up to now, there have only been a few investigations into how dominance
55 affects sweep signatures. In a simulation study, Teshima and Przeworski (2006)
56 explored how recessive mutations spend long periods of time at low frequencies,
57 increasing the amount of recombination that acts on derived haplotypes, weakening
58 signatures of hard sweeps. Fully recessive mutations may need a long time to reach
59 a significantly high frequency to be detectable by genome scans (Teshima *et al.*
60 2006). Ewing *et al.* (2011) have carried out a general mathematical analysis of
61 how dominance affects hard sweeps. Yet the impact of dominance on soft sweeps
62 has yet to be explored in depth.

63 In addition, existing models have so far focussed on randomly mating popu-

64 lations, with haplotypes freely mixing between individuals over generations. Dif-
65 ferent reproductive modes alter how alleles are inherited, affecting the hitchhiking
66 effect. Self-fertilisation, where male and female gametes produced from the same
67 individual can fertilise one another, can alter adaptation rates and selection signa-
68 tures (Hartfield *et al.* 2017). This mating system is prevalent amongst angiosperms
69 (Igic and Kohn 2006), some animals (Jarne and Auld 2006) and fungi (Billiard
70 *et al.* 2011). As the effects of dominance and self-fertilisation become strongly in-
71 tertwined, it is important to consider both together. Dominant mutations are more
72 likely to fix than recessive ones in outcrossers, as they have a higher initial selection
73 advantage (Haldane 1927). Yet recessive alleles can fix more easily in selfers than
74 in outcrossers as homozygote mutations are created more rapidly (Charlesworth
75 1992; Glémin 2012). Furthermore, a decrease in effective recombination rates
76 in selfers (Nordborg *et al.* 1996; Nordborg 2000; Charlesworth and Charlesworth
77 2010) can interfere with selection acting at linked sites, making it likelier that dele-
78 terious mutations hitchhike to fixation with adaptive alleles (Hartfield and Glémin
79 2014), or competition between adaptive mutations at closely-linked loci increases
80 the probability that rare mutations are lost by drift (Hartfield and Glémin 2016).

81 In a constant-sized population, beneficial mutations can be less likely to fix
82 from standing variation (either neutral or deleterious) in selfers as they maintain
83 lower diversity levels (Glémin and Ronfort 2013). Yet adaptation from standing
84 variation becomes likelier in selfers compared to outcrossers under ‘evolutionary
85 rescue’ scenarios, where swift adaptation is needed to prevent population extinc-
86 tion following environmental change. Here, rescue mutations are only present
87 in standing variation as the population size otherwise becomes too small (Glémin
88 and Ronfort 2013). Self-fertilisation further aids this process by creating beneficial

89 homozygotes more rapidly than in outcrossing populations (Uecker 2017).

90 Little data currently exists on the extent of soft sweeps in self-fertilisers. Many
91 selfing organisms exhibit sweep-like patterns, including *Arabidopsis thaliana* (Long
92 *et al.* 2013; Huber *et al.* 2014; Fulgione *et al.* 2018; Price *et al.* 2018); *Caenorhab-*
93 *ditis elegans* (Andersen *et al.* 2012); *Medicago truncatula* (Bonhomme *et al.* 2015);
94 and *Microbotryum* fungi (Badouin *et al.* 2017). Soft sweeps have also been reported
95 in soya bean (Zhong *et al.* 2017). Detailed analyses of these cases has been ham-
96 pered by a lack of theory on how hard and soft sweep signatures should manifest
97 themselves under different self-fertilisation and dominance levels. Previous studies
98 have only focussed on special cases; Hedrick (1980) analysed linkage disequilib-
99 rium caused by a hard sweep under self-fertilisation, while Schoen *et al.* (1996)
100 modelled sweep patterns caused by modifiers that altered the mating system in
101 different ways.

102 To this end, we develop a selective sweep model that accounts for dominance
103 and inbreeding via self-fertilisation. We determine the genetic diversity present
104 following a sweep from either a *de novo* mutation, or from standing variation. We
105 also determine the number of segregating sites and the site frequency spectrum,
106 while comparing results to an alternative soft-sweep model where adaptive alleles
107 arise via recurrent mutation. Note that we focus here on single sweep events, rather
108 than characterising how sweeps affect genome-wide diversity (Elyashiv *et al.* 2016;
109 Campos *et al.* 2017; Booker and Keightley 2018; Rettelbach *et al.* 2019).

110 Results

111 Model Outline

112 We consider a diploid population of size N (carrying $2N$ haplotypes in total).
113 Individuals reproduce by self-fertilisation with probability σ , and outcross with
114 probability $1 - \sigma$. A derived allele arises at a locus, and we are interested in de-
115 termining the population history of neutral regions that are linked to it, with a
116 recombination rate r between them. We principally look at the case where the ben-
117 efcial allele arises from previously-neutral standing variation, and subsequently
118 look at a sweep arising from recurrent mutation. The derived allele initially seg-
119 regates neutrally for a period of time, then becomes advantageous with selective
120 advantage $1 + hs$ when heterozygous and $1 + s$ when homozygous, with $0 < h < 1$
121 and $s > 0$. We further assume that the population size is large and selection is
122 large enough so that the beneficial allele's change in frequency can be modelled
123 deterministically (i.e., $N_ehs \gg 1$ and $1/N_e \ll s \ll 1$). Table 1 lists the notation
124 used in the analysis.

125 Our goal is to determine how the spread of the derived, adaptive allele affects
126 genealogies at linked neutral regions. For a sweep originating from standing vari-
127 ation, we follow the approach of Berg and Coop (2015) and, looking backwards
128 in time, break down the selected allele history into two phases. The first phase
129 (the 'sweep phase') considers the derived allele being selectively favoured from an
130 initial frequency p_0 and spreading through the population. The second phase (the
131 'standing phase') assumes that the derived allele is present at a fixed frequency
132 p_0 . During both phases, a pair of haplotypes can either coalesce, or one of them
133 recombines onto the ancestral background. A schematic is shown in Figure 1.

Symbol	Usage
N	Population size (with $2N$ haplotypes)
σ	Proportion of matings that are self-fertilising
F	Wright's inbreeding coefficient, probability of identity-by-descent at a single gene, equal to $\sigma/(2 - \sigma)$ at steady-state
Φ	Joint probability of identity-by-descent at two loci (Equation 1)
N_e	Effective population size, equal to $N/(1 + F)$ with selfing
r	Recombination rate between loci A and B
r_{eff}	'Effective' recombination rate, approximately equal to $r(1 - 2F + \Phi)$ with selfing
R	$2Nr$, the population-level recombination rate
p_0	Frequency at which the derived allele at B becomes advantageous
$p_{0,A}$	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_{p_0}	Time in the past when derived locus became beneficial
$p(t)$	Frequency of beneficial allele at time t
P_c	Probability of coalescence at time t
P_r	Probability of recombination at time t
P_m	Probability of mutation at time t
P_{NE}	Probability that neutral marker does not coalesce or recombine during sweep phase
$P_{R,Sw}$	Probability that neutral marker recombines during sweep phase
$P_{R,Sd}$	Probability that neutral marker recombines during standing phase
$P_{M,Sw}$	Probability that a lineage mutates during sweep phase
$P_{M,Sd}$	Probability that a lineage mutates during standing phase
H_l, H_h	'Effective' dominance coefficient for allele at low, high frequency
π	Pairwise diversity at site (π_0 is expected value without a sweep)
π_{SV}	Pairwise diversity following sweep from standing variation
π_M	Pairwise diversity following sweep from recurrent mutation
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.

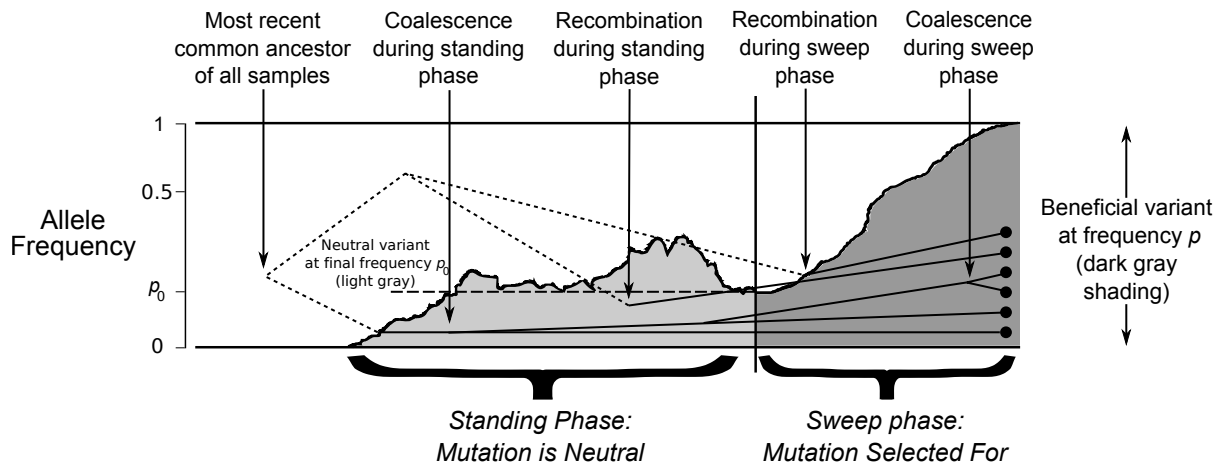


Figure 1. A schematic of the model. The history of the derived variant is separated into two phases; the ‘standing phase’ (shown in light gray), and the ‘sweep phase’ (shown in dark gray). Axis on the left-hand side show allele frequency on a log-scale. Dots on the right-hand side represent a sample of haplotypes taken at the present day, with lines representing their genetic histories. Solid lines represent coalescent histories for the derived genetic background; dotted lines represent coalescent histories for the ancestral, neutral background.

