

## Modeling the Growth and Decline of Pathogen Effective Population Size Provides Insight into Epidemic Dynamics and Drivers of Antimicrobial Resistance

ERIK M. VOLZ\* AND XAVIER DIDELOT

Department of Infectious Disease Epidemiology, Imperial College London, Norfolk Place, W2 1PG, UK

\*Correspondence to be sent to: Department of Infectious Disease Epidemiology, Imperial College London, Norfolk Place W2 1PG, UK; Email: [e.volz@imperial.ac.uk](mailto:e.volz@imperial.ac.uk).

Received 19 October 2017; reviews returned 1 February 2018; accepted 4 February 2018

Associate Editor: Jeffrey Townsend

**Abstract.**—Nonparametric population genetic modeling provides a simple and flexible approach for studying demographic history and epidemic dynamics using pathogen sequence data. Existing Bayesian approaches are premised on stochastic processes with stationary increments which may provide an unrealistic prior for epidemic histories which feature extended period of exponential growth or decline. We show that nonparametric models defined in terms of the growth rate of the effective population size can provide a more realistic prior for epidemic history. We propose a nonparametric autoregressive model on the growth rate as a prior for effective population size, which corresponds to the dynamics expected under many epidemic situations. We demonstrate the use of this model within a Bayesian phylodynamic inference framework. Our method correctly reconstructs trends of epidemic growth and decline from pathogen genealogies even when genealogical data are sparse and conventional skyline estimators erroneously predict stable population size. We also propose a regression approach for relating growth rates of pathogen effective population size and time-varying variables that may impact the replicative fitness of a pathogen. The model is applied to real data from rabies virus and *Staphylococcus aureus* epidemics. We find a close correspondence between the estimated growth rates of a lineage of methicillin-resistant *S. aureus* and population-level prescription rates of  $\beta$ -lactam antibiotics. The new models are implemented in an open source R package called *skygrowth* which is available at <https://github.com/mrc-ide/skygrowth>. [Antimicrobial resistance; effective population size; growth rate; MRSA; phylodynamics; *skygrowth*.]

Nonparametric population genetic modeling has emerged as a simple, flexible, popular, and powerful tool for interrogating genetic sequence data to reveal demographic history (Ho and Shapiro 2011). This approach has proved especially useful for analysis of pathogen sequence data to reconstruct epidemic history, and such models are increasingly incorporated into surveillance systems for infectious diseases (Volz et al. 2013). The most commonly used techniques are derivatives of the original *skyline* coalescent model, which describes the evolution of effective population size as a piecewise constant function of time (Pybus et al. 2000). The basic *skyline* model is prone to overfitting and estimating drastic fluctuations in effective population size, so that numerous approaches were subsequently developed for smoothing population size trajectories. Initial approaches to smoothing *skyline* estimators were based on aggregating adjacent coalescent intervals within a maximum likelihood framework (Strimmer and Pybus 2001). Subsequent development has largely focused on Bayesian approaches where a more complex stochastic diffusion process provides a prior for the evolution of a piecewise constant function of effective population size (Drummond et al. 2005). Nonparametric Bayesian approaches are now the most popular approach for phylodynamic inference, and such approaches have illuminated the epidemic history of numerous pathogens in humans and animals (Ho and Shapiro 2011).

To date, all Bayesian nonparametric models have assumed that the effective population size (or its logarithm) follows a stochastic process such as a Brownian motion (BM) (Minin et al. 2008; Palacios and

Minin 2013). The choice of a process with stationary increments as prior can have large influence on size estimates especially when genealogical data are sparse and uninformative. Genealogies often provide very little information about effective population size near the present (or most recent sample), especially in exponentially increasing populations (de Silva et al. 2012). In such cases, *skyline* estimators with BM priors on the effective population size may produce estimates which stabilize at a constant level even when the true size is increasing or decreasing exponentially. We argue that in many situations, a more realistic prior can be defined in terms of the growth rate of the effective population size. Below, we describe such a prior based on a simple autoregressive stochastic process defined on the growth rate of effective population size. We show how this prior can lead to substantially different estimates and argue that these estimates are more accurate in many situations. When genealogical data are sparse, our model will retain the growth rate learned from other parts of the genealogy and will correctly capture trends of exponential growth or decline. Even though our approach is nonparametric, we consider its relationship with parametric models of epidemic population genetics to show that our estimates of growth rates of pathogen effective population size are often likely to correspond to growth rates of an infectious disease epidemic.

Smoothing effective population size trajectories using a prior on growth rates also have important advantages when incorporating nongenetic covariate data into phylodynamic inference (Baele et al. 2016). Recent work has focused on refining effective population size estimates using both the times of sequencing sampling

(Karcher et al. 2016) or using environmental data which are expected to correlate with size estimates, such as independent epidemic size estimates based on nongenetic data (Gill et al. 2016). Existing statistical models have assumed that the effective population size has a linear or log-linear relationship with temporal covariates. However in many cases, a more realistic model would specify that the growth rate of effective population size is correlated with covariates, as when for example an environmental variable impacts the replicative fitness of a pathogen. We provide a similar extension of previous *skyride* models with covariate data (Gill et al. 2016) to show how such data can be used to test hypotheses concerning their effect and, when a significant effect exists, to refine estimates of both the growth rates and the effective population sizes.

We illustrate the potential advantages of our growth rate model using a rabies virus data set that has been thoroughly studied using previous phylodynamic methods (Biek et al. 2007; Gill et al. 2016). In particular, we show how our model correctly estimates a recent decline in epidemic size whereas previous models mistakenly predict a stabilization of the epidemic prevalence. We also apply our methodology to a genomic data set of methicillin-resistant *S. aureus* (MRSA) that had not formally been analyzed using phylodynamic methods (Uhlemann et al. 2014). We show how time series on prescription rates of  $\beta$ -lactam antibiotics correlate strongly with growth and decline of the effective population size, revealing the impact of antibiotic use on the emergence and spread of resistant bacterial pathogens.

## METHODS AND MATERIALS

We model effective population size through time as a first order autoregressive stochastic process on the growth rate. This provides an intuitive link between the growth rate of effective population size of pathogens and epidemic size as well as the reproduction number of the epidemic. We further show how to incorporate time-varying environmental covariates into phylodynamic inference.

### Previous Bayesian Nonparametric Phylodynamic Models

Several nonparametric phylodynamic models have been proposed based on BM processes and the Kingman coalescent genealogical model (Kingman 1982). In particular, the Bayesian nonparametric *skyride* model uses a BM prior to smooth trajectories of the logarithm of the effective population size (Minin et al. 2008). Let  $\gamma(t) = \log(\text{Ne}(t))$  denote the logarithm of the effective population size as a function of time. The BM prior is defined as:

$$\gamma(t+dt) \sim \gamma(t) + \mathcal{N}(0, dt/\tau), \quad (1)$$

where  $\tau$  is an estimated precision parameter, for which an uninformative Gamma prior is typically used.

