

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

A major research goal in evolutionary genetics is to uncover loci experiencing positive selection. One approach involves finding ‘selective sweeps’, either formed by *de novo* mutation, or ‘soft sweeps’ arising from recurrent mutation or existing standing variation. Existing theory generally assumes outcrossing populations, and it is unclear how dominance affects soft sweeps. We consider how arbitrary dominance and inbreeding via self-fertilisation affect hard and soft sweep signatures. With increased self-fertilisation, they are maintained over longer map distances due to reduced effective recombination and faster beneficial allele fixation times. Dominance can affect sweep patterns in outcrossers if the derived variant originates from either a single novel allele, or from recurrent mutation. These models highlight the challenges in distinguishing hard and 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
 35 common ancestor of a beneficial allele appeared before it became selected for (re-
 36 viewed by Barrett and Schluter (2008); Messer and Petrov (2013); Hermisson and
 37 Pennings (2017)). Soft sweeps can originate from beneficial mutations being intro-
 38 duced by recurrent mutation at the target locus (Pennings and Hermisson 2006a,b),

39 or originating from existing standing variation that was either neutral or deleterious (Orr and Betancourt 2001; Innan and Kim 2004; Przeworski *et al.* 2005; 40 Hermisson and Pennings 2005; Wilson *et al.* 2014; Berg and Coop 2015; Wilson 41 *et al.* 2017). A key property of soft sweeps is that the beneficial variant is present 42 on multiple genetic backgrounds as it sweeps to fixation, so different haplotypes 43 may be present around the derived allele. This property is often used to detect 44 soft sweeps in genetic data (Peter *et al.* 2012; Vitti *et al.* 2013; Garud *et al.* 2015; 45 Garud and Petrov 2016; Schrider and Kern 2016; Sheehan and Song 2016; Harris 46 *et al.* 2018a; Kern and Schrider 2018; Harris and DeGiorgio 2018, 2019). Soft 47 sweeps have been inferred in *Drosophila* (Karasov *et al.* 2010; Garud *et al.* 2015; 48 Garud and Petrov 2016; Vy *et al.* 2017), humans (Peter *et al.* 2012; Schrider and 49 Kern 2017), maize (Fustier *et al.* 2017), cattle (Qanbari *et al.* 2014) and pathogens 50 including *Plasmodium falciparum* (Anderson *et al.* 2016) and HIV (Pennings *et al.* 51 2014; Williams and Pennings 2019). Yet determining how extensive soft sweeps 52 are in 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). Detailed analyses of these cases
 95 has been hampered by a lack of theory on how hard and soft sweep signatures
 96 should manifest themselves under different self-fertilisation and dominance levels.
 97 Previous studies have only focussed on special cases; Hedrick (1980) analysed link-
 98 age disequilibrium caused by a hard sweep under self-fertilisation, while Schoen
 99 *et al.* (1996) modelled sweep patterns caused by modifiers that altered the mating
 100 system in different ways.

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

Results

Model Outline

We consider a diploid population of size N (carrying $2N$ haplotypes in total). Individuals reproduce by self-fertilisation with probability σ , and outcross with probability $1 - \sigma$. There are two biallelic loci A , B with a recombination rate r between them. Locus A represents a region where neutral polymorphism accumulates under an infinite-sites model (Kimura 1971). Locus B determines fitness differences, carrying an allele that initially segregates neutrally for a period of time. After a period of time the allele becomes advantageous with selective advantage $1 + hs$ when heterozygous and $1 + s$ when homozygous, with $0 < h < 1$ and $s > 0$. We further assume that selection is strong (i.e., $N_e hs \gg 1$) so the sweep trajectory can be modelled deterministically. p_0 is the frequency at which the allele becomes selected for. Table 1 lists the notation used in the analysis.

Our overall goal is to determine how the spread of an adaptive allele at locus B affects genealogies underlying polymorphism at locus A , by considering whether neutral alleles at A are associated with the selected derived allele or ancestral neutral allele at locus B . A schematic is shown in Figure 1. We follow the approach of Berg and Coop (2015) and, looking backwards in time, break down the allele history of B into two phases. The first phase (the ‘sweep phase’) considers the derived allele at B being selectively favoured from an initial frequency p_0 and spreading through the population. The second phase (the ‘standing phase’) assumes that the derived allele is present at an frequency with mean p_0 . During both phases, a pair of haplotypes can either coalesce, or one of them recombines onto the ancestral background.

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
A, B	Loci carrying neutral, selected alleles
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
$P_c(t)$	Probability of coalescence at time t
$P_r(t)$	Probability of recombination at time t
P_{NE}	Probability that neutral marker does not coalesce or recombine during sweep phase
$P_{NE,SL}$	P_{NE} using 'star-like' approximation (no coalescence 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
\tilde{s}	'Effective' selection coefficient to map hard sweep onto standing variation cases
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.

