

# BAYESIAN INFERENCE OF NATURAL SELECTION FROM ALLELE FREQUENCY TIME SERIES

JOSHUA G. SCHRAIBER, STEVEN N. EVANS, AND MONTGOMERY SLATKIN

**ABSTRACT.** The advent of accessible ancient DNA technology now allows the direct ascertainment of allele frequencies in ancestral populations, thereby enabling the use of allele frequency time series to detect and estimate natural selection. Such direct observations of allele frequency dynamics are expected to be more powerful than inferences made using patterns of linked neutral variation obtained from modern individuals. We developed a Bayesian method to make use of allele frequency time series data and infer the parameters of general diploid selection, along with allele age, in non-equilibrium populations. We introduce a novel path augmentation approach, in which we use Markov chain Monte Carlo to integrate over the space of allele frequency trajectories consistent with the observed data. Using simulations, we show that this approach has good power to estimate selection coefficients and allele age. Moreover, when applying our approach to data on horse coat color, we find that ignoring a relevant demographic history can significantly bias the results of inference. Our approach is made available in a C++ software package.

## 1. INTRODUCTION

The ability to obtain high-quality genetic data from ancient samples is revolutionizing the way that we understand the evolutionary history of populations. One of the most powerful applications of ancient DNA (aDNA) is to study the action of natural selection. While methods making use of only modern DNA sequences have successfully identified loci evolving subject to natural selection [Nielsen et al., 2005, Voight et al., 2006, Pickrell et al., 2009], they are inherently limited because they look indirectly for selection, finding its signature in nearby neutral variation. In contrast, by sequencing ancient individuals, it is possible to directly track the change in allele frequency that is characteristic of the action of natural selection. This approach has been exploited recently using whole genome data to identify candidate loci under selection in European humans [Mathieson et al., 2015].

To infer the action of natural selection rigorously, several methods have been developed to explicitly fit a population genetic model to a time series of allele frequencies obtained via aDNA. Initially, Bollback et al. [2008] extended an approach devised by Williamson and Slatkin [1999] to estimate the population-scaled selection coefficient,  $\alpha = 2N_e s$ , along with the effective size,  $N_e$ . To incorporate natural selection, Bollback et al. [2008] used

---

*Date:* Started on December 10, 2013. Compiled on March 3, 2016.

JGS supported by NSF grant DBI-1402120 to JGS, NIH grant R01-GM40282 to MS, SNE supported in part by NSF grant DMS-0907630, NSF grant DMS-1512933, and NIH grant 1R01GM109454-01, MS supported by NIH grant R01-GM40282 to MS.

20 the continuous diffusion approximation to the discrete Wright-Fisher model. This required  
21 them to use numerical techniques to solve the partial differential equation (PDE) associated  
22 with transition densities of the diffusion approximation to calculate the probabilities of the  
23 population allele frequencies at each time point. Ludwig et al. [2009] obtained an aDNA  
24 time series from 6 coat-color-related loci in horses and applied the method of Bollback  
25 et al. [2008] to find that 2 of them, ASIP and MC1R, showed evidence of strong positive  
26 selection.

27 Recently, a number of methods have been proposed to extend the generality of the  
28 Bollback et al. [2008] framework. To define the hidden Markov model they use, Bollback  
29 et al. [2008] were required to posit a prior distribution on the allele frequency at the first  
30 time point. They chose to use a uniform prior on the initial frequency; however, in truth  
31 the initial allele frequency is dictated by the fact that the allele at some point arose as a  
32 new mutation. Using this information, Malaspinas et al. [2012] developed a method that  
33 also infers allele age. They also extended the selection model of Bollback et al. [2008] to  
34 include fully recessive fitness effects. A more general selective model was implemented by  
35 Steinrücken et al. [2014], who model general diploid selection, and hence they are able to fit  
36 data where selection acts in an over- or under-dominant fashion; however, Steinrücken et al.  
37 [2014] assumed a model with recurrent mutation and hence could not estimate allele age.  
38 The work of Mathieson and McVean [2013] is designed for inference of metapopulations  
39 over short time scales and so it is computationally feasible for them to use a discrete time,  
40 finite population Wright-Fisher model. Finally, the approach of Feder et al. [2014] is ideally  
41 suited to experimental evolution studies because they work in a strong selection, weak drift  
42 limit that is common in evolving microbial populations.

43 One key way that these methods differ from each other is in how they compute the  
44 probability of the underlying allele frequency changes. For instance, Malaspinas et al.  
45 [2012] approximated the diffusion with a birth-death type Markov chain, while Steinrücken  
46 et al. [2014] approximate the likelihood analytically using a spectral representation of  
47 the diffusion discovered by Song and Steinrücken [2012]. These different computational  
48 strategies are necessary because of the inherent difficulty in solving the Wright-Fisher  
49 partial differential equation. A different approach, used by Mathieson and McVean [2013]  
50 in the context of a densely-sampled discrete Wright-Fisher model, is to instead compute  
51 the probability of the entire allele frequency trajectory in between sampling times.

52 In this work, we develop a novel approach for inference of general diploid selection and  
53 allele age from allele frequency time series obtained from aDNA. The key innovation of  
54 our approach is that we impute the allele frequency trajectory between sampled points  
55 when they are sparsely-sampled. Moreover, by working with a diffusion approximation,  
56 we are able to easily incorporate general diploid selection and changing population size.  
57 This approach to inferring parameters from a sparsely-sampled diffusion is known as high-  
58 frequency path augmentation, and has been successfully applied in a number of contexts  
59 [Roberts and Stramer, 2001, Golightly and Wilkinson, 2005, 2008, Sørensen, 2009, Fuchs,  
60 2013]. The diffusion approximation to the Wright-Fisher model, however, has several  
61 features that are atypical in the context of high-frequency path augmentation, including  
62 a time-dependent diffusion coefficient and a bounded state-space. We test this approach

63 with simulation, showing that it's important to accurately model demography history, then  
64 apply it to several datasets and find that we have power to estimate parameters of interest  
65 from real data.

## 66 2. MODEL AND METHODS

67 **2.1. Overview.** We begin by first reviewing the Wright-Fisher model, presenting its dif-  
68 fusion approximation as a stochastic differential equation (SDE). We then describe our  
69 inferential strategy using a path augmentation approach, in which we model the under-  
70 lying allele frequency trajectory as an additional (infinite dimensional) parameter. This  
71 requires us to derive an expression for the likelihood of an allele frequency trajectory, in-  
72 cluding accounting for the fact that we model alleles that start from low frequency as new  
73 mutants. Finally, we describe a Markov chain Monte Carlo algorithm for obtaining a pos-  
74 terior distribution of the parameters of natural selection, as well as the allele frequency  
75 trajectory.

76 **2.2. Generative model.** We assume a randomly mating diploid population that is size  
77  $N(t)$  at time  $t$ , where  $t$  is measured in units of  $2N_0$  generations for some arbitrary, constant  
78  $N_0$ . At the locus of interest, the ancestral allele,  $A_0$ , was fixed until some time  $t_0$  when the  
79 derived allele,  $A_1$ , arose with diploid fitnesses as given in Table 1.

80 [Table 1 about here.]

81 Given that an allele arises at some finite population frequency  $0 < x_0 < 1$  at some time  
82  $t_0$ , the trajectory of population frequencies of  $A_1$  at times  $t \geq t_0$ ,  $(X_t)_{t \geq t_0}$ , is modeled  
83 by the usual diffusion approximation to the Wright-Fisher model (and many other models  
84 such as the Moran model), which we will henceforth call the Wright-Fisher diffusion. While  
85 many treatments of the Wright-Fisher diffusion define it in terms of the partial differential  
86 equation that characterizes its transition densities (e.g. Ewens [2004]), we instead describe  
87 it as the solution of a stochastic differential equation (SDE). Specifically,  $(X_t)_{t \geq t_0}$  satisfies  
88 the SDE

$$(1) \quad \begin{aligned} dX_t &= X_t(1 - X_t)(\alpha_1(2X_t - 1) - \alpha_2 X_t) dt + \sqrt{\frac{X_t(1 - X_t)}{\rho(t)}} dB_t \\ X_{t_0} &= x_0, \end{aligned}$$

89 where  $B$  is a standard Brownian motion,  $\alpha_1 = 2N_0s_1$ ,  $\alpha_2 = 2N_0s_2$ , and  $\rho(t) = N(t)/N_0$ .  
90 If  $X_{t_*} = 0$  (resp.  $X_{t_*} = 1$ ) at some time  $t_* > t_0$ , then  $X_t = 0$  (resp.  $X_t = 1$ ) for all  $t \geq t_*$ .

91 In order to make this description of the dynamics of the population allele frequency  
92 trajectory  $(X_t)_{t \geq t_0}$  complete, we need to specify an initial condition at time  $t_0$ . In a finite  
93 population Wright-Fisher model we would take the allele  $A_1$  to have frequency  $\frac{1}{2N(t_0)}$  at  
94 the time  $t_0$  when it first arose in a single chromosome. This frequency converges to 0 when  
95 we pass to the diffusion limit, but we cannot start the Wright-Fisher diffusion at 0 at time  
96  $t_0$  because the diffusion started at 0 remains at 0. Instead, we take the value of  $X_{t_0}$  to  
97 be some small, but arbitrary, frequency  $x_0$ . This arbitrariness in the choice of  $x_0$  may  
98 seem unsatisfactory, but we will see that, in the context of a Bayesian inference procedure,

99 the resulting posterior distribution for the parameters  $\alpha_1, \alpha_2, t_0$  converges as  $x_0 \downarrow 0$  to a  
100 limit which can be thought of as the posterior corresponding to a certain improper prior  
101 distribution, and so, in the end, there is actually no need to specify  $x_0$ .

102 Finally, we require a model for how alleles arise. We assume that mutations at time  $t$   
103 occur at a rate proportional to  $2N(t)$ , and that a mutant allele arises exactly once. Further  
104 constraining alleles to have arisen more recently than some time,  $T$ , in the past, this implies  
105 that the prior density of allele ages is

$$\pi(t_0) = \frac{\rho(t_0)}{\int_T^0 \rho(t_0) ds}.$$

106 Taking the limit as  $T \downarrow -\infty$  results in an improper distribution on allele age, which, in the  
107 context of our Bayesian inference algorithm, implies an improper prior distribution on  $t_0$   
108 that is proportional to  $\rho$ . However, we emphasize that this still produces a proper posterior  
109 distribution on allele age (see also Slatkin [2001]).

110 Finally, we model the data assuming that at known times  $t_1, t_2, \dots, t_k$  samples of known  
111 sizes  $n_1, n_2, \dots, n_k$  chromosomes are taken and  $c_1, c_2, \dots, c_k$  copies of the derived allele are  
112 found at the successive time points (Figure 1). Note that it is possible that some of the  
113 sampling times are more ancient than  $t_0$ , the age of the allele.

114

[Figure 1 about here.]

115 **2.3. Bayesian path augmentation.** We are interested in devising a Bayesian method  
116 to obtain the posterior distribution on the parameters,  $\alpha_1, \alpha_2$ , and  $t_0$  given the sampled  
117 allele frequencies and sample times – data which we denote collectively as  $D$ . Because  
118 we are dealing with objects that don't necessarily have distributions which have densities  
119 with respect to canonical reference measures, it will be convenient in the beginning to  
120 treat priors and posteriors as probability measures rather than as density functions. For  
121 example, the posterior is the probability measure

$$(2) \quad P(d\alpha_1, d\alpha_2, dt_0 | D) = \frac{P(dD | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0)}{P(dD)},$$

122 where  $\pi$  is a joint prior on the model parameters. However, computing the likelihood  
123  $P(dD | \alpha_1, \alpha_2, t_0)$  is computationally challenging because, implicitly,

$$P(dD | \alpha_1, \alpha_2, t_0) = \int P(dD | X) P(dX | \alpha_1, \alpha_2, t_0),$$

124 where the integral is over the (unobserved, infinite-dimensional) allele frequency path  
125  $X = (X_t)_{t \geq t_0}$ ,  $P(\cdot | \alpha_1, \alpha_2, t_0)$  is the distribution of a Wright-Fisher diffusion with selection  
126 parameters  $\alpha_1, \alpha_2$  started at time  $t_0$  at the small but arbitrary frequency  $x_0$ , and

$$P(dD | X) = \prod_{i=1}^k \binom{n_i}{c_i} X_{t_i}^{c_i} (1 - X_{t_i})^{n_i - c_i}$$

127 because we assume that sampled allele frequencies at the times  $t_1, \dots, t_k$  are independent  
128 binomial draws governed by underlying population allele frequencies at the these times.

129 Integrating over the infinite-dimensional path  $(X_t)_{t \geq t_0}$  involves either solving partial dif-  
130 ferential equations numerically or using Monte Carlo methods to find the joint distribution  
131 of population allele frequency path at the times  $t_1, \dots, t_k$ .

132 To address this computational difficulty, we introduce a path augmentation method that  
133 treats the underlying allele frequency path  $(X_t)_{t \geq t_0}$  as an additional parameter. Observe  
134 that the posterior may be expanded out to

$$P(d\alpha_1, d\alpha_2, dt_0 | D) = \frac{\int P(dD | X') P(dX' | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0)}{\int P(dD | X') P(dX' | \alpha'_1, \alpha'_2, t'_0) \pi(d\alpha'_1, d\alpha'_2, dt'_0)},$$

135 where we use primes to designate dummy variables over which we integrate. Thinking of  
136 the path  $(X_t)_{t \geq t_0}$  as another parameter and taking the prior distribution for the augmented  
137 family of parameters to be

$$P(dX | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0),$$

138 the posterior for the augmented family of parameters is

$$(3) \quad P(d\alpha_1, d\alpha_2, dt_0; dX | D) = \frac{P(dD | X) P(dX | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0)}{\int P(dD | X') P(dX' | \alpha'_1, \alpha'_2, t'_0) \pi(d\alpha'_1, d\alpha'_2, dt'_0)}.$$

139 We thus see that treating the allele frequency path as a parameter is consistent with  
140 the initial “naive” Bayesian approach in that if we integrate the path variable out of the  
141 posterior (3) for the augmented family of parameters, then we recover the posterior (2)  
142 for the original family of parameters. In practice, this means that marginalizing out the  
143 path variable from a Monte Carlo approximation of the augmented posterior gives a Monte  
144 Carlo approximation of the original posterior.

145 Implicit in our set-up is the initial frequency  $x_0$  at time  $t_0$ . Under the probability  
146 measure governing the Wright-Fisher diffusion, any process started from  $x_0 = 0$  will stay  
147 there forever. Thus, we would be forced to make an arbitrary choice of some  $x_0 > 0$   
148 as the initial frequency of our allele. However, we argue in the Appendix that in the limit  
149 as  $x_0 \downarrow 0$ , we can achieve an improper prior distribution on the space of allele frequency  
150 trajectories. We stress that our inference using such an improper prior is not one that arises  
151 directly from a generative probability model for the allele frequency path. However, it does  
152 arise as a limit as the initial allele frequency  $x_0$  goes to zero of inferential procedures based  
153 on generative probability models and the limiting posterior distributions are probability  
154 distributions. Therefore, the parameters  $\alpha_1, \alpha_2, t_0$  retain their meaning, our conclusions  
155 can be thought of approximations to those that we would arrive at for all sufficiently small  
156 values of  $x_0$ , and we are spared the necessity of making an arbitrary choice of  $x_0$ .

