

1 Estimating the fitness cost and benefit of antimicrobial resistance from  
2 pathogen genomic data

3 David Helekal<sup>1</sup>, Matt Keeling<sup>2</sup>, Yonatan H Grad<sup>3</sup>, Xavier Didelot<sup>4,\*</sup>

4 <sup>1</sup> Centre for Doctoral Training in Mathematics for Real-World Systems, University of Warwick, UK

5

6 <sup>2</sup> Mathematics Institute and School of Life Sciences, University of Warwick, UK

7

8 <sup>3</sup> Department of Immunology and Infectious Diseases, TH Chan School of Public Health, Harvard  
9 University, USA

10

11 <sup>4</sup> School of Life Sciences and Department of Statistics, University of Warwick, UK

12

13 \* Corresponding author. Tel: 0044 (0)2476 572827. Email: [xavier.didelot@warwick.ac.uk](mailto:xavier.didelot@warwick.ac.uk)

14 Running title: Phylodynamics of pathogen antimicrobial resistance

15

16 Keywords: genomic epidemiology, phylodynamics, antimicrobial resistance, resistance fitness cost

## 17 ABSTRACT

18 Increasing levels of antibiotic resistance in many bacterial pathogen populations is a major threat  
19 to public health. Resistance to an antibiotic provides a fitness benefit when the bacteria is exposed  
20 to this antibiotic, but resistance also often comes at a cost to the resistant pathogen relative to  
21 susceptible counterparts. We lack a good understanding of these benefits and costs of resistance for  
22 many bacterial pathogens and antibiotics, but estimating them could lead to better use of antibiotics  
23 in a way that reduces or prevents the spread of resistance. Here, we propose a new model for the  
24 joint epidemiology of susceptible and resistant variants, which includes explicit parameters for the cost  
25 and benefit of resistance. We show how Bayesian inference can be performed under this model using  
26 phylogenetic data from susceptible and resistant lineages and that by combining data from both we are  
27 able to disentangle and estimate the resistance cost and benefit parameters separately. We applied our  
28 inferential methodology to several simulated datasets to demonstrate good scalability and accuracy.  
29 We analysed a dataset of *Neisseria gonorrhoeae* genomes collected between 2000 and 2013 in the USA.  
30 We found that two unrelated lineages resistant to fluoroquinolones shared similar epidemic dynamics  
31 and resistance parameters. Fluoroquinolones were abandoned for the treatment of gonorrhoea due to  
32 increasing levels of resistance, but our results suggest that they could be used to treat a minority of  
33 around 10% of cases without causing resistance to grow again.

## 34 INTRODUCTION

35 The levels of antimicrobial resistance of many pathogens have risen worryingly over the past few  
36 decades. In a report on the threat posed by antibiotic resistance published by the CDC (Centres  
37 for Disease Control and Protection), three microorganisms including *N. gonorrhoeae* are classified as  
38 posing an urgent threat level, and twelve more represent a serious threat to public health [1]. A  
39 review on antimicrobial resistance estimated that resistance claims at least 700,000 lives per year  
40 worldwide and that the death toll could go up to 10 million per year by 2050 if current trends  
41 are allowed to continue [2], and a recent study estimated that there were almost 5 million deaths  
42 associated with resistance in 2019 [3]. Few new antimicrobials have been developed and deployed  
43 since the 1970s, whereas resistance to new drugs often emerges soon after initial introduction [4],  
44 so that several pathogens are dangerously close to becoming completely untreatable. Effectively  
45 tackling antimicrobial resistance requires greater understanding of epidemiological and evolutionary  
46 factors leading to emergence of resistance and the spread of resistance through pathogen populations.  
47 Achieving this goal requires development of mathematical models of antimicrobial resistance and robust  
48 statistical analysis of epidemiological models with informative observations. This modelling approach  
49 to resistance was initiated in the late 1990s [5, 6] and has led to the development of many models,  
50 appropriate for different organisms, mode of spread, study scale and context [7].

51 Resistance brings a clear fitness benefit to pathogens acquiring it in the presence of antimicrobials.  
52 The net value of this fitness benefit therefore increases with the frequency with which the specific  
53 antimicrobial is employed, either against the pathogen itself or more generally in the case of a pathogen  
54 that can be carried asymptotically. However, resistance also typically comes with a fitness cost to the  
55 pathogen [8]. The simplest demonstration of this effect is when discontinued use of an antimicrobial  
56 leads to reductions in resistance rates. The fitness costs and benefits of resistance remain poorly  
57 understood for many pathogens and antimicrobials [9]. A better quantification of resistance benefits  
58 and costs is required to provide a solid basis for evaluating the potential effectiveness of public health  
59 intervention measures proposed to exploit fitness costs in the hope of stopping or even reversing the  
60 spread of resistance [9]. For example, the numbers of gonorrhoea cases sensitive and resistant to  
61 cefixime in England over a decade was recently analysed to quantify the cost and benefit associated  
62 with resistance to this antibiotic [10]. These estimates were used to predict that cefixime could  
63 be reintroduced to treat a minority ( $\sim 25\%$ ) of gonorrhoea cases without causing an increase in  
64 cefixime resistance levels, which would reduce the risk of emergence of resistance to the currently used  
65 antibiotics. Moreover, the extent of the fitness cost of resistance can vary by genomic background [11],  
66 such that the effect of interventions that seek to capitalize on the fitness costs of resistance may be  
67 lineage dependent. Therefore, it is necessary to estimate fitness costs at the per lineage level.

68 Pathogen genomic data has great potential to help us understand the evolutionary and epidemiological  
69 dynamics of infectious disease [12]. An important advantage of this phylodynamic approach is that  
70 analysis of genomic data is less sensitive to sampling biases, especially when using a coalescent  
71 framework which describes the ancestry process conditional on sampling [13]. A few studies have used  
72 this approach to shed light on the fitness cost associated with antimicrobial resistance. For example, a  
73 study showed the association between the growth rate of a methicillin-resistant *Staphylococcus aureus*  
74 lineage and consumption of beta-lactams [14]. Other studies quantified the relative transmission fitness  
75 of resistance mutations in HIV [15] and *Mycobacterium tuberculosis* [16]. Here, we take a different  
76 approach by modelling explicitly the phylodynamic trajectories of the sensitive and resistant lineages  
77 as a function of the fitness cost, which is constant, and the fitness benefit, which depends on the  
78 antimicrobial consumption. Our method therefore requires three inputs: the amount of antimicrobial  
79 being used over time, genomic data from a sensitive lineage, and genomic data from a resistant lineage.  
80 From this we disentangle the fitness cost and benefit of resistance, thereby providing the parameters

81 needed to predict phylodynamic trajectories and inform recommendations on how to use antimicrobials  
82 without worsening the resistance threat.

## 83 METHODS

### 84 Overall approach

85 Pathogen phylogenetic data contains information about past population size dynamics of the pathogen  
86 under study [12, 17]. Under assumptions of the epidemic process being characterised well enough by  
87 a simple compartmental epidemic model, this information about population size dynamics can be  
88 translated into epidemic trajectories [18, 19]. These epidemic trajectories can be described using an  
89 epidemic model which accounts for the effects of a fitness cost and benefit of resistance to a specific  
90 antimicrobial. As the use of this antimicrobial changes through time, so will the net fitness of the  
91 particular lineage in consideration. This will in turn lead to changes in the behaviour of the epidemic  
92 trajectory. However, not all changes in the behaviour of the epidemic trajectory will be due to changes  
93 in the fitness of the resistant phenotype. Confounding factors, such as as depletion of susceptibles or  
94 changes in host behaviour will also affect the epidemic trajectory. Under relatively mild assumptions  
95 detailed below changes in these confounding factors will affect other strains equally. We can therefore  
96 use as “control” some data from a susceptible lineage, ideally closely related and with the same  
97 resistance profile to other antimicrobials used in significant amounts as primary treatment. Differences  
98 between the trajectories of the sensitive and resistant lineages can then be ascribed specifically to  
99 resistance, allowing us to estimate the associated fitness cost and benefit parameters.

100 Let us consider a pathogen causing infections that are or were treated with a certain antimicrobial  
101 compound. We assume that at some point in the past one or several strains with resistance to this  
102 antimicrobial compound have arisen. Our aim is to quantify the fitness cost and benefit of the resistance  
103 to this antimicrobial for a given lineage as a function of use of the antimicrobial of interest through time.  
104 To this end we need data that quantify the use over time of the given antimicrobial to treat infections  
105 caused by this pathogen, as well as a reasonable sample of sequenced case isolates from infections caused  
106 by the pathogen over time. Furthermore, we need information that characterises the resistance profiles  
107 of the individual isolates, which can be obtained either by resistance screening *in vitro*, or predicted  
108 from the sequences *in silico* [20]. A dated phylogeny of these samples is estimated, for example using  
109 BEAST [21], BEAST2 [22] or BactDating [23]. This phylogeny is then used as the starting point for  
110 analysis [24], to identify which samples belong to resistant and susceptible lineages and to select related  
111 lineages for further study that are wholly resistant or susceptible to the antimicrobial of interest, but  
112 otherwise similar in their resistance profiles. Note that for simplicity resistance is treated as a binary  
113 trait, with samples being either resistant or susceptible to antimicrobials, as is usually the case in  
114 resistance modelling studies [7].