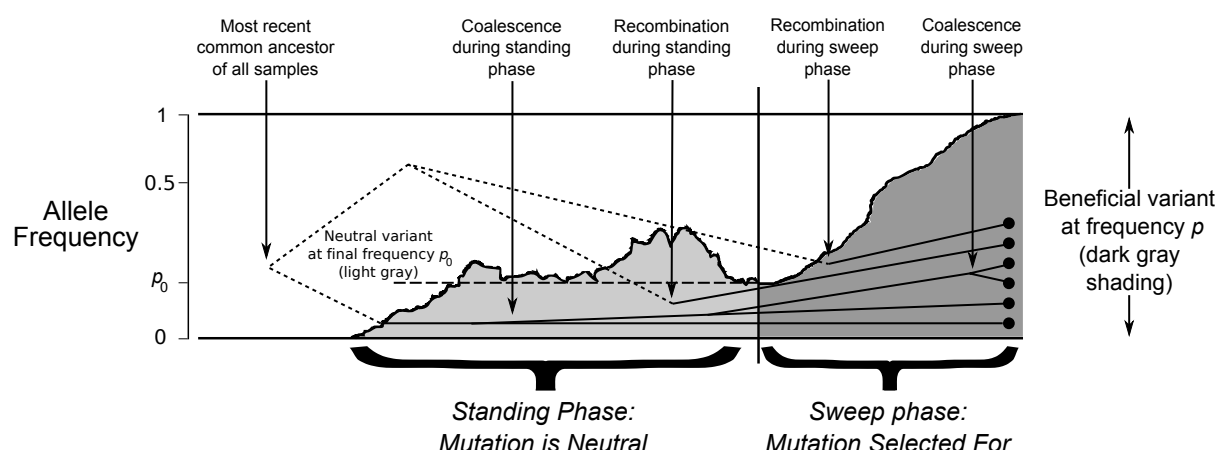


Figure 1. A schematic of the model. The history of the derived variant is separated into two phases. The ‘standing phase’ (shown in light gray), is when the derived variant is segregating at a frequency close to p_0 for a long period of time. The ‘sweep phase’ (shown in dark gray) is when the variant becomes selected for and starts increasing in frequency. Scale on the left-hand side show allele frequency on an arbitrary log-scale. Dots on the right-hand side represent a sample of haplotypes taken at the present day, with lines representing their genetic histories. Samples can either coalesce or recombine onto the ancestral background during either phase. Solid lines represent coalescent histories for the derived genetic background; dotted lines represent coalescent histories for the ancestral background.

For tightly linked loci (r close to 0), the relatively rapid fixation time of the derived variant makes it unlikely for a given neutral variant to be present on different backgrounds (with respect to the selected locus), reducing neutral diversity. Further from the target locus, recombination can transfer allele copies at A away from the selected background to the ancestral background, so diversity reaches the initial level. Self-fertilisation creates two key differences compared to outcrossing. First, the adaptive allele trajectory, which underlies expected diversity patterns, is affected by the levels of self-fertilisation (σ) and dominance (h). Second, the effective population size (which determines the coalescence rate) and recombina-

tion frequency are scaled by factors $1/(1 + F)$ and $1 - 2F + \Phi$ respectively, for $F = \sigma/(2 - \sigma)$ the inbreeding coefficient (Wright 1951; Pollak 1987; Charlesworth 1992; Caballero and Hill 1992; Nordborg *et al.* 1996; Nordborg and Donnelly 1997; Nordborg 2000) and Φ the joint probability of identity-by-descent at the two loci (Roze 2009, 2016; Hartfield and Glémin 2016). At equilibrium, it equals:

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

Note that $1 - 2F + \Phi$ approximates to $1 - F$ (as $\Phi \approx F$) for most combinations of recombination and selfing fractions, unless σ is close to one and r is high. The N_e scaling factor $1/(1 + F)$ can also be a good approximation if there is non-Poisson variation in offspring, unless female fitness strongly affects reproduction number (Laporte and Charlesworth 2002). Although we focus on inbreeding via self-fertilisation, the scalings $N_e = N/(1 + F)$ and $r_e \approx r(1 - F)$ should also hold under other systems of regular inbreeding (Caballero and Hill 1992; Charlesworth and Charlesworth 2010, Box 8.4).

We will outline how both coalescence and recombination act during both of these phases, and use these calculations to determine selective sweep properties. Although previous models tended to exclude coalescence during the sweep phase, including it is important for producing accurate matches with simulation results (Barton 1998; Charlesworth, in prep.).

Simulation Procedures

Throughout, analytical solutions are compared to results from Wright-Fisher forward-in-time stochastic simulations. Simulations were ran using SLiM version 3.3 (Haller

and Messer 2019), with simulation scripts available from GitHub (<https://github.com/MattHartfield/SweepDomSelf>). There exists N diploid individuals of length 100,000 nucleotides. The target locus carrying the derived allele is present at the left-hand end of the haplotype, while the remaining loci can carry neutral mutations only. Selected alleles have a homozygous selective advantage s and dominance coefficient h . Mutation and recombination parameters are input as population-level rates, which are subsequently scaled down to obtain per-locus mutation rates, or a per-inter-base-pair recombination probability.

A ‘burn-in’ phase is first run to generate background neutral diversity, where the population evolves without any beneficial alleles present for $20N$ generations; this population was subsequently saved. The second phase acts differently depending on whether the beneficial allele is instantly selected for (a hard sweep), or whether it went through a neutral phase. If a hard sweep was simulated, then the beneficial allele was introduced into a single individual as a heterozygote, and tracked until it is fixed or lost. If the latter, the burn-in population is reloaded, the random seed changed and the beneficial mutation reintroduced. The procedure is repeated until the mutation has fixed.

If the derived mutation was initially neutral, then following the burn-in a neutral allele is introduced into a random individual as a heterozygote, and tracked until it is lost or it reaches a frequency p_0 . If it is lost then the burn-in population is reloaded, the random seed changed and the derived allele is reintroduced. If the mutation reaches the target frequency p_0 then it is then converted into a selected mutation, and is tracked until fixation or loss. If the beneficial mutation is subsequently lost then the simulation is stopped and restarted from scratch (i.e., the burn-in population is also regenerated).

188 100 burn-in populations were generated for each parameter set. After the
189 beneficial allele has gone to fixation, we sampled 10 haplotypes 10 times from each
190 burn-in population to create 1,000 replicate simulations. Mutations are placed
191 in one of 10 bins depending on the distance from the sweep. Relevant statistics
192 (pairwise diversity, relative to neutral expectations; number of segregating sites;
193 site frequency spectrum) were calculated per bin. Mean values are calculated over
194 all 1,000 outputs. 95% confidence intervals were calculated by bootstrapping the
195 data 1,000 times.

196 **Data Availability.** File S1 is a *Mathematica* notebook of analytical deriva-
197 tions and simulation results. File S2 contains additional results and figures. File
198 S3 contains copies of the simulation scripts, which are also available from [https:](https://github.com/MattHartfield/SweepDomSelf)
199 [//github.com/MattHartfield/SweepDomSelf](https://github.com/MattHartfield/SweepDomSelf). Supplemental material has also
200 been uploaded to Figshare.

