



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Featherstone, LA;Di Giallonardo, F;Holmes, EC;Vaughan, TG;Duchêne, S

Title:

Infectious disease phylodynamics with occurrence data

Date:

2021-08-01

Citation:

Featherstone, L. A., Di Giallonardo, F., Holmes, E. C., Vaughan, T. G. & Duchêne, S. (2021). Infectious disease phylodynamics with occurrence data. *Methods in Ecology and Evolution*, 12 (8), pp.1498-1507. <https://doi.org/10.1111/2041-210X.13620>.

Persistent Link:

<https://hdl.handle.net/11343/337970>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

LEO FEATHERSTONE (Orcid ID : 0000-0002-8878-1758)

Article type : Research Article

Editor : Tiago Quental

Methods in Ecology and Evolution – Research Article

Infectious disease phylodynamics with occurrence data

Leo A. Featherstone^{1*}, Francesca Di Giallonardo², Edward C. Holmes^{3,5,4}, Timothy G. Vaughan^{6,7}, Sebastián Duchêne¹

¹Department of Microbiology and Immunology, Peter Doherty Institute for Infection and Immunity, University of Melbourne, Melbourne, VIC, Australia.

²The Kirby Institute, UNSW Sydney, Sydney, NSW, Australia.

³Marie Bashir Institute for Infectious Diseases and Biosecurity, The University of Sydney, Sydney, NSW, Australia.

⁴Charles Perkins Centre, School of Life and Environmental Sciences, The University of Sydney, Sydney, NSW, Australia.

⁵School of Medical Sciences, The University of Sydney, Sydney, NSW, Australia.

⁶Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland.

⁷Swiss Institute of Bioinformatics (SIB), Switzerland.

* Contact

Email: leo.featherstone@unimelb.edu.au

Abstract (350 words max. currently 173)

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/2041-210X.13620](https://doi.org/10.1111/2041-210X.13620)

This article is protected by copyright. All rights reserved

30 1: Phylodynamic models use pathogen genome sequence data to infer epidemiological
31 dynamics. With the increasing genomic surveillance of pathogens, especially during the
32 SARS-CoV-2 pandemic, new practical questions about their use are emerging.

33

34 2: One such question focuses on the inclusion of un-sequenced case occurrence data
35 alongside sequenced data to improve phylodynamic analyses. This approach can be
36 particularly valuable if sequencing efforts vary over time.

37

38 3: Using simulations, we demonstrate that birth-death phylodynamic models can employ
39 occurrence data to eliminate bias in estimates of the basic reproductive number due to
40 misspecification of the sampling process. In contrast, the coalescent exponential model is
41 robust to such sampling biases, but in the absence of a sampling model it cannot exploit
42 occurrence data. Subsequent analysis of the SARS-CoV-2 epidemic in the northwest USA
43 supports these results.

44

45 4: We conclude that occurrence data are a valuable source of information in combination
46 with birth-death models. These data should be used to bolster phylodynamic analyses of
47 infectious diseases and other rapidly spreading species in the future.

48

49 **Key Words:** Phylodynamics, pathogens, coalescent, birth-death, Bayesian statistics

50

51 **Introduction**

52 Outbreak investigations increasingly rely on genome sequencing of causative pathogens.
53 Phylodynamic methods take advantage of these data to infer epidemiological dynamics (Rife
54 et al., 2017). New sequencing technologies generate these data rapidly, such that
55 phylodynamic inferences can be conducted in actionable time frames (Gardy & Loman,
56 2018; Grubaugh et al., 2019; Hadfield et al., 2018). In this context, the main appeal of
57 phylodynamics is that it uses sequence data to infer epidemiological dynamics preceding the
58 earliest collected sample, or during periods when sequences have not been collected, and
59 offers insight into transmission chains.

60

61 Phylodynamic models describe a branching process, modelling both how a branching
62 transmission chain and phylogenetic tree of the underlying pathogen evolve. These are
63 central to linking epidemiological dynamics to the evolution of a pathogen. In Bayesian
64 phylogenetic implementations the particular model of a branching process is part of the prior
65 and is sometimes referred to as the 'tree prior', such as the birth-death or coalescent

66 exponential. Internal nodes in the tree are associated with transmission events while the tips
67 of the tree represent sampling events (du Plessis & Stadler, 2015). Internal nodes may also
68 represent divergence distinct from transmission events in the case of free-living pathogens.
69 The basic reproductive number, R_0 , is a key parameter that reflects the average number of
70 secondary infections in a fully susceptible population. The simplest tree priors that can infer
71 R_0 posit that the number of infected individuals increases exponentially over time. Although
72 more sophisticated methods now exist (Kühnert et al., 2014; Poppinga et al., 2015;
73 Rasmussen et al., 2017; Vaughan et al., 2019; Volz & Siveroni, 2018), we focus here on tree
74 priors assuming simple exponential growth since they are appropriate for the early stages of
75 an outbreak and are increasingly used to assess the efficacy of public health interventions
76 (Geoghegan et al., 2020; Vasylyeva et al., 2019).

77
78 Two commonly used phylodynamic tree priors are the coalescent exponential and the birth-
79 death, both of which assume that the infected population size, N , grows at a rate r , $N(t)=e^{rt}$,
80 where t is time after the origin and noting that the starting population size is 1 ($N(t=0) = 1$).
81 From an epidemiological perspective, r is the difference between the transmission rate, λ ,
82 and the become uninfected rate, δ , ($r = \lambda - \delta$). An individual is infectious for on average $1/\delta$
83 units of time. R_0 is estimated as $R_0 = \lambda/\delta$. The coalescent exponential is a generalisation of
84 the Kingman n -coalescent where population size is a deterministic function of time (Griffiths
85 & Tavaré, 1994; Volz et al., 2009, 2013). In contrast, the birth-death tree prior assumes
86 stochastic population growth with sampling through time (Stadler, 2010; Stadler et al., 2012;
87 Stadler & Yang, 2013). This is captured in the death rate $\delta = \psi + \mu$, where μ is the recovery rate
88 and ψ is the sequencing rate such that the sampling proportion, p , can be calculated as $p =$
89 $\frac{\psi}{\psi + \mu}$. Note that we assume individuals become non-infectious after sequencing.

