

# Jointly Inferring the Dynamics of Population Size and Sampling Intensity from Molecular Sequences

Kris V. Parag <sup>\*,†,1,2</sup>, Louis du Plessis,<sup>\*,†,1</sup> and Oliver G. Pybus<sup>\*,1</sup>

<sup>1</sup>Department of Zoology, University of Oxford, Oxford, United Kingdom

<sup>2</sup>Department of Infectious Disease Epidemiology, MRC Centre for Global Infectious Disease Analysis, Imperial College London, London, United Kingdom

<sup>†</sup>These authors contributed equally to this work.

\*Corresponding authors: E-mails: k.parag@imperial.ac.uk; louis.duplessis@zoo.ox.ac.uk; oliver.pybus@zoo.ox.ac.uk.

Associate editor: Keith Crandall

## Abstract

Estimating past population dynamics from molecular sequences that have been sampled longitudinally through time is an important problem in infectious disease epidemiology, molecular ecology, and macroevolution. Popular solutions, such as the skyline and skygrid methods, infer past effective population sizes from the coalescent event times of phylogenies reconstructed from sampled sequences but assume that sequence sampling times are uninformative about population size changes. Recent work has started to question this assumption by exploring how sampling time information can aid coalescent inference. Here, we develop, investigate, and implement a new skyline method, termed the epoch sampling skyline plot (ESP), to jointly estimate the dynamics of population size and sampling rate through time. The ESP is inspired by real-world data collection practices and comprises a flexible model in which the sequence sampling rate is proportional to the population size within an epoch but can change discontinuously between epochs. We show that the ESP is accurate under several realistic sampling protocols and we prove analytically that it can at least double the best precision achievable by standard approaches. We generalize the ESP to incorporate phylogenetic uncertainty in a new Bayesian package (BESP) in BEAST2. We re-examine two well-studied empirical data sets from virus epidemiology and molecular evolution and find that the BESP improves upon previous coalescent estimators and generates new, biologically useful insights into the sampling protocols underpinning these data sets. Sequence sampling times provide a rich source of information for coalescent inference that will become increasingly important as sequence collection intensifies and becomes more formalized.

**Key words:** coalescent processes, sampling models, skyline plots, demographic inference, influenza, bison, Bayesian phylogenetics.

## Introduction

The coalescent process describes how the size of a population influences the genealogical relationships of individuals randomly sampled from that population (Kingman 1982). Coalescent-based models are widely used in molecular epidemiology and ecology as null models of ancestry, and of the diversity of observed gene or genome sequences. In many instances, these sequences are sampled longitudinally through time from a study population, for example, when individual infections are sampled across an epidemic caused by a rapidly evolving virus or bacterium (Pybus and Rambaut 2009), or when ancient DNA is extracted from preserved animal tissue that may be tens of thousands of years old (Shapiro and Hofreiter 2014). If sequences accrue measurable amounts of genetic divergence between sampling times, then the data set is termed heterochronous (Drummond et al. 2003; Biek et al. 2015). A common problem in molecular evolution is the estimation of effective population size history from these heterochronous sequences or from time-scaled

genealogies (trees) that are reconstructed from those sequences.

Several coalescent-based approaches have been developed to solve this problem, including the popular and prevalent skyline and skygrid families of inference methods (Pybus et al. 2000; Strimmer and Pybus 2001; Drummond et al. 2005; Minin et al. 2008; Gill et al. 2013). These approaches, which originated with the classic skyline plot of Pybus et al. (2000), estimate population size history as a piecewise-constant function using only the coalescent event times (i.e., the tree branching times) of the reconstructed genealogy. For heterochronous data sets, these methods typically assume that the sequence sampling times (i.e., the tree tips) are defined by extrinsic factors such as sample availability or operational capacity (Ho and Shapiro 2011), and are thus uninformative about, and independent of, population size (Drummond et al. 2005; Parag and Pybus 2019).

Recent work has started to challenge this assumption and assess its consequences. Volz and Frost (2014) showed, for a coalescent process with exponentially growing population

size, that including sequence sampling time information can notably improve the precision of demographic parameter estimates, if the sampling process is correctly specified. They recommended an augmented coalescent sequence sampling model, and defined a *proportional sampling* process, in which the rate of sampling sequences at any time from a population is linearly dependent on its effective size at that time. Karcher et al. (2016) generalized this to include non-linear dependence, which they termed *preferential sampling*, and to allow for piecewise-constant effective population size changes. Karcher et al. (2016) cautioned that misleading inferences can result when the relationship between population size and the sequence sampling rate is misspecified.

Although these works have brought attention to the benefits of exploiting sampling time information for population size inference, further progress is needed. Previous studies have treated the sampling model as a statistical addition to the coalescent process (Karcher et al. 2016) and have not explicitly considered the types of sampling designs and surveillance protocols that are commonly implemented by epidemiologists and ecologists in the field. Moreover, there are to date few provable or general analytical insights into the joint inference of sampling and population size using coalescent models. A flexible model that can accurately assess the role of experimental and surveillance design is warranted as there is still uncertainty about what constitutes good rules for sequence sampling, and about the relative benefits and pitfalls of different sampling protocols (Stack et al. 2010; Hall et al. 2016; Parag and Pybus 2019). These issues will only increase in importance as sequence sampling intensifies and heterochronous data sets become more common (Ho and Shapiro 2011; Baele et al. 2017).

Here, we aim to advance the field by developing a new sampling-aware coalescent skyline model, which we term the “epoch sampling skyline plot” (ESP). The ESP extends the classic skyline plot to include a flexible epoch-based sampling model that can represent biologically realistic sampling scenarios. Its formulation also renders it amenable to theoretical exploration and straightforward implementation within a Bayesian phylogenetic MCMC framework. The ESP assumes that sampling occurs in epochs, which are defined as periods of time during which the sampling rate per capita is deemed constant. In practical applications, an epoch might, for example, represent weekly or monthly surveillance windows, epidemic seasons, archeological periods, or geological strata. The boundaries of each epoch are delineated by the sequence sampling times of the heterochronous genealogy. This guarantees model identifiability and helps guard against unsupported inferences by ensuring that the number of allowed per-capita sampling rate changes are fewer than the count of sampling events.

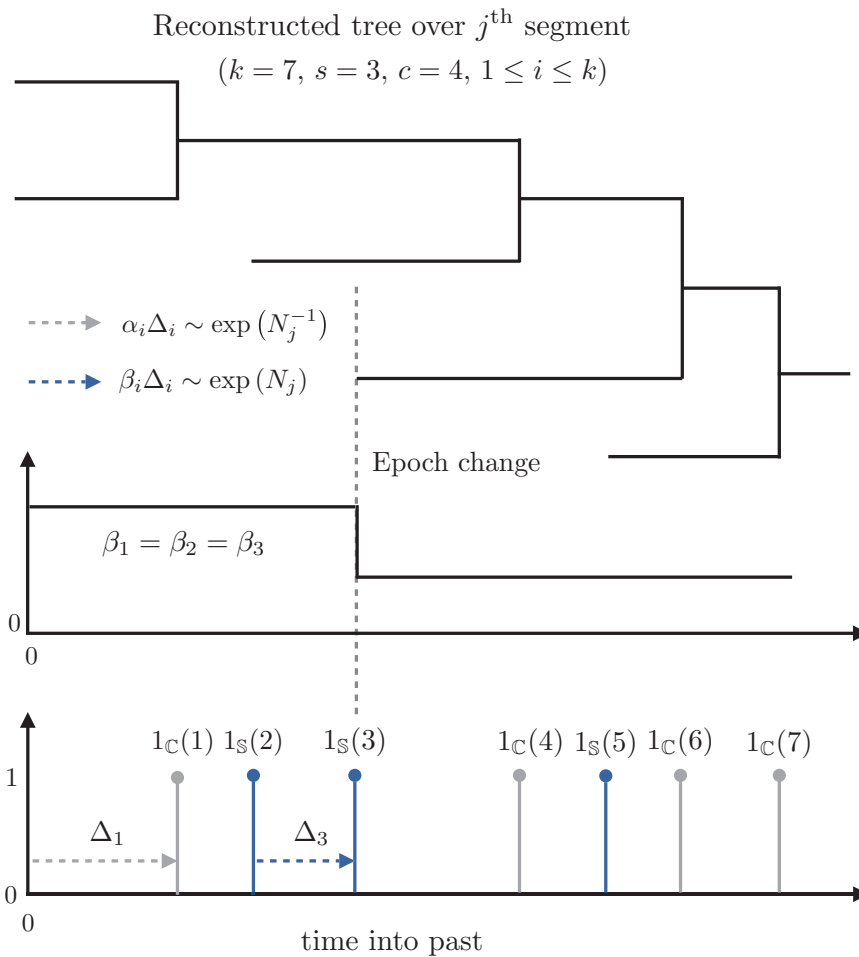
Within an epoch, the ESP assumes that tree tips are sampled in proportion to population size, with a constant of proportionality that we call the sampling intensity. This intensity measures the average sampling effort over the epoch per capita, with larger values corresponding to faster rates of sequence sample accumulation. We allow the sampling intensity to change discontinuously between epochs, resulting

in a flexible piecewise-constant sampling process. Within each epoch, the ESP locally models *density-defined sampling*, in which the sampling rate directly correlates with effective population size. Consequently, the ESP can describe a wide range of *time-varying density-defined* sampling protocols. This allows the ESP to: (1) account for external, population-independent fluctuations in sampling effort and (2) provide a means to quantify sampling effort through the testing of competing sequence collection hypotheses. For example, we may be interested in whether sampling intensity increases or decreases through time, or whether known historical or experimental events are associated with a change in sampling intensity.

Although the flexibility of the ESP means it can model a wide range of sampling models, we here give attention to two specific sampling models, inspired by real-world collection practices. The first is density-defined sampling, which embodies the assumption that the availability of sequences depends on the size of the study population, and leads to a fixed proportion of the population being sampled across the sampling time frame of the study. It can be modeled in the ESP by forcing all epochs to have identical sampling intensities or, equivalently, by defining a single epoch that spans the entire sampling period. Density-defined sampling is a simple sequence collection protocol and can be obtained directly from the proportional and preferential models of Volz and Frost (2014) and Karcher et al. (2016).

The second sampling model arises when studies aim to collect an approximately constant number of samples per unit time (e.g., week, epidemic season, or geological era), irrespective of the size of the study population. This protocol is called *frequency-defined sampling* and is modeled within the ESP by allocating epochs uniformly over time, and allowing their individual sampling intensities to vary such that the process samples an (approximately) deterministic number of samples per epoch. Frequency-defined sampling is often undertaken in molecular epidemiology when resources for surveillance are limited or predefined; or when the primary research aim is to diagnose and classify infections or to provide snapshots of the genetic diversity of pathogen populations (Ho and Shapiro 2011). Frequency-defined sampling cannot be described within previous frameworks and can represent the impact of external factors on the rate of sample collection.