### 115 Transmission model derivation

In order to estimate the fitness cost and benefit of antimicrobial resistance, a transmission model  
needs to be specified. We focus on estimating the fitness parameters of a particular lineage harbouring  
a certain treatment resistant phenotype when previous infection does not confer immunity against  
reinfection. Under the simplifying assumptions that the host population is unstructured and that  
past infections do not confer any immunity, the multi-strain Susceptible-Infected-Susceptible (SIS)

is a reasonable model [25, 26]. Fluctuations in the carriage levels of different strains can also be due to external factors, such as changes in host demography or behaviours. Left unaccounted, such fluctuations would bias estimates of the fitness cost and benefit of resistance to a given antimicrobial. Therefore, we modify the model with time-varying transmission rate  $\beta(t)$  and population size  $N(t)$ . This leads to an  $n$ -strain model described by a system of the following  $n$ -coupled ordinary differential equations (ODEs):

$$\begin{aligned} \frac{dI_1(t)}{dt} &= \frac{\beta(t)S(t)I_1(t)}{N(t)} - \gamma_1(t)I_1(t) \\ \frac{dI_2(t)}{dt} &= \frac{\beta(t)S(t)I_2(t)}{N(t)} - \gamma_2(t)I_2(t) \\ &\vdots \\ \frac{dI_n(t)}{dt} &= \frac{\beta(t)S(t)I_n(t)}{N(t)} - \gamma_n(t)I_n(t) \end{aligned} \tag{1}$$

Where  $I_j(t)$  denotes the number of people infected with the  $j$ -th strain at time  $t$ .  $\beta(t)$  is the transmission rate that varies with time due for example to changes that are not specific to any strain, for example host behaviour.  $N(t)$  is the host population size which may also change with time due to demographic factors.  $\gamma_j(t)$  is the recovery rate of the  $j$ -th strain at time  $t$ . These may or may not vary with time through their dependency on the antimicrobial usage which changes with time. Finally  $S(t)$  denotes the number of susceptible hosts

$$S(t) = \left( N(t) - \sum_{j=1}^n I_j(t) \right) \tag{2}$$

116 Typically this model could simply be reduced to a two strain model, averaging over all lineages that  
 117 are phenotypically similar in their resistance profiles. However, this is undesirable, as some of the  
 118 lineages with the same resistance phenotype could differ in fitness due to different genomic background  
 119 which would confound our estimates. Furthermore this sort of model would not be readily tractable  
 120 in a genomic framework, because phylogenetic data is generally going to be informative about the  
 121 dynamics of a particular lineage only.

We therefore need to focus on the resolution of individual lineages. We note that environmental effects such as fluctuations in host population size or behaviour affect all lineages equally, if the population is well mixed. We denote the combination of these effects as  $b(t) = \beta(t)S(t)/N(t)$ . Conditional on the knowledge trajectory of  $b(t)$  the ODEs in Equation become uncoupled, and this allows us to reduce the system to uncoupled equations corresponding to the strains we will be focusing on. As such we will treat  $b(t)$  as a random object that needs to be inferred. We further assume that for the susceptible lineages the average recovery rate denoted  $\gamma_s$  does not change over time, whereas for the resistant strains it takes one of two values:  $q_U\gamma_s$  if a given patient is treated with the antimicrobial of interest, or  $q_T\gamma_s$  otherwise. If we also consider the known proportion  $u(t)$  with which the antimicrobial is used as primary treatment at time  $t$ , this fully determines the average recovery rate of the resistant lineages as:

$$\gamma_r(t) = u(t)q_T\gamma_s + (1 - u(t))q_U\gamma_s \tag{3}$$

We can now fully write down the equations of the model we will be using for the sensitive and resistant

lineages, respectively:

$$\begin{aligned}\frac{dI_s(t)}{dt} &= b(t)I_s(t) - \gamma_s I_s(t) \\ \frac{dI_r(t)}{dt} &= b(t)I_r(t) - [u(t)q_T + (1 - u(t))q_U] \gamma_s I_r(t)\end{aligned}\tag{4}$$

122 This model can be applied to any number of resistant and sensitive lineages, simply by adding lineages  
123 associated terms to the likelihood and adding required parameters. This is straightforward as the  
124 individual lineages are independent conditional on  $b(t)$ . but for simplicity the remainder of methods  
125 description focuses on the case of a single sensitive and a single resistant strain, with the general case  
126 being a straightforward extension.

## 127 Link to phylogenies

Having defined the epidemiological model, we can now link it to the phylogenetic process. Based on [18, 27], the instantaneous coalescent rates for a single pair of lineages can be derived as

$$\lambda_s(t) = \frac{2b(t)}{I_s(t)} \text{ and } \lambda_r(t) = \frac{2b(t)}{I_r(t)}\tag{5}$$

in the susceptible and resistant populations, respectively. The likelihood of a dated phylogeny  $\mathbf{g}$  with  $n$  leaves at times  $s_1 < \dots < s_n$  and  $n - 1$  coalescent events at times  $c_1 < \dots < c_{n-1}$  and  $A(t)$  lineages at time  $t$  is therefore given by [28]:

$$p(\mathbf{g}|\lambda(t)) = \exp\left(-\int_{-\infty}^{\infty} \mathbb{1}[A(t) \geq 2] \binom{A(t)}{2} \lambda(t) dt\right) \prod_{i=1}^{n-1} \lambda(c_i)\tag{6}$$

128 Where  $\lambda(t) = \lambda_s(t)$  and  $\lambda(t) = \lambda_r(t)$  for the susceptible and resistant phylogenies, respectively.  
129 However, in most cases, and indeed in our case, the integral in Equation 6 is not analytically intractable.  
130 Furthermore, the antibiotic use data is unlikely to span the entire phylogeny. Therefore, we define the  
131 approximate likelihood for the phylogeny truncated to  $[t_{\min}, t_{\max}]$ , which is the intersection interval  
132 spanned by the antibiotic use data and the phylogenies under study.

As such we resort to the standard way of approximating coalescent likelihoods [29], partitioning the interval  $[t_{\min}, t_{\max}]$  into a fine mesh  $t_{\min} = t_1 < t_2 < t_3 < \dots < t_N = t_{\max}$  such that  $t_i - t_{i-1} < \Delta_t$  and that all sampling and coalescent times between  $t_{\min}$  and  $t_{\max}$  are included in the mesh:

$$p(\mathbf{g}|\lambda(t)) = \exp\left(-\sum_{i=2}^N (t_i - t_{i-1}) \binom{A(t_{i-1})}{2} \lambda(t_{i-1})\right) \prod_{i=1}^{n-1} \mathbb{1}[c_i \in [t_{\min}, t_{\max}]] \lambda(c_i)\tag{7}$$

## 133 Bayesian inference

We first re-scale time from the interval  $[t_{\min}, t_{\max}]$  to  $[-1, 1]$ . Denoting the scale factor  $D = (t_{\max} - t_{\min})/2$  associated with this re-scaling, we account for this in the model by defining  $\tilde{\gamma}_s = \gamma_s D$ . The model consists of independent first-order linear homogeneous ODEs for each strain with time-varying coefficients. The solutions at time  $t$  subject to initial conditions  $I_s(0) = I_{s0}$  and  $I_r(0) = I_{r0}$

can be obtained in terms of the integral of the instantaneous rates up to time  $t$ :

$$\begin{aligned} I_s(t) &= I_{s0} \exp \left\{ \int_0^t b(\tau) - \gamma_s d\tau \right\} \\ I_r(t) &= I_{r0} \exp \left\{ \int_0^t b(\tau) - [u(\tau)q_T + (1 - u(\tau))q_U]\gamma_s d\tau \right\} \end{aligned} \quad (8)$$

As it stands, this model would not be well-suited for performing Bayesian inference, primarily due to the difficulty in choosing a sensible prior on  $b(t)$ , and a very complicated dependency structure between the initial conditions and  $b(t)$ . As such we re-parameterise the model by directly modelling the logarithm of  $I_s(t)$  as a Gaussian Process:

$$C(t) = \log I_s(t) - \mu_s \quad (9)$$

Where  $C(t)$  is an appropriately chosen zero mean Gaussian Process, and  $\mu_s$  is the susceptible intercept which relates to the susceptible initial condition  $I_{s0}$  as follows:

$$\mu_s = \log I_{s0} - C(0) \quad (10)$$

We use this formulation principally to loosen the coupling between the intercept parameter and the Gaussian Process in order to speed up sampling. From this we can compute  $b(t)$  and  $\log I_r(t)$  as

$$b(t) = \frac{d}{dt}C(t) + \gamma_s \quad (11)$$

and

$$\begin{aligned} \log I_r(t) &= C(t) + \mu_r + \int_0^t \gamma_s d\tau - \int_0^t u(\tau)q_T\gamma_s d\tau - \int_0^t (1 - u(\tau))q_U\gamma_s d\tau \\ &= C(t) + \mu_r + \gamma_s \int_0^t 1 - u(\tau)(q_T - q_U) - q_U d\tau \\ &= C(t) + \mu_r + (1 - q_U)t - (q_T - q_U) \int_0^t u(\tau) d\tau \end{aligned} \quad (12)$$

