

How does date-rounding affect phylodynamic inference for public health?

Leo A. Featherstone^{*,1}, Danielle J. Ingle¹, Wytamma Wirth¹,
Sebastian Duchene^{1,2}

September 11, 2024

¹ Peter Doherty Institute for Infection and Immunity, University of Melbourne, Australia.

² Department of Computational Biology, Institut Pasteur, Paris, France.

* leo.featherstone@unimelb.edu.au

1 Abstract

2 Phylodynamic analyses enable the inference of epidemiological parameters from pathogen
3 genome sequences for enhanced genomic surveillance in public health. Pathogen
4 genome sequences and their associated sampling times are the essential data in every
5 analysis. However, sampling times are usually associated with hospitalisation or test-
6 ing dates and can sometimes be used to identify individual patients, posing a threat
7 to patient confidentiality. To lower this risk, sampling times are often given with
8 reduced date-resolution to the month or year, which can potentially bias inference
9 of epidemiological parameters. Here, we characterise the extent to which reduced
10 date-resolution biases phylodynamic analyses across a diverse range of empirical and
11 simulated datasets. We develop a practical guideline on when date-rounding biases
12 phylodynamic inference and we show that this bias is both unpredictable in its direc-
13 tion and compounds with decreasing date-resolution, higher substitution rates, and
14 shorter sampling intervals. We conclude by discussing future solutions that prioritise
15 patient confidentiality and propose a method for safer sharing of sampling dates by
16 translating them uniformly by a random number.

17 Introduction

18 Phylodynamics is commonly used to estimate the parameters of viral spread with
19 increasing application to bacteria. It allows estimation of important epidemiological
20 quantities including rates of transmission, the age of outbreaks, rates of spatial ad-
21 vance, and the prevalence of variants of concern ([Attwood *et al.*, 2022](#), [du Plessis and
22 Stadler, 2015](#), [Featherstone *et al.*, 2022](#), [Volz, 2023](#)). It is applicable across the scales
23 of transmission from the pandemic and epidemic scales, such as for SARS-CoV-2
24 and Ebola virus ([Lancet, 2021](#), [Mbala-Kingebeni *et al.*, 2019](#)), to long-term bacterial
25 transmission such as in *Salmonella enterica* and *Klebsiella pneumoniae*. Phylody-
26 namic analyses are most useful where temporal and spatial records of transmission
27 are sparse, using genomic information to help fill in the gaps.

28 The basis of all phylodynamic inference is that epidemiological spread leaves a
29 trace in the form of substitutions in pathogen genomes that can be used to recon-
30 struct transmission histories. Pathogen populations meeting this assumption are said
31 to be ‘measurably evolving populations’ ([Biek *et al.*, 2015](#), [Drummond *et al.*, 2003](#)).
32 In accordance, phylodynamics uses a combination of genome sequences and associ-
33 ated sampling times to leverage measurable evolution and infer temporally explicit
34 parameters of transmission and pathogen demography.

35 Ideal phylodynamic datasets should include precise sampling dates alongside genome
36 sequences ([Black *et al.*, 2020](#)), but sampling times necessarily carry over sensitive in-

37 formation about times of hospitalisation, testing, or treatment than can be used to
38 identify individual patients. This can pose an unacceptable risk for patient confiden-
39 tiality. In some cases, sampling times or dates of admission are even available for
40 purchase or have allowed identification for a majority of patients in a given record
41 (Sweeney, 2013). In acknowledgement of this risk, Shean and Greninger (2018) sug-
42 gest that Expert Determination govern whether sampling times be released alongside
43 genome sequences, and the resolution to which they are disclosed (day, month, year).
44 Essentially, this approach involves an expert opinion on whether information is safe
45 to release on a case-by-case basis.

46 From a phylodynamic point of view, sampling times with reduced resolution are
47 usable. Uncertainty in sampling times can be accommodated in Bayesian inference
48 (Shapiro *et al.*, 2011), but such an approach is only effective when samples with
49 uncertain dates comprise a small proportion of the total data (Rieux and Khatchikian,
50 2017).

51 The most common technique for incorporating data with a majority of uncertain
52 sampling times is to assume that sampling occurred at the middle of the uncertainty
53 range, such as all samples from 2020 being assigned 15 June 2020. Other approaches
54 would include sampling a random day within 2020 using a probability distribution
55 over the duration of 2020 for each sample. Both approaches introduce a degree of
56 error, which may cause bias because sampling times can drive phylodynamic infer-

57 ence (Featherstone *et al.*, 2021, 2023, Volz and Frost, 2014). Understanding this
58 bias has practical significance, as there are many examples of phylodynamic analyses
59 conducted with reduced date resolution for a diverse array of pathogens. These in-
60 clude viral pathogens such as Rabies Virus, Enterovirus, SARS-CoV-2, Dengue virus
61 (Bennett *et al.*, 2010, Talbi *et al.*, 2010, Wolf *et al.*, 2022, Xiao *et al.*, 2022), and
62 bacterial pathogens, such as *Klebsiella pneumoniae*, *Streptococcus pneumoniae*, and
63 *Mycobacterium tuberculosis* (Azarian *et al.*, 2018, Cella *et al.*, 2017, Merker *et al.*,
64 2015).

65 Precision in sampling dates is also relevant to the design and curation of pathogen
66 sequence databases because sampling dates are often considered as metadata, and
67 thus recorded inconsistently throughout repositories (Raza and Luheshi, 2016). For
68 example, as of early September 2024, there were roughly 19.9M SARS-CoV-2 genome
69 sequences available on GISAID with roughly 2.4% (382K) of these having incomplete
70 date information, where sampling dates are absent or only given to the month or year.
71 In other words, roughly 1 in 50 sequences lacked clear date resolution, reflecting global
72 inconsistency in SARS-CoV-2 sampling time records.

73 In recognition of this issue, we characterised the conditions under which biases
74 arise from reduced date resolution in phylodynamic inference. We analysed four
75 empirical datasets of SARS-CoV-2, H1N1 Influenza, *M. tuberculosis*, *Staphylococcus*
76 *aureus*, and conducted a simulation study with parameters corresponding to each em-

77 pirical dataset. These pathogens are key examples of candidates for genome surveil-
78 lance, with SARS-CoV-2 and H1N1 having caused pandemics and *S. aureus* and
79 *M. tuberculosis* being global priority pathogens (WHO, 2024). These data also have
80 diverse infectious periods and molecular evolutionary rates, thus providing a broad
81 representation of phylodynamics' applicability to pathogens presenting human-health
82 threats. For each empirical and simulated dataset, we studied the bias in estimated
83 epidemiological parameters across treatments with sampling times rounded to the
84 day, month, or year. For example, 2021-10-11 would be specified as 2021-10-15 when
85 rounding to the month and 2021-06-15 when the month and day are not provided.

86 We focused on inference of the reproductive number (R_0 or R_e for the basic and
87 effective reproductive number, respectively), defined as the average number of sec-
88 ondary infections stemming from an individual case (reviewed by du Plessis and
89 Stadler (2015), Featherstone *et al.* (2022), Kühnert *et al.* (2011)), the time to the
90 most recent common ancestor (tMRCA), and the substitution rate (substitutions per
91 site per year) in each dataset. Together, these parameters span much of the insight
92 that phylodynamics offers through inferring when an outbreak started and how fast it
93 proceeded. The evolutionary rate is also the central parameter relating evolutionary
94 time to epidemiological time, so any resulting bias in this parameter is expected to
95 have a pervasive effect throughout each phylodynamic model.

96 We hypothesised that reduced date resolution causes bias that compounds where

97 the uncertainty in dates exceeds the average time for a substitution to arise in a given
98 pathogen. We visualise the relationship between date resolution and average substi-
99 tution time in Fig 1. For example, H1N1 influenza virus accumulates substitutions at
100 a rate of about 4×10^{-3} subs/site/year (Hedge *et al.*, 2013). With a genome length
101 of 13,158bp, we then expect roughly one substitution to accrue per week. Therefore,
102 rounding dates to the month or year conflates molecular evolution in time and bi-
103 ases inferences. Based on this, we expected the SARS-CoV-2 and H1N1 datasets to
104 exhibit bias from month resolution onwards, the *S. aureus* dataset to exhibit bias
105 at year resolution, and the *M. tuberculosis* dataset to not display bias up to and
106 including year resolution (See Table 2 for average substitution times).

107 Our results across the simulation study and analyses of empirical data support
108 using the average substitution time as a rough threshold for when date-rounding
109 causes compounding bias. We also discuss factors that modulate the extent of bias,
110 such as duration of sampling intervals and the choice of phylodynamic model. We
111 finish by discussing future solutions that prioritise both patient confidentiality and
112 accurate data sharing for routine phylodynamic analyses for public health.

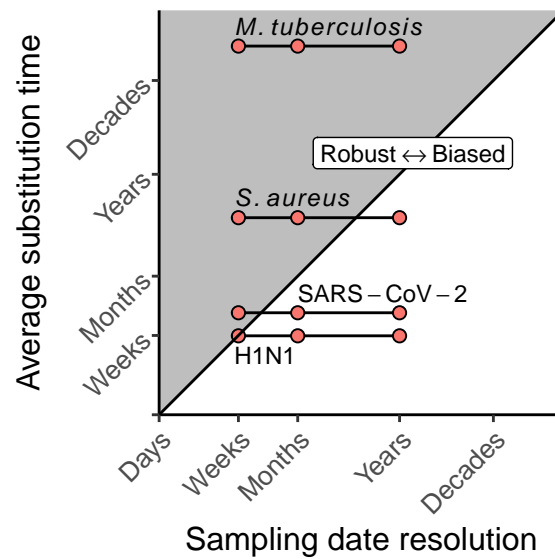


Figure 1: The average time to accrue one substitution based on a fixed genome size and evolutionary rate, $T_s = [\text{Genome Length (sites)} \times \text{Evolutionary rate (subs/site/yr)}]^{-1}$ against the temporal resolution lost by date-rounding. We hypothesised and showed that when analyses for a given pathogen round dates to an extent nearing or crossing the diagonal from left to right, biases is induced in R_e , tMRCA, and substitution rate. substitution rates are taken from each source for the empirical data. We do not report the numerical axis as this figure is designed to illustrate a concept rather than serve as a reference, in the same spirit as is inspiration in Figure 2 of [Biek *et al.* \(2015\)](#).