157 **2.4. Path likelihoods.** Most instances of Bayesian inference in population genetics have  
158 hitherto involved finite-dimensional parameters. Recall that for continuous, finite-dimensional  
159 parameters, one simply includes the prior *density* of the parameter value in place of the  
160 prior *probability*. Finite dimensional parameters usually have densities defined with respect  
161 to Lebesgue measure in an appropriate dimension; however, there is no infinite-dimensional  
162 Lebesgue measure against which to define a density for our infinite-dimensional augmented  
163 path. We thus require a reference measure on the infinite-dimensional space of paths that

164 will play a role analogous to that of Lebesgue measure in the finite-dimensional case, al-  
 165 lowing us to write down the probability density for each sampled path.

166 To see what is involved, suppose we have a diffusion process  $(Z_t)_{t \geq t_0}$  that satisfies the  
 167 SDE

$$(4) \quad \begin{aligned} dZ_t &= a(Z_t, t) dt + dB_t \\ Z_{t_0} &= z_0, \end{aligned}$$

168 where  $B$  is a standard Brownian motion (the Wright-Fisher diffusion is not of this form  
 169 but, as we shall soon see, it can be reduced to it after suitable transformations of time  
 170 and space). Let  $\mathbb{P}$  be the distribution of  $(Z_t)_{t \geq t_0}$  – this is a probability distribution on the  
 171 space of continuous paths that start from position  $z_0$  at time  $t_0$ . While the probability  
 172 assigned by  $\mathbb{P}$  to any particular path is zero, we can, under appropriate conditions, make  
 173 sense of the probability of a path under  $\mathbb{P}$  relative to its probability under the distribution  
 174 of Brownian motion. If we denote by  $\mathbb{W}$  the distribution of Brownian motion starting from  
 175 position  $z_0$  at time  $t_0$ , then Girsanov’s theorem [Girsanov, 1960] gives the density of the  
 176 path segment  $(Z_s)_{t_0 \leq s \leq t}$  under  $\mathbb{P}$  relative to  $\mathbb{W}$  as

$$(5) \quad \frac{d\mathbb{P}}{d\mathbb{W}}((Z_s)_{t_0 \leq s \leq t}) = \exp \left\{ \int_{t_0}^t a(Z_s, s) dZ_s - \frac{1}{2} \int_{t_0}^t a^2(Z_s, s) ds \right\},$$

177 where the first integral in the exponent is an Itô integral. In order for (5) to hold, the  
 178 integral  $\int_{t_0}^t a^2(Z_s, s) ds$  must be finite, in which case the Itô integral  $\int_{t_0}^t a(Z_s, s) dZ_s$  is also  
 179 well-defined and finite.

180 However, the Wright-Fisher SDE (1) is not of the form (4). In particular, the factor  
 181 multiplying the infinitesimal Brownian increment  $dB_t$  (the so-called diffusion coefficient)  
 182 depends on both space and time. To deal with this issue, we first apply a well-known time  
 183 transformation (see e.g. Slatkin and Hudson [1991] and Griffiths and Tavaré [1994]) and  
 184 consider the process  $(\tilde{X}_\tau)_{\tau \geq 0}$  given by  $\tilde{X}_\tau = X_{f^{-1}(\tau)}$ , where

$$(6) \quad f(t) = \int_{t_0}^t \frac{1}{\rho(s)} ds, \quad t \geq t_0.$$

185 It is not hard to see that  $(\tilde{X}_\tau)_{\tau \geq 0}$  satisfies the following SDE with a time-independent  
 186 diffusion coefficient,

$$\begin{aligned} d\tilde{X}_\tau &= \rho(f^{-1}(\tau)) \tilde{X}_\tau (1 - \tilde{X}_\tau) (\alpha_1 (2\tilde{X}_\tau - 1) - \alpha_2 \tilde{X}_\tau) d\tau + \sqrt{\tilde{X}_\tau (1 - \tilde{X}_\tau)} d\tilde{B}_\tau \\ \tilde{X}_0 &= x_0, \end{aligned}$$

187 where  $\tilde{B}$  is a standard Brownian motion. Next, we employ an angular space transformation  
 188 first suggested by Fisher [1922],  $Y_\tau = \arccos(1 - 2\tilde{X}_\tau)$ . Applying Itô’s lemma [Itô, 1944]  
 189 shows that  $(Y_\tau)_{\tau \geq 0}$  is a diffusion that satisfies the SDE

$$(7) \quad \begin{aligned} dY_\tau &= \frac{1}{4} (\rho(f^{-1}(\tau)) \sin(Y_\tau) (\alpha_2 + (2\alpha_1 - \alpha_2) \cos(Y_\tau)) - 2 \cot(Y_\tau)) d\tau + dW_\tau \\ Y_0 &= y_0 = \arccos(1 - 2x_0), \end{aligned}$$



190 where  $W$  is a standard Brownian motion. If the process  $X$  hits either of the boundary  
 191 points  $0, 1$ , then it stays there, and the same is true of the time and space transformed  
 192 process  $Y$  for its boundary points  $0, \pi$ .

193 The restriction of the distribution of the time and space transformed process  $Y$  to some  
 194 set of paths that don't hit the boundary is absolutely continuous with respect to the dis-  
 195 tribution of standard Brownian motion restricted to the same set; that is, the distribution  
 196 of  $Y$  restricted to such a set of paths has a density with respect to the distribution of  
 197 Brownian motion restricted to the same set. However, the infinitesimal mean in (7) (that  
 198 is, the term multiplying  $d\tau$ ) becomes singular as  $Y_\tau$  approaches the boundary points  $0$  and  
 199  $\pi$ , corresponding to the boundary points  $0$  and  $1$  for allele frequencies. These singularities  
 200 prevent the process  $Y$  from re-entering the interior of its state space and ensure that a  
 201 Wright-Fisher path will be absorbed when the allele is either fixed or lost. A consequence  
 202 is that the density of the distribution of  $Y$  relative to that of a Brownian motion blows up  
 203 as the path approaches the boundary. We are modeling the appearance of a new mutation  
 204 in terms of a Wright-Fisher diffusion starting at some small initial frequency  $x_0$  at time  
 205  $t_0$  and we want to perform our parameter inference in such a way that we get meaning-  
 206 ful answers as  $x_0 \downarrow 0$ . This suggests that rather than working with the distribution  $\mathbb{W}$   
 207 of Brownian motion as a reference measure it may be more appropriate to work with a  
 208 tractable diffusion process that exhibits similar behavior near the boundary point  $0$ .

209 To start making this idea of matching singularities more precise, consider a diffusion  
 210 process  $(\bar{Z}_t)_{t \geq t_0}$  that satisfies the SDE

$$(8) \quad \begin{aligned} d\bar{Z}_t &= b(\bar{Z}_t, t) dt + d\bar{B}_t \\ \bar{Z}_0 &= z_0, \end{aligned}$$

where  $\bar{B}$  is a standard Brownian motion. Write  $\mathbb{Q}$  for the distribution of the diffusion  
 process  $(\bar{Z}_t)_{t \geq t_0}$  and recall that  $\mathbb{P}$  is the distribution of a solution of (4). If  $(Z_s)_{t_0 \leq s \leq t}$  is a  
 segment of path such that both  $\int_{t_0}^t a^2(Z_s, s) ds < \infty$  and  $\int_{t_0}^t b^2(Z_s, s) ds < \infty$ , then

$$(9) \quad \begin{aligned} \frac{d\mathbb{P}}{d\mathbb{Q}}((Z_s)_{t_0 \leq s \leq t}) &= \frac{d\mathbb{P}}{d\mathbb{W}}((Z_s)_{t_0 \leq s \leq t}) / \frac{d\mathbb{Q}}{d\mathbb{W}}((Z_s)_{t_0 \leq s \leq t}) \\ &= \exp \left\{ \int_{t_0}^t (a(Z_s, s) - b(Z_s, s)) dZ_s - \frac{1}{2} \int_{t_0}^t (a^2(Z_s, s) - b^2(Z_s, s)) ds \right\}. \end{aligned}$$

211 Note that the right-hand side will stay bounded if one considers a sequence of paths,  
 212 indexed by  $\eta$ ,  $(Z_s^\eta)_{t_0 \leq s \leq t}$ , with  $\int_{t_0}^t a^2(Z_s^\eta, s) ds < \infty$  and  $\int_{t_0}^t b^2(Z_s^\eta, s) ds < \infty$ , provided  
 213 that  $\int_{t_0}^t (a^2(Z_s^\eta, s) - b^2(Z_s^\eta, s)) ds$  stays bounded. These manipulations with densities may  
 214 seem somewhat heuristic, but they can be made rigorous and, moreover, the form of  $\frac{d\mathbb{P}}{d\mathbb{Q}}$   
 215 follows from an extension of Girsanov's theorem that gives the density of  $\mathbb{P}$  with respect  
 216 to  $\mathbb{Q}$  directly without using the densities with respect to  $\mathbb{W}$  as intermediaries (see, for  
 217 example, [Kallenberg, 2002, Theorem 18.10]).

218 We wish to apply this observation to the time and space transformed Wright-Fisher  
219 diffusion of (7). Because

$$-\frac{1}{2} \cot(y) + \frac{1}{4} \rho(f^{-1}(t)) \sin(y) ((2\alpha_1 - \alpha_2) \cos(y) + \alpha_2) = -\frac{1}{2y} + O(y)$$

220 when  $y$  is small, an appropriate reference process should have infinitesimal mean  $b(y, t) \approx$   
221  $-1/(2y)$  as  $y \downarrow 0$ . Following suggestions by Schraiber et al. [2013] and Jenkins [2013], we  
222 compute path densities relative to the distribution  $\mathbb{Q}$  of the Bessel(0) process, a process  
223 which is the solution of the SDE

$$(10) \quad \begin{aligned} d\bar{Y}_t &= -\frac{1}{2\bar{Y}_t} dt + d\bar{B}_t, \\ \bar{Y}_0 &= y_0 = \arccos(1 - 2x_0). \end{aligned}$$

224 up until the first time that  $\bar{Y}_t$  hits 0, after which time  $\bar{Y}_t$  stays at 0 [Revuz and Yor, 1999,  
225 Chapter XI].

226 As we show more explicitly in the Appendix, this choice of dominating measure allows us  
227 to arrive at a proper posterior distribution as we send the initial frequency of the allele down  
228 to 0. In brief, if we write  $\mathbb{P}^{y_0}$  and  $\mathbb{Q}^{y_0}$  for the respective distributions of the solutions of (7)  
229 and (10) to emphasize the dependence on  $y_0$  (equivalently, on the initial allele frequency  
230  $x_0$ ), then there are  $\sigma$ -finite measures  $\mathbb{P}^0$  and  $\mathbb{Q}^0$  with infinite total mass such that for each  
231  $\epsilon > 0$

$$\lim_{y_0 \downarrow 0} \mathbb{P}^{y_0}((Y_t)_{t \geq \epsilon} \in \cdot \mid Y_\epsilon > 0) = \mathbb{P}^0((Y_t)_{t \geq \epsilon} \in \cdot) / \mathbb{P}^0(Y_\epsilon > 0)$$

232 and

$$\lim_{y_0 \downarrow 0} \mathbb{Q}^{y_0}((\bar{Y}_t)_{t \geq \epsilon} \in \cdot \mid \bar{Y}_\epsilon > 0) = \mathbb{Q}^0((\bar{Y}_t)_{t \geq \epsilon} \in \cdot) / \mathbb{Q}^0(\bar{Y}_\epsilon > 0),$$

233 where the numerators and denominators in the last two equations are all finite. Moreover,  
234  $\mathbb{P}^0$  has a density with respect to  $\mathbb{Q}^0$  that arises by naively taking limits as  $y_0 \downarrow 0$  in the  
235 functional form of the density of  $\mathbb{P}^{y_0}$  with respect to  $\mathbb{Q}^{y_0}$  (we say “naively” because  $\mathbb{P}^{y_0}$  and  
236  $\mathbb{Q}^{y_0}$  assign all of their mass to paths that start at position  $y_0 = \arccos(1 - 2x_0)$  at time 0,  
237 whereas  $\mathbb{P}^0$  and  $\mathbb{Q}^0$  assign all of their mass to paths that start at position 0 at time 0, and  
238 so the set of paths at which it is relevant to compute the density changes as  $y_0 \downarrow 0$ ). As  
239 we have already remarked, the limit of our Bayesian inferential procedure may be thought  
240 of as Bayesian inference with an improper prior, but we stress that the resulting posterior  
241 is proper.

242 The notion of the infinite measure  $\mathbb{Q}^0$  may seem somewhat forbidding, but this measure  
243 is characterized by the following simple properties:

$$\mathbb{Q}^0(\bar{Y}_\epsilon \in dy) = \frac{y^2}{\epsilon^2} \exp\left\{-\frac{y^2}{2\epsilon}\right\} dy, \quad y > 0,$$

244 so that  $\mathbb{Q}^0(\bar{Y}_\epsilon > 0) = \sqrt{\frac{\pi}{2}} \frac{1}{\sqrt{\epsilon}}$ , and conditional on the event  $\{\bar{Y}_\epsilon = y\}$  the evolution of  $(\bar{Y}_t)_{t \geq \epsilon}$   
245 is exactly that of the Bessel(0) process started at position  $y$  at time  $\epsilon$ . In the Appendix,  
246 we provide a more explicit construction of the measure  $\mathbb{Q}^0$  as part of our derivation of  
247 the proposal ratios in our MCMC algorithm. Moreover, conditional on the event  $\{\bar{Y}_s =$



248  $a, \bar{Y}_u = b\}$  for  $0 \leq s < u$  and  $a, b > 0$ , the evolution of the “bridge”  $(\bar{Y}_u)_{s \leq t \leq u}$  is the same  
 249 as that of the corresponding bridge for a Bessel(4) process; a Bessel(4) process satisfies the  
 250 SDE

$$d\hat{Y}_t = \frac{3}{2\hat{Y}_t} dt + d\hat{B}_t.$$

251 Very importantly for the sake of simulations, the Bessel(4) process is just the radial part of  
 252 a 4-dimensional standard Brownian motion – in particular, this process started at 0 leaves  
 253 immediately and never returns.

254 Note that the Bessel(0) process arises naturally because our space transformation  $x \mapsto$   
 255  $\arccos(1 - 2x) = \int_0^x \frac{1}{\sqrt{w(1-w)}} dw$  is approximately  $x \mapsto \int_0^x \frac{1}{\sqrt{w}} dw = 2\sqrt{x}$  when  $x > 0$  is  
 256 small. Interestingly, a multiple of the square of Bessel(0) process, sometimes called Feller’s  
 257 continuous state branching processes, arises naturally as an approximation to the Wright-  
 258 Fisher diffusion for low frequencies and has a long history in population genetics [Haldane,  
 259 1927, Feller, 1951].

260 **2.5. The joint likelihood of the data and the path.** To write down down the full  
 261 likelihood of the observations and the path, we make the assumption that the population  
 262 size function  $\rho(t)$  is continuously differentiable except at a finite set of times  $d_1 < d_2 <$   
 263  $\dots < d_M$ . Further, we require that that  $\rho(d_i^+) = \lim_{t \downarrow d_i} \rho(t)$  exists and is equal to  $\rho(d_i)$   
 264 while  $\rho(d_i^-) = \lim_{t \uparrow d_i} \rho(t)$  also exists (though it may not necessarily equal  $\rho(d_i)$ ).

