

MOLECULAR ECOLOGY

RESOURCES

Inferring the timing and strength of natural selection and gene migration in the evolution of chicken from ancient DNA data

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID	MER-21-0253
Manuscript Type:	Resource Article
Date Submitted by the Author:	22-May-2021
Complete List of Authors:	Lyu, Wenyang; University of Bristol, School of Mathematics Dai, Xiaoyang; University of Bristol, School of Biological Sciences Beaumont, Mark; University of Bristol, School of Biological Sciences Yu, Feng; University of Bristol, School of Mathematics He, Zhangyi; University of Cambridge, MRC Toxicology Unit
Keywords:	Natural selection, Gene migration, Continent-island model, Wright-Fisher diffusion, Hidden Markov model, Blockwise particle marginal Metropolis-Hastings

SCHOLARONE™
Manuscripts

Major comments

1) Impact of dominance not explored in simulation study

Inferring the timing and strength of natural selection and gene migration in the evolution of chicken from ancient DNA data

Wenyang Lyu^a, Xiaoyang Dai^{b,1}, Mark Beaumont^b, Feng Yu^{a,*}, Zhangyi He^{c,2,*}

^a*School of Mathematics, University of Bristol, Bristol BS8 1UG, United Kingdom*

^b*School of Biological Sciences, University of Bristol, Bristol BS8 1TQ, United Kingdom*

^c*MRC Toxicology Unit, University of Cambridge, Cambridge CB2 1QR, United Kingdom*

Abstract

With the rapid growth of the number of sequenced ancient genomes, there has been increasing interest in using this new information to study past and present adaptation. Such an additional temporal component has the promise of providing improved power for the estimation of natural selection. Over the last decade, statistical approaches for detection and quantification of natural selection from ancient DNA (aDNA) data have been developed. However, most of the existing methods do not allow us to estimate the timing of natural selection along with its strength, which is key to understanding the evolution and persistence of organismal diversity. Additionally, most methods ignore the fact that natural populations are almost always structured, which can result in overestimation of the effect of natural selection. To address these issues, we propose a novel Bayesian framework for the inference of natural selection and gene migration from aDNA data with Markov chain Monte Carlo techniques, co-estimating both timing and strength of natural selection and gene migration. Such an advance enables us to infer drivers of natural selection and gene migration by correlating genetic evolution with potential causes such as the changes in the ecological context in which an organism has evolved. The performance of our procedure is evaluated through extensive simulations, with its utility shown with an application to ancient chicken samples.

Keywords: Natural selection, Gene migration, Continent-island model, Wright-Fisher diffusion, Hidden Markov model, Blockwise particle marginal Metropolis-Hastings

*Corresponding author.

Email addresses: feng.yu@bristol.ac.uk (Feng Yu), z.he@imperial.ac.uk (Zhangyi He)

¹Present address: The Blizard Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London E1 2AT, United Kingdom

²Present address: Department of Epidemiology and Biostatistics, School of Public Health, Imperial College London, London W2 1PG, United Kingdom

1 1. Introduction

With modern advances in ancient DNA (aDNA) techniques, there has been a rapid increase in the availability of time serial samples of segregating alleles across one or more related populations. The temporal aspect of such samples reflects the combined evolutionary forces acting within and among populations such as genetic drift, natural selection and gene migration, which can contribute to our understanding of how these evolutionary forces are responsible for the patterns observed in contemporaneous samples. One of the most powerful applications of such genetic time series is to study the action of natural selection since the expected changes in allele frequencies over time are closely related to the timing and strength of natural selection.

Over the past fifteen years, there has been a growing literature on the statistical inference of natural selection from time series data of allele frequencies, especially in aDNA (see Malaspinas, 2016; Dehasque et al., 2020, for excellent reviews). Typically, estimating natural selection from genetic time series is built on the hidden Markov model (HMM) framework proposed by Bollback et al. (2008), where the allele frequency trajectory of the underlying population through time was modelled as a latent variable following the Wright-Fisher model introduced by Fisher (1922) and Wright (1931), and the allele frequency of the sample drawn from the underlying population at each sampling time point was treated as a noisy observation of the underlying population allele frequency. In their likelihood computation, the Wright-Fisher model was approximated through its standard diffusion limit, called the Wright-Fisher diffusion, which was then discretised for numerical integration with a finite difference scheme. Their method was used to analyse aDNA data associated with horse coat colouration in Ludwig et al. (2009) and extended to more complex evolutionary scenarios (see e.g., Malaspinas et al., 2012; Steinrücken et al., 2014; Ferrer-Admetlla et al., 2016; Schraiber et al., 2016; He et al., 2020b,c).

Natural populations are almost always structured, which affects the relative effect of natural selection and genetic drift on the changes in allele frequencies over time. This can cause overestimation of the selection coefficient (Mathieson et al., 2015). However, all existing methods based on the Wright-Fisher model for the inference of natural selection from time series data of allele frequencies lack the ability to account for the confounding effect of gene migration, with the exception of Mathieson & McVean (2013), which could model population structure. Mathieson & McVean (2013) is an extension of Bollback et al. (2008) for the inference of metapopulations,

clarify what exactly
is the undesirable
feature

31 which enables joint estimation of the selection coefficient and the migration rate from genetic
32 time series. However, their method could become computationally cumbersome for large popu-
33 lation sizes and evolutionary timescales since their likelihood computation was carried out with
34 the Wright-Fisher model, which is an undesirable feature in aDNA.

35 More recently, Loog et al. (2017) developed a Bayesian statistical framework for estimating
36 the timing and strength of natural selection from genetic time series while explicitly modelling
37 gene migration from external sources. Their approach also allowed joint estimation of the allele
38 frequency trajectory of the underlying population through time, which was important for un-
39 derstanding the drivers of natural selection. However, the population size in their approach was
40 assumed to be large enough to ignore genetic drift, which simplifies their likelihood computation
41 but limits the application of their method to aDNA data.

42 In this work, we develop a novel HMM-based approach for the Bayesian inference of natural
43 selection and gene migration to re-analyse the temporally-spaced ancient chicken samples from
44 Loog et al. (2017). Our approach is built upon the HMM framework of Bollback et al. (2008),
45 but unlike most existing methods, it enables the joint estimation of the timing and strength of
46 natural selection and gene migration. Such an advance allows us to infer the drivers of natural
47 selection and gene migration by correlating genetic evolution with ecological and cultural shifts.
48 Our key innovation is to propose a multi-allele Wright-Fisher diffusion for a single locus evolving
49 under natural selection and gene migration, including the timing of natural selection and gene
50 migration. This diffusion process characterises the allele frequency trajectories of the underlying
51 population through time, where the alleles that migrate from external sources are no longer
52 treated as the same as those that originate in the underlying population. Such a setup allows us
53 to take full advantage of known quantities like the proportion of the modern European chicken
54 that have Asian origin as a direct result of gene migration. Our Bayesian inference of natural
55 selection and gene migration is carried out through the particle marginal Metropolis-Hastings
56 (PMMH) algorithm of Andrieu et al. (2010) with blockwise sampling, which permits a joint
57 update of the underlying population allele frequency trajectories. We evaluate the performance
58 of our procedure through extensive simulations, with its utility shown with an application to
59 the time serial samples of segregating alleles from ancient chicken reported in Loog et al. (2017).

N finite and fixed over time

60 2. Materials and Methods

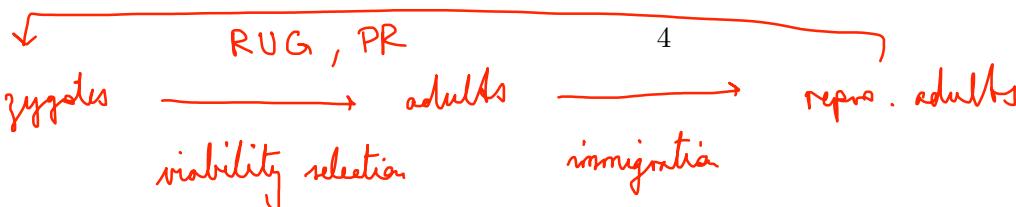
61 In this section, we first introduce the multi-allele Wright-Fisher diffusion for a single locus
 62 evolving under natural selection and gene migration and then present our Bayesian method for
 63 the joint inference of natural selection and gene migration from time series allele frequency data.

64 2.1. Wright-Fisher diffusion

65 Let us consider a population of N randomly mating diploid individuals at a single locus \mathcal{A}
 66 with discrete and nonoverlapping generations, where the population size N is finite and fixed
 67 over time. Suppose that at locus \mathcal{A} there are two allele types, labelled \mathcal{A}_1 and \mathcal{A}_2 , respectively.
 68 We attach the symbol \mathcal{A}_1 to the mutant allele, which arises only once in the population and
 69 is positively selected once the evolution starts to act through natural selection. We attach the
 70 symbol \mathcal{A}_2 to the ancestral allele, which originally exists in the population.

71 According to Loog et al. (2017), we characterise the population structure with the continent-
 72 island model (see, e.g., Hamilton, 2011, for an introduction). More specifically, the population is
 73 subdivided into two demes, the continental population and the island population. To distinguish
 74 between the alleles found on the island but emigrated from the continent or were originally on
 75 the island, the mutant and ancestral alleles that originated on the island are labelled \mathcal{A}_1^i and
 76 \mathcal{A}_2^i , respectively, and the mutant and ancestral alleles that were results of emigration from the
 77 continent are labelled \mathcal{A}_1^c and \mathcal{A}_2^c , respectively. We assume that the continent population is large
 78 enough such that the gene migration between the continent and the island does not influence the
 79 genetic composition of the continent population. Our interests in this work primarily focus on
 80 the island population dynamics. In what follows, the population refers to the island population
 81 unless noted otherwise.

82 To investigate the island population dynamics under natural selection and gene migration,
 83 we need to specify the life cycle of the island population, which starts with zygotes that natural
 84 selection acts on. Suppose that natural selection takes the form of viability selection, and the
 85 relative viabilities of all possible genotypes are shown in Table 1, where $s \in [0, 1]$ is the selection
 86 coefficient, and $h \in [0, 1]$ is the dominance parameter. After natural selection, a fraction m of
 87 the adults on the continent migrate into the population of mating adults on the island, which
 88 results in the change of the genetic composition of the island population, i.e., fraction m of
 89 the adults on the island are immigrants from the continent, and fraction $1 - m$ of adults were



Why distinguish between A_1^i and A_1^c and between A_2^i and A_2^c ?

originally already on the island. The Wright-Fisher reproduction introduced by Fisher (1922) and Wright (1931) finally completes the life cycle, which corresponds to randomly sampling $2N$ gametes with replacement from an effectively infinite gamete pool to form new zygotes in the next generation through random union of gametes.

We let $\mathbf{X}^{(N)}(k) = (X_1^{(N)}(k), X_2^{(N)}(k), X_3^{(N)}(k), X_4^{(N)}(k))$ be frequencies of the A_1^i , A_2^i , A_1^c and A_2^c alleles in N zygotes of generation $k \in \mathbb{N}$ on the island, which follows the multi-allele Wright-Fisher model with selection and migration described in Supplemental Material, File S1.