This BM prior has been adapted and applied in a variety of ways to enable statistical inference. In the *skygrid* model (Gill et al. 2013), time is discretized, and  $\gamma$  is defined to be a piecewise constant function of time over a grid with time increments  $h$ , and the value  $\gamma_i$  is estimated for each interval  $i$ . Time intervals do not in general correspond to coalescent times in the genealogy. In this case, the BM prior is computed over increments of  $\gamma$ :

$$p(\gamma_{1:m}|\tau) \propto \prod_{i=1}^{m-1} p(\gamma_{i+1} - \gamma_i|\tau), \quad (2)$$

where

$$p(\gamma_{i+1} - \gamma_i|\tau) = \sqrt{\frac{\tau}{2\pi h}} e^{-\frac{\tau}{2h}(\gamma_{i+1} - \gamma_i)^2}.$$

The genealogical data take the form  $\mathcal{G} = (c_{1:(n-1)}, s_{1:n})$ , where  $c$  and  $s$  are respectively ordered coalescent times (internal nodes of the genealogy) and sampling times (terminal nodes of the genealogy). In the coalescent framework, the sampling times are usually considered to be fixed, so that  $p(s) = 1$  and  $p(\mathcal{G}) = p(c|s)$ . Alternatively, in some variations of this model, a prior  $p(s|\text{Ne})$  is also provided for the sequence of sampling times, making this approach similar to but more flexible than sampling-birth-death-models (Volz and Frost 2014; Karcher et al. 2016).

Given a genealogy, the posterior distribution of the parameters  $\tau$  and  $\gamma_{1:m}$  is decomposed as:

$$p(\gamma_{1:m}, \tau|\mathcal{G}) \propto p(\mathcal{G}|\gamma_{1:m})p(\gamma_{1:m}|\tau)p(\tau). \quad (3)$$

The second term is given by Equation 2 and the last term by the prior on  $\tau$ . To assist with the definition of the first term, we first denote  $A(t)$  to be the number of extant lineages at time  $t$ :

$$A(t) = \sum_{i=1}^n I(s_i > t) - \sum_{i=1}^{n-1} I(c_i > t), \quad (4)$$

where  $I(x)$  is an indicator function equal to one when  $x$  is true and equal to zero otherwise. The probability density of the genealogical data given the population size history  $\gamma_{1:m}$  is then equal to (Griffiths and Tavaré 1994):

$$p(\mathcal{G}|\gamma_{1:m}) = \prod_{i=1}^{2n-2} \left( I(t_i \in c_i) \frac{\binom{A(t_i)}{2}}{\text{Ne}(t_{i+1})} e^{-\int_{t_i}^{t_{i+1}} -\frac{A(t)}{2\text{Ne}(t)} dt} + (1 - I(t_i \in c_i)) e^{-\int_{t_i}^{t_{i+1}} -\frac{A(t)}{2\text{Ne}(t)} dt} \right), \quad (5)$$

where  $t_{1:(2n-1)} = c_{1:(n-1)} \cup s_{1:n}$  is the set union of sample and coalescent times in descending order.

*Relationship Between the Growth Rate of Effective Population Size and Epidemic Properties*

Several recent studies have investigated the relationship between the effective population size of a pathogen and the number of infected hosts (Rosenberg and Nordborg 2002; Koelle et al. 2011; Dearlove and Wilson 2013). A simple link between these quantities does not exist, since the relationship depends on how incidence and epidemic size change through time (Volz et al. 2009), population structure (Volz 2012), and complex evolution of the pathogen within hosts (Didelot et al. 2016; Volz et al. 2017). Under idealized situations, there is however a simple relationship between the growth rate of effective population size and the growth rate of an epidemic (Frost and Volz 2010; Volz et al. 2013).

Let  $Y(t)$  and  $\beta(t)$  denote the number of infected hosts and per-capita transmission rate, respectively, as functions of time. Note that  $\beta(t)$  may depend on the density of susceptible individuals in the population, as in the common susceptible-infected-removed (SIR) model, in which case  $\beta(t) \propto S(t)/N$  (Allen 2008). The coalescent rate for an infectious disease epidemic was previously derived under the assumption that within-host effective population size is negligible and that superinfection does not occur (Volz et al. 2009; Frost and Volz 2010):

$$\lambda(t) = \binom{A(t)}{2} \frac{2\beta(t)}{Y(t)}. \tag{6}$$

Equating this rate with the coalescent rate under the coalescent model  $\lambda(t) = \binom{A(t)}{2} / Ne(t)$  (Kingman 1982) yields the following formula for the effective population size:

$$Ne(t) = \frac{Y(t)}{2\beta(t)}. \tag{7}$$

Differentiating with respect to time (denoting with a dot superscript) yields:

$$\dot{Ne}(t) = \frac{\dot{Y}(t)}{2\beta(t)} - \frac{\dot{\beta}(t)Y(t)}{2(\beta(t))^2}. \tag{8}$$

Note that, in general the growth rate of the effective population size does not correspond to the growth rate of  $Y$ , however if the per-capita transmission rate is constant ( $\dot{\beta} = 0$ ), we have  $\dot{Ne} = \dot{Y} / (2\beta) \propto \dot{Y}$ . Thus, we expect that over phases of the epidemic where per-capita transmission rates are nearly constant there will be close correspondence between the growth or decline of the effective population size and the growth or decline of the unobserved number of infected hosts. This condition is often satisfied near the beginning of an outbreak which has an exponential phase. It is also often satisfied towards the end of epidemics when the epidemic size is decreasing at a constant exponential rate.

The basic reproduction number  $R_0$  describes the expected number of transmission events caused by a single infected individual in an otherwise susceptible population. By extension, we can define  $R(t)$  as the

expected number of transmissions by an infected host infected at time  $t$  (Fraser 2007). Assuming that all infected individuals are equally infectious (as is the case e.g., in the SIR model), we have that during periods when the epidemic growth rate is constant, each infected individual transmits at rate  $\beta(t) = R(t) / \psi$ , where  $\psi$  is the mean duration of infections. With these definitions, the number of infections  $Y(t)$  varies according to the following differential equation:

$$\dot{Y}(t) = Y(t) \frac{R(t) - 1}{\psi} \tag{9}$$

Combining Equations 7, 8, and 9 leads to the following approximate estimator for the reproduction number through time:

$$\hat{R}(t) = 1 + \psi \frac{\dot{Ne}(t)}{Ne(t)} \tag{10}$$

This estimator makes use of the quantity  $\dot{Ne}(t) / Ne(t)$  which will be estimated in our model below. Equation 10 is likely to be a good estimator over periods of the epidemic where per-capita transmission rates are invariant. A special case of this occurs at the start of an epidemic, in which case Equation 10 can be used to estimate the basic reproduction number  $R_0$ , as previously noted (Pybus 2001).

