1 **Inferring demographic parameters in bacterial genomic data using Bayesian and hybrid phylogenetic methods**
2
3 Sebastian Duchene[1]*, David A Duchene[2], Jemma L Geoghegan[3], Zoe A Dyson[1], Jane Hawkey[1], Kathryn E Holt[1]
4
5 [1] Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, University of Melbourne, Parkville,
6 Victoria 3020, Australia
7
8 [2] School of Life and Environmental Sciences, University of Sydney, Sydney, NSW 2006, Australia
9
10 [3] Department of Biological Sciences, Macquarie University, Sydney, NSW 2109, Australia
11
12 * To whom correspondence should be addressed
13
14
15 Sebastian Duchene
16 sebastian.duchene@unimelb.edu.au
17
18 David A Duchene
19 david.duchene@sydney.edu.au
20
21 Jemma L Geoghegan
22 jemma Geoghegan@mq.edu.au
23
24 Zoe A Dyson
25 zoe.dyson@unimelb.edu.au
26
27 Jane Hawkey
28 jane.hawkey@unimelb.edu.au
29
30 Kathryn E Holt
31 kholt@unimelb.edu.au
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

**Abstract**

**Background:** Recent developments in sequencing technologies make it possible to obtain genome sequences from a large number of isolates in a very short time. Bayesian phylogenetic approaches can take advantage of these data by simultaneously inferring the phylogenetic tree, evolutionary timescale, and demographic parameters (such as population growth rates), while naturally integrating uncertainty in all parameters. Despite their desirable properties, Bayesian approaches can be computationally intensive, hindering their use for outbreak investigations involving genome data for a large numbers of pathogen isolates. An alternative to using full Bayesian inference is to use a hybrid approach, where the phylogenetic tree and evolutionary timescale are estimated first using maximum likelihood. Under this hybrid approach, demographic parameters are inferred from estimated trees instead of the sequence data, using maximum likelihood, Bayesian inference, or approximate Bayesian computation. This can vastly reduce the computational burden, but has the disadvantage of ignoring the uncertainty in the phylogenetic tree and evolutionary timescale.

**Results:** We compared the performance of a fully Bayesian and a hybrid method by analysing six whole-genome SNP data sets from a range of bacteria and simulations. The estimates from the two methods were very similar, suggesting that the hybrid method is a valid alternative for very large datasets. However, we also found that congruence between these methods is contingent on the presence of strong temporal structure in the data (i.e. clocklike behaviour), which is typically verified using a date-randomisation test in a Bayesian framework. To reduce the computational burden of this Bayesian test we implemented a date-randomisation test using a rapid maximum likelihood method, which has similar performance to its Bayesian counterpart.

**Conclusions:** Hybrid approaches can produce reliable inferences of evolutionary timescales and phylodynamic parameters in a fraction of the time required for fully Bayesian analyses. As such, they are a valuable alternative in outbreak studies involving a large number of isolates.

**Keywords**

Bayesian phylogenetics, phylodynamics, molecular clock, bacterial evolution

**Background**

Genomic data are increasingly used to investigate infectious disease outbreaks caused by microbial pathogens. Recent developments in sequencing technologies have made it possible to obtain data for a very large number of samples, at low cost and within a very short timeframe. Phylogenetic methods can make use of these data to infer their evolutionary dynamics, known as phylodynamic inference. For example, genome data obtained during the first months of the 2013-2016 Ebola virus epidemic were used to determine the time of origin of the outbreak and the basic reproductive number ($R_0$) of the circulating strains [1,2]. Some of the key requirements for these inferences are that the data must have sufficient genetic diversity and that they should be a representative sample of the circulating strains.

87 Serially sampled data are particularly useful because their sampling times can be used to calibrate the molecular
88 clock. This consists of calculating the rate of evolution, which is the amount of genetic change that has accumulated
89 per unit of time. The rate of evolution is key to infer an evolutionary timescale, typically represented by a
90 phylogenetic tree where the branch lengths correspond to time, known as a chronogram. Some methods assume
91 that the rate of evolution is constant over time, known as a strict molecular clock, but popular Bayesian
92 implementations, such as that in BEAST [3,4], include relaxed-clock models that use a statistical distribution to
93 describe rate variation across time and lineages (reviewed in [5]). Phylodynamic models can be used to estimate the
94 epidemic growth rate ($r$), $R_o$, and other parameters [6,7]. Importantly, these models describe the expectation of the
95 distribution of node times in the chronogram. As such, inferences drawn from phylodynamic models rely on accurate
96 estimates of evolutionary rates and timescales. A number of statistical methods are available to assess the
97 robustness of inferences of evolutionary rates and timescales; those that are most widely used are implemented
98 under a Bayesian framework (reviewed in [8]).
99
100 Bayesian phylogenetic approaches allow sophisticated evolutionary models to be specified. For example, the
101 evolution of a pathogen during an outbreak can be defined as an exponentially growing population with considerable
102 evolutionary rate variation among lineages; which can be modelled by specifying a nucleotide substitution model, a
103 relaxed-clock model and an exponential-growth tree prior. The parameters for all these models are obtained
104 simultaneously and their estimates correspond to posterior probability distributions, such that their uncertainty is a
105 natural by-product of the analysis. Bayesian methods require specifying a prior distribution for all parameters.
106 Although specifying a prior distribution is not trivial for some parameters, their influence can be assessed by
107 comparing them to the posterior. An advantage of specifying prior distributions is that it is possible to include
108 previous knowledge about the data. As a case in point, a known probability of sampling can be represented with a
109 prior distribution in birth-death models [9].
110
111 Whilst Bayesian phylogenetic methods have many desirable properties, analysing large genomic data sets under
112 complex models is often computationally prohibitive (e.g. [10,11]). An alternative to full Bayesian methods is to
113 conduct the analysis in several steps. In this hybrid approach the phylogenetic tree, evolutionary rates and
114 timescales, and demographic parameters are estimated separately.
115
116 Phylogenetic trees can be rapidly estimated using various maximum likelihood implementations [12–15]. These
117 methods assume a substitution model, but not a molecular-clock or demographic model, such that the branch
118 lengths of the trees represent the expected number of substitutions per site, and are known as phylograms.
119
120 Next, phylograms can be used to estimate evolutionary rates and chronograms, for example, using a recently
121 developed molecular clock method based on least-squares optimisation, called LSD (Least Squares Dating) [16]. LSD
122 is more computationally tractable than Bayesian molecular-clock methods, such that it is feasible to analyse genomic
123 data sets with thousands of samples. Although LSD assumes a strict molecular clock, its accuracy is frequently similar
124 to that obtained using more sophisticated Bayesian clock models [17]. Other non-Bayesian molecular-clock methods
125 have also been developed recently with the purpose of analysing large genomic data sets [18–20].
126