Once again we follow the same reasoning for the resistant trajectory intercept  $\mu_r$ , relating it to  $I_{r0}$  as:

$$\mu_r = \log I_{r0} - C(0) \quad (13)$$

Note that  $\frac{d}{dt}C(t)$  exists as long as the associated covariance kernel is sufficiently smooth such as in the case of the radial basis function (RBF) kernel [30] which we used. Evaluating a full-rank, Gaussian process with differentiable trajectories on the entirety of the mesh would be prohibitively expensive due to the  $O(n^3)$  computational complexity. Instead, we work with a low-rank representation of  $C(t)$  based on the framework introduced in [31]. This leads to the representation of the low-rank projection of  $C(t)$ , denoted by  $\hat{C}(t)$

$$\hat{C}(t) = \sum_{j=1}^m S_{\text{RBF}} \left( \sqrt{\frac{j\pi}{2L}}; \rho, \alpha \right) \sqrt{\frac{1}{L}} \sin \left( \frac{j\pi}{2L}(t + L) \right) f_j \quad (14)$$

and

$$\frac{d}{dt}\hat{C}(t) = \sum_{j=1}^m S_{\text{RBF}} \left( \sqrt{\frac{j\pi}{2L}}; \rho, \alpha \right) \sqrt{\frac{1}{L}} \frac{j\pi}{2L} \cos \left( \frac{j\pi}{2L}(t + L) \right) f_j \quad (15)$$

134 Where  $f_j$  are independent and identically distributed random variables following the standard Gaussian  
 135 distribution,  $S_{\text{RBF}}(\cdot; \cdot, \cdot)$  is the appropriate spectral density for the RBF kernel,  $\rho$  is the kernel length  
 136 scale and  $\alpha$  is the marginal standard deviation of the kernel [31].

Denote by  $\boldsymbol{\theta} = (\gamma_s, q_U, q_T, I_{s0}, I_{r0}, \hat{C}(t))$  the parameters of the pathogen dynamics model. We can now factorise the model posterior  $\pi(\boldsymbol{\theta}, \alpha, \rho, f_{1:m} \mid \mathbf{g}_s, \mathbf{g}_r)$ , suppressing dependency on  $t$  where appropriate:

$$\pi(\boldsymbol{\theta}, \alpha, \rho, f_{1:m} \mid \mathbf{g}_s, \mathbf{g}_r) \propto \pi(\mathbf{g}_s \mid \lambda_s) \pi(\mathbf{g}_r \mid \lambda_r) \pi(\lambda_s \mid \boldsymbol{\theta}) \pi(\lambda_r \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}, \alpha, \rho, f_{1:m}) \quad (16)$$

The first two terms are computed using the coalescent likelihood in Equation 6. The third term is given by combining Equations 5, 9 and 11. The fourth term is obtained by combining Equations 5, 11 and 12. Finally, the last term is given by:

$$\pi(\boldsymbol{\theta}, \alpha, \rho, f_{1:m}) = \pi(\hat{C}(t) \mid \alpha, \rho, f_{1:m}) \pi(\gamma_s) \pi(q_U) \pi(q_T) \pi(I_{s0}) \pi(I_{r0}) \pi(\alpha) \pi(\rho) \pi(f_{1:m}) \quad (17)$$

137 where the first term is given by the Gaussian process (Equations 14 and 15) and the remaining terms  
 138 correspond to the prior distributions listed below.

## 139 Choice of prior and parameterisation

The model is parameterised with the following prior distributions on the  $[-1, 1]$  time scale:

$$\begin{aligned} \gamma_s &\sim \text{log-normal}(\log \gamma^*, \sigma) \\ q_U, q_T &\sim \text{log-normal}(0, 0.5) \\ I_{r0}, I_{s0} &\sim \text{log-normal}(6, 2) \\ \alpha &\sim \text{gamma}(4, 4) \\ \rho &\sim \text{inverse-gamma}(4.63, 2.21) \\ f_{1:m} &\sim \mathcal{N}(0, 1) \end{aligned} \quad (18)$$

140 The data is not expected to be very informative about the value of  $\gamma_s$ . As such, we impose a fairly  
 141 informative prior on this parameter, centred around a guess  $\gamma^*$  which must be known and supplied  
 142 *a priori*.  $\sigma$  then governs how informative the prior is. We typically use a value of  $\sigma = 0.15$ , which  
 143 includes relative fluctuations of around 10% in its 95% interval. The higher the value of  $\sigma$ , the more  
 144 complicated the geometry and subsequently sampling of the posterior becomes.  $q_U$  and  $q_T$  represent  
 145 relative changes in the recovery rate of the resistant strain corresponding to the fitness cost and benefit  
 146 of resistance. A log-normal prior is a relatively natural choice here, with mean=1, so that there is  
 147 no assumption made on the significance of the effects. A log-standard deviation of 0.5 represents a  
 148 weakly informative choice, with fluctuations over 40% being included in the 95% quantile. Fitness  
 149 costs and benefits that exceed this value are hardly of interest here since they would lead to a very  
 150 rapid selective sweep or extinction. The prior on  $\rho$  was chosen so that approximately 1% of mass lies  
 151 on values of  $\rho < 0.2$  and approximately 1% of mass lies on  $\rho > 2$ . The lower bound was chosen to  
 152 avoid over-fitting, and the upper bound to suppress length scales that exceed the range of data and  
 153 thus cannot be informed about by the data. For all results, we used Hilbert Space Gaussian Process  
 154 (HSGP) approximation parameters of  $L = 6.5$  and  $H = 60$ . This approximation is valid for the 99%  
 155 interval of the length-scale prior used as per [31].

In practice, we encounter very tight correlation between  $q_U$  and  $q_T$ . Heuristically, this is due to the linear structure of the dynamics the magnitude of the cost and benefit of resistance is effectively



observed through their sum weighted by the antimicrobial use over time. This complicates sampling from the posterior and is amenable to re-parameterisation that loosens this coupling. By recognising that the data will in general be much more informative about the overall change in prevalence of the resistant strain rather the instantaneous rates, we reparameterise as follows. Denoting by  $\bar{u}$  the average of the antimicrobial use  $u(t)$  across the time interval we are working on, we introduce the following parameters  $\tilde{q}_1$  and  $\tilde{q}_2$  that relate to  $q_U$  and  $q_T$  via:

$$\begin{aligned} q_U &= \frac{e^{\tilde{q}_1} - \bar{u}e^{\tilde{q}_2}}{1 - \bar{u}} \\ q_T &= e^{\tilde{q}_2} \end{aligned} \tag{19}$$

The Jacobian adjustment to the likelihood associated with this re-parameterisation is proportional to

$$|\det J_q| \propto e^{\tilde{q}_1 + \tilde{q}_2} \tag{20}$$

## 156 Computational implementation

157 The posterior in Equation 16 is a high dimensional distribution and we expect many parameters  
158 to have a high degree of interdependency. In order to sample from this distribution, we use  
159 Dynamic Hamiltonian Monte Carlo, a Hamiltonian Monte Carlo (HMC) sampler available in Stan  
160 [32]. We implemented the model and inference method in a R package which is available at  
161 <https://github.com/dhelekal/ResistPhy/>. All results shown used 4 chains with 2000 iterations  
162 for warmup and 2000 iterations for sampling. For all model parameters and all analysis the bulk  
163 effective sample size (bulk-ESS) was always greater than 500, and all  $\hat{R}$  statistics were lower than  
164 1.05 [33], values that indicate no issues with mixing. We also checked that there were no divergent  
165 transitions at least during the sampling phase.

## 166 RESULTS

### 167 Detailed analysis of a single simulated dataset

168 To validate the performance of this model we first resort to simulation from a 3-strain stochastic SIS  
169 with population size  $N(t)$ , transmission rate  $\beta(t)$  and antimicrobial usage function  $u(t)$  varying over  
170 the past 20 years, as illustrated in Figure 1. The first two strains are susceptible and thus unaffected by  
171 fluctuations in antimicrobial usage, whereas the third strain is resistant and therefore affected. The first  
172 strain represents the bulk of the susceptible lineages and is thus left unobserved. The remaining two  
173 strains represent the observed lineages, susceptible and resistant, respectively. The per-day recovery  
174 rate of the sensitive strain was set to  $\gamma_s = 1/60$ , the fitness cost of resistance to  $q_U = 1.25$  and the  
175 fitness benefit of resistance to  $q_T = 0.45$ . From each of these two observed strains, a dated phylogeny  
176 with 200 leaves was observed, which were sampled over the past 6 years with density proportional to  
177 prevalence.

178 We performed inference on this simulated dataset; the traces are shown in Figure S1 and the posterior  
179 distribution of the kernel parameters in Figure S2. The prevalence and reproduction number  $R(t)$   
180 of both the susceptible and resistant strains are shown in Figure 2. As expected, the inferred values  
181 followed the correct values used in the simulation. The inferred values of the susceptible strain recovery

182 rate  $\gamma_s$  and the cost and benefit of resistance  $q_U$  and  $q_T$  were also found to be close to their correct  
183 values, as shown in Figure 3. The posterior distribution of  $\gamma_s$  was almost identical to the prior, which  
184 was centered on the correct value  $1/60$ , reflecting the fact that the data is uninformative about this  
185 parameter and stressing the importance of using an informative prior. There was a strong negative  
186 correlation between the inferred values of  $q_U$  and  $q_T$ , as expected since these two parameters play  
187 opposite roles in the overall fitness of the resistant strain relative to the sensitive strain. Nevertheless,  
188 we detected both the cost and the benefit associated with resistance, since the ranges of inferred values  
189 for  $q_U$  and  $q_T$  were respectively above and below one, contrary to their log-normal priors with mean one  
190 (Figure 3). Finally, we computed the posterior predictive distribution [34] for the number of ancestral  
191 lineages through time  $A(t)$  and compared this with the input phylogenetic data (Figure S3). The data  
192 and posterior predictive trajectories were similar, indicating a good fit of the model to the data as  
193 indeed would be expected here since the same model was used for simulation and inference.