We assume that the selection coefficient and the migration rate are both of order $1/(2N)$ and fixed from the time of the onset of natural selection and gene migration up to present. We run time at rate $2N$, i.e., $t = k/(2N)$, and scale the selection coefficient and the migration rate as

$$\alpha(t) = \begin{cases} 0, & \text{if } t < t_s \\ 2Ns, & \text{otherwise} \end{cases} \quad \text{and} \quad \beta(t) = \begin{cases} 0, & \text{if } t < t_m \\ 4Nm, & \text{otherwise,} \end{cases}$$

*Multivariate difference factor
2N vs. 4N*

where t_s and t_m denote the starting times of natural selection and gene migration on the island measured in the unites of $2N$ generations. With the population size N approaching infinity, the Wright-Fisher model $\mathbf{X}^{(N)}$ converges to a diffusion process, denoted by $\mathbf{X} = \{\mathbf{X}(t), t \geq t_0\}$, evolving in the state space (i.e., a three-simplex)

$$\Omega_{\mathbf{X}} = \left\{ \mathbf{x} \in [0, 1]^4 : \sum_{i=1}^4 x_i = 1 \right\}$$

and satisfying the stochastic differential equation (SDE) of the form

$$d\mathbf{X}(t) = \underbrace{\boldsymbol{\mu}(\mathbf{X}(t), t)dt}_{\text{drift term}} + \underbrace{\boldsymbol{\sigma}(\mathbf{X}(t), t)d\mathbf{W}(t)}_{\text{diffusion term}}, \quad t \geq t_0 \quad (1)$$

with initial condition $\mathbf{X}(t_0) = \mathbf{x}_0$. In Eq. (1), $\boldsymbol{\mu}(\mathbf{x}, t)$ is the drift term

$$\begin{aligned} \mu_1(\mathbf{x}, t) &= \alpha(t)x_1(x_2 + x_4)[(x_1 + x_3)h + (x_2 + x_4)(1 - h)] - \frac{1}{2}\beta(t)x_1 \\ \mu_2(\mathbf{x}, t) &= -\alpha(t)x_2(x_1 + x_3)[(x_1 + x_3)h + (x_2 + x_4)(1 - h)] - \frac{1}{2}\beta(t)x_2 \\ \mu_3(\mathbf{x}, t) &= \alpha(t)x_3(x_2 + x_4)[(x_1 + x_3)h + (x_2 + x_4)(1 - h)] - \frac{1}{2}\beta(t)(x_3 - x_c) \\ \mu_4(\mathbf{x}, t) &= -\alpha(t)x_4(x_1 + x_3)[(x_1 + x_3)h + (x_2 + x_4)(1 - h)] - \frac{1}{2}\beta(t)(x_4 + x_c - 1), \end{aligned} \quad (2)$$

why not incorporate this to $\beta(t)$?

$$\binom{4}{2} = \frac{4!}{(4-2)! 2!} = \frac{4 \cdot 3 \cdot 2}{2} = 6$$

106 where x_c is the frequency of the A_1^c allele in the continent population, which is fixed over time,

107 $\sigma(\mathbf{x}, t)$ is the diffusion term

$$\sigma(\mathbf{x}, t) = \begin{pmatrix} \sqrt{x_1 x_2} & \sqrt{x_1 x_3} & \sqrt{x_1 x_4} & 0 & 0 & 0 \\ -\sqrt{x_2 x_1} & 0 & 0 & \sqrt{x_2 x_3} & \sqrt{x_2 x_4} & 0 \\ 0 & -\sqrt{x_3 x_1} & 0 & -\sqrt{x_3 x_2} & 0 & \sqrt{x_3 x_4} \\ 0 & 0 & -\sqrt{x_4 x_1} & 0 & -\sqrt{x_4 x_2} & -\sqrt{x_4 x_3} \end{pmatrix}, \quad (3)$$

108 and $\mathbf{W}(t)$ is a six-dimensional standard Brownian motion. Notice that the explicit formula for
 109 the diffusion term $\sigma(\mathbf{x}, t)$ in Eq. (3) is obtained by following He et al. (2020a). The proof of
 110 the convergence follows in the similar manner to that employed for the neutral two-locus case
 111 in Durrett (2008, p. 323). We refer to the stochastic process $\mathbf{X} = \{\mathbf{X}(t), t \geq t_0\}$ that satisfies
 112 the SDE in Eq. (1) as the multi-allele Wright-Fisher diffusion with selection and migration.

113 2.2. Bayesian inference of natural selection and gene migration

114 Suppose that the available data are sampled from the underlying island population at time
 115 points $t_1 < t_2 < \dots < t_K$, which are measured in units of $2N$ generations to be consistent with
 116 the Wright-Fisher diffusion timescale. At the sampling time point t_k , there are c_k mutant alleles
 117 (i.e., the A_1^i and A_2^i alleles) and d_k continent alleles (i.e., the A_1^c and A_2^c alleles) observed in
 118 the sample of n_k chromosomes drawn from the underlying island population. Note that in real
 119 data, the continent allele count of the sample may not be available at each sampling time point,
 120 e.g., the proportion of European chicken that have Asian origin is only available in the modern
 121 sample (Loog et al., 2017). The population genetic quantities of interest are the scaled selection
 122 coefficient $\alpha = 2Ns$, the dominance parameter h , the selection time t_s , the scaled migration
 123 rate $\beta = 4Nm$, and the migration time t_m , as well as the allele frequency trajectories of the
 124 underlying island population. For simplicity, in the sequel we let $\vartheta_s = (\alpha, h, t_s)$ be the selection-
 125 related parameters and $\vartheta_m = (\beta, t_m)$ be the migration-related parameters, respectively.

126 2.2.1. Hidden Markov model

127 We apply an HMM framework similar to the one proposed in Bollback et al. (2008). We as-
 128 sume that the underlying population evolves according to the Wright-Fisher diffusion in Eq. (1)
 129 and the observations are modelled through independent sampling from the underlying popula-

sampling : $c_k = x_1 + x_3 , \quad 6 \quad \# \text{ mutant alleles}$
 $d_k = x_3 + x_4 , \quad \# \text{ continental alleles}$

tion at each given time point. Unlike Loog et al. (2017), we co-estimate the timing and strength of natural selection and gene migration, including the allele frequency trajectories of the underlying population. Our Wright-Fisher diffusion can trace the changes over time in the frequency of the allele in the island population that results from emigrants from the continent population, which allows us to make the most of other available information such as the proportion of the modern European chicken with Asian ancestry in the most recent sample reported in Loog et al. (2017). This provides valuable information regarding the timing and strength of gene migration.

Let $\mathbf{x}_{1:K} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)$ be the allele frequency trajectories of the underlying population at the sampling time points $t_{1:K}$. Under our HMM framework, the joint posterior probability distribution for the population genetic quantities of interest and the allele frequency trajectories of the underlying population is

$$p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m, \mathbf{x}_{1:K} | \mathbf{c}_{1:K}, \mathbf{d}_{1:K}) \propto p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) p(\mathbf{x}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \mathbf{x}_{1:K}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m), \quad (4)$$

where $p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the prior probability distribution for the population genetic quantities of interest and can be taken to be a uniform prior over the parameter space if their prior knowledge is poor, $p(\mathbf{x}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the probability distribution for the allele frequency trajectories of the underlying population at the sampling time points $t_{1:K}$, and $p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \mathbf{x}_{1:K}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the probability of the observations at the sampling time points $t_{1:K}$ conditional on the underlying population allele frequency trajectories.

With the Markov property of the Wright-Fisher diffusion, we have

$$p(\mathbf{x}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = p(\mathbf{x}_1 | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) \prod_{k=1}^{K-1} p(\mathbf{x}_{k+1} | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m), \quad (5)$$

where $p(\mathbf{x}_1 | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the prior probability distribution for the allele frequencies of the underlying population at the initial sampling time point, commonly taken to be non-informative (e.g., flat over the entire state space) if the prior knowledge is poor, and $p(\mathbf{x}_{k+1} | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the transition probability density function of the Wright-Fisher diffusion between two consecutive sampling time points for $k = 1, 2, \dots, K - 1$, which satisfies the Kolmogorov backward equation (or its adjoint) resulting from the Wright-Fisher diffusion in Eq. (1). Unless otherwise specified, in this work we take the prior $p(\mathbf{x}_1 | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ to be a uniform distribution over the state space

155 Ω_X , known as the flat Dirichlet distribution, if gene migration starts before the first sampling
 156 time point, i.e., $t_m \leq t_1$; otherwise, the prior $p(\mathbf{x}_1 | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is set to be a uniform distribution
 157 over the state space Ω_X restricted to the line satisfying the condition that $x_3 = x_4 = 0$, i.e.,
 158 there is no continent allele in the island population.

159 Given the allele frequency trajectories of the underlying population, the observations at each
 160 sampling time point are independent. Therefore, we have

$$p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \mathbf{x}_{1:K}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = \prod_{k=1}^K p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m), \quad (6)$$

161 where $p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ is the probability for the observations at the k -th sampling time
 162 point given its relevant allele frequencies of the underlying population for $k = 1, 2, \dots, K$. If the
 163 sample continent allele count d_k is available, we introduce $\mathbf{z}_k = (z_{1,k}, z_{2,k}, z_{3,k}, z_{4,k})$ to denote
 164 the counts of the \mathcal{A}_1^i , \mathcal{A}_2^i , \mathcal{A}_1^c and \mathcal{A}_2^c alleles in the sample at the k -th sampling time point, and
 165 the emission probability $p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ can be expressed as

$$p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = \sum_{\mathbf{z}_k \in \Omega_{\mathbf{Z}_k}} \frac{n_k!}{\prod_{i=1}^4 z_{i,k}!} \prod_{i=1}^4 x_{i,k}^{z_{i,k}} \mathbb{1}_{\{z_{1,k}+z_{3,k}=c_k, z_{3,k}+z_{4,k}=d_k\}}(\mathbf{z}_k),$$

166 where

$$\Omega_{\mathbf{Z}_k} = \left\{ \mathbf{z}_k \in \mathbb{N}^4 : \sum_{i=1}^4 z_{i,k} = n_k \right\}$$

167 and $\mathbb{1}_A$ is the indicator function that equals to 1 if condition A holds and 0 otherwise. Otherwise,
 168 the emission probability $p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ can be reduced to

$$p(c_k, d_k | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = \frac{n_k!}{c_k!(n_k - c_k)!} (x_{1,k} + x_{3,k})^{c_k} (x_{2,k} + x_{4,k})^{n_k - c_k}.$$

169 2.2.2. Particle marginal Metropolis-Hastings

170 The most challenging part in the computation of the posterior $p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m, \mathbf{x}_{1:K} | \mathbf{c}_{1:K}, \mathbf{d}_{1:K})$ is

171 obtaining the ~~the~~ transition probability density function $p(\mathbf{x}_{k+1} | \mathbf{x}_k; \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ for $k = 1, 2, \dots, K-1$.

In principle, it can be obtained by numerically solving the Kolmogorov backward equation (or its

adjoint) associated with the Wright-Fisher diffusion in Eq. (1) like Bollback et al. (2008) and He

et al. (2020c). However, it typically requires a fine discretisation of the

state space Ω_X to guarantee numerically stable computation of the solution. Moreover, how fine

long sentence → split

176 the discretisation needs to be strongly depends on the underlying population genetic quantities
 177 that we aim to estimate (He et al., 2020a). We thus resort to the PMMH algorithm developed
 178 by Andrieu et al. (2010) in this work, ~~that~~^{The PMMH algo.} only involves simulating the Wright-Fisher SDE in
 179 Eq. (1), which permits a joint update of the population genetic parameters of interests and the
 180 allele frequency trajectories of the underlying population. Full details of the PMMH algorithm
 181 can be found in Andrieu et al. (2010). Fearnhead & Künsch (2018) provided a detailed review
 182 of Monte Carlo methods for parameter estimation in the HMM based on the particle filter.