265 We can write the joint likelihood of the data and the path as

$$L(D, (Y_t)_{t \geq 0} \mid \alpha_1, \alpha_2, t_0) = \mathbb{F}(D \mid (Y_t)_{t \geq 0}, t_0) \frac{d\mathbb{P}}{d\mathbb{Q}}((Y_t)_{t \geq 0}; \alpha_1, \alpha_2, t_0)$$

where  $\mathbb{F}(\cdot)$  is the binomial sampling probability of the observed allele frequencies,  $\mathbb{P}$  is the  
 distribution of transformed Wright-Fisher paths, and  $\mathbb{Q}$  is the distribution of Bessel(0)  
 paths. In the Appendix, we show that

$$\begin{aligned} & L(D, (Y_s)_{0 \leq s \leq t_k} \mid \alpha_1, \alpha_2, t_0) \\ &= \exp \left\{ A(Y_{f(t_k)}, t_k^-) + A(Y_{f(d_m)}, d_m^-) - (A(Y_{f(d_K)}, d_K) + A(Y_{f(t_0)}, t_0)) \right. \\ & \quad + \sum_{i=m}^K [A(Y_{f(d_{i+1})}, d_{i+1}^-) - A(Y_{f(d_i)}, d_i)] \\ & \quad \left. - \int_{t_0}^{t_k} B(Y_{f(s)}, s) ds - \frac{1}{2} \int_{t_0}^{t_k} C(Y_{f(s)}, s) ds - \frac{1}{2} \int_{t_0}^{t_k} D(Y_{f(s)}, s) ds \right\} \\ & \quad \times \prod_{i=1}^k \binom{n_i}{c_i} \left( \frac{1 - \cos(Y_{f(t_i)})}{2} \right)^{c_i} \left( \frac{1 + \cos(Y_{f(t_i)})}{2} \right)^{n_i - c_i}, \end{aligned} \tag{11}$$

where  $f$  is as in (6),  $m = \min\{i : d_i > t_0\}$  and  $K = \max\{i : d_i > t_k\}$ , and

$$\begin{aligned} A(y, t) &= \frac{\log(y)}{2} - \frac{1}{8} (\rho(t) \cos(y)(2\alpha_2 + (2\alpha_1 - \alpha_2) \cos(y)) + 4 \log(\sin(y))) \\ B(y, t) &= -\frac{1}{8} \frac{d\rho}{dt}(t) \cos(y)(2\alpha_2 + (2\alpha_1 - \alpha_2) \cos(y)) \\ C(y, t) &= \frac{1}{2} \left( \alpha_1 \cos(y) + \frac{\csc(y)^2}{\rho(t)} \right) - \frac{1}{2y^2 \rho(t)} \\ D(y, t) &= \frac{1}{16\rho(t)} (\rho(t) \sin(y)(\alpha_2 + (2\alpha_1 - \alpha_2) \cos(y)) - 2 \cot(y))^2 - \frac{1}{4y^2 \rho(t)}. \end{aligned}$$

266 While this expression may appear complicated, it has the important feature that, unlike  
267 the form of the likelihood that would arise by simply applying Girsanov's theorem, it only  
268 involves Lebesgue (indeed Riemann) integrals and not Itô integrals, which, as we recall  
269 in the Appendix, are known from the literature to be potentially difficult to compute  
270 numerically.

271 **2.6. Metropolis-Hastings algorithm.** We now describe a Markov chain Monte Carlo  
272 method for Bayesian inference of the parameters  $\alpha_1$ ,  $\alpha_2$  and  $t_0$ , along with the allele  
273 frequency path  $(X_t)_{t \geq t_0}$  (equivalently, the transformed path  $(Y_t)_{t \geq 0}$ ). While updates to  
274 the selection parameters  $\alpha_1$  and  $\alpha_2$  do not require updating the path, updating the time  $t_0$   
275 at which the derived allele arose requires proposing updates to the segment of path from  $t_0$   
276 up to the time of the first sample with a non-zero number of derived alleles. Additionally,  
277 we require proposals to update small sections of the path without updating any parameters  
278 and proposals to update the allele frequency at the most recent sample time.

279 [Figure 2 about here.]

280 **2.6.1. Interior path updates.** To update a section of the allele frequency, we first choose a  
281 time  $s_1 \in (t_0, t_k)$  uniformly at random, and then choose a time  $s_2$  that is a fixed fraction of  
282 the path length subsequent to  $s_1$ . We prefer this approach of updating a fixed fraction of  
283 the path to an alternative strategy of holding  $s_2 - s_1$  constant because paths for very strong  
284 selection may be quite short. Recalling the definition of  $f$  from (6), we subsequently propose  
285 a new segment of transformed path between the times  $f(s_1)$  and  $f(s_2)$  while keeping the  
286 values  $Y_{f(s_1)}$  and  $Y_{f(s_2)}$  fixed (Figure 2a). Such a path that is conditioned to take specified  
287 values at both end-points of the interval over which it is defined is called a bridge, and by  
288 updating small portions of the path instead of the whole path at once, we are able to obtain  
289 the desirable behavior that our Metropolis-Hastings algorithm is able to stay in regions of  
290 path space with high posterior probability. If we instead drew the whole path each time,  
291 we would much less efficiently target the posterior distribution.

292 Noting that bridges must be sampled against the *transformed* time scale, the best bridges  
293 for the allele frequency path would be realizations of Wright-Fisher bridges themselves.  
294 However, sampling Wright-Fisher bridges is challenging (but see Schraiber et al. [2013],  
295 Jenkins and Spano [2015]), so we instead opt to sample bridges for the transformed path

296 from the Bessel(0) process. Sampling Bessel(0) bridges can be accomplished by first sam-  
 297 pling Bessel(4) bridges (as described in Schraiber et al. [2013]) and then recognizing that  
 298 a Bessel(4) process is the same as a Bessel(0) process conditioned to never hit 0 and hence  
 299 has the same bridges – in the language of the general theory of Markov processes, the  
 300 Bessel(0) and Bessel(4) processes are Doob  $h$ -transforms of each other and it is well-known  
 301 that processes related in this way share the same bridges. We denote by  $(Y'_\tau)_{\tau \geq 0}$  the path  
 302 that has the proposed bridge spliced in between times  $f(s_1)$  and  $f(s_2)$  and coincides with  
 303  $(Y_\tau)_{\tau \geq 0}$  outside the interval  $[f(s_1), f(s_2)]$ .

304 In the Appendix, we show that the acceptance probability in this case is simply

$$(12) \quad \min \left\{ 1, \frac{L(D, (Y'_\tau)_{f(s_1) \leq \tau \leq f(s_2)} \mid \alpha_1, \alpha_2, t_0)}{L(D, (Y_\tau)_{f(s_1) \leq \tau \leq f(s_2)} \mid \alpha_1, \alpha_2, t_0)} \right\}.$$

305 Note that we only need to compute the likelihood ratio for the segment of transformed  
 306 path that changed between the times  $f(s_1)$  and  $f(s_2)$ .

307 **2.6.2. Allele age updates.** The first sample time with a non-zero count of the derived allele  
 308 (Figure 2b) is  $t_s$ , where  $s = \min\{i : c_i > 0\}$ . We must have  $t_0 < t_s$ . Along with proposing  
 309 a new value  $t'_0$  of the allele age  $t_0$  we will propose a new segment of the allele frequency  
 310 path from time  $t'_0$  to time  $t_s$ . Changing the allele age  $t_0$  to some new proposed value  $t'_0$   
 311 changes the definition of the function  $f$  in (6). Write  $f'(t) = \int_{t'_0}^t \frac{1}{\rho(s)} ds$ , where we stress  
 312 that the prime does not denote a derivative. The proposed transformed path  $Y'$  consists  
 313 of a new piece of path that goes from location 0 at time 0 to location  $Y_{f(t_s)}$  at time  $f'(t_s)$   
 314 and then has  $Y'_{f'(t)} = Y_{f(t)}$  for  $t \geq t_s$ . Recall that we use the improper prior  $\rho(t_0)$  for  $t_0$ ,  
 315 which reflects the fact that an allele is more likely to arise during times of large population  
 316 size [Slatkin, 2001]. In the Appendix, we show that the acceptance probability is

$$(13) \quad \min \left\{ 1, \frac{L(D, (Y'_\tau)_{0 \leq \tau \leq f'(t_s)} \mid \alpha_1, \alpha_2, t'_0) \psi(Y'_{f'(t_s)}; f'(t_s)) q(t_0|t'_0) \rho(t'_0)}{L(D, (Y_\tau)_{0 \leq \tau \leq f(t_s)} \mid \alpha_1, \alpha_2, t_0) \psi(Y_{f(t_s)}; f(t_s)) q(t'_0|t_0) \rho(t_0)} \right\}$$

317 where, in the notation of Subsection 2.4,

$$(14) \quad \psi(y; \epsilon) = \frac{y^2}{\epsilon^2} \exp \left\{ -\frac{y^2}{2\epsilon} \right\} = \frac{\mathbb{Q}^0(\bar{Y}_\epsilon \in dy)}{dy}$$

318 is the density of the so-called entrance law for the Bessel(0) process that appears in the  
 319 characterization of the  $\sigma$ -finite measure  $\mathbb{Q}^0$  and  $q(t'_0|t_0)$  is the proposal distribution of  $t'_0$   
 320 (in practice, we use a half-truncated normal distribution centered at  $t_0$ , with the upper  
 321 truncation occurring at the first time of non-zero observed allele frequency).

322 **2.6.3. Most recent allele frequency update.** While the allele frequency at sample times  
 323  $t_1, t_2, \dots, t_{k-1}$  are updated implicitly by the interior path update, we update the allele  
 324 frequency at the most recent sample time  $t_k$  separately (note that the most recent allele  
 325 frequency is not an additional parameter, but simply a random variable with a distribution  
 326 implied by the Wright-Fisher model on paths). We do this by first proposing a new allele  
 327 frequency  $Y'_{f(t_k)}$  and then proposing a new bridge from  $Y_{f(t_f)}$  to  $Y'_{f(t_k)}$  where  $t_f \in (t_{k-1}, t_k)$

328 is a fixed time (Figure 2c). If  $q(Y'_{f(t_k)} | Y_{f(t_k)})$  is the proposal density for  $Y'_{f(t_k)}$  given  $Y_{f(t_k)}$   
 329 (in practice, we use a truncated normal distribution centered at  $Y_{f(t_k)}$  and truncated at 0  
 330 and  $\pi$ ), then, arguing along the same lines as the interior path update and the allele age  
 331 update, we accept this update with probability

$$(15) \quad \min \left\{ 1, \frac{L(D, (Y'_\tau)_{f(t_f) \leq \tau \leq f(t_k)} | \alpha_1, \alpha_2, t_0) q(Y_{f(t_k)} | Y'_{f(t_k)}) Q(Y_{f(t_f)}, Y_{f(t_k)}; f(t_k) - f(t_f))}{L(D, (Y_\tau)_{f(t_f) \leq \tau \leq f(t_k)} | \alpha_1, \alpha_2, t_0) q(Y'_{f(t_k)} | Y_{f(t_k)}) Q(Y_{f(t_f)}, Y'_{f(t_k)}; f(t_k) - f(t_f))} \right\},$$

332 where

$$(16) \quad Q(x, y; t) = \frac{y}{t} \exp \left\{ -\frac{x^2 + y^2}{2t} \right\} I_1 \left( \frac{xy}{t} \right)$$

333 is the transition density of the Bessel(0) process (with  $I_1(\cdot)$  being the Bessel function of  
 334 the first kind with index 1) – see Knight [1981, Section 4.3.6]. Again, it is only necessary  
 335 to compute the likelihood ratio for the segment of transformed path that changed between  
 336 the times  $f(t_f)$  and  $f(t_k)$ .

337 **2.7. Updates to  $\alpha_1$  and  $\alpha_2$ .** Updates to  $\alpha_1$  and  $\alpha_2$  are conventional scalar parameter  
 338 updates. For example, letting  $q(\alpha'_1 | \alpha_1)$  be the proposal density for the new value of  $\alpha_1$ ,  
 339 we accept the new proposal with probability

$$\min \left\{ 1, \frac{L(D, (Y_\tau)_{\tau \geq 0} | \alpha'_1, \alpha_2, t_0) q(\alpha_1 | \alpha'_1) \pi(\alpha'_1, \alpha_2, t_0)}{L(D, (Y_\tau)_{\tau \geq 0} | \alpha_1, \alpha_2, t_0) q(\alpha'_1 | \alpha_1) \pi(\alpha_1, \alpha_2, t_0)} \right\}.$$

340 The acceptance probability for  $\alpha_2$  is similar. For both  $\alpha_1$  and  $\alpha_2$ , we use a heavy-tailed  
 341 Cauchy prior with median 0 and scale parameter 100, and we take the parameters  $\alpha_1, \alpha_2, t_0$   
 342 to be independent under the prior distribution. In addition, we use a normal proposal  
 343 distribution, centered around the current value of the parameter. Here, it is necessary to  
 344 compute the likelihood across the whole path.

345

### 3. RESULTS

346 We first test our method using simulated data to assess its performance and then apply  
 347 it to two real datasets from horses.

348 **3.1. Simulation performance.** To test the accuracy of our MCMC approach, we per-  
 349 formed two sets of simulations. First, we simulated data under a constant demographic  
 350 history to assess the quality of parameter inference under a simple model. Second, we  
 351 simulated data under the horse demographic history of Der Sarkissian et al. [2015] and  
 352 compared inferences performed with and without accounting for the demographic history.

353 In the constant demography simulations, we simulated allele frequency trajectories with  
 354 ages uniformly distributed between 0.1 and 0.3 diffusion time units ago, evolving with  $\alpha_1$   
 355 and  $\alpha_2$  uniformly distributed between 0 and 100. We simulate allele frequency trajectories  
 356 using an Euler approximation to the Wright-Fisher SDE (1) with  $\rho(t) \equiv 1$ . At each time  
 357 point between  $-0.4$  and  $0.0$  in steps of  $0.05$ , we simulated the sampling of 20 chromosomes.

358 We then ran the MCMC algorithm for 1,000,000 generations, sampling every 1000  
359 generations to obtain 1000 MCMC samples for each simulation. After discarding the first  
360 500 samples from each MCMC run as burn-in, we computed the effective sample size of  
361 the allele age estimate using the R package `coda` [Plummer et al., 2006]. For the analysis  
362 of the simulations, we only included simulations that had an effective sample size greater  
363 than 150 for the allele age, resulting in retaining 744 out of 1000 simulations.

364 Because our MCMC analysis provides a full posterior distribution on parameter val-  
365 ues, we summarized the results by computing the maximum *a posteriori* estimate of each  
366 parameter. We find that across the range of simulated  $\alpha_1$  values, estimation is quite ac-  
367 curate (Figure 3A). There is some downward bias for large true values of  $\alpha_1$ , indicating  
368 the influence of the prior. On the other hand, the strength of selection in favor of the ho-  
369 mozygote,  $\alpha_2$ , is less well estimated, with a more pronounced downward bias (Figure 3B).  
370 This is largely because most simulated alleles do not reach sufficiently high frequency for  
371 homozygotes to be common. Hence, there is very little information regarding the fitness of  
372 the homozygote. Allele age is estimated accurately, although there is a slight bias toward  
373 estimating a more recent age than the truth (Figure 3C).