127 Finally, a range of tools are available to infer phylodynamic parameters from a chronogram, such as that obtained
128 using LSD. For example: TreePar uses maximum likelihood to fit birth-death and skyline models [21]; BEAST2 [4] and
129 RevBayes [22] can fit a range of birth-death, coalescent, and Skyline models using Bayesian inference [7]; and
130 approximate Bayesian computation (ABC) approaches that use tree summary statistics have recently been
131 developed to fit phylogenetic epidemiological models [23,24]. The main disadvantage of these approaches over
132 those that are fully Bayesian is that the estimates are based on a single tree, such that uncertainties in tree topology,
133 branch lengths, and evolutionary rates are ignored. A potential solution is to repeat the analysis using non-
134 parametric bootstrap replicates, but combining the different sources of uncertainty under this framework is not
135 trivial.
136
137 Here, we compare the following two methods to infer evolutionary rates and timescales, and demographic
138 parameters:
139     (i)        The fully Bayesian method, implemented in BEAST2, to simultaneously infer the phylogenetic tree,
140                 evolutionary timescales and phylodynamic parameters;
141     (ii)       The hybrid method: phylogram inference using maximum likelihood in PhyML v3.1 [14], chronogram
142                 inference using LSD v0.3, and estimation of phylodynamics parameters in BEAST2 using Bayesian
143                 inference.
144
145 To compare the performance of these two methods, we analysed previously published whole genome SNP bacterial
146 data sets of *Mycobacterium tuberculosis* Lineage 2 [25], *Vibrio cholerae* [26], *Shigella dysenteriae* type 1 [11], and
147 *Staphylococcus aureus* ST239 [27]. Because these data sets have small numbers of samples (n=63 for *M. tuberculosis*,
148 n=122 for *V. cholerae*, n=121 for *S. dysenteriae*, and n=74 for *S. aureus*) their analyses are computationally tractable
149 using both approaches. We also demonstrate the unique potential of the hybrid approach by analysing two genomic
150 data sets with larger numbers of sequences, which have been difficult to analyse using a fully Bayesian approach; a
151 global sample of *S. dysenteriae* type 1 (n=329) and *S. dysenteriae* type 1 lineage IV (n = 208) [11]. Finally, we validated
152 the performance of the hybrid approach using a simulation experiment.
153
154 **Results**
155
156 *Estimates of evolutionary rates and timescales*
157
158 We compared estimates of rates and evolutionary timescales using the full Bayesian approach in BEAST2 and LSD.
159 Because our data consist of SNPs, we used ascertainment bias correction by specifying the number of constant sites
160 from the core genome. In BEAST2 we used both the strict and the uncorrelated lognormal (UCLN [28]) clock models.
161 We investigated the degree of rate variation among lineages by inspecting the coefficient of rate variation, estimated
162 in the UCLN model. This parameter is the standard deviation of branch rates divided by the mean rate. The data are
163 considered to display clocklike behaviour if the distribution for this parameter abuts zero. Therefore, we used this
164 parameter to select the clock model in BEAST2 for each data set, as suggested in previous studies [29,30]. The *M.*
165 *tuberculosis* data set was the only data set to support a strict clock over the UCLN model, whereas the remaining data
166 sets favoured the UCLN model (Fig.1). We set uniform prior distributions for the clock rate, the growth rate (*r*) and

167 the scaled population size ($\Phi$). In the context of pathogen evolution, $r$ determines the speed of spread of the

168 pathogen in the host population, while $\Phi$ is proportional to the infected host population size at present.

169

170 The estimates of evolutionary rates and timescales from these different methods were largely congruent (Fig.1). In all

171 four cases, the 95% credible intervals for the evolutionary rate and age of the root node obtained with BEAST2

172 overlapped with the 95% confidence intervals obtained for the same parameters with LSD (Fig.1). However, we

173 observed some differences in the mean evolutionary rate estimates, with the estimates from BEAST2 consistently

174 producing higher values than those from LSD. The largest difference in mean rate estimates was observed in *M.*

175 *tuberculosis*, with a mean rate of $9.37 \times 10^{-8}$ (95% credible interval: $4.25 \times 10^{-8} - 1.73 \times 10^{-7}$) using BEAST2, and $1.10 \times 10^{-8}$

176 (95% confidence interval: $1.00 \times 10^{-10} - 2.02 \times 10^{-7}$) in LSD (see Fig.1). In contrast we found more congruent mean rate

177 estimates in the *V. cholerae* data set, with estimates of $7.20 \times 10^{-7}$ (95% credible interval: $5.87 \times 10^{-7} - 8.65 \times 10^{-7}$) for the

178 BEAST2 and $6.76 \times 10^{-7}$ (95% confidence interval: $5.76 \times 10^{-7} - 8.89 \times 10^{-7}$) for LSD. The differences in estimates of the

179 root-node age were similar, with the largest difference in the mean root-node age found in *S. aureus* ST239 (mean

180 root-node age of 1958 for BEAST2 and 1949 for LSD) (Fig.1). In most cases, the estimates from BEAST2 were more

181 uncertain with credible intervals that were wider than the confidence intervals from LSD. We investigated two

182 aspects of phylogenetic data that can affect estimates of evolutionary rates; the topological uncertainty and the

183 degree of clocklike variation. We found that the maximum likelihood trees were highly supported, according to local

184 likelihood ratio tests (aLRT) [31] (which ranges from 0 to 1, for low to high branch support, respectively). The median

185 aLRT values across nodes were 0.9 for *M. tuberculosis*, 0.83 for *V. cholerae*, 0.99 for *S. dysenteriae* type 1, and 0.92 for

186 *S. aureus*.