201 Probability of events during sweep phase

202 We first look at the probability of events (coalescence or recombination) acting
203 throughout the sweep phase for a pair of alleles. Looking back in time following
204 a sweep, sites linked to the beneficial allele can either coalesce or recombine onto
205 the ancestral genetic background. Let $p(t)$ be the adaptive mutation frequency
206 at time t , defined as the number of generations prior to the present day. Further
207 define $p(0) = 1$ (i.e., the allele is fixed at the present day), and τ_{p_0} the time in the
208 past when the derived variant became beneficial (i.e., $p(\tau_{p_0}) = p_0$).

209 For a pair of haplotype samples carrying the derived allele, if it is at frequency
210 $p(t)$ at time t , this lineage pair can either coalesce or one of the haplotypes recom-

bine onto the ancestral background. Each event occurs with probability:

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

Equation 2 is based on those obtained by Kaplan *et al.* (1989), assuming that N_e is reduced by a factor $1+F$ due to self-fertilisation (Pollak 1987; Charlesworth 1992; Caballero and Hill 1992; Nordborg and Donnelly 1997), and $r_{eff} = r(1-2F+\Phi)$ is the ‘effective’ recombination rate after correcting for increased homozygosity due to self-fertilisation (Nordborg *et al.* 1996; Nordborg 2000; Charlesworth and Charlesworth 2010; Roze 2009, 2016; Hartfield and Glémin 2016).

We are interested in calculating (i) the probability P_{NE} that no coalescence or recombination occurs; (ii) the probability $P_{R,Sw}$ that recombination acts on a lineage to transfer it to the neutral background carrying the ancestral allele, assuming that no more than one recombination event occurs per generation (see Campos and Charlesworth (2019) for derivations assuming multiple recombination events). We will go through these probabilities in turn to determine expected pairwise diversity. For P_{NE} , the total probability that the two lineages do not coalesce or recombine over τ_{p_0} generations equals:

$$\begin{aligned}
 P_{NE} &= \prod_{t=0}^{\tau_{p0}} [1 - P_c(t) - P_r(t)] \\
 &\approx \exp \left(- \int_{t=0}^{\tau_{p0}} [P_c(t) + P_r(t)] dt \right) && \text{assuming } P_c, P_r \ll 1 \\
 &\approx \exp \left(- \int_{t=0}^{\tau_{p0}} \left[\frac{1+F}{2Np(t)} + 2r(1-2F+\Phi)(1-p(t)) \right] dt \right) \\
 &\approx \exp \left(- \int_{p=1-\epsilon}^{p_0} \left[\frac{\frac{1+F}{2Np} + 2r(1-2F+\Phi)(1-p)}{dp/dt} \right] dp \right) && \text{taking the integral over } p
 \end{aligned}
 \tag{3}$$

Here ϵ is a small term and $1 - \epsilon$ is the upper limit of the deterministic spread of the beneficial allele. We will discuss in the section ‘Effective starting frequency from a de novo mutation’ what a reasonable value for ϵ should be. We can calculate P_{NE} for a general self-fertilisation level if the selection coefficient is not too weak (i.e., $1/N_e \ll s \ll 1$). Here the rate of change of the allele frequency is given by (Glémin 2012):

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

Note the negative factor in Equation 4 since we are looking back in time. By substituting Equation 4 into Equation 3, we obtain an analytical solution for P_{NE} , although the resulting expression is complicated (Section A of Supplementary File S1).

To calculate $P_{R,Sw}$, the probability that recombination acts during the sweep, we first calculate the probability that recombination occurs when the beneficial allele is at frequency p' . Here, no events occur in the time leading up to p' , then

a recombination event occurs with probability $2r(1 - 2F + \Phi)(1 - p')$. $P_{R,Sw}$ is obtained by summing this probability over the entire sweep from time 0 to τ_{p_0} , which can be approximated in continuous time by integration:

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

where:

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

dp/dt is Equation 4 but with p' instead of p . We evaluate Equation 5 numerically.

‘Star-like’ approximation (no coalescence during the sweep phase)

It is possible to obtain more tractable analytical solutions by using a ‘separation-of-timescales’ approximation, and assume that no coalescence occurs during the sweep phase (Pennings and Hermisson 2006b; Berg and Coop 2015). Here we only have to calculate the probability that no recombination occurs during the sweep phase, which for a single lineage equals:

$$\begin{aligned} P_{NE,SL} &= \exp \left(- \int_{p=1}^{p_0} \frac{r_{eff}(1 - p)}{dp/dt} dp \right) \\ &= \exp \left(- \frac{r_{eff}}{H_l s} \log \left[\frac{H_l}{H_h} \left(\frac{1}{p_0} + 1 \right) - 1 \right] \right) \\ &= \left[\frac{H_l}{H_h} \left(\frac{1}{p_0} + 1 \right) - 1 \right]^{-r_{eff}/(H_l s)} \end{aligned} \quad (6)$$

Here, $H_l = F + h - Fh$, $H_h = 1 - h + Fh$ are the ‘effective’ dominance coefficients when the beneficial variant is at a low or high frequency (Glémin 2012). Note that for the special case $\sigma = 0$ and $h = 1/2$, $H_l = H_h = 1/2$ and Equation 6 reduces to $(1/p_0)^{-(2r/s)}$, which is equivalent to Equation 2 of Berg and Coop (2015) after scaling the selection coefficient by $1/2$ to include semidominance.