In this article, we develop and define the ESP and show how it facilitates the joint inference of effective population sizes and sampling intensities, within maximum likelihood (ML) and Bayesian frameworks. We validate its performance using simulated data, before exploring its improvements over existing skyline-based methods on empirical data sets (H3N2 influenza A virus sequences from New York state, and ancient mtDNA sequences from Beringian steppe bison). We focus on biologically inspired sampling protocols (see above) and demonstrate how the epoch sampling model facilitates the testing and exploration of different data collection hypotheses. We highlight how the inverse relationship between the rates of sampling and coalescence can substantially improve



**Fig. 1.** Illustration of the epoch sampling skyline plot model. A temporally sampled (heterochronous) tree consists of sampled tips and coalescing lineages. A portion of this tree is shown (top). This portion covers the  $j$ th segment, during which the effective population size is assumed to be fixed at  $N_j$ . Our epoch model assumes a piecewise-constant sampling intensity function which, in this illustration, comprises two epochs over this tree segment (middle). The sampling times (blue, bottom) provide information about the sampling intensities in each epoch and also determine the epoch boundaries. The coalescent event times (gray, bottom) allow inference of  $N_j$  and also delimit the segment boundaries. See New Approaches for definitions of the mathematical notation used.

population size estimation bias, especially when coalescent events are sparse, and prove that the information available for inferring population size (and hence the precision of those estimates) could more than double by using the ESP. Finally, we describe in detail both ML and Bayesian ESP implementations. The latter is available as an integrated package called BEBP in the popular Bayesian phylogenetics platform BEAST2 (Bouckaert et al. 2019).

### New Approaches

Consider a coalescent tree reconstructed from sequences sampled longitudinally through time. Let the effective population size underlying this process at time  $t$ , into the past, be  $N(t)$ . Standard coalescent skyline-based approaches to estimating  $N(t)$  assume that sequence sample times are uninformative (Drummond et al. 2005) and therefore draw all of their inferential power from the reconstructed coalescent event times. These methods approximate  $N(t)$

with a piecewise-constant function comprising  $p$  segments:  $\sum_{j=1}^p N_j 1_{[t_{j-1}, t_j)}(t)$ , where  $t_j - t_{j-1}$  is the duration of the  $j$ th segment and  $1_{\mathbb{A}}(x)$  is an indicator variable, which equals 1 if  $x \in \mathbb{A}$  and is 0 otherwise, for some set  $\mathbb{A}$ . Here,  $t_0 = 0$  is the present. Figure 1 illustrates a coalescent subtree spanning the  $j$ th segment, during which the effective population size is  $N_j$ . Two epochs with distinct sampling intensities occur within this segment. The coalescent event times (gray) form the branching points of the reconstructed tree, whereas sampling events (blue) determine when new tips are introduced.

We use  $\Delta_i$  to denote the duration of the  $i$ th interevent period or interval within a given segment, and define the lineage count in this interval as  $\ell_i$ . If there are  $k$  intervals in the  $j$ th segment, then  $t_j - t_{j-1} = \sum_{i=1}^k \Delta_i$ . We use the sets  $\mathbb{S}$  and  $\mathbb{C}$  to indicate whether an interval ends with a sampling or coalescent event, respectively. Then  $s = \sum_{i=1}^k 1_{\mathbb{S}}(i)$  and  $c = \sum_{i=1}^k 1_{\mathbb{C}}(i)$  count the number of sampling and coalescent

events in a given interval, and  $k = s + c$ . Note that,  $s, c$ , and  $k$  are not fixed, and can have different values for all  $p$  segments. Events which occur at a change-point belong to the interval that precedes that change-point, that is, the one closer to the present. Hence the sampling events  $1_{\mathbb{S}}(2)$  and  $1_{\mathbb{S}}(3)$  in [figure 1](#) belong to the first epoch, and the starting two lineages are included in the likelihood of the  $(j - 1)$ th segment.

Coalescent events falling within the  $j$ th segment follow a Poisson process with rate  $\alpha_i N_j^{-1}$ , with  $\alpha_i := \binom{\ell_i}{2}$ , and  $N_j$ , as the unknown effective population size during that segment ([Kingman 1982](#)). As a result,  $\alpha_i \Delta_i \sim \exp(N_j^{-1})$  describes the key informative relationship in coalescent processes. Standard skyline methods capitalize on this dependence, but assume that intervals ending in sampling events (i.e., those satisfying  $\{\Delta_i : i \in \mathbb{S}\}$ ) are uninformative. Under this assumption, the maximum Fisher information about  $N_j$  that can be extracted by these methods is  $cN_j^{-2}$  ([Parag and Pybus 2017](#)).

The ESP instead posits that the sample times within the  $i$ th interval of the  $j$ th segment derive from a Poisson process of rate  $\beta_i N_j$ . Here,  $\beta_i$  is the sampling intensity governing the average sampling effort (per capita or unit of  $N_j$ ) made across  $\Delta_i$ . This encodes the extra informative relationship:  $\beta_i \Delta_i \sim \exp(N_j)$ , and is the most complex sampling model admissible within the skyline framework (i.e., it is maximally parametrized). We remove unnecessary complexity by defining an epoch as a grouping of consecutive intervals (which may span multiple segment boundaries) over which the sampling intensity is constant. Thus, within an epoch, all  $\beta_i$  take the same value (in [fig. 1](#) there are two epochs). We force epoch change times to coincide with sequence sampling times and assume that no sampling effort was made before the most ancient sample, that is, we set  $\beta_i = 0$  for all intervals from the most ancient sample to the last coalescent event time (the time of the most recent common ancestor of the tree).

This description guarantees that the ESP is maximally flexible yet statistically identifiable ([Parag and Pybus 2019](#)), because every skyline segment and epoch has at least one coalescent and one sampling event, respectively (see Results). Our epochal model, unlike previous attempts at incorporating sample times ([Volz and Frost 2014](#); [Karcher et al. 2016](#)), can account for the temporal heterogeneity of sampling protocols undertaken in real-world studies. For example, sampling often occurs in bursts with discontinuous sampling effort that changes between collection periods. In the ESP, the sampling intensities of the epochs are independent of one other. Using this framework, we construct the ESP log-likelihood for the  $j$ th segment,  $\mathcal{L}_j = \log P(\mathcal{T} | N_j)$ , as in [equation \(1\)](#), with  $\mathcal{T}$  as the reconstructed tree.

$$\mathcal{L}_j = \sum_{i=1}^k 1_{\mathbb{S}}(i) \log(\beta_i N_j) + 1_{\mathbb{C}}(i) \log(\alpha_i N_j^{-1}) - \Delta_i(\beta_i N_j + \alpha_i N_j^{-1}) \quad (1)$$

The complete tree log-likelihood is  $\mathcal{L} = \sum_{j=1}^p \mathcal{L}_j$ . The waiting time until the end of any interval contributes the

$-\Delta_i(\beta_i N_j + \alpha_i N_j^{-1})$  term, whereas sampling and coalescent events introduce terms  $1_{\mathbb{S}}(i) \log(\beta_i N_j)$  and  $1_{\mathbb{C}}(i) \log(\alpha_i N_j^{-1})$ , respectively. If we define  $p'$  epochs over  $\mathcal{T}$ , then there are  $p + p'$  unknown parameters in our log-likelihood (the set of  $N_j$  and distinct, nonzero  $\beta_i$ ). [Equation \(1\)](#) is related to the augmented log-likelihood from [Karcher et al. \(2016\)](#) but differs in both the population size and sampling models used.

The ESP is obtained from [equation \(1\)](#) by computing the grouped ML estimate (MLE),  $\hat{N}_j$ , for each segment. This involves solving a pair of quadratic equations that depend on the relative number of sampling and coalescent events in that segment,  $s - c$ . Defining  $a = \sum_{i=1}^k \alpha_i \Delta_i$  and  $b = \sum_{i=1}^k \beta_i \Delta_i$ , we obtain [equation \(2\)](#), from the roots of these quadratics (see [eqs. 11 and 12](#) in Materials and Methods).

$$\hat{N}_j = \begin{cases} \frac{s-c}{2b} + \sqrt{\left(\frac{s-c}{2b}\right)^2 + \frac{a}{b}} & \text{if } s \geq c \\ \left(\frac{c-s}{2a} + \sqrt{\left(\frac{c-s}{2a}\right)^2 + \frac{b}{a}}\right)^{-1} & \text{if } s < c \end{cases} \quad (2)$$

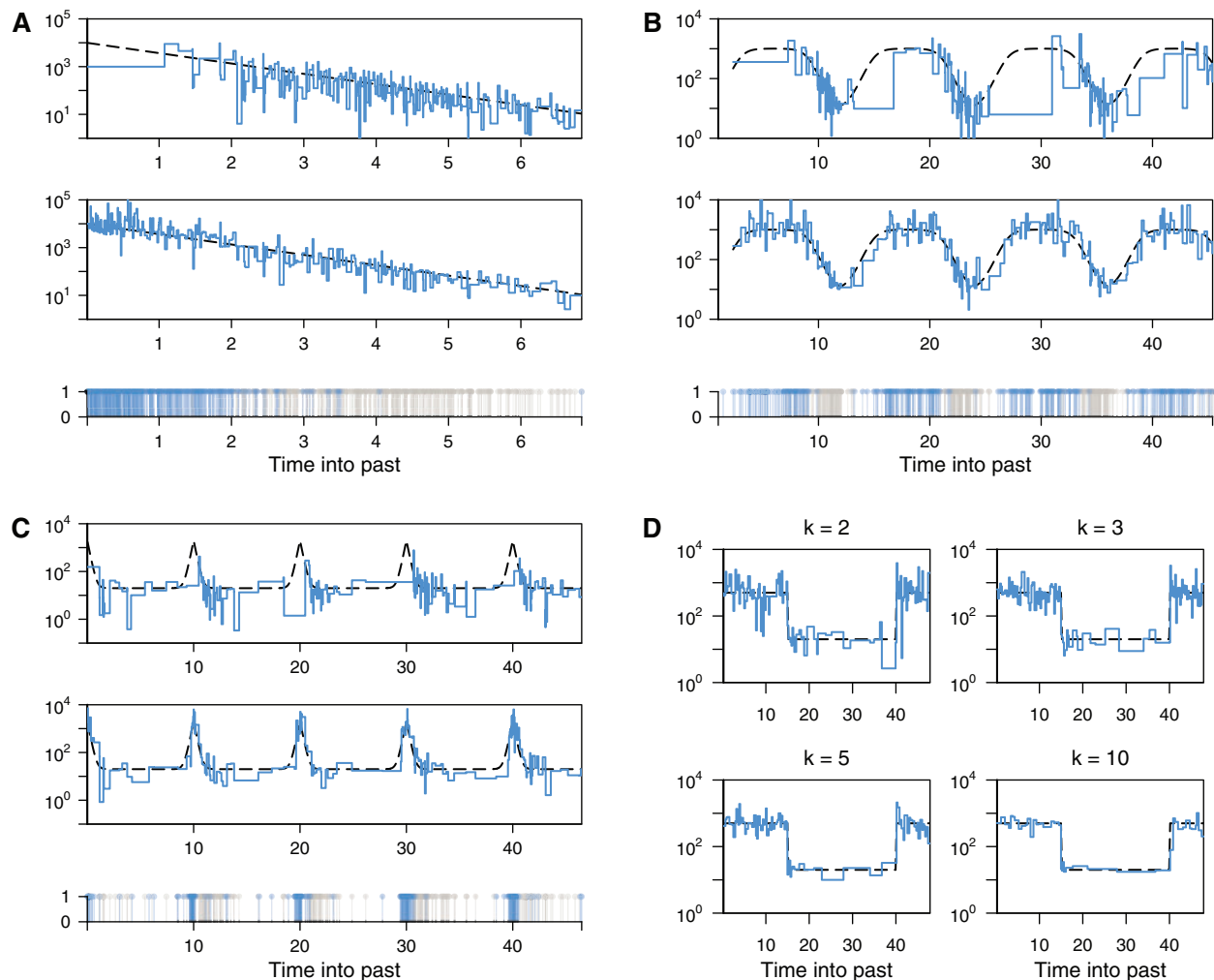
[Equation \(2\)](#) forms our main result and requires the MLE of each  $\beta_i, \beta_i$ , to be jointly estimated (see Materials and Methods for appropriate algorithms). If  $s = c$ , both parts of [equation \(2\)](#) converge to the simple square root estimator,  $\hat{N}_j = \sqrt{ab^{-1}}$ . Grouping over  $k$  adjacent intervals in our skyline leads to smoother population size estimates that are quick to compute and easy to generalize. Note that, if  $s = 0$ , all  $\beta_i = 0$  and  $c = 1$ , then [equation \(2\)](#) simplifies to the classic skyline plot estimator of  $N_j$  ([Pybus et al. 2000](#)).