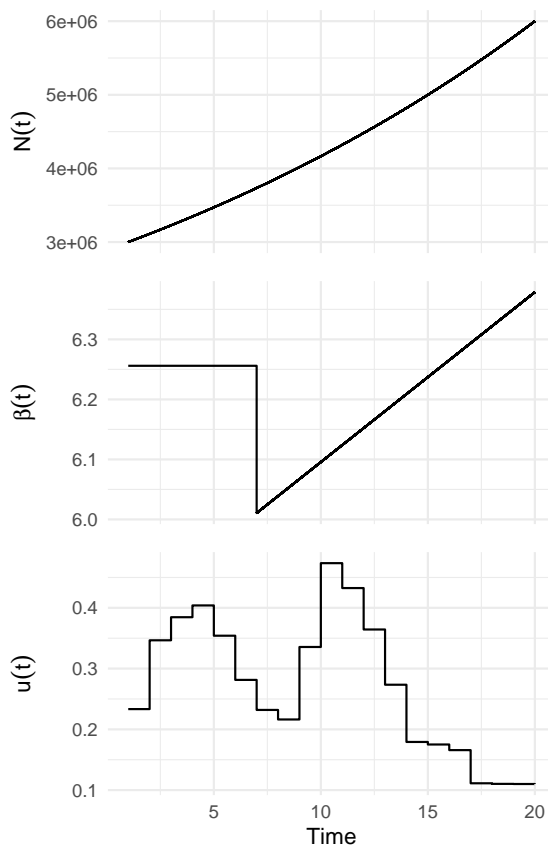


Figure 1: Host population size function  $N(t)$ , transmission rate over time  $\beta(t)$  and antibiotic usage function  $u(t)$  used in the simulated datasets.

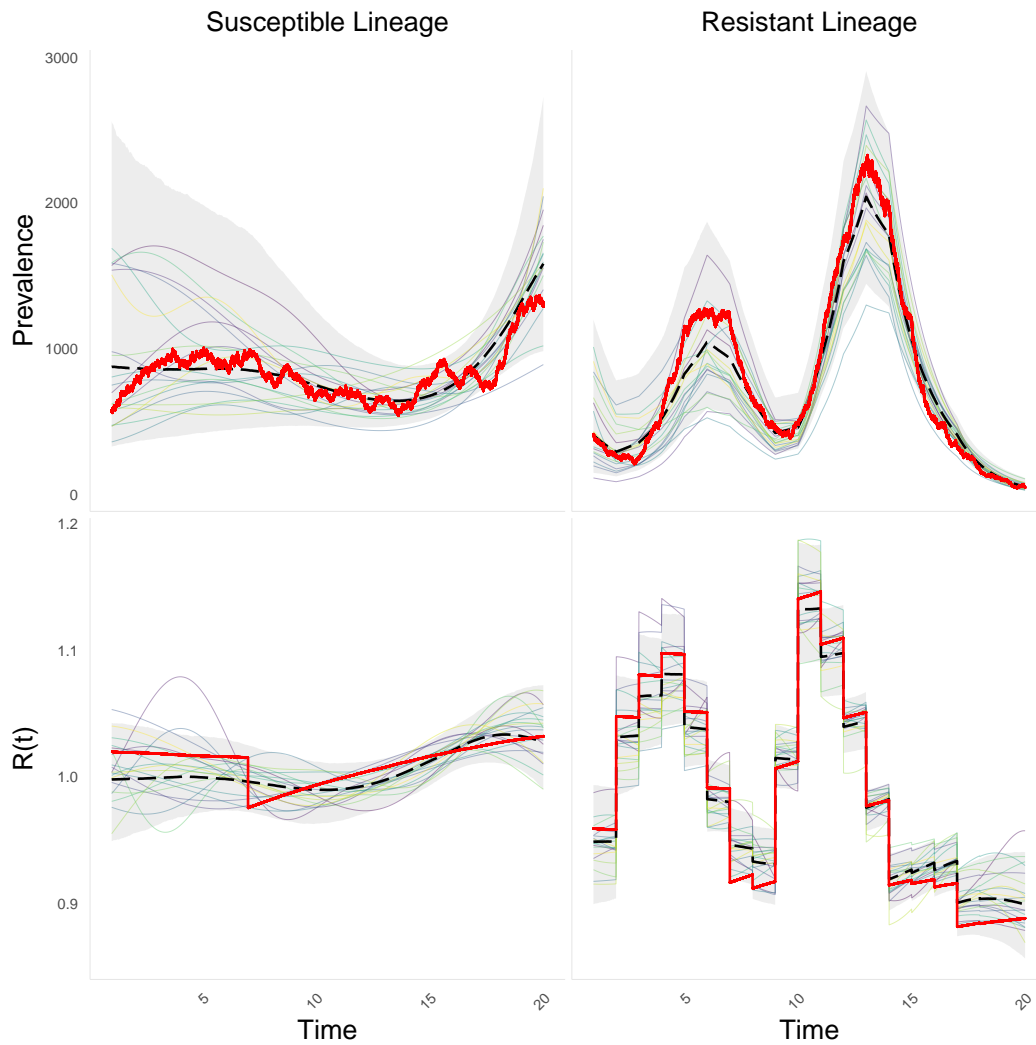


Figure 2: Posterior summary of dynamics for the sensitive (left) and resistant (right) lineages, showing prevalence (top) and reproduction number (bottom). Bold solid red lines indicates simulated values. Posterior median in bold dashed black line. Shaded bands indicate 95% posterior credible intervals. Solid light lines represent posterior draws.

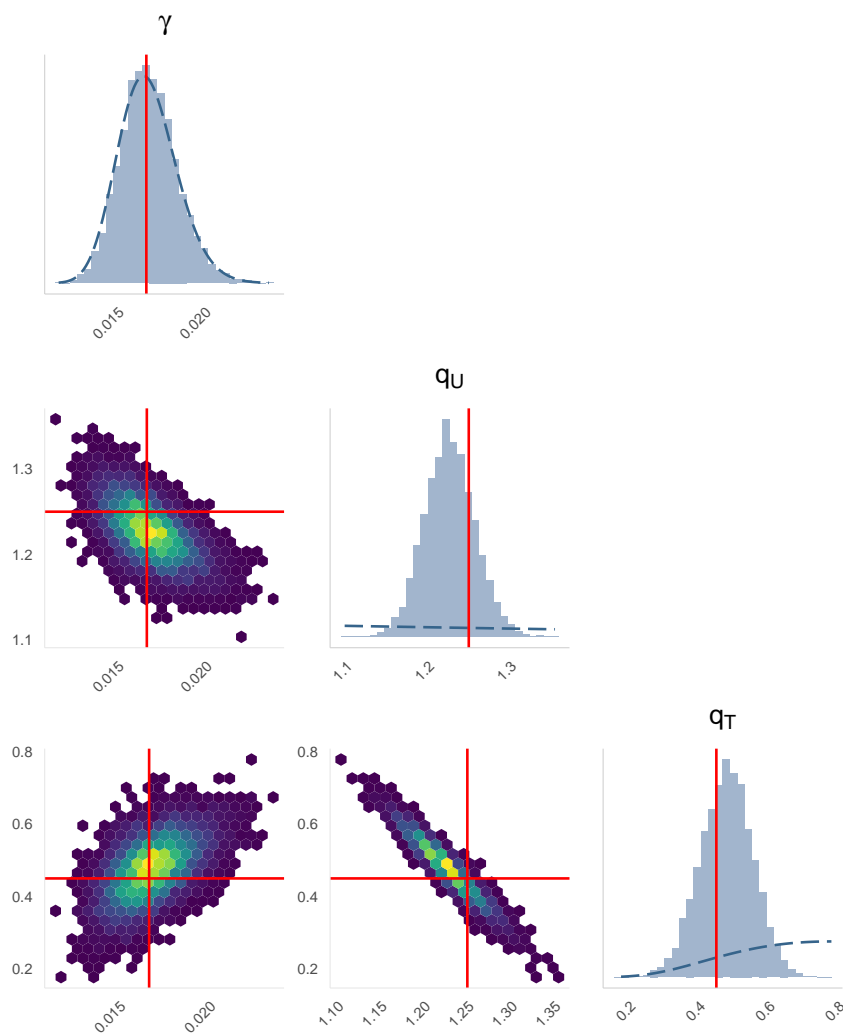


Figure 3: Marginal and joint posterior distributions for the recovery rate of the sensitive lineage ( $\gamma_s$ ), fitness cost ( $q_U$ ) and fitness benefit ( $q_T$ ) of resistance. Bold red solid lines indicate simulation values. Bold blue dashed lines indicate prior density values.