183 In our PMMH-based procedure, the marginal likelihood

$$p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) = \int_{\Omega_{\mathbf{X}}^K} p(\mathbf{x}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \mathbf{x}_{1:K}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) d\mathbf{x}_{1:K}$$

184 is estimated with the bootstrap particle filter introduced by Gordon et al. (1993), where the
 185 particles are generated by simulating the Wright-Fisher diffusion in Eq. (1) through the Euler-
 186 Maruyama scheme. The product of average weights of the set of particles at the sampling time
 187 points $\mathbf{t}_{1:K}$ yields an unbiased estimate of the marginal likelihood $p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \mathbf{x}_{1:K}, \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$,
 188 and the underlying population allele frequency trajectories $\mathbf{x}_{1:K}$ are sampled once from the final
 189 set of particles with their corresponding weights. Given that the strength of natural selection and
 190 gene migration can be strongly correlated with their timing, we resort to a blockwise updating
 191 scheme to avoid the small acceptance ratio of the PMMH with full dimensional updates. We first
 192 split the population genetic quantities of interest into two disjoint blocks, the selection-related
 193 parameters $\boldsymbol{\vartheta}_s$ and the migration-related parameters $\boldsymbol{\vartheta}_m$, respectively, and then we iteratively
 194 update one block at a time through the PMMH.

195 More specifically, we first generate a set of the initial candidates of the parameters $(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$
 196 from the prior $p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$. We then run a bootstrap particle filter with the proposed parameters
 197 $(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ to obtain an initial candidate of the underlying population allele frequency trajectories
 198 $\mathbf{x}_{1:K}$ and a bootstrap particle filter's estimate of the marginal likelihood $p(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$.
 199 Repeat the following steps until a sufficient number of the samples of the parameters $(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$
 200 and the underlying population allele frequency trajectories $\mathbf{x}_{1:K}$ have been obtained:

201 Step 1: Update the selection-related parameters $\boldsymbol{\vartheta}_s$.

202 Step 1a: Draw a sample of new candidates of the selection-related parameters $\boldsymbol{\vartheta}_s^*$ from the
 203 proposal $q_s(\cdot | \boldsymbol{\vartheta}_s)$.

204 Step 1b: Run a bootstrap particle filter with the parameters $(\boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m)$ to yield the underlying
 205 population allele frequency trajectories $\mathbf{x}_{1:K}^*$ and the marginal likelihood estimate
 206 $\hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m)$.

207 Step 1c: Accept the parameters $\boldsymbol{\vartheta}_s^*$ and the underlying population allele frequency trajec-
 208 tories $\mathbf{x}_{1:K}^*$ with probability equal to the Metropolis-Hastings ratio

$$A = \frac{p(\boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m) \hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m)}{p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m) \hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)} \frac{q_s(\boldsymbol{\vartheta}_s | \boldsymbol{\vartheta}_s^*)}{q_s(\boldsymbol{\vartheta}_s^* | \boldsymbol{\vartheta}_s)}.$$

209 If the new candidates are rejected, put the parameters $\boldsymbol{\vartheta}_s^* = \boldsymbol{\vartheta}_s$ and the underlying
 210 population allele frequency trajectories $\mathbf{x}_{1:K}^* = \mathbf{x}_{1:K}$.

211 Step 2: Update the migration-related parameters $\boldsymbol{\vartheta}_m$.

212 Step 2a: Draw a sample of new candidates of the migration-related parameters $\boldsymbol{\vartheta}_m^*$ from
 213 the proposal $q_m(\cdot | \boldsymbol{\vartheta}_m)$.

214 Step 2b: Run a bootstrap particle filter with the parameters $(\boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m^*)$ to yield the underlying
 215 population allele frequency trajectories $\mathbf{x}_{1:K}^*$ and the marginal likelihood estimate
 216 $\hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m^*)$.

217 Step 2c: Accept the parameters $\boldsymbol{\vartheta}_m^*$ and the underlying population allele frequency trajec-
 218 tories $\mathbf{x}_{1:K}^*$ with probability equal to the Metropolis-Hastings ratio

$$A = \frac{p(\boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m^*) \hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m^*)}{p(\boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m) \hat{p}(\mathbf{c}_{1:K}, \mathbf{d}_{1:K} | \boldsymbol{\vartheta}_s^*, \boldsymbol{\vartheta}_m)} \frac{q_m(\boldsymbol{\vartheta}_m | \boldsymbol{\vartheta}_m^*)}{q_m(\boldsymbol{\vartheta}_m^* | \boldsymbol{\vartheta}_m)},$$

219 If the new candidates are rejected, put the parameters $\boldsymbol{\vartheta}_m^* = \boldsymbol{\vartheta}_m$ and the under-
 220 lying population allele frequency trajectories $\mathbf{x}_{1:K}^* = \mathbf{x}_{1:K}$. corret?

221 In this work we use random walk proposals for both selection- and migration-related parameters
 222 in our blockwise PMMH algorithm unless otherwise specified.

223 Once enough samples of the parameters $(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ and the underlying population allele fre-
 224 quency trajectories $\mathbf{x}_{1:K}$ have been yielded, we can compute the posterior $p(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m | \mathbf{c}_{1:K}, \mathbf{d}_{1:K})$
 225 from the samples of the parameters $(\boldsymbol{\vartheta}_s, \boldsymbol{\vartheta}_m)$ using nonparametric density estimation techniques
 226 (see Izenman, 1991, for a detailed review) and achieve the maximum a posteriori probability
 227 (MAP) estimates for the population genetic quantities of interest. Our estimates for the under-
 228 lying population allele frequency trajectories $\mathbf{x}_{1:K}$ are the posterior mean of the stored samples

229 of the underlying population allele frequency trajectories. Our method can be readily extended
230 to the analysis of multiple (independent) loci. Given that the migration-related parameters are
231 shared by all loci, in each iteration of our procedure we only need to replicate Step 1 once to
232 update selection-related parameters specified for each locus and then update migration-related
233 parameters with shared by all loci in Step 2, where the likelihood is replaced by the product of
234 the likelihoods for each locus.

235 3. Results

236 In this section, we first evaluate the performance of our approach using simulated datasets
237 with various population genetic parameter values and then apply it to re-analyse the time series
238 allele frequency data from ancient chicken in Loog et al. (2017) genotyped at the locus encoding
239 for the thyroid-stimulating hormone receptor (*TSHR*), which is hypothesised to have undergone
240 strong and recent natural selection in domestic chicken.

241 3.1. Robustness and performance

242 To test our procedure, we run forward-in-time simulations of the multi-allele Wright-Fisher
243 model with selection and migration described in Supplemental Material, File S1 and examine
244 the bias and the root mean square error (RMSE) of our estimates obtained from these replicate
245 simulations. We vary the selection coefficient $s \in \{0.003, 0.006, 0.009\}$, and fix the selection time
246 $k_s = 180$ and the dominance parameter $h = 0.5$. We set the migration rate $m = 0.005$ and vary
247 the migration time $k_m \in \{90, 360\}$. The starting times of natural selection and gene migration,
248 k_s and k_m , respectively, are measured in generations. Additionally, we vary the population size
249 $N \in \{5000, 50000, 500000\}$. In principle, the conclusions we draw in this section hold for other
250 values of the population genetic parameters in similar ranges.

251 For each of the 18 possible combinations of the selection coefficient, the migration time and
252 the population size, we perform 300 replicated runs. For each run, we take the starting allele
253 frequencies of the underlying island population at generation 0 (*i.e.*, the first sampling time
254 point) to be $x_1 = (0.4, 0.6, 0, 0)$ and the mutant allele frequency of the underlying continent
255 population to be $x_c = 0.9$. These values are similar to the allele frequencies of ancient chicken
256 samples reported in Loog et al. (2017). We simulate a total of 500 generations under the multi-
257 allele Wright-Fisher model with selection and migration and generate a multinomial sample of

van

17
entire
section
should
be in
past
line

In practice, how can you distinguish between continental and island alleles? By sequence differences in the flanking regions?

258 100 chromosomes from the underlying island population every 50 generations from generation 0,
 259 11 sampling time points in total. Given that in real data only mutant allele counts and continent
 260 allele counts are available, in our simulation studies we generate the mutant allele count of the
 261 sample by summing the first and third components of the simulated sample allele counts, and
 262 the continent allele count by summing the third and fourth components at each sampling time
 263 point. Moreover, in real data the continent allele count of the sample may not be available at
 264 each sampling time point (e.g., Loog et al., 2017), thereby assuming that the continent allele
 265 counts of the sample are unavailable at first three and seven sampling time points, respectively,
 266 for each run in our simulation studies (as shown in simulated dataset A and B, respectively, in
 267 Figure 1) *To study the impact of missing continent allele counts, we assume*

268 In our procedure, we choose a uniform prior over the interval $[-1, 1]$ for the selection coefficient s and a uniform prior over the interval $[0, 1]$ for the migration rate m . We assume that the
 269 starting times of natural selection and gene migration k_s and k_m are uniformly distributed over
 270 the set of all possible time points, i.e., $k_s, k_m \in \{k \in \mathbb{Z} : k \leq 500\}$. We run 10000 iterations of
 271 the blockwise PMMH with 1000 particles, and in the Euler-Maruyama scheme each generation
 272 is divided into five subintervals. We discard the first half of the iterations as the burn-in period
 273 and then thin the remaining PMMH output by selecting every fifth value. See Figures 2 and 3
 274 for our posteriors of the timing and strength of natural selection and gene migration obtained
 275 from the simulated datasets shown in Figure 1, including our estimates for the mutant and con-
 276 tinent allele frequency trajectories of the underlying island population. Evidently our approach
 277 is capable of identifying the signature of natural selection and gene migration and accurately
 278 estimate their timing and strength in these two examples. Also, the mutant and continent allele
 279 frequency trajectories of the underlying island population are well matched with our estimates,
 280 i.e., the mutant and continent allele frequency trajectories of the underlying island population
 281 fluctuate slightly around our estimates and are completely covered by our 95% highest posterior
 282 density (HPD) intervals.

284 In Figure 4, we present the boxplots of our estimates for the scenario where continent allele
 285 counts are not available at the first three sampling time points. These boxplots illustrate the
 286 relative bias of (a) the selection coefficient estimates, (b) the selection time estimates, (c) the
 287 migration rate estimates and (d) the migration time estimates across 18 different combinations

288 of the selection coefficient, the migration time and the population size. The tips of the whiskers
289 represent the 2.5%-quantile and the 97.5%-quantile, and the boxes denote the first and third
290 quartiles with the median in the middle. We summarise the bias and the RMSE of the estimates
291 in Supplemental Material, Tables S1 and S2.

292 As shown in Figure 4, our estimates for the selection coefficient and time are approximately
293 median-unbiased across 18 different parameter combinations, but the migration rate and time
294 are both slightly overestimated (i.e., our estimates show a small positive bias). An increase in
295 the population size results in better overall performance of our estimates (i.e., smaller bias with
296 smaller variance). In particular, the average proportion of the replicates that the signature of
297 natural selection can be identified (i.e., the 95% HPD interval does not contain the value of 0)
298 increases from 17.17% to 59.33% and then to 80.67% as the population size increases. Such an
299 improvement in the performance of our estimation is to be expected since large population sizes
300 reduce the magnitude of the stochastic effect on the changes in allele frequencies due to genetic
301 drift, which degrades evidence of natural selection and gene migration. for which I didn't

302 Compared to the case of natural selection starting after gene migration (i.e., $k_m = 90$), our
303 estimates for the case of natural selection starting before gene migration (i.e., $k_m = 360$) reveal
304 smaller bias and variances in both the selection coefficient and time, with the average proportion
305 of the replicates that the signature of natural selection can be identified increasing by 15.96%.
306 This might be because if natural selection begins before gene migration, there is a period of time
307 that the underlying trajectory of allele frequencies is only under natural selection. In contrast,
308 our estimates perform better for the migration rate when gene migration begins before natural
309 selection, but the performance for the migration time deteriorates somewhat unexpectedly. This
310 might be due to our parameter setting where the starting time of gene migration is within the
311 period of availability of continent allele counts for $k_m = 360$, but not for $k_m = 90$. ? ↘

312 In addition, we see from Figure 4 that the bias and variance of our estimates for the selection
313 coefficient and time are largely reduced as the selection coefficient increases, especially in terms
314 of outliers. The average proportion of the replicates where the signature of natural selection can
315 be identified increases from 27.56% to 63.11% and then to 66.50% as the selection coefficient
316 increases, with 97.17% for the case of large population size ($N = 500000$) and selection coefficient
317 ($s = 0.009$). For weak natural selection, the underlying trajectory of allele frequencies is

318 extremely stochastic so that it is difficult to disentangle the effects of genetic drift and natural
319 selection (Schraiber et al., 2013). An increase in the strength of natural selection leads to more
320 pronounced changes through time in allele frequencies, making the signature of natural selection
321 more identifiable. In contrast, an increase in the selection coefficient demonstrates little effect
322 on our estimates of the migration rate and time.