The ESP has several desirable properties. Its counteracting proportional and inverse dependence on  $N(t)$  means that it has more informative intervals during time periods when coalescent events are infrequent, which otherwise hinders standard skyline inference. This property spreads the information about  $N(t)$  more uniformly through time and reduces estimator bias. The ESP can also significantly improve overall estimate precision. The Fisher information that the ESP extracts from the  $j$ th segment of the reconstructed tree is now at least  $(s + c)N_j^{-2}$  (see Results for analytic derivation).

## Results

### Simulated Performance

We start by comparing the estimates from [equation \(2\)](#) to those of the classic skyline plot ([Pybus et al. 2000](#)), which ignores the information in sequence sampling times and is the basis of several popular skyline methods. We keep the number of piecewise-constant segments (parameters) inferred in the ESP (model dimensionality) approximately the same as that of the classic skyline plot by fixing  $k = 2$ . For clarity, we assume a single, known sampling intensity and examine only the period more recent than the most ancient sampling time. We compare the abilities of the ESP and classic skyline plot methods to recover a variety of population size dynamics in [figure 2A–C](#). In each panel ([fig. 2A–C](#)), the top graph gives the classic skyline plot estimate, the middle one shows the ESP estimate (for the same fixed tree), and the



**Fig. 2.** ESP and classic skyline plot estimates. Panels (A–C) compare the performance of the classic skyline plot (top graph) to the ESP at  $k = 2$  (middle graph) for a range of demographic models: (A) exponential growth, (B) cyclical logistic growth, and (C) steep periodic dynamics. Estimates of  $N(t)$  are shown in blue and the true demographic functions are in dashed black, on a logarithmic scale. The classic skyline plot performs poorly near the present in (A), or when there are notable fluctuations between large and small population sizes in (B and C). This results from the uneven temporal distribution of coalescent events (gray lines in the bottom graph of each panel). The sampling events (blue lines in the bottom graph of each panel) are inversely distributed to coalescent events. Consequently, the ESP tracks changes in population size more accurately, for the same number of population size segments. Panel (D) shows how increasing the grouping of adjacent intervals,  $k$ , can improve the smoothing of the ESP, in the context of a stepwise demographic function. All trees were simulated using the phylodyn R package (Karcher et al. 2017) with  $\sim 300$  coalescent and sampling events.

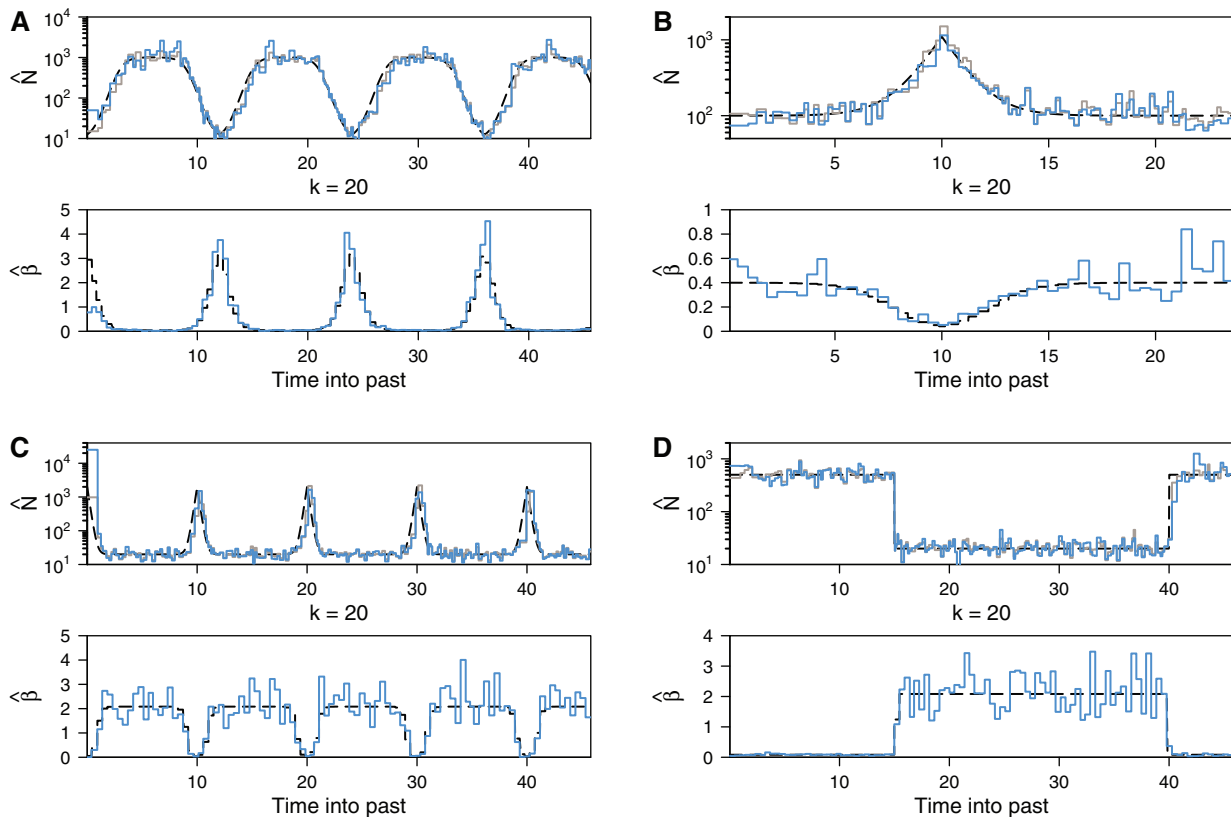
bottom one plots the distribution of sampling (blue) and coalescent (gray) event times.

The ESP significantly improves demographic inference, relative to skyline plot methods, when population size is large (fig. 2A) and in periods featuring sharp demographic changes (fig. 2C). In these scenarios, standard skyline or skygrid approaches are known to perform poorly because coalescent events, due to their inverse dependence on population size, are sparse and hence unable to capture these population dynamics. Accordingly, coalescent events also tend to cluster around bottlenecks (fig. 2B), and so cause standard methods to lose fidelity across cyclic epidemics. Sampling events, however, fall in periods of sparse coalescence, allowing the ESP to circumvent these problematic conditions.

The generalized skyline plot was introduced in Strimmer and Pybus (2001) to ameliorate the noisy nature

of the classic skyline plot. It grouped adjacent intervals to achieve a bias-variance trade-off that led to smoother estimates of  $N(t)$ . This grouping is used in some popular skyline approaches, notably the Bayesian skyline plot (BSP) (Drummond et al. 2005). We achieve a similar smoothing effect in the ESP by increasing the grouping parameter,  $k$  (see fig. 2D). This extends the generalized skyline plot approach in two ways; first by incorporating sampling time information and second by including the specific times of events within a grouped interval.

Having clarified the attributes of the ESP, we now investigate examples in which the sampling intensities are unknown and can vary through time. We assume that the times corresponding to all sampling events are available for analysis. We consider two realistic, and widely used sampling protocols, which we, respectively, refer to as *density-defined* and



**Fig. 3.** Joint inference of effective population size and sampling intensity using the ESP. Panels (A–D) show estimates of population size (upper) and sampling intensity (lower) through time. A single, fixed tree was simulated under a frequency-defined sampling model for demographic scenarios featuring: (A) cycles of logistic growth and decline, (B) exponential growth and decline (the boom-bust model), (C) steep periodic cycles, and (D) a stepwise population size change. Simulations in (A) and (C) comprised 2,000 sampled tree tips over four population cycles, with  $p' = 100$  and  $k = 20$ . Simulations in (B) and (D) comprised 1,000 sampled tree tips, with  $p' = 50$ . The upper graph in each panel compares the true  $N(t)$  demographic function (dashed black) to  $\hat{N}$  when  $\beta$  is known without error (gray), and  $\hat{N}$  (blue) when it is coestimated with  $\beta$  within the ML framework (see eq. 2 and Materials and Methods). The lower graphs show the corresponding plots of  $\hat{\beta}$  (blue) against the true sampling intensity  $\beta$  (dashed black).

*frequency-defined* samplings. In the first, there is a direct correlation between the time-varying effective population size and the rate of sampling, and a single sampling intensity persists throughout the complete sampling period. Density-defined sampling is the simplest model described within the ESP framework. It represents the process of proportional sampling (i.e., more samples are taken if the population to be sampled is larger).

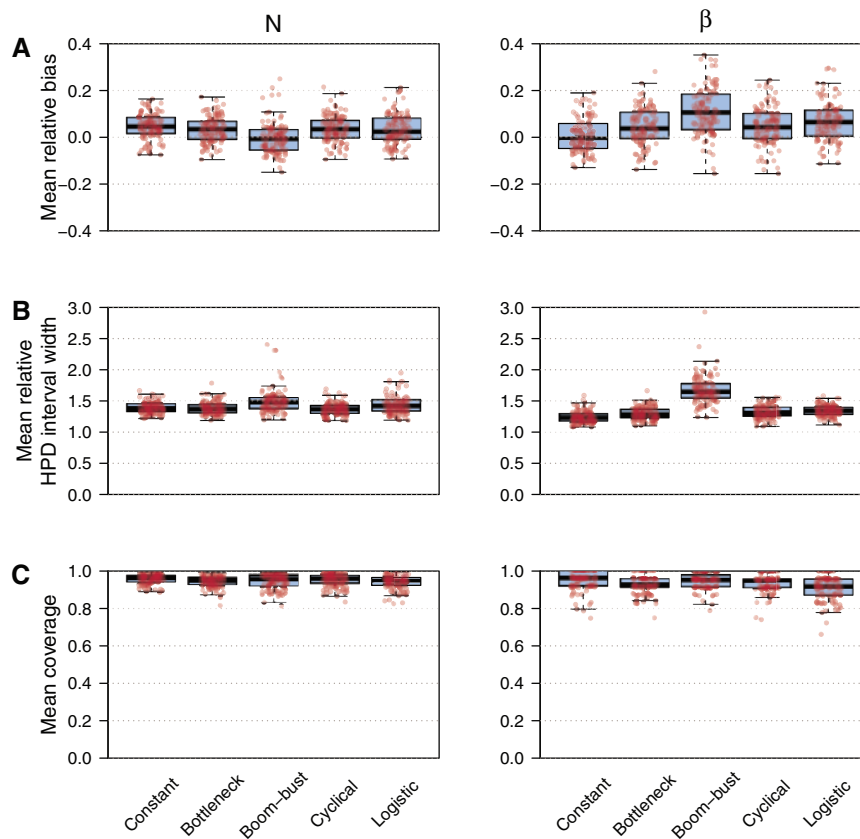
However, in many epidemiological scenarios, surveillance organizations or treatment centers will often examine a relatively fixed number of samples per unit time (e.g., per month or epidemic season). This number may be constrained by extrinsic factors such as funding or operational capacity. Similar constraints may control the availability of ancient DNA sequences generated by molecular evolutionary studies. In such circumstances, frequency-defined sampling results and the sampling intensity temporally fluctuates due to underlying changes in population size. As this sampling scheme is more complex (it is a *time-varying density-defined* model), we use it to validate ESP performance. For clarity, in this section, we restrict our analysis to fixed, time-scaled trees that are assumed to be known without error and apply our ML approach (see Materials and Methods for details). In later

sections, we examine both sampling models using a Bayesian implementation of the ESP that incorporates phylogenetic uncertainty.