134 During the sweep phase, the derived allele will also cause the spread of linked
 135 haplotypes that it appeared on. Over the course of the sweep, haplotypes are bro-
 136 ken down by recombination; the total number of recombination events is propor-
 137 tional to $r\tau_{p_0}$, where τ_{p_0} is the fixation time of the beneficial allele, given an initial
 138 frequency p_0 (Maynard Smith and Haigh 1974). Dominance and self-fertilisation
 139 have different effects on τ_{p_0} , and therefore the number of fixing haplotypes. If p_0
 140 is low ($\sim 1/2N$) then highly recessive or dominant mutations take longer to go to
 141 fixation (Glémin 2012), which can increase the number of recombination events.
 142 Dominance also affects the nature of the sweep trajectory. For example, recessive
 143 mutations spend more time at a low frequency compared to dominant mutations.
 144 These different sweep trajectories can also affect the final sweep profile (Teshima

145 and Przeworski 2006). Self-fertilisation leads to decreased fixation time of adap-
146 tive mutations through converting heterozygotes to homozygotes (Glémin 2012).
147 Recombination is likelier to act between homozygotes under self-fertilisation, so its
148 effective rate is reduced by a factor $1 - 2F + \Phi$, for $F = \sigma/(2 - \sigma)$ the inbreeding
149 coefficient (Nordborg *et al.* 1996; Nordborg 2000) and Φ the joint probability of
150 identity-by-descent at the two loci (Roze 2009, 2016; Hartfield and Glémin 2016),
151 defined as:

$$\Phi = \frac{\sigma(2 - \sigma - 2(1 - r)r(2 - 3\sigma))}{(2 - \sigma)(2 - (1 - 2(1 - r)r)\sigma)} \quad (1)$$

152 Note that $1 - 2F + \Phi$ approximates to $1 - F$ (as $\Phi \approx F$), unless σ is close to one
153 and r is high (approximately greater than 0.1).

154 During the standing phase, the amount of initial recombinant haplotypes that
155 are swept to fixation depend on the relative rates of recombination and coalescence.
156 The latter occurs with probability proportional to $1/2N_e$ for N_e the effective pop-
157 ulation size. Under self-fertilisation $N_e = N/(1 + F)$ (Wright 1951; Pollak 1987;
158 Charlesworth 1992; Caballero and Hill 1992; Nordborg and Donnelly 1997), so
159 self-fertilisation increases the coalescence probability. This scaling factor remains
160 a good approximation if there is non-Poisson variation in offspring, unless female
161 fitness strongly affects reproduction number (Laporte and Charlesworth 2002).
162 Although we focus on inbreeding via self-fertilisation, the scalings $N_e = N/(1 + F)$
163 and $r_e \approx r(1 - F)$ should also hold under other systems of regular inbreeding
164 (Caballero and Hill 1992; Charlesworth and Charlesworth 2010, Box 8.4).

165 We will outline how both coalescence and recombination act during both of
166 these phases, and use these calculations to determine selective sweep properties.

167 Previous models tended to only determine how lineages recombine away from the
168 derived background during the sweep phase, without considering how two lineages
169 coalesce during the sweep phase. If lineages coalesce during the sweep, then the
170 total number of unique recombination events, and hence the number of linked
171 haplotypes, are reduced. Barton (1998) showed that these coalescent events are
172 negligible only for very strong selection ($\log(Ns) \gg 1$; and B. Charlesworth, un-
173 published results). Hence, accounting for these coalescent events is important for
174 producing accurate matches with simulation results.

175 Throughout, analytical solutions are compared to results from Wright-Fisher
176 forward-in-time stochastic simulations that were ran using SLiM version 3.3 (Haller
177 and Messer 2019). Results for outcrossing populations were also tested using coa-
178 lescent simulations ran with *msms* (Ewing and Hermisson 2010). The simulation
179 methods are outlined in Supplementary File S2.

180 **Data Availability.** File S1 is a *Mathematica* notebook of analytical deriva-
181 tions and simulation results. File S2 contains additional methods, results and
182 figures. File S3 contains copies of the simulation scripts, which are also available
183 from <https://github.com/MattHartfield/SweepDomSelf>. Supplemental mate-
184 rial has also been uploaded to Figshare.

185 **Probability of events during sweep phase**

186 We first look at the probability of events (coalescence or recombination) acting
187 during the sweep phase for the simplest case of two alleles. Looking back in time
188 following the fixation of the derived mutation, sites linked to the beneficial allele
189 can either coalesce or recombine onto the ancestral genetic background. Let $p(t)$

190 be the adaptive mutation frequency at time t , defined as the number of genera-
191 tions prior to the present day. Further define $p(0) = 1$ (i.e., the allele is fixed at
192 the present day), and τ_{p_0} the time in the past when the derived variant became
193 beneficial (i.e., $p(\tau_{p_0}) = p_0$).

194 For a pair of haplotype samples carrying the derived allele, if it is at frequency
195 $p(t)$ at time t , this lineage pair can either coalesce or one of the haplotypes recom-
196 bine onto the ancestral background. Each event occurs with probability:

$$\begin{aligned} P_c(t) &= \frac{1}{2N_e p(t)} = \frac{(1+F)}{2N p(t)} \\ P_r(t) &= 2r_{eff}(1-p(t)) = 2r(1-2F+\Phi)(1-p(t)) \end{aligned} \tag{2}$$

197 Equation 2 is based on those obtained by Kaplan *et al.* (1989), assuming that
198 $N_e = N/(1+F)$ due to self-fertilisation (Pollak 1987; Charlesworth 1992; Ca-
199 ballero and Hill 1992; Nordborg and Donnelly 1997), and $r_{eff} = r(1-2F+\Phi)$
200 is the ‘effective’ recombination rate after correcting for increased homozygosity
201 due to self-fertilisation (Nordborg *et al.* 1996; Nordborg 2000; Charlesworth and
202 Charlesworth 2010; Roze 2009, 2016; Hartfield and Glémin 2016). Equation 2
203 demonstrates how each event is differently influenced by p . In particular, the per-
204 generation coalescence probability P_c can be small unless p is close to $1/2N$. The
205 total probability that coalescence occurs during the sweep phase increases if the
206 beneficial allele spends a sizeable time at low frequency, e.g., when it is recessive.
207 The terms in Equation 2 can also be defined as functions of p .

208 We are interested in calculating (i) the probability P_{NE} that no coalescence or
209 recombination occurs in the sweep phase; (ii) the probability $P_{R,Sw}$ that recombi-

210 nation acts on a lineage to transfer it to the neutral background that is linked to
 211 the ancestral allele, assuming that no more than one recombination event occurs
 212 per generation (see Campos and Charlesworth (2019) for derivations assuming
 213 multiple recombination events). We will go through these probabilities in turn to
 214 determine expected pairwise diversity. For P_{NE} , the total probability that the two
 215 lineages do not coalesce or recombine over τ_{p_0} generations equals:

$$\begin{aligned}
 P_{NE} &= \prod_{t=0}^{\tau_{p_0}} [1 - P_c(t) - P_r(t)] \\
 &\approx \exp\left(-\int_{t=0}^{\tau_{p_0}} [P_c(t) + P_r(t)] dt\right) && \text{assuming } P_c, P_r \ll 1 \\
 &\approx \exp\left(-\int_{t=0}^{\tau_{p_0}} \left[\frac{1+F}{2Np(t)} + 2r(1-2F+\Phi)(1-p(t))\right] dt\right) \\
 &\approx \exp\left(-\int_{p=1-\epsilon}^{p_0} \left[\frac{1+F}{2Np} + 2r(1-2F+\Phi)(1-p)\right] dp\right) && \text{taking the integral over } p
 \end{aligned}
 \tag{3}$$

216 Here ϵ is a small term and $1 - \epsilon$ is the upper limit of the deterministic spread
 217 of the beneficial allele. We will discuss in the section ‘Effective starting frequency
 218 from a *de novo* mutation’ what a reasonable value for ϵ should be. Also note that
 219 we switch from a discrete-time calculation to a continuous-time calculation, which
 220 can give simplifying results. To calculate P_{NE} we insert the deterministic change
 221 in allele frequency p (Glémin 2012):

$$\frac{dp}{dt} = -sp(1-p)(F+h-Fh+(1-F)(1-2h)p)
 \tag{4}$$

222 Note the negative factor in Equation 4 since we are looking back in time. By

223 substituting Equation 4 into Equation 3, we obtain an analytical solution for P_{NE} ,
 224 although the resulting expression is complicated (Section A of Supplementary File
 225 S1).

226 To calculate $P_{R,Sw}$, the probability that recombination acts during the sweep,
 227 we first calculate the probability that recombination occurs when the beneficial
 228 allele is at frequency p' . Here, no events occur in the time leading up to p' , then
 229 a recombination event occurs with probability $P_r(p') = 2r(1 - 2F + \Phi)(1 - p')$.
 230 $P_{R,Sw}$ is obtained by integrating this probability over the entire sweep from time
 231 0 to τ_{p_0} :

$$P_{R,Sw} \approx \int_{p'=1-\epsilon}^{p_0} \frac{P_{R,p'}}{dp'/dt} dp' \quad (5)$$

where:

$$\begin{aligned} P_{R,p'} &= \exp \left[- \int_{p=1-\epsilon}^{p'} \frac{P_c(p) + P_r(p)}{dp/dt} dp \right] \cdot P_r(p') \\ &= \exp \left[- \int_{p=1-\epsilon}^{p'} \frac{\frac{1+F}{2Np} + 2r(1 - 2F + \Phi)(1 - p)}{dp/dt} dp \right] \cdot [2r(1 - 2F + \Phi)(1 - p')] \end{aligned} \quad (6)$$

232 Note that the exponential term of $P_{R,p'}$ is different from P_{NE} (Equation 3) since
 233 the upper integral limit is to p' rather than p_0 . That is, it only covers part of the
 234 sweep phase. Equation 5 is evaluated numerically. In Supplementary File S2, we
 235 provide a ‘star-like’ analytical approximation to P_{NE} that assumes no coalescence
 236 during the sweep phase.