113 Methods

114 Overview

115 Our study is based on four empirical datasets including with two viruses, H1N1
 116 influenza and SARS-CoV-2, and two bacterial species, *Staphylococcus aureus* and
 117 *Mycobacterium tuberculosis*. We also conducted a simulation study with parameters
 118 tailored to each dataset. These data were chosen to span the usual parameter space for

119 substitution rate and sampling duration in phylodynamics for epidemiology (roughly
120 10^{-3} -to- 10^{-8} (subs/site/yr) for substitution rate and months-to-decades for duration
121 of sampling).

122 To assess the effects of date-rounding, we conducted phylodynamic analyses for
123 both the empirical and simulated datasets with sampling dates rounded to the day,
124 month, or year. For example, two samples from 2000-05-29 and 2000-05-02 would
125 become 2000-05-15 if rounded to the month. We then measured the resulting bias
126 in epidemiologically- or phylodynamically-important parameters: the reproductive
127 number (R_0 or R_e), substitution rate (subs/site/year), and the tMRCA. The tMRCA
128 gives a measure of the age of the pathogen population driving the outbreak and is
129 often interpreted as the age of the outbreak. We also consider the tMRCA to facilitate
130 comparison, because there is variability in which phylodynamic models include the
131 length of the root branch in the age of the outbreak (Stadler *et al.*, 2012).

132 The two viral datasets consist of samples from the 2009 H1N1 pandemic (n=161)
133 from Hedge *et al.* (2013), and a cluster of early SARS-CoV-2 cases from Victoria,
134 Australia in 2020 (n = 112) (Lane *et al.*, 2021). The bacterial datasets consist of
135 *S. aureus*, with 104 samples from New York sampled over ≈ 2 years (Duchêne *et al.*,
136 2016, Uhlemann *et al.*, 2014, Volz and Didelot, 2018), and 30 *M. tuberculosis* sam-
137 ples from an ≈ 25 year outbreak studied by Kühnert *et al.* (2018). These data were
138 chosen because they encompass a diversity of epidemiological dynamics, timescales,

139 and variable substitution rates.

140 **Simulation Study**

141 We simulated outbreaks as birth-death sampling processes using the ReMaster pack-
142 age in BEAST v2.7.6 (Bouckaert *et al.*, 2019, Vaughan, 2024). Simulations consisted
143 of four parameter settings corresponding to each empirical dataset (Table 1), with 100
144 replicates of each. All parameter sets include a proportion of sequenced cases (p), out-
145 break duration (T), and a ‘becoming un-infectious’ rate ($\delta = \frac{1}{\text{Duration of infection}}$).
146 For simulations corresponding the viral datasets, transmission was modelled via R_0 ,
147 the average number of secondary infections (assuming a fully susceptible population).
148 For those corresponding to the bacterial datasets, we allowed the effective reproduc-
149 tive numbers to vary over two intervals (R_{e_1} and R_{e_2} respectively). For the *S. aureus*
150 setting, the change time for R_e was set at $t = 22$ with the sequencing proportion
151 (p) also set to zero before this time to replicate the sampling effort in the empiri-
152 cal dataset. For the *M. tuberculosis* dataset, the change time was fixed at halfway
153 through simulations ($t = 12.5$) with a fixed sequencing proportion throughout.

Table 1: Parameter sets for the simulation study corresponding to each empirical dataset. δ is the ‘becoming un-infectious’ rate, which is the reciprocal of the duration of infection in units of years⁻¹. R_0 is the basic reproductive number, describing the average number of secondary infections arising at the beginning of an outbreak where the susceptible population is greatest. $R_{e\bullet}$ refer to the effective reproductive number over two successive intervals of an outbreak as the susceptible population varies. p is the proportion of sequenced cases. T is the duration of the outbreak.

| Microbe | $\delta(\text{yrs})^{-1}$ | R_0 | R_{e_1} | R_{e_2} | p | $T(\text{yrs})$ | Source |
|------------------------|---------------------------|-------|-----------|-----------|------------------|-----------------|---------------------------------------|
| H1N1 | 91.31 | 1.3 | - | - | 0.015 | 0.25 | Hedge et al. (2013) |
| SARS-CoV-2 | 36.56 | 2.5 | - | - | 0.80 | 0.16 | Lane et al. (2021) |
| <i>S. aureus</i> | 0.93 | - | 2.0 | 1.0 | 0.2 [†] | 25 | Duchêne et al. (2016) |
| <i>M. tuberculosis</i> | 0.125 | - | 2.0 | 1.10 | 0.08 | 25.0 | Kühnert et al. (2018) |

154 [†] p was set to zero before $T = 22$

155 Simulations generated a total of 400 outbreaks which we then used to simulate
156 sequences data under a Jukes-Cantor model using Seq-Gen v1.3.4 ([Rambaut and](#)
157 [Grass, 1997](#)) with fixed substitution rates (Table 2). We chose a simple substitution
158 model to reduce parameter space and because substitution model mismatch has been
159 widely explored elsewhere (e.g. [Lemmon and Moriarty \(2004\)](#)).

Table 2: Substitution rates and genome length for sequence simulation.

| Microbe | Substitution Rate (subs/site/yr) | Genome Length | Time/Sub/Genome (yrs) |
|------------------------|----------------------------------|---------------|-----------------------|
| H1N1 | 4×10^{-3} | 13158 | 0.0190 |
| SARS-CoV-2 | 1×10^{-3} | 29903 | 0.0334 |
| <i>S. aureus</i> | 1×10^{-6} | 2900000 | 0.3458 |
| <i>M. tuberculosis</i> | 1×10^{-7} | 4300000 | 2.3256 |

160 We then analysed each of the 400 simulated datasets under each tree prior and
161 three date resolutions (day, month, and year), yielding 1800 analyses (1200 for the
162 birth-death and 600 for coalescent with exponential growth, referred to hereon as
163 the ‘coalescent exponential’ or CE). We used identical model specifications and prior
164 distributions as for the corresponding empirical datasets. We ran each MCMC chain
165 for 5×18^8 steps, sampling every $10^{4\text{th}}$ step and discarding the first 50% as burnin.
166 We then discarded all analyses that did not have effective sample sizes of the MCMC
167 (ESS) of at least 200 ($ESS \geq 200$), leaving a total of 1670 replicates incorporated
168 in our results.

169 Empirical Data

170 We conducted Bayesian phylodynamic analyses using a birth-death skyline tree prior
171 in BEAST v2.7.6 for all datasets ([Bouckaert *et al.*, 2019](#), [Stadler *et al.*, 2012](#)). We also

172 fit a coalescent tree prior with exponential population growth for the viral datasets
173 (Kingman, 1982). We sampled from the posterior distribution using Markov chain
174 Monte Carlo (MCMC), with 5×10^7 steps (1×10^7 for SARS-CoV-2 data), sampling
175 every 10^4 steps, and discarding the initial 10% as burnin. We assessed sufficient
176 sampling from the stationary distribution by ensuring $ESS \geq 200$ for all parameters
177 and likelihoods.

178 H1N1

179 The H1N1 data consist of 161 samples from North America during the 2009 H1N1
180 influenza virus pandemic, previously analysed by Hedge *et al.* (2013). Samples orig-
181 inate from April to September 2009 and provide an example of a rapidly evolving
182 pathogen sparsely sequenced during an emerging outbreak.

183 Under the birth-death model, we placed a Lognormal($\mu = 0, \sigma = 1$) prior on R_0 ,
184 $\beta(1, 1)$ prior on p , and fixed the becoming-uninfectious rate to ($\delta = 91 \text{ years}^{-1}$), cor-
185 responding to a four-day duration of infection. We also placed an improper ($U(0, \infty)$)
186 prior on the age of the outbreak and a Gamma(shape = 2, rate = 400) prior on the
187 substitution rate.

188 Under the coalescent exponential, we placed a Laplace($\mu = 0, \text{scale} = 100$) prior
189 on the growth rate, which was later transformed to R_0 via $R_0 = rD + 1$ where r is
190 the growth rate and D is the duration of infection. We also placed an improper prior

191 on the effective population size, and otherwise included the same priors as for the
192 birth-death.

193 SARS-CoV-2

194 The SARS-CoV-2 data consist of 112 samples from a densely sequenced transmission
195 cluster from Victoria, Australia over late July to mid September 2020 [Lane et al.](#)
196 [\(2021\)](#). These data are similar to the H1N1 datasets in presenting a quickly evolving
197 viral pathogen, but differ in that a high proportion of cases were sequenced.

198 Under the birth-death, we placed a Lognormal(mean = 1, sd = 1.25) prior on R_0
199 and an Inv-Gamma($\alpha = 5.807$, $\beta = 346.020$) prior on the becoming-uninfectious rate
200 (δ). The sampling proportion was fixed to $p = 0.8$ since every the target was to
201 sequence every known SARS-CoV-2 case in Victoria at this stage of the pandemic,
202 with a roughly 20% sequencing failure rate. We also placed an Exp(mean = 0.019)
203 prior on the origin, corresponding to a lag of up to one week between the index
204 case and the first putative transmission event. Lastly, we placed a Gamma(shape =
205 2, rate = 2000) prior on the substitution rate.