We assume  $p'$  epochs, so there are  $p'$  unknown sets of  $\beta_i$  values to infer (within each epoch all  $\beta_i$  take the same value). We use  $\beta$  to represent this vector of unknowns, and let its MLE be  $\hat{\beta}$ . Note that, epoch and population size change-points are not synchronized (i.e., they are generally nonoverlapping), and we are jointly estimating a total of  $p + p'$  parameters. Figure 3A–D presents our joint estimates of  $N$  and  $\beta$  at  $k = 20$  for heterochronous genealogies simulated under four different demographic scenarios with frequency-defined sampling at  $p' = 100$  (fig. 3A and C) or  $p' = 50$  (fig. 3B and D) (see figure legend for details). As the sample count in each epoch is approximately the same, the  $\hat{\beta}$  estimates (lower graphs of fig. 3) take a complementary form to the  $\hat{N}$  ones (upper graphs). These results show that the ESP has the ability to faithfully reproduce changes in both population size and sampling intensity through time.

### Bayesian Implementation Simulation Study

Having explored the ML performance of the ESP, we now investigate and validate a Bayesian implementation of the



**FIG. 4.** Boxplots and stripcharts showing measures of statistical performance for the BEBP, evaluated on trees simulated under five different demographic models (constant, bottleneck, boom-bust, cyclical boom-bust, logistic growth, and decline). We simulated 100 replicate trees for each scenario. Three measures of performance are shown: (A) mean relative bias, (B) mean relative HPD interval size, and (C) mean coverage. Left and right columns illustrate estimation performance for effective population size ( $N$ ) and sampling intensity ( $\beta$ ), respectively.

ESP, which we call the BEBP (see Materials and Methods). The BEBP incorporates the ESP log-likelihood within the computational framework of BEAST2. In this section, we benchmark the ability of the BEBP to recover accurate and unbiased parameter estimates. We simulated 100 replicate coalescent genealogies (using the *phylodyn* R package; Karcher et al. 2017) under five demographic scenarios: 1) constant size, 2) bottleneck, 3) boom-bust, 4) cyclical boom-bust, and 5) logistic growth and decline. In all simulations, we used frequency-defined sampling with approximately equal numbers of samples split over 24 equidistant epochs. We jointly inferred  $N$  and  $\beta$  from each simulated tree using the BEBP and assumed that trees were known without error (to render the simulations computationally feasible, and to distinguish uncertainty in the coalescent model from phylogenetic noise). Estimation of  $N$  and  $\beta$  directly from sets of empirical gene sequences is demonstrated in the next section.

We grouped coalescent and sampling events into  $p = 100$  equally informed population size segments (i.e.,  $k$  is equal for all segments) to estimate  $N$  and used  $p' = 24$  approximately equidistant sampling epochs for  $\beta$ . To quantify the bias and precision of the BEBP method, we computed the relative bias, the relative highest posterior density (HPD) interval width and the coverage of estimates of  $N$  and  $\beta$ , averaged across the time between the most recent and most ancient samples.

Further details on the simulations, inferences, and summary statistics can be found in the [Supplementary Material](#) online. The results of our simulation study are summarized in [figure 4](#). Example simulated trees and inferred parameter trajectories are shown in [supplementary figures S1–S5](#), [Supplementary Material](#) online (see [https://github.com/laduplessis/BESP\\_paper-analyses/](https://github.com/laduplessis/BESP_paper-analyses/); last accessed November 22, 2019, for simulated trees and inferred parameter trajectories for all replicates).

Both  $N$  and  $\beta$  appear to be slightly overestimated with a larger bias in the  $\beta$  estimates. Nonetheless, the boxplots for the mean relative bias intersect 0 for all five demographic scenarios, verifying acceptable accuracy. The mean relative HPD interval widths of the population size estimates are  $< 2$  for all replicate cases, with only a few outliers. Relative HPD intervals  $< 2$  indicate that estimates are at least twice as precise as a standard Gaussian approximation (the width under a Gaussian distribution with SD equal to the absolute value of the parameter is  $\approx 3.92$ ). Estimates of  $\beta$  under the boom-bust scenario occasionally have relative HPD interval widths  $> 2$ . We found this to be a consequence of the BEBP not having sufficient power to precisely estimate  $\beta$  during the most recent sampling epoch (see [supplementary fig. S3](#), [Supplementary Material](#) online, for an illustrative example of this effect). Lastly, the mean coverage is always close to 1, indicating that the true  $N$  and  $\beta$  values are included within

the HPD intervals for the majority of the sampling period. These results verify that the BEBP exhibits comparatively low bias and high precision.

### Case Study 1: Seasonal Human Influenza

Human influenza A virus (IAV) is a leading threat to global public health, causing an estimated 290,000–650,000 deaths per year (WHO 2018). Two subtypes of IAV currently cocirculate worldwide (H3N2 and H1N1-pdm) which, in temperate regions, cause annual winter epidemics. Strong immune pressure on the virus surface glycoprotein *hemagglutinin* (HA) drives a continuous replacement of circulating strains with new variants, termed antigenic drift (Ferguson et al. 2003). Rambaut et al. (2008) reported 1,302 complete genomes of A/H3N2 and A/H1N1 viruses that were sampled longitudinally through time from temperate regions (specifically New York state and New Zealand) and analyzed the dynamics of IAV genetic diversity using the BSP (Drummond et al. 2005).

Rambaut et al. (2008) found that the BSP could recover cyclical evolutionary dynamics from these sequences, with an increase in genetic diversity at the start of each winter influenza season, followed by a bottleneck at the end of that season, although the cycles were not sharply defined. Subsequently, Karcher et al. (2016) showed that estimates of IAV effective population size could be improved by incorporating sequence sampling time information within a preferential sampling model. However, that analysis assumed density-defined sampling and conditioned on the tree being known without error, thus eliminating phylogenetic noise. Here, we extend the analysis of this data set by using our BEBP approach to coestimate the effective population size history and sampling intensity across epidemic seasons of A/H3N2 HA genes sampled from New York state.

As with the BSP, the population size parameter of the BEBP,  $N$ , is proportional to the effective population size in the absence of natural selection ( $N_e$ ), that is,  $N = N_e\tau$  where  $\tau$  is the average generation time. This assumption does not hold for human IAV HA genes, which are subjected to strong directional selection. We follow previous practice and instead interpret  $N$  as a measure of relative genetic diversity (Rambaut et al. 2008). Our data set comprises an alignment of 637 HA gene sequences (1,698 nt long) sampled across 12 complete influenza seasons, from 1993/1994 to 2004/2005 (fig. 5A). Our estimates are inferred directly from the heterochronous sequence alignment using MCMC sampling and therefore incorporate phylogenetic uncertainty. Substitution and clock models are similar to those in Rambaut et al. (2008) (see Supplementary Material online for model details). We estimate a BEBP with  $p = 40$  population size segments and  $p' = 12$  sampling epochs, so that each epoch corresponds approximately to the duration of one influenza season.

As figure 5A shows, considerably fewer sequences were sampled during the 1995/1996, 2000/2001, and 2002/2003 influenza seasons. The inferred dynamics of A/H3N2 genetic diversity (fig. 5B) are strongly cyclical, with peaks coinciding with the midpoint of each epidemic season, except for 2000/2001 and 2002/2003. These results agree with epidemiological

surveillance data for New York and New Jersey states, which show that nearly all infections during the 2000/2001 season were caused by A/H1N1 and influenza B viruses, and that the 2002/2003 season was dominated by A/H1N1 infections (CDC 2019). We do infer a clear peak for A/H3N2 in the 1995/1996 season (fig. 5B), reflecting the fact that influenza cases during the 1995/1996 season were a mixture of A/H1N1 and A/H3N2 infections (Ferguson et al. 2003), which resulted in an intermediate number of sequences being sampled that year (fig. 5A).

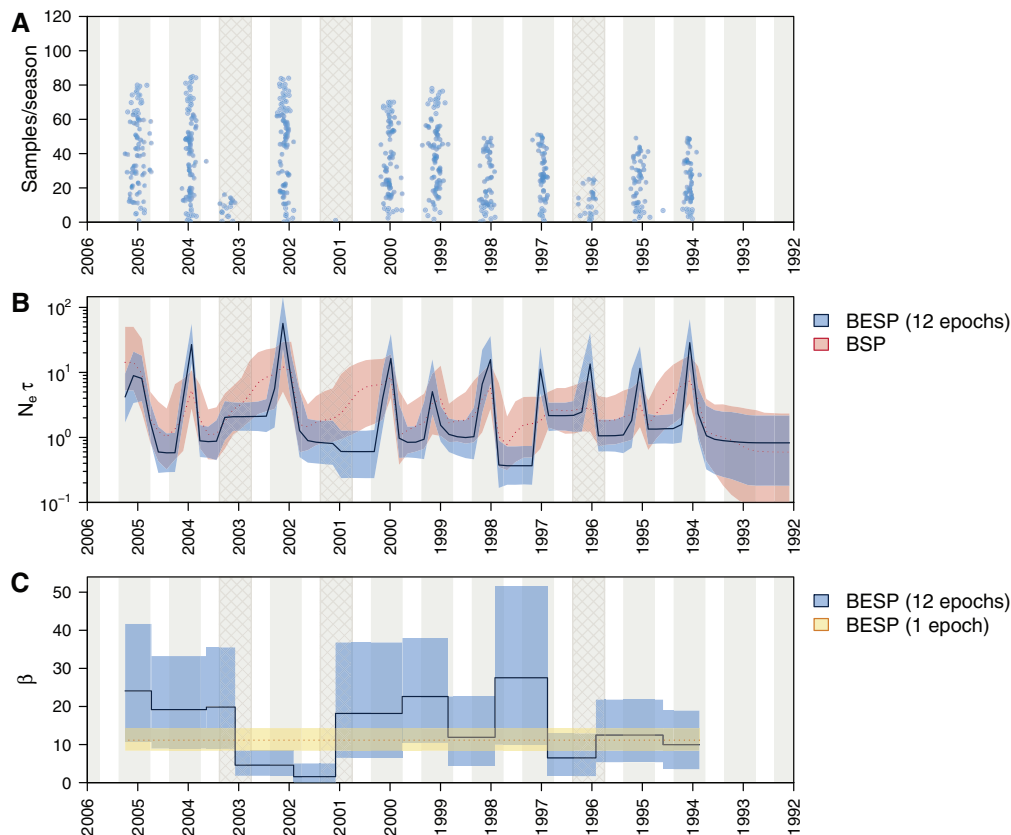
A comparison of the BEBP and BSP estimates of  $N_e\tau$  (fig. 5B) on the same data set shows that the BEBP infers an epidemic peak for 1996/1997, whereas no such peak was revealed by the BSP. This indicates that the BEBP has greater inferential power. Further, the peaks in the BEBP are typically more defined than those in the BSP and have narrower 95% HPD intervals. Specifically, in the BEBP, genetic diversity drops more sharply at the end of each season. This agrees well with our simulation results (see fig. 2B), as coalescent events tend to be sparse when population sizes are large (e.g., at the start of a bottleneck), but sampling events are plentiful. Unlike the BEBP, the BSP cannot exploit these informative sampling events and fails to efficiently track the fall in the number of infections.