*A Growth Rate Prior for Effective Population Size*

We propose a model in which the growth rate of the effective population size, as opposed to effective population size itself, is an autoregressive process with stationary increments. This growth rate is defined as:

$$\rho(t) = \frac{\dot{Ne}(t)}{Ne(t)}. \tag{11}$$

Note that  $\rho(t)$  is a real-valued quantity, with negative and positive values respectively indicating an increase and decrease in the effective population size. In particular, if the population is exponentially growing or declining from  $t=0$ , then we would have  $Ne(t) = Ne(0)\exp(\rho t)$  so that  $\rho(t) = \rho$  at every time  $t \geq 0$ . More generally, we model  $\rho(t)$  using a BM process:  $\rho(t) \sim BM(\tau)$  (cf Equation 1). To facilitate statistical inference, we work with a discretized time axis with  $m$  intervals of length  $h$  as in the *skygrid* model (Gill et al. 2013). We define the growth rate in time interval  $i$  as:

$$\rho_i = \frac{Ne_{i+1} - Ne_i}{hNe_i}. \tag{12}$$

We use the following approximate model for  $p(\rho_{i+1} | \rho_i)$ :

$$\rho_{i+1} \sim \rho_i + \mathcal{N}(0, h/\tau). \tag{13}$$

Note that Equation 12 implies  $\rho_i \in (-1/h, \infty)$  since  $Ne$  cannot decline below zero, whereas the approximate model in Equation 13 assumes support on the entire real line. We have found performance with this approximate

model to be superior to exact models on the log transformation of  $N_e$  provided that  $h$  is small.

With the above definitions, the prior density of a sequence  $\rho_{1:m}$  is defined in terms of the increments:

$$p(\rho_{1:m}|\tau) \propto \prod_{i=1}^{m-2} p(\rho_{i+1} - \rho_i|\tau), \quad (14)$$

where

$$p(\rho_{i+1} - \rho_i|\tau) = \sqrt{\frac{\tau}{2\pi h}} e^{-\frac{\tau}{2h}(\rho_{i+1} - \rho_i)^2}.$$

This equation can be compared with the *skygrid* density, Equation 2.

### Incorporating Covariates into Phylodynamic Inference

A simple model was recently proposed for incorporating time-varying covariates into phylodynamic inference with *skygrid* models (Gill et al. 2016). Suppose we observe  $q$  covariates at  $m$  time points denoted  $X = (X_{1:m,1:q})$ , and such that observation times correspond to the grid used in the phylodynamic model. The following linear model for the marginal distribution of  $\gamma$  with covariate vector  $\alpha_{1:q}$  was proposed:

$$p(\gamma_i|X, \alpha_{1:q}, \epsilon) \sim \mathcal{N}(\alpha_0 + X_{i,1:q}\alpha_{1:q}, \epsilon), \quad (15)$$

where  $\alpha_0$  is the expected mean of  $\gamma$  without covariate effects.

This implies, along with the BM model, the following marginal distribution of the increments:

$$p(\gamma_{i+1} - \gamma_i|X, \alpha_{1:q}, \tau, \epsilon) \sim \mathcal{N}(X_{i+1,1:q}\alpha_{1:q} - X_{i,1:q}\alpha_{1:q}, h/\tau + 2\epsilon). \quad (16)$$

When covariates are likely to be associated with growth rates of the effective population size instead of the logarithm of the effective population size, we can analogously define the density of increments of  $\rho$ :

$$p(\rho_{i+1} - \rho_i|X, \alpha_{1:q}, \tau, \epsilon) \sim \mathcal{N}(X_{i+1,1:q}\alpha_{1:q} - X_{i,1:q}\alpha_{1:q}, h/\tau + 2\epsilon). \quad (17)$$

When fitting this model, we drop  $\epsilon$  for simplicity (as in Gill et al. 2016), and estimate a single variance parameter  $\tau$  along with the regression coefficients  $\alpha$ .

### Inference and Software Implementation

Our growth rate model is implemented in an open source R package called *skygrowth*, available from <https://mrc-ide.github.io/skygrowth/>, and which includes both maximum *a posteriori* (MAP) and Bayesian Markov Chain Monte Carlo (MCMC) methods for model fitting.

The MCMC procedure uses a Gibbs-within-Metropolis algorithm that alternates between sampling the growth rate vector  $\rho_{1:m}$  and sampling of the

precision parameter  $\tau$ . Metropolis-Hastings sampling is also performed for regression coefficients  $\alpha_{1:q}$  if covariate data are provided with univariate normal proposals. The elements of  $\rho_{1:m}$  are sampled in sequence (from past to present), and multiple Gibbs iterations (by default 100) are performed before updating other parameters using Metropolis-Hastings steps.

MAP is used as a starting point for the MCMC. The MAP estimator alternates between optimization of  $\gamma_{1:m}$  using gradient descent (*BFGS* in R, Goldfarb 1970) and univariate optimization of  $\tau$  until convergence in the posterior is observed. Approximate credible intervals are provided for the MAP estimator based on curvature of the posterior around the optimum.

## RESULTS

### Simulations

We evaluated the ability of the *skygrowth* model to infer epidemic trends by simulating partially-sampled genealogies from a stochastic individual-based susceptible-infected-removed (SIR) model. Simulated data were generated using the BEAST2 package MASTER (Vaughan and Drummond 2013), and code to reproduce simulated results is available at <https://github.com/emvolz/skygrowth-experiments>.

The *skygrowth* model was also compared to *skygrid* as implemented in the *phylodyn* R package (Karcher et al. 2016, 2017) which estimates effective population size using a fast approximate Bayesian nonparametric reconstruction (BNPR). The SIR model was density dependent with a reaction rate  $\beta S(t)I(t)$  of generating new infections. Figure 1 shows results of a single simulation with  $R_0=1.3$  and 10,000 initial susceptible individuals. Additional simulations are shown in Supplementary Fig. S1 available on Dryad at <http://dx.doi.org/10.5061/dryad.9qh7t9t>. Estimates with *skygrowth* were obtained using the MCMC algorithm and an Exponential(0.1) prior on the precision parameter. We report the posterior means from both *skygrowth* and *skygrid* BNPR. Genealogies were reconstructed by sampling 200 or 1000 infected individuals at random from the entire history of the epidemic. In this scenario, both the *skygrowth* and *skygrid* models reproduce the true epidemic trend, capturing both the rate of initial exponential increase, the time of peak prevalence, and the rate of epidemic decline. However, when sampling only 200 lineages (Fig. 1B), the genealogy contains relatively little information about later epidemic dynamics, and the *skygrid* estimates an unrealistic levelling-off of  $N_e$ . Estimates using the *skygrid* BNPR model were highly similar to results using an exact MCMC algorithm for sampling the posterior also included in the *phylodyn* package.