237 **Probability of coalescence from standing variation**

238 The variant becomes advantageous at frequency p_0 . We assume that p_0 , and hence
239 event probabilities, remain fixed over time. Berg and Coop (2015) have shown this
240 assumption provides a good approximation to coalescent rates during the standing
241 phase. The outcome during the standing phase is thus determined by competing
242 Poisson processes. The two haplotypes could coalesce, with an exponentially-
243 distributed waiting time with rate $P_c(p_0) = (1 + F)/(2Np_0)$. Alternatively, one
244 of the two haplotypes could recombine onto the ancestral background with mean
245 waiting time $P_r(p_0) = 2r_{eff}(1 - p_0)$. For two competing exponential distribu-
246 tions with rates λ_1 and λ_2 , the probability of the first event occurring given an
247 event happens equals $\lambda_1/(\lambda_1 + \lambda_2)$ (Wakeley 2009). Hence the probability that
248 recombination occurs instead of coalescence equals:

$$\begin{aligned} P_{R,sd} &= \frac{P_r(p_0)}{P_c(p_0) + P_r(p_0)} \\ &= \frac{2r_{eff}(1 - p_0)}{\frac{1+F}{2Np_0} + 2r_{eff}(1 - p_0)} \\ &= \frac{2R(1 - 2F + \Phi)p_0(1 - p_0)/(1 + F)}{1 + 2R(1 - 2F + \Phi)p_0(1 - p_0)/(1 + F)} \\ &\approx \frac{2R(1 - \sigma)p_0(1 - p_0)}{1 + 2R(1 - \sigma)p_0(1 - p_0)} \end{aligned} \quad (7)$$

249 The probability of coalescence rather than recombination is $P_{C,sd} = 1 - P_{R,sd}$.
250 Here $R = 2Nr$ is the population-scaled recombination rate. The final approxima-
251 tion arises as $(1 - 2F + \Phi)/(1 + F) \approx (1 - F)/(1 + F) = (1 - \sigma)$ if $\Phi \approx F$. This term
252 reflects how increased homozygosity reduces both effective recombination and N_e ,
253 with the latter making coalescence more likely. In addition, it also highlights how

254 the signature of a sweep from standing variation, as characterised by the spread
255 of different initial recombinant haplotypes, is spread over an increased distance of
256 $1/(1 - \sigma)$ under self-fertilisation.

257 **Effective starting frequency for a *de novo* mutation**

258 When a new beneficial mutation appears at a single copy, it is highly likely to
259 go extinct by chance (Fisher 1922; Haldane 1927). Beneficial mutations that in-
260 crease in frequency faster than expected when rare are more able to overcome this
261 stochastic loss and reach fixation. These beneficial mutations will hence display
262 an apparent ‘acceleration’ in their logistic growth, equivalent to having a starting
263 frequency that is greater than $1/(2N)$ (Maynard Smith 1976; Barton 1998; Desai
264 and Fisher 2007; Martin and Lambert 2015). Correcting for this acceleration is
265 important to accurately model hard sweep signatures, and inform on the mini-
266 mum level of standing variation needed to differentiate a hard sweep from one
267 originating from standing variation.

268 In Section B of Supplementary File S1, we determine that hard sweeps that go
269 to fixation have the following effective starting frequency:

$$p_{0,A} = \frac{1 + F}{4NsH_l} \quad (8)$$

270 where $H_l = F + h - Fh$ is the effective dominance coefficient for mutations at a low
271 frequency. This result is consistent with those of Martin and Lambert (2015), who
272 obtained a distribution of effective starting frequencies using stochastic differential
273 equations. This acceleration effect can create substantial increases in the effective
274 p_0 , especially for recessive mutations (Figure 2).

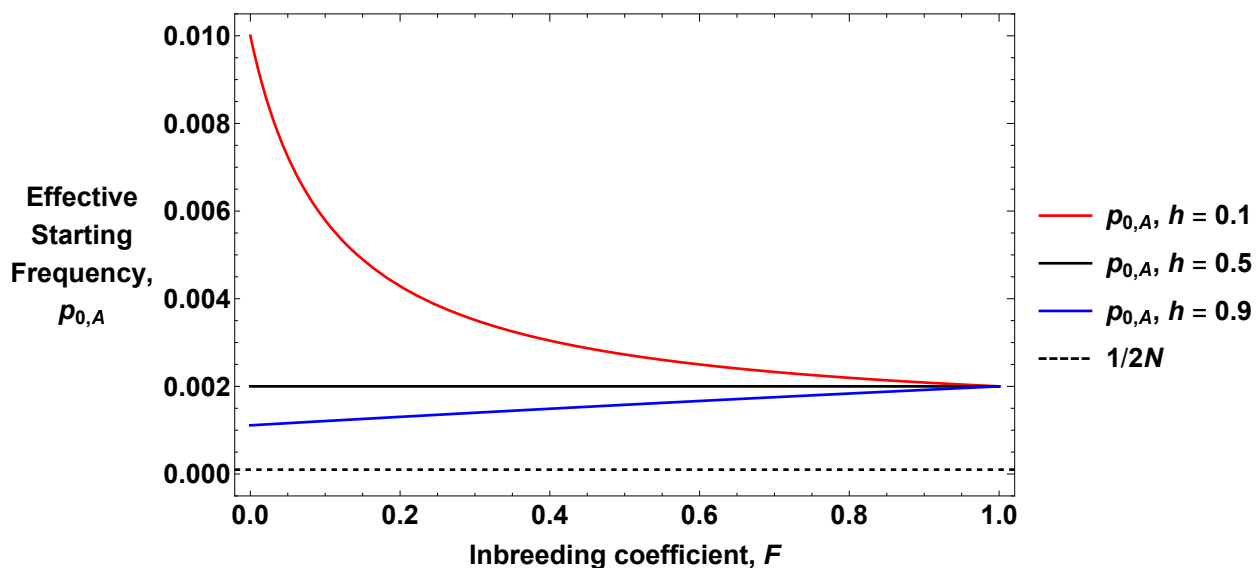


Figure 2. Examples of the effective starting frequency. Equation 8 is plotted as a function of F for different dominance values, as shown in the legend. Other parameters are $N = 5,000$, $s = 0.05$. The dashed line shows the actual starting frequency, $1/2N$.

275 **Effective final frequency:** The effective final frequency of the derived allele
276 $1 - \epsilon$, at which its spread is no longer deterministic, can be obtained by setting
277 $\epsilon = p_{0,A}(1 - h)$; that is, by substituting H_l to $H_h = 1 - h + Fh$ in Equation 8.
278 Van Herwaarden and Van der Wal (2002) determined that the sojourn time for
279 an allele with dominance coefficient h that is increasing in frequency, is the same
280 for an allele decreasing in frequency with dominance $1 - h$. Glémin (2012) showed
281 that this result also holds under any inbreeding value F (and B. Charlesworth,
282 unpublished results).

283 **Expected Pairwise Diversity**

284 We use P_{NE} , $P_{R,sw}$ and $P_{R,sd}$ to calculate the expected pairwise diversity (denoted
285 π) present around a sweep. During the sweep phase, the two neutral sites could
286 either coalesce, or one of them recombines onto the ancestral background. If
287 coalescence occurs, since it does so in the recent past then it is assumed that no
288 diversity exist between samples, i.e., $\pi \approx 0$ for π the average number of differences
289 between two alleles (Tajima 1983). In reality there may be some residual diversity
290 caused by appearance of mutations during the sweep phase; we do not account
291 for these mutations while calculating π but will do so when calculating the site-
292 frequency spectrum. Alternatively, if one of the two samples recombines onto the
293 neutral background, they will have the same pairwise diversity between them as
294 the background population (π_0). If the two samples trace back to the standing
295 phase (with probability P_{NE}) then the same logic applies. Hence the expected
296 diversity following a sweep π_{SV} , relative to the background value π_0 , equals:

$$\mathbb{E}\left(\frac{\pi_{SV}}{\pi_0}\right) = P_{R,sw} + (P_{NE} \cdot P_{R,sd}) \quad (9)$$

297 The full solution to Equation 9 can be obtained by plugging in the relevant
298 parts from Equations 3, 5 and 7, which we evaluate numerically. Equation 9 is
299 undefined for $h = 0$ or 1 with $\sigma = 0$; these cases can be derived separately.

300 Figure 3 plots Equation 9 with different dominance, self-fertilisation, and stand-
301 ing frequency values. The analytical solution fits well compared to forward-in-time
302 simulations, yet slightly overestimates them for high self-fertilisation frequencies.
303 It is unclear why this mismatch arises. One explanation could be that drift effects
304 are magnified under self-fertilisation, which causes a quicker sweep fixation time

305 than expected from deterministic spread, if conditioning on a sweep going to fixa-
306 tion. Although $p_{0,A}$ (Equation 8) captures these drift effects for rare alleles, there
307 may be additional effects that are not accounted for. Under complete outcross-
308 ing, baseline diversity is restored (i.e., $\mathbb{E}(\pi_{SV}/\pi_0)$ goes to 1) closer to the sweep
309 origin for recessive mutations ($h = 0.1$), compared to semidominant ($h = 0.5$)
310 or dominant ($h = 0.9$) mutations. Sweeps caused by dominant and semidomi-
311 nant mutations result in a similar genetic diversity, so these cases may be hard to
312 differentiate from diversity data alone.

313 These results can be better understood by examining the underlying allele
314 trajectories, using logic described by Teshima and Przeworski (2006) (Figure 4).
315 For outcrossing populations, recessive mutations spend most of the sojourn time at
316 low frequencies, maximising recombination events and restoring neutral variation.
317 These trajectories mimic sweeps from standing variation, which spend extended
318 periods of time at low frequencies in the standing phase. Conversely, dominant
319 mutations spend most of their time at high frequencies, reducing the chance for
320 neutral markers to recombine onto the ancestral background.

321 As self-fertilisation increases, sweep signatures become similar to the co-dominant
322 case as the derived allele is more likely to spread as a homozygote, weakening the
323 influence that dominance exerts over beneficial allele trajectories. Increasing p_0
324 also causes sweeps with different dominance coefficients to produce comparable
325 signatures, as beneficial mutation trajectories become similar after conditioning
326 on starting at an elevated frequency.