The relative genetic diversity at the epidemic trough varies little among years, although it appears higher during 2002 and lower during 1997. It is possible that the bottleneck level largely depends on the availability of data since, in the absence of coalescent and sampling events, the smoothing prior maintains a roughly constant population size estimate (Volz and Frost 2014). As the informative events in a given season mostly stem from sequences sampled during that season (see supplementary figs. S7 and S8, Supplementary Material online), the BEBP reveals no information about population dynamics prior to the first sampled season (1993/1994).

The inferred sampling intensities,  $\beta$ , for each season, are given in figure 5C. Except for the period from 2001 to 2003 (which includes both of the seasons without an inferred epidemic peak), the 95% HPD intervals of the estimated  $\beta$  values for each season are overlapping. The estimated  $\beta$  for 1996 also appears lower, however, the 95% HPD interval still overlaps with other seasons. Although there is some variation in the median estimates, the uncertainty in these estimates is large, especially when  $\beta$  is high.

We also analyzed the same data set using a simpler single-epoch model (i.e., density-defined sampling with a constant  $\beta$  through time). We found that the  $N_e\tau$  dynamics estimated using this simpler model (supplementary fig. S6B, Supplementary Material online) closely matches those inferred using the more complex 12-epoch model. The estimated sampling intensities obtained under the single- and 12-epoch models are also congruent (fig. 5C and supplementary fig. S6, Supplementary Material online). The density-defined model estimates a median sampling intensity of 11.16 (95% HPD 8.38–14.32), whereas the mean–median estimate of the 12-epoch model is 14.87 (mean 95% HPD 6.03–27.19). We conclude that variation in sampling intensity through time is comparatively weak. Thus, if the aim of the





**FIG. 5.** (A) Density of sequence sampling dates through time for the alignment of 637 A/H3N2 HA sequences from New York state that we analyzed. Blue dots indicate stripcharts of individual samples for each season. The stripchart heights give the number of samples in each season. Gray shading indicates the approximate period of influenza observation in New York state during each season (epidemiological week 40, to week 20 in the next year). Cross-hatched seasons are those where A/H3N2 was not the dominant influenza virus subtype. (B) Median (solid/dotted line) and 95% highest posterior density (HPD) intervals (shaded areas) for the genetic diversity estimates ( $N_e \tau$ ) through time. The BESP estimate is shown in blue and the BSP estimate is in red. (C) Median (solid line/dotted line) and 95% HPD intervals (shaded areas) of the estimated sampling intensities ( $\beta$ ) for each sampling epoch. The 12-epoch BESP estimates are shown in blue and a single-epoch (density-defined) estimate is in yellow.

original study authors was to undertake a density-defined sampling protocol, then our  $\beta$  estimates provide an independent validation that this aim was, at least approximately, achieved.

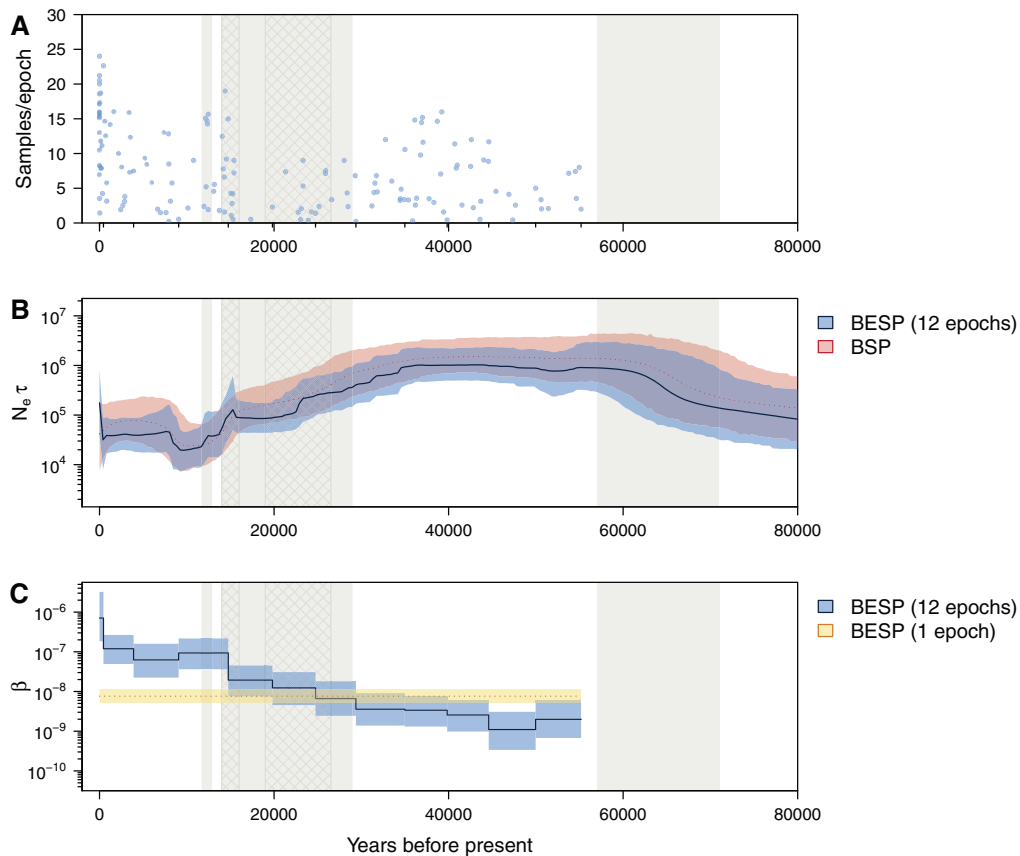
### Case Study 2: Steppe Bison

To illustrate the application of the ESP model to nonvirus data sets, we now analyze a heterochronous alignment of mtDNA genomes from modern and ancient bison that has previously been used to evaluate the performance of skyline-based methods (Shapiro et al. 2004; Drummond et al. 2005; Gill et al. 2013; Faulkner et al. forthcoming). During the Late Pleistocene, Beringia (eastern Siberia, the Bering land bridge, Alaska, and northwestern Canada) supported a large diversity of megafauna including bison, horses, and mammoths. A favorable climate for specimen preservation means that bison fossils suitable for ancient DNA extraction are abundant across the region (Shapiro et al. 2004). Sequences tens of thousands of years old can be recovered and dated with high confidence using radiocarbon dating (Shapiro and Hofreiter 2014). Reconstructing the past population dynamics of bison in this region can help clarify, and improve our

understanding of, the contributions of climate change and human presence to megafaunal population decline.

The data set we use is the same as that in Gill et al. (2013) and consists of mtDNA control region sequences from 135 ancient and 17 modern bison samples, with the oldest sample dated 55,182 years before present (BP). We treat sequence sampling dates as known and use the BESP to jointly infer the effective population size trajectory and sampling intensity through time, with  $p = 20$  segments, and  $p' = 12$  epochs. Each epoch lasts  $\sim 5,000$  years, except for the most recent, which stretches from the present to 450 years BP. We compared our population size estimates to a BSP with 20 population size segments. Both analyses used an HKY+ $\Gamma$  substitution model and a strict molecular clock (see Supplementary Material online for further model details).

Figure 6A shows that sequence sampling has been approximately constant through time, except for the most recent epoch (0–450 years BP), which contains the most samples, and the period 17–22 ka BP, which contains only three samples. This period coincides with the end of the last glacial maximum (LGM), whence fossil material is sparse (Shapiro et al. 2004). The BESP estimates of  $N_e \tau$  and  $\beta$  through time



**FIG. 6.** (A) Density of sequence sampling dates through time for the alignment of 152 bison mtDNA sequences that we used. Blue dots indicate stripcharts of individual samples for each sampling epoch. The heights of the stripcharts are equal to the number of samples in each epoch. Small tick marks on the  $x$ -axis represent epoch times. Grey shading indicates cool periods in the Earth's climate (from the present: Younger Dryas, Marine Isotope Stages (MIS) 2, MIS 4). The two cross hatched areas delimit the time of the last glacial maximum ( $\approx 26.5$ – $19$  ka BP) and approximate time of substantial human settlement of the Americas ( $\approx 16$ – $14$  ka BP). (B) Median (solid/dotted line) and 95% highest posterior density (HPD) intervals (shaded areas) for the genetic diversity estimates ( $N_e\tau$ ) through time. The BESP estimate is shown in blue and the BSP estimate in red. (C) Median (solid line/dotted line) and 95% HPD intervals (shaded areas) of the estimated sampling intensities ( $\beta$ ) for each sampling epoch. The 12-epoch BESP estimates are in blue and a single-epoch (density-dened) estimate is in yellow.

are shown in blue in [figure 6B and C](#), respectively. Estimated effective population size exhibits sustained growth until a population peak  $\sim 45$  ka BP. This is followed by a population size decline and a population bottleneck  $\sim 12$  ka BP, with a slight recovery in the recent past.

Both the BESP and BSP infer similar  $N_e\tau$  dynamics, with largely overlapping HPD intervals. However, the BESP shows a more rapid and less smooth decline. The BESP recovers a period of stable effective population size around  $\approx 20$  ka BP that coincides with the low number of sequences sampled during the LGM. HPD intervals are not notably narrower under the BESP model, likely because phylogenetic uncertainty in this data set masks any dramatic gains in precision from using the sampling date information.

Estimates of  $\beta$ , vary substantially, increasing over four orders of magnitude as time moves from the oldest sample (55 ka BP) toward the present. This contrasts with the limited variation in  $\beta$  that was observed in the IAV data set ([fig. 5C](#)). Thus, this data set demonstrates how the BESP can be used to detect a strong temporal trend in sampling intensity that requires further exploration.

It is likely that this remarkable increase in sampling intensity is caused by a combination of two factors: 1) sample preservation and successful ancient DNA recovery increases toward the present and 2) bison effective population sizes were substantially larger in the past, hence the likelihood of sampling “per-capita” in the past was smaller. There are two notable discontinuous increases in estimated  $\beta$ , one at the present (0–450 BP), and one as  $N_e\tau$  declines sharply  $\sim 15$  ka BP. The first is due to the 17 modern sequences in the data set. The second increase coincides with the period of substantial human settlement of the Americas.

We also investigated a simpler BESP with a single epoch (i.e., density-defined sampling with constant  $\beta$  across time). Comparison of the single- and 12-epoch models highlights the significant rise in  $\beta$  through time in the latter and demonstrates that multiple sampling epochs are needed to properly characterize this data set ([fig. 6C](#) and [supplementary fig. S9, Supplementary Material](#) online). The single-epoch model generates  $N_e\tau$  estimates that are unrealistically high between 15 ka BP and the present, and implies rapid exponential growth in the bison population after the LGM

(supplementary fig. S9B, Supplementary Material online). This result is an artifact of misspecification of the sampling model: enforcing a constant  $\beta$  means that sampling effort in the recent past is greatly underestimated, whereas sampling effort in the distant past is correspondingly overestimated. As a consequence, the  $N_e\tau$  estimates are biased upward (downward) during periods when the sampling intensity is underestimated (overestimated). We conclude that a BSP with a constant  $\beta$  is inadequate for this data set and would promote misleading inferences.