While the results in Figure 1A and B suggest that  $N_e(t)$  can serve as a very effective proxy for epidemic size, the degree of correspondence will depend on details of the epidemic model as discussed in the Methods section.

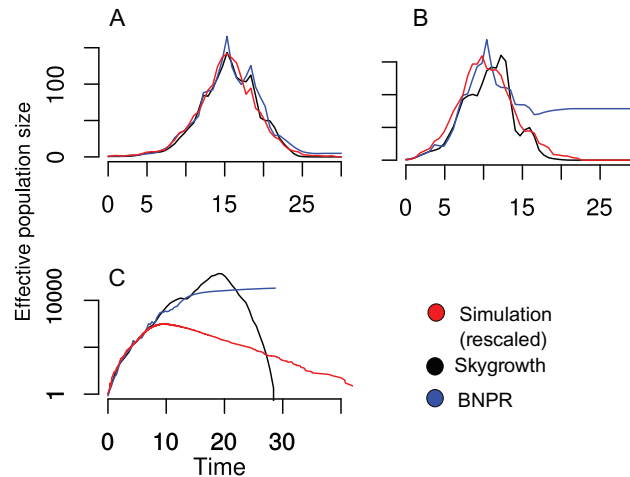


FIGURE 1. Comparison of effective population size estimates using the *skygrowth* and *skygrid* models applied to data from a susceptible-infected-recovered simulated epidemic. Effective population size estimates are also compared to the number of infected hosts through time under a linear rescaling (red). a) Estimates using a SIR model and simulated genealogy with 1000 sampled lineages and  $R_0=1.3$ . b) Estimates using a SIR model and simulated genealogy with 200 sampled lineages and  $R_0=1.3$ . c) Estimates using a SIR model and simulated genealogy with 200 sampled lineages and  $R_0=5$ .

Figure 1C and Supplementary Fig. S2 available on Dryad show a scenario where estimates of  $N_e(t)$  capture the initial rate of exponential growth but fail to estimate the time of peak epidemic prevalence, and the *skygrid* model also fails to detect that the epidemic ever decreases. This scenario was based on a higher  $R_0=5$  and only 2000 initially susceptible individuals, such that almost all hosts are eventually infected and the rate of epidemic decline predominantly reflects the host recovery rate. This is easily understood using the formula  $N_e(t) \propto I(t)/S(t)$  (cf. Equation 7). When  $R_0$  is large,  $S(t)$  will change drastically over the course of the epidemic. In the later stages, almost all hosts have been infected so that  $1/S(t)$  is large, producing correspondingly large effective population sizes. There is a very slight signal of decreasing growth rate which is detected shortly following epidemic peak using the *skygrowth* model. In the absence of other information, the *skygrowth* model retains this growth rate which produces estimates of decreasing  $N_e(t)$ .

#### Rabies Virus

An epidemic of rabies broke out in the late 1970s in the North American raccoon population, following the emergence of a host-adapted variant of the virus called raccoon rabies virus (RRV). By the end of the 1990s, this outbreak had spread to a vast geographical area including all Northeast and mid-Atlantic US states (Childs et al. 2000). A sample of 47 RRV isolates has been sequenced in a previous study (Biek et al. 2007), and BEAST (Drummond et al. 2012) was used to reconstruct a dated phylogenetic tree. A standard *skyline* analysis (Drummond et al. 2005) was performed, which visually suggested a correlation between the inferred effective population size ( $N_e$ ) and the monthly area newly affected by RRV (hereafter denoted  $V$ ), but without attempting to quantify the strength or significance of this association.

These data were recently reanalyzed using the *skygrid* model with covariates (Gill et al. 2016). No significant association was found between  $N_e$  and  $V$ , but the authors noted that since  $V$  is the newly affected area,  $V$  would be expected to be associated with a change in  $N_e$  rather than  $N_e$  itself. Since the *skyride* method is focused on  $N_e$ , like all previous phylodynamic methods, the authors considered the cumulative distribution of  $V$  and showed that this is slightly associated with  $N_e$  [with a 95% credible interval of (0.18–2.86) on the covariate effect size, Gill et al. 2016]. However, this approach is not fully satisfactory. In particular, since  $V$  is always positive, the cumulative distribution of  $V$  is always increasing, whereas  $N_e$  is in principle equally likely to increase or decrease over time. Furthermore both  $V$  and its cumulative distribution were considered on a logarithm scale, so that the latter flattens over time by definition.

A more natural solution is to keep the covariate  $V$  untransformed and investigate its association with the growth rate  $\rho(t)$  rather than  $N_e(t)$  as implemented in our methodology (Fig. 2). For this analysis, we used exactly the same dated phylogeny as previously published (Biek et al. 2007) (reproduced in Supplementary Fig. S3 available on Dryad). When the covariate was not used (red results in Fig. 2), the growth rate was inferred to be positive but declining progressively to zero from 1973 to ~1983, then stable around zero up to ~1990, followed by a period of positive growth until ~2000, after which the growth rate decreased below zero. This implies that the effective population size increased from 1973 to ~1983, then was stable until ~1990, increased to a peak in ~1997 and afterwards decreased. Two waves of spread have therefore been inferred as in previous analyses (Biek et al. 2007; Gill et al. 2016), with the first one starting in the 1970s and ending in ~1983 and the second one lasting from ~1990 to ~1997.