206 Under the coalescent exponential, we placed an improper prior ($\frac{1}{x}$) on the effective
207 population size and a Laplace($\mu = 0.01$, scale = 0.5) prior on the growth rate. Other
208 parameters were given the priors as under the birth-death. Note that we fit the
209 coalescent exponential tree prior for completeness here, but in practice it would not be

210 a reliable model choice due to the high sequencing proportion violating the assumption
211 of a low sequencing proportion under the coalescent. This poor fit is reflected later
212 in the results.

213 *Staphylococcus aureus*

214 The *S. aureus* dataset originates from [Uhlemann *et al.* \(2014\)](#) and we analysed a
215 subset of the data later analysed in [Duchêne *et al.* \(2016\)](#) and [Volz and Didelot \(2018\)](#).
216 It consists of a single nucleotide polymorphism (SNP) alignment of 104 sequenced
217 isolates sampled in New York from 2009 to 2011. Populations growth is understood
218 to have been driven by β -lactam antibiotic use beginning in the 1980s. These data
219 therefore provide a comparison to the *M. tuberculosis* dataset in a briefer sampling
220 span from an outbreak of similar duration.

221 To accommodate changing transmission dynamics, we included two intervals for
222 R_e with a Lognormal($\mu = 0, \sigma = 1$) prior on each. We also placed a $\beta(1, 1)$ prior
223 on the sampling proportion, which was otherwise fixed to 0 before the first sample
224 to capture the lag in sampling. We also placed a $U(0, 1000)$ prior on the origin, and
225 fixed the becoming un-infectious rate at $\delta = 0.93$, corresponding to a nearly year-long
226 duration of infection following [Volz and Didelot \(2018\)](#).

227 *Mycobacterium tuberculosis*

228 The *M. tuberculosis* dataset consists of 36 sequenced isolates from a retrospectively
229 recognised outbreak in California, USA, that originated in the Wat Tham Krabok
230 refugee camp in Thailand. The data were originally analysed using the birth-tree
231 prior by Kühnert *et al.* (2018). We applied the same prior configurations as Kühnert
232 *et al.* (2018), with the exception of including two intervals for R_e and fitting a strict
233 molecular clock with a Gamma(shape = 0.001, rate = 1000.0) prior.

234 Results

235 Simulation study

236 The viral simulation conditions (i.e. SARS-CoV-2 and H1N1) display the greatest
237 bias in mean posterior estimates of substitution rate, tMRCA, and reproductive num-
238 ber with decreasing date resolution (Figure 2 A-C). The *S. aureus* simulations exhibit
239 similar trend with lesser bias in response to decreasing date resolution when rounding
240 dates to the year. The *M. tuberculosis* condition is effectively inert to decreasing date
241 resolution, with mean posterior estimates for each parameter of interest remaining
242 consistent across date resolution (day to year). The *S. aureus* data provide an impor-
243 tant intermediate case in that estimates of each parameter change when transitioning
244 from month to year resolution (see crossing of lines from month to year resolution in

245 the *S. aureus* column of Figure 2). These trends are in agreement with the hypothesis
246 of decreasing date resolution causing increased bias where the resolution lost exceeds
247 the average time for a substitution to arise. This occurs because date-rounding com-
248 presses divergent sequences in time, driving a signal for higher rates of substitution
249 and transmission locally to each temporal cluster of sampling times. This effect is
250 less pronounced in the bacterial simulation conditions relative to the viral conditions,
251 because the date resolution lost is a smaller fraction of the effective substitution time
252 (average time to until substitution is ≈ 4 months and ≈ 28 months for *S. aureus*
253 and *M. tuberculosis* conditions respectively, Table 2). In other words, the bacterial
254 sequences clustered in time were on average less divergent than for the viral data,
255 which is biologically realistic given that bacteria tend to accrue substitutions more
256 slowly than viruses. There are also notable deviations from these general trends
257 across date resolutions, simulation conditions, and tree priors that we attribute to
258 the duration of the sampling intervals below.

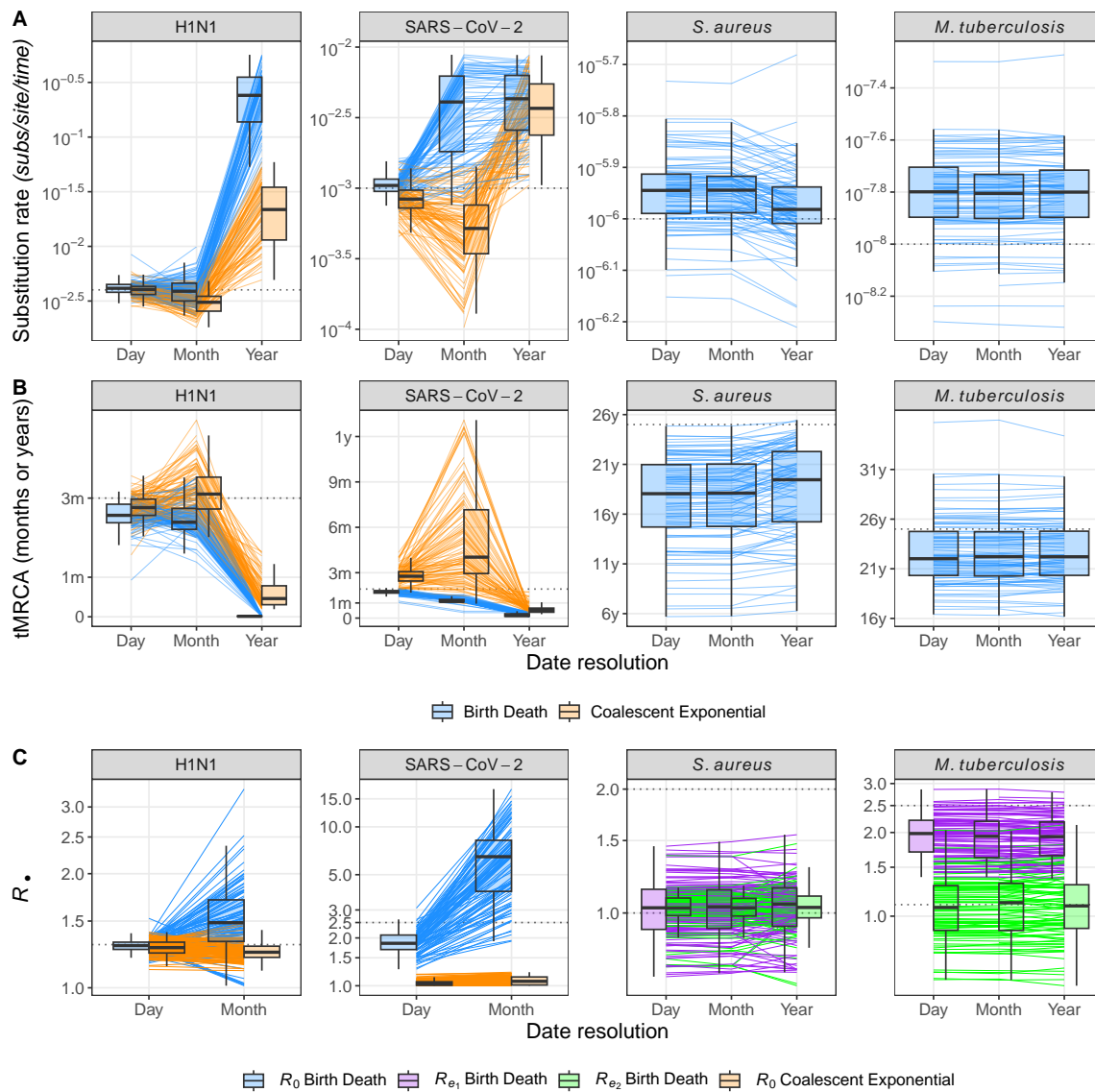


Figure 2: Mean posterior estimates for parameters of interest for each simulated dataset varying across date resolution. Individual lines track mean posterior estimates for each simulated dataset and boxplots are given to summarise the spread and direction of bias across all simulated datasets at each date resolution. Rows correspond to individual parameters, columns correspond to simulation conditions (underlying parameters matching each empirical dataset), and colour corresponds to tree prior or reproductive number interval. Dashed horizontal lines correspond to the true value under which each dataset was simulated. (A) Mean posterior substitution rate across simulation scenarios. (B) Mean posterior tMRCA, a measure of the age of the population driving the outbreak

. (C) mean posterior reproductive number.

259 The coalescent exponential shows overall downwards bias in the substitution rate
260 for the SARS-CoV-2 and H1N1 treatments at month resolution, while the birth-
261 death exhibits upwards bias. Since the sampling times for each viral dataset are
262 distributed over three months, date-rounding compresses samples within a month
263 to one time, simultaneously increasing the time between samples across months and
264 driving a signal for lower transmission and substitution rates between months. The
265 different phylodynamic likelihood functions for each tree prior respond differently
266 to this warped distribution of diversity over time with the coalescent exponential
267 placing more weight on decreased rates of substitution rates while the birth-death
268 favoured an increase. This can be explained by the birth-death drawing signal for
269 increased transmission among coincident sampling times within each month, while the
270 coalescent exponential instead conditions on sampling times (Volz and Frost, 2014).
271 At the year resolution there is there is also lower bias in estimates of substitution rate
272 for the coalescent exponential than the birth-death, however both models estimate
273 upwards-biased substitution rates as year resolution. This is probably because year
274 resolution clusters all sampling times to a single time, meaning a highly inflated rate
275 of substitution is needed to model the artificial burst in diversity at one time for
276 both tree priors - See Figure S2 to see sampling times compressed in time across date
277 resolution for posterior trees. For all viral simulation conditions, the mean posterior
278 tMRCA of each outbreak shifts inversely to the substitution rate. This is the result

279 of a well understood relationship among phylodynamic models where higher rates of
280 evolution suggest shorter periods of evolution.