90
91 Phylodynamic analyses use sequence data and sequencing times as information (Biek et
92 al., 2015; Drummond et al., 2002, 2003; Rambaut, 2000; Rieux & Balloux, 2016). In contrast,
93 occurrence data are confirmed infectious that are sampled but not sequenced. They are
94 equivalent to a sequencing times without sequence data. In the coalescent exponential,
95 sequencing times are useful insofar as they influence the distribution of coalescent events
96 through time, influencing R_0 in turn. Coalescent models typically condition upon sequencing
97 times instead of using them to infer sequencing rates. Some ‘augmented likelihood’
98 approaches combine the coalescent with a sequencing process (Karcher et al., 2020; Parag
99 et al., 2020; Volz & Frost, 2014), but they are so-far not standard practice. For the birth-
100 death tree prior, the number of sequences and their times are naturally informative because

101 they are explicitly modelled through the sequencing rate (i.e. they inform p) (Boskova et al.,
102 2018). This difference between the two tree priors is a well understood (Stadler et al., 2015;
103 Volz & Frost, 2014), but its consequences remain to be explored in the context of
104 occurrence data. Although the amount of sequence data in outbreak investigations has
105 increased, a key consideration is that sequencing efforts are often conducted only after a
106 relatively large number of cases are reported. This latency in sequencing can bias estimates
107 of epidemiological parameters. To visualise this, the trees in Fig 1 were simulated under an
108 R_0 of 2, a constant sampling effort, and over the course of 1 year. If sequencing were only
109 conducted for samples collected after 0.75 years, samples from the deep sections of the tree
110 would be missed (*late sampling* in Fig 1). Such sequencing bias can mislead inferences of
111 epidemiological dynamics because there is no sequence data and very few branching
112 events to inform inferences of the early stages of the outbreak.

113
114 Here we consider bias in epidemiological parameters due to sequencing heterogeneity and
115 present two approaches to reduce such bias using occurrence data. The **first** approach
116 involves using a birth-death skyline tree prior that requires an understanding of the
117 sequencing effort (Stadler et al., 2013). If it is known that there was no attempt to collect
118 sequences early in the outbreak, one can set two intervals for the p parameter where one is
119 zero. However, without knowledge of sequencing effort this scenario is indistinguishable
120 from a constant sequencing effort where initial prevalence was so low as to preclude
121 obtaining any sequence data early in an outbreak. The **second** approach consists of
122 including early case occurrences in analyses, where an occurrence is a laboratory confirmed
123 case that was sampled but not sequenced (*occurrences* scenario in Fig 1). Occurrence data
124 are less expensive and more readily available than sequence data because they are
125 traditionally used in epidemiology and accurately identified via contact tracing and testing.
126 Our approach and others have been modelled, but not applied in phylodynamics hitherto
127 (Gupta et al., 2020; Manceau et al., 2020; Zarebski et al., 2020).

128
129 In a Bayesian phylogenetic framework, topological uncertainty due to occurrence data is
130 naturally incorporated into the analysis through the posterior. An analogous approach is
131 used in the fossilised birth death process in which fossils perform analogously to
132 occurrences for molecular clock calibration (Heath et al., 2014; Heath & Moore, 2014). Our
133 results further suggest that occurrence data are useful for reducing bias when inferring
134 diversification rates. However, careful consideration must be given to diversification rates
135 derived from fossil data as the data collection may violate key assumptions about the

136 sampling process, which may be more difficult to detect than in an epidemiological setting
137 and bias estimates (Matschiner, 2019).

138

139 More broadly, our inclusion of occurrences follows macroevolutionary studies employing
140 'total evidence dating' approaches combining molecular, morphological, and fossil data
141 (Zhang et al., 2016). Recent work has begun to unify such concepts across phylodynamics
142 and macroevolution through the development new tree priors which accommodate
143 occurrences and or fossil data in either context (Andréoletti et al., 2020). Our use of
144 occurrence data further affirms that additional sources of information, such as drug-
145 resistance profiles may help inform epidemiological inference in the future. It can be
146 expected that Bayesian implementations such as BEAST and RevBayes will continue to
147 provide platforms for these techniques in the future (Bouckaert et al., 2019; Höhna et al.,
148 2016; Suchard et al., 2018).

149

150 **Materials and Methods**

151 *Simulation study*

152 We simulated phylogenetic trees under a birth-death process in MASTER v6.1 (Vaughan &
153 Drummond, 2013), with the following parameterisation; $R_0 = 2$, $\delta = 91$, $p = 0.05$, and an
154 outbreak duration of one year ($\frac{1}{\delta} = 0.011$ years corresponding to an expected infectious
155 period of about 4 days). We chose the birth-death model over the coalescent exponential,
156 because the sampling process is a natural by-product of the process, whereas it would need
157 to be conditioned upon *a priori* under the coalescent exponential. This approach simulates
158 an outbreak from which on average one in every twenty cases is randomly sampled (equiv.
159 $p = 0.05$), with these samples becoming tips in the tree. This is representative of empirical
160 data sets where identical, epidemiologically related sequences are often removed resulting
161 in data sets that approach a well-mixed population with homogeneity in the sequencing
162 proportion (Rambaut, 2020). We filtered for trees between 100 to 150 tips, otherwise
163 allowing their number and ages to vary naturally.

164

165 We then assumed a strict molecular clock with an evolutionary rate of 0.01 substitutions per
166 site per year (subs/site/year) and the HKY+ Γ substitution model to produce alignments of
167 13,000 nucleotides corresponding to each tree using NELSI (Ho et al., 2015) and Phangorn
168 v2.4 (Schliep, 2011). These settings are broadly similar to an influenza virus outbreak
169 (Hedge et al., 2013), but a rescaling of the epidemiological parameters could apply to many
170 other pathogens. We then assumed three scenarios: (i) *constant* sequencing with all

171 sequences from the simulation included (e.g. the sequence for every sample in the tree in
172 Fig 1 is included), (ii) *late sequencing* only with sequences after time T_s (e.g. only sequences
173 for samples after the dashed line in the tree in Fig 1), and (iii) *occurrences* in which
174 sequence data are available only after time T_s with those preceding recorded as
175 occurrences. We set T_s to 0.75. For each parameter configuration we simulated 100
176 sequence data sets which were subsampled according to the three scenarios above.
177 Occurrences were emulated by replacing simulated DNA sequences with 'n' (i.e. missing
178 data) in the alignment. This resulted in a median of 556.5 polymorphic sites in constant
179 sequencing alignments, and 535.5 in the late sequencing and occurrences alignments.