187

188 *Assessing temporal structure using a date-randomisation test*

189

190 We assessed the reliability of our estimates of evolutionary rate and timescales by conducting a date-randomisation

191 test [32,33]. The motivation of this test is similar to that of root-to-tip regressions implemented in TempEst [34]. That

192 is, to determine whether there is sufficient sampling in the data. However, root-to-tip regressions should be

193 interpreted for visual inspection, as opposed to date-randomisations, which are a formal statistical test. The date

194 randomisation test consists in repeating the analysis several times after randomising the sampling dates. The

195 resulting rate estimates correspond to the expected values if there is no association between sampling times and

196 genetic divergence. The data are considered to have strong temporal structure if the rate estimate obtained using

197 the correct sampling times is not contained within the range of values from the randomisations. In a Bayesian

198 context, 10 to 20 randomisations appear to be sufficient [33,35]. We conducted this test in BEAST2 using 20

199 randomisations and in LSD using 100 randomisations (Fig.2). Interestingly, the results from both tests were

200 congruent, and consistent with visualisations of clock-like behaviour of the data using root-to-tip regressions

201 (Fig.S1). The *M. tuberculosis* data set had no temporal structure with either method (Fig.2): the credible interval of the

202 Bayesian estimate with the correct sampling times overlapped with those from all of the randomisations; using LSD,

203 the estimate with the correct sampling times was around the lower threshold in the program, at $1.00 \times 10^{-10}$

204 subs/site/year, which also corresponds to the value obtained for most of the randomisations. The other data sets

205 showed strong temporal structure with both date-randomisation tests: the Bayesian credible intervals using the

206 correct sampling times did not overlap with those from any of the randomisations, and the estimates from LSD using
207 the correct sampling times were not contained within the distributions of the 100 randomisations (Fig. 2).
208
209 *Inference of phylodynamic parameters*
210
211 We analysed the data sets using the exponential-growth coalescent model in BEAST2, which has two parameters, $r$
212 and $\Phi$. Because these are compound parameters, they cannot be interpreted in an absolute scale without additional
213 information about the size of the infected host population at present [36]. In most cases, the posterior distributions
214 of both parameters were very similar when using either BEAST2 or the hybrid approach, with similar means and
215 uncertainties (Fig.3). Although the intervals overlapped in *V. cholerae*, *S. dysenteriae*, and *S. aureus*, the mode of the
216 posterior distribution of $\Phi$ was higher when using the hybrid approach. The posterior distributions of $r$ were almost
217 identical across methods for the three data sets with temporal signal (Fig.3). The uncertainty in estimates of this
218 parameter did not include 0, except in the case of *V. cholerae*, suggesting that most of these bacterial data sets were
219 undergoing population growth. Interestingly, the *M. tuberculosis* data set, which had no temporal structure, was the
220 only data set to display large differences in estimates among the methods (Fig.3).
221
222 *Application: analysing large data sets using the hybrid approach*
223
224 Having demonstrated good performance of the hybrid approach on small data sets with strong temporal signal, we
225 applied it to analyse two published genome-wide SNP data sets whose sample size was prohibitively large to analyse
226 under a full Bayesian framework in the original publication. These data sets consisted of: (i) 329 samples of *S.*
227 *dysenteriae* type 1 from [11], which included BEAST2 analysis of a subset of 125 samples; and (ii) 208 samples of
228 lineage IV of *S. dysenteriae* type 1, which was represented by 61 samples in the BEAST2 analysis in the same study
229 [11]. These three data sets displayed strong temporal structure according to the date-randomisation test in LSD, with
230 rate estimates that were not contained within the range of estimates from 100 date-randomisations (Fig.4). The
231 evolutionary rate estimates from LSD were $5.93 \times 10^{-7}$ (95% confidence interval: $3.65 \times 10^{-7}$ - $1.65 \times 10^{-6}$) subs/site/year for
232 *S. dysenteriae* type 1, and $7.04 \times 10^{-7}$ (95% confidence interval: $3.92 \times 10^{-7}$ - $1.54 \times 10^{-6}$) subs/site/year for *S. dysenteriae*
233 type 1 Lineage IV (Fig.4). Interestingly, the estimate of $r$ for *S. dysenteriae* type 1 lineage IV was over an order of
234 magnitude higher than that for the global data set of this bacterium, with a mean of $2.00 \times 10^{-2}$ for lineage IV
235 compared with $3.40 \times 10^{-3}$ for the global data set. Importantly, the posterior distributions of $r$ for these three data sets
236 did not include zero, indicating epidemic growth (Fig.4).
237
238 *Validation using simulations*
239
240 Although our empirical analyses suggest that the hybrid and the full Bayesian method can produce largely congruent
241 results, it is unclear whether the methods are accurate. That is, whether they can recover the true parameter
242 estimates. To investigate this, we conducted a simulation experiment. We simulated 100 whole genome data sets
243 using similar parameters to those we inferred for our *S. dysenteriae* data set. We extracted the SNPs from the
244 synthetic genomes and analysed them using the hybrid and full Bayesian methods, with the same settings that we
245 used for the empirical data. Our date-randomisations in LSD indicated that all of these data sets had temporal

6

246     structure, with *p*-values of 0.00. The estimates for the age of the root-node from both methods were very similar.

247     However, it is important to note that our hybrid method uses a single tree, such that the age of the root-node is a

248     point value, whereas the full Bayesian analyses include uncertainty in this parameter. Accordingly, the estimates

249     from LSD were very close to those used to generate the data (within 5 years of the true value), and those from the full

250     Bayesian method always included the true value within their credible interval. The estimates for the demographic

251     parameters, *r* and *Φ*, had credible intervals that always included the true value for both methods, with mean values

252     that often matched those used to generate the data (Fig. 5a). Interestingly, in 10 randomly selected simulation

253     replicates, we found that the credible intervals for the demographic parameters were very similar for both methods,

254     with the hybrid approach sometimes producing more precise estimates. We found no estimation biases in any of the

255     methods (Fig. 5a).

256

257     We conducted a second set of simulations of data with no temporal structure. To do this, we generated similar

258     sequence alignments as described above, but we assigned random sampling times in our analyses in LSD and in

259     BEAST2. This means that the molecular clock calibration is effectively uninformative. The age of the root-node was

260     over estimated by both methods. In LSD this bias was of over three orders of magnitude, whereas in BEAST2 it

261     ranged between half and three orders of magnitude. The value of *Φ* was similarly overestimated in both methods.