281 The reproductive number for each viral dataset (R_0) also changes markedly with
282 decreasing date resolution under the birth-death, but not under the coalescent. For
283 the birth-death, this is in agreement with temporal clustering of samples driving
284 a signal for higher transmission rates. Conversely, estimates under the coalescent
285 exponential remain near-identical at month resolution, which is again due to its con-
286 ditioning on sampling times. Estimates of R_0 for the SARS-CoV-2 settings under
287 the coalescent exponential are also heavily biased downwards. This is probably due
288 to high sequencing proportions violating the assumption of low sampling under the
289 coalescent, thus leading to poorly fitting model in the first place.

290 The *S. aureus* condition yields consistent estimates of substitution rate, tMRCA,
291 and reproductive number (R_e in this case) when days are rounded to the month (Fig-
292 ure 2 *S. aureus* column). At year resolution the posterior substitution rate appears
293 biased downwards. This can be explained by the two year sampling duration of the
294 *S. aureus* condition, such that samples rounded to the year will be on average fur-
295 ther apart in time than if dates are given to the month or day (Figure S2). This
296 spacing of diversity in time likely drives the signal for lower substitution rates and
297 an older outbreak in turn. There is no clear pattern in the direction of bias for R_{e_1}
298 and R_{e_2} at year resolution, though estimates deviate from those at month and day

299 resolution. Estimates for R_{e1} are also overall lower than their true value of 2.0, and
300 this is attributable to inconsistent sampling over the duration of the outbreak which
301 was previously demonstrated for other datasets with late sampling in [Featherstone](#)
302 *et al.* (2021).

303 The *M. tuberculosis* simulation condition effectively acts as a control, since it
304 appears inert to date-rounding. This is expected because this dataset reflects longer
305 simulation time, with temporal clustering less likely to inflate R_e , and an average
306 substitution time is longer than a year. As such, even rounding to the year is unlikely
307 to drive a signal for increased evolutionary rate or a more recent origin time.

308 Phylodynamic and phylogenetic terms from the total posterior likelihood also vary
309 with decreasing date resolution (Figure S3). deviation also increases with lesser date
310 resolution from month and year. This verifies that altered date resolution affects the
311 likelihood manifold of each analysis, which is reflected in the different trends of bias
312 in each parameter of interest.

313 **Empirical Results**

314 Broadly, analyses of the empirical datasets reproduce the patterns of bias in the sim-
315 ulation study(Figure 3). That is, the reproductive number increases with decreasing
316 date resolution along with an increase in the substitution rate and corresponding
317 decrease in the tMRCA. There are a few exceptions to this trend that we consider

318 below and which we again attribute to the difference between simulated and empirical
319 sampling time distributions.

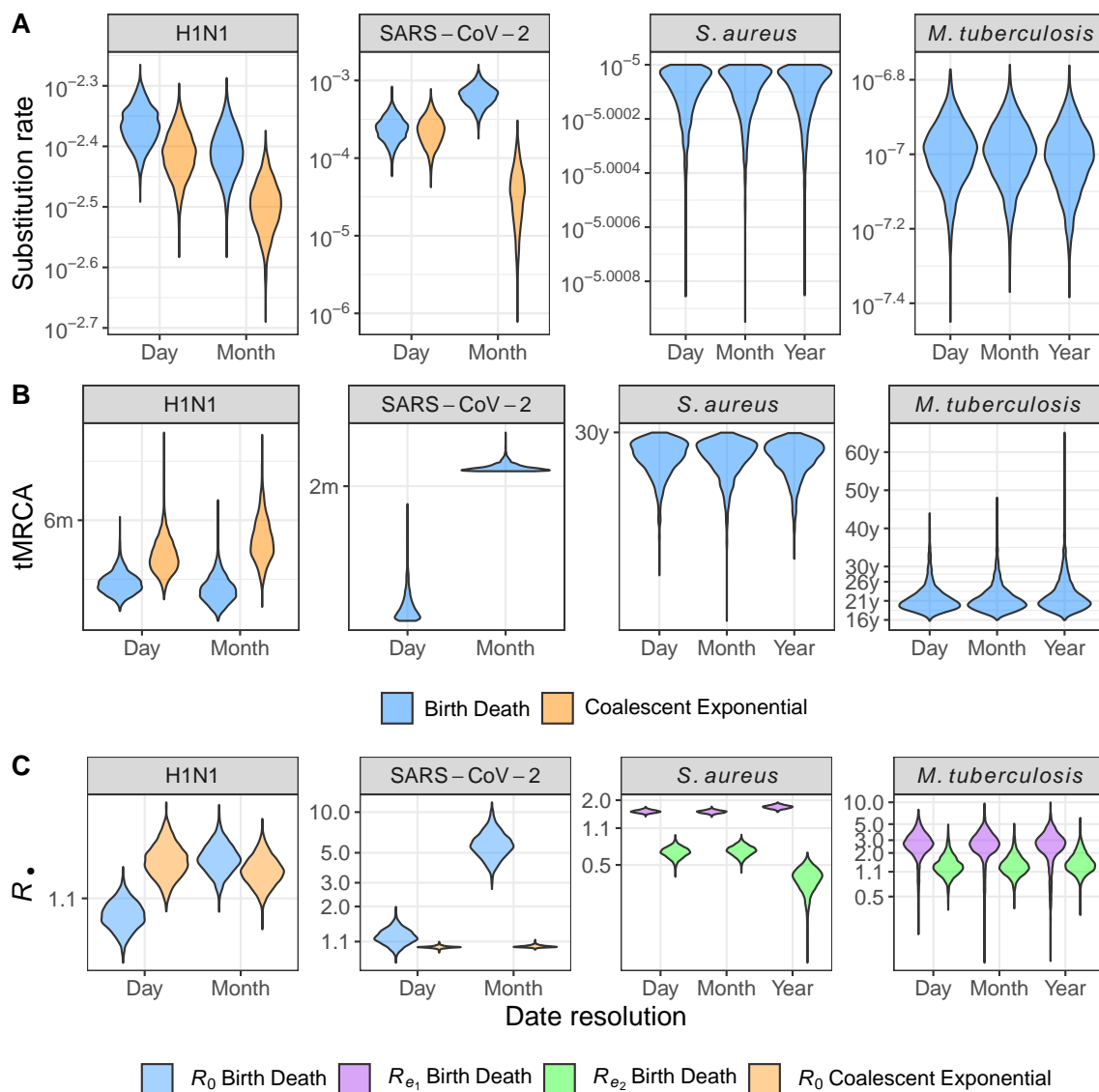


Figure 3: Posterior distributions for parameters of interest estimated for each empirical dataset. Date resolution is given on the horizontal axis and colour denotes tree prior. Estimates for viral datasets at year-resolution are omitted because results deviate by implausible orders of magnitude due to sampling times rounded to identical dates. (A) Posterior substitution rate across date resolutions. (B) Posterior tMRCA in units of months (m) or years (y). (C) Posterior reproductive number on a log-transformed axis.

320

Phylogenetic and phylogenetic likelihoods also diverge where the loss in date

321 resolution exceed the average time for a substitution to arise (Figure 4). For both
322 viral datasets, month and day posterior distributions of phylodynamic and phyloge-
323 netic likelihood are diverged, while likelihoods overlap at all date resolutions for the
324 *M. tuberculosis* data. The *S. aureus* data provide an intermediate case where only
325 the posterior likelihoods for year-resolution differ. Together, these likelihood distri-
326 butions support the hypothesis that date-uncertainty that is wider (in time) than the
327 average time to one substitution causes a qualitative shift in the likelihood manifold
328 for analyses under both birth-death based and coalescent tree priors.

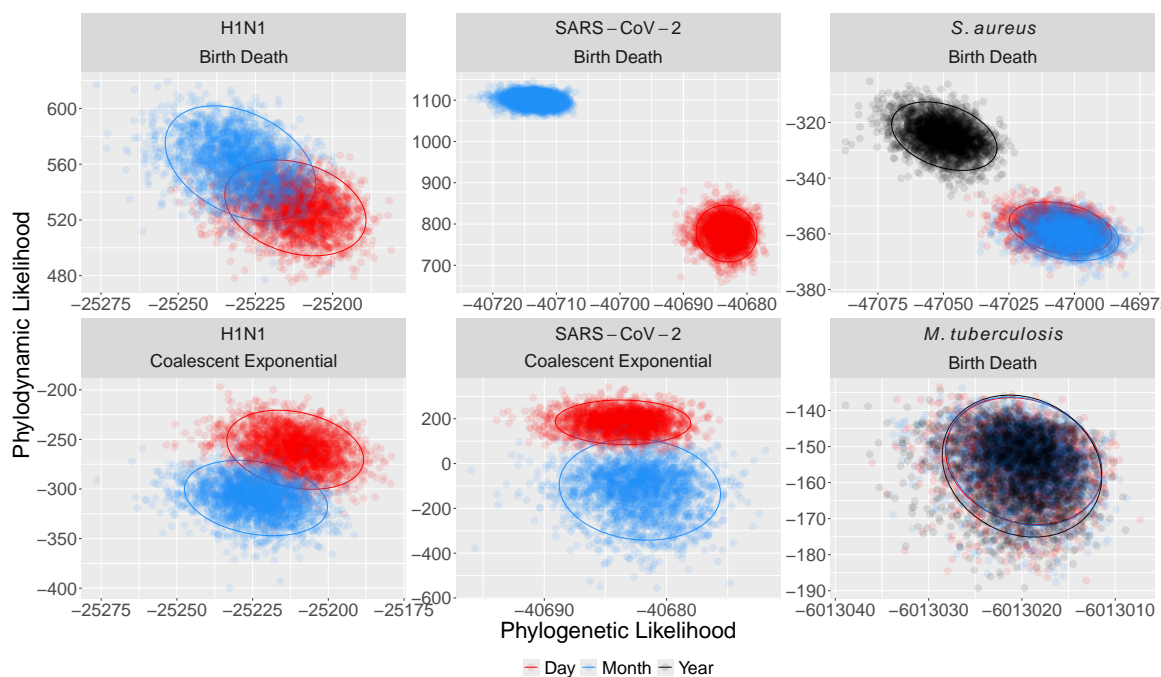


Figure 4: Posterior phylodynamic likelihood against phylogenetic likelihood for each combination of empirical dataset and with colour corresponding to date resolution. Ellipses surround the 95% highest posterior density region for each posterior. Both Phylodynamic and phylogenetic likelihoods diverge between day and month resolution for the viral datasets, while year resolution differs from month and day for the *S. aureus* data. Posterior likelihoods all coincide for *M. tuberculosis*.