180
181 We analysed the data in BEAST v2.5 (Bouckaert et al., 2019) with coalescent exponential,
182 birth-death, and birth-death skyline tree priors (Table 1). Our results focus on the birth-death,
183 but the coalescent exponential forms a valuable point of comparison through its robustness
184 to variation in sampling. For the birth-death skyline, we set two intervals for the p parameter,
185 with the interval time fixed at T_s . We matched the substitution and clock model to the values
186 used to generate the data and we used an informative Γ prior distribution with mean set to
187 the true value of 91 and standard deviation of 1 for δ . For R_0 we used a lognormal prior with
188 $\mu = 0$ and $\sigma = 1$. For the coalescent exponential, we placed a one-on-x prior on the scaled
189 effective population size and a Laplace prior on the growth rate with $\mu = 0.001$ and scale =
190 30.7.

191
192 We assessed the effectiveness of each analysis treatment using three statistics. First, we
193 considered the coverage as a measure of accuracy: the number of times the 95% highest
194 posterior density (HPD) intervals covered the true value of a given parameter. Second, we
195 consider average absolute error, which is the absolute difference between the posterior
196 mean and true mean for a given parameter averaged across the 100 simulations for each
197 sampling treatment. Third, we consider average 95% HPD width for each treatment, as a
198 measure of precision.

199 200 *Empirical case study*

201 To illustrate the accuracy of occurrence data relative to completely sequenced data sets we
202 analysed 821 whole genome sequences sampled from the SARS-CoV-2 pandemic from
203 Washington State, USA, and the adjacent Washington County, Oregon. These were
204 downloaded from GISAID (Supplementary material) and partially documented by (Bedford et
205 al., 2020). This data set was chosen because it presented a well-sequenced outbreak from
206 which we could take a set of sequences that grew exponentially in time. This best fits the

207 coalescent exponential and birth-death tree priors such that we could test the impact of
208 occurrence data clearly. To arrive at the data set of 821 sequences, we downloaded 2,164
209 high-coverage genome sequences collected between January 18th and June 30th 2020, but
210 selected the 821 sequences taken up to March 21st 2020 (~50% of confirmed cases) to
211 capture exponential growth in cases and sequences, matching our simulated epidemic
212 trajectories and lending itself to a consistent sequencing proportion (Fig S1). We then
213 corroborated exponential growth in the underlying population using an Epoch Sampling
214 Proportion Skyline Plot (Parag et al., 2020). We further divided this data set into five subsets
215 as per our simulation study: (i) 'complete sequencing including all 821 sequences; (ii) late
216 sequencing post March 6th 2020 (decimal date 2020.18) including 637 sequences; (iii) late
217 sequencing post March 14th 2020 (2020.20) including 340 sequences; (iv) late sequencing
218 post March 6th 2020 (2020.18) including 637 sequences and 184 occurrences; and (v) late
219 sequencing post 2020.2 including 340 sequences and 481 occurrences. Including two late
220 sequencing data subsets offers information about how inflation in R_0 varies with latency in
221 sequencing.

222
223 We then analysed each data set with each tree prior used in the simulation study with
224 BEASTv2.5. We first employed a birth-death model with serial sequencing. We placed a
225 lognormal prior on R_0 with mean 0 and standard deviation of 1; fixed δ at 36.5 (i.e. 10-day
226 duration of infection as estimated recently (Price et al., 2020)); a β prior on sampling
227 proportion with shape and scale equal to 2 to penalise extreme values. We did not use a
228 hyperprior on the age of the root node, such that it was determined by the birth-death
229 process. Second, we used a birth-death skyline model with the same priors as the birth-
230 death, but with two sequencing rate parameters. The first pertained to after the 2020.18 or
231 2020.2 cut-off, and the second to before the cut-off. Both used the same beta prior for
232 sampling proportion as for the birth-death. Third, a coalescent exponential tree prior was
233 used with a Laplace prior on growth rate with mean 0 and scale 100 and an exponential prior
234 with mean 100 on the coalescent exponential effective population size (ϕ). All priors are
235 listed in Table 1. For both tree priors, we assumed HKY+ Γ substitution model with a strict
236 molecular clock rate fixed to 10^{-3} subs/site/year, following recent estimates (Duchene et al.,
237 2020). We ran a Markov chain Monte Carlo of 5×10^8 steps, sampling at every 1000th step.
238 We determined sufficient sampling from the posterior by verifying that the effective sample
239 size all parameters of interest was above 200.

240

241 **Results**

242 *Simulation study*

243 Analyses of data sets with late sequencing using the birth-death model were least accurate
244 in estimating R_0 . In only 12 of 100 simulations did the 95% HPD include 2 (Table 2 and Fig
245 2b). The birth-death skyline was more accurate with 95 and 92 of 100 simulations covering
246 $R_0 = 2$. The coalescent exponential was also more accurate with all 100 HPD intervals
247 covering $R_0 = 2$. However, this came at the cost of low precision as HPD width was the
248 largest for the coalescent out of all treatments.

249

250 In general, we observed that the birth-death model tended to overestimate R_0 while the
251 coalescent exponential underestimated it for data sets with late sequencing (Fig 2).

252 Estimates of the evolutionary rate displayed an identical pattern to those of R_0 , with the
253 coalescent exponential and the birth-death model being the most and least accurate
254 respectively at the expense of precision. However, the evolutionary rate appeared overall
255 robust to the choice of the tree prior, with the only treatment producing a less than 90%
256 coverage being the birth-death model with late sequencing. This is a valuable consideration
257 for analyses of future outbreaks as considerable attention is initially devoted to estimating a
258 reliable evolutionary rate for a given pathogen because this is key to phylodynamic inference
259 (Duchene et al., 2020).

260

261 As expected, analyses of the data with constant sampling were accurate in a majority of
262 cases, with 94, 95, and 89 out of 100 simulations covering R_0 , alongside 94 and 92 for the
263 evolutionary rate under the birth-death, birth-death skyline, and coalescent exponential
264 models, respectively. The true model is the birth-death, and as such it is expected to perform
265 better than the coalescent. Estimates of R_0 including occurrence data were similar in
266 accuracy to those with complete sampling. A total of 94 analyses correctly estimated this
267 parameter under the birth-death model, 97 under the birth-death skyline, and 96 analyses
268 included the true value for the coalescent exponential. Evolutionary rate estimates with
269 occurrence data were similar, with 95 accurate estimates using the birth-death model, 93
270 under the birth-death skyline, and 91 using the coalescent exponential (Table 2, Fig 2).