262     The growth rate, *r*, was underestimated by several orders of magnitude with the hybrid approach, but it tended to be

263     overestimated with the full Bayesian method (Fig. 5b). A key result about the simulations with no temporal structure

264     is that *Φ* was always incorrectly estimated, and the true value of *r* was only contained within the 95% credible interval

265     in about 14% of the analyses using the full Bayesian method. Moreover, the estimates with the hybrid approach often

266     displayed larger discrepancies with the correct values.

267

268     *Computational demands of the Bayesian and the hybrid methods*

269

270     The hybrid approach was several times faster than the full Bayesian approach. For example, the computation time for

271     each randomisation of the *V. cholerae* data set each was about 2 hours using BEAST2, where as those in LSD took

272     1.23 seconds (sec). However, a key aspect of the date-randomisation test in LSD is that the tree topology and branch

273     lengths are fixed for all randomisations, where as they are re-estimated for each randomisation in BEAST2. For the *V.*

274     *cholerae* data set, a complete analysis using the hybrid approach took: 10.06 minutes (min) to infer a maximum

275     likelihood tree in PhyML, 1.23 sec to estimate the evolutionary rate and timescale in LSD, and 5 min to infer *r* and *Φ* in

276     BEAST2 to obtain effective sample sizes (ESS) of over 200 for all parameters (drawing $1 \times 10^7$ steps, with 1 minutes per

277     $10^6$ steps), for a total of about 15 min, and $1/12^{th}$ of the time required in BEAST2. Analysis of the full *S. dysenteriae*

278     dataset from [11], the largest data set in our study, took 10.6 sec to analyse in LSD and 1 hour infer *r* and *Φ* BEAST2

279     (drawing $5 \times 10^7$ steps, with 1.2 minutes per $10^6$ steps), for the 329 sampled sequences.

280

281     **Discussion**

282

283     Our results demonstrate that, as long as a strong temporal signal is present, the hybrid and fully Bayesian methods

284     can produce congruent estimates of evolutionary parameters, even in cases where the data display substantial rate

285     variation among lineages. These methods also yielded similar estimates of demographic parameters in data sets with

7

286  strong temporal signal, indicating the hybrid approach is a reliable alternative to full Bayesian analyses. However, $r$

287  appears to be more robust than $\Phi$ to mild differences in estimates of the rate and timescale. This probably occurs

288  because the age of the root-node plays an important role in the population size under the coalescent. In particular,

289  the effective population size, and therefore $\Phi$, are known to scale positively with the age of the root-node [37].

290

291  Obtaining congruent estimates between the two methods depends on whether the data meet certain criteria. In

292  practice, it is important to verify that the trees have high branch support and that the data have strong temporal

293  structure. The trees inferred here were highly supported, but it is likely that the hybrid approach will produce

294  misleadingly precise estimates (i.e. with narrow confidence intervals) if branch support is low, because the

295  demographic parameters will still be conditioned on a single, and possibly incorrect, tree obtained in step 1 that does

296  not capture uncertainty in the topology. In contrast, in such circumstances the Bayesian method will simply integrate

297  over phylogenetic uncertainty and yield wider credible intervals. Our simulations illustrate ideal conditions, in which

298  the data evolve under the correct model and have strong temporal structure. In this case, we find that both methods

299  produce accurate estimates with similar precision.

300

301  Our simulations of data with no temporal structure demonstrate, not only that the hybrid and full Bayesian methods

302  will produce different estimates, but that they both tend to be inaccurate. In the absence of temporal structure, LSD

303  often produces rate estimates at the lower threshold of the program, which was $10^{-10}$ here. This means that the

304  timescale of the chronogram is overestimated. The value of $\Phi$ is also overestimated, which occurs because this

305  parameter scales positively with the age of the root-node [37]. Although, we found that $r$ was also overestimated, this

306  parameter is determined by the distribution of branches in the tree, such that its error is less predictable. The full

307  Bayesian method produced estimates with smaller bias. We used uniform priors for $\Phi$ and $r$, and the prior for the age

308  of the root was determined by the coalescent prior. It is likely that these parameters, especially $\Phi$, will be affected by

309  different choice of priors. For empirical data with low temporal structure, the hybrid approach will likely be

310  misleading because it is conditioned on a single tree which is probably incorrect. In such cases, it may be necessary to

311  use the full Bayesian method approach because it is possible to include sources of molecular clock calibration via prior

312  parametric distributions, at the expense of much higher computational demands. For instance, a reasonable

313  calibration on the age of the root-node might be sufficient to overcome low temporal structure and to obtain reliable

314  estimates for $\Phi$ and $r$. To investigate this, it is important to verify that there exists a difference between the prior and

315  posterior for parameters of interest (see Boskova et al. [38] for an investigation of the prior and posterior in Bayesian

316  phylodynamics).

317

318  Our results show that the date-randomisation test in LSD appears to be as effective as it is in BEAST2, with the

319  advantage of being much less computationally demanding. As a result, it is possible to use a larger number of

320  replicates, which can improve the power of the test. Moreover, the sampling times under a Bayesian analysis of

321  sequentially sampled data are informative about the tree topology. That is, they impose a high prior probability on

322  trees that cluster sequences with similar sampling times, which can render the date-randomisation test unreliable,

323  with an increase in type I error [39]. Moreover, in some phylodynamic models, the estimate of the age of the root-

324  node and the evolutionary rate are determined by a combination of the sequence data and their sampling times [38],

325  such that assessing temporal structure via the date randomisation test is not trivial. The date-randomisation test in