327 *Star-like approximation.* An analytical approximation can be obtained by using
328 the ‘star-like’ result for P_{NE} (described in Supplementary Files S1, S2). In this
329 case the expected pairwise diversity approximates to:

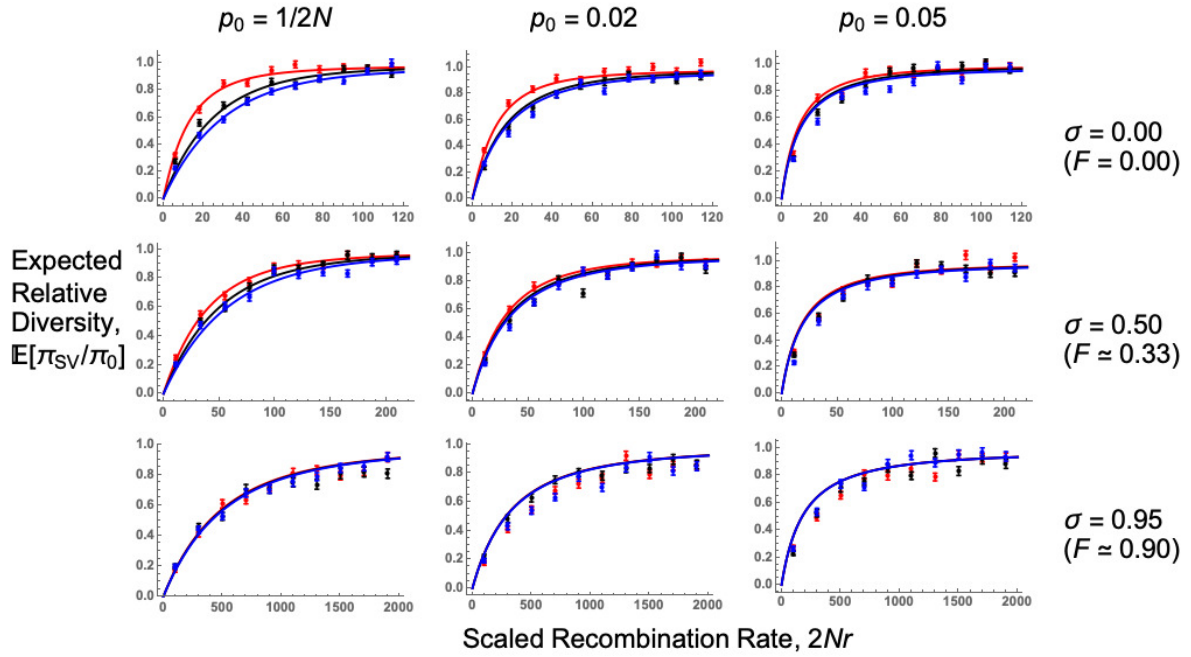


Figure 3. Expected relative pairwise diversity following a selective sweep. Plots of $\mathbb{E}(\pi_{SV}/\pi_0)$ as a function of the recombination rate scaled to population size $2Nr$. Lines are analytical solutions (Equation 9), points are forward-in-time simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$ (note μ is scaled by N , not N_e), and dominance coefficient $h = 0.1$ (red lines, points), 0.5 (black lines, points), or 0.9 (blue lines, points). Values of p_0 and self-fertilisation rates σ used are shown for the relevant row and column; note the x -axis range changes with the self-fertilisation rate. For $p_0 = 1/2N$ we use $p_{0,A}$ in our model, as given by Equation 8. Further results are plotted in Section C of Supplementary File S1.

$$\begin{aligned}
 \mathbb{E}_{SL} \left(\frac{\pi_{SV}}{\pi_0} \right) &= 1 - (P_{NE} \cdot P_{C,sd}) \\
 &= 1 - \left[\frac{1}{1 + 2R(1 - 2F + \Phi)p_0(1 - p_0)/(1 + F)} \right] \cdot \left[\frac{H_l}{H_h} \left(\frac{1}{p_0} + 1 \right) - 1 \right]^{-2r(1-2F+\Phi)/(H_l s)}
 \end{aligned}
 \tag{10}$$

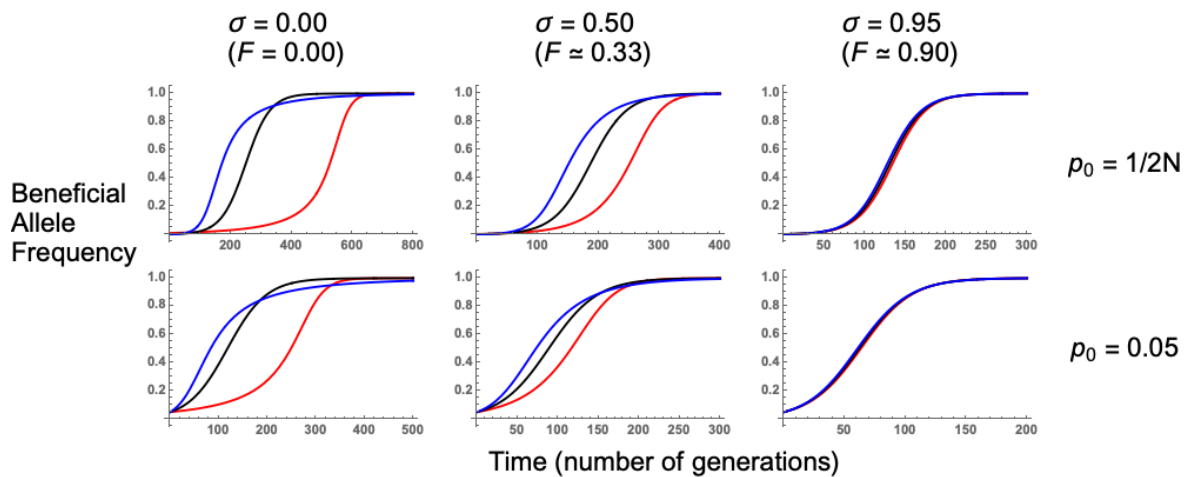


Figure 4. Beneficial allele trajectories. These were obtained by numerically evaluating the negative of Equation 4 forward in time. $N = 5,000$, $s = 0.05$, and h equals either 0.1 (red lines), 0.5 (black lines), or 0.9 (blue lines). Values of p_0 and self-fertilisation rates σ used are shown for the relevant row and column. Note the different x -axis scales used in each panel. Further results are plotted in Section C of Supplementary File S1.

330 Note that Equation 10 instead uses the probability of coalescence during the
 331 standing phase, $P_{C,sd} = 1 - P_{R,sd}$. This approximation reflects similar formulas
 332 for diversity following soft sweeps in haploid outcrossing populations (Pennings
 333 and Hermisson 2006b; Berg and Coop 2015). There is a factor of two in the
 334 power term to account for two lineages. In Supplementary File S2 we demonstrate
 335 that this equation overestimates the relative diversity following a selective sweep.
 336 This mismatch arises since the star-like assumption of no coalescence during the
 337 sweep phase is only accurate for very strongly selected mutations (Barton 1998; B.
 338 Charlesworth, unpublished results). Hence it is important to consider coalescence
 339 during the sweep phase to accurately model selective sweeps that do not have an
 340 extremely high selection coefficient.

341 Site Frequency Spectrum

342 The star-like approximation can be used to obtain analytical solutions for the
343 number of segregating sites and the site frequency spectrum (i.e., the probability
344 that $l = 1, 2 \dots n - 1$ of n alleles carry derived variants). The full derivation
345 for these statistics are outlined in Supplementary File S2, which uses the star-like
346 approximation. Figure 5 plots the SFS (Equation A12 in Supplementary File S2)
347 alongside simulation results. Analytical results fit the simulation data well after
348 including an adjusted singleton class, which accounts for recent mutations that
349 arise on the derived background during both the standing and sweep phases (Berg
350 and Coop 2015). Including this new singleton class improves the model fit, but
351 there remains a tendency for analytical results to underestimate the proportion of
352 low- and high-frequency classes ($l = 1$ and 9 in Figure 5), and overestimate the
353 proportion of intermediate-frequency classes. Additional inaccuracies could have
354 arisen due to the use of the star-like approximation, which assumes that there is
355 no coalescence during the sweep phase.

356 Hard sweeps in either outcrossers or partial selfers are characterised by a large
357 number of singletons and highly-derived variants (Figure 5), which is a typical
358 selective sweep signature (Braverman *et al.* 1995; Barton 1998; Kim and Stephan
359 2002). As the initial frequency p_0 increases, so does the number of intermediate-
360 frequency variants (Figure 5). This signature is often seen as a characteristic of
361 soft sweeps (Pennings and Hermisson 2006b; Berg and Coop 2015). Recessive
362 hard sweeps ($h = 0.1$ and $p_0 = 1/2N$) can produce SFS profiles that are similar to
363 sweeps from standing variation, as there are an increased number of recombination
364 events occurring since the allele is at a low frequency for long time periods (Fig-

365 ure 4). With increased self-fertilisation, both hard and soft sweep signatures (e.g.,
 366 increased number of intermediate-frequency alleles) are recovered when measuring
 367 the SFS at a longer recombination distance than in outcrossers (Figure 5, bottom
 368 row). This is an example of how signatures of sweeps from standing variation
 369 are extended over an increased recombination distance of around $1/(1 - \sigma)$, as
 370 demonstrated by Equation 7.

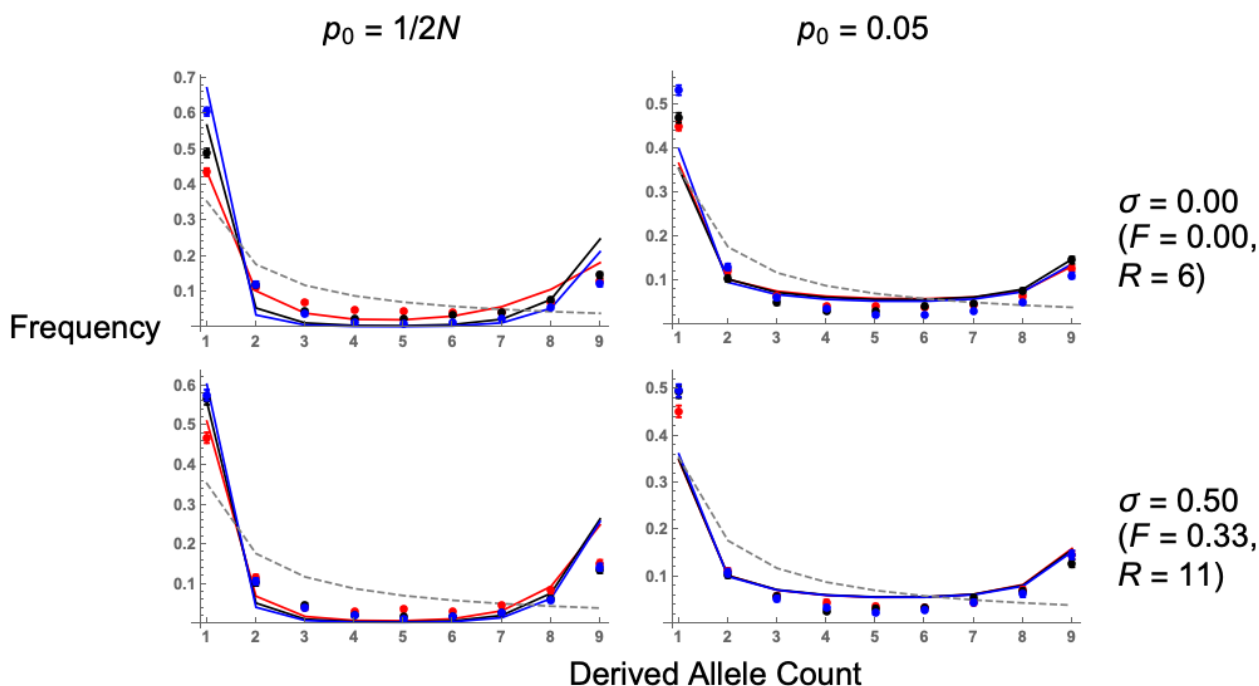


Figure 5. Expected site frequency spectrum, in flanking regions to the adaptive mutation, following a selective sweep. Lines are analytical solutions (Equation A12 in Supplementary File S2), points are simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$, and dominance coefficient $h = 0.1$ (red lines, points), 0.5 (black lines, points), or 0.9 (blue lines, points). The neutral SFS is also included for comparisons (grey dashed line). Values of p_0 , self-fertilisation rates σ and recombination distances R are shown for the relevant row and column. Results for other recombination distances are in Section E of Supplementary File S1.