329 **H1N1**

330 Mean posterior R_0 increases from day to month resolution for the birth-death (1.08 to
331 1.14), yet remains near-identical for the coalescent exponential (1.14 to 1.13) (Table
332 S1). The mean posterior substitution rate also decreases for both tree priors across
333 day to month resolution (4.3×10^{-3} to 3.9×10^{-3} and 3.9×10^{-3} to 3.2×10^{-3} for the
334 birth-death and coalescent exponential respectively) (Table S1). The posterior tM-
335 RCA also differs between tree priors mirroring substitution rate, with a decrease from
336 day to month resolution for the birth-death and an increase for the coalescent expo-
337 nential. For the coalescent exponential, we can attribute the decrease in reproductive
338 number and substitution rate from day to month resolution to samples being spread
339 further in time (Figure S1), which drives the signal for an older outbreak and
340 lower transmission rates. While the same is true for the sampling distribution under
341 the birth-death, the additional information it draws from identical sampling times as
342 month-resolution likely inflates the mean posterior reproductive number despite the
343 signal for a lower substitution rate and older outbreak.

344 **SARS-CoV-2**

345 Under the birth-death, the SARS-CoV-2 dataset behaves as expected, with an in-
346 crease in posterior R_0 from day to month rounding. In particular, rounding to the
347 month results in a high, yet plausible value of $R_0 = 5.972$ (Table S1). Under the

348 coalescent exponential, the mean posterior R_0 remains near identical across day to
349 month treatments (1.00 to 1.01 respectively). We again note that the coalescent ex-
350 ponential is included for completeness for the SARS-CoV-2 dataset, but is not an
351 appropriate choice of model in practice due to the near complete-sequencing of the
352 original transmission cluster. Thus, poor model-fit is probably the cause of unrealistic
353 estimates of R_0 .

354 The mean posterior substitution rate under the birth-death increases over two-
355 fold when rounding to the month (2.47×10^{-4} to 6.56×10^{-4} , Table S1). Mean
356 posterior tMRCA also increases from 0.15 years to 0.17 years from day to month,
357 which contradicts the expectation of a decreased estimate of tMRCA under date-
358 rounding. We again attribute these differences to the distribution of the empirical
359 sampling times under date-rounding. Sampling for the SARS-CoV-2 dataset mainly
360 occurred over August to September 2020, with most August samples originating later
361 in the month (Figure S1). Rounding to 15th of August therefore made these samples
362 appear older in time and likely contributed to the older origin under month-rounding.

363 *S. aureus*

364 For R_{e1} , the *S. aureus* dataset recapitulated the simulation study with month round-
365 ing having a minimal effect, but year rounding inducing an upwards bias (mean
366 values of 1.57, 1.56, 1.73 respectively)(Figure 3, Table S1). R_{e2} displays a similar

367 pattern with consistent estimates at day and month-rounding before a reduction at
368 year rounding (0.66, 0.67, and 0.37 respectively). This result is consistent with the
369 estimates an initial increase in growth rate in previous analyses of the dataset [Volz](#)
370 [and Didelot \(2018\)](#).

371 Mean posterior substitution rate and tMRCA remain identical across date resolu-
372 tions (10^{-5} subs/site/year and a tMRCA of 30 years), despite the change in reproduc-
373 tive numbers at year rounding. This is surprising given the change in posterior phy-
374 lodynamic and phylogenetic likelihoods (Figure 4), and highlights that date-rounding
375 can perturb the likelihood without predictable changes in parameters of epidemiolog-
376 ical significance.

377 *M. tuberculosis*

378 The *M. tuberculosis* data recapitulate the outcome of the simulation study in be-
379 ing robust to date-rounding. Posterior substitution rates and outbreak ages remain
380 consistent across decreasing date resolution (1.02×10^{-7} , 1.02×10^{-7} , 9.86×10^{-8}
381 (subs/site/time) and 21.7, 21.7, and 22.5 years respectively) (Table S1, Figure 3).
382 We also infer that $R_{e_1} > R_{e_s}$ across date-rounding conditions, coinciding with an
383 earlier burst of transmission in agreement with [Kühnert et al. \(2018\)](#). However, R_{e_1}
384 decreases slightly date date-rounding (mean posterior estimates of 2.77, 2.74, 2.66 for
385 day, month and year rounding)(Table S1), while R_{e_2} increases (1.4, 1.41, 1.53 from

386 day to year rounding). This was likely caused by the higher number of samples in the
387 second sampling interval, from roughly 2002 to 2010, such that compressing sampling
388 times drive drove a signal for higher transmission in the second interval with longer
389 periods between sampling in the first interval at year resolution. Again, this shows
390 that distribution of sampling times for empirical data, which is largely unpredictable,
391 modulate the effects of date-rounding.

392 Discussion

393 The results of the simulation study and analyses of empirical data support our hypoth-
394 esis that phylodynamic inference is most biased where the temporal resolution lost in
395 date rounding exceeds the average time for one substitution to arise. In the both the
396 simulation study and empirical analyses, the viral datasets (H1N1 and SARS-CoV-2)
397 display the greatest bias in mean posterior reproductive number, substitution rate,
398 and tMRCA when rounding to the month or year, with the average substitution time
399 being less than one month in both simulation conditions. The *S. aureus* data pro-
400 vide an intermediate case, with estimated parameters displaying bias when rounding
401 dates to the year (average substitution time between the order of months to a year).
402 Lastly, the *M. tuberculosis* data also provide supporting evidence in not displaying
403 any notable bias between estimates at day, month, or year date-rounding. This is ex-
404 pected because the average substitution time longer than a year in all *M. tuberculosis*

405 analyses.

406 We therefore propose the average substitution time as a rough practical threshold
407 after which genomic epidemiologists can invariably expect date-rounding to distort in-
408 ference. Genomic epidemiologists can make this assessment by calculating the average
409 substitution time, T_s , as $T_s = [\text{Genome Length (sites)} \times \text{Evolutionary rate (subs/site/yr)}]^{-1}$
410 and checking whether $T_s < \frac{1}{12}$ (indicating substitutions arising faster than monthly)
411 when justifying rounding to the day, or $T_s < 1$ (substitutions arising more than yearly)
412 when justifying rounding to the month. In the more general terms, we propose that
413 date rounding is problematic for fast-evolving RNA viruses, such as in the H1N1 and
414 SARS-CoV-2 datasets. We urge others uploading data to repositories such as GISAID
415 to include dates to the day where possible, and support the practice of including dates
416 to the day on pathoplexus (pathoplexus.org). This will increase the added-value of
417 phylodynamic analysis for future infectious disease threats. Rounding to the year is
418 sufficient for slowly evolving bacteria such as *M. tuberculosis*. We suggest case-by-case
419 assessment for pathogens with intermediate average substitution times, such as the *S.*
420 *aureus* herein and other faster-evolving bacteria including *Streptococcus*, multi-drug
421 resistant *Escherichia coli*, or *Klebsiella pneumoniae* ([Gorrie et al., 2018](#), [Sherry et al.,](#)
422 [2022](#), [Xie et al., 2024](#)). In the specific cases of *S. aureus* and other high disease-burden
423 bacteria with asymptomatic and/or community carriage, we suggest preserving dates
424 as much as possible to recover maximal information given the additional work that

425 is often dedicated to screening samples. Finally, we note that genome samples with
426 or without rounded dates reflect considerable efforts in the field to collect and pro-
427 cess samples. In the case where only low-resolution dates are available, we do not
428 discourage phylodynamic analyses, but instead encourage additional analyses to test
429 the effects of rounded dates, such as by including priors on sampling ranges.

430 We also strongly emphasise that this proposal is a rough guideline lacking rigor-
431 ous mathematical derivation. Any degree of date-rounding may alter likelihood and
432 parameter estimation in phylodynamic analyses. Other factors such as as the length
433 of the sampling window, distribution of sampling times over this interval, and choice
434 of tree prior also affect the direction and severity of bias when rounding dates.

435 Shorter sampling intervals can also exacerbate the bias due to date-rounding.
436 For example, in the SARS-CoV-2 data and simulation conditions, most sequences
437 originated over one month with the remainder towards the end of each of the previous
438 two months. Bias for these data was greater for each parameter compared to otherwise
439 similar H1N1 data, which had a more even distribution of sampling over three full
440 months. This result is in line with previous results for ancient DNA data showing
441 that date-rounding has negligible effects for timescales of millennia or longer, which
442 we expect to span the average substitution time several-fold (Molak *et al.*, 2013). This
443 emphasises the importance of accurate dates for phylodynamic datasets of emerging
444 pathogens sequenced over shorter timescales, where results are likely to be the most

445 urgent and reflect shorter sampling intervals.

446 The choice of tree prior also affects bias when rounding dates. For example, the
447 coalescent exponential tended to infer decreased substitution rates while the birth-
448 death favoured increased substitution rates across simulated and empirical viral data.
449 The inverse trend also arose for the tMRCA. This is because the birth-death draws
450 additional information from clustered sampling times, which serves to elevate rates
451 of substitution and transmission, while the coalescent conditions on these and relies
452 more on the longer duration between sampling times at month resolution for both
453 datasets.