271 These results are attributable to the fact that the birth-death model treats sequencing times
272 as data, whereas the coalescent exponential model conditions on the number of samples
273 and their ages (Boskova et al., 2018; Stadler et al., 2015). In other words, while sequencing
274 times partially inform the root height for each, the birth-death and birth-death skyline go on to
275 use sequencing times as information while the coalescent draws no further information from
276 them. In the birth-death and birth-death skyline models, occurrence data improve accuracy
277 for R_0 and are also informative about the age of the tree height under this tree prior, which

278 can also improve the accuracy of the evolutionary rate relative to the coalescent exponential
279 model. But these estimates are unlikely to be as accurate as those with complete sequence
280 data because they necessarily include less information.

281

282 The coalescent exponential model appears to be more robust to the sampling treatment,
283 with greater accuracy than the birth-death and birth-death skyline model across late
284 sequencing and occurrence treatments. Our simulations suggest that this comes at the
285 expense of less precise estimates than those from birth-death models (Table 2). In turn,
286 birth-death and birth-death skyline models tend to produce more precise estimates with less
287 error (Table 2, Fig 2, Fig S2). Together these results suggest that in a genomic-reporting
288 scenario, the coalescent exponential is suitable when sequencing proportion is assumed to
289 be low, when the sequencing process is otherwise poorly understood, or when reliable
290 occurrence data are not available. However, when increased precision is desirable and
291 occurrence data are available, birth-death tree priors may provide the sharper estimates with
292 comparable accuracy. The choice of tree prior could be optimised depending on prioritisation
293 of precision and error based on the ordering of bars in Figure S2.

294

295 *Empirical case study: SARS-CoV-2 from the northwest USA*

296 Mirroring trends in our simulated data sets, the coalescent exponential returned consistent
297 estimates of R_0 across treatments which were generally lower than those inferred by the
298 birth-death tree prior (Fig 3a, Table 3). Coalescent exponential treatments again produced
299 wider HPD intervals than birth-death and birth-death skyline treatments, with the exception
300 of late sequencing under the birth-death as expected from simulations. Uncertainty in
301 posterior R_0 does not appear to change when substituting sequence data for occurrence
302 data (Fig 3A), indicating that late samples are highly informative while occurrence data
303 contribute relatively little additional information to coalescent analyses. Moreover, we
304 observed a near perfect match between estimates from analyses with only late sequencing
305 and those that included occurrences. This pattern can be explained because occurrence
306 data have no influence on marginal posterior estimates under the coalescent. By contrast,
307 our simulations show small differences in performance between coalescent analyses with
308 late sequencing and those with occurrence data, which we attribute to noise in the
309 simulation study.

310

311 The results of the birth-death analyses recapitulate our observation from simulations that
312 later sequencing inflates estimates of R_0 , and that occurrence data rectify this (Figure 3B
313 table 3). Complete sequencing estimated a mean R_0 of 1.96 (95% HPD: [1.85, 2.07]) and

314 late sequencing with occurrence data gave means of 1.95 and 2.00 (95% HPDs: [1.8, 2.11]
315 and [1.9, 2.12] for post - 2020.18 and 2020.2, respectively). These estimates are slightly
316 lower than those from earlier work to estimate R_0 in the Washington state epidemic
317 (Vaughan et al., 2020). This discrepancy may be due to the former being conducted earlier
318 when the virus may have been spreading more rapidly. Late sequencing alone inferred a
319 mean R_0 of 2.44 and 3.53 for post 2020.18 and 2020.2 ([2.31, 2.58] and [3.24, 3.82] 95%
320 HPDs respectively). The way in which the latest sequencing dataset inferred the highest
321 values of R_0 further suggests that upward bias increases with lateness in sequencing.
322 Interestingly, the birth-death skyline produced slightly higher estimates than constant
323 sequencing when occurrences were included (Fig 3). This may be attributable to
324 occurrences with similar sampling times clustering phylogenetically, therefore inflating
325 posterior R_0 in lieu of a genomic signal to break them apart.

326

327 In both late sequencing treatments, the birth-death skyline posterior R_0 distributions were
328 lower than their equivalents under the standard birth-death model, with later sequencing
329 corresponding to lower estimates (Fig 3). This is consistent with the simulated data (Fig 2)
330 and suggests that including occurrence data is a preferential strategy to rectify posterior R_0
331 estimates amid late genome sequence sampling. Furthermore, the entropy of each birth-
332 death based posterior R_0 distribution, a measure of uncertainty, is comparable at 3.68-3.78
333 as calculated with the mlf R package (Peterson, 2018). This further suggests that the
334 topological uncertainty induced by occurrence data does not considerably increase
335 uncertainty in posterior R_0 (Fig 3).

336

337 **Discussion**

338 *Occurrence data in empirical phylodynamic studies*

339 Occurrence data present an extreme case of there being no genome coverage in samples
340 as opposed to complete or partial coverage. Our results therefore extend to show that low-
341 coverage or partially overlapping sequences can be useful in phylodynamics so long as they
342 are accurate, with their informative capacity lying somewhere between that of the
343 occurrence and complete sequence. An outstanding task is to characterise an upper-bound
344 on the relative proportion of occurrence to genomic samples from which genomic samples
345 can still inform tree topology for epidemiological dynamics. To this end, we caution against
346 over-inflating occurrences among genomic data sets without comparing to results obtained
347 with genomic samples alone.

348

349 The simulated and empirical data sets we consider present scenarios where the balance of
350 occurrences to sequences is at most slightly above equal (~59% in late sequencing post
351 2020.2). This is accurate for many sequencing scenarios, such as during the SARS-CoV-2
352 where a high proportion of cases have been sequenced in settings such as New Zealand,
353 Iceland, and Australia (Duchene et al., 2020; Geoghegan et al., 2020; Seemann et al.,
354 2020). Of course, sequencing proportion may also be significantly lower than this, such that
355 the extent to which sequences inform estimates in data sets largely composed of
356 occurrences is questionable. This is a complex question incorporating the relative
357 proportions of occurrences and sequences, their temporal diversity, and genetic diversity of
358 the latter. In such a situation, we suggest running analyses with sequence information
359 removed (i.e. converting every sequence to an occurrence) to make a decision as to
360 whether sequences drive a signal in the data and phylodynamic analyses are justified.