371 **Soft sweeps from recurrent mutation**

372 So far, we have only focussed on a soft sweep that arises from standing variation.
373 An alternative type of soft sweep is one where recurrent mutation at the selected
374 locus introduces the beneficial allele onto different genetic backgrounds. We can
375 examine this case by modifying existing results. Below we derive the expected
376 relative diversity between two alleles following this type of soft sweep, and outline
377 the SFS for more than two samples in Supplementary File S2.

378 In this model, derived alleles arise from recurrent mutation and are instan-
379 taneously beneficial (i.e., there is no ‘standing phase’). During the sweep phase,
380 lineages can escape the derived background by recombination, or if they are derived
381 from a mutation event. If the beneficial allele is at frequency p then the probability
382 of being descended from an ancestral allele by mutation is $P_m(p) = 2\mu_b(1 - p)/p$,
383 for μ_b the mutation probability (Pennings and Hermisson 2006b). Denote the
384 probability of a lineage experiencing recombination or mutation during this sweep
385 phase by $P_{R,sw}$, $P_{M,sw}$ respectively. In both these cases the expected diversity
386 present at linked sites is π_0 . If none of these events arise with probability P_{NE} ,
387 then remaining lineages can either coalesce, or they arise from independent muta-
388 tion events. If they coalesce then they have approximately zero pairwise diversity
389 between them; alternatively, they have different origins and thus exhibit the same
390 pairwise diversity π_0 as the neutral background. Let $P_{M,sd}$ denote the probability
391 that mutation occurs at the sweep origin, as opposed to coalescence.

392 Following this logic, the expected relative diversity for a sweep arising from

393 recurrent mutation equals (with additional details in Supplementary File S1):

$$\mathbb{E}\left(\frac{\pi_M}{\pi_0}\right) = P_{R,sw} + P_{M,sw} + (P_{NE} \cdot P_{M,sd}) \quad (11)$$

394 π_M denotes the diversity around a soft sweep from recurrent mutation. $P_{R,sw}$,
 395 P_{NE} are similar to the equations used when modelling a sweep from standing
 396 variation. They are both modified to account for additional beneficial mutation
 397 arising during the sweep phase:

$$P_{R,sw} \approx \int_{p'=1-\epsilon}^{p_0} \frac{P_{R,p'}}{dp'/dt} dp' \quad (12)$$

where:

$$\begin{aligned} P_{R,p'} &= \exp\left[-\int_{p=1-\epsilon}^{p'} \frac{P_c(p) + P_r(p) + P_m(p)}{dp/dt} dp\right] \cdot P_r(p') \\ &= \exp\left[-\int_{p=1-\epsilon}^p \frac{\frac{1+F}{2Np} + 2r(1-2F+\Phi)(1-p) + \frac{2\mu_b(1-p)}{p}}{dp/dt} dp\right] \cdot [2r(1-2F+\Phi)(1-p')] \end{aligned} \quad (13)$$

398 and:

$$\begin{aligned} P_{NE} &\approx \exp\left(-\int_{p=1-\epsilon}^{p_{0,A}} \left[\frac{P_c(p) + P_r(p) + P_m(p)}{dp/dt}\right] dp\right) \\ &= \exp\left(-\int_{p=1-\epsilon}^{p_{0,A}} \left[\frac{\frac{1+F}{2Np} + 2r(1-2F+\Phi)(1-p) + \frac{2\mu_b(1-p)}{p}}{dp/dt}\right] dp\right) \end{aligned} \quad (14)$$

399 Note that Equation 14 has an upper integral limit of $p_{0,A}$, as opposed to a
 400 general p_0 used in the sweep from standing variation model, reflecting that there

401 is no standing phase.

402 $P_{M,sw}$ is the mutation probability during the sweep phase, and is similar to
 403 Equation 13 except that $2r(1 - 2F + \Phi)(1 - p')$ is replaced by $2\mu_b(1 - p')/p'$, for p'
 404 is the derived allele frequency when the event occurs. $P_{M,sd}$ is the probability that,
 405 at the sweep origin, the derived allele appears by mutation instead of coalescing,
 406 and is defined in a similar manner to $P_{R,sd}$ (Equation 7):

$$\begin{aligned}
 P_{M,sd} &= \frac{P_m(p_{0,A})}{P_c(p_{0,A}) + P_m(p_{0,A})} \\
 &= \frac{\frac{2\mu_b(1-p_{0,A})}{p_{0,A}}}{\frac{1+F}{2Np_{0,A}} + \frac{2\mu_b(1-p_{0,A})}{p_{0,A}}} \\
 &= \frac{2\Theta_b(1 - p_{0,A})}{1 + F + 2\Theta_b(1 - p_{0,A})} \tag{15}
 \end{aligned}$$

407 where $\Theta_b = 2N\mu_b$. The coalescence probability is $1 - P_{M,sd}$. Equation 15 implies
 408 that self-fertilisation makes it more likely for beneficial mutations to coalesce at the
 409 start of a sweep, rather than arising from independent mutation events. Hence the
 410 signatures of soft sweeps via recurrent mutation will be weakened under inbreeding.

411 Figure 6 compares $\mathbb{E}(\pi_{SV}/\pi_0)$ in the standing variation case, and $\mathbb{E}(\pi_M/\pi_0)$ for
 412 the recurrent mutation case, under different levels of self-fertilisation. While dom-
 413 inance only weakly affects sweep signatures arising from standing variation under
 414 outcrossing, it more strongly affects sweeps from recurrent mutation in outcrossing
 415 populations, as each variant arises from an initial frequency close to $1/(2N)$ (Fig-
 416 ure 4). Second, the two models exhibit different behaviour close to the selected
 417 locus (R close to zero). The recurrent mutation model has non-zero diversity
 418 levels, while the standing variation model exhibits zero diversity. As R increases,

419 diversity eventually becomes higher for the standing variation case compared to
420 the recurrent mutation case. We can heuristically determine when this transition
421 occurs as follows. Assume a large population size but weak recombination and mu-
422 tation rates. Hence, it is unlikely that any events occur during the sweep phase, so
423 $P_{R,sw}$, $P_{M,sw} \approx 0$ and $P_{NE} \approx 1$. Then the expected relative diversity (Equation 11)
424 equals $P_{R,sd}$ for a sweep from standing variation, and $P_{M,sd}$ for one from recurrent
425 mutation. To find the recombination rate R_{Lim} at which a sweep from recurrent
426 mutation yields higher diversity than one from standing variation, we find the R
427 value needed to equate the two probabilities, giving:

$$\begin{aligned} R_{Lim} &= \frac{\Theta_b}{p_0(1 - 2F + \Phi)} \\ &\approx \frac{\Theta_b}{p_0(1 - F)} \end{aligned} \quad (16)$$

428 The last approximation arises as $\Phi \approx F$. Hence for a fixed Θ_b , the window
429 where recurrent mutations create higher diversity near the selected locus increases
430 for lower p_0 or higher F , since both these factors reduces the potential for re-
431 combination to create new haplotypes during the standing phase. Equation 16 is
432 generally accurate when sweeps from standing variation have higher diversity than
433 sweeps with recurrent mutations (Figure 6, bottom row), but becomes inaccurate
434 for $h = 0.1$ in outcrossing populations, as some events are likely to occur during
435 the sweep phase. In Supplementary File S2 we show how similar results apply to
436 the SFS.