Probability of coalescence from standing variation

When the variant becomes advantageous at frequency p_0 , we expect $\sim 2Np_0$ haplotypes to carry it. We assume that p_0 , and hence event probabilities, remain invariant over time. Berg and Coop (2015) have shown this assumption provides a good approximation to coalescent rates during the standing phase. The outcome during the standing phase is thus determined by competing Poisson processes. The two haplotypes could coalesce, with a waiting time being exponentially distributed with rate $(1 + F)/(2Np_0)$. Alternatively, one of the two haplotypes could recombine onto the ancestral background with mean waiting time $2r_{eff}(1 - p_0)$. For two competing exponential distributions with rates λ_1 and λ_2 , the probability of the first event occurring given an event happens equals $\lambda_1/(\lambda_1 + \lambda_2)$ (Wakeley 2009). Hence the probability that recombination occurs instead of coalescence equals:

$$P_{R,Sd} = \frac{2r_{eff}(1 - p_0)}{\frac{1+F}{2Np_0} + 2r_{eff}(1 - p_0)} = \frac{2R(1 - 2F + \Phi)p_0(1 - p_0)/(1 + F)}{1 + 2R(1 - 2F + \Phi)p_0(1 - p_0)/(1 + F)} \quad (7)$$

The probability of coalescence rather than recombination is $1 - P_{R,Sd}$. Here $R = 2Nr$ is the population-scaled recombination rate. The $(1 - 2F + \Phi)/(1 + F)$ term in Equation 7 approximates to $1 - \sigma$ if $\Phi \approx F$; this approximation holds unless

both self-fertilisation and recombination are high. This term reflects how increased homozygosity reduces both effective recombination and N_e ; with the latter making coalescence more likely. Sweeps from standing variation are characterised by recombination during this standing phase, so different background haplotypes go to fixation alongside the selected mutation. Equation 7 shows that for a fixed recombination rate R , haplotypes are more likely to coalesce with increased self-fertilisation, limiting the creation of different background haplotypes. However, the same coalescent probability can be recovered by increasing the recombination distance by a factor $\sim 1/(1 - \sigma)$, so the various background haplotypes can be captured by a genetic scan if a longer genetic region is analysed.

Effective starting frequency for a de novo mutation

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

In Section B of Supplementary File S1, we determine that hard sweeps that go

to fixation have the following effective starting frequency:

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

where $H_l = F + h - Fh$ is the effective dominance coefficient for mutations at a low frequency. This result is consistent with those of Martin and Lambert (2015), who obtained a distribution of effective starting frequencies using stochastic differential equations. This acceleration effect can create substantial increases in the apparent p_0 , especially for recessive mutations. For example, with $N = 5,000$, $s = 0.05$, $h = 0.1$ and $F = 0$, $p_{0,A} = 0.01$, an 100-fold increase above $p_0 = 1/(2N) = 0.0001$.

Effective final frequency: The effective final frequency of the derived allele, at which its spread is no longer deterministic, can be obtained by changing H_l to $H_h = 1 - h + Fh$ in Equation 8. Van Herwaarden and Van der Wal (2002) determined that the sojourn time for an allele with dominance coefficient h that is increasing in frequency, is the same for an allele decreasing in frequency with dominance $1 - h$. Glémin (2012) showed that this result also holds under any inbreeding value F (see also Charlesworth, in prep.).

Expected Pairwise Diversity

We use P_{NE} , $P_{R,sw}$ and $P_{R,sd}$ to calculate the expected pairwise diversity (denoted π) present around a sweep. During the sweep phase, then the two neutral sites could either coalesce, or one of them recombines onto the ancestral background. If coalescence occurs, then since it occurred in the recent past then no diversity exist between samples (this assumption is later relaxed when calculating the site-frequency spectrum). Alternatively, if one of the two samples recombines onto the

neutral background, they will have the same pairwise diversity between them as the background population (π_0). If the two samples trace back to the standing phase (with probability P_{NE}) then the same logic applies. Hence the expected diversity following a sweep equals:

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

The full solution to Equation 9 can be obtained by plugging in the relevant parts from Equations 3, 5 and 7, which we evaluate numerically. An analytical approximation can be obtained by using the ‘star-like’ result for P_{NE} (Equation 6). In this case we are interested in calculating the probability of coalescence during the standing phase $P_{C,sd} = 1 - P_{R,sd}$, and the expected pairwise diversity approximates to:

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

Equation 10 reflects similar formulas for diversity following soft sweeps in haploid outcrossing populations (Pennings and Hermisson 2006b; Berg and Coop 2015). There is a factor of two in the power term to account for two lineages. Note that both Equations 9 and 10 are undefined for $h = 0$ or 1 with $\sigma = 0$; these cases can be derived separately.

Figure 2 plots Equation 9 with different dominance, self-fertilisation, and stand-

ing frequency values. The analytical solution fits well compared to forward-in-time simulations, yet slightly overestimates them for high self-fertilisation frequencies. Under complete outcrossing, baseline diversity is restored (i.e., $\mathbb{E}(\pi/\pi_0)$ goes to 1) closer to the sweep origin for recessive mutations ($h = 0.1$), compared to semidominant ($h = 0.5$) or dominant ($h = 0.9$) mutations. Dominant and semidominant mutations produce similar reductions in genetic diversity, so these cases may be hard to differentiate from diversity data alone.