### The Information in Sample Timing

We now provide some theoretical basis for why the ESP improves upon the estimates of standard skyline approaches. Although sample times are known to provide additional information for demographic inference (Volz and Frost 2014), their exact contribution has not been quantified. We apply the Fisher information approach from Parag and Pybus (2019) to investigate the benefits of integrating sampling and coalescent events. As in New Approaches, we consider the subtree of  $\mathcal{T}$  that spans the  $j$ th population size,  $N_j$ , and contains  $s$  sampling and  $c$  coalescent events (see fig. 1). We use the Fisher information because it delimits the maximum asymptotic precision attainable by any unbiased estimator of  $N_j$  (Kay 1993). This precision defines the inverse of the variance (uncertainty) around that estimator. The Fisher information is computed as the expected second derivative of the log-likelihood (see Materials and Methods).

Popular skyline-based inference methods such as the BSP (Drummond et al. 2005), the skyride (Minin et al. 2008), and the skygrid (Gill et al. 2013), are founded on the coalescent log-likelihood  $\mathcal{L}_{j,c}$ , of equation (3).

$$\mathcal{L}_{j,c} = \sum_{i=1}^k 1_{\mathcal{C}}(i) \log(\alpha_i N_j^{-1}) - \Delta_i(\alpha_i N_j^{-1}) \quad (3)$$

This considers only the  $c$  coalescent events to be informative about  $N_j$ . The log-likelihoods specific to each method can be obtained from equation (3) by simply altering its population size grouping procedure. The estimates of these approaches are the MLEs of equation (3) or some related Bayesian variant. This gives the left side of equation (4), which modifies the grouped generalized skyline plot of Strimmer and Pybus (2001) to incorporate the exact times of individual events within that group.

$$\hat{N}_{j,c} = \frac{1}{c} \sum_{i=1}^k \alpha_i \Delta_i = \frac{a}{c}, \quad \mathcal{I}_c(N_j) = cN_j^{-2}. \quad (4)$$

The Fisher information available about  $N_j$  from these various skyline-based methods is identical and given by the right side of equation (4) (Parag and Pybus 2019). The maximum precision (minimum variance), around  $\hat{N}_{j,c}$ , achievable by these approaches is therefore  $\mathcal{I}_c(N_j)^{-1}$  (Kay 1993; Parag and Pybus 2017).

Next, we define an equivalent log-likelihood for sequence sampling events in equation (5). This assumes that only the  $s$

epochal sampling times are informative and ignores the coalescent events.

$$\mathcal{L}_{j,s} = \sum_{i=1}^k 1_{\mathcal{S}}(i) \log(\beta_i N_j) - \Delta_i(\beta_i N_j). \quad (5)$$

The MLE and Fisher information for this likelihood follow in equation (6).

$$\hat{N}_{j,s} = s \left( \sum_{i=1}^k \beta_i \Delta_i \right)^{-1} = \frac{s}{b}, \quad \mathcal{I}_s(N_j) = sN_j^{-2}. \quad (6)$$

Interestingly, the per event Fisher information ( $\frac{1}{s} \mathcal{I}_s(N_j)$ ) attained by this sampling-event only model is the same as that from any standard skyline method ( $\frac{1}{c} \mathcal{I}_c(N_j)$ ). This result formalizes and quantifies the assertion in Volz and Frost (2014) that  $N(t)$  can in theory be estimated using only the sampling event times.

Having considered the two information sources separately, we now examine the ESP, which deems both the  $s$  sampling and  $c$  coalescent events to be informative. Using equation (1), we compute the Fisher information of the  $j$ th segment,  $\mathcal{I}(N_j)$  (see Materials and Methods). This results in equation (7), with  $\zeta_j = \sum_{i=1}^k 1_{\mathcal{S}}(i) \alpha_i \beta_i^{-1} \geq 0$  as a grouping factor.

$$\mathcal{I}(N_j) = (s + c)N_j^{-2} + 2\zeta_j N_j^{-4}. \quad (7)$$

Intriguingly,  $\mathcal{I}(N_j) \geq \mathcal{I}_s(N_j) + \mathcal{I}_c(N_j)$ . This means that we gain additional precision by integrating both sampling and coalescent models (the per event Fisher information  $\frac{1}{s+c} \mathcal{I}(N_j)$  has increased). This extra information comes from the counteracting proportional and inverse dependencies of the two event types. Further, any segment with equal numbers of sampling and coalescent events can now be estimated with at least twice the precision of any standard skyline approach, for the same reconstructed tree  $\mathcal{T}$ . As  $n$  sampled sequences lead to  $n-1$  coalescent events, and the total Fisher information is  $\mathcal{I}(N) = \sum_{j=1}^p \mathcal{I}(N_j)$ , then the overall asymptotic precision across  $\mathcal{T}$  is also roughly, at minimum, doubled.

Equation (7) explains the improvements in population size inference that the ESP can achieve. However, this improvement may sometimes be clouded by other sources of uncertainty, such as phylogenetic error, and disappears if the sampling times contain no information about population size (in which case, the ESP converges to a standard skyline plot). Estimation precision for a given segment depends explicitly on the number of events informing that estimate, that is,  $c$  for standard skylines (eq. 3),  $s$  for sampling-events only (eq. 5), and  $s + c$  for the ESP (eq. 1). This suggests that estimates of  $N_j$  should be disregarded when the number of events falling in the  $j$ th segment is small (if this number is 0 the skyline is unidentifiable as the Fisher information matrix becomes singular; Rothenberg 1971). We recommend identifying and excluding such regions from population size estimates as a precaution against overconfident inference.

The log-likelihood of equation (1) also provides insight into the statistical power available to infer sampling intensities

across time (the  $\beta_i$  parameters). The MLE and Fisher information provided by  $\mathcal{T}$  about  $\beta_i$  over the duration of the  $j$ th population segment are given in [equation \(8\)](#).

$$\hat{\beta}_i = 1_{\mathcal{S}}(i)(\Delta_i N_j)^{-1}, \quad \mathcal{I}(\beta_i) = 1_{\mathcal{S}}(i)\beta_i^{-2}. \quad (8)$$

This MLE depends on  $N_j$ , and thus, the two parameters must be jointly estimated (see Materials and Methods for the algorithms that we used to solve this). The Fisher information shows that only intervals ending with sampling events offer the power to estimate a sampling intensity parameter. In our implementation, we group the  $\beta_i$  into a smaller number of epochs, so that the power for estimating the sampling intensity during an epoch depends on the total number of sampling events within that epoch. As, by definition, each epoch contains at least one sampling event, statistical identifiability is guaranteed ([Parag and Pybus 2019](#)). As with population size (discussed above), we recommend ignoring inferences from epochs that contain small numbers of sampling events.

## Discussion

The ESP and its Bayesian implementation (BESP) infer population size history from heterochronous phylogenies and longitudinally sampled genetic sequences. These methods generalize the skyline approach to include flexible yet tractable models of sequence sampling through time that can more accurately reflect and characterize real-world data collection protocols. This flexible formulation allows the ESP and BESP to serve as tools for the exploration and selection of appropriate time-varying sampling models. This is analogous to how the BSP can be used to select among suitable parametric demographic models for a given data set.

The improvement in population size inference exhibited by the ESP results from two factors. First, by incorporating sampling time information within an epochal framework, we essentially double the number of data points available for inference. As sampling and coalescent events are equally informative ([eqs. 4 and 6](#)) about population size, we also at least double our best asymptotic estimate precision.

Second, the bias of any coalescent inference method depends on the temporal distribution of its informative events. In standard skyline methods the rate of informative events is inversely dependent on population size, such that periods of large population size possess few coalescent events (resulting in long tree branches), whereas population bottlenecks feature high event densities. Such skewed distributions can promote inconsistent estimation ([Gattepaille et al. 2016](#)). By including sampling events, which cluster in a contrasting way to coalescent events, the ESP achieves more uniform distributions of informative events through time ([fig. 2](#)). This not only reduces bias but also increases temporal resolution, which in turn improves its power to detect and infer rapid population size changes, as seen in both simulated and empirical examples ([figs. 3–6](#)).

The ESP was partly inspired by the surveillance and data collection protocols often employed in infectious disease epidemiology. Our assumption that local sampling intensity within an epoch is proportional to population size reflects

situations in which sampling is based on availability or convenience, and hence often correlated with the number of infections in an epidemic ([Stack et al. 2010](#)). Our inclusion of epochs embodies the expectation that sequence collection rates will likely change discontinuously over time due to fluctuations in funding, resources, and timelines of individual research projects or patient cohorts. Our formulation also allows for external and unpredictable factors that may dramatically alter the sampling effort over an epidemic, such as “fog of war” effects ([Viboud et al. 2018](#)).

An analogous situation exists for studies that generate ancient DNA sequences from preserved biological material of different archeological and geological ages. Specimen preservation and the rate of DNA decay are not only highly dependent on sample age but also on moisture, temperature, and other conditions ([Shapiro and Hofreiter 2014](#)). Thus, although the number of specimens sampled from a given time period might be expected to vary proportionally with species abundance, the constant of proportionality is likely to shift through time. The epoch-based sampling model is sufficiently flexible to capture and extract these types of trends.

Although this flexibility is a benefit of the ESP, we find that biases can result when sampling intensities are defined too rigidly. When an epoch spans a long period of substantial variation in sampling effort,  $\beta$  is an estimate of the average sampling intensity over that epoch. If  $N$  also changes across this epoch, then parameter correlations mean that the ESP can overestimate population size in periods where the sampling intensity is underestimated, and vice versa. This effect is apparent when using the single-epoch model to analyze the Beringian steppe bison data set ([fig. 6C](#) and [supplementary fig. S9, Supplementary Material](#) online). This issue possibly underlies the reported biases in previous sampling-aware methods, which all effectively use a single epoch and are based around density-defined models ([Karcher et al. 2016](#)).

However, when multi-epoch models are used, the ESP is able to compensate for this bias and expose the vastly different ancient sampling dynamics which underlie this data set and corroborate previous investigations [Shapiro et al. \(2004\)](#). Analysis of the New York influenza epidemic ([fig. 5](#)) highlighted an opposite trend. Here, we found that the multi-epoch model offered little advantage over single-epoch formulations, hence providing evidence for a simpler, density-defined description. These results showcase how the ESP can serve as a tool for selecting among various sampling hypotheses and for avoiding model misspecification.

In spite of these benefits, our method has some known limitations. The ESP does not model spatial structure and hence assumes that samples are randomly drawn from a single well-mixed population. Parameter estimates may therefore be biased, if sampling efforts vary across geographic regions. Further, our analysis has relied on having some basic, prior knowledge of how to specify epoch change-points (e.g., knowing epidemic seasons or understanding practical constraints, as in frequency-defined sampling). If good *a priori* information is unavailable and epoch times are set arbitrarily, biases can result as we may have periods over which the ESP is too rigidly formulated (akin to the single-epoch bias). In these

cases, we recommend distributing sampling events evenly among epochs to guard against this type of misspecification.