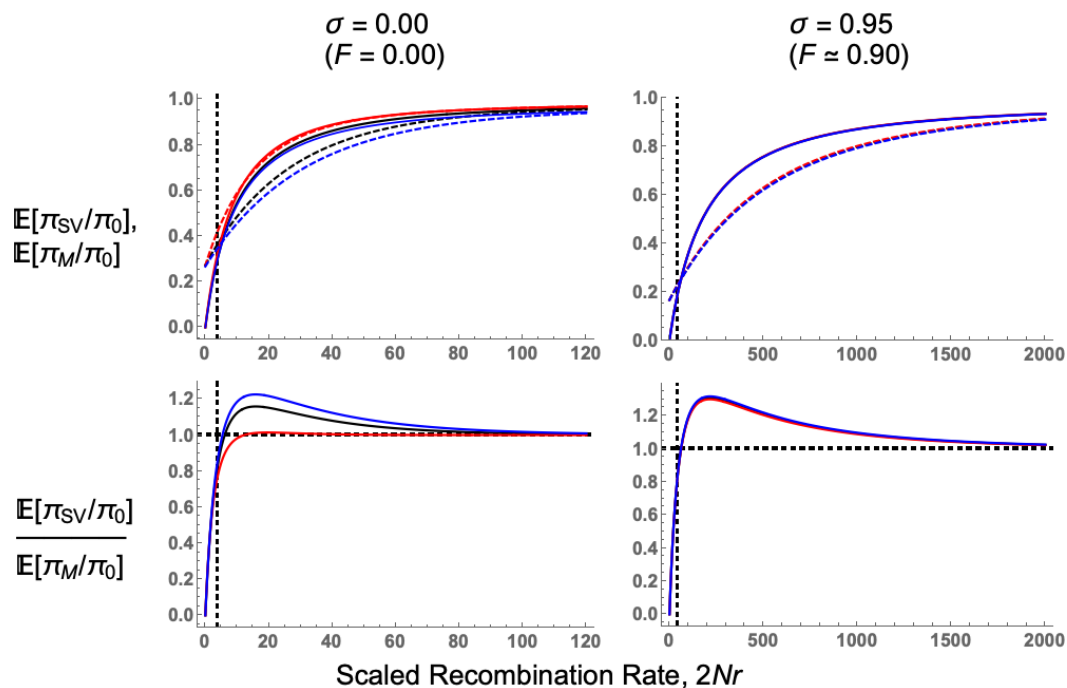


Figure 6. Comparing sweeps from recurrent mutation to those from standing variation. Top row: comparing relative diversity following a soft sweep, from either standing variation (Equation 9 with $p_0 = 0.05$, solid lines) or recurrent mutation (using Equation 11 with $\Theta_b = 0.2$, dashed lines). $N = 5,000$, $s = 0.05$, and dominance coefficient $h = 0.1$ (red lines), 0.5 (black lines), or 0.9 (blue lines). Bottom row: the ratio of the diversity following a sweep from standing variation to one from recurrent mutation. Parameters for each panel are as in the respective plot for the top row. Vertical dashed black line indicates R_{Lim} (the approximate form of Equation 16); horizontal dashed line in the bottom-row plots show when the ratio equals 1. Note the different x -axis between left- and right-hand panels. Results are also plotted in Section F of Supplementary File S1.

437 Discussion

438 Summary of Theoretical Findings

439 While there has been many investigations into how different sweep processes can
 440 be detected from next-generation sequence data (Pritchard and Di Rienzo 2010;

441 Messer and Petrov 2013; Stephan 2016; Hermisson and Pennings 2017), these
442 models generally assumed idealised randomly mating populations and beneficial
443 mutations that are semidominant ($h = 0.5$). Here we have created a more general
444 selective sweep model, with arbitrary self-fertilisation and dominance levels. Our
445 principal focus is on comparing a hard sweep arising from a single allele copy
446 to a soft sweep arising from standing variation, but we also consider the case of
447 recurrent mutation (Figure 6).

448 We find that the qualitative patterns of different selective sweeps under selfing
449 remain similar to expectations from outcrossing models. In particular, a sweep
450 from standing variation still creates an elevated number of intermediate-frequency
451 variants compared to a sweep from *de novo* mutation (Figures 5, 6). This pattern is
452 standard for soft sweeps (Pennings and Hermisson 2006b; Messer and Petrov 2013;
453 Berg and Coop 2015; Hermisson and Pennings 2017) so existing statistical methods
454 for detecting them (e.g., observing an higher than expected number of haplotypes;
455 Vitti *et al.* (2013); Garud *et al.* (2015)) can, in principle, also be applied to self-
456 ing organisms. Under self-fertilisation, these signatures are stretched over longer
457 physical regions than in outcrossers. These extensions arise as self-fertilisation
458 affects gene genealogies during both the sweep and standing phases in different
459 ways. During the sweep phase, beneficial alleles fix more rapidly under higher
460 self-fertilisation as homozygous mutations are created more rapidly (Charlesworth
461 1992; Glémin 2012). In addition, the effective recombination rate is reduced by
462 approximately $1 - F$ (Nordborg *et al.* 1996; Nordborg 2000; Charlesworth and
463 Charlesworth 2010), and slightly more for highly inbred populations (Roze 2009,
464 2016). These two effects mean that neutral variants linked to an adaptive allele are
465 less likely to recombine onto the neutral background during the sweep phase, as re-

466 flected in Equation 3 for P_{NE} . During the standing phase, two haplotypes are more
467 likely to coalesce under high levels of self-fertilisation since N_e is decreased by a fac-
468 tor $1/(1+F)$ (Pollak 1987; Charlesworth 1992; Caballero and Hill 1992; Nordborg
469 and Donnelly 1997). This effect, combined with a reduced effective recombination
470 rate, means that the overall recombination probability during the standing phase
471 is reduced by a factor $(1-\sigma)$ (Equation 7). Hence intermediate-frequency variants,
472 which could provide evidence of adaptation from standing variation, will be spread
473 out over longer genomic regions (this result can be seen in the site-frequency spec-
474 trum results, Figure 5). The elongation of sweep signatures means sweeps from
475 standing variation can be easier to detect in selfing organisms than in outcrossers.
476 Conversely, sweeps from recurrent mutation will have weakened signatures under
477 self-fertilisation. This result is due to a reduced effective population size, making
478 it likelier that lineages trace back to a common ancestor rather than independent
479 mutation events.

480 We have also investigated how dominance affects soft sweep signatures, since
481 previous analyses have only focussed on how dominance affects hard sweeps (Teshima
482 and Przeworski 2006; Teshima *et al.* 2006; Ewing *et al.* 2011). In outcrossing or-
483 ganisms, recessive mutations leave weaker sweep signatures than additive or domi-
484 nant mutations as they spend more time at low frequencies, increasing the amount
485 of recombination that restores neutral variation (Figures 3, 4). With increased
486 self-fertilisation, dominance has a weaker impact on sweep signatures as most mu-
487 tations are homozygous (Figure 4). We also show that the SFS for recessive alleles
488 can resemble a soft sweep, with a higher number of intermediate-frequency vari-
489 ants than for other hard sweeps (Figure 5). Dominance only weakly affects sweeps
490 from standing variation, as trajectories of beneficial alleles become similar once

491 the variant's initial frequency exceeds $1/(2N)$ (Figures 3, 4). Yet different domi-
492 nance levels can affect sweep signatures if the beneficial allele is reintroduced by
493 recurrent mutation (Figure 6). Hence if one wishes to understand how dominance
494 affects sweep signatures, it is also important to consider which processes underlie
495 observed patterns of genetic diversity.

496 These results also demonstrate that the effects of dominance on sweeps are
497 not necessarily intuitive. For example, both highly dominant and recessive muta-
498 tions have elongated fixation times compared to co-dominant mutations (Glémin
499 2012). Based on this intuition, one could expect both dominant and recessive
500 mutations to both produce weaker sweep signatures than co-dominant ones. In
501 practice, dominant mutations have similar sweep signatures to co-dominant mu-
502 tations (Figures 3, 5), and recessive sweeps could produce similar signatures to
503 sweeps to standing variation (Figure 5). Dominance also has a weaker impact on
504 sweeps on standing variation (Figures 3, 5).

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

506 These theoretical results shed light onto how to distinguish between soft sweeps
507 that arise either from standing variation, or from recurrent mutation. Both mod-
508 els are characterised by an elevated number of intermediate-frequency variants,
509 in comparison to a hard sweep. Yet sweeps arising from recurrent mutation have
510 non-zero diversity at the selected locus, whereas a sweep from standing variation
511 exhibits approximately zero diversity. Hence a sweep from recurrent mutation
512 shows intermediate-frequency variants closer to the beneficial locus, compared to
513 sweeps from standing variation (Figures 6 and C in Supplementary File S2). Fur-
514 ther from the selected locus, a sweep from standing variation exhibits greater

515 variation than one from recurrent mutation, due to recombinant haplotypes being
516 created during the standing phase. Equation 16 provides a simple condition for
517 R_{Lim} , the recombination distance needed for a sweep from standing variation to
518 exhibit higher diversity than one from recurrent mutation; from this equation, we
519 see that the size of this region increases under higher self-fertilisation. Hence it
520 may be easier to differentiate between these two sweep scenarios in self-fertilising
521 organisms.

522 Differences in haplotype structure between sweeps from either standing varia-
523 tion or recurrent mutation should be more pronounced in self-fertilising organisms,
524 due to the reduction in effective recombination rates. However, when investigating
525 sweep patterns over broad genetic regions, it becomes likelier that genetic diversity
526 will be affected by multiple beneficial mutations spreading throughout the genome.
527 Competing selective sweeps can lead to elevated diversity near a target locus for
528 two reasons. First, selection interference increases the fixation time of individual
529 mutations, allowing more recombination that can restore neutral diversity (Kim
530 and Stephan 2003). In addition, competing selective sweeps can drag different sets
531 of neutral variation to fixation, creating asymmetric diversity levels around a sub-
532 stitution (Chevin *et al.* 2008). Further investigations of selective sweep patterns
533 across long genetic distances will prove to be a rich area of future research.

534 Finally, we have assumed a fixed population size, and that sweeps from standing
535 variation arose from neutral variation. The resulting signatures could differ if the
536 population size has changed over time, or if the beneficial allele was previously
537 deleterious. Both issues could also affect our ability to discriminate between soft
538 and hard sweeps.

539 **Potential applications to self-fertilising organisms**

540 Existing methods for finding sweep signatures in nucleotide polymorphism data
541 are commonly based on finding regions with a site-frequency spectrum matching
542 what is expected under a selective sweep (Nielsen *et al.* 2005; Boitard *et al.* 2009;
543 Pavlidis *et al.* 2013; DeGiorgio *et al.* 2016; Huber *et al.* 2016). The more general
544 models developed here can be used to create more specific sweep-detection methods
545 that include self-fertilisation. However, a recent analysis found that soft-sweep
546 signatures can be incorrectly inferred if analysing genetic regions that flank hard
547 sweeps, which was named the ‘soft shoulder’ effect (Schridder *et al.* 2015). Due to
548 the reduction in recombination in selfers, these model results indicate that ‘soft-
549 shoulder’ footprints can arise over long genetic distances and should be taken into
550 account. One remedy to this problem is to not just classify genetic regions as being
551 subject to either a hard or soft sweep, but also as being linked to a region subject
552 to one of these sweeps (Schridder and Kern 2016). These more general calculations
553 can also be extended to quantify to what extent background selection and sweeps
554 jointly shape genome-wide diversity in self-fertilising organisms (Elyashiv *et al.*
555 2016; Campos *et al.* 2017; Booker and Keightley 2018; Rettelbach *et al.* 2019), or
556 detect patterns of introgression (Setter *et al.* 2019).