↙ has

323 In Figure 5, we show the boxplots of our estimates for the case where continent allele counts
324 are not available at the first seven sampling time points, with their bias and RMSE summarised
325 in Supplemental Material, Tables S3 and S4. They reveal similar behaviour in estimation bias
326 and variance, although the estimates for the migration-related parameters illustrate significantly
327 larger bias and variances, probably caused by the increased length of time when continent allele
328 counts are unavailable. This, however, has little effect on our estimation of the selection-related
329 parameters, with similar average proportions of the replicates where the signature of natural
330 selection can be identified (52.39% vs. 52.17%).

331 In conclusion, our approach can produce reasonably accurate joint estimates of the timing
332 and strength of natural selection and gene migration from time series data of allele frequencies
333 across different parameter combinations. Our estimates for the selection coefficient and time are
334 approximately median-unbiased, with smaller variances as the population size or the selection
335 coefficient (or both) increases. Our estimates for the migration rate and time both show little
336 positive bias. Their performance improves with an increase in population size or the number of
337 the sampling time points when continent allele counts are available (or both).

338 3.2. Application to ancient chicken samples

339 We re-analyse aDNA data of 452 European chicken genotyped at the *TSHR* locus (position
340 43250347 on chromosome 5) from earlier studies of Flink et al. (2014) and Loog et al. (2017).
341 The time from which the data come ranges from approximately 2200 years ago to the present.
342 The data shown in Table 2 come from grouping the raw chicken samples by their sampling time
343 points. The raw sample information and genotyping results can be found in Loog et al. (2017).
344 The derived *TSHR* allele has been associated with reduced aggression to conspecifics and faster
345 onset of egg laying (Belyaev, 1979; Rubin et al., 2010; Karlsson et al., 2015, 2016), which was
346 hypothesised to have undergone strong and recent selection in domestic chicken (Rubin et al.,
347 2010; Karlsson et al., 2015) from the period of time when changes in Medieval dietary preferences

part
please

provide a generic
motivation earlier
or in the methods
for >2 alleles

348 and husbandry practices across northwestern Europe occurred (Loog et al., 2017).

349 To avoid overestimating the effect of natural selection on allele frequency changes, we need
350 to account for recent gene migration in domestic chicken from Asia to Europe in our analyses.

351 More specifically, in our approach the European chicken population is represented as the island
352 population while the Asian chicken population is represented as the continent population with
353 a mutant allele frequency of $x_c = 0.99$ fixed from the time of the onset of gene migration, which
354 is a conservative estimate chosen in Loog et al. (2017). Gene migration from Asia in domestic
355 chicken, beginning around 250 years ago and continuing until the present, was historically well
356 documented (Dana et al., 2011; Flink et al., 2014; Lyimo et al., 2015). Unlike Loog et al. (2017),
357 we co-estimate the migration rate along with the selection coefficient and time by incorporating
358 the estimate reported in Loog et al. (2017) that approximately 15% of modern European chicken
359 have Asia origin. This allows us to obtain the sample frequency of the allele in European chicken
360 at the most recent sampling time point that resulted from immigration from Asia. We take the
361 average length of a generation of chicken to be one year.

✓ power analyses above should
cover the case of $h = 1$

362 In our analyses, we adopt the dominance parameter $h = 1$ since the derived *TSHR* allele is
363 recessive, and set the population size $N = 180000$ (95% HPD 26000-460000) estimated by Loog
364 et al. (2017). We pick a uniform prior over the interval $[-1, 1]$ for the selection coefficient s and
365 a uniform prior over the set $\{-9000, -8999, \dots, 0\}$ for the selection time k_s , which covers the
366 time period of chicken domestication dated to about 8000 years ago (95% CI 7014–8768 years)
367 (Lawal et al., 2020). We choose a uniform prior over the interval $[0, 1]$ for the migration rate m
368 and take the migration time to be $k_m = -250$. All settings in the Euler-Maruyama scheme and
369 the blockwise PMMH algorithm are the same as we applied in Section 3.1. The posteriors of
370 the selection coefficient, the selection time and the migration rate are shown in Figure 6, as well
371 as estimates for the mutant and Asian allele frequency trajectories of the underlying European
372 population. The MAP estimates, as well as 95% HPD intervals, are summarised in Table 3.

373 From Table 3, we observe that our estimate of the selection coefficient for the mutant allele
374 is 0.005120 with 95% HPD interval [0.003591, 0.007064], strong evidence to support the derived
375 *TSHR* allele being selectively advantageous in the European chicken population. Such positive
376 selection results in an increase over time in the mutant allele frequency, starting from 975 AD
377 with 95% HPD interval [611, 1174] AD (see Figure 6e). The starting frequency of the mutant

minus sign
consuming
here

↓
remind
reader
of limit
being
measured
rel. to

O BC;

cf. Table 2

378 allele in 128 BC is 0.454200 with 95% HPD interval [0.349024, 0.562094], which is similar to that
379 estimated in a red junglefowl captive zoo population in Rubin et al. (2010). Our estimate of
380 the migration rate for the Asian allele is 0.000659 with 95% HPD interval [0.000483, 0.000861].
381 This gene migration, starting about 250 years ago, leads to 15.2848% of European chicken with
382 Asian ancestry in 1995 AD, with 95% HPD interval [0.116412, 0.191382] (see Figure 6f). Our
383 findings are consistent with those reported in Loog et al. (2017). This is further confirmed by the
384 results obtained with different values of the population size (*i.e.*, $N = 26000$ and $N = 460000$,
385 the lower and upper bounds of 95% HPD interval for the population size given in Loog et al.
386 (2017), respectively). These results are shown in Supplemental Material Figures S3 and S4 and
387 summarised in Table 3.

simulations to test the effect of sparse sampling

388 To further evaluate the performance of our approach when samples are sparsely distributed
389 with small uneven sizes, such as the European chicken samples at the *TSHR* locus we have
390 studied above, we generate 300 simulated datasets that mimic the *TSHR* data, *i.e.*, we use the
391 sample times and sizes as given in Table 2, the timing and strength of natural selection and gene
392 migration as given by MAP estimates found in Table 3, and population size $N = 180000$. From
393 Figure 7, we find that our simulation studies based on the *TSHR* data yield median-unbiased
394 estimates for the selection coefficient, the selection time and the migration rate, similar to the
395 performance of our procedure applied in the simulation studies in Section 3.1. Additionally, the
396 signature of natural selection can be identified in all these 300 replicates. This further shows
397 that our approach can achieve a good performance with time serial samples that are sparsely
398 distributed with small uneven sizes, which is highly desirable for aDNA data.

399 In summary, our approach works well on the ancient chicken samples, even though they are
400 sparsely distributed with small uneven sizes. Our estimates show strong evidence for the derived
401 *TSHR* allele being positively selected between the 7th and 12th centuries AD, which coincides
402 exactly with the time period of changes in dietary preferences and husbandry practices across
403 northwestern Europe. This again demonstrates possible links established by Loog et al. (2017)
404 between the selective advantage of the derived *TSHR* allele and a historically attested cultural
405 shift in food preference in Medieval Europe.

406 4. Discussion

407 In this work, we introduced a novel MCMC-based approach for the joint inference of the
408 timing and strength of natural selection and gene migration from aDNA data. To our knowledge,
409 Mathieson & McVean (2013) and Loog et al. (2017) described the only existing approaches that
410 can jointly infer natural selection and gene migration from time series data of allele frequencies.
411 However, the method of Mathieson & McVean (2013) is unable to estimate the time of the onset
412 of natural selection and gene migration. Loog et al. (2017) only showed the applicability of their
413 approach in the case where timing and strength of gene migration were both prespecified. In
414 addition, their method is restricted by the assumption of infinite population size, which highly
415 limits the application of their approach to aDNA data.

416 Our Bayesian inference procedure was built on an HMM framework incorporating a multi-
417 allele Wright-Fisher diffusion with selection and migration. Our estimates for the timing and
418 strength of natural selection and gene migration were achieved through the PMMH algorithm
419 with blockwise sampling, which enables joint estimation of the underlying trajectories of allele
420 frequencies as well. This is a highly desirable feature for aDNA because it allows us to infer the
421 drivers of natural selection and gene migration by correlating patterns of genetic variation with
422 potential evolutionary events such as changes in the ecological context in which an organism
423 has evolved.

424 We showed through extensive simulation studies that our approach could deliver reasonably
425 accurate estimates for the timing and strength of natural selection and gene migration, including
426 the estimates for the underlying trajectories of allele frequencies through time. In our simulation
427 studies, the estimates for the selection rate and time were largely unbiased, while the estimates
428 for the migration rate and time showed a slight positive bias. We applied our Bayesian inference
429 procedure to re-analyse ancient European chicken samples genotyped at the *TSHR* locus from
430 previous studies of Flink et al. (2014) and Loog et al. (2017). We observed that the derived
431 *TSHR* allele became selectively advantageous from 975 AD (95% HPD 611-1174 AD), which
432 was similar to that reported in Loog et al. (2017). Our results further confirmed the findings
433 of Loog et al. (2017) that positive selection acting on the *TSHR* locus in European chicken
434 could be driven by chicken intensification and egg production in Medieval Europe as a result
435 of Christian fasting practices (i.e., the consumption of birds, eggs and fish became allowed

!

discussed
well
enough?

436 (Venarde, 2011). Except for religiously inspired dietary preferences, this could also result from
437 changes in Medieval husbandry practices along with population growth and urbanisation in the
438 High Middle Ages (around 1000–250 AD). See Loog et al. (2017) and references cited therein
439 for more details.

440 Unlike Loog et al. (2017), our method takes the influence of genetic drift into account. From
441 Table 3, we see that our estimates from aDNA data for *TSHR* are close to each other regardless
442 of what population size we pick from the 95% HPD interval for European chicken population
443 size reported in Loog et al. (2017). This implies that ignoring genetic drift might have little
444 effect on the inference of natural selection from aDNA data such as those in Loog et al. (2017).
445 To explore the effect of genetic drift, we further simulate 300 datasets based on the aDNA data
446 for *TSHR*, where the timing and strength of natural selection and gene migration are set to our
447 estimates found in Table 3 but the true population size is picked to be $N = 4500$. We run our
448 method with an incorrect population size $N = 180000$ for these 300 replicates and find that this
449 incorrect (larger) population size results in significant overestimation of the selection coefficient
450 and time (see Figure 8). Moreover, the misspecified population size causes a reduction of 9.67%
451 in the proportion of the replicates that the signature of natural selection can be identified, which
452 implies the necessity of modelling genetic drift in the inference of natural selection from aDNA
453 data.