326    LSD does not suffer from these problems because sequence data alone, not tip dates, are used to infer the tree
327    topology in maximum likelihood.
328
329    Critically, the rates estimated using the date-randomisation in test in LSD are not necessarily unimodal in their
330    distribution. This occurs because a lack of temporal structure usually leads to very low rate estimates, which affects
331    randomisations in LSD and in BEAST2. In the case of LSD, very low values for the rate will correspond to the lower
332    threshold set in the program [17], which we arbitrarily set at $10^{-10}$ subs/site/year, such that most randomisations will
333    have this value. As such, a reasonable approach to interpret the date-randomisation test in LSD is to ensure that the
334    rate estimate with the correct sampling times is higher than those from at least 95% of the randomisations, following
335    the frequentist one-tailed $p$-value of $\alpha=0.05$.
336
337    **Conclusions**
338
339    As shown here, hybrid methods offer an attractive alternative to full Bayesian approaches for genome-scale data sets
340    with very large numbers of samples. The accuracy and precision of both methods are comparable, but hybrid
341    methods can perform an analysis in a about an eighth of the time required for full Bayesian analyses. Nevertheless,
342    some steps of the hybrid method used here require oversimplifications of the evolutionary process. For example, LSD
343    always assumes a strict molecular clock, such that it is impossible to assess among-lineage rate variation or to
344    pinpoint potential biological causes for why lineages have different rates. The choice of whether to use a hybrid
345    method should be made based on what parameters a user wishes to interrogate. In the context of molecular
346    epidemiology, demographic parameters ($r$ and $\Phi$) and divergence time information are of primary interest, all of
347    which appear robust to some among-lineage rate variation.
348
349    In this study, we used a simple demographic model, the exponential-growth coalescent. This model appears to be
350    well suited when outbreak data are sampled at an early stage, but it makes several assumptions, including that the
351    population of susceptible hosts is constant and that there is no population structure [6]. A better understanding of
352    the data used here requires more sophisticated phylodynamic models, such as those that include changes in
353    diversification parameters over time [40], and migration [41]. To this end, our results suggest that harnessing the
354    power of such models and large-scale genome sequencing can be done through hybrid approaches.
355
356    **Materials and Methods**
357
358    *Data collection*
359
360    Our bacterial data sets consisted of publically available genome data. We obtained all of our genome-wide SNP
361    alignments from a previous studies [11, 25, 27, 35]. These data sets are freely available online
362    (github.com/sebastianduchene/bacteria_genomic_rates_data). These data have had regions with evidence of
363    recombination removed using Gubbins v2 [42].
364
365    *Phylogenetic analyses under the fully Bayesian approach*

366

367  We analysed the sequence alignments in BEAST v2.4 using the sampling times for calibration, the GTR+Γ

368  substitution model, the exponential-growth coalescent tree prior, and two clock models; the strict and the UCLN. We

369  used the default priors for all parameters. Our Markov chain Monte Carlo (MCMC) sampling scheme consisted of a

370  chain length of $5 \times 10^8$ steps, sampling every $10^4$ steps. We verified that the ESS for all parameters was at least 200. To

371  determine whether the data had temporal structure, we conducted a date-randomisation test by randomising the

372  sampling dates 20 times and repeating the analyses [33].

373

374  *Phylogenetic analyses using the hybrid approach*

375

376  We inferred phylogenetic trees using maximum likelihood in PhyML v3.1. We used the GTR+Γ substitution model,

377  and a search strategy that combines the nearest-neighbour interchange and subtree prune and regraft algorithms.

378  To assess branch support, we calculated the aLRT score for each branch. To visually assess temporal structure, we

379  conducted a regression of the root-to-tip distances as a function of the sampling times using TempEst v1.5 [34]. To

380  determine the optimal root in this program we selected the position that maximised $R^2$.

381

382  We analysed the maximum likelihood trees (i.e. phylograms) in LSD v0.3 to infer the evolutionary rate and timescale.

383  We set the sampling times as calibrations and allowed the program to determine the optimal position of the root. We

384  constrained the branching times of the estimated chronograms such that daughter nodes must be younger than their

385  parent nodes. To obtain an uncertainty around estimates of times and rates, we conducted 100 parametric bootstrap

386  replicates of the branch lengths, as implemented in the program. Therefore, the uncertainty corresponds to the 95%

387  confidence interval of the parametric bootstrap values. We conducted a date-randomisation test 100 times by

388  randomising the sampling times in the 'date' file and running LSD each time. In this version of the test, the

389  phylogenetic tree topology and branch lengths are fixed.

390

391  We used the chronograms estimated in LSD to infer demographic parameters in BEAST2. This consists in setting the

392  input file to calculate the posterior as the likelihood of the tree given the model parameters multiplied by the priors

393  on the parameters. In the exponential growth coalescent there are two parameters; $\Phi$ and $r$. We used an MCMC chain

394  length of $1 \times 10^7$ sampling every $10^4$ steps, and we verified that all parameters had ESS values of at least 200.

395

396  *Simulations*

397  We simulated whole genome sequence alignments using the parameters from our *S. dysenteriae* data set. To do this,

398  we took the highest clade credibility tree from this data set inferred in BEAST2 and simulated the evolutionary rate

399  using NELSI [29], according to an UCLN clock model. We used a mean rate of $10^{-6}$ subs/site/year and a standard

400  deviation of $10^{-7}$. We used Seq-Gen v1.3 [43] to simulate genome sequence alignments of 3,750,125 nucleotides using

401  the GTR+Γ substitution model with the mean parameter estimates for the empirical *S. dysenteriae* data. Finally, we

402  extracted the SNPs from these alignments and analysed using the same method as for our empirical data. For our

403  simulations with no temporal structure we set random sampling times for our analyses in LSD and BEAST2. In all

404  cases, we conducted a date-randomisation test in LSD, as used in our empirical data analysis.

405

406 **Abbreviations**

407 LSD, Least-squares dating; ABC, Approximate Bayesian Computation; UCLN, uncorrelated lognormal clock; aLRT,

408 local Likelihood ratio test for branch support; MCMC, Markov chain Monte Carlo.

409

410

411 **Availability of supporting data**

412 The datasets generated and/or analysed during the current study are available in the github repository,

413 github.com/sebastianduchene/bacteria_genomic_rates_data

414

415 **Declarations**

416 *Ethics and consent to participate*

417 Not applicable.

418

419 *Consent to publish*

420 Not applicable.

421

422 *Availability of data and materials*

423 All the software used in this study is freely available and open source. The data are all available online

424 github.com/sebastianduchene/bacteria_genomic_rates_data

425

426 *Competing interest*

427 The authors declare no competing interests.

428

429 *Funding*

430 SD was supported by a McKenzie fellowship from the University of Melbourne. ZAD is funded by strategic award

431 #106158 from the Wellcome Trust of Great Britain. KEH is supported by fellowship #1061409 from the NHMRC of

432 Australia.

433

434 *Author contributions*

435 SD, DD, JLG, and KEH conceived and designed the experiments. SD, JLG, JH and ZD analysed the data. SD wrote the

436 manuscript with input from all the authors.