557 **Acknowledgments.** We would like to thank Sally Otto for providing infor-
558 mation on the elevated effective starting frequency of beneficial mutations; Brian
559 Charlesworth on providing advice on modelling selective sweeps, sharing unpub-
560 lished results, and providing comments on the manuscript; Ben Haller for answer-
561 ing questions about SLiM; Nick Barton and other anonymous referees for providing
562 feedback on the manuscript. MH was supported by a Marie Curie International

563 Outgoing Fellowship (MC-IOF-622936) and a NERC Independent Research Fel-
564 lowship (NE/R015686/1). MH and TB also acknowledge financial support from
565 the European Research Council under the European Union’s Seventh Framework
566 Program (FP7/20072013, ERC Grant 311341).

References

- Andersen, E. C., J. P. Gerke, J. A. Shapiro, J. R. Crissman, R. Ghosh, *et al.*,
2012 Chromosome-scale selective sweeps shape *Caenorhabditis elegans* genomic
diversity. *Nat. Genet.* **44**: 285–290.
- Anderson, T. J. C., S. Nair, M. McDew-White, I. H. Cheeseman, S. Nkhoma, *et al.*,
2016 Population parameters underlying an ongoing soft sweep in southeast asian
malaria parasites. *Mol. Biol. Evol.* **34**: 131–144.
- Badouin, H., P. Gladieux, J. Gouzy, S. Siguenza, G. Aguilera, *et al.*, 2017
Widespread selective sweeps throughout the genome of model plant pathogenic
fungi and identification of effector candidates. *Mol. Ecol.* **26**: 2041–2062.
- Barrett, R. D. H. and D. Schluter, 2008 Adaptation from standing genetic varia-
tion. *Trends Ecol. Evol.* **23**: 38–44.
- Barton, N. H., 1998 The effect of hitch-hiking on neutral genealogies. *Genet. Res.*
72: 123–133.
- Berg, J. J. and G. Coop, 2015 A coalescent model for a sweep of a unique standing
variant. *Genetics* **201**: 707–725.

- Billiard, S., M. López-Villavicencio, B. Devier, M. E. Hood, C. Fairhead, *et al.*, 2011 Having sex, yes, but with whom? Inferences from fungi on the evolution of anisogamy and mating types. *Biol. Rev. Camb. Philos. Soc.* **86**: 421–442.
- Boitard, S., C. Schlötterer, and A. Futschik, 2009 Detecting selective sweeps: A new approach based on hidden markov models. *Genetics* **181**: 1567–1578.
- Bonhomme, M., S. Boitard, H. San Clemente, B. Dumas, N. Young, *et al.*, 2015 Genomic signature of selective sweeps illuminates adaptation of *Medicago truncatula* to root-associated microorganisms. *Mol. Biol. Evol.* **32**: 2097–2110.
- Booker, T. R. and P. D. Keightley, 2018 Understanding the factors that shape patterns of nucleotide diversity in the house mouse genome. *Mol. Biol. Evol.* **35**: 2971–2988.
- Braverman, J. M., R. R. Hudson, N. L. Kaplan, C. H. Langley, and W. Stephan, 1995 The hitchhiking effect on the site frequency spectrum of DNA polymorphisms. *Genetics* **140**: 783–796.
- Caballero, A. and W. G. Hill, 1992 Effects of partial inbreeding on fixation rates and variation of mutant genes. *Genetics* **131**: 493–507.
- Campos, J. L. and B. Charlesworth, 2019 The effects on neutral variability of recurrent selective sweeps and background selection. *Genetics* **212**: 287–303.
- Campos, J. L., L. Zhao, and B. Charlesworth, 2017 Estimating the parameters of background selection and selective sweeps in drosophila in the presence of gene conversion. *Proc. Natl. Acad. Sci. USA* **114**: E4762–E4771.

- Charlesworth, B., 1992 Evolutionary rates in partially self-fertilizing species. *Am. Nat.* **140**: 126–148.
- Charlesworth, B. and D. Charlesworth, 2010 *Elements of Evolutionary Genetics*. Roberts & Company Publishers, Greenwood Village, Colo.
- Chevin, L.-M., S. Billiard, and F. Hospital, 2008 Hitchhiking both ways: Effect of two interfering selective sweeps on linked neutral variation. *Genetics* **180**: 301–316.
- DeGiorgio, M., C. D. Huber, M. J. Hubisz, I. Hellmann, and R. Nielsen, 2016 SweepFinder2: increased sensitivity, robustness and flexibility. *Bioinformatics* **32**: 1895–1897.
- Desai, M. M. and D. S. Fisher, 2007 Beneficial mutation-selection balance and the effect of linkage on positive selection. *Genetics* **176**: 1759–1798.
- Elyashiv, E., S. Sattath, T. T. Hu, A. Strutsosky, G. McVicker, *et al.*, 2016 A genomic map of the effects of linked selection in *Drosophila*. *PLoS Genet.* **12**: e1006130.
- Ewing, G. and J. Hermisson, 2010 MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics* **26**: 2064–2065.
- Ewing, G., J. Hermisson, P. Pfaffelhuber, and J. Rudolf, 2011 Selective sweeps for recessive alleles and for other modes of dominance. *J. Math. Biol.* **63**: 399–431.
- Fay, J. C. and C.-I. Wu, 2000 Hitchhiking Under Positive Darwinian Selection. *Genetics* **155**: 1405–1413.

- Ferrer-Admetlla, A., M. Liang, T. Korneliussen, and R. Nielsen, 2014 On detecting incomplete soft or hard selective sweeps using haplotype structure. *Mol. Biol. Evol.* **31**: 1275–1291.
- Fisher, R. A., 1922 On the dominance ratio. *Proc. R. Soc. Edinburgh* **42**: 321–341.
- Fujito, N. T., Y. Satta, T. Hayakawa, and N. Takahata, 2018 A new inference method for detecting an ongoing selective sweep. *Genes Genet. Syst.* **93**: 149–161.
- Fulgione, A., M. Koornneef, F. Roux, J. Hermisson, and A. M. Hancock, 2018 Madeiran *Arabidopsis thaliana* reveals ancient long-range colonization and clarifies demography in Eurasia. *Mol. Biol. Evol.* **35**: 564–574.
- Fustier, M. A., J. T. Brandenburg, S. Boitard, J. Lapeyronnie, L. E. Eguiarte, *et al.*, 2017 Signatures of local adaptation in lowland and highland teosintes from whole-genome sequencing of pooled samples. *Mol. Ecol.* **26**: 2738–2756.
- Garud, N. R., P. W. Messer, E. O. Buzbas, and D. A. Petrov, 2015 Recent Selective Sweeps in North American *Drosophila melanogaster* Show Signatures of Soft Sweeps. *PLoS Genet.* **11**: e1005004.
- Garud, N. R. and D. A. Petrov, 2016 Elevated linkage disequilibrium and signatures of soft sweeps are common in *Drosophila melanogaster*. *Genetics* **203**: 863–880.
- Glémin, S., 2012 Extinction and fixation times with dominance and inbreeding. *Theor. Popul. Biol.* **81**: 310–316.

- Glémin, S. and J. Ronfort, 2013 Adaptation and maladaptation in selfing and outcrossing species: New mutations versus standing variation. *Evolution* **67**: 225–240.
- Haldane, J. B. S., 1927 A mathematical theory of natural and artificial selection, part V: Selection and mutation. *Math. Proc. Cambridge Philos. Soc.* **23**: 838–844.
- Haller, B. C. and P. W. Messer, 2019 SLiM 3: Forward Genetic Simulations Beyond the Wright–Fisher Model. *Mol. Biol. Evol.* **36**: 632–637.
- Harris, A. M. and M. DeGiorgio, 2018 Identifying and classifying shared selective sweeps from multilocus data. bioRxiv p. 446005.
- Harris, A. M. and M. DeGiorgio, 2019 A likelihood approach for uncovering selective sweep signatures from haplotype data. bioRxiv .
- Harris, A. M., N. R. Garud, and M. DeGiorgio, 2018a Detection and classification of hard and soft sweeps from unphased genotypes by multilocus genotype identity. *Genetics* **210**: 1429–1452.
- Harris, R. B., A. Sackman, and J. D. Jensen, 2018b On the unfounded enthusiasm for soft selective sweeps II: Examining recent evidence from humans, flies, and viruses. *PLoS Genet.* **14**: e1007859.
- Hartfield, M., T. Bataillon, and S. Glémin, 2017 The evolutionary interplay between adaptation and self-fertilization. *Trends Genet.* **33**: 420–431.
- Hartfield, M. and S. Glémin, 2014 Hitchhiking of deleterious alleles and the cost of adaptation in partially selfing species. *Genetics* **196**: 281–293.

Hartfield, M. and S. Glémin, 2016 Limits to adaptation in partially selfing species.

Genetics **203**: 959–974.

Hedrick, P. W., 1980 Hitchhiking: A comparison of linkage and partial selection.

Genetics **94**: 791–808.

Hermisson, J. and P. S. Pennings, 2005 Soft sweeps: Molecular population genetics of adaptation from standing genetic variation. Genetics **169**: 2335–2352.

Hermisson, J. and P. S. Pennings, 2017 Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation.

Methods Ecol. Evol. **8**: 700–716.

Huber, C. D., M. DeGiorgio, I. Hellmann, and R. Nielsen, 2016 Detecting recent selective sweeps while controlling for mutation rate and background selection.

Mol. Ecol. **25**: 142–156.

Huber, C. D., M. Nordborg, J. Hermisson, and I. Hellmann, 2014 Keeping It Local: Evidence for Positive Selection in Swedish *Arabidopsis thaliana*. Mol. Biol. Evol. **31**: 3026–3039.

Igic, B. and J. R. Kohn, 2006 The distribution of plant mating systems: study bias against obligately outcrossing species. Evolution **60**: 1098–1103.

Innan, H. and Y. Kim, 2004 Pattern of polymorphism after strong artificial selection in a domestication event. Proc. Natl. Acad. Sci. USA **101**: 10667–10672.

Innan, H. and M. Nordborg, 2003 The extent of linkage disequilibrium and haplotype sharing around a polymorphic site. Genetics **165**: 437.