454 Taken together, the results from the simulation study and empirical data show that
455 although date-rounding biases epidemiological estimates in a theoretically predictable
456 directions (upwards for transmission and substitution rates, downwards for tMRCA),
457 the intensity of the bias is difficult to predict and varies with the distribution and
458 span of sampling times as well as tree prior. We conclude that sufficiently accurate
459 sampling times are essential where phylodynamic insight is needed to understand
460 infectious disease epidemiology and evolution. There does not appear to be an clear
461 way to adjust for the bias caused otherwise. Accurate sampling times will be essential
462 for employing phylodynamics amid future infectious disease threats.

463 We also acknowledge that while accurate sampling times are essential for reliable
464 phylodynamic results, it may pose an unacceptable level of risk to patient confiden-

465 tiality to release sampling times. We therefore emphasise the importance of methods
466 that prioritise both patient confidentiality and data transparency and finish by dis-
467 cussing potential future solutions.

468 **Translating dates by random seeds**

469 The functional component of phylodynamic data are the differences among genome
470 sequences and among dates, rather than their absolute values. It may therefore be
471 possible to protect patient confidentiality while sharing accurate dates by translating
472 dates uniformly by a random number. This would protect the true sampling dates
473 while preserving the relative times between them. For example, if the sampling times
474 in a dataset of 3 genomes are 2000, 2001 and 2002, then data providers may randomly
475 draw a number of 1000, which is kept secret, to shift dates. The genome-associated
476 dates 2000, 2001 and 2002 are then shared as 3000, 3001 and 3002. While currently
477 implausible, these translated dates are usable in phylodynamic analyses and preserve
478 the distance between sampling times. Once results are returned the data provider
479 can internally account for the translation in any estimated ages, such as node ages or
480 the tMRCA, by subtracting 1000. For example if the estimated age of the outbreak
481 (taken as the tMRCA) was 5 years before the most recent sample, then the data
482 provider can privately estimate the outbreak's onset as 1997 (2002 - 5), while those
483 conducting the analysis externally can only estimate the relative age of 5 years. In

484 the same way, intervals of transmission parameters such as R_e can be placed with
485 respect to the true sampling times. Rates, such as growth or infection rates can also
486 be accurately estimated via this method since these are not biased by shifting dates
487 uniformly in time.

488 **Distributed computing**

489 Approaches based on distributed computing, where data are analysed across remote
490 servers, also offer promise for maximising data transparency and patient confiden-
491 tiality. For example, Santos *et al.* (2022) recently developed a method to estimate
492 phylogenetic trees from private genome data using distributed computing and quan-
493 tum cryptographic protocols. Routine phylodynamic analysis for genomic surveillance
494 may also benefit from adopting protocols from so-called swarm learning approaches
495 that allow artificial intelligence models in precision medicine to be trained across
496 distributed datasets (together comprising a swarm) (Warnat-Herresthal *et al.*, 2021).
497 Such approaches are in general complementary with hub-and-spoke networks, which
498 are commonly used for storing sensitive pathogen genome data in national reposi-
499 ries (Hoang *et al.*, 2022). We remain optimistic that future advances in distributed
500 computing can eliminate the need for date-rounding in phylodynamic analysis.

501 Data Availability

502 All analyses and data used in this work can be accessed and run as a Snakemake
503 pipeline at <https://github.com/LeoFeatherstone/pdp>.

504 Authors' Contributions

505 LAF designed the study, performed all analyses, and wrote the paper. DJI provided
506 initial empirical data and contributed to writing of the manuscript. WW and SDG
507 provided supervision and contributed to writing of the manuscript.

508 References

- 509 Attwood, S. W. *et al.* (2022). Phylogenetic and phylodynamic approaches to under-
510 standing and combating the early sars-cov-2 pandemic. *Nature Reviews Genetics*,
511 **23**(9), 547–562.
- 512 Azarian, T. *et al.* (2018). The impact of serotype-specific vaccination on phylo-
513 dynamic parameters of *Streptococcus pneumoniae* and the pneumococcal pan-
514 genome. *PLOS Pathogens*, **14**(4), e1006966. Publisher: Public Library of Science.
- 515 Bennett, S. *et al.* (2010). Epidemic Dynamics Revealed in Dengue Evolution. *Molec-
516 ular Biology and Evolution*, **27**(4), 811–818.
- 517 Biek, R. *et al.* (2015). Measurably evolving pathogens in the genomic era. *Trends in
518 Ecology & Evolution*, **30**(6), 306–313. Publisher: Elsevier.
- 519 Black, A. *et al.* (2020). Ten recommendations for supporting open pathogen genomic
520 analysis in public health. *Nature medicine*, **26**(6), 832–841.
- 521 Bouckaert, R. *et al.* (2019). BEAST 2.5: An advanced software platform for bayesian
522 evolutionary analysis. *PLOS Computational Biology*, **15**(4), e1006650. Publisher:
523 Public Library of Science.

- 524 Cella, E. *et al.* (2017). Multi-drug resistant *Klebsiella pneumoniae* strains circulating
525 in hospital setting: whole-genome sequencing and Bayesian phylogenetic analysis
526 for outbreak investigations. *Scientific Reports*, **7**(1), 3534. Number: 1 Publisher:
527 Nature Publishing Group.
- 528 Drummond, A. J. *et al.* (2003). Measurably evolving populations. *Trends in ecology*
529 *& evolution*, **18**(9), 481–488.
- 530 du Plessis, L. and Stadler, T. (2015). Getting to the root of epidemic spread with
531 phylodynamic analysis of genomic data. *Trends in Microbiology*, **23**(7), 383–386.
- 532 Duchêne, S. *et al.* (2016). Genome-scale rates of evolutionary change in bacteria.
533 *Microbial Genomics*, **2**(11).
- 534 Featherstone, L. A. *et al.* (2021). Infectious disease phylodynamics with occur-
535 rence data. *Methods in Ecology and Evolution*, **12**(8), 1498–1507. eprint:
536 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13620>.
- 537 Featherstone, L. A. *et al.* (2022). Epidemiological inference from pathogen genomes:
538 A review of phylodynamic models and applications. *Virus Evolution*, **8**(1), veac045.
- 539 Featherstone, L. A. *et al.* (2023). Decoding the Fundamental Drivers of Phylodynamic
540 Inference. *Molecular Biology and Evolution*, **40**(6), msad132.
- 541 Gorrie, C. L. *et al.* (2018). Antimicrobial-Resistant *Klebsiella pneumoniae* Carriage
542 and Infection in Specialized Geriatric Care Wards Linked to Acquisition in the
543 Referring Hospital. *Clinical Infectious Diseases*, **67**(2), 161–170.
- 544 Hedge, J. *et al.* (2013). Real-time characterization of the molecular epidemiology of
545 an influenza pandemic. *Biology Letters*, **9**(5), 20130331.
- 546 Hoang, T. *et al.* (2022). AusTrakka: Fast-tracking nationalized genomics surveillance
547 in response to the COVID-19 pandemic. *Nature Communications*, **13**(1), 865.
548 Publisher: Nature Publishing Group.
- 549 Kingman, J. F. C. (1982). The coalescent. *Stochastic Processes and their Applications*,
550 **13**(3), 235–248.
- 551 Kühnert, D. *et al.* (2011). Phylogenetic and epidemic modeling of rapidly evolving
552 infectious diseases. *Infection, genetics and evolution*, **11**(8), 1825–1841.
- 553 Kühnert, D. *et al.* (2018). Tuberculosis outbreak investigation using phylodynamic
554 analysis. *Epidemics*, **25**, 47–53.

- 555 Lancet, T. (2021). Genomic sequencing in pandemics. *Lancet (London, England)*,
556 **397**(10273), 445.
- 557 Lane, C. R. *et al.* (2021). Genomics-informed responses in the elimination of covid-19
558 in victoria, australia: an observational, genomic epidemiological study. *The Lancet*
559 *Public Health*, **6**(8), e547–e556.
- 560 Lemmon, A. R. and Moriarty, E. C. (2004). The importance of proper model assump-
561 tion in bayesian phylogenetics. *Systematic Biology*, pages 265–277.
- 562 Mbala-Kingebeni, P. *et al.* (2019). Medical countermeasures during the 2018 ebola
563 virus disease outbreak in the north kivu and ituri provinces of the democratic
564 republic of the congo: a rapid genomic assessment. *The Lancet infectious diseases*,
565 **19**(6), 648–657.
- 566 Merker, M. *et al.* (2015). Evolutionary history and global spread of the Mycobac-
567 terium tuberculosis Beijing lineage. *Nature Genetics*, **47**(3), 242–249. Number: 3
568 Publisher: Nature Publishing Group.
- 569 Molak, M. *et al.* (2013). Phylogenetic Estimation of Timescales Using Ancient DNA:
570 The Effects of Temporal Sampling Scheme and Uncertainty in Sample Ages. *Molec-
571 ular Biology and Evolution*, **30**(2), 253–262.
- 572 Rambaut, A. and Grass, N. C. (1997). Seq-gen: an application for the monte carlo
573 simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics*,
574 **13**(3), 235–238.
- 575 Raza, S. and Luheshi, L. (2016). Big data or bust: realizing the microbial genomics
576 revolution. *Microbial Genomics*, **2**(2).
- 577 Rieux, A. and Khatchikian, C. E. (2017). Tipdatingbeast: An r package to assist
578 the implementation of phylogenetic tip-dating tests using beast. *Molecular Ecology*
579 *Resources*, **17**(4), 608–613.
- 580 Santos, M. B. *et al.* (2022). Private Computation of Phylogenetic Trees Based on
581 Quantum Technologies. *IEEE Access*, **10**, 38065–38088. Conference Name: IEEE
582 Access.
- 583 Shapiro, B. *et al.* (2011). A bayesian phylogenetic method to estimate unknown
584 sequence ages. *Molecular biology and evolution*, **28**(2), 879–887.
- 585 Shean, R. C. and Greninger, A. L. (2018). Private collection: high correlation of sam-
586 ple collection and patient admission date in clinical microbiological testing compli-
587 cates sharing of phylodynamic metadata. *Virus Evolution*, **4**(1), vey005.