374 [Figure 3 about here.]

375 When simulating under the horse demographic history, we drew 1000 allele ages with  
376 probability proportional to  $\rho(t)$  for  $t$  between 0.1 and 0.3 diffusion time units ago. Similarly  
377 to the simulations with constant demography, we drew  $\alpha_1$  and  $\alpha_2$  uniformly between 0 and  
378 100), and then simulated allele frequency trajectories using an Euler approximation to (1)  
379 with  $\rho(t)$  given by the history inferred by Der Sarkissian et al. [2015]. The sampling scheme  
380 is identical to the constant demography simulations.

381 We ran our simulated data through two separate MCMC pipelines, one accounting for  
382 the true simulated demographic history, and the other assuming a constant population  
383 size. All other settings were identical to the analysis of the data simulated under constant  
384 demography. We retained MCMC runs where the sampling likelihood, path likelihood,  $\alpha_1$   
385 estimate,  $\alpha_2$  estimate, and allele age estimate all had effective sample sizes greater than  
386 50, resulting in 561 analyses retained from the inference with variable demography, 647  
387 analyses retained from the inference with constant demography, and 454 analyses that were  
388 retained in both.

To quantify the overall impact of demographic model misspecification on parameter  
inference, we approximated the posterior root mean square error of a parameter (generically  
 $\theta$ ) by averaging over the posterior distribution,

$$\begin{aligned} RMSE(\theta) &= \left( \int (\hat{\theta} - \theta)^2 P(\hat{\theta}|D) d\hat{\theta} \right)^{\frac{1}{2}} \\ &\approx \left( \frac{1}{N} \sum_i (\hat{\theta}_i - \theta)^2 \right)^{\frac{1}{2}}, \end{aligned}$$

389 where the sum is over retained MCMC samples.

390 We found substantially smaller RMSE for inference of  $\alpha_1$  when demography is properly  
391 modeled (Figure 4). While inference of  $\alpha_2$  was similar between the two models, there is  
392 somewhat larger RMSE when demography is incorrectly assumed to be constant. Interest-  
393 ingly, there seem to be two regimes of error in allele age estimation: for the most recent  
394 allele ages, modeling demography results in higher RMSE, while for more ancient ages,  
395 inferences with constant population size result in larger RMSE. These are likely caused by  
396 a particular feature of this demographic model, which is a very strong bottleneck inferred  
397 in the recent past. Because alleles are more likely to arise during periods of larger popula-  
398 tion size, accounting for demographic history extends the tail of the posterior distribution  
399 further into the past, when the population was larger.

400 [Figure 4 about here.]

401 **3.2. Application to ancient DNA.** We applied our approach to real data by reanalyzing  
402 the MC1R and ASIP data from Ludwig et al. [2009]. In contrast to earlier analyses of these  
403 data, we explicitly incorporated the demography of the domesticated horse, as inferred  
404 by Der Sarkissian et al. [2015], using a generation time of 8 years. Table 2 shows the  
405 sample configurations and sampling times corresponding to each locus, where diffusion  
406 units are scaled to  $2N_0$ , with  $N_0 = 16000$  being the most recent effective size reported  
407 by Der Sarkissian et al. [2015]. For comparison, we also analyzed the data assuming the  
408 population size has been constant at  $N_0$ .

409 [Table 2 about here.]

410 [Figure 5 about here.]

411 With the MC1R locus, we found that posterior inferences about selection coefficients  
412 can be strongly influenced by whether or not demographic information is included in the  
413 analysis (Figure 5). Marginally, we see that incorporating demographic information results  
414 in an inference that  $\alpha_1$  is larger than the constant-size model (MAP estimates of 267.6 and  
415 74.1, with and without demography, respectively; Figure 5A), while  $\alpha_2$  is inferred to be  
416 smaller (MAP estimates of 59.1 and 176.2, with and without demography, respectively;  
417 Figure 5B). This has very interesting implications for the mode of selection inferred on the  
418 MC1R locus. Recall that  $\alpha_2 > \alpha_1 > 0$  is directional selection, in which the derived allele  
419 is always beneficial,  $\alpha_2 < \alpha_1 > 0$  is overdominant selection, in which the heterozygote  
420 is favored, and  $\alpha_2 > \alpha_1 < 0$  is underdominant selection, in which the heterozygote is  
421 disfavored. With constant demography, the trajectory of the allele is estimated to be shaped  
422 by positive directional selection (joint MAP,  $\alpha_1 = 87.6$ ,  $\alpha_2 = 394.8$ ; Figure 5C), while when  
423 demographic information is included, selection is inferred to act in an overdominant fashion  
424 (joint MAP,  $\alpha_1 = 262.5$ ,  $\alpha_2 = 128.1$ ; Figure 5D).

425 [Figure 6 about here.]

426 Incorporation of demographic history also has substantial impacts on the inferred distri-  
427 bution of allele ages (Figure 6). Most notably, the distribution of the allele age for MC1R  
428 is significantly truncated when demography is incorporated, in a way that correlates to  
429 the demographic events (Figure S1). While both the constant-size history and the more  
430 complicated history result in a posterior mode at approximately the same value of the



431 allele age, the domestication bottleneck inferred by Der Sarkissian et al. [2015] makes it  
432 far less likely that the allele rose more anciently than the recent population expansion.  
433 Because the allele is inferred to be younger under the model incorporating demography,  
434 the strength of selection in favor of the homozygote must be higher to allow it to escape  
435 low frequency quickly and reach the observed allele frequencies. Hence,  $\alpha_1$  is inferred to  
436 be much higher when demographic history is explicitly modeled.

437 [Figure 7 about here.]

438 Incorporation of demographic history has an even more significant impact on inferences  
439 made about the ASIP locus (Figure 7). Most strikingly, while  $\alpha_1$  is inferred to be very  
440 large without demography, it is inferred to be close to 0 when demography is incorporated  
441 (MAP estimates of 16.3 and 159.9 with and without demography, respectively; Figure  
442 7A). On the other hand, inference of  $\alpha_2$  is largely unaffected (MAP estimates of 34.7  
443 and 39.8 with and without demography, respectively; Figure 7B). Interestingly, this has  
444 an opposite implication for the mode of selection compared to the results for the MC1R  
445 locus. With a constant-size demographic history, the allele is inferred to have evolved  
446 under overdominance (joint MAP,  $\alpha_1 = 153.3$ ,  $\alpha_2 = 47$ ; Figure 7C), but when the more  
447 complicated demography is modeled, the allele frequency trajectory is inferred to be shaped  
448 by positive, nearly additive, selection (joint MAP,  $\alpha_1 = 16.4$ ,  $\alpha_2 = 46.8$ ; Figure 7D).

449 [Figure 8 about here.]

450 Incorporating demography has a similarly opposite effect on inference of allele age (Fig-  
451 ure 8). In particular, the allele is inferred to be much older when demography is modeled,  
452 and features a multi-modal posterior distribution on allele age, with each mode corre-  
453 sponding to a period of historically larger population size (Figure S2). Because the allele  
454 is inferred to be substantially older when demography is modeled, selection in favor of the  
455 heterozygote must have been weaker than would be inferred with the much younger age.  
456 Hence, the mode of selection switches from one of overdominance in a constant demography  
457 to one in which the homozygote is more fit than the heterozygote.

#### 458 4. DISCUSSION

459 Using DNA from ancient specimens, we have obtained a number of insights into evolu-  
460 tionary processes that were previously inaccessible. One of the most interesting aspects of  
461 ancient DNA is that it can provide a *temporal* component to evolution that has long been  
462 impossible to study. In particular, instead of making inferences about the allele frequencies,  
463 we can directly measure these quantities. To take advantage of this new data, we developed  
464 a novel Bayesian method for inferring the intensity and direction of natural selection from  
465 allele frequency time series. In order to circumvent the difficulties inherent in calculat-  
466 ing the transition probabilities under the standard Wright-Fisher process of selection and  
467 drift, we used a data augmentation approach in which we learn the posterior distribution  
468 on allele frequency paths. Doing this not only allows us to efficiently calculate likelihoods,  
469 but provides an unprecedented glimpse at the historical allele frequency dynamics.

470 The key innovation of our method is to apply high-frequency path augmentation meth-  
471 ods [Roberts and Stramer, 2001] to analyze genetic time series. The logic of the method is

472 similar to the logic of a path integral, in which we average over all possible allele frequency  
473 trajectories that are consistent with the data [Schraiber, 2014]. By choosing a suitable  
474 reference probability distribution against which to compute likelihood ratios, we were able  
475 to adapt these methods to infer the age of alleles and properly account for variable popu-  
476 lation sizes through time. Moreover, because of the computational advantages of the path  
477 augmentation approach, we were able to infer a model of general diploid selection. To  
478 our knowledge, ours is the first work that can estimate both allele age and general diploid  
479 selection while accounting for demography.

480 Using simulations, we showed that our method performs well for strong selection and  
481 densely sampled time series. However, it is worth considering the work of Watterson [1979],  
482 who showed that even knowledge of the full trajectory results in very flat likelihood surfaces  
483 when selection is not strong. This is because for weak selection, the trajectory is extremely  
484 stochastic and it is difficult to disentangle the effects of drift and selection [Schraiber et al.,  
485 2013].

486 We also used simulations to test how misspecification of demographic history impacts  
487 inference. We saw substantially increased posterior root mean square error in inference  
488 of selection parameters if demographic history is misspecified. To examine the impact of  
489 demographic history in the context of real data, we then applied our method to a classic  
490 dataset from horses. We found that our inference of both the strength and mode of natural  
491 selection depended strongly on whether or not we incorporated demography. For the MC1R  
492 locus, a constant-size demographic model results in an inference of positive selection, while  
493 the more complicated demographic model inferred by Der Sarkissian et al. [2015] causes the  
494 inference to tilt toward overdominance, as well as a much younger allele age. In contrast,  
495 the ASIP locus is inferred to be overdominant under a constant-size demography, but the  
496 complicated demographic history results in an inference of positive selection, and a much  
497 older allele age.

498 These results stand in contrast to those of Steinrücken et al. [2014], who found that  
499 the most likely mode of evolution for both loci under a constant demographic history  
500 is one of overdominance. There are a several reasons for this discrepancy. First, we  
501 computed the diffusion time units differently, using  $N_0 = 16000$  and a generation time of  
502 8 years, as inferred by Der Sarkissian et al. [2015], while Steinrücken et al. [2014] used  
503  $N_0 = 2500$  (consistent with the bottleneck size found by Der Sarkissian et al. [2015]) and  
504 a generation time of 5 years. Hence, our constant-size model has far less genetic drift  
505 than the constant-size model assumed by Steinrücken et al. [2014]. This emphasizes the  
506 importance of inferring appropriate demographic scaling parameters, even when a constant  
507 population size is assumed. Secondly, we use MCMC to integrate over the distribution of  
508 allele ages, which can have a very long tail going into the past, while Steinrücken et al.  
509 [2014] assume a fixed allele age.

510 One key limitation of this method is that it assumes that the aDNA samples all come  
511 from the same, continuous population. If there is in fact a discontinuity in the populations  
512 from which alleles have been sampled, this could cause rapid allele frequency change and  
513 create spurious signals of natural selection. Several methods have been devised to test this  
514 hypothesis [Sjödín et al., 2014], and one possibility would be to apply these methods to

515 putatively neutral loci sampled from the same individuals, thus determining which samples  
516 form a continuous population. Alternatively, if our method is applied to a number of loci  
517 throughout the genome and an extremely large portion of the genome is determined to  
518 be evolving under selection, this could be evidence for model misspecification and suggest  
519 that the samples do not come from a continuous population.

520 An advantage of the method that we introduced is that it may be possible to extend it to  
521 incorporate information from linked neutral diversity. In general, computing the likelihood  
522 of neutral diversity linked to a selected site is difficult and many have used Monte Carlo  
523 simulation and importance sampling [Slatkin, 2001, Coop and Griffiths, 2004, Chen and  
524 Slatkin, 2013]. These approaches average over allele frequency trajectories in much same  
525 way as our method; however, each trajectory is drawn completely independently of the  
526 previous trajectories. Using a Markov chain Monte Carlo approach, as we do here, has the  
527 potential to ensure that only trajectories with a high posterior probability are explored  
528 and hence greatly increase the efficiency of such approaches.

## 529 5. ACKNOWLEDGMENTS

530 We are grateful for helpful comments and discussion with Yun Song, Matthias Stein-  
531 rucken, and Anand Bhaskar during the conception and implementation of this work. We  
532 would also like thank 2 anonymous reviewers for their helpful comments that improved the  
533 clarity and thoroughness of this manuscript.

## 534 6. SOFTWARE AVAILABILITY

535 C++ software implementing the method described in this manuscript is freely available  
536 under a GNU Public License at <https://github.com/Schraiber/selection>.

## 537 7. APPENDIX

### 538 7.1. A proper posterior in the limit as the initial allele frequency approaches 0.

539 For reasons that we explain in Subsection 2.4, we re-parametrize our model by replacing  
540 the path variable  $(X_t)_{t \geq t_0}$  with a deterministic time and space transformation of it  $(Y_t)_{t \geq 0}$   
541 that takes values in the interval  $[0, \pi]$  with the boundary point 0 (resp.  $\pi$ ) for  $(Y_t)_{t \geq 0}$   
542 corresponding to the boundary point 0 (resp. 1) for  $(X_t)_{t \geq t_0}$ . The transformation producing  
543  $(Y_t)_{t \geq 0}$  is such that  $(X_t)_{t \geq t_0}$  can be recovered from  $(Y_t)_{t \geq 0}$  and  $t_0$ .

544 Implicit in our set-up is the initial frequency  $x_0$  at time  $t_0$  which corresponds to an  
545 initial value  $y_0$  at time 0 of the transformed process  $(Y_t)_{t \geq 0}$ . For the moment, let us make  
546 the dependence on  $y_0$  explicit by including it in relevant notation as a superscript. For  
547 example,  $\mathbb{P}^{y_0}(\cdot | \alpha_1, \alpha_2, t_0)$  is the prior distribution of  $(Y_t)_{t \geq 0}$  given the specified values of the  
548 other parameters  $\alpha_1, \alpha_2, t_0$ . We will construct a tractable “reference” process  $(\tilde{Y}_t)_{t \geq 0}$  with  
549 distribution  $\mathbb{Q}^{y_0}(\cdot)$  such that the probability distribution  $\mathbb{P}^{y_0}(\cdot | \alpha_1, \alpha_2, t_0)$  has a density  
550 with respect to  $\mathbb{Q}^{y_0}(\cdot)$  – explicitly,  $\mathbb{Q}^{y_0}(\cdot)$  is the distribution of a Bessel(0) process started  
551 at location  $y_0$  at time 0. That is, there is a function  $\Phi^{y_0}(\cdot; \alpha_1, \alpha_2, t_0)$  on path space such  
552 that

$$(17) \quad \mathbb{P}^{y_0}(dy | \alpha_1, \alpha_2, t_0) = \Phi^{y_0}(y; \alpha_1, \alpha_2, t_0) \mathbb{Q}^{y_0}(dy)$$

553 for a path  $(y_t)_{t \geq 0}$ . Assuming that  $\pi$  has a density with respect to Lebesgue measure which,  
 554 with a slight abuse of notation, we also denote by  $\pi$ , the outcome of our Bayesian inferential  
 555 procedure is determined by the ratios

$$(18) \quad \frac{\mathbb{P}(dD | y^{**}, t_0^{**}) \Phi^{y_0}(y^{**}; \alpha_1^{**}, \alpha_2^{**}, t_0^{**}) \pi(\alpha_1^{**}, \alpha_2^{**}, t_0^{**})}{\mathbb{P}(dD | y^*, t_0^*) \Phi^{y_0}(y^*; \alpha_1^*, \alpha_2^*, t_0^*) \pi(\alpha_1^*, \alpha_2^*, t_0^*)}$$

556 for pairs of augmented parameter values  $(y^*, \alpha_1^*, \alpha_2^*, t_0^*)$  and  $(y^{**}, \alpha_1^{**}, \alpha_2^{**}, t_0^{**})$  (*i.e.* the  
 557 Metropolis-Hastings ratio).