437

438 *Acknowledgements*

439 Not applicable.

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457 **References**

458 1. Gire SK, Goba A, Andersen KG, Sealfon RSG, Park DJ, Kanneh L, et al. Genomic surveillance elucidates Ebola virus
459 origin and transmission during the 2014 outbreak. Science; 2014;345:1369–72.

460 2. Holmes EC, Dudas G, Rambaut A, Andersen KG. The evolution of Ebola virus: Insights from the 2013–2016
461 epidemic. Nature. 2016;538:193–200.

462 3. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. Mol Biol
463 Evol. 2012;29:1969–73.

464 4. Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, et al. BEAST 2: a software platform for Bayesian
465 evolutionary analysis. PLOS Comput Biol. 2014;10:e1003537.

466 5. Ho SYW, Duchêne S. Molecular-clock methods for estimating evolutionary rates and time scales. Mol Ecol.
467 2014;23:5947–75.

468 6. Volz EM, Koelle K, Bedford T. Viral phylodynamics. PLOS Comput Biol. 2013;9:e1002947.

469 7. du Plessis L, Stadler T. Getting to the root of epidemic spread with phylodynamic analysis of genomic data. Trends
470 Microbiol. 2015;23:383–6.

471 8. Rieux A, Balloux F. Inferences from tip-calibrated phylogenies: a review and a practical guide. Mol Ecol.
472 2016;25:1911–24.

473 9. Stadler T, Kouyos R, von Wyl V, Yerly S, Böni J, Bürgisser P, et al. Estimating the basic reproductive number from
474 viral sequence data. Mol Biol Evol. 2012;29:347–57.

475 10. Wong VK, Baker S, Pickard DJ, Parkhill J, Page AJ, Feasey NA, et al. Phylogeographical analysis of the dominant
476 multidrug-resistant H58 clade of Salmonella Typhi identifies inter-and intracontinental transmission events. Nat
477 Genet. 2015;47:632–9.

478 11. Njamkepo E, Fawal N, Tran-Dien A, Hawkey J, Strockbine N, Jenkins C, et al. Global phylogeography and
479 evolutionary history of Shigella dysenteriae type 1. Nat Microbiol. 2016;1:16027.

480 12. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.

481   Bioinformatics. 2014;30:1312–3.

482   13. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for
483   estimating maximum-likelihood phylogenies. Mol Biol Evol. SMBE; 2015;32:268–74.

484   14. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods to estimate
485   maximum likelihood phylogenies: assessing the performance of PhyML 3.0. Syst Biol. 2010;59:307–21.

486   15. Price MN, Dehal PS, Arkin AP. FastTree 2–approximately maximum-likelihood trees for large alignments. PLoS
487   One. 2010;5:e9490.

488   16. To T-H, Jung M, Lycett S, Gascuel O. Fast dating using least-squares criteria and algorithms. Syst Biol.
489   2016;65:82–97.

490   17. Duchêne S, Geoghegan JL, Holmes EC, Ho SYW. Estimating evolutionary rates using time-structured data: a
491   general comparison of phylogenetic methods. Bioinformatics. 2016;32:3375–9.

492   18. Kumar S, Hedges SB. Advances in time estimation methods for molecular data. Mol Biol Evol. 2016;

493   19. Volz EM, Frost SDW. Scalable relaxed clock phylogenetic dating. Virus Evol. 2017;3.

494   20. Sagulenko P, Puller V, Neher R. TreeTime: maximum likelihood phylodynamic analysis. bioRxiv. 2017;153494.

495   21. Stadler T. Mammalian phylogeny reveals recent diversification rate shifts. Proc Natl Acad Sci. 2011;108:6187–92.

496   22. Höhna S, Landis MJ, Heath TA, Boussau B, Lartillot N, Moore BR, et al. RevBayes: Bayesian phylogenetic
497   inference using graphical models and an interactive model-specification language. Syst Biol. 2016;65:726–36.

498   23. Poon AFY. Phylodynamic inference with kernel ABC and its application to HIV epidemiology. Mol Biol Evol.
499   2015;32:2483–95.

500   24. Saulnier E, Alizon S, Gascuel O. Assessing the accuracy of Approximate Bayesian Computation approaches to
501   infer epidemiological parameters from phylogenies. PLOS Comput Biol. 2017;13:e1005416.

502   25. Merker M, Blin C, Mona S, Duforet-Frebourg N, Lecher S, Willery E, et al. Evolutionary history and global spread
503   of the Mycobacterium tuberculosis Beijing lineage. Nat Genet. 2015;47:242–9.

504   26. Devault AM, Golding GB, Waglechner N, Enk JM, Kuch M, Tien JH, et al. Second-pandemic strain of Vibrio
505   cholerae from the Philadelphia cholera outbreak of 1849. N Engl J Med. 2014;370:334–40.

506   27. Baines SL, Holt KE, Schultz MB, Seemann T, Howden BO, Jensen SO, et al. Convergent adaptation in the
507   dominant global hospital clone ST239 of methicillin-Resistant Staphylococcus aureus. MBio. 2015;6:e00080-15.

508   28. Drummond AJ, Ho SYW, Phillips MJ, Rambaut A. Relaxed phylogenetics and dating with confidence. PLOS Biol.
509   2006;4:699–710.

510   29. Ho SYW, Duchêne S, Duchêne D. Simulating and detecting autocorrelation of molecular evolutionary rates
511   among lineages. Mol Ecol Resour. 2015;15.

512   30. Duchêne S, Duchêne DA, Di Giallonardo F, Eden J-S, Geoghegan JL, Holt KE, et al. Cross-validation to select
513   Bayesian hierarchical models in phylogenetics. BMC Evol Biol. 2016;16.

514    31. Anisimova M, Gascuel O. Approximate likelihood-ratio test for branches: A fast, accurate, and powerful
515    alternative. Syst Biol. 2006;55:539–52.

516    32. Ramsden C, Holmes EC, Charleston MA. Hantavirus evolution in relation to its rodent and insectivore hosts: no
517    evidence for codivergence. Mol Biol Evol. 2009;26:143–53.

518    33. Duchêne S, Duchêne DA, Holmes EC, Ho SYW. The performance of the date-randomization test in phylogenetic
519    analyses of time-structured virus data. Mol Biol Evol. 2015;32:1895–906.