## 194 Benchmark using multiple simulated datasets

195 We repeated the same application of our inference method to data simulated in the same conditions  
196 as described above and illustrated in Figure 1, except the values of the fitness cost and benefit of  
197 resistance were varied. A total of 50 simulated datasets were generated and analysed, with the fitness  
198 cost  $q_U$  increasing linearly from 1 to 1.2, and the fitness benefit  $q_T$  decreasing linearly from 1 to 0.5.  
199 The prevalence of the susceptible and resistant strains in these simulations are shown in Figure S4.  
200 The results of inference are illustrated in Figure 4 and show that in almost all cases the posterior 95%  
201 credible intervals covered the correct values of the fitness cost and benefit of resistance used in the  
202 simulations.

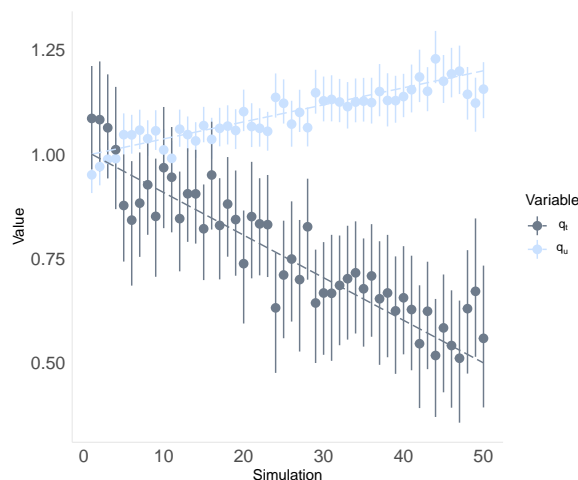


Figure 4: Posterior recovery rate summaries versus ground truth.

## 203 Application to fluoroquinolone resistant *N. gonorrhoeae* in USA

204 We demonstrate the use of our model and inferential framework by estimating the cost and benefit  
205 of fluoroquinolone resistance in *N. gonorrhoeae*. Based on the 1102 genomes collected between 2000  
206 and 2013 by the CDC Gonococcal Isolate Surveillance Project [35], a recombination-corrected tree  
207 was constructed using ClonalFrameML [36] and dated using BactDating [23]. As there are two major  
208 fluoroquinolone resistant lineages present in this phylogeny [35], we decided to do a comparative study.  
209 The two fluoroquinolone resistant lineages and one fluoroquinolone susceptible lineage were selected  
210 based on similar resistance profiles against other relevant antibiotics. By inspecting the antibiotic usage  
211 data and the resistance profiles for the the three lineages (Figure 5) we can see that the resistance  
212 profiles match for antimicrobials that were in use as primary treatment at significant levels after 1995.  
213 As such this is the year we set as the analysis start date ( $t_{\min} = 1995$ ) and the end date is the date  
214 when the last genomes were collected ( $t_{\max} = 2013$ ). Note that a subclade within the susceptible  
215 lineage that displayed a *de novo* gain of resistance to cefixime has been removed. The prior mean for  
216 the per-day recovery rate for the susceptible strain was set to  $\gamma^* = 1/90$  based on previous gonorrhoea  
217 modelling studies [10, 37, 38].

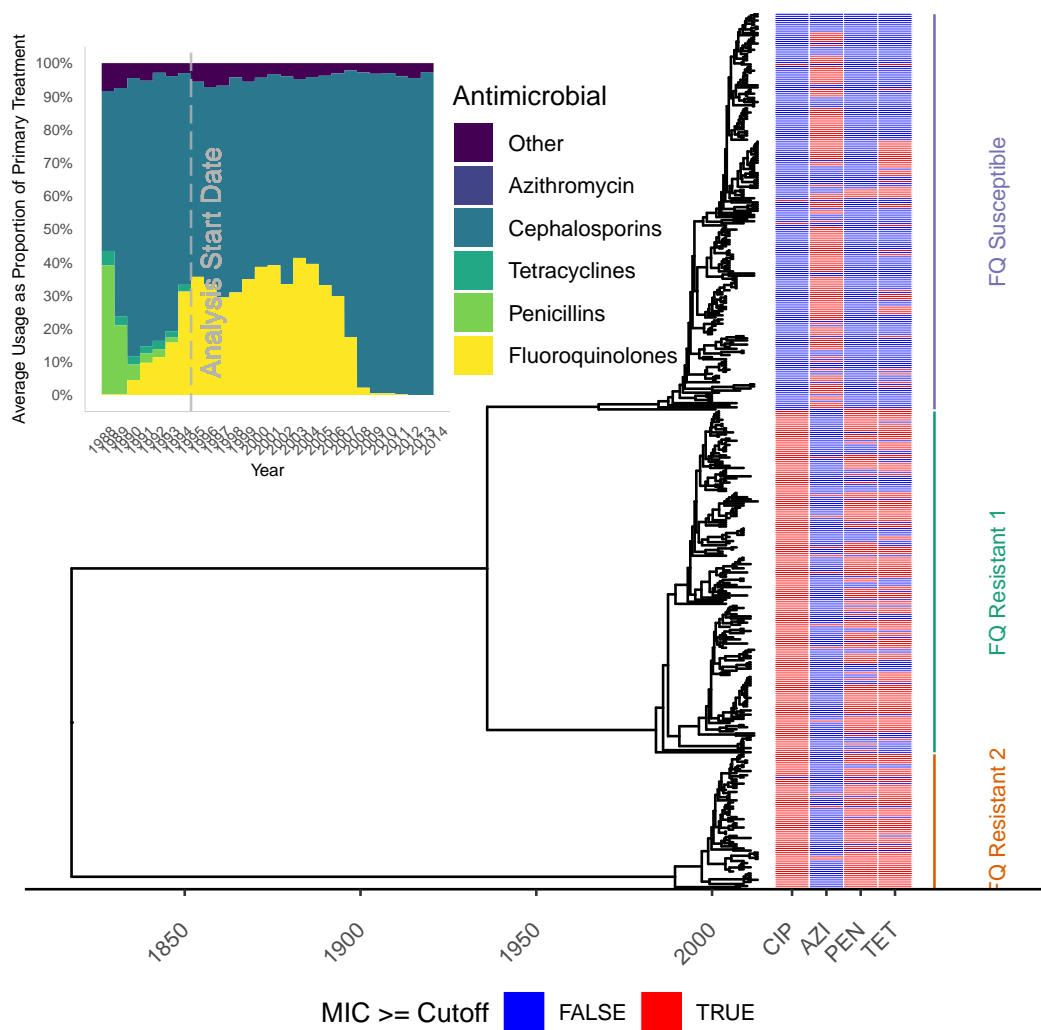


Figure 5: Antibiotic usage data and phylogeny used for the application to fluoroquinolone resistant *N. gonorrhoeae*.

218 We performed inference for this dataset; the traces are shown in Figure S5 and the posterior distribution  
 219 of kernel parameters in Figure S6. Figure 6 depicts the summary of posterior latent transmission  
 220 dynamics for the two resistant lineages, whereas Figure S7 shows the same for the susceptible lineage.  
 221 The two resistant lineages have similar dynamics, with a peak in prevalence around 2007, which  
 222 corresponds to the moment when fluoroquinolone use dropped (Figure 5). Figure 7 depicts the  
 223 marginal and joint posterior distributions for the resistance parameters  $q_U$  and  $q_T$  for both resistant  
 224 lineages. This is consistent with there being both a cost and benefit to fluoroquinolone resistance  
 225 for both lineages, since both  $q_T$  and  $q_U$  are respectively localised below 1 and above 1, with high  
 226 posterior probability. It is noteworthy that while both of these lineages come from distinct genetic  
 227 background, their resistance profile is qualitatively very similar, indicating both of these lineages faced  
 228 similar selective pressures and neither seems to have successfully adapted to overcome the fitness cost  
 229 associated with fluoroquinolone resistance. We used a posterior predictive approach to ensure that  
 230 the model can explain the data appropriately [34]. Posterior predictive trajectories for the function of

231 ancestral lineages through time  $A(t)$  were simulated and found to be very similar to the ones implied  
232 by the phylogenetic data (Figure S8).