- 588 Sherry, N. L. *et al.* (2022). Multi-site implementation of whole genome sequencing
589 for hospital infection control: A prospective genomic epidemiological analysis. *The*
590 *Lancet Regional Health - Western Pacific*, **23**, 100446.
- 591 Stadler, T. *et al.* (2012). Estimating the basic reproductive number from viral se-
592 quence data. *Molecular biology and evolution*, **29**(1), 347–357.
- 593 Sweeney, L. (2013). Matching Known Patients to Health Records in Washington
594 State Data. *SSRN Electronic Journal*.
- 595 Talbi, C. *et al.* (2010). Phylodynamics and Human-Mediated Dispersal of a Zoonotic
596 Virus. *PLOS Pathogens*, **6**(10), e1001166. Publisher: Public Library of Science.
- 597 Uhlemann, A.-C. *et al.* (2014). Molecular tracing of the emergence, diversification, and
598 transmission of *S. aureus* sequence type 8 in a New York community. *Proceedings*
599 *of the National Academy of Sciences*, **111**(18), 6738–6743. Publisher: Proceedings
600 of the National Academy of Sciences.
- 601 Vaughan, T. G. (2024). ReMASTER: improved phylodynamic simulation for BEAST
602 2.7. *Bioinformatics*, **40**(1), btae015.
- 603 Volz, E. (2023). Fitness, growth and transmissibility of SARS-CoV-2 genetic variants.
604 *Nature Reviews Genetics*, pages 1–11. Publisher: Nature Publishing Group.
- 605 Volz, E. M. and Didelot, X. (2018). Modeling the Growth and Decline of Pathogen
606 Effective Population Size Provides Insight into Epidemic Dynamics and Drivers of
607 Antimicrobial Resistance. *Systematic Biology*, **67**(4), 719–728.
- 608 Volz, E. M. and Frost, S. D. W. (2014). Sampling through time and phylodynamic
609 inference with coalescent and birth-death models. *Journal of the Royal Society,*
610 *Interface*, **11**(101), 20140945.
- 611 Warnat-Herresthal, S. *et al.* (2021). Swarm Learning for decentralized and confi-
612 dential clinical machine learning. *Nature*, **594**(7862), 265–270. Publisher: Nature
613 Publishing Group.
- 614 WHO (2024). WHO bacterial priority pathogens list, 2024: Bacterial pathogens of
615 public health importance to guide research, development and strategies to prevent
616 and control antimicrobial resistance.
- 617 Wolf, J. M. *et al.* (2022). Temporal spread and evolution of SARS-CoV-2 in the
618 second pandemic wave in Brazil. *Journal of Medical Virology*, **94**(3), 926–936.
619 eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jmv.27371>.

- 620 Xiao, J. *et al.* (2022). Genomic Epidemiology and Phylodynamic Analysis of En-
621 terovirus A71 Reveal Its Transmission Dynamics in Asia. *Microbiology Spectrum*,
622 **10**(5), e01958–22. Publisher: American Society for Microbiology.
- 623 Xie, O. *et al.* (2024). Temporal and Geographic Strain Dynamics of Invasive Strepto-
624 coccus Pyogenes In Australia: A Multi-Centre Clinical and Genomic Epidemiology
625 Study 2011-2023. *Preprints with The Lancet*.

626 **Supplementary Material**

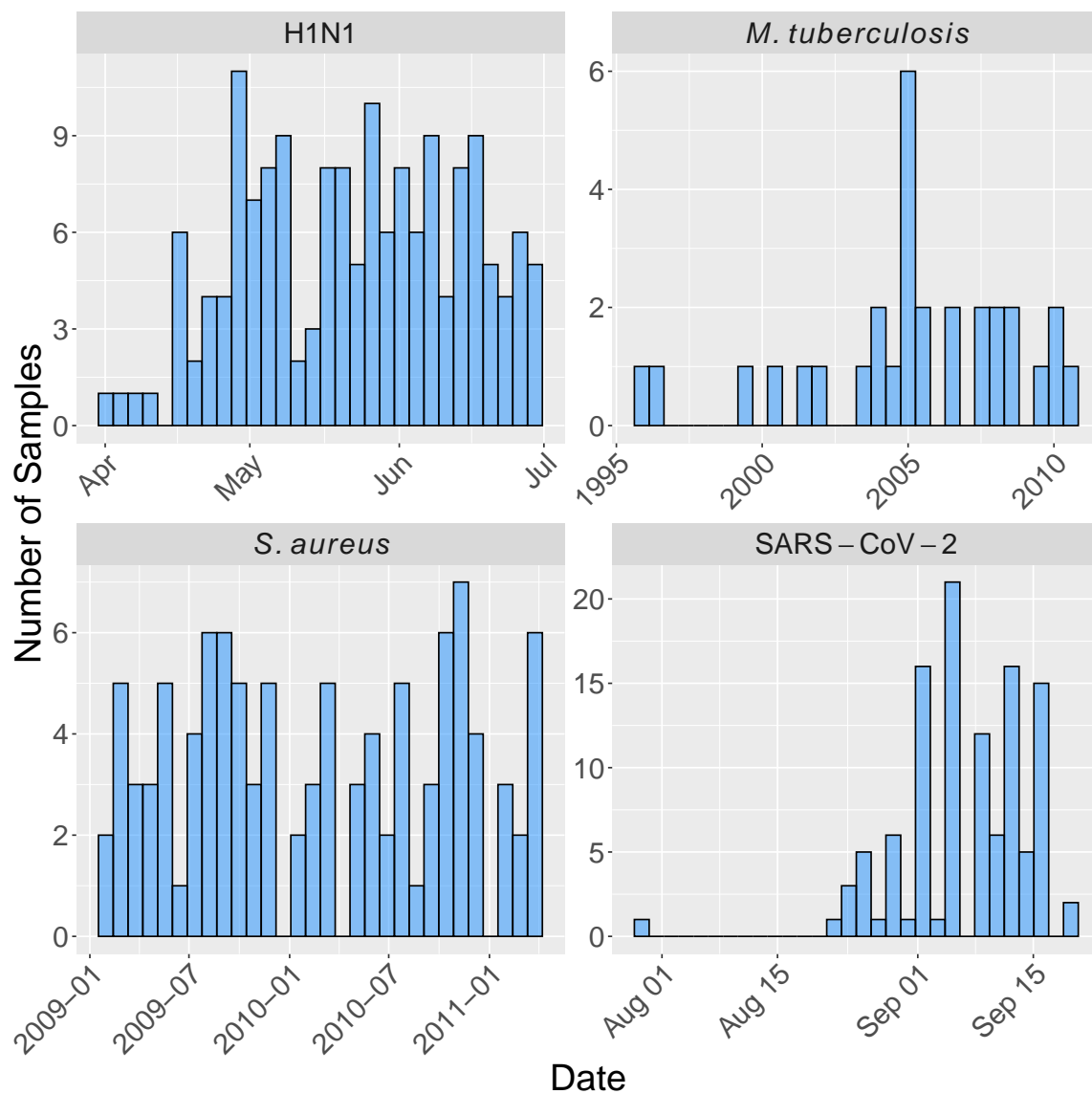


Figure S1: The number of samples over time for each empirical dataset. Date-rounding has the effect of moving each sampling within a month or year to the middle of that month or year (15th of the month or June 15th of the year).

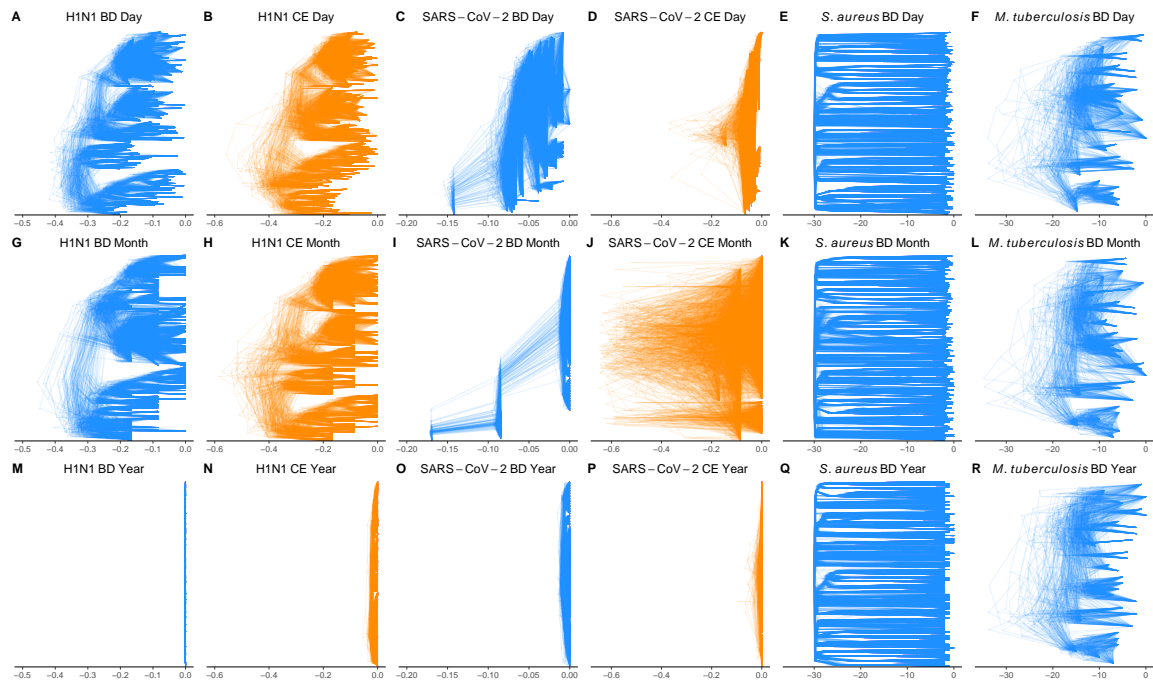


Figure S2: Desnitrees (overlaid posterior trees) for empirical data with columns corresponding to pathogen under each combination of date resolution and tree prior. For the H1N1 and SARS-CoV-2 treatments, Year resolution causes trees to collapse to instantaneous bursts.