520    34. Rambaut A, Lam TT, Carvalho LM, Pybus OG. Exploring the temporal structure of heterochronous sequences
521    using TempEst (formerly Path-O-Gen). Virus Evol. 2016;2:vew007.

522    35. Duchêne S, Holt KE, Weill F-X, Le Hello S, Hawkey J, Edwards DJ, et al. Genome-scale rates of evolutionary
523    change in bacteria. Microb Genomics. Microbiology Society; 2016;2:e000094.

524    36. Boskova V, Bonhoeffer S, Stadler T. Inference of epidemiological dynamics based on simulated phylogenies using
525    birth-death and coalescent models. PLOS Comput Biol. 2014;10:e1003913.

526    37. Rosenberg NA, Nordborg M. Genealogical trees, coalescent theory and the analysis of genetic polymorphisms.
527    Nat Rev Genet. 2002;3:380.

528    38. Boskova V, Stadler T, Magnus C. The influence of phylodynamic model specifications on parameter estimates of
529    the Zika virus epidemic. Virus Evol. 2018;4:vex044.

530    39. Murray GGR, Wang F, Harrison EM, Paterson GK, Mather AE, Harris SR, et al. The effect of genetic structure on
531    molecular dating and tests for temporal signal. Methods Ecol Evol. 2015;7:80–9.

532    40. Stadler T, Kühnert D, Bonhoeffer S, Drummond AJ. Birth–death skyline plot reveals temporal changes of
533    epidemic spread in HIV and hepatitis C virus (HCV). Proc Natl Acad Sci. 2013;110:228–33.

534    41. Kühnert D, Stadler T, Vaughan TG, Drummond AJ. Phylodynamics with migration: A computational framework to
535    quantify population structure from genomic data. Mol Biol Evol. 2016;33:2102–16.

536    42. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid phylogenetic analysis of large
537    samples of recombinant bacterial whole genome sequences using Gubbins. Nucleic Acids Res. 2014;43:e15.

538    43. Rambaut A, Grass NC. Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along
539    phylogenetic trees. Bioinformatics. 1997;13:235–8.

540

541

542

543    **Figure legends**

544

545    **Figure 1. Estimates of evolutionary rate, time to the most recent common ancestor, and the coefficient of rate**
546    **variation of the UCLN.** The histograms correspond to the posterior distribution in BEAST2 using the full Bayesian
547    approach. With the exception of the *Mycobacterium tuberculosis* data set, we used the UCLN clock model because the
548    coefficient of rate variation was not abutting zero. The red solid line is the estimate from LSD, and the dashed lines

549     correspond to the 95% confidence interval. Note that the coefficient of rate variation is not computed for LSD, which

550     assumes a strict molecular clock.

551

552     **Figure 2. Date randomisation test using LSD and BEAST2.** The left column shows histograms of the rate estimates

553     with randomised sampling times in LSD (grey). The red line corresponds to the estimate using the correct sampling

554     times. The right column shows the date randomisation test in BEAST2. The grey bars denote the 95% credible

555     intervals of substitution rate estimates from the randomisations. The red lines correspond to the 95% credible

556     interval of the rate estimates using the correct sampling times. The circles denote the mean value. The $x$-axis in the

557     left column and the $y$-axis in the right column are in logarithmic scale.

558

559     **Figure 3. Posterior estimates of demographic parameters, $\Phi$ and $r$ using the full Bayesian and hybrid**

560     **approaches.** The red histograms correspond to the estimates from the hybrid approach, where the coalescent

561     likelihood is calculated on a fixed tree. The grey histograms correspond to the posterior estimates using the full

562     Bayesian method.

563

564     **Figure 4. Date randomisation test in LSD and estimates of demographic parameters for large data sets using the**

565     **hybrid approach.** The grey histograms correspond to rate estimates from the randomisations, while the red lines

566     correspond to the estimates using the correct sampling times. The red histograms correspond to the posterior

567     distribution of parameters $\Phi$ and $r$.

568

569     **Figure 5. Parameter estimates for 10 randomly selected simulations (from a total of 100).** Simulations with strong

570     temporal structure (a) had a $p$-value for the date randomisations test of 0.00, where as those with no temporal

571     structure (b) had a $p$-value of 1. Each row within each panel is for a simulated genome analysis. Estimates in red were

572     obtained using the hybrid method, while those in grey are for the full Bayesian approach. The circles correspond to

573     the mean value, except for the age of the root-node for the hybrid approach (LSD), where it is the point estimate.

574     The bars denote the 95% credible interval. The dashed lines are the value used to generate the data. Note that the $x$-

575     axes in (b) are in $\log_{10}$ scale.

576

577     **Supplementary material legends**

578

579     **Fig.S1. Root-to-tip regression for all data sets.** The blue points correspond to tips in the tree. The black line

580     represents the linear regression of root-to-tip distance as a function of the sampling time. The root-to-tip distance is

581     measured by fitting the root of the tree that maximises $R^2$.
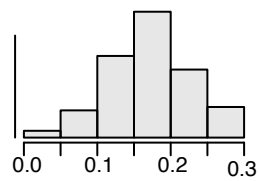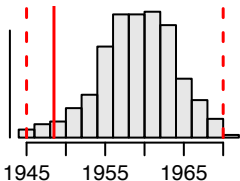
*Mycobacterium tuberculosis* Lineage 2

*Vibrio cholerae*

*Shigella dysenteriae* type I

*Staphylococcus aureus* ST239

Posterior density

Substitution rate
(subs/site/year)

Age of root-node
(year)