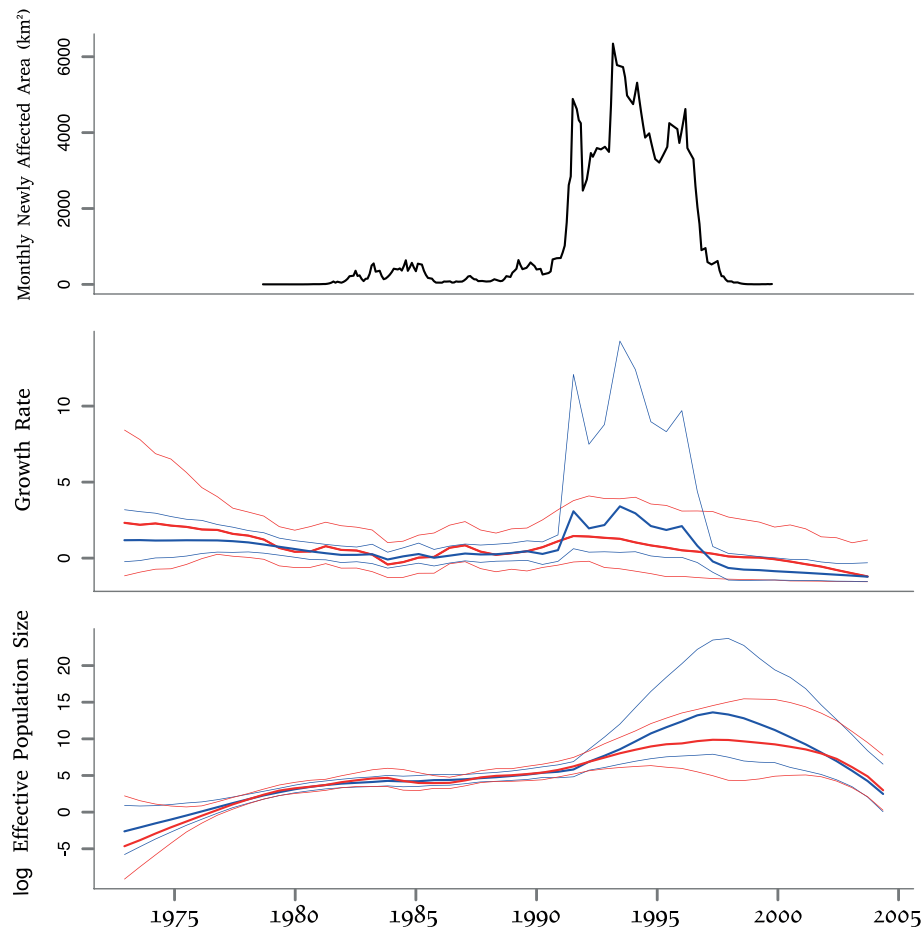


FIGURE 2. Results on the rabies application. Top: covariate data, representing the area in  $\text{km}^2$  newly affected by rabies recorded monthly between September 1978 and October 1999. Middle: growth rate estimates. Bottom: log effective population size estimates. The middle and bottom plots show results without (red) and with (blue) the use of the covariate data, and with a solid line indicating posterior means and shaded areas indicating the 95% credible regions.

Unfortunately the covariate data  $V$  start in September 1978 and therefore do not cover the first wave. However, the covariate data show that the epidemic was spreading very quickly between 1992 and 1997, much faster than before or after these dates, and this timing corresponds fairly precisely to the second wave of spread. When the covariate data were integrated into phylodynamic inference, the covariate effect size was found to be statistically significant but only slightly so, with a large 95% credible interval for the covariate effect size of (0.03–4.61) and posterior mean of 1.09. The reconstructed growth rate and effective population size when using the covariate data (blue results in Fig. 2) were compatible with results without covariate data. Using additional informative data tighten the credible interval as would be expected, except in the second wave during which the covariate data suggest higher values for both the growth rate and effective population size. The mean posterior growth rate reached a value of about 2.5 per year in the 1990s (Fig. 2) and the average generation time of raccoon rabies has previously been estimated to be around 2 months (Biek et al. 2007). We can use Equation 10 to infer a reproduction number of  $R = 1.4$ , slightly higher than a

previous estimate around  $R = 1.1$  based on the same data (Biek et al. 2007).

One of the main novel findings of our analysis is that we found a significant decline of the effective population size of raccoon rabies post-2000, whereas previous phylodynamic studies based on the same data found this to be constant (Biek et al. 2007; Gill et al. 2016). Previous methods consider a BM on the logarithm of  $N_e$ , which results in a strong prior that  $N_e$  is constant in recent time. In contrast, our model results in the growth rate being *a priori* constant, so that the clear decline in growth rate started in the mid-1990s is likely to have continued to the point that the growth rate became negative and  $N_e$  declined. Our result is in good agreement with Centers for Disease Control surveillance that shows a clear decline in rabid raccoons after the peak in the mid-1990s (Monroe et al. 2016).

#### *Staphylococcus aureus* USA300

*Staphylococcus aureus* is a bacterium that causes infections ranging from mild skin infections to life-threatening septicemia. In the 1980s and 1990s,

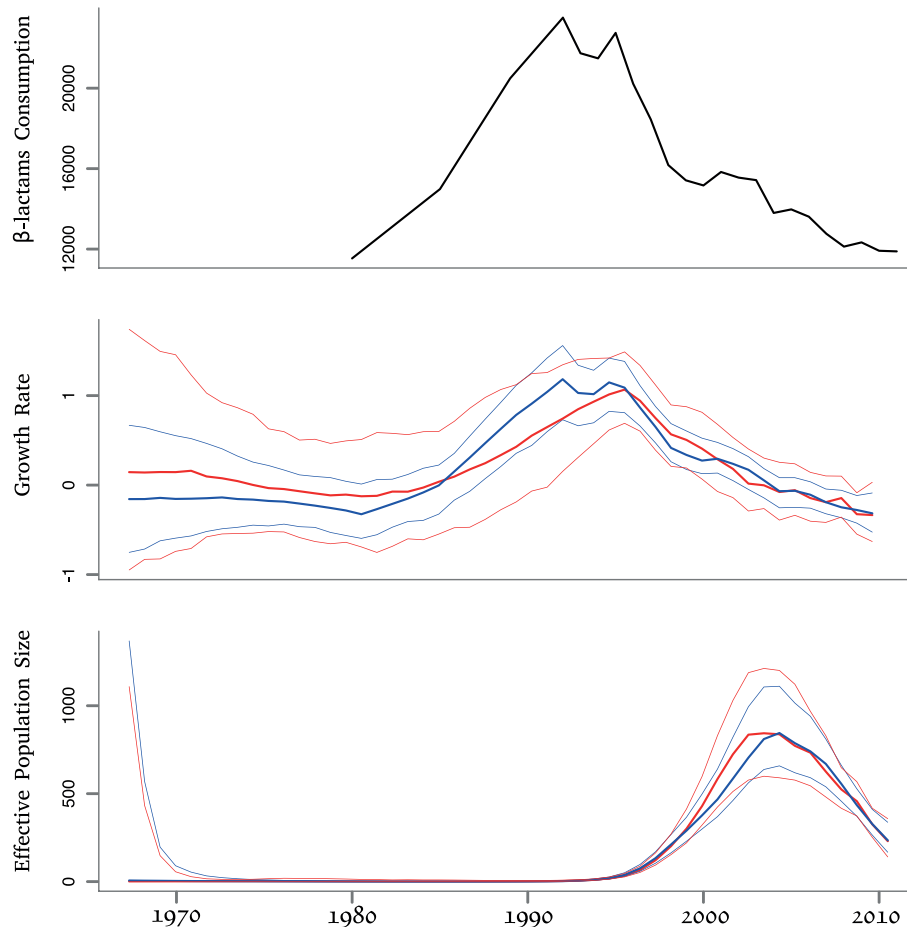


FIGURE 3. Results on the USA300 application. Top: covariate data, representing the consumption of  $\beta$ -lactams between 1980 to 2012 in the USA, measured in standard units per 1000 population. Middle: growth rate estimates. Bottom: effective population size estimates. The middle and bottom plots show results without (red) and with (blue) the use of the covariate data, and with a solid line indicating posterior means and shaded areas indicating the 95% credible regions.

several variants of *S. aureus* have emerged that are resistant to methicilin and other  $\beta$ -lactam antibiotics, and collectively called MRSA (Chambers and Deleo 2009). MRSA are well known as a leading cause of hospital infections worldwide, but the MRSA variant called USA300 differs from most others by causing infections mostly in communities rather than hospitals. USA300 was first reported in 2000 and has since spread throughout the USA and internationally (Tenover and Goering 2009; Challagundla et al. 2018). A recent study sequenced the genomes from 387 isolates of USA300 sampled from New York between 2009 and 2011, and reconstructed phylogeographic spread that frequently involved transmission within households (Uhlemann et al. 2014).