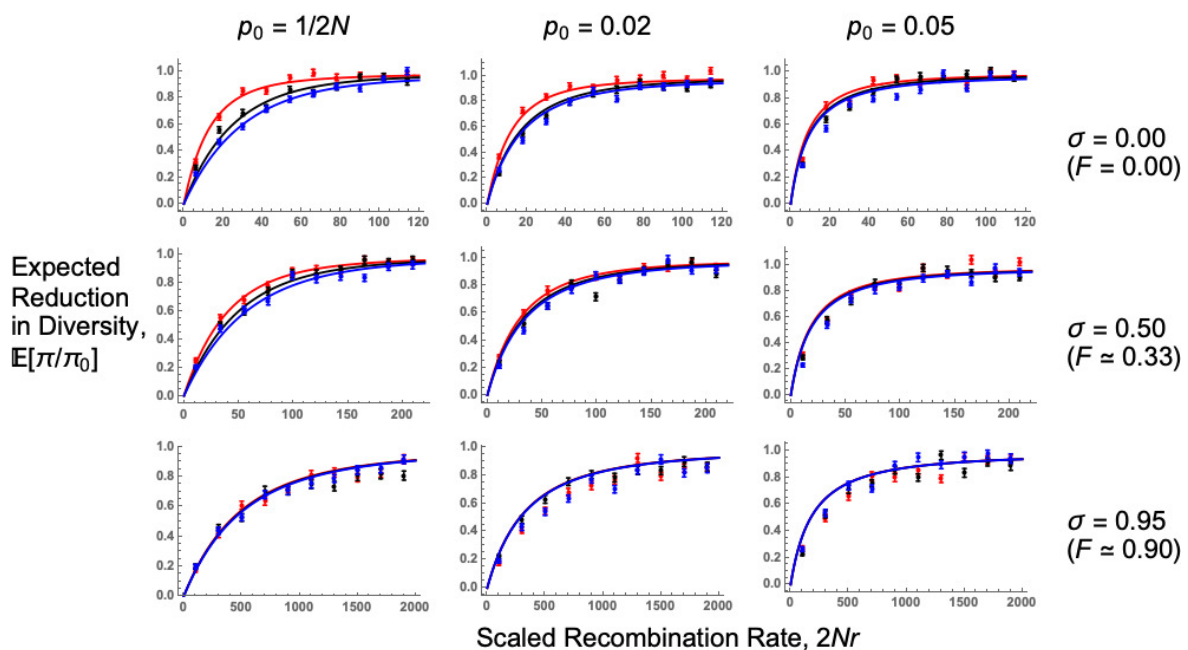


Figure 2. Expected pairwise diversity following a selective sweep. Plots of $\mathbb{E}(\pi/\pi_0)$ as a function of the recombination rate scaled to population size $2Nr$. Lines are analytical solutions (Equation 9), points are forward-in-time simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$ (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). 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.

These patterns can be understood by examining the underlying allele trajectories, using logic described by Teshima and Przeworski (2006) (Figure 3). For outcrossing populations, recessive mutations spend most of the sojourn time at low frequencies, maximising recombination events and restoring neutral variation. These trajectories mimic sweeps from standing variation, which spend extended periods of time at low frequencies in the standing phase. Conversely, dominant mutations spend most of their time at high frequencies, reducing the chance for neutral markers to recombine onto the ancestral background.

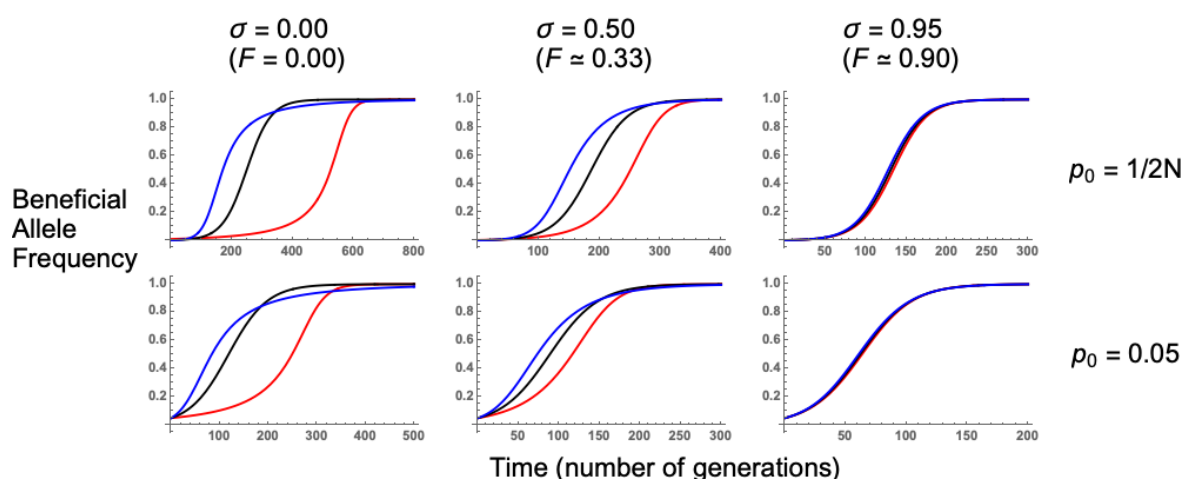


Figure 3. Beneficial allele trajectories. These were obtained by numerically evaluating the negative of Equation 4 forward in time. $N = 5,000$, $s = 0.05$, and h equals either 0.1 (red lines), 0.5 (black lines), or 0.9 (blue lines). Values of p_0 and self-fertilisation rates σ used are shown for the 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.

As self-fertilisation increases, sweep signatures become similar to the co-dominant case as the derived allele is more likely to spread as a homozygote, weakening the influence that dominance exerts over beneficial allele trajectories. Increasing p_0 also causes sweeps with different dominance coefficients to produce comparable

signatures, as beneficial mutation trajectories become similar after conditioning on starting at an elevated frequency.

In Supplementary File S2, we show that the star-like approximation (Equation 10) systematically predicts higher diversity levels compared to the full solution (Equation 9). This is because Equation 10 assumes that no coalescence occurs during the sweep phase, which is only accurate for very strongly selected mutations (Barton 1998; Charlesworth, in prep.). We also compare forward-in-time simulations to coalescent simulations for the outcrossing case. Both methods yield similar results, although forward-in-time simulations produce slightly higher diversity estimates compared to coalescent simulations.