The ESP differs from previous approaches that use parametric sampling models (Volz and Frost 2014; Karcher et al. 2016). This mirrors the distinction between skyline plot methods and coalescent estimators of parametric demographic functions (Parag and Pybus 2017). Karcher et al. (2016), for example, used a nonlinear sampling rate model of form  $e^{\gamma_0 N(t)^{\gamma_1}}$ , with  $\gamma_0$  and  $\gamma_1$  as parameters to be inferred. Although such formulations do not model the same range of sampling behaviors as the ESP, they can provide specific biological insights (e.g.,  $\gamma_1$  informs about sample clustering) if the true (unknown) sampling rate lies within their functional class. The ESP, by providing insight into what types of parametric hypotheses might be supported by a given data set, can complement these approaches.

The sampling intensity,  $\beta$ , inferred by the ESP can be used to reconstruct the absolute sampling rate,  $\beta N$ . Practically,  $\beta$  measures how quickly new sequences accumulate relative to the effective population size (i.e., “per capita”). It has units of [time<sup>-2</sup>]. As  $N$  has dimensions of [time] (measured in the units of the time-scaled genealogy) then the ESP directly infers changes in the rate of collecting samples per genealogical time unit. The separation of  $\beta$  and  $N$  is important, as it disaggregates the relative contributions of each time-varying unknown. Further, as  $\beta$  modulates a Poisson process, then over an infinitesimal period it defines a piecewise-constant sampling probability that is analogous to (but not equal to) the sampling model used in phylogenetic birth–death skyline methods (Stadler et al. 2013).

As sequence data become more prevalent, heterochronous sampling design will play an increasingly important role in phylodynamics (Ho and Shapiro 2011; Parag and Pybus 2019). Continuing improvements in infectious disease monitoring and sequencing will result in richer and more diverse epidemiological data (Baele et al. 2017), whereas ongoing advances in techniques for isolating and generating ancient DNA will lead to strengthened molecular evolution data sets. We hope that the ESP will prove useful in exploiting and exploring such data and help inform future debates surrounding sequence sampling protocol design and misspecification.

## Materials and Methods

### Deriving the ESP

Here, we construct the log-likelihood for the ESP (eq. 1), and derive its population size MLE (eq. 2) and Fisher information (eq. 7). Let the  $j$ th piecewise-constant segment of a sampled-coalescent process have unknown population size  $N_j$ , and duration  $t_j - t_{j-1} = \sum_{i=1}^k \Delta_i$ . We assume that this segment consists of  $k \geq 1$  event intervals, the  $i$ th of which has duration  $\Delta_i$ . If this interval ends in a sampling (coalescent) event, then  $1_S(i) = 1(0)$ , and  $1_C(i) = 0(1)$ . The coalescent lineage factors, and sampling intensities, for the  $i$ th interval are, respectively,  $\alpha_i$  and  $\beta_i$ . Figure 1 clarifies this notation for a simple reconstructed coalescent genealogy (tree),  $\mathcal{T}$ , over this segment. Standard

skyline and skygrid approaches model coalescent events as the outputs of a Poisson process with rate (over each interval)  $\sum_{i=1}^k 1_C(i)\alpha_i N_j^{-1}$ , but ignore sampling events. The ESP assumes that sampling events are also produced by a Poisson point process, with rate  $\sum_{i=1}^k 1_S(i)\beta_i N_j$ . The result is a piecewise-constant dual-type Poisson process, with combined event rate  $\lambda(t)$  as in equation (9).

$$\lambda(t) = \sum_{i=1}^k 1_S(i)\beta_i N_j + 1_C(i)\alpha_i N_j^{-1}. \quad (9)$$

Note that  $\lambda_t$  changes as time  $t$  traverses the intervals  $\Delta_i$ . We can construct the combined Poisson log-likelihood function for the  $j$ th segment,  $\mathcal{L}_j := \log P(\mathcal{T}|N_j, \{\beta_i\})$ , as in equation (10) with  $\lambda_0(t) = \sum_{i=1}^k 1_{\Delta_i}(t)(\beta_i N_j + \alpha_i N_j^{-1})$  and  $1_{\Delta_i}(t)$  indicating when  $t$  is in  $\Delta_i$  (Snyder and Miller 1991; Parag and Pybus 2018).

$$\mathcal{L}_j = - \int_{t_{j-1}}^{t_j} \lambda_0(t) dt + \int_{t_{j-1}}^{t_j} \log \lambda(t) du_t. \quad (10)$$

Here the  $\lambda_0(t)$  integral accounts for no events occurring within the intervals of the  $j$ th segment while the second term indicates the events that transpire at the interval end-points, since  $du_t = 1$  at event times and 0 otherwise.

The total log-likelihood over all  $p$  segments of  $\mathcal{T}$  is  $\mathcal{L} = \sum_{j=1}^p \mathcal{L}_j$ . For now, we only focus on the set of  $N_j$  unknowns in this log-likelihood (we discuss the power to estimate  $\{\beta_i\}$  in the next section). Equation (1) is derived by splitting the integrals in equation (10) over the  $k$  intervals. Note that,  $\mathcal{L}$  defines population size change-points at (irregular) event times. This contrasts with the approach of Karcher et al. (2016), where change-point times are regular, predefined, and do not depend on the temporal event distribution. One advantage of our formulation is that we always have at least one event informing on each  $N_j$  parameter. This results in a nonsingular Fisher information matrix, which guarantees the statistical identifiability of the ESP (Rothenberg 1971; Parag and Pybus 2019).

The skyline estimator that we propose is the grouped MLE of equation (1). This solves  $\nabla_{N_j} \mathcal{L}_j = 0$  when  $s \geq c$ , and leads to the quadratic expression in  $N_j$  given in equation (11).

$$N_j^2 - (s - c)b^{-1}N_j - ab^{-1} = 0. \quad (11)$$

Here,  $\nabla_x$  is the first partial derivative with respect to  $x$ , while  $s = \sum_{i=1}^k 1_S(i)$ , and  $c = \sum_{i=1}^k 1_C(i)$  count the total number of sampling and coalescent events falling in the  $j$ th segment of  $\mathcal{T}$ . If  $s < c$ , then  $\nabla_{N_j} \mathcal{L}_j = 0$  must be computed, and then inverted. This gives equation (12), which is a quadratic in  $N_j^{-1}$ .

$$N_j^{-2} - (c - s)a^{-1}N_j^{-1} - ba^{-1} = 0. \quad (12)$$

This conditional MLE approach is needed to avoid singularities in cases when either  $s = 0$ , or  $c = 0$ , and to keep population sizes positive. The roots of these quadratics form equation (2).

The Fisher information of the ESP with respect to  $N_j$ , is defined as  $\mathcal{I}(N_j) := -\mathbb{E}[\nabla_{N_j}^2 \mathcal{L}_j]$ , with  $\nabla_x^2$  as the second partial derivative (Kay 1993). The expectation is taken across the event intervals,  $\Delta_i$ . Applying this to equation (1), we obtain equation (13).

$$\mathcal{I}(N_j) = (s - c)N_j^{-2} + 2N_j^{-3} \sum_{i=1}^k \alpha_i \mathbb{E}[\Delta_i]. \quad (13)$$

Note that, we can replace  $\mathcal{L}_j$  with  $\mathcal{L}_{c,j}$  or  $\mathcal{L}_{s,j}$  in the above definition, to also recover equations (4) and (6), the Fisher information stemming from only the coalescent and sampling events, respectively. The expectation in equation (13) conditions on the type of event in each interval, that is,  $\mathbb{E}[\Delta_i] = \frac{1_{\mathbb{S}(i)}}{\beta_i N_j} + \frac{1_{\mathbb{C}(i)} N_j}{\alpha_i}$ . Expanding  $\sum_{i=1}^k \alpha_i \mathbb{E}[\Delta_i]$  we get  $cN_j + N_j^{-1} \sum_{i=1}^k 1_{\mathbb{S}(i)} \alpha_i \beta_i^{-1}$ . Substituting this into equation (13) simplifies to equation (7), which when  $s \approx c$  reveals a minimum  $\mathcal{I}(N_j)$  of  $2cN_j^{-2}$ . This is twice the value obtained in equation (4) and shows the marked improvement in estimate precision that results from including sampling events.

Lastly, we comment on how ESP population size estimates relate to those in equations (4) and (6). We group our skyline over the entire tree so that there is only a single population size to estimate,  $N_1$ . This is equivalent to a Kingman coalescent assumption (i.e., constant population size). As the number of coalescent and sampling events are always roughly the same then, we can use the  $s \approx c$  solution of equation (2), and the MLEs from equations (4) and (6) to derive  $\hat{N}_1 = \sqrt{\hat{N}_{1,s} \hat{N}_{1,c}}$ . If we think of the true population size,  $N(t)$ , as being continuously time-varying, then standard sky-lines estimate its harmonic mean with  $\hat{N}_{1,c}$  (Pybus et al. 2000). Similarly,  $\hat{N}_{1,s}$  estimates the arithmetic mean of  $N(t)$ . The ESP is then the geometric mean of these two mean estimators, and hence smooths the individual population size estimates from equation (4) and (6).

### Estimating the Epoch Sampling Intensities

We now define our epochal sampling model, characterize the power of the ESP for estimating sampling intensities, and present algorithms to compute the ML estimates of these sampling intensities. We assume a total of  $p'$  epochs, spanning the duration of the first (most recent) to last (most ancient) sampling event (time increases into the past). This is the period over which nonzero sampling effort is assumed. Within each epoch, the sampling intensities of each interval are the same, and epoch times coincide with sample times. This results in a piecewise-constant, time delimited, longitudinal sampling intensity. We first consider the most flexible, naïve epoch model, in which, each interval is treated as a new epoch. For the  $j$ th segment, this means there are  $k$  sampling unknowns,  $\{\beta_i\}$ . The MLE,  $\hat{\beta}_i$ , is the solution to  $\nabla_{\beta_i} \mathcal{L}_j = 0$ . The Fisher information that  $\mathcal{T}$  contains about  $\beta_i$  is  $\mathcal{I}(\beta_i) := -\mathbb{E}[\nabla_{\beta_i}^2 \mathcal{L}_j]$ .

Applying these to equation (1) gives equation (8), the MLE and Fisher information of  $\beta_i$  during the  $j$ th population segment. Two key observations emerge: (1)  $\{\hat{\beta}_i\}$  depends on  $N_j$

and (2) we only have power to estimate sampling intensities in intervals that contain sampling events ( $\mathcal{I}(\beta_i) = 0 | i \in \mathbb{C}$ ). Point (2) suggests that if  $i' \in \mathbb{S}$  and  $i' + 1 \in \mathbb{C}$  then we should assume either  $\beta_{i'+1} = 0$  or  $\beta_{i'+1} = \beta_{i'}$ , to ensure identifiability. We can resolve point (2) by grouping our sampling intensities (similar to how we group over  $N_j$ ) so that there are only  $p'$  distinct epochs. Within these epochs, there is only one sampling intensity parameter, and there is always at least one sampling event, guaranteeing identifiability (the Fisher information with respect to grouped  $\beta_i$  is nonsingular; Rothenberg 1971). The minimum variance around these per-epoch estimates of sampling intensity is then related to the sum of the  $\mathcal{I}(\beta_i)$  comprising the epoch. For example, if there is 1 epoch over the  $j$ th segment, with unknown intensity  $\beta_j$ , then  $\mathcal{I}(\beta_j) = s\beta_j^{-2}$ .