558 Under the probability measure  $\mathbb{P}^{y_0}(\cdot | \alpha_1, \alpha_2, t_0)$ , the process  $(Y_t)_{t \geq 0}$  converges in distri-  
 559 bution as  $y_0 \downarrow 0$  (equivalently,  $x_0 \downarrow 0$ ) to the trivial process that starts at location 0 at time  
 560 0 and stays there. However, for all  $\epsilon > 0$  the conditional distribution of  $(Y_t)_{t \geq \epsilon}$  under the  
 561 probability measure  $\mathbb{P}^{y_0}(\cdot | \alpha_1, \alpha_2, t_0)$  given the event  $\{Y_\epsilon > 0\}$  converges to a non-trivial  
 562 probability measure as  $y_0 \downarrow 0$ . Similarly, the conditional distribution of the reference  
 563 diffusion process  $(\bar{Y}_t)_{t \geq \epsilon}$  under the probability measure  $\mathbb{Q}^{y_0}(\cdot)$  given the event  $\{\bar{Y}_\epsilon > 0\}$   
 564 converges as  $y_0 \downarrow 0$  to a non-trivial limit. There are  $\sigma$ -finite measures  $\mathbb{P}^0(\cdot | \alpha_1, \alpha_2, t_0)$  and  
 565  $\mathbb{Q}^0(\cdot)$  on path space that both have infinite total mass, are such that for any  $\epsilon > 0$  both of  
 566 these measures assign finite, non-zero mass to the set of paths that are strictly positive at  
 567 the time  $\epsilon$ , and the corresponding conditional probability measures are the limits as  $y_0 \downarrow 0$   
 568 of the conditional probability measures described above. Moreover, there is a function  
 569  $\Phi^0(\cdot; \alpha_1, \alpha_2, t_0)$  on path space such that

$$(19) \quad \mathbb{P}^0(dy | \alpha_1, \alpha_2, t_0) = \Phi^0(y; \alpha_1, \alpha_2, t_0) \mathbb{Q}^0(dy).$$

570 The posterior distribution (3) converges to

$$(20) \quad \mathbb{P}^0(d\alpha_1, d\alpha_2, dt_0; dY | D) = \frac{\mathbb{P}(dD | Y, t_0) \mathbb{P}^0(dY | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0)}{\int \mathbb{P}(dD | Y') \mathbb{P}^0(dY' | \alpha'_1, \alpha'_2, t'_0) \pi(d\alpha'_1, d\alpha'_2, dt'_0)}.$$

571 Thus, the limit as  $y_0 \downarrow 0$  of a Bayesian inferential procedure for the augmented set of  
 572 parameters can be viewed as a Bayesian inferential procedure with the improper prior  
 573  $\mathbb{P}^0(dY | \alpha_1, \alpha_2, t_0) \pi(d\alpha_1, d\alpha_2, dt_0)$  for the parameters  $Y, \alpha_1, \alpha_2, t_0$ . In particular, the limit-  
 574 ing Bayesian inferential procedure is determined by the ratios

$$(21) \quad \frac{\mathbb{P}(dD | y^{**}, t_0^{**}) \Phi^0(h^{**}; \alpha_1^{**}, \alpha_2^{**}, t_0^{**}) \pi(\alpha_1^{**}, \alpha_2^{**}, t_0^{**})}{\mathbb{P}(dD | y^*, t_0^*) \Phi^0(y^*; \alpha_1^*, \alpha_2^*, t_0^*) \pi(\alpha_1^*, \alpha_2^*, t_0^*)}$$

575 for pairs of augmented parameter values  $(y^*, \alpha_1^*, \alpha_2^*, t_0^*)$  and  $(y^{**}, \alpha_1^{**}, \alpha_2^{**}, t_0^{**})$ .

576 **7.2. The likelihood of the data and the path.** Write  $\tau_i = f(t_i)$ . Note that  $\tau_0 =$   
 577  $f(t_0) = 0$ . Using equation (9), the density of the distribution of the transformed allele  
 578 frequency process  $(Y_t)_{0 \leq s \leq \tau_k}$  against the reference distribution of the Bessel(0) process  
 579  $(\bar{Y}_s)_{0 \leq s \leq \tau_k}$  when  $Y_0 = \bar{Y}_0 = y_0$  can be written

$$(22) \quad \exp \left\{ \int_0^{\tau_k} (a(Y_r, r) - b(Y_r)) dY_r - \frac{1}{2} \int_0^{\tau_k} (a^2(Y_r, r) - b^2(Y_r)) dr \right\}$$

580 where

$$a(y, \tau) = -\frac{1}{2} \cot(Y_\tau) + \frac{1}{4} (\rho(f^{-1}(\tau)) \sin(y) (\alpha_2 + (2\alpha_1 - \alpha_2) \cos(y)))$$

581 is the infinitesimal mean of the transformed Wright-Fisher process and

$$b(y) = -\frac{1}{2y}$$

is the infinitesimal mean of the Bessel(0) process. However, as shown by Sermaidis et al. [2013], attempting to approximate the Itô integral in (22) using a discrete representation of the path can lead to biased estimates of the posterior distribution. Instead, consider the potential functions

$$\begin{aligned} H_1(y, \tau) &= \int^y a(\xi, \tau) d\xi \\ &= -\frac{1}{8} (\rho(f^{-1}(\tau)) \cos^2(y)(2\alpha_1 - \alpha_2) + 4 \log(\sin(y))) \end{aligned}$$

and

$$\begin{aligned} H_2(y) &= \int^y b(\xi, \tau) d\xi \\ &= -\frac{\log(y)}{2}. \end{aligned}$$

If we assume that  $\rho$  is continuous (not merely right continuous with left limits), then Itô's lemma shows that we can write

$$\begin{aligned} \int_0^{\tau_k} (a(Y_r, r) - b(Y_r)) dY_r &= H_1(Y_{\tau_k}, \tau_k) - H_2(Y_{\tau_k}) - (H_1(Y_0, 0) - H_2(Y_0)) \\ &\quad - \int_0^{\tau_k} \left( \frac{\partial H_1}{\partial \tau}(Y_r, r) - \frac{\partial H_2}{\partial \tau}(Y_r) \right) dr \\ &\quad - \int_0^{\tau_k} \left( \frac{\partial^2 H_1}{\partial y^2}(Y_r, r) - \frac{\partial^2 H_2}{\partial y^2}(Y_r) \right) dr. \end{aligned}$$

582 To generalize this to the case where  $\rho$  is right continuous with left limits, write

$$\int_0^{\tau_k} (a(Y_r, r) - b(Y_r)) dY_r = I_0 + \sum_{i=m}^K I_i,$$

583 where  $m$  and  $K$  are defined in the main text,

$$I_0 = \lim_{\tau \uparrow f(d_m)} \int_0^{\tau} (a(Y_r, r) - b(Y_r)) dY_r,$$

584 for  $m < i < K$ ,

$$I_i = \lim_{\tau \uparrow f(d_{i+1})} \int_{f(d_i)}^{\tau} (a(Y_r, r) - b(Y_r)) dY_r,$$

585 and

$$I_K = \lim_{\tau \uparrow \tau_k} \int_{f(d_K)}^{\tau} (a(Y_r, r) - b(Y_r)) dY_r.$$

586 Itô's lemma can then be applied to each segment in turn. Following the conversion of the  
 587 Itô integrals into ordinary Lebesgue integrals, making the substitution  $s = f^{-1}(r)$  results  
 588 in the path likelihood displayed in (11).

589 **7.3. Acceptance probability for an interior path update.** When we propose a new  
 590 path  $(y'_t)_{0 \leq t \leq \tau_k}$  to update the current path  $(y_t)_{0 \leq t \leq \tau_k}$  which doesn't hit the boundary, the  
 591 new path agrees with the existing path outside some time interval  $[v_1, v_2]$ , and has a new  
 592 segment spliced in that goes from  $y_{v_1}$  at time  $v_1$  to  $y_{v_2}$  at time  $v_2$ . The proposed new path  
 593 segment comes from a Bessel(0) process over the time interval  $[v_1, v_2]$  that is pinned to  
 594 take the values  $y_{v_1}$  and  $y_{v_2}$  at the end-points; that is, the proposed new piece of path is a  
 595 bridge.

596 The ratio that determines the probability of accepting the proposed path is

$$(23) \quad \frac{P(dD | y', t_0)}{P(dD | y, t_0)} \times \frac{\mathbb{P}(dy') \kappa(dy | y')}{\mathbb{P}(dy) \kappa(dy' | y)},$$

597 where  $P(\cdot | y', t_0)$  and  $P(\cdot | y, t_0)$  give the probability of the observed allele counts given  
 598 the transformed allele frequency paths and initial time  $t_0$ ,  $\mathbb{P}(\cdot)$  is the distribution of the  
 599 transformed Wright-Fisher diffusion starting from  $y_0 > 0$  at time 0 (that is, the distribution  
 600 we have sometimes denoted by  $\mathbb{P}^{y_0}$ ), the probability kernel  $\kappa(\cdot | y)$  gives the distribution of  
 601 the proposed path when the current path is  $y$ , and  $\kappa(\cdot | y')$  is similar. To be completely  
 602 rigorous, the second term in the product in (23) should be interpreted as the Radon-  
 603 Nikodym derivative of two probability measures on the product of path space with itself.

604 Consider a finite set of times  $0 \equiv \tau_0 \equiv u_0 < u_1 < \dots < u_\ell \equiv \tau_k$ . Suppose that  
 605  $\{v_1, v_2\} \in \{u_0, \dots, u_\ell\}$ ,  $v_1 = u_m$  and  $v_2 = u_n$  for some  $m < n$ . Let  $(y_t)_{0 \leq t \leq \tau_k}$  and  $(y'_t)_{0 \leq t \leq \tau_k}$   
 606 be two paths that coincide on  $[0, v_1] \cup [v_2, \tau_k] = [u_0, u_m] \cup [u_n, u_\ell]$ . Write  $P(x, y; s, t)$  for  
 607 the transition density (with respect to Lebesgue measure) of the transformed Wright-  
 608 Fisher diffusion from time  $s$  to time  $t$  and  $Q(x, y; t)$  for the transition density (with respect  
 609 to Lebesgue measure) of the Bessel(0) process. Suppose that  $(\xi, \zeta)$  is a pair of random  
 610 paths with  $P((\xi, \zeta) \in (dy, dy')) = \mathbb{P}(dy) \kappa(dy' | y)$ . Then, writing  $z_t = y_t = y'_t$  for  $t \in$   
 611  $[0, v_1] \cup [v_2, \tau_k] = [u_0, u_m] \cup [u_n, u_\ell]$ , we have

$$\begin{aligned} & P(\xi_{u_1} \in dy_{u_1}, \dots, \xi_{u_\ell} \in dy_{u_\ell}, \zeta_{u_1} \in dy'_{u_1}, \dots, \zeta_{u_\ell} \in dy'_{u_\ell}) \\ &= P(z_{u_0}, z_{u_1}; u_0, u_1) dz_{u_1} \times \dots \times P(z_{u_{m-1}}, z_{u_m}; u_{m-1}, u_m) dz_{u_m} \\ & \quad \times P(z_{u_m}, y_{u_{m+1}}; u_m, u_{m+1}) dy_{u_{m+1}} \times \dots \times P(y_{u_{n-1}}, z_{u_n}; u_{n-1}, u_n) dz_{u_n} \\ & \quad \times P(z_{u_n}, z_{u_{n+1}}; u_n, u_{n+1}) dz_{u_{n+1}} \times \dots \times P(z_{u_{\ell-1}}, z_{u_\ell}; u_{\ell-1}, u_\ell) dz_{u_\ell} \\ & \quad \times Q(z_{u_m}, y'_{u_{m+1}}; u_{m+1} - u_m) dy_{u_{m+1}} \times \dots \times Q(y_{u_{n-1}}, z_{u_n}; u_n - u_{n-1}) \\ & \quad \Big/ Q(z_{u_m}, z_{u_n}; u_n - u_m), \end{aligned}$$

612 where the factor in the denominator arises because we are proposing *bridges* and hence  
 613 conditioning on going from a fixed location at  $v_1 = u_m$  to another fixed location at  $v_2 = u_n$ .



614 Thus,

$$\begin{aligned} & \frac{P(\xi_{u_1} \in dy'_{u_1}, \dots, \xi_{u_\ell} \in dy'_{u_\ell}, \zeta_{u_1} \in dy_{u_1}, \dots, \zeta_{u_\ell} \in dy_{u_\ell})}{P(\xi_{u_1} \in dy_{u_1}, \dots, \xi_{u_\ell} \in dy_{u_\ell}, \zeta_{u_1} \in dy'_{u_1}, \dots, \zeta_{u_\ell} \in dy'_{u_\ell})} \\ &= \frac{\prod_{j=m}^{n-1} P(y'_{u_j}, y'_{u_{j+1}}; u_j, u_{j+1})/Q(y'_{u_j}, y'_{u_{j+1}}; u_{j+1} - u_j)}{\prod_{j=m}^{n-1} P(y_{u_j}, y_{u_{j+1}}; u_j, u_{j+1})/Q(y_{u_j}, y_{u_{j+1}}; u_{j+1} - u_j)}. \end{aligned}$$

615 Therefore, the Radon-Nikodym derivative appearing in (23) is the ratio of Radon-Nikodym  
616 derivatives

$$\frac{\frac{d\tilde{\mathbb{P}}}{d\tilde{\mathbb{Q}}}(y')}{\frac{d\tilde{\mathbb{P}}}{d\tilde{\mathbb{Q}}}(y)},$$

617 where  $\tilde{\mathbb{P}}$  (resp.  $\tilde{\mathbb{Q}}$ ) is the distribution of the transformed Wright-Fisher diffusion (resp. the  
618 Bessel(0) process) started at location  $y_{v_1} = y'_{v_1}$  at time  $v_1$  and run until time  $v_2$ . The  
619 formula (12) for the acceptance probability associated with an interior path update follows  
620 immediately.