The USA300 phylogenetic tree (Uhlemann et al. 2014) was dated using a previously described method (Didelot et al. 2012) and a clock rate of  $\sim 3$  substitutions per year for USA300 (Uhlemann et al. 2014; Alam et al. 2015). We analyzed the resulting dated phylogeny (Supplementary Fig. S4 available on Dryad) using our phylodynamic methodology (Fig. 3). We initially performed this analysis without the use of any covariate

data (red results in Fig. 3) and found that the growth rate had been around zero up until 1985, after which it steadily increased until  $\sim 1995$ , and subsequently decreased almost linearly, becoming negative in  $\sim 2002$  and continuing to decrease afterwards. The effective population size was accordingly found to have been very small until the mid-1990s, to have peaked in  $\sim 2002$  and to have declined since. These results are in very good agreement with a phylodynamic analysis of USA300 performed using a traditional *skyline* plot on a different genomic data set (Glaser et al. 2016), and epidemiological data also suggest that USA300 may be declining (Planet 2017). The overall MRSA incidence has declined by 31.2% in the USA between 2005 and 2011 (Dantes et al. 2013), with some MRSA lineages showing encouraging signs of reverting to methicilin susceptibility (Ledda et al. 2017).

We hypothesized that the dynamics of USA300 may be driven by the consumption of  $\beta$ -lactams in the USA, and we therefore gathered data on this from three different sources covering respectively the periods between 1980 and 1992 (McCaig and Hughes 1995), between 1992 and 2000 (McCaig et al. 2003), and between 2000 and

2012 (CDDEP 2017). There first and second sources overlapped in the year 1992, and the second and third sources overlapped in the year 2000. We used these 2 years of overlap to scale the data for consistency between the three sources. Specifically, the values from the second source (McCaig et al. 2003) were scaled so that the 2000 value is equal to the one in the third source (CDDEP 2017), and values from the first source (McCaig and Hughes 1995) were then scaled so that the 1992 value is equal to the one in the previously rescaled second source. The final rescaled data are therefore measured as in the third source, namely in standard units of  $\beta$ -lactams (i.e., narrow-spectrum and broad spectrum penicilins plus cephalosporins) consumed per 1000 population in the USA (CDDEP 2017). These data show that the consumption of  $\beta$ -lactams almost doubled between 1980 and 1991 and subsequently decreased to reach around 2010 levels comparable to the early 1980s (Fig. 3). These trends on  $\beta$ -lactams consumption therefore appear to be very similar to the ones observed for the USA300 growth rate without the use of covariates (red results in Fig. 3). To confirm this observation, we repeated our phylodynamic analysis with integration of the  $\beta$ -lactam use as a covariate (blue results in Fig. 3). We found that the covariate was significantly associated with growth rate, with a mean posterior effect of 0.48% and 95% credible interval (0.18–0.71). The growth rate dynamics inferred when using covariate data were almost identical to those inferred without the use of covariate data, except for a clear reduction of the width of the intervals which reflects the gain in information when combining two independent types of data. USA300 is also partly resistant to macrolides and quinolones, but the consumption of these antibiotics increased in the USA throughout the 1990s (McCaig et al. 2003), and the subset of sensitive genomes (31.9% for ciprofloxacin and 6.3% for erythromycin) did not exhibit different phylodynamic properties compared to resistant genomes (Uhlemann et al. 2014), so that these antibiotics could not explain the USA300 growth rate dynamics.

Our analysis therefore suggests that the rise in  $\beta$ -lactams consumption in the 1980s was responsible for the emergence of the highly successful USA300 lineage. From the mid-1990s, the use of  $\beta$ -lactams has declined, both due to an overall reduction in antibiotic use and a diversification of the type of antibiotics prescribed (McCaig et al. 2003; CDDEP 2017), and the growth rate of USA300 has consequently decreased. Importantly, the consumption of antibiotics is expected to be associated with the growth rates of resistant bacterial pathogens, rather than with their effective population sizes, which here is not at all correlated with the covariate (Fig. 3). Amongst pairs of genomes sampled from the same hosts, the distribution of genomic distance had a mean of 1.4 substitution (Uhlemann et al. 2014). If we assume that sampling occurred on average in the middle of the carriage duration (i.e.,  $\psi/2$  time after infection), the evolutionary time separating the two genomes is between 0  $\psi$  depending on the level of within-host genetic drift (Didelot et al. 2016). Given that

the molecular clock rate of USA300 is approximately 3 substitutions per year (Uhlemann et al. 2014; Alam et al. 2015), the average duration of infections in this outbreak is therefore around  $\psi = 1.4 \times 2/3 = 0.93$  year. In the first half of the 1990s, the growth rate peaked around 1 per year (Fig. 3) and using Equation 10 we estimate that the reproduction number was around  $R = 1.93$ , which is in good agreement with the recent estimate  $R = 1.5$  for MRSA in the US population (Hogea et al. 2014). The fact that this estimate is only modestly above the minimum threshold of  $R = 1$  required for outbreaks to take place could help explain why the USA300 is declining, even though  $\beta$ -lactams are still widely used. The consumption level may have lowered below the threshold caused by the fitness cost of resistance, as previously discussed for other resistant bacteria (Dingle et al. 2017; Whittles et al. 2017).

## DISCUSSION

Many environmental covariates, particularly those with a mechanistic influence on replicative fitness of pathogens, are closely related to the growth rate of epidemic size but not necessarily related to absolute epidemic size. We have found that these relationships can be inferred from random samples of pathogen genetic sequences by relating environmental covariates to the growth rate of the effective population size. This enables the estimation of the fitness effect of environmental covariates as well as the prediction of future epidemic dynamics should conditions change. We have found a clear and highly significant relationship between the growth and decline of community-associated MRSA USA300 and the population-level prescription rates of  $\beta$ -lactam antibiotics (Fig. 3). This relationship is not apparent when comparing antibiotic usage directly with the effective population size of MRSA USA300. Our methodology focused on growth rate is therefore well suited to investigate the drivers of antibiotic resistance, compared to previous phylodynamic methods focused on the effective population size.