Thus, the ESP contains power to estimate (sensibly) flexible sampling intensity changes through time. Computing these estimates, and hence resolving point (1), requires joint inference of the population size and sampling intensity parameters. For ML inference, we achieve this with a simple iterative algorithm. Let  $\beta$  and  $N$  be the  $p'$  and  $p$  element vectors of unknowns that we want to estimate. We draw an initial  $\hat{\beta}(1)$  from a wide uniform distribution and then compute the conditional estimate  $\hat{N}(1) | \hat{\beta}(1)$  using equation (2). We substitute this into equation (8) to get  $\hat{\beta}(2) | \hat{N}(1)$ . Repeating this procedure iteratively yields the desired joint MLEs,  $\hat{\beta}$  and  $\hat{N}$ , usually within 100 steps (it does not require tuning and is robust to the initial  $\hat{\beta}(1)$ ). This algorithm, and the above ML solutions are all implemented in Matlab and are available at <https://github.com/kpzoo/epoch-sampling-skyline-plot> (last accessed August 7, 2019).

### The Bayesian ESP

Here, we extend the BSP (Drummond et al. 2005) to incorporate the epochal sampling model defined in the previous section. Given a genealogy  $\mathcal{T}$ , a set of  $p$  segment sizes,  $K = \{k_1, k_2, \dots, k_p\}$ , counting the numbers of events (coalescent/sampling) in each piecewise population size segment, and a set of  $p'$  epoch sizes,  $K' = \{k'_1, k'_2, \dots, k'_{p'}\}$ , counting the sampling events in each epoch, we can compute the likelihood  $f(\mathcal{T} | N, K, \beta, K')$  from equation (1). Applying Bayes' theorem yields the joint posterior distribution of  $N$ ,  $\beta$ , and  $K$  given in equation (14).

$$f(N, K, \beta | \mathcal{T}, K') \propto f(\mathcal{T} | N, K, \beta, K') \times f(N) f(K) f(\beta). \quad (14)$$

We obtain the Bayesian ESP (i.e., BEBP) by sampling from this posterior using standard MCMC proposal distributions. Equation (14) features priors on the population size vector,  $N$ , its grouping parameter (the number of events in each population size segment),  $K$ , and the sampling intensity vector,  $\beta$ . We have assumed that  $p$ ,  $p'$  and the epoch grouping parameter,  $K'$ , are all specified *a priori*, which reflects the belief that we generally have a reasonable idea of the timescale over which sampling intensities vary. This assumption could in theory be relaxed by sampling epoch sizes ( $K'$ ) within BEAST2.

We impose the same smoothing prior on  $N$  as in the BSP. This assumes that neighboring effective population size segments are autocorrelated, and implements this by drawing  $N_j$  from an exponential distribution with a mean equal to  $N_{j-1}$  (i.e.,  $N_j \sim \exp(N_{j-1}^{-1})$  for  $2 \leq j \leq p$ ) and a Jeffreys prior on  $N_1$  (Drummond et al. 2005). As we expect sampling efforts to change discontinuously, we do not assume that neighboring sampling intensities are autocorrelated, and place independent and identical priors on each  $\beta_i$ . It is trivial to relax this assumption and apply different priors to each  $\beta_i$ , for example, if *a priori* information is available about changes in sampling effort through time. This is analogous to the recent approach in Karcher et al. (2019), which embeds time-varying external covariates within the sampling process.

Our BEBP implementation also contains some practical adjustments. We constrain the minimum segment duration for both population size segments and sampling epochs to be above some threshold  $\epsilon$ . This guards against zero-length segments or epochs, which can result if too many sampling events coincide in time or if phylogenies contain bursts of branching events. Further, we constrain segments and epochs to contain at least two informative events each, to safely ensure identifiability. The BEBP is implemented as a BEAST2.6 (Bouckaert et al. 2019) package and uses  $f(\mathcal{T}|N, K, \beta, K')$  as a tree-prior for Bayesian phylogenetic analysis. This allows the BEBP, in conjunction with existing substitution and clock models, to jointly infer changes in effective population size and sampling intensity directly from sequence data, while incorporating phylogenetic uncertainty. The BEBP package is available at <https://github.com/laduplessis/besp> (last accessed November 25, 2019) and raw data, workflows, XML files, and additional figures for the simulations and empirical analyses presented above are available at [https://github.com/laduplessis/BESP\\_paper-analyses](https://github.com/laduplessis/BESP_paper-analyses) (last accessed November 22, 2019) and <https://doi.org/10.5281/zenodo.3649734>.

## Supplementary Material

Supplementary data are available at *Molecular Biology and Evolution* online.

## Acknowledgments

We thank Edward Holmes, Beth Shapiro, Cécile Viboud, and Amanda Perofsky for access to and advice on the empirical data sets. We acknowledge funding from (i) the European Research Council under the European Commission Seventh Framework Programme (FP7/2007-2013)/European Research Council (Grant Agreement 614725-PATHPHYLODYN), (ii) the Oxford Martin School, and (iii) the UK Medical Research Council (MRC) and the UK Department for International Development (DFID) under the MRC/DFID Concordat agreement, which is also part of the EDCTP2 programme supported by the European Union (Grant Reference MR/R015600/1).

## References

- Baele G, Suchard MA, Rambaut A, Lemey P. 2017. Emerging concepts of data integration in pathogen phylodynamics. *Syst Biol*. 66(1):e47–e65.
- Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015. Measurably evolving pathogens in the genomic era. *Trends Ecol Evol*. 30(6):306–313.
- Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, Heled J, Jones G, Kühnert D, De Maio N, et al. 2019. BEAST 2.5: an advanced software platform for Bayesian evolutionary analysis. *PLoS Comput Biol*. 15(4):e1006650.
- CDC. 2019. Overview of influenza surveillance in the United States [Internet]. [cited 2019 Jul 9]. Atlanta, Georgia: Centers for Disease Control and Prevention. Available from: <https://www.cdc.gov/flu/weekly/overview.htm>.
- Drummond AJ, Rambaut A, Shapiro B, Pybus OG. 2005. Bayesian coalescent inference of past population dynamics from molecular sequences. *Mol Biol Evol*. 22(5):1185–1192.
- Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. 2003. Measurably evolving populations. *Trends Ecol Evol*. 18(9):481–488.
- Faulkner JR, Magee AF, Shapiro B, Minin VN. Forthcoming. Horseshoe-based Bayesian nonparametric estimation of effective population size trajectories. *Biometrics*.
- Ferguson NM, Galvani AP, Bush RM. 2003. Ecological and immunological determinants of influenza evolution. *Nature* 422(6930):428–433.
- Gattepaille L, Torsten G, Jakobsson M. 2016. Inferring past effective population size from distributions of coalescent times. *Genetics* 204(3):1191–1206.
- Gill MS, Lemey P, Faria NR, Rambaut A, Shapiro B, Suchard MA. 2013. Improving Bayesian population dynamics inference: a coalescent-based model for multiple loci. *Mol Biol Evol*. 30(3):713–724.
- Hall MD, Woolhouse ME, Rambaut A. 2016. The effects of sampling strategy on the quality of reconstruction of viral population dynamics using Bayesian skyline family coalescent methods: a simulation study. *Virus Evol*. 2(1):vew003.
- Ho SYW, Shapiro B. 2011. Skyline-plot methods for estimating demographic history from nucleotide sequences. *Mol Ecol Res*. 11(3):423–434.
- Karcher MD, Palacios JA, Bedford T, Suchard MA, Minin VN. 2016. Quantifying and mitigating the effect of preferential sampling on phylodynamic inference. *PLoS Comput Biol*. 12(3):e1004789.
- Karcher MD, Suchard MA, Dudas G, Minin VN. 2019. Estimating effective population size changes from preferentially sampled genetic sequences. *arXiv e-Prints, Page arXiv*. 1903.11797.
- Karcher MD, Palacios JA, Lan S, Minin VN. 2017. PHYLODYN: an R package for phylodynamic simulation and inference. *Mol Ecol Resour*. 17(1):96–100.
- Kay SM. 1993. Fundamentals of statistical signal processing: estimation theory. New Jersey: Prentice Hall.
- Kingman JFC. 1982. On the genealogy of large populations. *J Appl Probab*. 19(A):27–43.
- Minin VN, Bloomquist EW, Suchard MA. 2008. Smooth skyride through a rough skyline: Bayesian coalescent-based inference of population dynamics. *Mol Biol Evol*. 25(7):1459–1471.
- Parag KV, Pybus OG. 2017. Optimal point process filtering and estimation of the coalescent process. *J Theor Biol*. 421:153–167.
- Parag KV, Pybus OG. 2018. Exact Bayesian inference for phylogenetic birth-death models. *Bioinformatics* 34(21):3638–3645.
- Parag KV, Pybus OG. 2019. Robust design for coalescent model inference. *Syst Biol*. 68(5):730–743.
- Pybus OG, Rambaut A. 2009. Evolutionary analysis of the dynamics of viral infectious disease. *Nat Rev Genet*. 10:240–250.
- Pybus OG, Rambaut A, Harvey PH. 2000. An integrated framework for the inference of viral population history from reconstructed genealogies. *Genetics* 155(3):1429–1437.
- Rambaut A, Pybus OG, Nelson MI, Viboud C, Taubenberger JK, Holmes EC. 2008. The genomic and epidemiological dynamics of human influenza A virus. *Nature* 453(7195):615–619.

- Rothenberg TJ. 1971. Identification in parametric models. *Econometrica* 39(3):577–591.
- Shapiro B, Drummond AJ, Rambaut A, Wilson MC, Matheus PE, Sher AV, Pybus OG, Gilbert MTP, Barnes I, Binladen J, et al. 2004. Rise and fall of the Beringian steppe bison. *Science* 306(5701):1561–1565.
- Shapiro B, Hofreiter M. 2014. A paleogenomic perspective on evolution and gene function: new insights from ancient DNA. *Science* 343(6169):1236573–1236573.
- Snyder DL, Miller MI. 1991. Random point processes in time and space. 2nd ed. New York: Springer-Verlag.
- Stack JC, Welch JD, Ferrari MJ, Shapiro BU, Grenfell BT. 2010. Protocols for sampling viral sequences to study epidemic dynamics. *J R Soc Interface*. 7(48):1119–1127.
- Stadler T, Kuhnert D, Bonhoeffer S, Drummond AJ. 2013. Birth-death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc Natl Acad Sci U S A*. 110(1):228–233.
- Strimmer K, Pybus OG. 2001. Exploring the demographic history of DNA sequences using the generalized skyline plot. *Mol Biol Evol*. 18(12):2298–2305.
- Viboud C, Sun K, Gaffey R, Ajelli M, Fumanelli L, Merler S, Zhang Q, Chowell G, Simonsen L, Vespignani A, et al. 2018. The RAPIDD ebola forecasting challenge: synthesis and lessons learnt. *Epidemics* 22:13–21.
- Volz EM, Frost SDW. 2014. Sampling through time and phylodynamic inference with coalescent and birth–death models. *J R Soc Interface*. 11:20140945.
- WHO. 2018. Fact sheet on seasonal influenza [Internet]. [cited 2019 Jul 25]. Geneva, Switzerland: World Health Organization. Available from: [https://www.who.int/en/news-room/fact-sheets/detail/influenza-\(seasonal\)](https://www.who.int/en/news-room/fact-sheets/detail/influenza-(seasonal)).