361

362 A related question is how heterogeneity in occurrence sampling skews results. In our
363 simulations we created a scenario of homogeneity – a constant rate – of sampling as this is
364 pertinent to many sequencing scenarios such as for individual transmission clusters that are
365 sequenced relatively consistently. Noting this, we stress that it is not appropriate to include
366 all known occurrences alongside a set of sequences to infer epidemiological parameters.
367 Rather, occurrences must be curated in a similar way to sequences in that they must come
368 from an epidemiologically related source to avoid skewing estimates due extraneous
369 epidemiological dynamics.

370

371 Our simulations and empirical data analyses reveal that occurrence data are a rich source of
372 information for birth-death tree priors that can dramatically improve the accuracy and
373 precision in estimates of epidemiological parameters. A key consideration is that
374 occurrences should represent confirmed cases that would have been sequenced if
375 sequencing effort had been constant, and which are known to belong to a particular
376 outbreak, such as via contact tracing. Combining occurrence and sequence data can be
377 particularly useful in situations where it is unknown if sequence sampling has been constant
378 over time or where there exist several confirmed cases but a smaller number of sequences.
379 This is valuable amid recently emerging outbreaks where combining both sources of data
380 can provide sharper and more timely insight into the recent evolution of the pathogen in
381 question.

382

383 **Acknowledgements**

384 We thank the Editor and three reviewers for thoughtful comments on previous versions of
385 this manuscript. We are also grateful to Trevor Bedford, authors and groups, originating
386 laboratories, and submitting laboratories of sequences downloaded from GISAID. We
387 provide a full acknowledgement table in the supplementary information. SD and LAF were
388 supported by a Discovery Early Career Fellowship from the Australian Research Council
389 (DE190100805), awarded to SD. ECH is supported by an Australian Research Council
390 Laureate Fellowship (FL170100022).

391

392 **Authors' contributions:**

393 All authors contributed to the design of experiments and writing of the manuscript. LF
394 conducted analyses of empirical data and lead writing of the manuscript. FG contributed
395 initial datasets, writing, and guidance with figures. TV contributed to writing the manuscript
396 and mathematical concepts. EH contributed to writing of the manuscript and original ideas.
397 SD conceived of fundamental concepts in the manuscript, conducted simulations, and
398 contributed to writing.

399

400 **Data availability:**

401 Input files to generate trees in MASTER and to analyse sequence data in BEAST according
402 to the birth-death skyline, birth-death, and the coalescent exponential tree priors, and
403 accession numbers for empirical SARS-CoV-2 virus data. Available at:
404 github.com/LeoFeatherstone/occurrences (<http://doi.org/10.5281/zenodo.4655304>)
405 (Featherstone, 2021). Accession numbers for empirical SARS-CoV-2 virus data and the
406 GISAID acknowledgements table are available as supplementary data online.

407

408 **References**

- 409 Andréoletti, J., Zwaans, A., Warnock, R. C. M., Aguirre-Fernández, G., Barido-Sottani, J.,
410 Gupta, A., Stadler, T., & Manceau, M. (2020). A skyline birth-death process for
411 inferring the population size from a reconstructed tree with occurrences. *BioRxiv*,
412 2020.10.27.356758. <https://doi.org/10.1101/2020.10.27.356758>
- 413 Bedford, T., Greninger, A. L., Roychoudhury, P., Starita, L. M., Famulare, M., Huang, M.-L.,
414 Nalla, A., Pepper, G., Reinhardt, A., Xie, H., Shrestha, L., Nguyen, T. N., Adler, A.,
415 Brandstetter, E., Cho, S., Giroux, D., Han, P. D., Fay, K., Frazar, C. D., ... Jerome, K. R.
416 (2020). *Cryptic transmission of SARS-CoV-2 in Washington State* [Preprint].
417 *Epidemiology*. <https://doi.org/10.1101/2020.04.02.20051417>

- 418 Biek, R., Pybus, O. G., Lloyd-Smith, J. O., & Didelot, X. (2015). Measurably evolving
419 pathogens in the genomic era. *Trends in Ecology and Evolution*, 30(6), 306–313.
- 420 Boskova, V., Stadler, T., & Magnus, C. (2018). The influence of phylodynamic model
421 specifications on parameter estimates of the Zika virus epidemic. *Virus Evolution*,
422 4(1), vex044.
- 423 Bouckaert, R., Vaughan, T. G., Barido-Sottani, J., Duchêne, S., Fourment, M., Gavryushkina,
424 A., Heled, J., Jones, G., Kühnert, D., Maio, N. D., Matschiner, M., Mendes, F. K.,
425 Müller, N. F., Ogilvie, H. A., Plessis, L. du, Poppinga, A., Rambaut, A., Rasmussen, D.,
426 Siveroni, I., ... Drummond, A. J. (2019). BEAST 2.5: An advanced software platform for
427 Bayesian evolutionary analysis. *PLOS Computational Biology*, 15(4), e1006650.
428 <https://doi.org/10.1371/journal.pcbi.1006650>
- 429 Drummond, A. J., Nicholls, G. K., Rodrigo, A. G., & Solomon, W. (2002). Estimating mutation
430 parameters, population history and genealogy simultaneously from temporally
431 spaced sequence data. *Genetics*, 161(3), 1307–1320.
- 432 Drummond, A. J., Pybus, O. G., Rambaut, A., Forsberg, R., & Rodrigo, A. G. (2003).
433 Measurably evolving populations. *Trends in Ecology & Evolution*, 18(9), 481–488.
- 434 du Plessis, L., & Stadler, T. (2015). Getting to the root of epidemic spread with phylodynamic
435 analysis of genomic data. *Trends in Microbiology*, 23(7), 383–386.
- 436 Duchene, S., Featherstone, L., Blasio, B. F. de, Holmes, E. C., Bohlin, J., & Pettersson, J. H.-O.
437 (2020). The impact of early public health interventions on SARS-CoV-2 transmission
438 and evolution. *MedRxiv*, 2020.11.18.20233767.
439 <https://doi.org/10.1101/2020.11.18.20233767>
- 440 Duchene, S., Featherstone, L., Haritopoulou-Sinanidou, M., Rambaut, A., Lemey, P., & Baele,
441 G. (2020). Temporal signal and the phylodynamic threshold of SARS-CoV-2. *BioRxiv*,
442 2020.05.04.077735. <https://doi.org/10.1101/2020.05.04.077735>
- 443 Featherstone, L. A. (2021). *LeoFeatherstone/occurrences: First release*.
444 <https://doi.org/10.5281/zenodo.4655304>
- 445 Gardy, J. L., & Loman, N. J. (2018). Towards a genomics-informed, real-time, global
446 pathogen surveillance system. *Nature Reviews Genetics*, 19(1), 9.
- 447 Geoghegan, J. L., Ren, X., Storey, M., Hadfield, J., Jelley, L., Jefferies, S., Sherwood, J., Paine,
448 S., Huang, S., Douglas, J., Mendes, F. K. L., Sporle, A., Baker, M. G., Murdoch, D. R.,
449 French, N., Simpson, C. R., Welch, D., Drummond, A. J., Holmes, E. C., ... de Ligt, J.