454 Furthermore, we explore the effect of gene migration in a similar manner. We simulate 300
455 datasets based on the aDNA data for *TSHR* with the migration rate $m = 0.00001$ and $m = 0.01$,
456 respectively, but run our procedure with a misspecified migration rate $m = 0.000659$ (i.e., the
457 migration rate estimated with the population size $N = 180000$ shown in Table 3). We can find
458 from Figure 9 that an incorrect migration rate does not dramatically alter the posterior median
459 of the selection-related parameters, but can change the overall shape of the posterior surface for
460 the selection time. This possibly results from incorrect information used for gene migration also
461 causing incorrect information for natural selection, therefore disturbing the resulting posterior
462 surface.

463 Although we have focused on the continent-island model in this work, our Bayesian statis-
464 tical framework lends itself to being extended to more complex models of gene migration, e.g.,
465 multiple islands. With an increase in the number of demes, our approach will be more compu-

466 tationally demanding, but improvements to exact-approximate particle filtering techniques like
467 the PMMH algorithm continue to be developed (See e.g. Yıldırım et al. 2018). Our approach
468 is also readily applicable to the analysis of multiple (independent) loci, where updating the
469 selection-related parameters for each locus can proceed in parallel on different cores. Moreover,
470 it is possible to extend our procedure to handle the case of non-constant demographic histories
471 like Schraiber et al. (2016) and He et al. (2020c). To achieve accurate estimation of relevant
472 population genetic quantities of interest, it is important to account for local linkage among loci,
473 which has been illustrated to be capable of further improving the inference of natural selection
474 (He et al., 2020b). Our Bayesian statistical framework can be readily extended to the scenario
475 of two linked loci by incorporating the method of He et al. (2020b) but such an extension will
476 probably be computationally prohibitive in the case of multiple linked loci. As a tractable al-
477 ternative for multiple linked loci, we can apply our two-locus method in a pairwise manner by
478 adding additional blocks in blockwise sampling.

479 References

- 480 Andrieu, C., Doucet, A., & Holenstein, R. (2010). Particle Markov chain Monte Carlo methods.
481 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72, 269–342.
- 482 Belyaev, D. K. (1979). Destabilizing selection as a factor in domestication. *Journal of Heredity*,
483 70, 301–308.
- 484 Bollback, J. P., York, T. L., & Nielsen, R. (2008). Estimation of $2N_e s$ from temporal allele
485 frequency data. *Genetics*, 179, 497–502.
- 486 Dana, N., Megens, H.-J., Crooijmans, R. P. M. A., Hanotte, O., Mwacharo, J. et al. (2011).
487 East Asian contributions to Dutch traditional and western commercial chickens inferred from
488 mtDNA analysis. *Animal Genetics*, 42, 125–133.
- 489 Dehasque, M., Ávila-Arcos, M. C., Díez-del Molino, D., Fumagalli, M., Guschanski, K. et al.
490 (2020). Inference of natural selection from ancient DNA. *Evolution Letters*, 4, 94–108.
- 491 Durrett, R. (2008). *Probability Models for DNA Sequence Evolution*. New York: Springer-
492 Verlag.

- 493 Fearnhead, P., & Künsch, H. R. (2018). Particle filters and data assimilation. *Annual Review
494 of Statistics and Its Application*, 5, 421–449.
- 495 Ferrer-Admetlla, A., Leuenberger, C., Jensen, J. D., & Wegmann, D. (2016). An approximate
496 Markov model for the Wright-Fisher diffusion and its application to time series data. *Genetics*,
497 203, 831–846.
- 498 Fisher, R. A. (1922). On the dominance ratio. *Proceedings of the Royal Society of Edinburgh*,
499 42, 321–341.
- 500 Flink, L. G., Allen, R., Barnett, R., Malmström, H., Peters, J. et al. (2014). Establishing the
501 validity of domestication genes using DNA from ancient chickens. *Proceedings of the National
502 Academy of Sciences*, 111, 6184–6189.
- 503 Gordon, N. J., Salmond, D. J., & Smith, A. F. M. (1993). Novel approach to nonlinear/non-
504 Gaussian Bayesian state estimation. *IEE Proceedings F (Radar and Signal Processing)*, 140,
505 107–113.
- 506 Hamilton, M. (2011). *Population Genetics*. Chichester: Wiley-Blackwell.
- 507 He, Z., Beaumont, M. A., & Yu, F. (2020a). Numerical simulation of the two-locus Wright-Fisher
508 stochastic differential equation with application to approximating transition probability den-
509 sities. *bioRxiv*, (p. 213769).
- 510 He, Z., Dai, X., Beaumont, M. A., & Yu, F. (2020b). Detecting and quantifying natural selection
511 at two linked loci from time series data of allele frequencies with forward-in-time simulations.
512 *Genetics*, 216, 521–541.
- 513 He, Z., Dai, X., Beaumont, M. A., & Yu, F. (2020c). Estimation of natural selection and allele
514 age from time series allele frequency data using a novel likelihood-based approach. *Genetics*,
515 216, 463–480.
- 516 Izenman, A. J. (1991). Recent developments in nonparametric density estimation. *Journal of
517 the American Statistical Association*, 86, 205–224.
- 518 Karlsson, A.-C., Fallahshahroudi, A., Johnsen, H., Hagenblad, J., Wright, D. et al. (2016). A
519 domestication related mutation in the thyroid stimulating hormone receptor gene (*TSHR*)

- 520 modulates photoperiodic response and reproduction in chickens. *General and Comparative*
521 *Endocrinology*, 228, 69–78.
- 522 Karlsson, A.-C., Sveyer, F., Eriksson, J., Darras, V. M., Andersson, L. et al. (2015). The
523 effect of a mutation in the thyroid stimulating hormone receptor (TSHR) on development,
524 behaviour and TH levels in domesticated chickens. *PLoS One*, 10, e0129040.
- 525 Lawal, R. A., Martin, S. H., Vanmechelen, K., Vereijken, A., Silva, P. et al. (2020). The wild
526 species genome ancestry of domestic chickens. *BMC Biology*, 18, 1–18.
- 527 Loog, L., Thomas, M. G., Barnett, R., Allen, R., Sykes, N. et al. (2017). Inferring allele
528 frequency trajectories from ancient DNA indicates that selection on a chicken gene coincided
529 with changes in medieval husbandry practices. *Molecular Biology and Evolution*, 34, 1981–
530 1990.
- 531 Ludwig, A., Pruvost, M., Reissmann, M., Benecke, N., Brockmann, G. A. et al. (2009). Coat
532 color variation at the beginning of horse domestication. *Science*, 324, 485–485.
- 533 Lyimo, C. M., Weigend, A., Mssoffe, P. L., Hocking, P. M., Simianer, H. et al. (2015). Maternal
534 genealogical patterns of chicken breeds sampled in Europe. *Animal Genetics*, 46, 447–451.
- 535 Malaspina, A.-S. (2016). Methods to characterize selective sweeps using time serial samples:
536 an ancient DNA perspective. *Molecular Ecology*, 25, 24–41.
- 537 Malaspina, A.-S., Malaspina, O., Evans, S. N., & Slatkin, M. (2012). Estimating allele age
538 and selection coefficient from time-serial data. *Genetics*, 192, 599–607.
- 539 Mathieson, I., Lazaridis, I., Rohland, N., Mallick, S., Patterson, N. et al. (2015). Genome-wide
540 patterns of selection in 230 ancient Eurasians. *Nature*, 528, 499–503.
- 541 Mathieson, I., & McVean, G. (2013). Estimating selection coefficients in spatially structured
542 populations from time series data of allele frequencies. *Genetics*, 193, 973–984.
- 543 Rubin, C.-J., Zody, M. C., Eriksson, J., Meadows, J. R. S., Sherwood, E. et al. (2010). Whole-
544 genome resequencing reveals loci under selection during chicken domestication. *Nature*, 464,
545 587–591.

- 546 Schraiber, J. G., Evans, S. N., & Slatkin, M. (2016). Bayesian inference of natural selection
547 from allele frequency time series. *Genetics*, 203, 493–511.
- 548 Schraiber, J. G., Griffiths, R. C., & Evans, S. N. (2013). Analysis and rejection sampling of
549 Wright-Fisher diffusion bridges. *Theoretical Population Biology*, 89, 64–74.
- 550 Steinrücken, M., Bhaskar, A., & Song, Y. S. (2014). A novel spectral method for inferring
551 general diploid selection from time series genetic data. *The Annals of Applied Statistics*, 8,
552 2203–2222.
- 553 Venarde, B. L. (2011). *The Rule of Saint Benedict*. Cambridge, Massachusetts: Harvard
554 University Press.
- 555 Wright, S. (1931). Evolution in Mendelian populations. *Genetics*, 16, 97–159.
- 556 Yıldırım, S., Andrieu, C., & Doucet, A. (2018). Scalable Monte Carlo inference for state-space
557 models. [arXiv:1809.02527](https://arxiv.org/abs/1809.02527).

558 **Data Accessibility Statement**

559 The authors state that all data necessary for confirming the conclusions of this work are
560 represented fully within the article. Source code implementing the method described in this work
561 is available at <https://github.com/zhangyi-he/WFM-1L-DiffusApprox-PMMH-Chicken>.

562 **Author Contributions**

563 F.Y. and Z.H. designed the project and developed the method; W.L. and Z.H. implemented
564 the method; W.L. and X.D. analysed the data under the supervision of M.B., F.Y. and Z.H.;
565 W.L., X.D. and Z.H. wrote the manuscript; M.B. and F.Y. reviewed the manuscript.

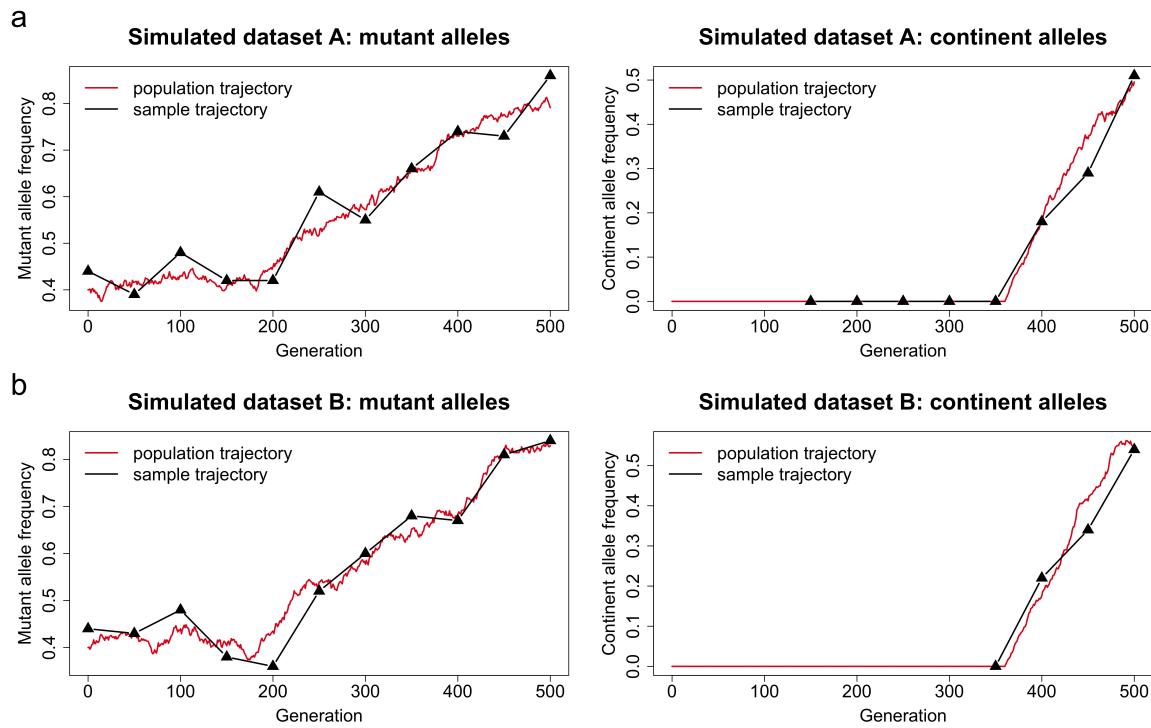


Figure 1: Representative examples of the datasets simulated using the Wright-Fisher model with selection and migration. We take the selection coefficient and time to be $s = 0.006$ and $k_s = 180$ and the migration rate and time to be $m = 0.005$ and $k_m = 360$, respectively. We set the dominance parameter $h = 0.5$ and the population size $N = 5000$. We adopt the starting allele frequencies of the underlying island population $\mathbf{x}_1 = (0.4, 0.6, 0, 0)$ and the mutant allele frequency of the underlying continent population $x_c = 0.9$. We sample 100 chromosomes at every 50 generations from generation 0 to 500. (a) simulated dataset A: continent allele counts are not available at the first three sampling time points. (b) simulated dataset B: continent allele counts of the sample are not available at the first seven sampling time points.