Site Frequency Spectrum

The star-like approximation can be used to obtain analytical solutions for the number of segregating sites and the site frequency spectrum (i.e., the probability that $l = 1, 2 \dots n - 1$ of n alleles carry derived variants). The full derivation for these statistics are outlined in Supplementary File S2. Figure 4 plots the SFS (Equation A12 in Supplementary File S2) alongside simulation results. Analytical results fit the simulation data well after including an inflated singleton class to account for new mutations that occur during the sweep phase (Berg and Coop 2015). There is a tendency for analytical results to underestimate the proportion of low- and high-frequency classes ($l = 1$ and 9 in Figure 4), and overestimate the proportion of intermediate-frequency classes. Hard sweeps in either outcrossers or partial selfers are characterised by a large number of singletons or highly-derived variants (Figure 4), which is a typical selective sweep signature (Braverman *et al.*

1995; Barton 1998; Kim and Stephan 2002). As the initial frequency p_0 increases, so does the number of intermediate-frequency variants (Figure 4). This signature is often seen as a characteristic of soft sweeps (Pennings and Hermisson 2006b; Berg and Coop 2015). Recessive hard sweeps ($h = 0.1$ and $p_0 = 1/2N$) can produce SFS profiles that are similar to sweeps from standing variation, as there are an increased number of recombination events occurring since the allele is at a low frequency for long time periods (Figure 3). With increased self-fertilisation, both hard and soft sweep signatures (e.g., increased number of intermediate-frequency alleles) are recovered when measuring the SFS at a longer recombination distance than in outcrossers (Figure 4, bottom row).

Soft sweeps from recurrent mutation

So far, we have only focussed on a soft sweep that arises from standing variation. An alternative type of soft sweep is one where recurrent mutation at the selected locus introduces the beneficial allele onto different genetic backgrounds. We can examine this case by modifying existing results. During the sweep phase, markers linked to the derived background can not only change state by recombination, but also via mutation. If the derived allele is at frequency p then the probability of a mutation event is $2\mu_b(1 - p)/p$, for μ_b the mutation probability (Pennings and Hermisson 2006b). In this case the expected reduction in diversity now equals:

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

where $P_{R,sw}$, P_{NE} are modified to include mutations arising during the sweep phase:

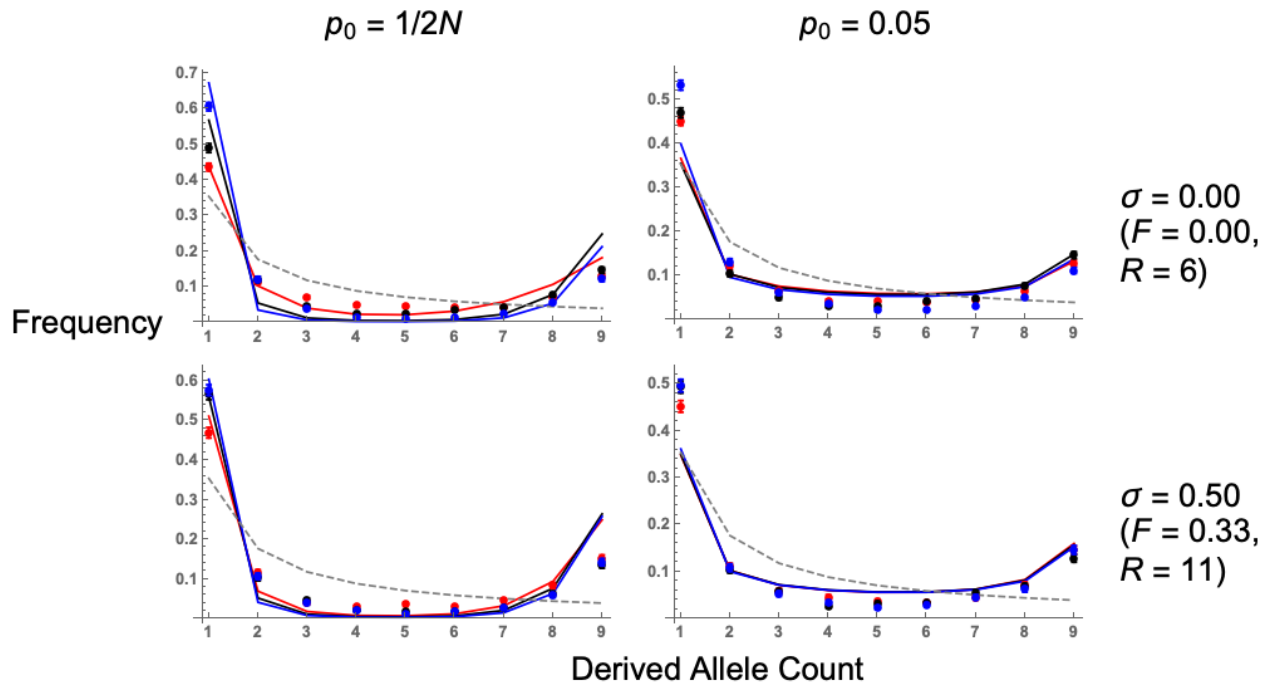


Figure 4. Expected site frequency spectrum, in flanking regions to the adaptive mutation, following a selective sweep. Lines are analytical solutions (Equation A12 in Supplementary File S2), points are simulation results. $N = 5,000$, $s = 0.05$, $4N\mu = 40$, and dominance coefficient $h = 0.1$ (red 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.

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

where:

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

and:

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

$P_{M,sw}$ is the mutation probability during the sweep phase, and is similar to Equation 12 except that $2r(1-2F+\Phi)(1-p')$ is replaced by $2\mu_b(1-p')/p'$, for p' is the derived allele frequency when the event occurs. $P_{M,sd}$ is the probability that, at the sweep origin, the derived allele appears by mutation instead of coalescing, and is defined in a similar manner to $P_{R,sd}$ (Equation 7):

$$P_{M,sd} = \frac{\frac{2\mu_b(1-p_{0,A})}{p_{0,A}}}{\frac{1+F}{2Np_{0,A}} + \frac{2\mu_b(1-p_{0,A})}{p_{0,A}}} = \frac{2\Theta_b(1-p_{0,A})}{1+F+2\Theta_b(1-p_{0,A})} \quad (15)$$

where $\Theta_b = 2N\mu_b$. The coalescence probability is one minus $P_{M,sd}$. We can also obtain an analytical solution using a ‘star-like’ approximation, which assumes that no coalescence or mutation events occur during the sweep phase:

$$\mathbb{E}_{SL} \left(\frac{\pi_M}{\pi_0} \right) = 1 - \left[\frac{1}{1 + 2(1-p_{0,A})\Theta_b/(1+F)} \right] \cdot \left[\frac{H_l}{H_h} \left(\frac{1}{p_{0,A}} + 1 \right) - 1 \right]^{-2r(1-2F+\Phi)/(H_l s)} \quad (16)$$