The *skygrowth* model can provide a more realistic prior for many infectious disease epidemics where the growth rate of epidemic size is likely to approach stationarity as opposed to the absolute effective population size. Conventional *skyride* and *skygrid* models are prone to erroneously estimating a stable effective population size when genealogical data are uninformative, as for example when estimating epidemic trends in the latter stages of SIR epidemics (Fig. 1). The *skygrowth* model will correctly predict epidemic decline in this situation. Moreover, under ideal conditions, the estimated growth rate can be related to the reproduction number of an epidemic, and the *skygrowth* model provides a simple nonparametric estimator of the reproduction number through time given additional information about the natural history of infection (Equation 10). Caution should be exercised when using the effective population



size as a proxy for epidemic size, as the relationship between the two is complex (cf. Simulation results). In general, there will be close correspondence between the growth of epidemic size and growth of effective population size during periods where the growth rate is relatively constant.

The methods presented here can be applied more generally to evaluate the role of antibiotic stewardship, vaccine campaigns, or other public health interventions on epidemic growth rates. Some environmental covariates, such as independent prevalence estimates, may be more closely related to effective population size rather than growth rates, and future work is indicated on the development of regression models in terms of both statistics. More complex stochastic models can also be considered, such as processes with both autoregressive and moving average components. A variety of mathematical models have been developed to explain *de novo* evolution of antimicrobial resistance as a function of population-level antimicrobial usage (Bonhoeffer et al. 1997; Austin et al. 1999; Spicknall et al. 2013; Whittles et al. 2017), and an important direction for future work will be the development of parametric and semiparametric structured coalescent models (Volz 2012) that can be applied to bacterial phylogenies featuring a mixture of antibiotic sensitive and resistant lineages. This methodology will allow us to estimate key evolutionary parameters, such as the fitness cost and benefit of resistance, or the rate of mutation from sensitive to resistant status, which are needed to make well-informed recommendations on resistance control strategies.

#### SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository: <http://dx.doi.org/10.5061/dryad.9qh7t9t>.

#### FUNDING

This work was supported by the National Institute of General Medical Sciences [U01 GM110749 to E.M.V.]; and the Medical Research Councils Centre for Outbreak Analysis and Modelling [MR/K010174].

#### REFERENCES

- Alam M.T., Read T.D., Petit R.A., Boyle-Vavra S., Miller L.G., Eells S.J., Daum R.S., David M.Z. 2015. Transmission and microevolution of USA300 MRSA in U.S. households: evidence from whole-genome sequencing. *MBio* 6:1–10.
- Allen L. 2008. An introduction to stochastic epidemic models. In: Jean-Michel Morel, Bernard Teissier, editors, *Mathematical epidemiology*. Lecture Notes in Mathematics. vol. 1945. Berlin: Springer. p. 81–130.
- Austin D.J., Kristinsson K.G., Anderson R.M. 1999. The relationship between the volume of antimicrobial consumption in human communities and the frequency of resistance. *Proc. Natl. Acad. Sci. USA* 96:1152–1156.
- Baele G., Suchard, M. A., Rambaut A., Lemey P. 2016. Emerging concepts of data integration in pathogen phylodynamics. *Syst. Biol.* 66:e47–e65.
- Biek R., Henderson J.C., Waller L.A., Rupprecht C.E., Real L.A. 2007. A high-resolution genetic signature of demographic and spatial expansion in epizootic rabies virus. *Proc. Natl. Acad. Sci. USA* 104:7993–8.
- Bonhoeffer S., Lipsitch M., Levin B.R. 1997. Evaluating treatment protocols to prevent antibiotic resistance. *Proc. Natl. Acad. Sci. USA* 94:12106–12111.
- CDDEP. 2017. The Center for Disease Dynamics Economics and Policy. ResistanceMap. Available from: URL <http://resistancemap.cddep.org/> (accessed July 2017).
- Challagundla L., Luo X., Tickler I.A., Didelot X., Coleman D.C., Shore A.C., Coombs G.W., Sordelli D.O., Brown E.L., Skov R., Larsen R., Reyes J., Robledo I.E., Vazquez G.J., Rivera R., Fey P.D., Stevenson K., Wang S.-H., Kreiswirth B.N., Mediavilla J.R., Arias C.A., Planet P.J., Nolan R.L., Tenover F.C., Goering R.V., Robinson D.A. 2018. Range expansion and the origin of USA300 North American epidemic methicillin-resistant *Staphylococcus aureus*. *MBio* 9:e02016–17.
- Chambers H.F., Deleo F.R. 2009. Waves of resistance: *Staphylococcus aureus* in the antibiotic era. *Nat. Rev. Microbiol.* 7:629–41.
- Childs J.E., Curns A.T., Dey M.E., Rev L.A., Feinstein L., Bjornstad O.N., Krebs J.W. 2000. Predicting the local dynamics of epizootic rabies among raccoons in the United States. *Proc. Natl. Acad. Sci. USA* 97:13666–13671.
- Dantes R., Mu Y., Belflower R., Aragon D., Dumyati G., Harrison L.H., Lessa F.C., Lynfield R., Nadle J., Petit S., Ray S.M., Schaffner W., Townes, J., Fridkin S. 2013. National burden of invasive methicillin-resistant *Staphylococcus aureus* infections, United States, 2011. *JAMA Intern. Med.* 173:1970–1979.
- de Silva E., Ferguson N.M., Fraser C. 2012. Inferring pandemic growth rates from sequence data. *J. R. Soc. Interface* 9:1797–1808.
- Dearlove B., Wilson D. 2013. Coalescent inference for infectious disease: meta-analysis of hepatitis C. *Philos. Trans. R. Soc. B Biol. Sci.* 368:20120314.
- Didelot X., Eyre D.W., Cule M., Ip C.L.C., Ansari M.A., Griffiths D., Vaughan A., O'Connor L., Golubchik T., Batty E.M., Piazza P., Wilson D.J., Bowden R., Donnelly P.J., Dingle K.E., Wilcox M., Walker A. S., Crook D. W., Peto T. E., Harding R.M. 2012. Microevolutionary analysis of *Clostridium difficile* genomes to investigate transmission. *Genome Biol.* 13:R118.
- Didelot X., Walker A.S., Peto T.E., Crook D.W., Wilson D.J. 2016. Within-host evolution of bacterial pathogens. *Nat. Rev. Microbiol.* 14:150–162.
- Dingle K.E., Didelot X., Quan T.P., Eyre D.W., Stoesser N., Golubchik T., Harding R.M., Wilson D.J., Griffiths D., Vaughan A., Others. 2017. Effects of control interventions on *Clostridium difficile* infection in England: an observational study. *Lancet Infect. Dis.* 17:411–421.
- Drummond A.J., Rambaut A., Shapiro B., and Pybus O.G. 2005. Bayesian coalescent inference of past population dynamics from molecular sequences. *Mol. Biol. Evol.* 22:1185–92.
- Drummond A.J., Suchard M.A., Xie D., Rambaut A. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* 29:1969–1973.
- Fraser C. 2007. Estimating individual and household reproduction numbers in an emerging epidemic. *PLoS One* 2:e758.
- Frost S.D.W., Volz E.M. 2010. Viral phylodynamics and the search for an 'effective number of infections'. *Philos. Trans. R. Soc. B* 365:1879–1890.
- Gill M.S., Lemey P., Bennett S.N., Biek R., Suchard M.A. 2016. Understanding past population dynamics: Bayesian coalescent-based modeling with covariates. *Syst. Biol.* 65:1041–1056.
- Gill M.S., Lemey P., Faria N.R., Rambaut A., Shapiro B., Suchard M.A. 2013. Improving Bayesian population dynamics inference: a coalescent-based model for multiple loci. *Mol. Biol. Evol.* 30:713–724.
- Glaser P., Martins-Simões P., Villain A., Barbier M., Tristan A., Bouchier C., Ma L., Bes M., Laurent F., Guillemot D., Wirth T., Vandenesch F. 2016. Demography and intercontinental spread of the USA300 community-acquired methicillin-resistant *Staphylococcus aureus* lineage. *MBio* 7:1–11.