The authors should also include panels for $k_m = 90$, keeping the number of generations for which continent allele counts are missing the same as above.

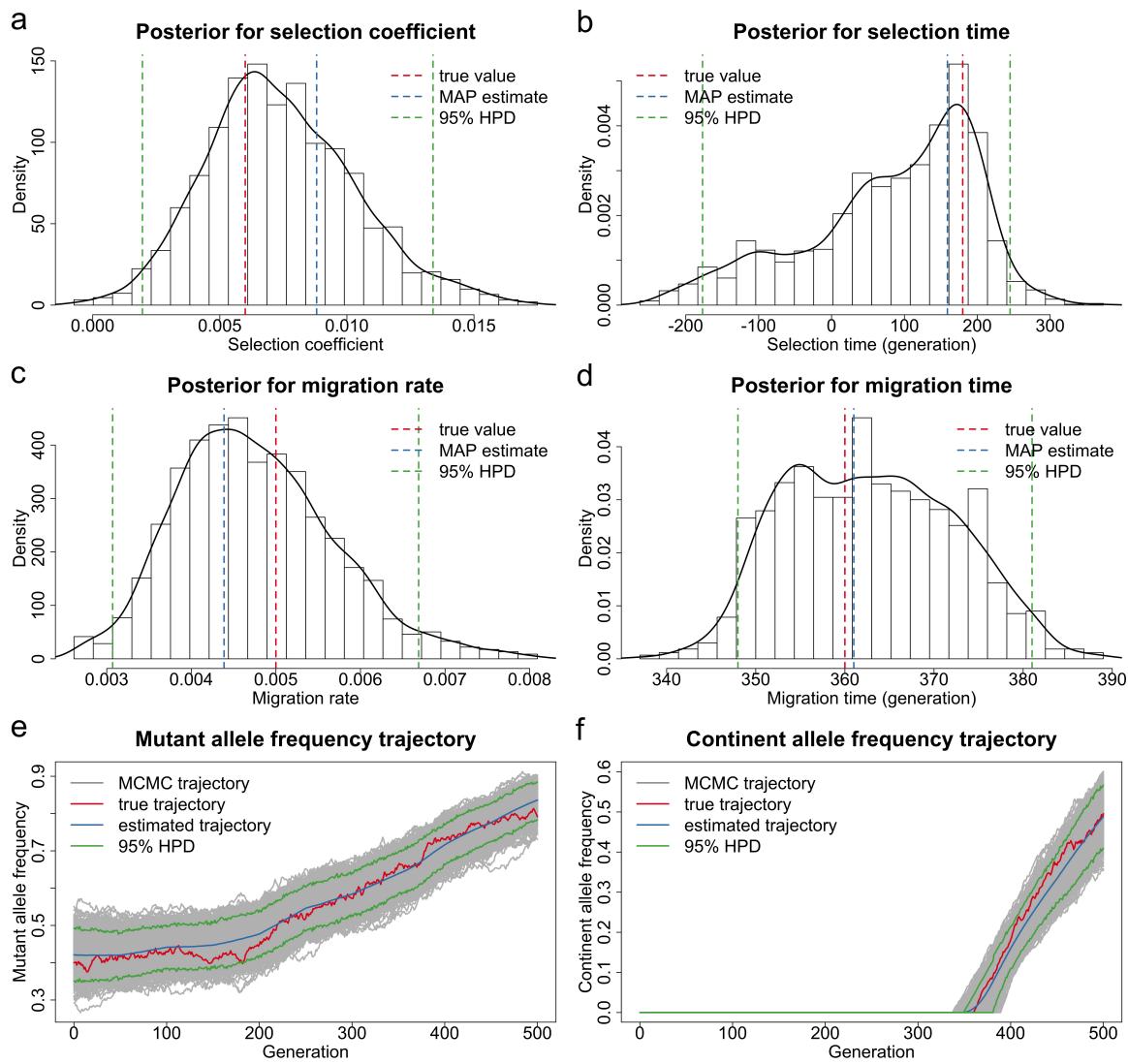


Figure 2: Bayesian estimates for the dataset shown in Figure 1a simulated for the case of the continent allele counts unavailable at the first three sampling time points. Posteriors for (a) the selection coefficient (b) the selection time (c) the migration rate and (d) the migration time. The MAP estimate is for the joint posterior, and may not correspond to the mode of the marginals. Estimated underlying trajectories of (e) the mutant allele frequency and (f) the continent allele frequency of the island population.

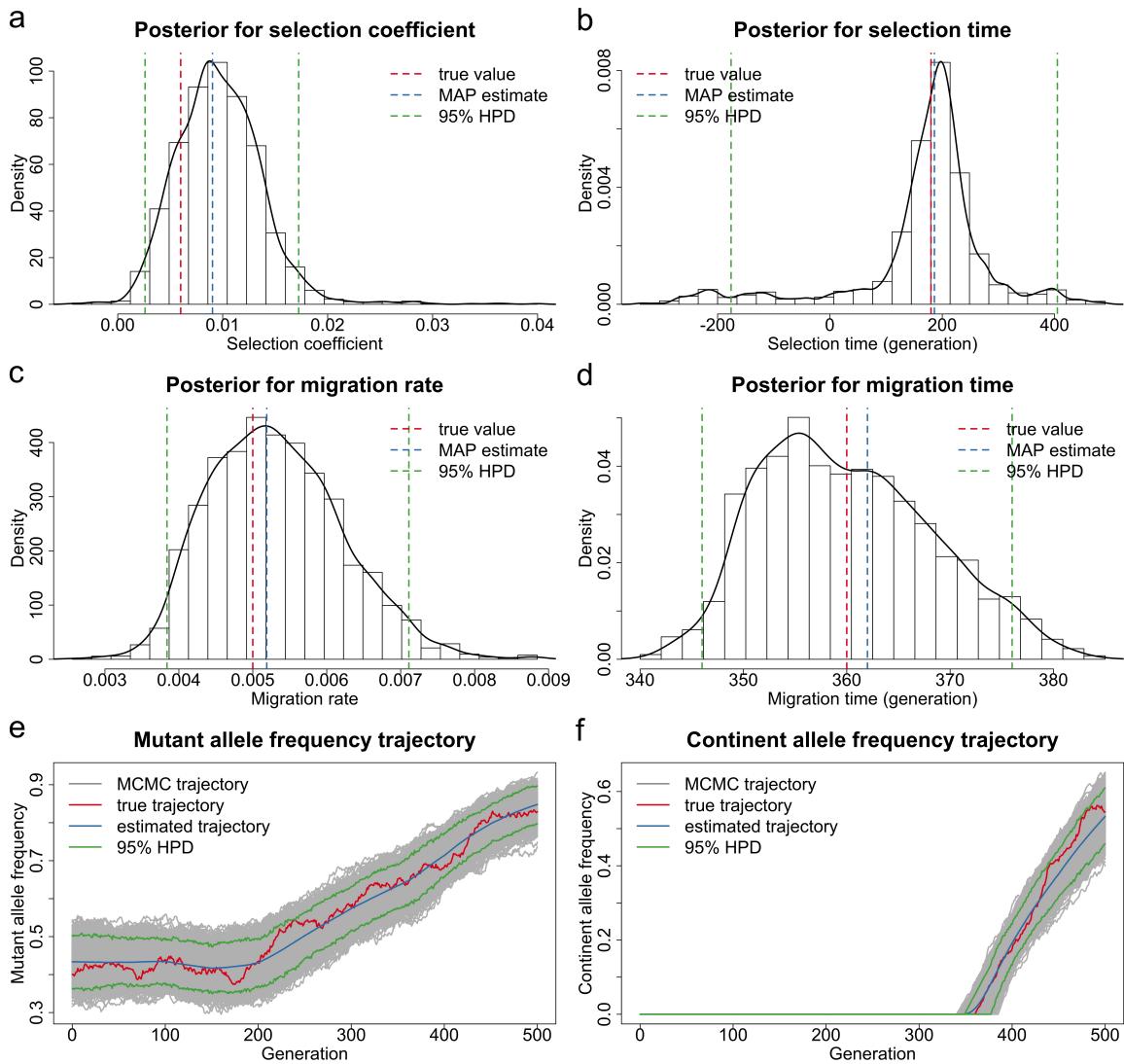


Figure 3: Bayesian estimates for the dataset shown in Figure 1b simulated for the case of continent allele counts unavailable at the first seven sampling time points. Posteriors for (a) the selection coefficient (b) the selection time (c) the migration rate and (d) the migration time. The MAP estimate is for the joint posterior, and may not correspond to the mode of the marginals. Estimated underlying trajectories of (e) the mutant allele frequency and (f) the continent allele frequency of the island population.