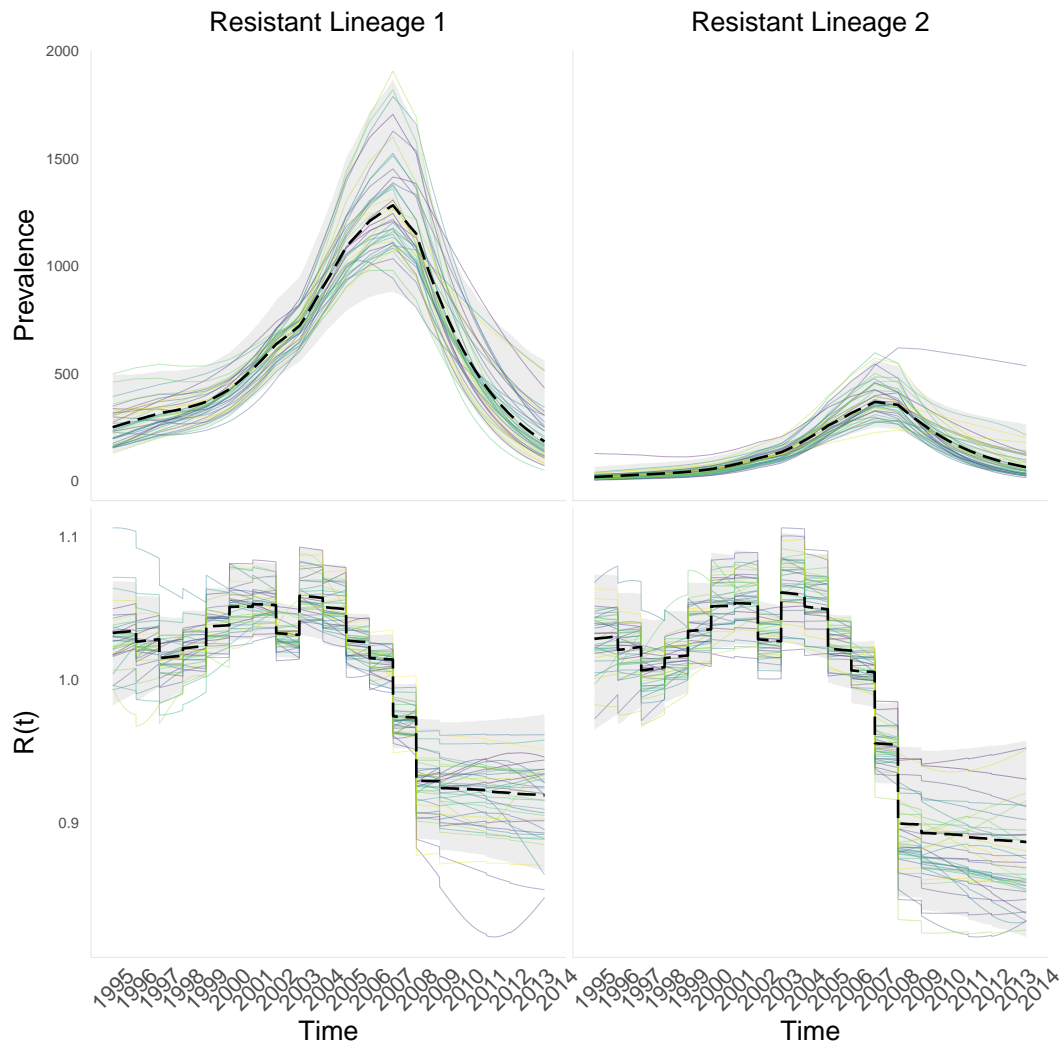


Figure 6: Posterior epidemic dynamics for both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

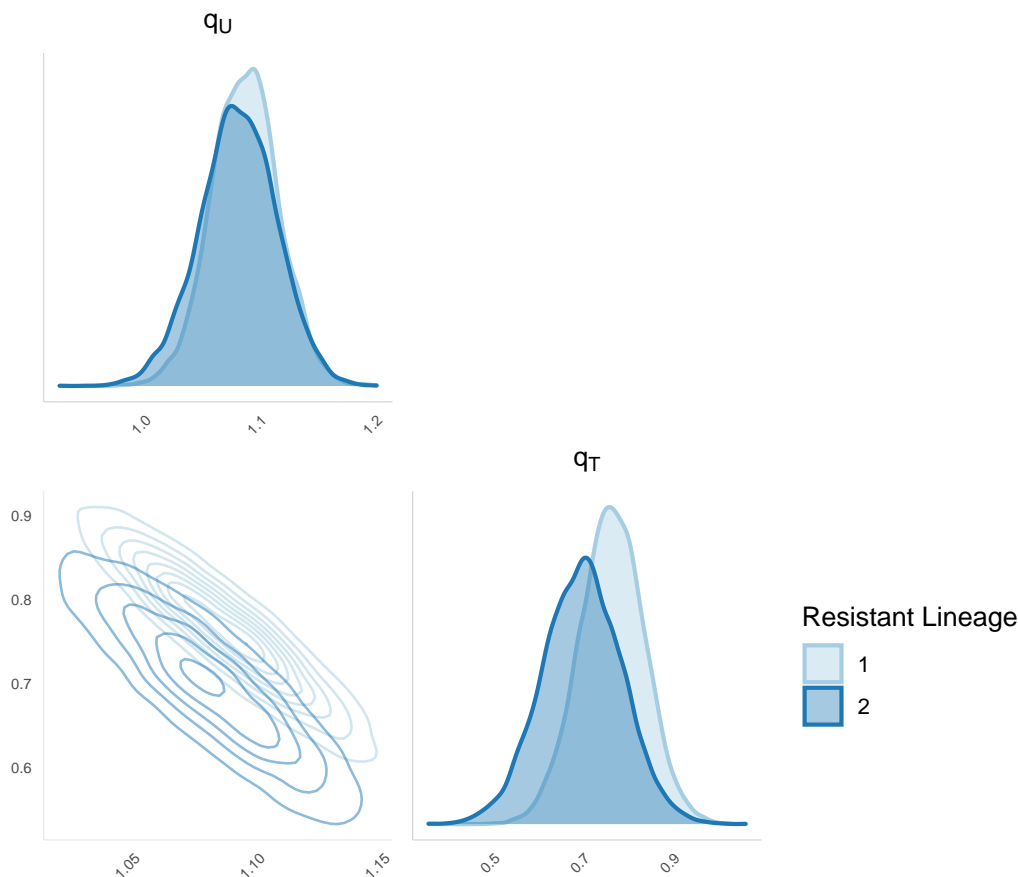


Figure 7: Marginal and joint posterior distribution for the cost ( $q_U$ ) and benefit ( $q_T$ ) of both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

233 Under the assumption of perfect between strain competition, if we want to ensure to that a resistant  
 234 strain cannot establish, and its proportion decays sufficiently fast, we fix a decay factor  $C < 1$  and  
 235 aim to ensure that the reproduction number of the resistant strain is at least  $C$  times lower than the  
 236 reproduction number of the sensitive strain, that is  $R_r(t)/R_s(t) < C$ . Given that the strains have  
 237 the same transmission rate function  $b(t)$ , this condition is equivalent to  $\gamma_r(t)/\gamma_s < C$ , and using the  
 238 definition of  $\gamma_r(t)$  from Equation 3, this is equivalent to  $u(t)q_T + (1 - u(t))q_U < C$ . We use this to  
 239 estimate posterior probabilities that the ratio of reproduction numbers is smaller than  $C$  for a given  
 240 value of fluoroquinolone usage  $u(t)$  for both resistant lineages, as shown in Figure 8. In order to be 95%  
 241 certain that the resistant lineages remain at a lower fitness than the susceptible lineage, fluoroquinolone  
 242 should not be prescribed to more than  $\sim 20\%$  and  $\sim 10\%$  of infected individuals, for resistant lineages  
 243 1 and 2, respectively.



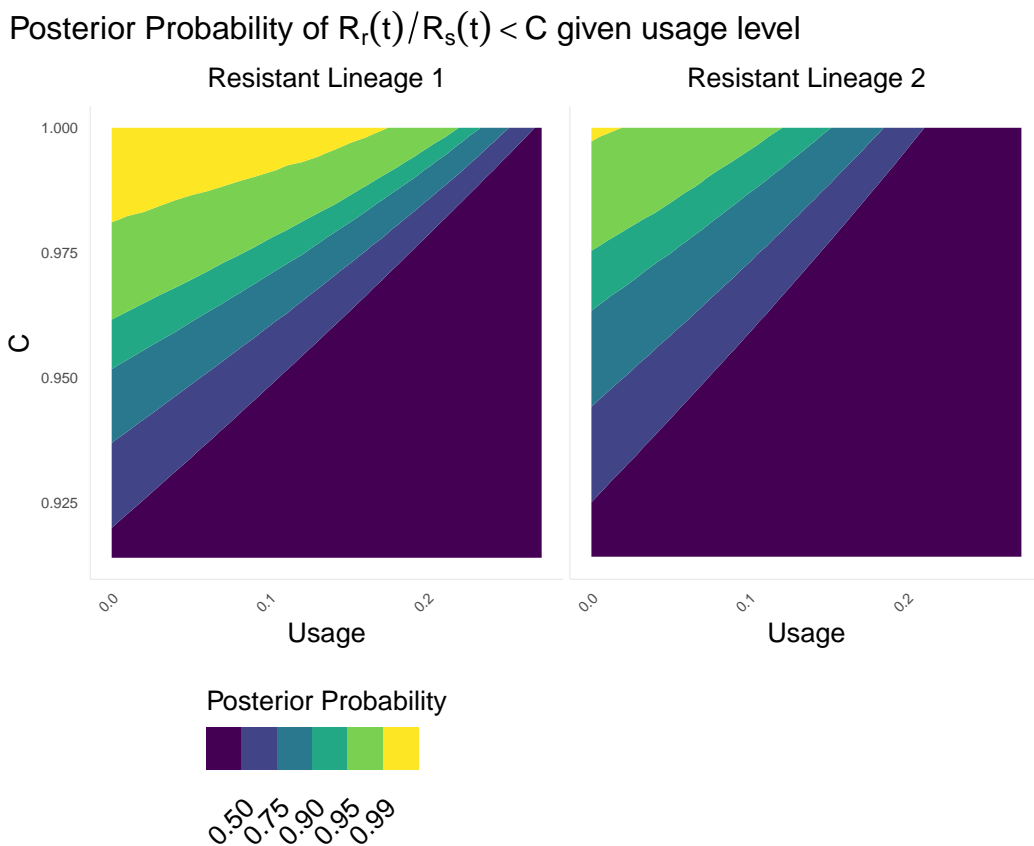


Figure 8: Posterior probabilities of  $R_r(t)/R_s(t) < C$  given usage  $u(t)$  in the x-axis and value of  $C$  in the y-axis, for both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