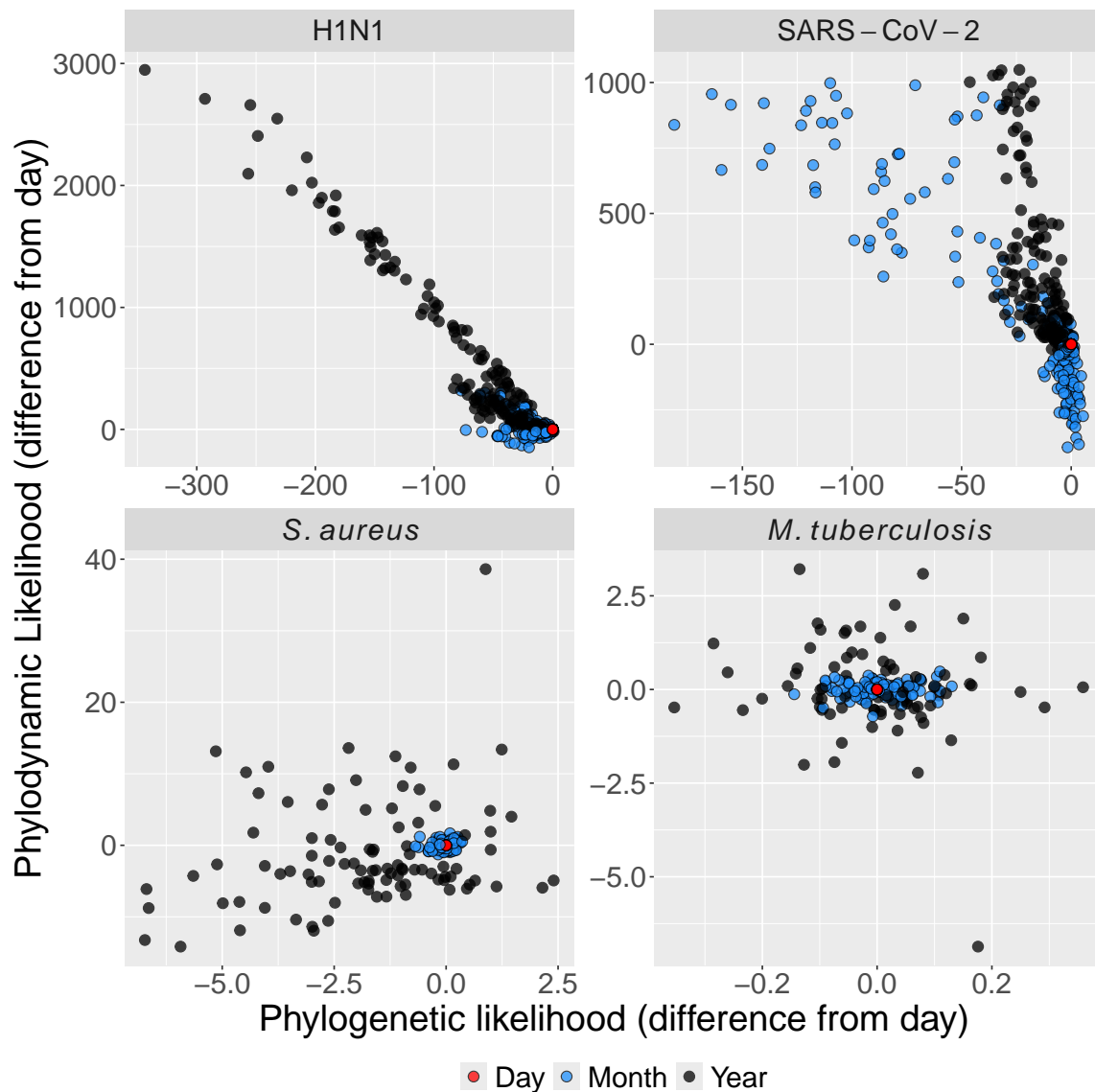


Figure S3: Adjusted phylodynamic likelihood against adjusted phylogenetic likelihood with panels corresponding to each simulation condition. Points correspond to mean posterior likelihood for each simulated dataset under each simulation condition. Colour corresponds to date resolution. Likelihoods are adjusted by subtracting the mean phylodynamic or phylogenetic likelihood at Day resolution from each the means under Month and year resolution. Resulting points therefore show the difference phylodynamic and phylogenetic likelihoods due to date-rounding with the point (0,0) representing likelihood at day resolution for each dataset. Month resolution generally results in smaller differences than Year resolution, suggesting coarser date resolution results in more perturbed likelihoods. There is also generally more error in phylodynamic likelihood than phylogenetic likelihood.

Table S1: Mean posterior estimates of substitution rate and tMRCA for empirical data with 95% HPD in brackets. The lower table gives mean posterior estimates of R_{\bullet} for empirical data with 95% HPD in brackets.

| | Tree Prior | Resolution | Substitution Rate (subs/site/yr) | tMRCA |
|------------------------|------------|------------|----------------------------------|-----------------------------|
| H1N1 | BD | Day | 4.31e-3 (3.7e-3, 4.9e-3) | 3.67e-1 (3.3e-1, 4.2e-1) |
| H1N1 | BD | Month | 3.9e-3 (3.2e-3, 4.6e-3) | 3.53e-1 (3.1e-1, 4.1e-1) |
| H1N1 | CE | Day | 3.87e-3 (3.2e-3, 4.5e-3) | 4.25e-1 (3.6e-1, 5.0e-1) |
| H1N1 | CE | Month | 3.17e-3 (2.6e-3, 3.8e-3) | 4.59e-1 (3.8e-1, 5.6e-1) |
| SARS-CoV-2 | BD | Day | 2.47e-4 (1.1e-4, 4.5e-4) | 1.45e-1 (1.4e-1, 1.5e-1) |
| SARS-CoV-2 | BD | Month | 6.56e-4 (3.3e-4, 1.1e-3) | 1.7e-1 (1.7e-1, 1.7e-1) |
| SARS-CoV-2 | CE | Day | 2.37e-4 (9.1e-5, 4.7e-4) | 2.03e-1 (1.4e-1, 3.6e-1) |
| SARS-CoV-2 | CE | Month | 4.34e-5 (4.4e-6, 1.4e-4) | 1.6 (2.9e-1, 5.9) |
| <i>S. aureus</i> | BD | Day | 1e-5 (1e-5, 1e-5) | 3e+01 (3e+01, 3e+01) |
| <i>S. aureus</i> | BD | Month | 1e-5 (1e-5, 1e-5) | 3e+01 (3e+01, 3e+01) |
| <i>S. aureus</i> | BD | Year | 1e-5 (1e-5, 1e-5) | 3e+01 (3e+01, 3e+01) |
| <i>M. tuberculosis</i> | BD | Day | 1.02e-7 (6.5e-8, 1.4e-7) | 2.17e+01 (1.7e+01, 3.2e+01) |
| <i>M. tuberculosis</i> | BD | Month | 1.02e-7 (6.6e-8, 1.4e-7) | 2.17e+01 (1.7e+01, 3.2e+01) |
| <i>M. tuberculosis</i> | BD | Year | 9.86e-8 (6.2e-8, 1.4e-7) | 2.25e+01 (1.8e+01, 3.4e+01) |

| | Tree Prior | Resolution | R_0 | R_{e1} | R_{e2} |
|------------------------|------------|------------|--------------------|--------------------|--------------------------|
| H1N1 | BD | Day | 1.08 (1.1, 1.1) | - | - |
| H1N1 | BD | Month | 1.14 (1.1, 1.2) | - | - |
| H1N1 | CE | Day | 1.14 (1.1, 1.2) | - | - |
| H1N1 | CE | Month | 1.13 (1.1, 1.2) | - | - |
| SARS-CoV-2 | BD | Day | 1.2 (9.3e-1, 1.6) | - | - |
| SARS-CoV-2 | BD | Month | 5.85 (3.7, 9.0) | - | - |
| SARS-CoV-2 | CE | Day | 1 (9.6e-1, 1.0) | - | - |
| SARS-CoV-2 | CE | Month | 1.01 (9.8e-1, 1.1) | - | - |
| <i>S. aureus</i> | BD | Day | - | 1.57 (1.5, 1.7) | 6.56e-1 (5.1e-1, 8.0e-1) |
| <i>S. aureus</i> | BD | Month | - | 1.56 (1.5, 1.7) | 6.78e-1 (5.4e-1, 8.3e-1) |
| <i>S. aureus</i> | BD | Year | - | 1.73 (1.6, 1.8) | 3.71e-1 (1.9e-1, 5.4e-1) |
| <i>M. tuberculosis</i> | BD | Day | - | 2.77 (5.8e-1, 5.3) | 1.4 (7.2e-1, 2.7) |
| <i>M. tuberculosis</i> | BD | Month | - | 2.74 (5.7e-1, 5.0) | 1.41 (7.4e-1, 2.7) |
| <i>M. tuberculosis</i> | BD | Year | - | 2.66 (4.6e-1, 5.1) | 1.53 (8.1e-1, 2.9) |