- Jarne, P. and J. R. Auld, 2006 Animals mix it up too: the distribution of self-fertilization among hermaphroditic animals. *Evolution* **60**: 1816–1824.
- Jensen, J. D., 2014 On the unfounded enthusiasm for soft selective sweeps. *Nat. Commun.* **5**.
- Kaplan, N. L., R. R. Hudson, and C. H. Langley, 1989 The “hitchhiking effect” revisited. *Genetics* **123**: 887–899.
- Karasov, T., P. W. Messer, and D. A. Petrov, 2010 Evidence that adaptation in *Drosophila* is not limited by mutation at single sites. *PLoS Genet.* **6**: e1000924.
- Kern, A. D. and D. R. Schrider, 2018 diploS/HIC: An updated approach to classifying selective sweeps. *G3* **8**: 1959–1970.
- Kim, Y. and R. Nielsen, 2004 Linkage disequilibrium as a signature of selective sweeps. *Genetics* **167**: 1513–1524.
- Kim, Y. and W. Stephan, 2002 Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics* **160**: 765–777.
- Kim, Y. and W. Stephan, 2003 Selective sweeps in the presence of interference among partially linked loci. *Genetics* **164**: 389–398.
- Laporte, V. and B. Charlesworth, 2002 Effective population size and population subdivision in demographically structured populations. *Genetics* **162**: 501–519.
- Long, Q., F. A. Rabanal, D. Meng, C. D. Huber, A. Farlow, *et al.*, 2013 Massive genomic variation and strong selection in *Arabidopsis thaliana* lines from Sweden. *Nat. Genet.* **45**: 884–890.

- Martin, G. and A. Lambert, 2015 A simple, semi-deterministic approximation to the distribution of selective sweeps in large populations. *Theor. Popul. Biol.* **101**: 40–46.
- Maynard Smith, J., 1976 What determines the rate of evolution? *Am. Nat.* **110**: 331–338.
- Maynard Smith, J. and J. Haigh, 1974 The hitch-hiking effect of a favourable gene. *Genet. Res.* **23**: 23–35.
- McVean, G. A. T., 2007 The structure of linkage disequilibrium around a selective sweep. *Genetics* **175**: 1395–1406.
- Messer, P. W. and D. A. Petrov, 2013 Population genomics of rapid adaptation by soft selective sweeps. *Trends Ecol. Evol.* **28**: 659–669.
- Nielsen, R., 2005 Molecular signals of natural selection. *Annu. Rev. Genet.* **39**: 197–218.
- Nielsen, R., S. Williamson, Y. Kim, M. J. Hubisz, A. G. Clark, *et al.*, 2005 Genomic scans for selective sweeps using SNP data. *Genome Res.* **15**: 1566–1575.
- Nordborg, M., 2000 Linkage disequilibrium, gene trees and selfing: An ancestral recombination graph with partial self-fertilization. *Genetics* **154**: 923–929.
- Nordborg, M., B. Charlesworth, and D. Charlesworth, 1996 Increased levels of polymorphism surrounding selectively maintained sites in highly selfing species. *Proc. R. Soc. B* **263**: 1033–1039.
- Nordborg, M. and P. Donnelly, 1997 The coalescent process with selfing. *Genetics* **146**: 1185–1195.

- Orr, H. A. and A. J. Betancourt, 2001 Haldane's sieve and adaptation from the standing genetic variation. *Genetics* **157**: 875–884.
- Pavlidis, P., D. Živković, A. Stamatakis, and N. Alachiotis, 2013 SweeD: Likelihood-Based Detection of Selective Sweeps in Thousands of Genomes. *Mol. Biol. Evol.* **30**: 2224–2234.
- Pennings, P. S. and J. Hermisson, 2006a Soft Sweeps II – Molecular Population Genetics of Adaptation from Recurrent Mutation or Migration. *Mol. Biol. Evol.* **23**: 1076–1084.
- Pennings, P. S. and J. Hermisson, 2006b Soft Sweeps III: The Signature of Positive Selection from Recurrent Mutation. *PLoS Genet.* **2**: e186.
- Pennings, P. S., S. Kryazhimskiy, and J. Wakeley, 2014 Loss and Recovery of Genetic Diversity in Adapting Populations of HIV. *PLoS Genet.* **10**: e1004000.
- Peter, B. M., E. Huerta-Sanchez, and R. Nielsen, 2012 Distinguishing between selective sweeps from standing variation and from a *De Novo* mutation. *PLoS Genet.* **8**: e1003011.
- Pollak, E., 1987 On the theory of partially inbreeding finite populations. I. Partial selfing. *Genetics* **117**: 353–360.
- Price, N., B. T. Moyers, L. Lopez, J. R. Lasky, J. G. Monroe, *et al.*, 2018 Combining population genomics and fitness QTLs to identify the genetics of local adaptation in *Arabidopsis thaliana*. *Proc. Natl. Acad. Sci. USA* **115**: 5028–5033.
- Pritchard, J. K. and A. Di Rienzo, 2010 Adaptation - not by sweeps alone. *Nat. Rev. Genet.* **11**: 665–667.

- Przeworski, M., G. Coop, and J. D. Wall, 2005 The signature of positive selection on standing genetic variation. *Evolution* **59**: 2312–2323.
- Rettelbach, A., A. Nater, and H. Ellegren, 2019 How linked selection shapes the diversity landscape in *Ficedula* flycatchers. *Genetics* **212**: 277–285.
- Roze, D., 2009 Diploidy, population structure, and the evolution of recombination. *Am. Nat.* **174**: S79–S94.
- Roze, D., 2016 Background selection in partially selfing populations. *Genetics* **203**: 937–957.
- Sabeti, P. C., D. E. Reich, J. M. Higgins, H. Z. P. Levine, D. J. Richter, *et al.*, 2002 Detecting recent positive selection in the human genome from haplotype structure. *Nature* **419**: 832–837.
- Schoen, D. J., M. T. Morgan, and T. Bataillon, 1996 How Does Self-Pollination Evolve? Inferences from Floral Ecology and Molecular Genetic Variation. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **351**: 1281–1290.
- Schrider, D. R. and A. D. Kern, 2016 S/HIC: Robust identification of soft and hard sweeps using machine learning. *PLoS Genet.* **12**: e1005928.
- Schrider, D. R. and A. D. Kern, 2017 Soft sweeps are the dominant mode of adaptation in the human genome. *Mol. Biol. Evol.* **34**: 1863–1877.
- Schrider, D. R., F. K. Mendes, M. W. Hahn, and A. D. Kern, 2015 Soft shoulders ahead: Spurious signatures of soft and partial selective sweeps result from linked hard sweeps. *Genetics* **200**: 267–284.

- Setter, D., S. Mousset, X. Cheng, R. Nielsen, M. DeGiorgio, *et al.*, 2019 Volcanofinder: genomic scans for adaptive introgression. bioRxiv p. 697987.
- Sheehan, S. and Y. S. Song, 2016 Deep learning for population genetic inference. PLoS Comput. Biol. **12**: e1004845.
- Stephan, W., 2016 Signatures of positive selection: from selective sweeps at individual loci to subtle allele frequency changes in polygenic adaptation. Mol. Ecol. **25**: 79–88.
- Stephan, W., 2019 Selective sweeps. Genetics **211**: 5–13.
- Tajima, F., 1983 Evolutionary Relationship of DNA Sequences in Finite Populations. Genetics **105**: 437–460.
- Teshima, K. M., G. Coop, and M. Przeworski, 2006 How reliable are empirical genomic scans for selective sweeps? Genome Res. **16**: 702–712.
- Teshima, K. M. and M. Przeworski, 2006 Directional positive selection on an allele of arbitrary dominance. Genetics **172**: 713–718.
- Thomson, G., 1977 The effect of a selected locus on linked neutral loci. Genetics **85**: 753–788.
- Uecker, H., 2017 Evolutionary rescue in randomly mating, selfing, and clonal populations. Evolution **71**: 845–858.
- Van Herwaarden, O. A. and N. J. Van der Wal, 2002 Extinction time and age of an allele in a large finite population. Theor. Popul. Biol. **61**: 311–318.

- Vatsiou, A. I., E. Bazin, and O. E. Gaggiotti, 2016 Detection of selective sweeps in structured populations: a comparison of recent methods. *Mol. Ecol.* **25**: 89–103.
- Vitti, J. J., S. R. Grossman, and P. C. Sabeti, 2013 Detecting natural selection in genomic data. *Annu. Rev. Genet.* **47**: 97–120.
- Voight, B. F., S. Kudaravalli, X. Wen, and J. K. Pritchard, 2006 A map of recent positive selection in the human genome. *PLoS Biol.* **4**: e72.
- Vy, H. M. T., Y.-J. Won, and Y. Kim, 2017 Multiple Modes of Positive Selection Shaping the Patterns of Incomplete Selective Sweeps over African Populations of *Drosophila melanogaster*. *Mol. Biol. Evol.* **34**: 2792–2807.
- Wakeley, J., 2009 *Coalescent theory: an introduction*, volume 1. Roberts & Company Publishers, Greenwood Village, Colorado.
- Williams, K.-A. and P. S. Pennings, 2019 Drug resistance evolution in HIV in the late 1990s: hard sweeps, soft sweeps, clonal interference and the accumulation of drug resistance mutations. bioRxiv p. 548198.
- Wilson, B. A., P. S. Pennings, and D. A. Petrov, 2017 Soft selective sweeps in evolutionary rescue. *Genetics* **205**: 1573–1586.
- Wilson, B. A., D. A. Petrov, and P. W. Messer, 2014 Soft Selective Sweeps in Complex Demographic Scenarios. *Genetics* **198**: 669–684.
- Wright, S., 1951 The genetical structure of populations. *Ann. Eugen.* **15**: 323–354.
- Xue, A. T., D. R. Schrider, A. D. Kern, and Ag1000G Consortium, 2019 Discovery of ongoing selective sweeps within *Anopheles* mosquito populations using deep learning. bioRxiv p. 589069.

Yang, Z., J. Li, T. Wiehe, and H. Li, 2018 Detecting recent positive selection with a single locus test bipartitioning the coalescent tree. *Genetics* **208**: 791–805.

Zhong, L., Q. Yang, X. Yan, C. Yu, L. Su, *et al.*, 2017 Signatures of soft sweeps across the Dt1 locus underlying determinate growth habit in soya bean [*Glycine max* (L.) Merr.]. *Mol. Ecol.* **26**: 4686–4699.