- Goldfarb D. 1970. A family of variable-metric methods derived by variational means. *Math. Comput.* 24:23–26.
- Griffiths R., Tavare S. 1994. Sampling theory for neutral alleles in a varying environment. *Philos. Trans. R. Soc. B Biol. Sci.* 344:403–410.
- Ho S.Y., Shapiro B. 2011. Skyline-plot methods for estimating demographic history from nucleotide sequences. *Mol. Ecol. Resour.* 11:423–434.
- Hogea C., Van Effelterre T., Acosta C.J. 2014. A basic dynamic transmission model of *Staphylococcus aureus* in the US population. *Epidemiol. Infect.* 142:468–478.
- Karcher M.D., Palacios J.A., Bedford T., Suchard M.A., Minin V.N. 2016. Quantifying and mitigating the effect of preferential sampling on phylodynamic inference. *PLoS Comput. Biol.* 12:e1004789.
- Karcher M.D., Palacios J.A., Lan S., Minin V.N. 2017. phylodyn: an R package for phylodynamic simulation and inference. *Mol. Ecol. Resour.* 17:96–100.
- Kingman J. 1982. The coalescent. *Stoch. Process. Appl.* 13:235–248.
- Koelle K., Ratmann O., Rasmussen D.A., Pasour V., Mattingly J. 2011. A dimensionless number for understanding the evolutionary dynamics of antigenically variable RNA viruses. *Proc. R. Soc. B Biol. Sci.* 278:3723–3730.
- Ledda A., Price J.R., Cole K., Llewelyn M.J., Kearns A.M., Crook D.W., Paul J., Didelot X. 2017. Re-emergence of methicillin susceptibility in a resistant lineage of *Staphylococcus aureus*. *J. Antimicrob. Chemother.* 72:1285–1288.
- McCaig L.F., Besser R.E., Hughes J.M. 2003. Antimicrobial drug prescription in ambulatory care settings, United States, 1992–2000. *Emerg. Infect. Dis.* 9:432–437.
- McCaig L.F., Hughes J.M. 1995. Trends in antimicrobial drug prescribing among office-based physicians in the United States. *J. Am. Med. Assoc.* 273:214–219.
- Minin V.N., Bloomquist E.W., Suchard M.A. 2008. Smooth skyride through a rough skyline: Bayesian coalescent-based inference of population dynamics. *Mol. Biol. Evol.* 25:1459–1471.
- Monroe B., Yager P., Blanton J., Birhane M., Wadhwa A., Orciari, L., Petersen B., Wallace R. 2016. Rabies surveillance in the United States during 2014. *J. Am. Vet. Med. Assoc.* 248:777–788.
- Palacios J.A., Minin V.N. 2013. Gaussian process-based Bayesian nonparametric inference of population size trajectories from gene genealogies. *Biometrics* 69:8–18.
- Planet P.J. 2017. Life after USA300: the rise and fall of a superbug. *J. Infect. Dis.* 215:S71–S77.
- Pybus O.G. 2001. The epidemic behavior of the hepatitis C virus. *Science* 292:2323–2325.
- Pybus O.G., Rambaut A., Harvey P.H. 2000. An integrated framework for the inference of viral population history from reconstructed genealogies. *Genetics* 155:1429–1437.
- Rosenberg N.A., Nordborg M. 2002. Genealogical trees, coalescent theory and the analysis of genetic polymorphisms. *Nat. Rev. Genet.* 3:380–90.
- Spicknall I.H., Foxman B., Marrs C.F., Eisenberg J.N.S. 2013. A modeling framework for the evolution and spread of antibiotic resistance: literature review and model categorization. *Am. J. Epidemiol.* 178:508–520.
- Strimmer K., Pybus O.G. 2001. Exploring the demographic history of DNA sequences using the generalized skyline plot. *Mol. Biol. Evol.* 18:2298–2305.
- Tenover F.C., Goering R.V. 2009. Methicillin-resistant *Staphylococcus aureus* strain USA300: origin and epidemiology. *J. Antimicrob. Chemother.* 64:441–446.
- Uhlemann A.-C., Dordel J., Knox J.R., Raven K.E., Parkhill J., Holden M.T.G., Peacock S.J., Lowy F.D. 2014. Molecular tracing of the emergence, diversification, and transmission of *S. aureus* sequence type 8 in a New York community. *Proc. Natl. Acad. Sci. USA* 111:6738–43.
- Vaughan T.G., Drummond A.J. 2013. A stochastic simulator of birth-death master equations with application to phylodynamics. *Mol. Biol. Evol.* 30:1480–1493.
- Volz E.M. 2012. Complex population dynamics and the coalescent under neutrality. *Genetics* 190:187–201.
- Volz E.M., Frost S.D.W. 2014. Sampling through time and phylodynamic inference with coalescent and birth–death models. *J. R. Soc. Interface* 11:20140945.
- Volz E.M., Koelle K., Bedford T. 2013. Viral phylodynamics. *PLoS Comput. Biol.* 9:e1002947.
- Volz E.M., Kosakovsky Pond S.L., Ward M.J., Leigh Brown A. J., Frost S.D.W. 2009. Phylodynamics of infectious disease epidemics. *Genetics* 183:1421–30.
- Volz E.M., Romero-Severson E., Leitner T. 2017. Phylodynamic inference across epidemic scales. *Mol. Biol. Evol.* 34:1276–1288.
- Whittles L., White P., Didelot X. 2017. Estimating the fitness cost and benefit of cefixime resistance in *Neisseria gonorrhoeae* to inform prescription policy: a modelling study. *PLoS Med.* 14:e1002416.