## 244 DISCUSSION

245 A bacterial pathogen lineage that is resistant to a given antibiotic incurs both a fitness cost and a  
246 fitness benefit compared to similar susceptible lineages [8]. When the antibiotic is used a lot, the  
247 benefit is likely to be greater than the cost so that the resistant lineage has an advantage and grows  
248 faster than susceptible lineages. Conversely, if the antibiotic is used rarely or not at all, the benefit is  
249 likely to become smaller than the cost, which will lead to the resistant lineage decreasing in frequency.  
250 Estimating these parameters is therefore of primary importance to determine how antibiotics should be  
251 prescribed without causing an increase in resistance [9]. Here, we have shown how genome sequencing  
252 data coupled with data on antibiotic prescriptions can be used for this purpose, following on previous  
253 work that demonstrated the link between epidemic dynamics and phylogenetics [12, 18, 19, 27]. By

254 comparing the phylodynamic trajectories of susceptible and resistant lineages, and relating them with  
255 a known function of antibiotic use, we show that it is possible to estimate separately the parameters  
256 corresponding to the fitness cost and benefit of resistance. In particular, we reanalysed a large published  
257 collection of *N. gonorrhoeae* genomes [35]. We were able to infer these parameters for two lineages of  
258 *N. gonorrhoeae* resistant to fluoroquinolones, and found similar estimates of cost and benefit in both  
259 (Figure 7). We were able to use this knowledge to make recommendations on antibiotic stewardship  
260 of fluoroquinolones (Figure 8).

261 Our inferential methodology is based on a well-defined and relatively simple epidemic model (Equation  
262 4) which means making a number of assumptions the validity of which was considered before performing  
263 our analysis. Our model assumes multiple-strain pathogen dynamics driven by person-to-person  
264 transmission in a well mixed host population in the absence of any significant population structure, so  
265 that there is perfect competition between strains. It also assumes that individuals become infectious  
266 as soon as they are infected, that their infectiousness remains constant until they recover, after which  
267 they become susceptible again without any immunity being gained. This list of relatively strong  
268 assumptions may seem to preclude application to any real infectious disease, but they are necessary  
269 to obtain a model under which inference can be performed. Furthermore, violation of some of these  
270 assumptions does not necessarily invalidate the results of inference. For example, if infection causes  
271 immunity, this will effectively reduce the number  $S(t)$  of susceptible individuals (Equation 2), but this  
272 number is not assumed to be constant in our model. In fact both the size  $N(t)$  of the host population  
273 and the number  $S(t)$  of susceptible individuals are integrated out as part of our parameterisation  
274 in terms of the function  $b(t)$  (cf Equation 4), so the inference is robust as long as the immunity  
275 conferred applies to all strains under study. Likewise the assumption of an unstructured population  
276 may seem problematic, including in our application to *N. gonorrhoeae* throughout the USA, but for  
277 anything other than small local outbreaks the genomes available for analysis are sparsely sampled from  
278 the whole infected population [39]. In these conditions, any effect of the host population structure  
279 on phylodynamics is likely to be insignificant as long as an effective rather than actual number of  
280 infections is considered [40, 41].

281 The compatibility of our model with the phylogenetic data under analysis can be tested using posterior  
282 predictive distribution checks (Figures S3 and S8). If these tests fail, or if the model assumptions are  
283 thought to be inappropriate, a solution may be to resort to other methods that postprocess a dated  
284 phylogeny [24] but make less assumptions, at the cost of not inferring directly the parameters of  
285 resistance. Alternative approaches includes non-parametric methods that detect differences in the  
286 branching patterns in different lineages [42, 43] as well as methods parameterised in terms of the  
287 pathogen population size growth rather than underlying epidemiological drivers [14, 44]. However, our  
288 model-based approach is both general and flexible, so that we expect it to be applicable in many settings  
289 using our software implementation which is available at <https://github.com/dhelekal/ResistPhy/>.  
290 We believe that this methodology, applied to the increasingly large genomic databases on many  
291 bacterial pathogens, will help quantify the exact link between antibiotic usage and resistance and  
292 therefore provide a much-needed evidence basis for the design of future antibiotic prescription strategies  
293 [9, 45, 46].

## 294 ACKNOWLEDGMENTS

295 We acknowledge funding from the National Institute for Health Research (NIHR) Health Protection  
296 Research Unit in Genomics and Enabling Data. This work was supported by the UK Engineering and  
297 Physical Sciences Research Council (EPSRC) grant EP/S022244/1 for the EPSRC Centre for Doctoral

## 299 References

- 300 [1] CDC, 2013 Antibiotic resistance threats in the United States, 2013. *Current* p. 114. ISSN  
301 10985530. (doi:CS239559-B).
- 302 [2] O’Neill, J., 2016 *Tackling drug-resistant infections globally: final report and recommendations*.  
303 London: Wellcome Trust & HM Government.
- 304 [3] Murray, C. J., Ikuta, K. S., Sharara, F., Swetschinski, L., Aguilar, G. R., Gray, A., Han, C.,  
305 Bisignano, C., Rao, P., Wool, E. *et al.*, 2022 Global burden of bacterial antimicrobial resistance  
306 in 2019: a systematic analysis. *The Lancet* **399**, 629–655.
- 307 [4] Clatworthy, A. E., Pierson, E. & Hung, D. T., 2007 Targeting virulence: a new paradigm for  
308 antimicrobial therapy. *Nat. Chem. Biol.* **3**, 541–548. ISSN 1552-4450. (doi:10.1038/nchembio.  
309 2007.24).
- 310 [5] Bonhoeffer, S., Lipsitch, M. & Levin, B. R., 1997 Evaluating treatment protocols to prevent  
311 antibiotic resistance. *Proc Natl Acad Sci U S A* **94**, 12106–12111. ISSN 0027-8424 (Print).  
312 (doi:10.1073/pnas.94.22.12106).
- 313 [6] Austin, D. J., Kristinsson, K. G. & Anderson, R. M., 1999 The relationship between the volume  
314 of antimicrobial consumption in human communities and the frequency of resistance. *PNAS* **96**,  
315 1152–1156. ISSN 0027-8424. (doi:10.1073/pnas.96.3.1152).
- 316 [7] Spicknall, I. H., Foxman, B., Marrs, C. F. & Eisenberg, J. N. S., 2013 A modeling framework  
317 for the evolution and spread of antibiotic resistance: Literature review and model categorization.  
318 *Am. J. Epidemiol.* **178**, 508–520. ISSN 00029262. (doi:10.1093/aje/kwt017).
- 319 [8] Andersson, D. I. & Levin, B. R., 1999 The biological cost of antibiotic resistance. *Curr. Opin.*  
320 *Microbiol.* **2**, 489–493. ISSN 13695274. (doi:10.1016/S1369-5274(99)00005-3).
- 321 [9] Andersson, D. I. & Hughes, D., 2010 Antibiotic resistance and its cost: is it possible to reverse  
322 resistance? *Nat. Rev. Microbiol.* **8**, 260–271. ISSN 1740-1526. (doi:10.1038/nrmicro2319).
- 323 [10] Whittles, L. K., White, P. J. & Didelot, X., 2017 Estimating the fitness benefit and cost of cefixime  
324 resistance in *Neisseria gonorrhoeae* to inform prescription policy: A modelling study. *PLoS Med.*  
325 **14**, e1002416. (doi:10.1371/journal.pmed.1002416).
- 326 [11] Rubin, D. H., Ma, K. C., Westervelt, K. A., Hullahalli, K., Waldor, M. K. & Grad, Y. H.,  
327 2022 Variation in supplemental carbon dioxide requirements defines lineage-specific antibiotic  
328 resistance acquisition in *neisseria gonorrhoeae*.  
329 *bioRxiv* (doi:10.1101/2022.02.24.481660). Publisher: Cold Spring Harbor Laboratory \_eprint:  
330 <https://www.biorxiv.org/content/early/2022/02/25/2022.02.24.481660.full.pdf>.
- 331 [12] Pybus, O. G. & Rambaut, A., 2009 Evolutionary analysis of the dynamics of viral infectious  
332 disease. *Nat. Rev. Genet.* **10**, 540–50. ISSN 1471-0064. (doi:10.1038/nrg2583).
- 333 [13] Volz, E. M. & Frost, S. D. W., 2014 Sampling through time and phylodynamic inference with  
334 coalescent and birth death models. *J. R. Soc. Interface* **11**, 20140945.
- 335 [14] Volz, E. M. & Didelot, X., 2018 Modeling the Growth and Decline of Pathogen Effective Population  
336 Size Provides Insight into Epidemic Dynamics and Drivers of Antimicrobial Resistance. *Syst. Biol.*  
337 **67**, 719–728. ISSN 1063-5157. (doi:10.1093/sysbio/syy007).