- 450 (2020). *Genomic epidemiology reveals transmission patterns and dynamics of SARS-*
451 *CoV-2 in Aotearoa New Zealand* [Preprint]. *Infectious Diseases (except HIV/AIDS)*.
452 <https://doi.org/10.1101/2020.08.05.20168930>
- 453 Griffiths, R. C., & Tavaré, S. (1994). Sampling theory for neutral alleles in a varying
454 environment. *Philosophical Transactions of the Royal Society of London. Series B:*
455 *Biological Sciences*, 344(1310), 403–410. <https://doi.org/10.1098/rstb.1994.0079>
- 456 Grubaugh, N. D., Ladner, J. T., Lemey, P., Pybus, O. G., Rambaut, A., Holmes, E. C., &
457 Andersen, K. G. (2019). Tracking virus outbreaks in the twenty-first century. *Nature*
458 *Microbiology*, 4(1), 10.
- 459 Gupta, A., Manceau, M., Vaughan, T., Khammash, M., & Stadler, T. (2020). The probability
460 distribution of the reconstructed phylogenetic tree with occurrence data. *Journal of*
461 *Theoretical Biology*, 488, 110115. <https://doi.org/10.1016/j.jtbi.2019.110115>
- 462 Hadfield, J., Megill, C., Bell, S. M., Huddleston, J., Potter, B., Callender, C., Sagulenko, P.,
463 Bedford, T., & Neher, R. A. (2018). Nextstrain: Real-time tracking of pathogen
464 evolution. *Bioinformatics*, 34(23), 4121–4123.
- 465 Heath, T. A., Huelsenbeck, J. P., & Stadler, T. (2014). The fossilized birth–death process for
466 coherent calibration of divergence-time estimates. *Proceedings of the National*
467 *Academy of Sciences*, 111(29), E2957–E2966.
- 468 Heath, T. A., & Moore, B. R. (2014). Bayesian inference of species divergence times. In M.-H.
469 Chen, L. Kuo, & P. O. Lewis (Eds.), *Bayesian Phylogenetics, Methods, Algorithms, and*
470 *Applications* (pp. 277–318). CRC Press.
- 471 Hedge, J., Lycett, S. J., & Rambaut, A. (2013). Real-time characterization of the molecular
472 epidemiology of an influenza pandemic. *Biology Letters*, 9(5), 20130331.
- 473 Ho, S. S., Duchêne, S., & Duchêne, D. A. (2015). Simulating and detecting autocorrelation of
474 molecular evolutionary rates among lineages. *Molecular Ecology Resources*, 15(4),
475 688–696. <https://doi.org/10.1111/1755-0998.12320>
- 476 Höhna, S., Landis, M. J., Heath, T. A., Boussau, B., Lartillot, N., Moore, B. R., Huelsenbeck, J.
477 P., & Ronquist, F. (2016). RevBayes: Bayesian Phylogenetic Inference Using Graphical
478 Models and an Interactive Model-Specification Language. *Systematic Biology*, 65(4),
479 726–736. <https://doi.org/10.1093/sysbio/syw021>
- 480 Karcher, M. D., Carvalho, L. M., Suchard, M. A., Dudas, G., & Minin, V. N. (2020). Estimating
481 effective population size changes from preferentially sampled genetic sequences.

- 482 *PLOS Computational Biology*, 16(10), e1007774.
483 <https://doi.org/10.1371/journal.pcbi.1007774>
- 484 Kühnert, D., Stadler, T., Vaughan, T. G., & Drummond, A. J. (2014). Simultaneous
485 reconstruction of evolutionary history and epidemiological dynamics from viral
486 sequences with the birth–death SIR model. *Journal of The Royal Society Interface*,
487 11(94), 20131106. <https://doi.org/10.1098/rsif.2013.1106>
- 488 Manceau, M., Gupta, A., Vaughan, T., & Stadler, T. (2020). The probability distribution of the
489 ancestral population size conditioned on the reconstructed phylogenetic tree with
490 occurrence data. *Journal of Theoretical Biology*, 110400.
491 <https://doi.org/10.1016/j.jtbi.2020.110400>
- 492 Matschiner, M. (2019). Selective Sampling of Species and Fossils Influences Age Estimates
493 Under the Fossilized Birth–Death Model. *Frontiers in Genetics*, 10.
494 <https://doi.org/10.3389/fgene.2019.01064>
- 495 Parag, K. V., du Plessis, L., & Pybus, O. G. (2020). Jointly Inferring the Dynamics of Population
496 Size and Sampling Intensity from Molecular Sequences. *Molecular Biology and*
497 *Evolution*, 37(8), 2414–2429. <https://doi.org/10.1093/molbev/msaa016>
- 498 Peterson, K. (2018). *mlf: Machine Learning Foundations* (1.2.1) [Computer software].
499 <https://CRAN.R-project.org/package=mlf>
- 500 Poppinga, A., Vaughan, T., Stadler, T., & Drummond, A. J. (2015). Inferring Epidemiological
501 Dynamics with Bayesian Coalescent Inference: The Merits of Deterministic and
502 Stochastic Models. *Genetics*, 199(2), 595–607.
503 <https://doi.org/10.1534/genetics.114.172791>
- 504 Price, D. J., Shearer, F. M., Meehan, M. T., McBryde, E., Moss, R., Golding, N., Conway, E. J.,
505 Dawson, P., Cromer, D., Wood, J., Abbott, S., McVernon, J., & McCaw, J. M. (2020).
506 *Early analysis of the Australian COVID-19 epidemic* [Preprint]. *Epidemiology*.
507 <https://doi.org/10.1101/2020.04.25.20080127>
- 508 Rambaut, A. (2020). *Phylodynamic Analysis | 176 genomes | 6 Mar 2020—SARS-CoV-2*
509 *coronavirus / nCoV-2019 Genomic Epidemiology*. *Virological*.
510 <https://virological.org/t/phylodynamic-analysis-176-genomes-6-mar-2020/356>
- 511 Rambaut, Andrew. (2000). Estimating the rate of molecular evolution: Incorporating non-
512 contemporaneous sequences into maximum likelihood phylogenies. *Bioinformatics*,
513 16(4), 395–399. <https://doi.org/10.1093/bioinformatics/16.4.395>