621 The above argument was carried out under the assumption that the transformed initial  
622 allele frequency  $y_0$  was strictly positive and so all the measures involved were probabil-  
623 ity measures. However, taking  $y_0 \downarrow 0$  we see that the formula (12) continues to hold.  
624 Alternatively, we could have worked directly with the measure  $\mathbb{P}^0$  in place of  $\mathbb{P}^{y_0}$ . The  
625 only difference is that we would have to replace  $P(y_0, y; 0, s)$  by the density  $\phi(y; 0, s)$  of an  
626 entrance law for  $\mathbb{P}^0$ . That is,  $\phi(y; 0, s)$  has the property that

$$\lim_{y_0 \downarrow 0} \frac{P(y_0, y'; 0, s')}{P(y_0, y''; 0, s'')} = \frac{\phi(y'; 0, s')}{\phi(y''; 0, s'')}$$

627 for all  $y', y'' > 0$  and  $s', s'' > 0$  so that

$$\int \phi(y; 0, s) P(y, z; s, t) dy = \phi(z; 0, t)$$

628 for  $0 < s < t$ . Such a density, and hence the corresponding entrance law, is unique up to  
629 a multiplicative constant. In any case, it is clear that the choice of entrance law in the  
630 definition of  $\mathbb{P}^0$  does not affect the formula (12) as the entrance law densities “cancels out”.

631 **7.4. Acceptance probability for an allele age update.** The argument justifying the  
632 formula (13) for the probability of accepting a proposed update to the allele age  $t_0$  is similar  
633 to the one just given for interior path updates. Now, however, we have to consider replacing  
634 a path  $y$  that starts from  $y_0$  at time 0 and runs until time  $f(t_k)$  with a path  $y'$  that starts  
635 from  $y_0$  at time 0 and runs until time  $f'(t_k)$ . Instead of removing an internal segment of  
636 path and replacing it by one of the same length with the same values at the endpoints, we  
637 replace the initial segment of path that runs from time 0 to  $f(t_s) = \int_{t_0}^{t_s} \frac{1}{\rho(s)} ds$  by one that  
638 runs from time 0 to time  $f'(t_s) = \int_{t'_0}^{t_s} \frac{1}{\rho(s)} ds$ , with  $y'_{f'(t_s)} = y_{f(t_s)}$ .

639 By analogy with the previous subsection, we need to consider

$$\frac{P(\xi \in dy', T_0^\xi \in dt', \zeta \in dy, T_0^\zeta \in dt)}{P(\xi \in dy, T_0^\xi \in dt, \zeta \in dy', T_0^\zeta \in dt')}$$

640 where  $\xi$  is a transformed Wright-Fisher process starting at  $y_0$  at time 0 and run to time  
641  $F^\xi = \int_{T_0^\xi}^{t_s} \frac{1}{\rho(s)} ds$ , where  $P(T_0^\xi \in dt) = \rho(t) dt$ , and conditional on  $\xi$ ,  $\zeta$  is a Bessel(0)  
642 bridge run from  $y_0$  at time 0 to  $\xi_{F^\xi}$  at time  $F^\zeta = \int_{T_0^\zeta}^{t_s} \frac{1}{\rho(s)} ds$ , where  $P(T_0^\zeta \in dt) = \rho(t) dt$   
643 independent of  $\xi$  and  $T_0^\xi$ .

644 Suppose that  $0 = u_0 < u_1 < \dots < u_m = \int_{t'}^{t_s} \frac{1}{\rho(s)} ds$  and  $0 = v_0 < v_1 < \dots < v_n =$   
 645  $\int_t^{t_s} \frac{1}{\rho(s)} ds$ . We have for  $y'_0, \dots, y'_m$  and  $y_0, \dots, y_n$  with  $y_0 = y'_0$  and  $y'_m = y_n$  that

$$\begin{aligned}
 & \frac{P(\xi_{u_i} \in dy'_i, 1 \leq i \leq m-1, T_0^\xi \in dt', \zeta_{v_j} \in dy_j, 1 \leq j \leq n, T_0^\zeta \in dt)}{P(\xi_{v_j} \in dy_j, 1 \leq j \leq n-1, T_0^\xi \in dt, \zeta_{u_i} \in dy'_i, 1 \leq i \leq m, T_0^\zeta \in dt')} \\
 &= \left\{ \prod_{i=0}^{m-1} P(y'_j, y'_{j+1}; u_i, u_{i+1}) dy'_{i+1} \times \rho(t') dt' \right. \\
 & \quad \times \left[ \prod_{j=0}^{n-2} Q(y_j, y_{j+1}; v_{j+1} - v_j) dy_{j+1} \times Q(y_{n-1}, y_n; v_n - v_{n-1}) / Q(y_0, y_n; v_n) \right] \times dt \left. \right\} \\
 & \quad / \left\{ \prod_{j=0}^{n-1} P(y_j, y_{j+1}; v_j, v_{j+1}) dy_{j+1} \times \rho(t) dt \right. \\
 & \quad \times \left[ \prod_{i=0}^{m-2} Q(y'_i, y'_{i+1}; u_{i+1} - u_i) dy'_{i+1} \times Q(y'_{m-1}, y'_m; u_m - u_{m-1}) / Q(y'_0, y'_m; u_m) \right] \times dt' \left. \right\} \\
 &= \left\{ \prod_{i=0}^{m-1} P(y'_j, y'_{j+1}; u_i, u_{i+1}) dy'_{i+1} \times \rho(t') dt' \right. \\
 & \quad \times \left[ \prod_{j=0}^{n-1} Q(y_j, y_{j+1}; v_{j+1} - v_j) dy_{j+1} / Q(y_0, y_n; v_n) \right] \times dt \left. \right\} \\
 & \quad / \left\{ \prod_{j=0}^{n-1} P(y_j, y_{j+1}; v_j, v_{j+1}) dy_{j+1} \times \rho(t) dt \right. \\
 & \quad \times \left[ \prod_{i=0}^{m-1} Q(y'_i, y'_{i+1}; u_{i+1} - u_i) dy'_{i+1} / Q(y'_0, y'_m; u_m) \right] \times dt' \left. \right\} \\
 &= \frac{\prod_{i=0}^{m-1} P(y'_j, y'_{j+1}; u_i, u_{i+1}) dy'_{i+1} / \left[ \prod_{i=0}^{m-1} Q(y'_i, y'_{i+1}; u_{i+1} - u_i) dy'_{i+1} \right]}{\prod_{j=0}^{n-1} P(y_j, y_{j+1}; v_j, v_{j+1}) dy_{j+1} / \left[ \prod_{j=0}^{n-1} Q(y_j, y_{j+1}; v_{j+1} - v_j) dy_{j+1} \right]} \\
 & \quad \times \frac{Q(y'_0, y'_m; u_m)}{Q(y_0, y_n; v_n)} \times \frac{\rho(t')}{\rho(t)},
 \end{aligned}$$

646 where the second equality follows from the fact that  $y_n = y'_m$ .

647 Thus,

$$\begin{aligned} & \frac{P(\xi \in dy', T_0^\xi \in dt', \zeta \in dy, T_0^\zeta \in dt)}{P(\xi \in dy, T_0^\xi \in dt, \zeta \in dy', T_0^\zeta \in dt')} \\ &= \frac{\frac{d\hat{\mathbb{P}}}{d\hat{\mathbb{Q}}}(y')}{\frac{d\hat{\mathbb{P}}}{d\hat{\mathbb{Q}}}(y)} \times \frac{Q(y_0, y_{T'}; T')}{Q(y_0, y_T; T)} \times \frac{\rho(t')}{\rho(t)}, \end{aligned}$$

648 where  $T = \int_t^{t_s} \frac{1}{\rho(s)} ds$  and  $T' = \int_{t'}^{t'_s} \frac{1}{\rho(s)} ds$ ,  $\hat{\mathbb{P}}$  (resp.  $\check{\mathbb{P}}$ ) is the distribution of the transformed  
649 Wright-Fisher diffusion starting at location  $y_0$  at time 0 and run until time  $T$  (resp.  $T'$ ),  
650 and  $\hat{\mathbb{Q}}$  (resp.  $\check{\mathbb{Q}}$ ) is the distribution of the Bessel(0) process starting at location  $y_0$  at time  
651 0 and run until time  $T$  (resp.  $T'$ ).

652 We have thusfar assumed that  $y_0$  is strictly positive. As in the previous subsection,  
653 we can let  $y_0 \downarrow 0$  to get an expression in terms of Radon-Nikodym derivatives of  $\sigma$ -finite  
654 measures and the density  $\psi(y; s)$  of an entrance law for  $\mathbb{Q}^0$ . That is,  $\psi(y; s)$  has the property  
655 that

$$\lim_{y_0 \downarrow 0} \frac{Q(y_0, y'; s')}{Q(y_0, y''; s'')} = \frac{\psi(y'; s')}{\psi(y''; s'')}$$

656 for all  $y', y'' > 0$  and  $s', s'' > 0$ , so that

$$\int \psi(y; s) Q(y, z; t) dy = \psi(z; s + t)$$

657 for  $s, t > 0$ . Up to an irrelevant multiplicative constant,  $\psi$  is given by the expression (14),  
658 and the formula (13) for the acceptance probability follows immediately.

659 **7.5. Acceptance probability for a most recent allele frequency update.** The deriva-  
660 tion of formula (15) for the probability of accepting a proposed update to the most recent  
661 allele frequency is similar to those for the other acceptance probabilities (12) and (13), so  
662 we omit the details.

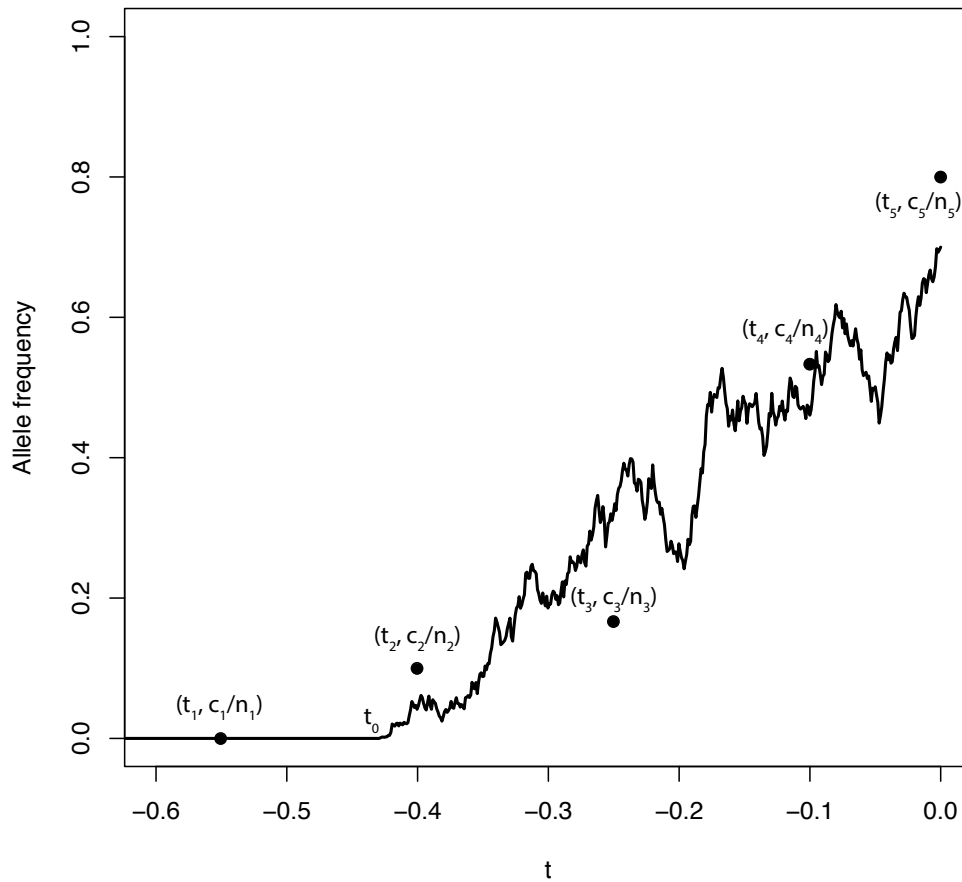


FIGURE 1. Taking samples from an allele frequency trajectory. An allele frequency trajectory is simulated from the Wright-Fisher diffusion (solid line). At each time,  $t_i$ , a sample of size  $n_i$  chromosomes is taken and  $c_i$  copies of the derived allele are observed. Each point corresponds to the observed allele frequency of sample  $i$ . Note that  $t_1$  is more ancient than the allele age,  $t_0$ .

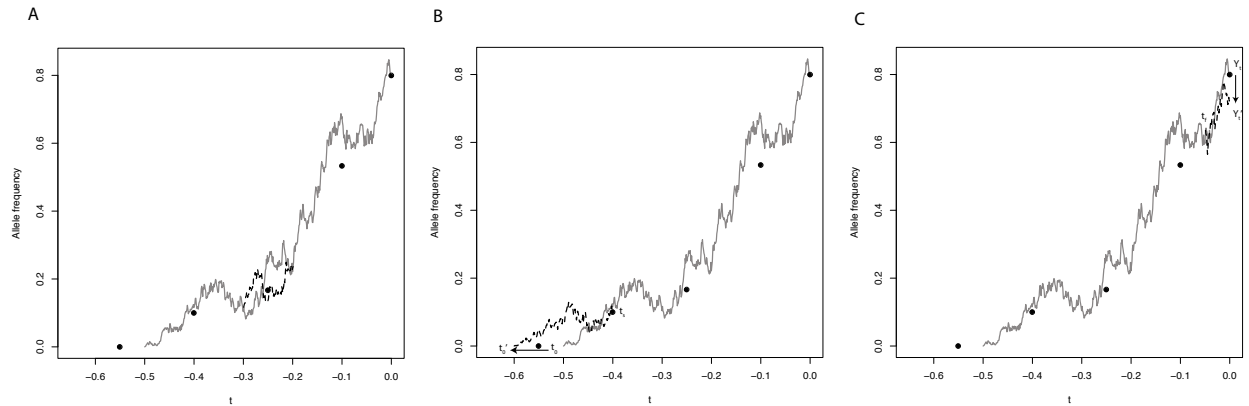


FIGURE 2. Illustration of path updates. Filled circles correspond to the same sample frequencies as in Figure 1. The solid gray line in each panel is the current allele frequency trajectory and the dashed black lines are the proposed updates. In panel a, an interior section of path is proposed between points  $s_1$  and  $s_2$ . In panel b, a new allele age,  $t'_0$  is proposed and a new path is drawn between  $t'_0$  and  $t_s$ . In panel c, a new most recent allele frequency  $Y'_{t_k}$  is proposed and a new path is drawn between  $t_f$  and  $t_k$ .