- 338 [15] Khnert, D., Kouyos, R., Shirreff, G., Peerska, J., Scherrer, A. U., Bni, J., Yerly, S., Klimkait, T.,  
339 Aubert, V., Gnthard, H. F. *et al.*, 2018 Quantifying the fitness cost of HIV-1 drug resistance  
340 mutations through phylodynamics. *PLOS Pathogens* **14**, e1006895. ISSN 1553-7374. (doi:  
341 10.1371/journal.ppat.1006895). Publisher: Public Library of Science.
- 342 [16] Peerska, J., Khnert, D., Meehan, C. J., Coscoll, M., de Jong, B. C., Gagneux, S. & Stadler, T.,  
343 2021 Quantifying transmission fitness costs of multi-drug resistant tuberculosis. *Epidemics* **36**,  
344 100471. ISSN 1755-4365. (doi:10.1016/j.epidem.2021.100471).
- 345 [17] Ho, S. Y. W. & Shapiro, B., 2011 Skyline-plot methods for estimating demographic history  
346 from nucleotide sequences. *Molecular Ecology Resources* **11**, 423–434. ISSN 1755098X. (doi:  
347 10.1111/j.1755-0998.2011.02988.x).
- 348 [18] Volz, E. M., Kosakovsky Pond, S. L., Ward, M. J., Leigh Brown, A. J. & Frost, S. D. W.,  
349 2009 Phylodynamics of infectious disease epidemics. *Genetics* **183**, 1421–30. ISSN 1943-2631.  
350 (doi:10.1534/genetics.109.106021).
- 351 [19] Dearlove, B. & Wilson, D., 2013 Coalescent inference for infectious disease: Meta-analysis of  
352 hepatitis C. *Philosophical Transactions of the Royal Society B* **368**, 20120314.
- 353 [20] Didelot, X., Bowden, R., Wilson, D. J., Peto, T. E. A. & Crook, D. W., 2012 Transforming  
354 clinical microbiology with bacterial genome sequencing. *Nature Reviews Genetics* **13**, 601–612.  
355 ISSN 1471-0056. (doi:10.1038/nrg3226).
- 356 [21] Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J. & Rambaut, A., 2018  
357 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol.* **4**,  
358 vey016. ISSN 2057-1577. (doi:10.1093/ve/vey016).
- 359 [22] Bouckaert, R., Vaughan, T. G., Fourment, M., Gavryushkina, A., Heled, J., Denise, K., Maio,  
360 N. D., Matschiner, M., Ogilvie, H., Plessis, L. *et al.*, 2019 BEAST 2.5 : An Advanced Software  
361 Platform for Bayesian Evolutionary Analysis. *PLoS Comput. Biol.* **15**, e1006650.
- 362 [23] Didelot, X., Croucher, N. J., Bentley, S. D., Harris, S. R. & Wilson, D. J., 2018 Bayesian inference  
363 of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Res.* **46**, e134–e134. ISSN 0305-  
364 1048. (doi:10.1093/nar/gky783).
- 365 [24] Didelot, X. & Parkhill, J., 2022 A scalable analytical approach from bacterial genomes to  
366 epidemiology. *Philosophical Transactions of the Royal Society B: Biological Sciences* **377**,  
367 20210246. (doi:10.1098/rstb.2021.0246).
- 368 [25] Allen, L. J., Kirupaharan, N. & Wilson, S. M., 2004 SIS epidemic models with multiple pathogen  
369 strains. *Journal of Difference Equations and Applications* **10**, 53–75.
- 370 [26] Keeling, M. J. & Rohani, P., 2008 *Modeling Infectious Diseases in Humans and Animals*. Princeton  
371 University Press. ISBN 9780691116174.
- 372 [27] Volz, E. M., 2012 Complex population dynamics and the coalescent under neutrality. *Genetics*  
373 **190**, 187–201. ISSN 1943-2631. (doi:10.1534/genetics.111.134627).
- 374 [28] Griffiths, R. & Tavaré, S., 1994 Sampling theory for neutral alleles in a varying environment.  
375 *Philos. Trans. R. Soc. B* **344**, 403–410.
- 376 [29] Gill, M. S., Lemey, P., Faria, N. R., Rambaut, A., Shapiro, B. & Suchard, M. A., 2013 Improving  
377 bayesian population dynamics inference: A coalescent-based model for multiple loci. *Mol. Biol.*  
378 *Evol.* **30**, 713–724. ISSN 07374038. (doi:10.1093/molbev/mss265).
- 379 [30] Rasmussen, C. E., 2004 *Gaussian Processes in Machine Learning*, pp. 63–71. Berlin, Heidelberg:  
380 Springer Berlin Heidelberg. ISBN 978-3-540-28650-9. (doi:10.1007/978-3-540-28650-9\_4).

- 381 [31] Riutort-Mayol, G., Brkner, P.-C., Andersen, M. R., Solin, A. & Vehtari, A., 2022 Practical  
382 hilbert space approximate bayesian gaussian processes for probabilistic programming. *arXiv* (doi:  
383 10.48550/arXiv.2004.11408).
- 384 [32] Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker,  
385 M. A., Guo, J., Li, P. & Riddell, A., 2017 Stan: A probabilistic programming language. *J. Stat.*  
386 *Softw.* **76**. ISSN 15487660. (doi:10.18637/jss.v076.i01).
- 387 [33] Vehtari, A., Gelman, A., Simpson, D., Carpenter, B. & Burkner, P. C., 2021 Rank-Normalization,  
388 Folding, and Localization: An Improved R hat for Assessing Convergence of MCMC. *Bayesian*  
389 *Anal.* **16**, 667–718. ISSN 19316690. (doi:10.1214/20-BA1221).
- 390 [34] Gelman, A., Meng, X. & Stern, H., 1996 Posterior predictive assessment of model fitness via  
391 realized discrepancies. *Stat Sinica* **6**, 733–807.
- 392 [35] Grad, Y. H., Harris, S. R., Kirkcaldy, R. D., Green, A. G., Marks, D. S., Bentley, S. D., Trees, D.,  
393 Lipsitch, M., Diseases, I., Health, P. *et al.*, 2016 Genomic epidemiology of gonococcal resistance  
394 to extended spectrum cephalosporins, macrolides, and fluoroquinolones in the US, 2000-2013. *J.*  
395 *Infect. Dis.* **214**, 1579–1587. ISSN 0022-1899. (doi:10.1093/infdis/jiw420).
- 396 [36] Didelot, X. & Wilson, D. J., 2015 ClonalFrameML: Efficient Inference of Recombination in  
397 Whole Bacterial Genomes. *PLoS Comput. Biol.* **11**, e1004041. ISSN 1553-7358. (doi:  
398 10.1371/journal.pcbi.1004041).
- 399 [37] Fingerhuth, S. M., Bonhoeffer, S., Low, N. & Althaus, C. L., 2016 Antibiotic-Resistant *Neisseria*  
400 *gonorrhoeae* Spread Faster with More Treatment, Not More Sexual Partners. *PLOS Pathog.* **12**,  
401 e1005611. ISSN 1553-7374. (doi:10.1371/journal.ppat.1005611).
- 402 [38] Whittles, L. K., White, P. J. & Didelot, X., 2019 A dynamic power-law sexual network  
403 model of gonorrhoea outbreaks. *PLoS Comput. Biol.* **15**, e1006748. ISSN 1553-7358. (doi:  
404 10.1371/journal.pcbi.1006748).
- 405 [39] Klinkenberg, D., Colijn, C. & Didelot, X., 2019 Methods for Outbreaks Using Genomic Data. In  
406 *Handbook of Infectious Disease Data Analysis*, pp. 245–263. CRC Press.
- 407 [40] Nordborg, M., 1997 Structured coalescent processes on different time scales. *Genetics* **146**, 1501–  
408 1514. ISSN 0016-6731.
- 409 [41] Frost, S. D. W. & Volz, E. M., 2010 Viral phylodynamics and the search for an 'effective number  
410 of infections'. *Philosophical Transactions of the Royal Society B* **365**, 1879–1890. ISSN 0962-8436.  
411 (doi:10.1098/rstb.2010.0060).
- 412 [42] Dearlove, B. L., Xiang, F. & Frost, S. D. W., 2017 Biased phylodynamic inferences from analysing  
413 clusters of viral sequences. *Virus Evolution* **3**, 1–10. (doi:10.1093/ve/vex020).
- 414 [43] Volz, E. M., Wiuf, C., Grad, Y. H., Frost, S. D. W., Dennis, A. M. & Didelot, X., 2020  
415 Identification of hidden population structure in time-scaled phylogenies. *Systematic Biology* **69**,  
416 884–896. (doi:10.1093/sysbio/syaa009).
- 417 [44] Helekal, D., Ledda, A., Volz, E., Wyllie, D. & Didelot, X., 2021 Bayesian inference of clonal  
418 expansions in a dated phylogeny. *Systematic Biology* p. syab095. ISSN 1063-5157. (doi:  
419 10.1093/sysbio/syab095).
- 420 [45] Michael, C. A., Dominey-Howes, D. & Labbate, M., 2014 The Antimicrobial Resistance Crisis:  
421 Causes, Consequences, and Management. *Frontiers in Public Health* **2**. ISSN 2296-2565. (doi:  
422 10.3389/fpubh.2014.00145).

- 423 [46] Holmes, A. H., Moore, L. S., Sundsfjord, A., Steinbakk, M., Regmi, S., Karkey, A., Guerin, P. J.  
424 & Piddock, L. J., 2016 Understanding the mechanisms and drivers of antimicrobial resistance.  
425 *Lancet* **387**, 176–187. ISSN 1474547X. (doi:10.1016/S0140-6736(15)00473-0).