$k_s = 180$

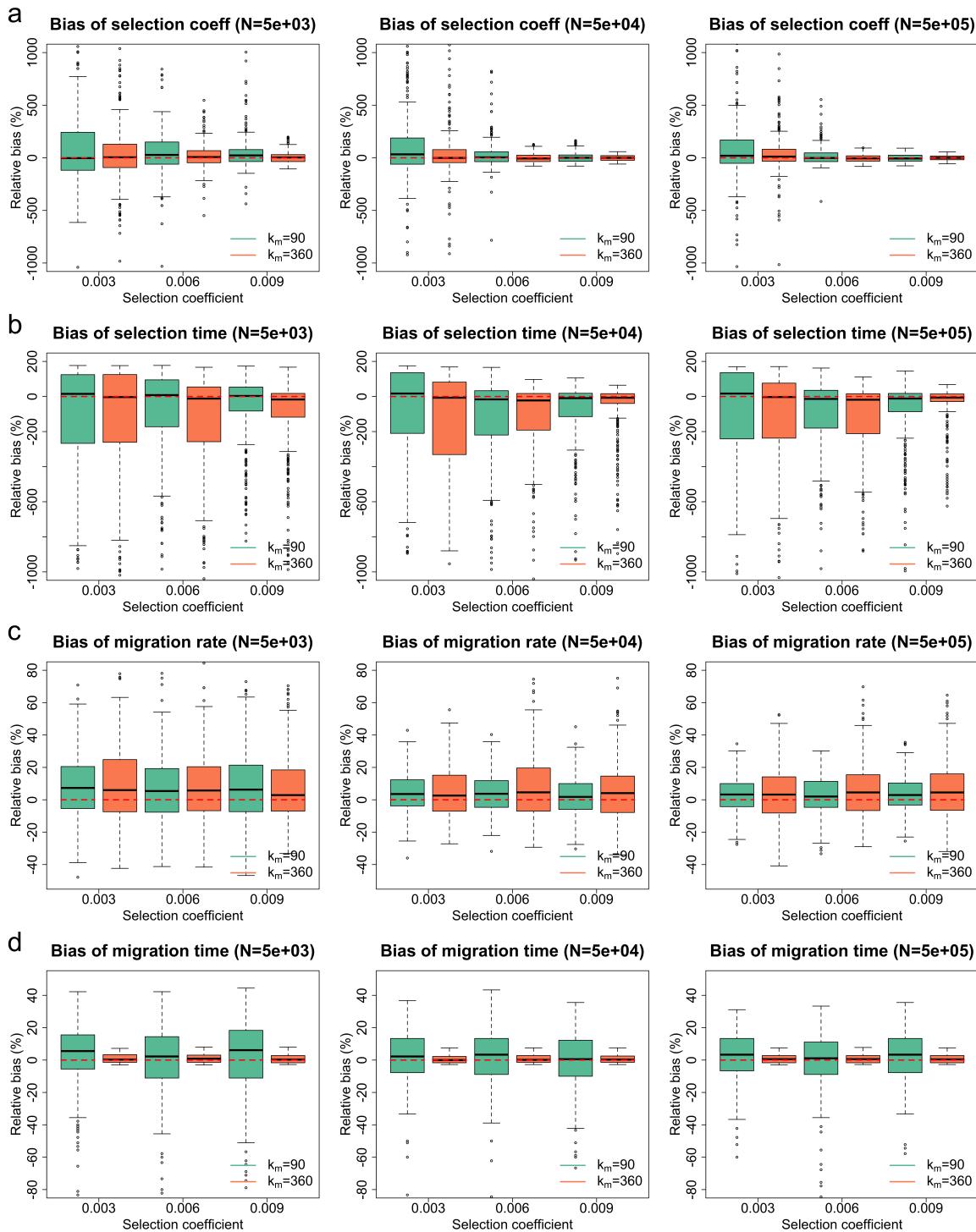


Figure 4: Empirical distributions of the estimates for 300 datasets simulated for the case of continent allele counts unavailable at the first three sampling time points. Aquamarine boxplots represent the estimates produced for the case of natural selection starting after gene migration, and coral boxplots represent the estimates produced for the case of natural selection starting before gene migration. Boxplots of the relative bias of (a) the selection coefficient estimates (b) the selection time estimates (c) the migration rate estimates and (d) the migration time estimates. To aid visual comparison, we have picked the y axes here so that boxes are of a relatively large size. This causes some outliers to lie outside the plots. Boxplots containing all outliers can be found in Supplemental Material, Figure S1.

just say 'green' and 'orange'?

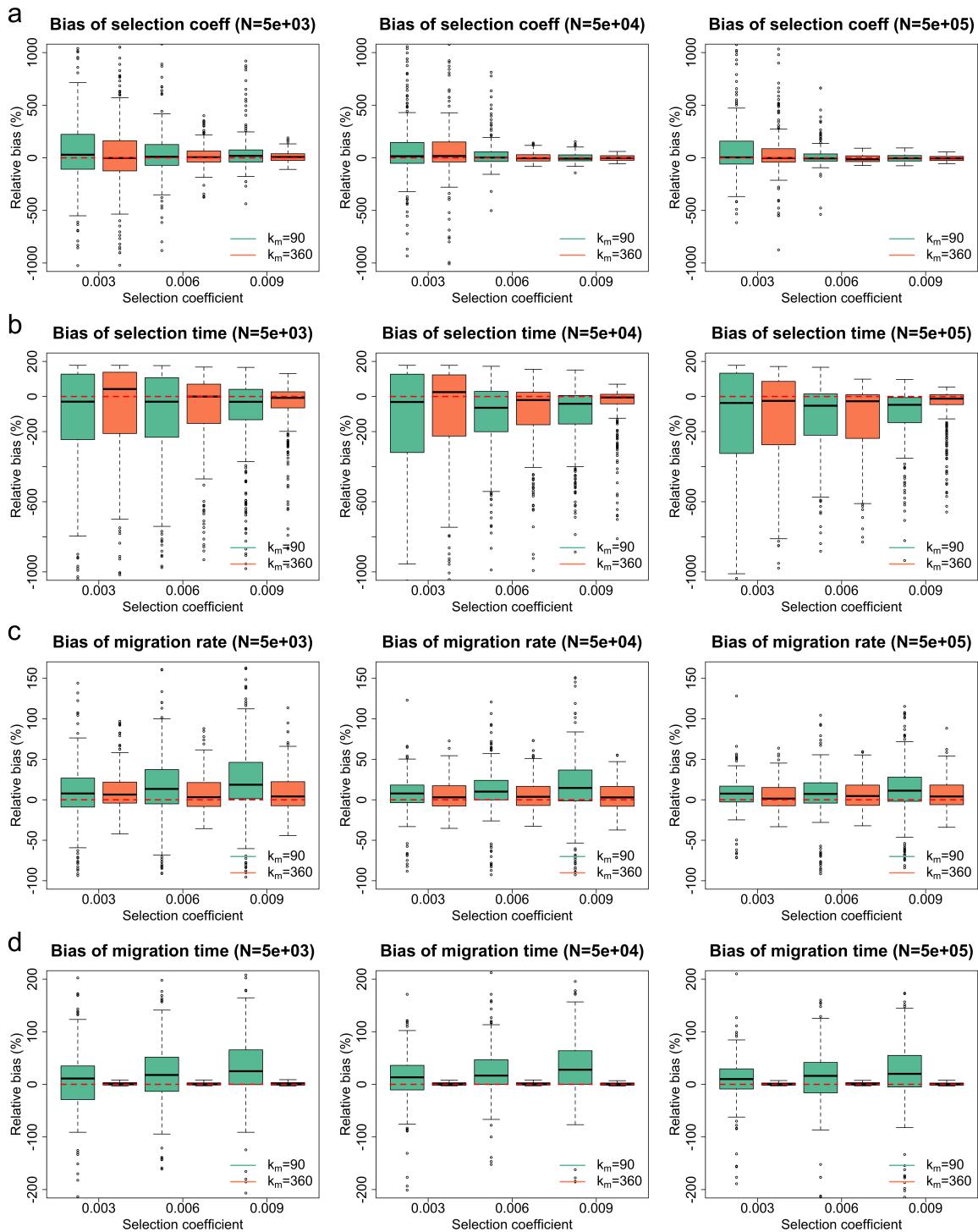


Figure 5: Empirical distributions of the estimates for 300 datasets simulated for the case of continent allele counts unavailable at the first seven sampling time points. Aquamarine boxplots represent the estimates produced for the case of natural selection starting after gene migration, and coral boxplots represent the estimates produced for the case of natural selection starting before gene migration. Boxplots of the relative bias of (a) the selection coefficient estimates (b) the selection time estimates (c) the migration rate estimates and (d) the migration time estimates. To aid visual comparison, we have picked the y axes here so that boxes are of a relatively large size. This causes some outliers to lie outside the plots. Boxplots containing all outliers can be found in Supplemental Material, Figure S2.

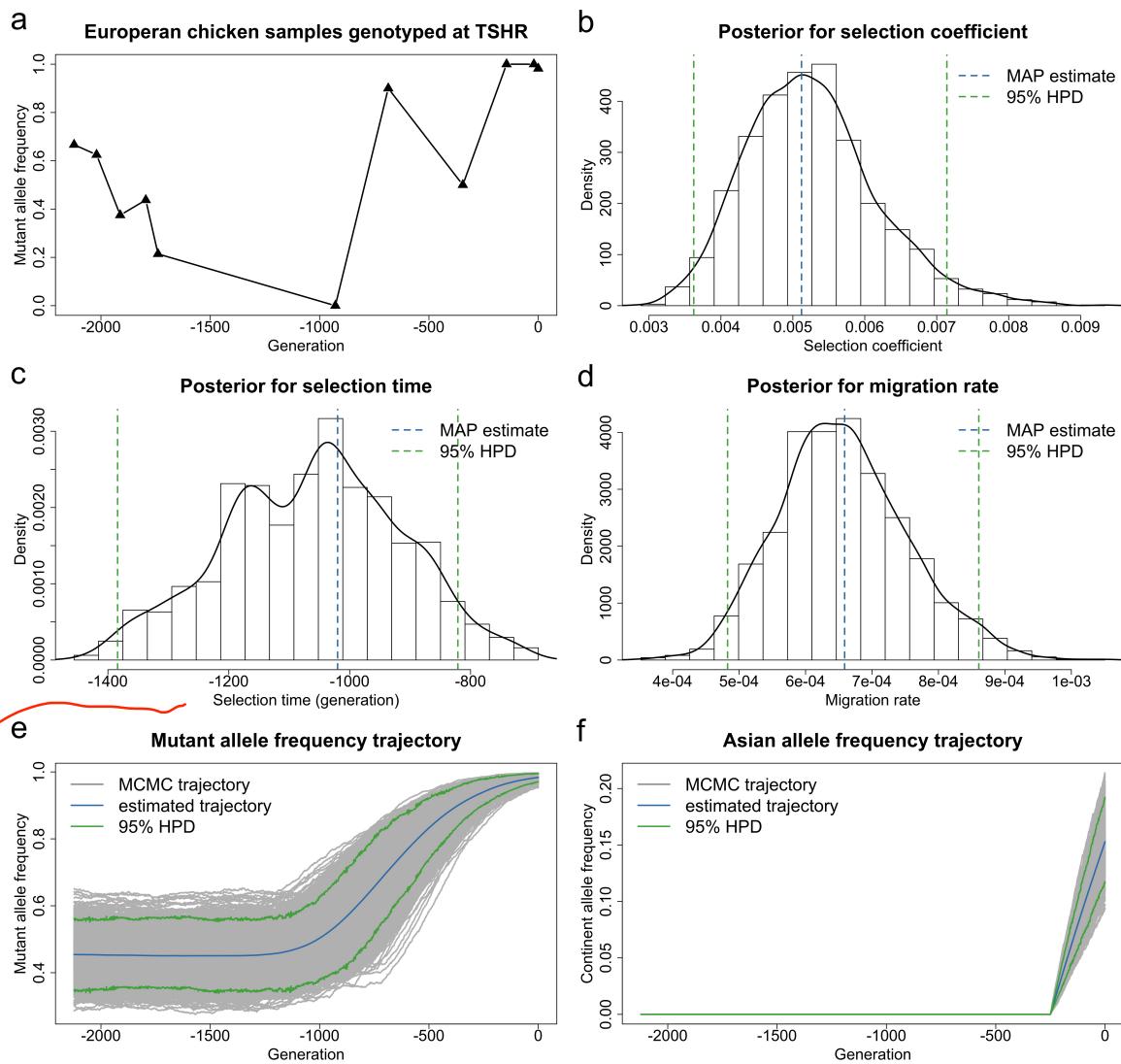


Figure 6: Bayesian estimates for aDNA data of European chicken genotyped at the *TSHR* locus from Loog et al. (2017) for the case of the population size $N = 180000$. (a) Temporal changes in the mutant allele frequencies of the sample, where the sampling time points have been offset so that the most recent sampling time point (1995 AD) is generation 0. Posteriors for (b) the selection coefficient (c) the selection time and (d) the migration rate. Estimated underlying trajectories of (e) the mutant allele frequency and (f) the Asian allele frequency in the European chicken population. The MAP estimate is for the joint posterior, and may not correspond to the mode of the marginals.