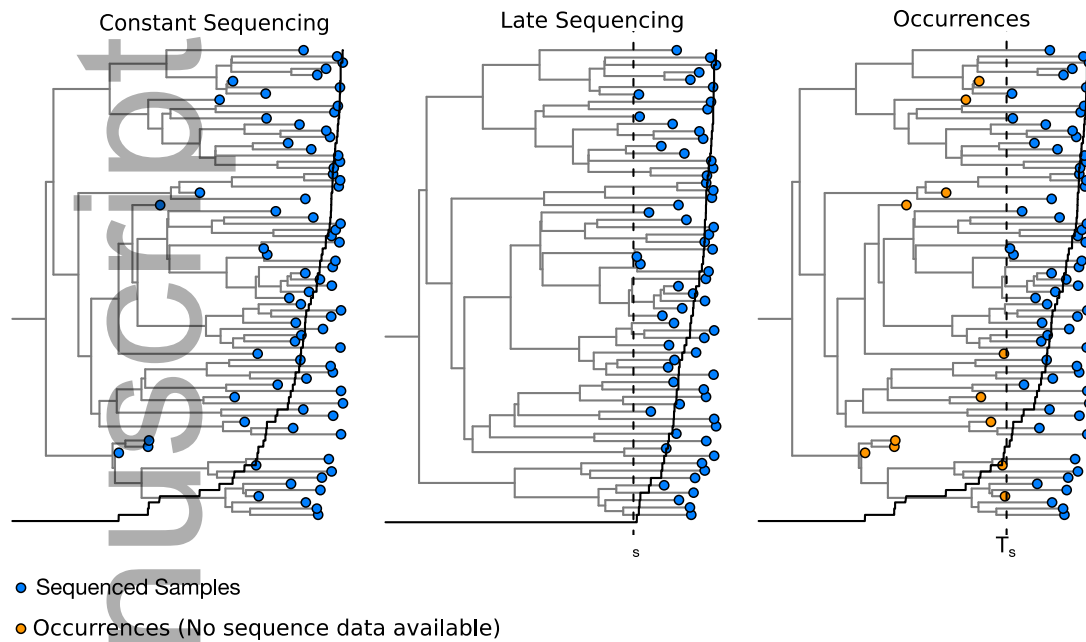
- 514 Rasmussen, D. A., Kouyos, R., Günthard, H. F., & Stadler, T. (2017). Phylodynamics on local
515 sexual contact networks. *PLOS Computational Biology*, *13*(3), e1005448.
516 <https://doi.org/10.1371/journal.pcbi.1005448>
- 517 Rieux, A., & Balloux, F. (2016). Inferences from tip-calibrated phylogenies: A review and a
518 practical guide. *Molecular Ecology*, *25*(9), 1911–1924.
- 519 Rife, B. D., Mavian, C., Chen, X., Ciccozzi, M., Salemi, M., Min, J., & Prosperi, M. C. (2017).
520 Phylodynamic applications in 21st century global infectious disease research. *Global*
521 *Health Research and Policy*, *2*(1), 13. <https://doi.org/10.1186/s41256-017-0034-y>
- 522 Schliep, K. P. (2011). phangorn: Phylogenetic analysis in R. *Bioinformatics*, *27*(4), 592–593.
- 523 Seemann, T., Lane, C. R., Sherry, N. L., Duchene, S., Gonçalves da Silva, A., Caly, L., Sait, M.,
524 Ballard, S. A., Horan, K., Schultz, M. B., Hoang, T., Easton, M., Dougall, S., Stinear, T.
525 P., Druce, J., Catton, M., Sutton, B., van Diemen, A., Alprent, C., ... Howden, B. P.
526 (2020). Tracking the COVID-19 pandemic in Australia using genomics. *Nature*
527 *Communications*, *11*(1), 4376. <https://doi.org/10.1038/s41467-020-18314-x>
- 528 Stadler, T. (2010). Sampling-through-time in birth-death trees. *Journal of Theoretical*
529 *Biology*, *167*(3), 696–404.
- 530 Stadler, T., Kouyos, R., von Wyl, V., Yerly, S., Böni, J., Bürgisser, P., Klimkait, T., Joos, B.,
531 Rieder, P., & Xie, D. (2012). Estimating the basic reproductive number from viral
532 sequence data. *Molecular Biology and Evolution*, *29*(1), 347–357.
- 533 Stadler, T., Kühnert, D., Bonhoeffer, S., & Drummond, A. J. (2013). Birth–death skyline plot
534 reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV).
535 *Proceedings of the National Academy of Sciences*, *110*(1), 228–233.
- 536 Stadler, T., Vaughan, T. G., Gavryushkin, A., Guindon, S., Kühnert, D., Leventhal, G. E., &
537 Drummond, A. J. (2015). How well can the exponential-growth coalescent
538 approximate constant-rate birth–death population dynamics? *Proceedings of the*
539 *Royal Society B: Biological Sciences*, *282*(1806), 20150420.
- 540 Stadler, T., & Yang, Z. (2013). Dating phylogenies with sequentially sampled tips. *Systematic*
541 *Biology*, *62*(5), 674–688.
- 542 Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J., & Rambaut, A. (2018).
543 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus*
544 *Evolution*, *4*(1). <https://doi.org/10.1093/ve/vey016>