Figures

27

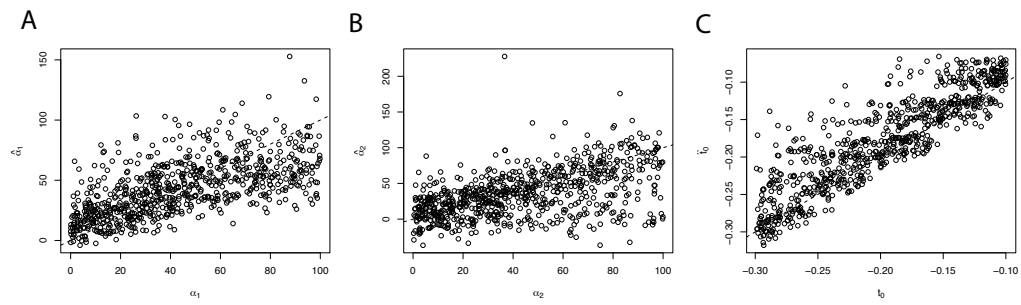


FIGURE 3. Maximum *a posteriori* estimates of different parameters. Each panel shows the true value of a parameter on the *x*-axis, while the inferred value is on the *y*-axis. Dashed line is  $y = x$ .

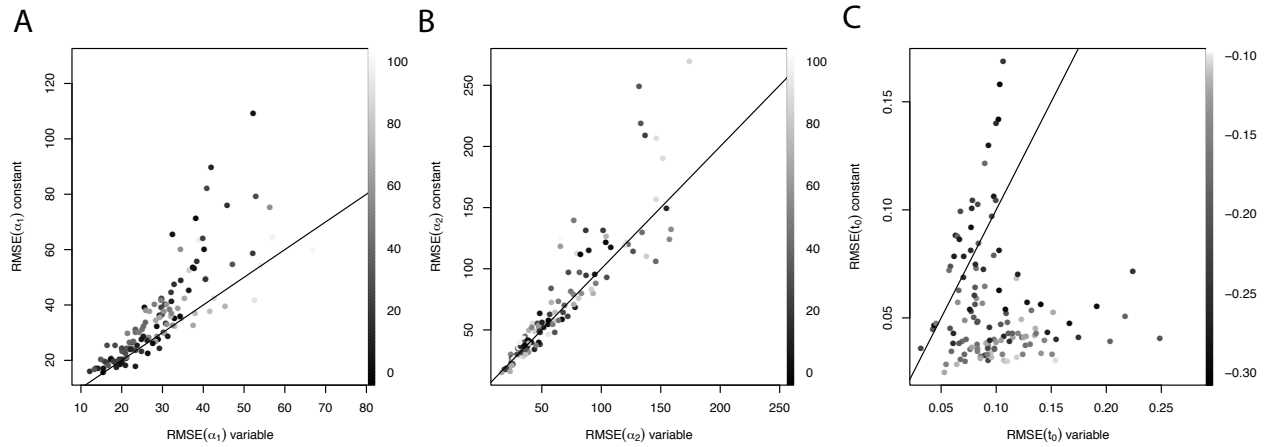


FIGURE 4. Comparison of root mean square error (RSME) when inference is performed with the proper (variable) demographic model on the  $x$ -axis compared to a misspecified constant demography model on the  $y$ -axis. Each point represents a single simulation, and points are colored according to simulated parameter value (scale on the right of each panel). Solid line is  $y = x$ .

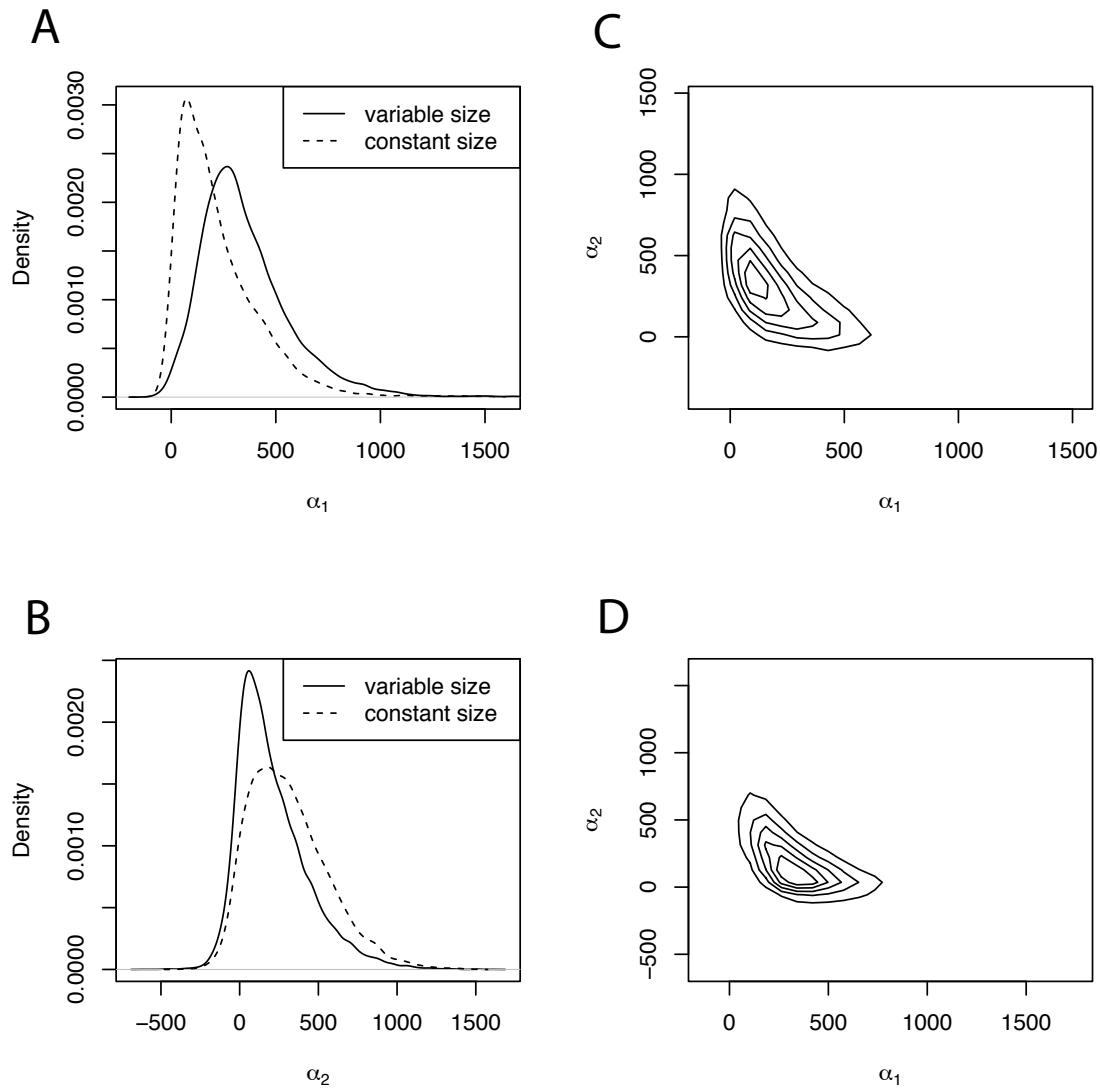


FIGURE 5. Posterior distributions of selection coefficients for the MC1R locus. Panels A and B show marginal distributions of  $\alpha_1$  and  $\alpha_2$ , respectively, with the solid line indicating the posterior obtained from an analysis including the full demographic history, and the dotted line showing what would be inferred in a constant size population. Panels C and D show contour plots of the joint distribution of  $\alpha_1$  and  $\alpha_2$  without and with demography, respectively.

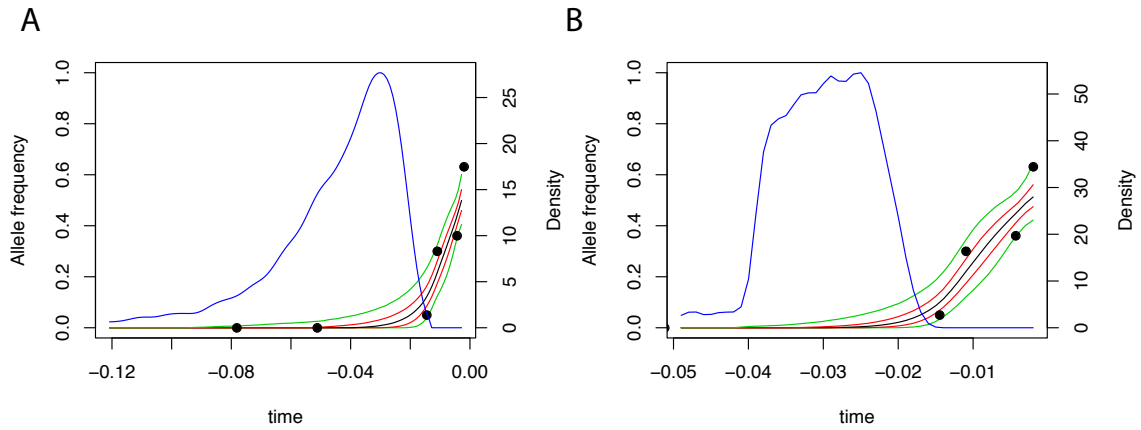


FIGURE 6. Posterior distribution on allele frequency paths for the MC1R locus. Each panel shows the sampled allele frequency data (filled circles), the point-wise median (black), 25 and 75% quantiles (red), and 5 and 95% quantiles (green) of the posterior distribution on paths, and the posterior distribution on allele age (blue). Panel A reports inference with constant demography, while panel B shows the result of inference with the full demographic history.

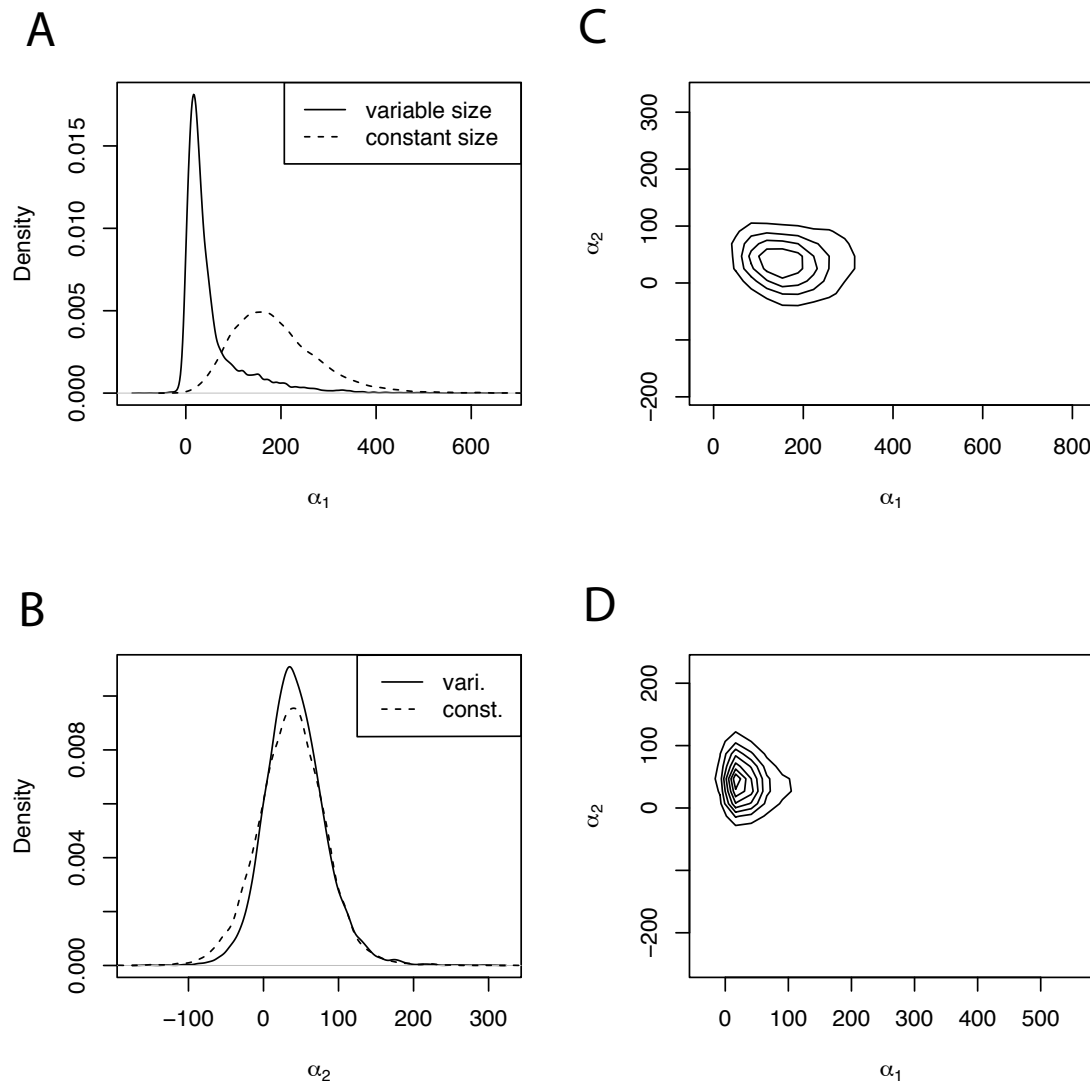


FIGURE 7. Posterior distributions of selection coefficients for the ASIP locus. Panels as in Figure 5

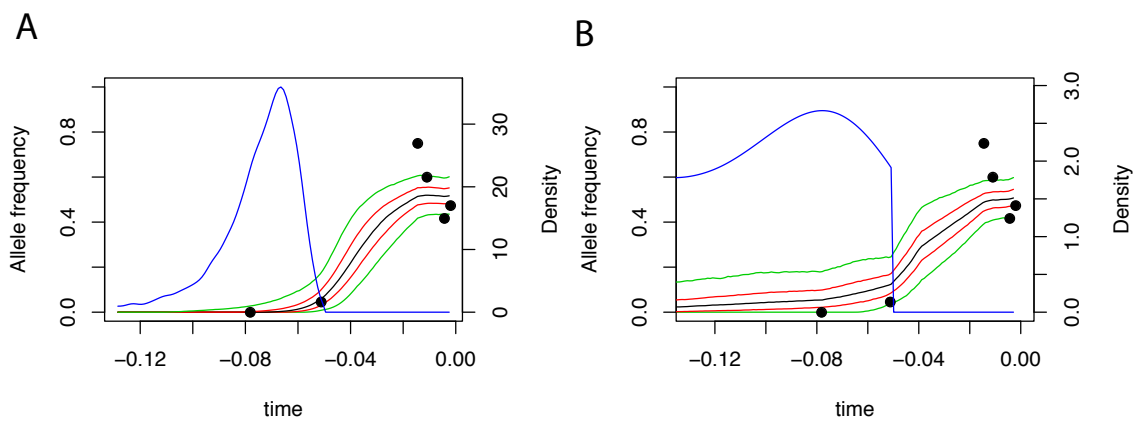


FIGURE 8. Posterior distribution on allele frequency paths for the ASIP locus. Panels are as in Figure 6.

Genotype	$A_1A_1$	$A_1A_0$	$A_0A_0$
Fitness	$1 + s_2$	$1 + s_1$	1

TABLE 1. Fitness scheme assumed in the text.



Sample time (years BCE)	20,000	13,100	3,700	2,800	1,100	500
Sample time (diffusion units)	0.078	0.051	0.014	0.011	0.004	0.002
Sample size	10	22	20	20	36	38
Count of ASIP alleles	0	1	15	12	15	18
Count of MC1R alleles	0	0	1	6	13	24

TABLE 2. Sample information for horse data. Diffusion time units are calculated assuming  $N_0 = 2500$  and a generation time of 5 years.

663

## 8. SUPPLEMENTARY FIGURES

664

[Figure S1 about here.]