Figure 5 compares $\mathbb{E}(\pi/\pi_0)$ in the standing variation case, and for the recurrent mutation case, under different levels of self-fertilisation. While dominance only weakly affects sweep signatures arising from standing variation under outcrossing, it more strongly affects sweeps from recurrent mutation in outcrossing populations, as each variant arises from an initial frequency close to $1/(2N)$ (Figure 3). Second, the two models exhibit different behaviour close to the selected

locus (R close to zero). The recurrent mutation model has diversity levels that are greater than zero, while the standing variation model exhibits little diversity. As R increases, diversity reaches higher levels in the standing variation case than for the recurrent mutation case. Assuming weak recombination (so that $P_{NE} \approx 1$), the recombination rate at which a sweep from recurrent mutation yields higher diversity than one from standing variation is given when the coalescence probability is higher for the mutation model than that for the standing variation case. This change occurs at:

$$R \leq R_{Lim} = \frac{\Theta_b}{p_0(1 - 2F + \Phi)} \quad (17)$$

$$\approx \frac{\Theta_b}{p_0(1 - F)} \quad (18)$$

The last approximation arises as $\Phi \approx F$ unless F is close to one, and recombination rates are high (r approximately greater than 0.1). Hence for a fixed Θ_b , the window where recurrent mutations create higher diversity near the selected locus increases for lower p_0 or higher F , since both these factors reduces the potential for recombination to create new haplotypes during the standing phase. Equation 18 is generally accurate when sweeps from standing variation have higher diversity than sweeps with recurrent mutations (Figure 5, bottom row), but becomes inaccurate for $h = 0.1$ in outcrossing populations, as some events are likely to occur during the sweep phase. In Supplementary File S2 we show how similar results apply to the SFS.

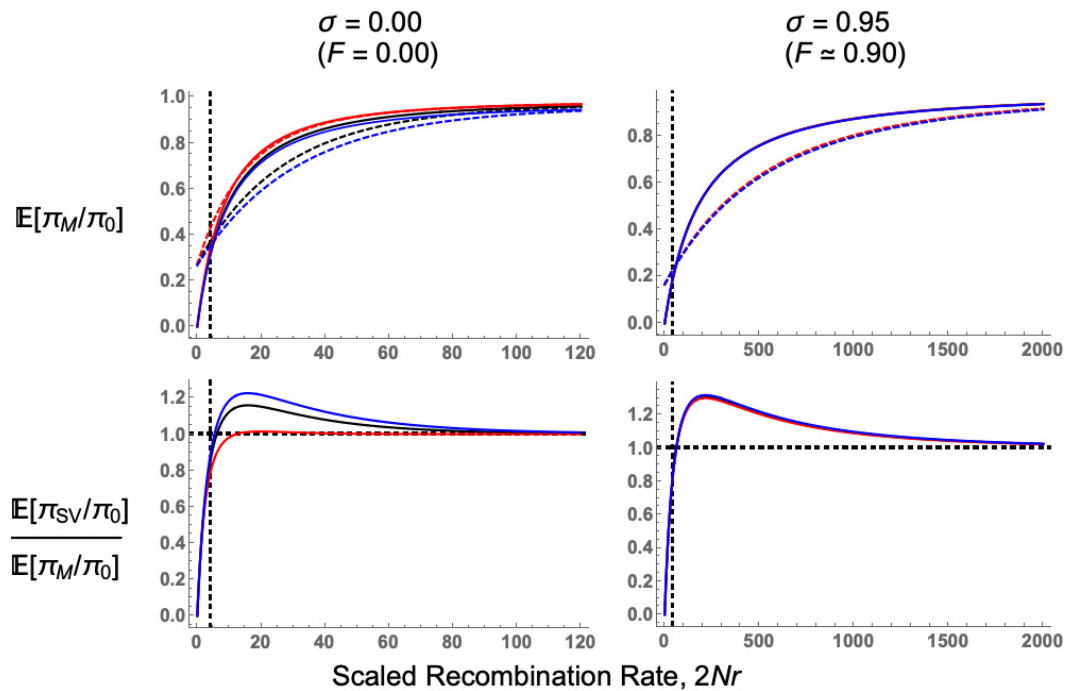


Figure 5. Comparing sweeps from recurrent mutation to those from standing variation. Top row: comparing the reduction in diversity following a soft sweep (Equation 11), from either standing variation ($p_0 = 0.05$, solid lines) or recurrent mutation (using $P_{coal,M}$ with $\Theta_b = 0.2$, dashed lines). $N = 5,000$, $s = 0.05$, and dominance coefficient $h = 0.1$ (red lines), 0.5 (black lines), or 0.9 (blue lines). Bottom row: 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 the respective plot for the top row. Vertical dashed black line indicates R_{Lim} (Equation 18), the predicted recombination rate where $\pi_{SV}/\pi_M = 1$ (horizontal dashed line in the bottom-row plots). Note the different x -axis lengths between panels. Results are also plotted in Section F of Supplementary File S1.

Discussion

Summary of Theoretical Findings

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

Messer and Petrov 2013; Stephan 2016; Hermisson and Pennings 2017), these models generally assumed idealised randomly mating populations and beneficial mutations that are semidominant ($h = 0.5$). Here we have created a more general selective sweep model, with arbitrary self-fertilisation and dominance levels. Our principal focus is on comparing a hard sweep arising from a single allele copy to a soft sweep arising from standing variation, but we also consider the case of recurrent mutation (Figure 5).