line scale continuing w.
to simulation study

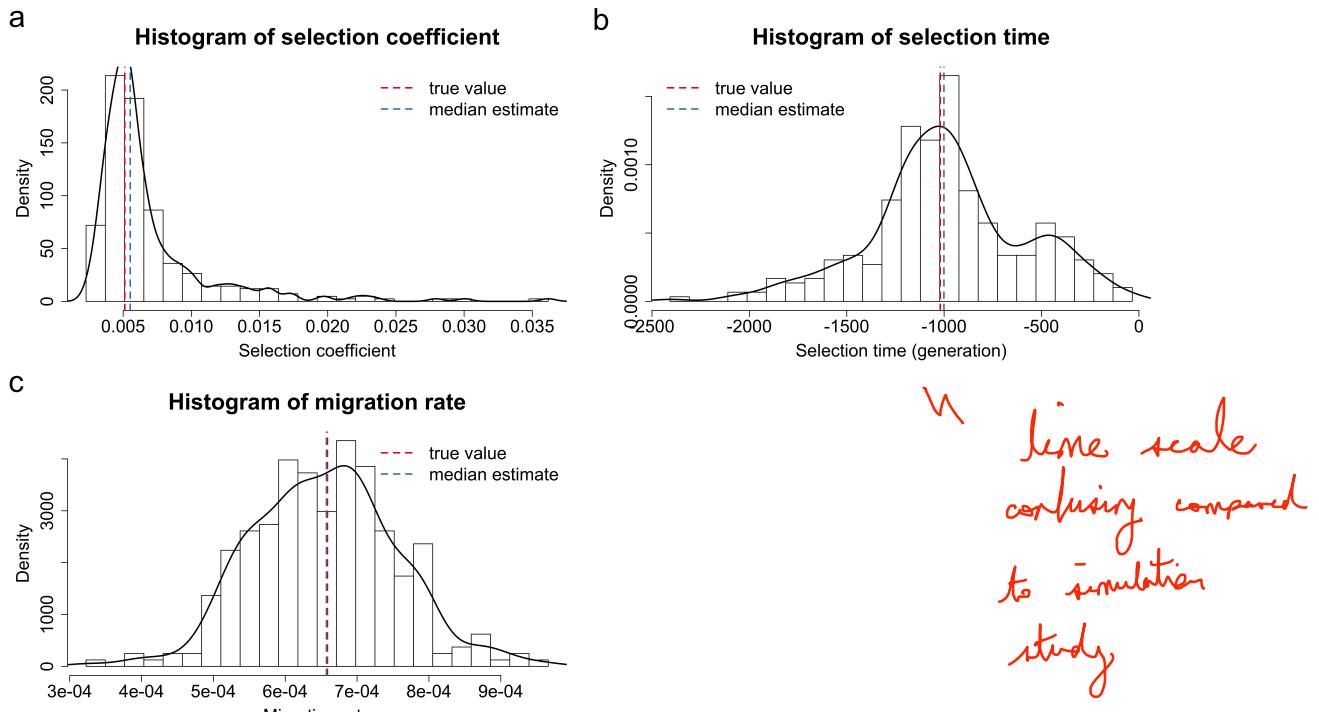


Figure 7: Empirical distributions of the estimates for 300 datasets simulated for *TSHR* based on the aDNA data shown in Table 2. We simulate the underlying population dynamics with the timing and strength of natural selection and gene migration estimated with the population size $N = 180000$ shown in Table 3. Histograms of (a) the selection coefficient estimates (b) the selection time estimates and (c) the migration rate estimates. To aid visual comparison, we have picked the x axis in (a) not to cover all 300 estimates. The histogram containing all 300 estimates can be found in Supplemental Material, Figure S5.

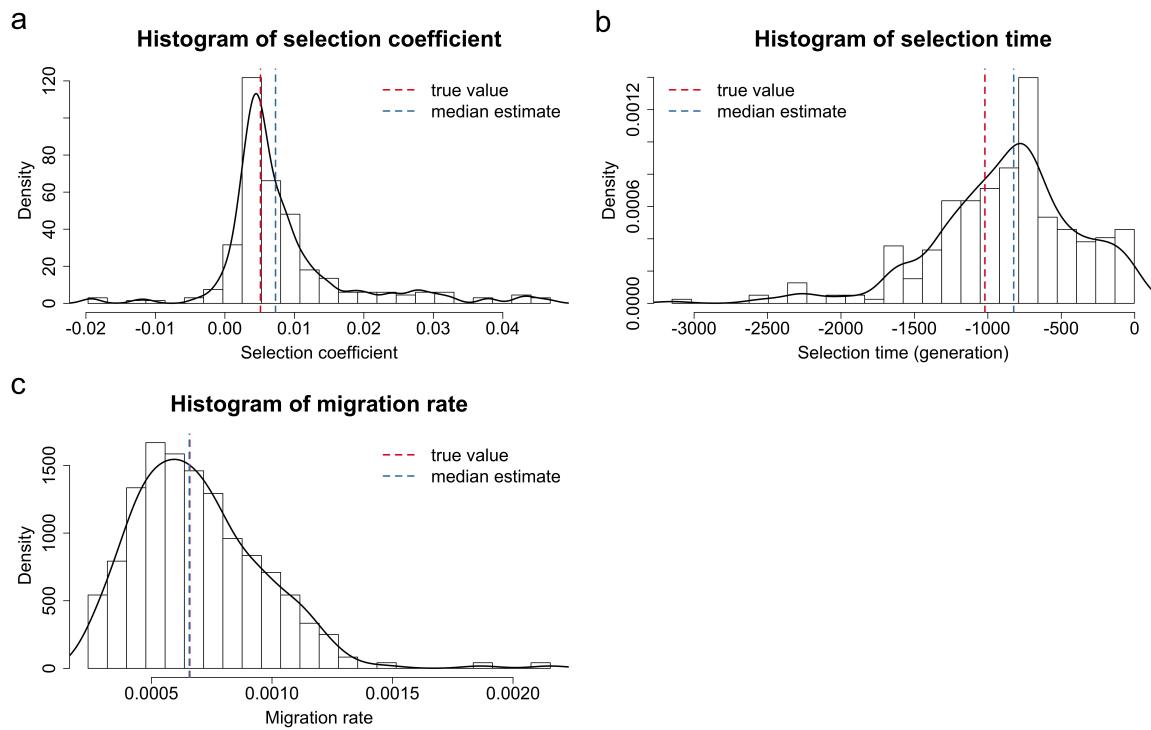


Figure 8: Empirical distributions of the estimates for 300 datasets simulated for *TSHR* based on the aDNA data shown in Table 2. We take the timing and strength of natural selection and gene migration to be those estimated with the population size $N = 180000$ shown in Table 3, but the true population size in the simulation is taken to be $N = 4500$. Histograms of (a) the selection coefficient estimates (b) the selection time estimates and (c) the migration rate estimates. To aid visual comparison, we have picked the x axis in (a) not to cover all 300 estimates. The histogram containing all 300 estimates can be found in Supplemental Material, Figure S6.

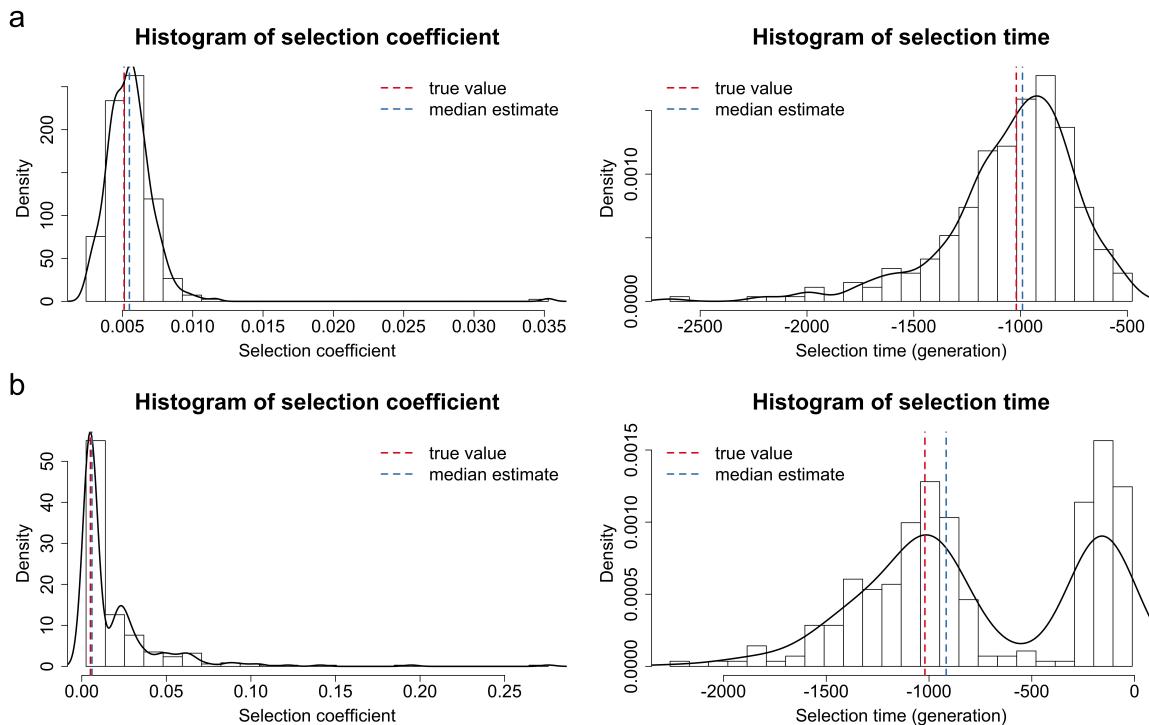


Figure 9: Empirical distributions of the estimates for 300 datasets simulated for *TSHR* based on the aDNA data shown in Table 2. We take the timing and strength of natural selection and gene migration to be those estimated with the population size $N = 180000$ shown in Table 3, but the true migration rate in the simulation is taken to be (a) $m = 0.0001$ and (b) $m = 0.001$. Histograms of the selection coefficient estimates and the selection time estimates for the migration rate (a) $m = 0.0001$ and (b) $m = 0.001$.

	\mathcal{A}_1^i	\mathcal{A}_2^i	\mathcal{A}_1^c	\mathcal{A}_2^c
\mathcal{A}_1^i	1	$1 - hs$	1	$1 - hs$
\mathcal{A}_2^i	$1 - hs$	$1 - s$	$1 - hs$	$1 - s$
\mathcal{A}_1^c	1	$1 - hs$	1	$1 - hs$
\mathcal{A}_2^c	$1 - hs$	$1 - s$	$1 - hs$	$1 - s$

Table 1: Relative viabilities of all possible genotypes at locus \mathcal{A} when we distinguish between the alleles that originate on the island and the alleles that emigrate from the continent.

↙
 not necessary ; can be written in the
 main text more compactly

For Review Only

Sample time	Sample size	Mutant allele
-128	12	8
-25	8	5
82	8	3
200	32	14
256	14	3
1067	6	0
1309	20	18
1650	2	1
1850	2	2
1975	14	14
1995	334	328

Table 2: Time serial European chicken samples of segregating alleles at the *TSHR* locus. The unit of the sampling time is the AD year.

	Population size	MAP	95% HPD
Selection coefficient	26000	0.005109	[0.003622, 0.007141]
	180000	0.005120	[0.003591, 0.007064]
	460000	0.005122	[0.003648, 0.006578]
Selection time	26000	-1047	[-1659, -857]
	180000	-1020	[-1384, -821]
	460000	-1047	[-1327, -893]
Migration rate	26000	0.000712	[0.000448, 0.000918]
	180000	0.000659	[0.000483, 0.000861]
	460000	0.000620	[0.000478, 0.000837]

Table 3: MAP estimates of the selection coefficient, the selection time and the migration rate, as well as their 95% HPD intervals, for *TSHR* achieved with the population size $N = 26000$, $N = 180000$ and $N = 460000$.