Coefficient of rate
variation

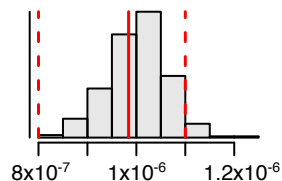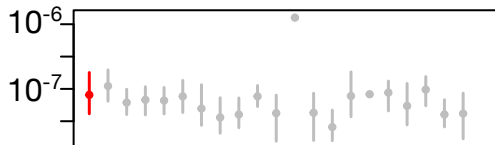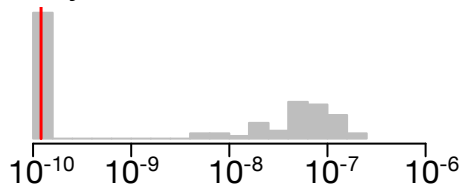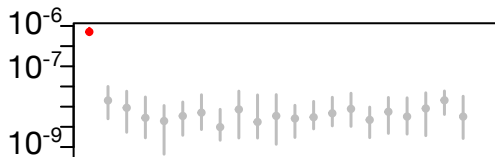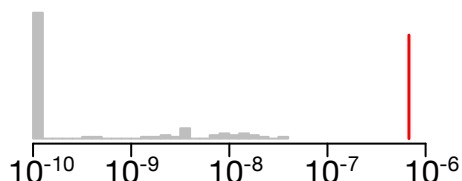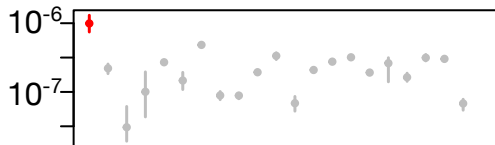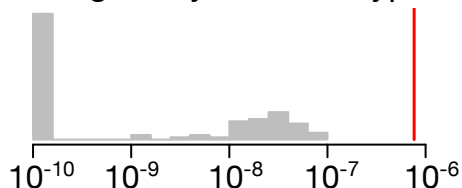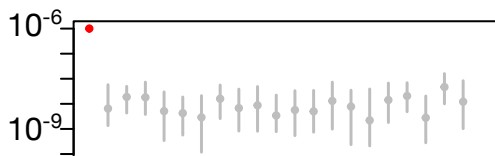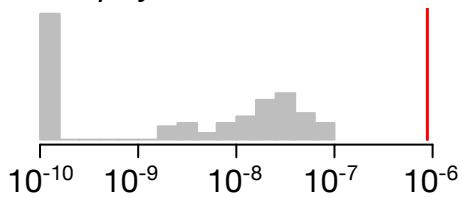*Mycobacterium tuberculosis* Lineage 2
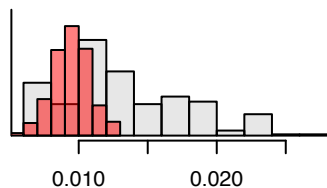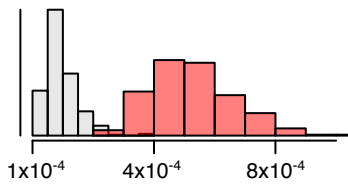
*Vibrio cholerae*

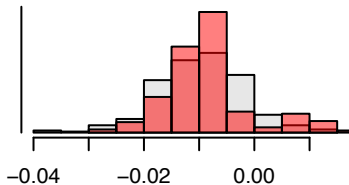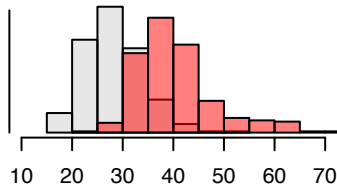*Shigella dysenteriae* type I

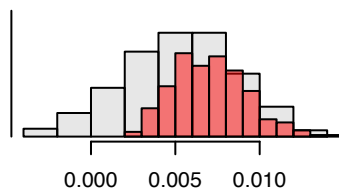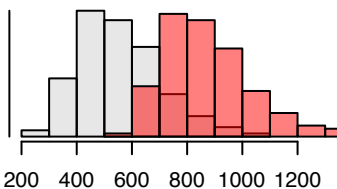*Staphylococcus aureus* ST239

Frequency

Substitution rate (subs/site/year)

*Mycobacterium tuberculosis* Lineage 2
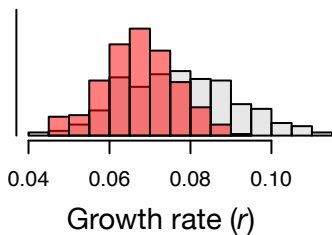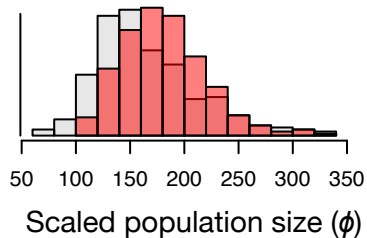
*Vibrio cholerae*

*Shigella dysenteriae*
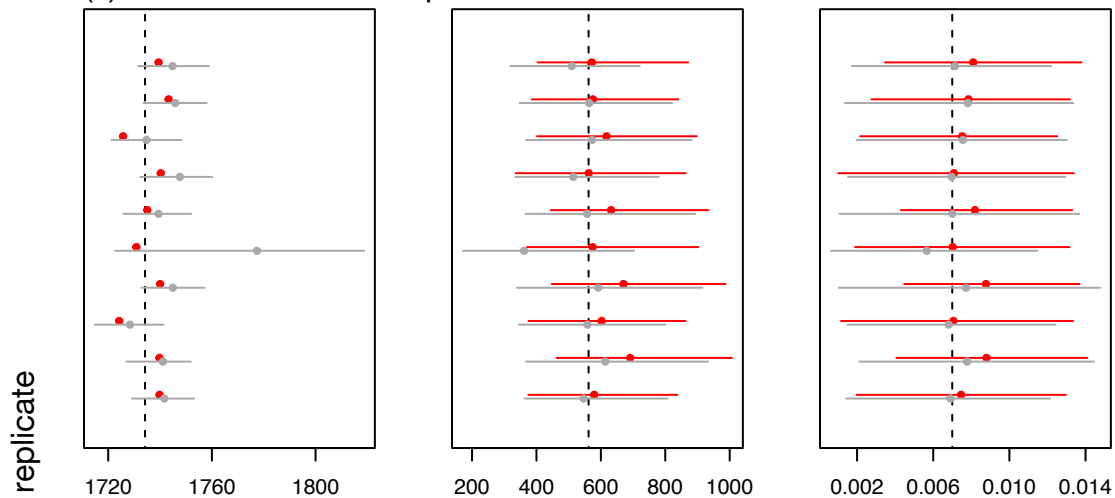
*Staphylococcus aureus* ST239

Scaled population size ($\phi$)

Growth rate ($r$)

*Shigella dysenteriae* type I (global data set)

*Shigella dysenteriae* type I (lineage IV)

Frequency

Posterior density

Substitution rate
(subs/site/year)

Scaled population
size (*Φ*)

Growth rate (*r*)

(a) Simulations with temporal structure

(b) Simulations with no temporal structure

Simulation replicate

Age of root-node (year)

Scaled population size ($\phi$)

Growth rate ($r$)