We find that the qualitative patterns of different selective sweeps under selfing remain similar to expectations from outcrossing models. In particular, a sweep from standing variation still creates an elevated number of intermediate-frequency variants compared to a sweep from *de novo* mutation (Figures 4, 5). This pattern is standard for soft sweeps (Pennings and Hermisson 2006b; Messer and Petrov 2013; Berg and Coop 2015; Hermisson and Pennings 2017) so existing statistical methods for detecting them (e.g., observing an higher than expected number of haplotypes; Vitti *et al.* (2013); Garud *et al.* (2015)) can, in principle, also be applied to selfing organisms. Under self-fertilisation, these signatures are stretched over longer physical regions than in outcrossers. These extensions arise as self-fertilisation affects gene genealogies during both the sweep and standing phases in different ways. During the sweep phase, beneficial alleles fix more rapidly under higher self-fertilisation as homozygous mutations are created more rapidly (Charlesworth 1992; Glémin 2012). In addition, the effective recombination rate is reduced by approximately $1 - F$ (Nordborg *et al.* 1996; Nordborg 2000; Charlesworth and Charlesworth 2010), and slightly further for highly inbred populations (Roze 2009, 2016). These two effects mean that neutral variants linked to an adaptive allele are less likely to recombine onto the neutral background during the sweep phase, as re-

flected in Equation 3 for P_{NE} . During the standing phase, two haplotypes are more likely to coalesce under high levels of self-fertilisation since N_e is decreased by a factor $1/(1+F)$ (Pollak 1987; Charlesworth 1992; Caballero and Hill 1992; Nordborg and Donnelly 1997). This effect, combined with a reduced effective recombination rate, means that the overall recombination probability during the standing phase is reduced by a factor $\sim(1-\sigma)$ (Equations 7). Hence intermediate-frequency variants, which could provide evidence of adaptation from standing variation, will be spread out over longer genomic regions. The elongation of sweep signatures means soft sweeps can be easier to detect in selfing organisms than in outcrossers.

We have also investigated how dominance affects soft sweep signatures, since previous analyses have only focussed on how dominance affects hard sweeps (Teshima and Przeworski 2006; Teshima *et al.* 2006; Ewing *et al.* 2011). In outcrossing organisms, recessive mutations leave weaker sweep signatures than additive or dominant mutations as they spend more time at low frequencies, increasing the amount of recombination that restores neutral variation (Figures 2, 3). With increased self-fertilisation, dominance has a weaker impact on sweep signatures as most mutations are homozygous (Figure 3). We also show that the SFS for recessive alleles can resemble a soft sweep, with a higher number of intermediate-frequency variants than for other hard sweeps (Figure 4). Dominance only weakly affects sweeps from standing variation, as trajectories of beneficial alleles become similar once the variant's initial frequency exceeds $1/(2N)$ (Figures 2, 3). Yet different dominance levels can affect sweep signatures if the beneficial allele is reintroduced by recurrent mutation (Figure 5). Hence if one wishes to understand how dominance affects sweep signatures, it is also important to consider which processes underlies observed patterns of genetic diversity.

479 Soft sweeps from recurrent mutation or standing variation?

480 These theoretical results shed light onto how to distinguish between soft sweeps
481 that arise either from standing variation, or from recurrent mutation. Both models
482 are characterised by an elevated number of intermediate-frequency variants, in
483 comparison to a hard sweep. Yet sweeps arising from recurrent mutation produce
484 intermediate-frequency variants closer to the beneficial locus, compared to sweeps
485 from standing variation (Figures 5 and C in Supplementary File S2). Equation 18
486 provides a simple condition for R_{Lim} , the recombination distance needed for a
487 sweep from standing variation to exhibit higher diversity than one from recurrent
488 mutation; the size of this region increases under higher self-fertilisation.

489 Differences in haplotype structure between sweeps from either standing varia-
490 tion or recurrent mutation should be more pronounced in self-fertilising organisms,
491 due to the reduction in effective recombination rates. However, when investigating
492 sweep patterns over broad genetic regions, it becomes likelier that genetic diversity
493 will be affected by multiple beneficial mutations spreading throughout the genome.
494 Competing selective sweeps can lead to elevated diversity near a target locus for
495 two reasons. First, selection interference increases the fixation time of individual
496 mutations, allowing more recombination that can restore neutral diversity (Kim
497 and Stephan 2003). In addition, competing selective sweeps can drag different
498 sets of neutral variation to fixation, creating asymmetric reductions in diversity
499 around a substitution (Chevin *et al.* 2008). Further investigations of selective
500 sweep patterns across long genetic distances will prove to be a rich area of future
501 research.

Potential applications to self-fertilising organisms

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

Acknowledgments. We would like to thank Sally Otto for providing information on the elevated effective starting frequency of beneficial mutations; Brian Charlesworth on providing advice on modelling selective sweeps, sharing unpublished results, and providing comments on the manuscript; and Ben Haller for answering questions about SLiM. MH was supported by a Marie Curie International Outgoing Fellowship (MC-IOF-622936) and a NERC Independent Research Fellowship (NE/R015686/1). MH and TB also acknowledge financial support from

526 the European Research Council under the European Union’s Seventh Framework
527 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. Aguileta, *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. Strutsovsky, G. McVicker, *et al.*, 2016 A genomic map of the effects of linked selection in *Drosophila*. *PLoS Genet.* **12**: e1006130.
- Ewing, G., J. Hermisson, P. Pfaffelhuber, and J. Rudolf, 2011 Selective sweeps for recessive alleles and for other modes of dominance. *J. Math. Bio.* **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* .
- 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.
- Kimura, M., 1971 Theoretical foundation of population genetics at the molecular level. Theor. Popul. Biol. **2**: 174–208.
- 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.
- Qanbari, S., H. Pausch, S. Jansen, M. Somel, T. M. Strom, *et al.*, 2014 Classic selective sweeps revealed by massive sequencing in cattle. *PLoS Genet.* **10**: e1004148.

- 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.
- 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.
- 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.
- 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.