665

[Figure S2 about here.]

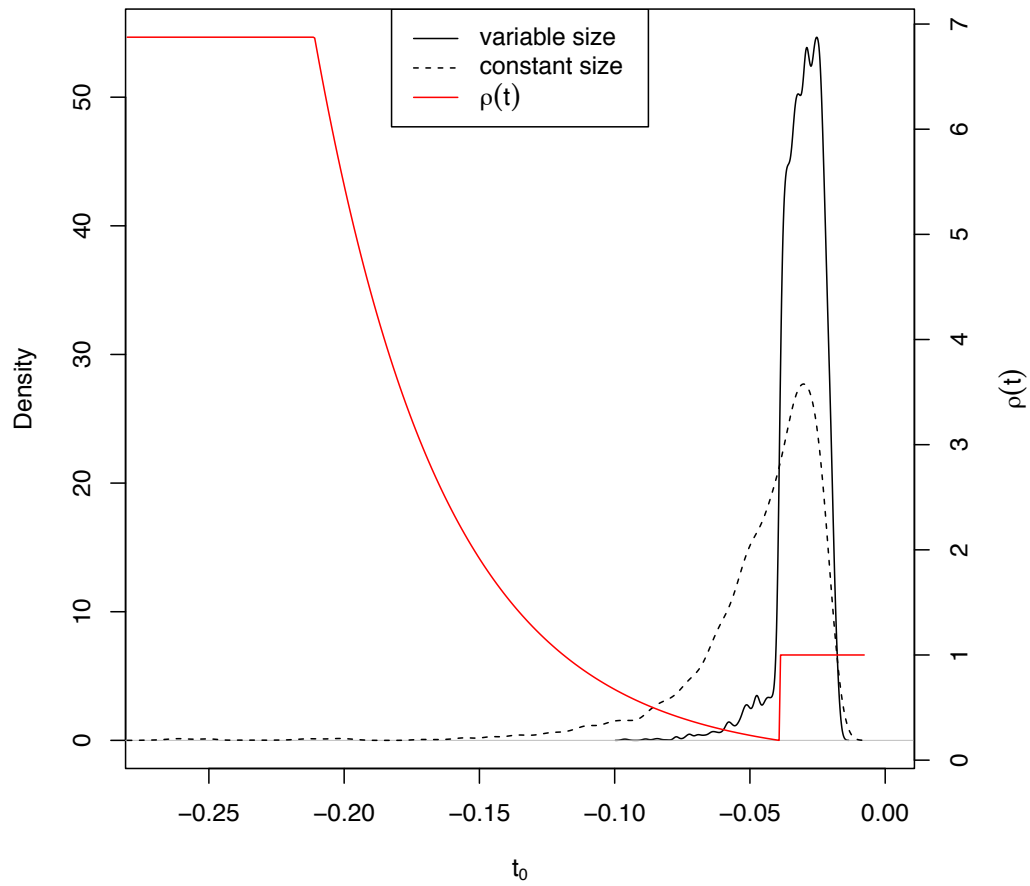


FIGURE S1. Influence of population size on age estimates of the MC1R locus. The solid and dashed lines show the posterior distribution on allele age with and without demography, respectively. In red, the demographic history inferred by Der Sarkissian et al. [2015].

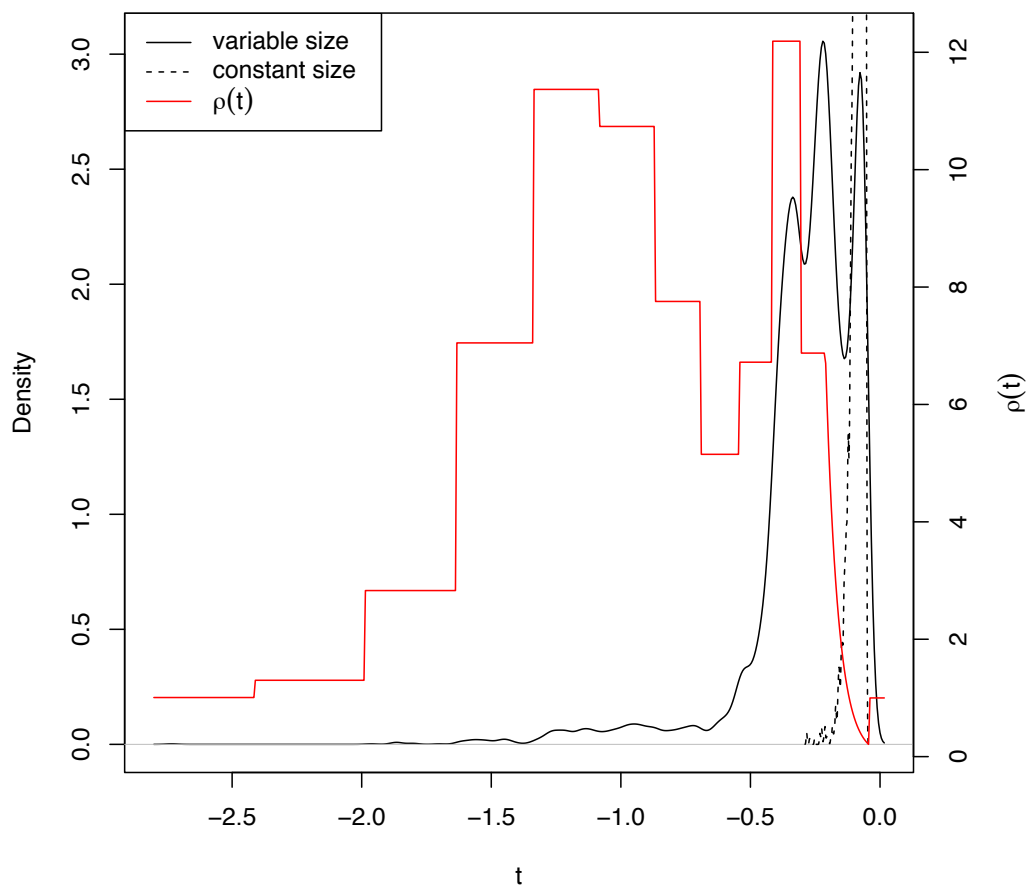


FIGURE S2. Influence of population size on age estimates of the ASIP locus. Data presented is as in Figure S1

666

REFERENCES

- 667 Jonathan P Bollback, Thomas L York, and Rasmus Nielsen. Estimation of  $2n\mu s$  from  
668 temporal allele frequency data. *Genetics*, 179(1):497–502, 2008.
- 669 Hua Chen and Montgomery Slatkin. Inferring selection intensity and allele age from mul-  
670 tilocus haplotype structure. *G3: Genes— Genomes— Genetics*, 3(8):1429–1442, 2013.
- 671 Graham Coop and Robert C Griffiths. Ancestral inference on gene trees under selection.  
672 *Theoretical population biology*, 66(3):219–232, 2004.
- 673 Clio Der Sarkissian, Luca Ermini, Mikkel Schubert, Melinda A Yang, Pablo Librado, Mat-  
674 teo Fumagalli, Hákon Jónsson, Gila Kahila Bar-Gal, Anders Albrechtsen, Filipe G Vieira,  
675 et al. Evolutionary genomics and conservation of the endangered przewalski’s horse.  
676 *Current Biology*, 25(19):2577–2583, 2015.
- 677 Warren J Ewens. *Mathematical population genetics: I. Theoretical introduction*, volume 27.  
678 Springer, 2004.
- 679 Alison F Feder, Sergey Kryazhimskiy, and Joshua B Plotkin. Identifying signatures of  
680 selection in genetic time series. *Genetics*, 196(2):509–522, 2014.
- 681 William Feller. Diffusion processes in genetics. In *Proc. Second Berkeley Symp. Math.*  
682 *Statist. Prob*, volume 227, page 246, 1951.
- 683 Ronald Aylmer Fisher. On the dominance ratio. *Proceedings of the royal society of Edin-*  
684 *burgh*, 42:321–341, 1922.
- 685 Christiane Fuchs. *Inference for Diffusion Processes: With Applications in Life Sciences*.  
686 Springer, 2013.
- 687 IV Girsanov. On transforming a certain class of stochastic processes by absolutely contin-  
688 uous substitution of measures. *Theory of Probability & Its Applications*, 5(3):285–301,  
689 1960.
- 690 Andrew Golightly and Darren J Wilkinson. Bayesian inference for stochastic kinetic models  
691 using a diffusion approximation. *Biometrics*, 61(3):781–788, 2005.
- 692 Andrew Golightly and Darren J Wilkinson. Bayesian inference for nonlinear multivariate  
693 diffusion models observed with error. *Computational Statistics & Data Analysis*, 52(3):  
694 1674–1693, 2008.
- 695 Robert C Griffiths and Simon Tavaré. Sampling theory for neutral alleles in a varying  
696 environment. *Philosophical Transactions of the Royal Society of London B: Biological*  
697 *Sciences*, 344(1310):403–410, 1994.
- 698 John Burdon Sanderson Haldane. A mathematical theory of natural and artificial selection,  
699 part v: selection and mutation. *Mathematical Proceedings of the Cambridge Philosophical*  
700 *Society*, 23(07):838–844, 1927.
- 701 Kiyosi Itô. Stochastic integral. *Proceedings of the Japan Academy, Series A, Mathematical*  
702 *Sciences*, 20(8):519–524, 1944.
- 703 Paul A Jenkins. Exact simulation of the sample paths of a diffusion with a finite entrance  
704 boundary. *arXiv preprint arXiv:1311.5777*, 2013.
- 705 Paul A Jenkins and Dario Spano. Exact simulation of the wright-fisher diffusion. *arXiv*  
706 *preprint arXiv:1506.06998*, 2015.

- 707 Olav Kallenberg. *Foundations of modern probability*. Probability and its Applications (New  
708 York). Springer-Verlag, New York, second edition, 2002. doi: 10.1007/978-1-4757-4015-8.  
709 URL <http://dx.doi.org/10.1007/978-1-4757-4015-8>.
- 710 Frank B Knight. *Essentials of Brownian motion and diffusion*, volume 18 of *Mathematical*  
711 *Surveys*. American Mathematical Society, Providence, R.I., 1981.
- 712 Arne Ludwig, Melanie Pruvost, Monika Reissmann, Norbert Benecke, Gudrun A Brock-  
713 mann, Pedro Castaños, Michael Cieslak, Sebastian Lippold, Laura Llorente, Anna-Sapfo  
714 Malaspinas, et al. Coat color variation at the beginning of horse domestication. *Science*,  
715 324(5926):485–485, 2009.
- 716 Anna-Sapfo Malaspinas, Orestis Malaspinas, Steven N Evans, and Montgomery Slatkin.  
717 Estimating allele age and selection coefficient from time-serial data. *Genetics*, 192(2):  
718 599–607, 2012.
- 719 Iain Mathieson and Gil McVean. Estimating selection coefficients in spatially structured  
720 populations from time series data of allele frequencies. *Genetics*, 193(3):973–984, 2013.
- 721 Iain Mathieson, Iosif Lazaridis, Nadin Rohland, Swapan Mallick, Nick Patterson,  
722 Songül Alpaslan Roodenberg, Eadaoin Harney, Kristin Stewardson, Daniel Fernandes,  
723 Mario Novak, et al. Genome-wide patterns of selection in 230 ancient eurasians. *Nature*,  
724 528(7583):499–503, 2015.
- 725 Rasmus Nielsen, Scott Williamson, Yuseob Kim, Melissa J Hubisz, Andrew G Clark, and  
726 Carlos Bustamante. Genomic scans for selective sweeps using snp data. *Genome research*,  
727 15(11):1566–1575, 2005.
- 728 Joseph K Pickrell, Graham Coop, John Novembre, Sridhar Kudaravalli, Jun Z Li, Devin  
729 Absher, Balaji S Srinivasan, Gregory S Barsh, Richard M Myers, Marcus W Feldman,  
730 et al. Signals of recent positive selection in a worldwide sample of human populations.  
731 *Genome research*, 19(5):826–837, 2009.
- 732 Martyn Plummer, Nicky Best, Kate Cowles, and Karen Vines. Coda: Conver-  
733 gence diagnosis and output analysis for mcmc. *R News*, 6(1):7–11, 2006. URL  
734 <http://CRAN.R-project.org/doc/Rnews/>.
- 735 Daniel Revuz and Marc Yor. *Continuous martingales and Brownian motion*, volume 293  
736 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathe-*  
737 *matical Sciences]*. Springer-Verlag, Berlin, third edition, 1999. ISBN 3-540-64325-7. doi:  
738 10.1007/978-3-662-06400-9. URL <http://dx.doi.org/10.1007/978-3-662-06400-9>.
- 739 Gareth O Roberts and Osnat Stramer. On inference for partially observed nonlinear diffu-  
740 sion models using the metropolis–hastings algorithm. *Biometrika*, 88(3):603–621, 2001.
- 741 Joshua G Schraiber. A path integral formulation of the wright–fisher process with genic  
742 selection. *Theoretical population biology*, 92:30–35, 2014.
- 743 Joshua G. Schraiber, Robert C. Griffiths, and Steven N. Evans. Analysis and rejection  
744 sampling of Wright-Fisher diffusion bridges. *Theoretical Population Biology*, 89(0):64–  
745 74, 2013.
- 746 Giorgos Sermaidis, Omiros Papaspiliopoulos, Gareth O Roberts, Alexandros Beskos, and  
747 Paul Fearnhead. Markov chain monte carlo for exact inference for diffusions. *Scandina-*  
748 *vian Journal of Statistics*, 2013.

- 749 Per Sjödin, Pontus Skoglund, and Mattias Jakobsson. Assessing the maximum contribution  
750 from ancient populations. *Molecular biology and evolution*, page msu059, 2014.
- 751 Montgomery Slatkin. Simulating genealogies of selected alleles in a population of variable  
752 size. *Genetical research*, 78(01):49–57, 2001.
- 753 Montgomery Slatkin and Richard R Hudson. Pairwise comparisons of mitochondrial dna  
754 sequences in stable and exponentially growing populations. *Genetics*, 129(2):555–562,  
755 1991.
- 756 Yun S Song and Matthias Steinrücken. A simple method for finding explicit analytic  
757 transition densities of diffusion processes with general diploid selection. *Genetics*, 190  
758 (3):1117–1129, 2012.
- 759 Michael Sørensen. Parametric inference for discretely sampled stochastic differential equa-  
760 tions. In *Handbook of financial time series*, pages 531–553. Springer, 2009.
- 761 Matthias Steinrücken, Anand Bhaskar, and Yun S Song. A novel spectral method for  
762 inferring general diploid selection from time series genetic data. *The annals of applied*  
763 *statistics*, 8(4):2203, 2014.
- 764 Benjamin F Voight, Sridhar Kudaravalli, Xiaoquan Wen, and Jonathan K Pritchard. A  
765 map of recent positive selection in the human genome. *PLoS biology*, 4(3):e72, 2006.
- 766 GA Watterson. Estimating and testing selection: the two-alleles, genic selection diffusion  
767 model. *Advances in Applied Probability*, pages 14–30, 1979.
- 768 Ellen G Williamson and Montgomery Slatkin. Using maximum likelihood to estimate  
769 population size from temporal changes in allele frequencies. *Genetics*, 152(2):755–761,  
770 1999.