- 545 Vasylyeva, T. I., du Plessis, L., Pineda-Peña, A. C., Kühnert, D., Lemey, P., Vandamme, A.-M.,
546 Gomes, P., Camacho, R. J., Pybus, O. G., Abecasis, A. B., & Faria, N. R. (2019). Tracing
547 the Impact of Public Health Interventions on HIV-1 Transmission in Portugal Using
548 Molecular Epidemiology. *The Journal of Infectious Diseases*, *220*(2), 233–243.
549 <https://doi.org/10.1093/infdis/jiz085>
- 550 Vaughan, T. G., & Drummond, A. J. (2013). A stochastic simulator of birth–death master
551 equations with application to phylodynamics. *Molecular Biology and Evolution*,
552 *30*(6), 1480–1493.
- 553 Vaughan, T. G., Leventhal, G. E., Rasmussen, D. A., Drummond, A. J., Welch, D., & Stadler, T.
554 (2019). Estimating Epidemic Incidence and Prevalence from Genomic Data.
555 *Molecular Biology and Evolution*, *36*(8), 1804–1816.
556 <https://doi.org/10.1093/molbev/msz106>
- 557 Vaughan, T. G., Nadeau, S. A., Sciré, J., & Stadler, T. (2020, March 13). Phylodynamic
558 Analyses of outbreaks in China, Italy, Washington State (USA), and the Diamond
559 Princess. *Virological.Org*. [https://virological.org/t/phylodynamic-analyses-of-](https://virological.org/t/phylodynamic-analyses-of-outbreaks-in-china-italy-washington-state-usa-and-the-diamond-princess/439)
560 [outbreaks-in-china-italy-washington-state-usa-and-the-diamond-princess/439](https://virological.org/t/phylodynamic-analyses-of-outbreaks-in-china-italy-washington-state-usa-and-the-diamond-princess/439)
- 561 Volz, E. M., & Frost, S. D. W. (2014). Sampling through time and phylodynamic inference
562 with coalescent and birth–death models. *Journal of the Royal Society Interface*,
563 *11*(101). <https://doi.org/10.1098/rsif.2014.0945>
- 564 Volz, E. M., Koelle, K., & Bedford, T. (2013). Viral phylodynamics. *PLOS Computational*
565 *Biology*, *9*(3), e1002947.
- 566 Volz, E. M., Pond, S. L. K., Ward, M. J., Brown, A. J. L., & Frost, S. D. W. (2009).
567 Phylodynamics of infectious disease epidemics. *Genetics*, *183*(4), 1421–1430.
- 568 Volz, E. M., & Siveroni, I. (2018). Bayesian phylodynamic inference with complex models.
569 *PLOS Computational Biology*, *14*(11), e1006546.
570 <https://doi.org/10.1371/journal.pcbi.1006546>
- 571 Zarebski, A. E., Plessis, L. du, Parag, K. V., & Pybus, O. G. (2020). A computationally tractable
572 birth–death model that combines phylogenetic and epidemiological data. *BioRxiv*,
573 2020.10.21.349068. <https://doi.org/10.1101/2020.10.21.349068>
- 574 Zhang, C., Stadler, T., Klopstein, S., Heath, T. A., & Ronquist, F. (2016). Total-Evidence
575 Dating under the Fossilized Birth–Death Process. *Systematic Biology*, *65*(2), 228–249.
576 <https://doi.org/10.1093/sysbio/syv080>

577

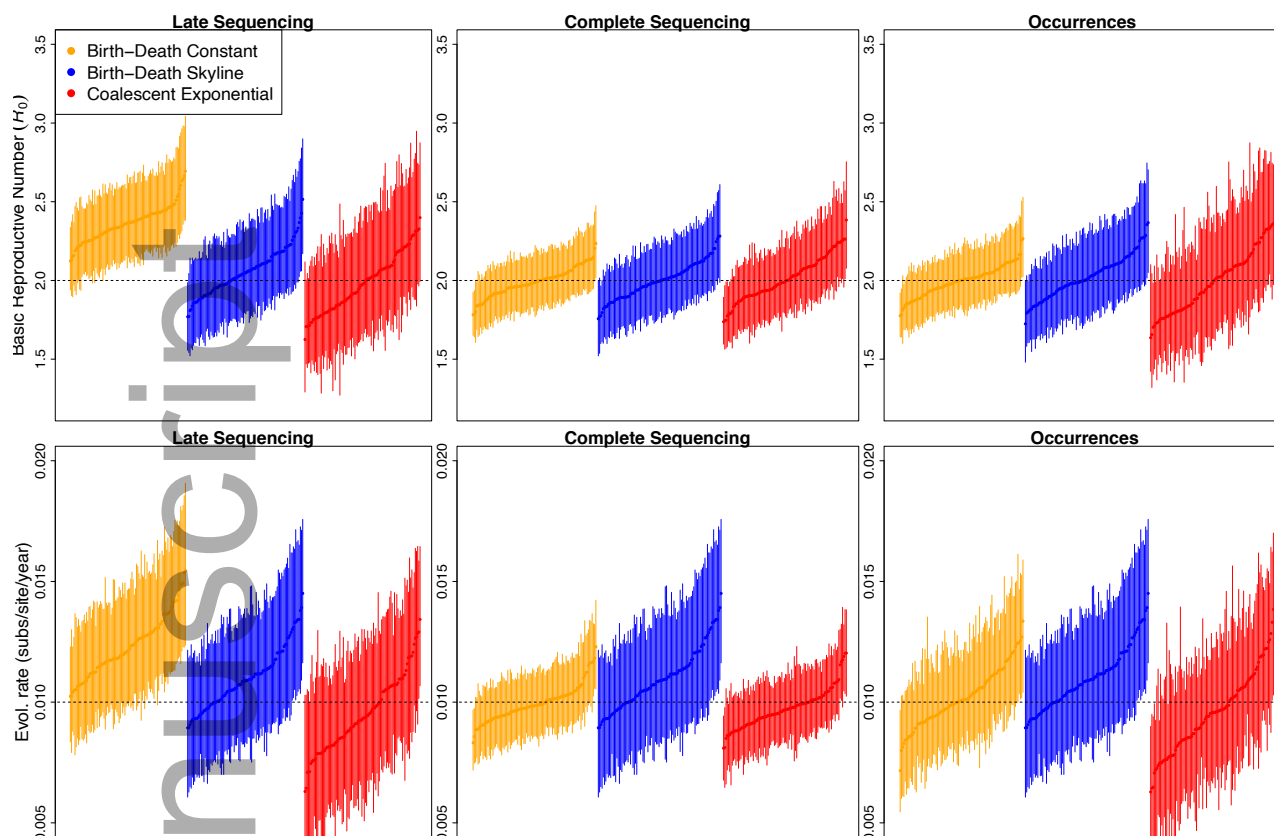
578

579 **Figure legends**



580

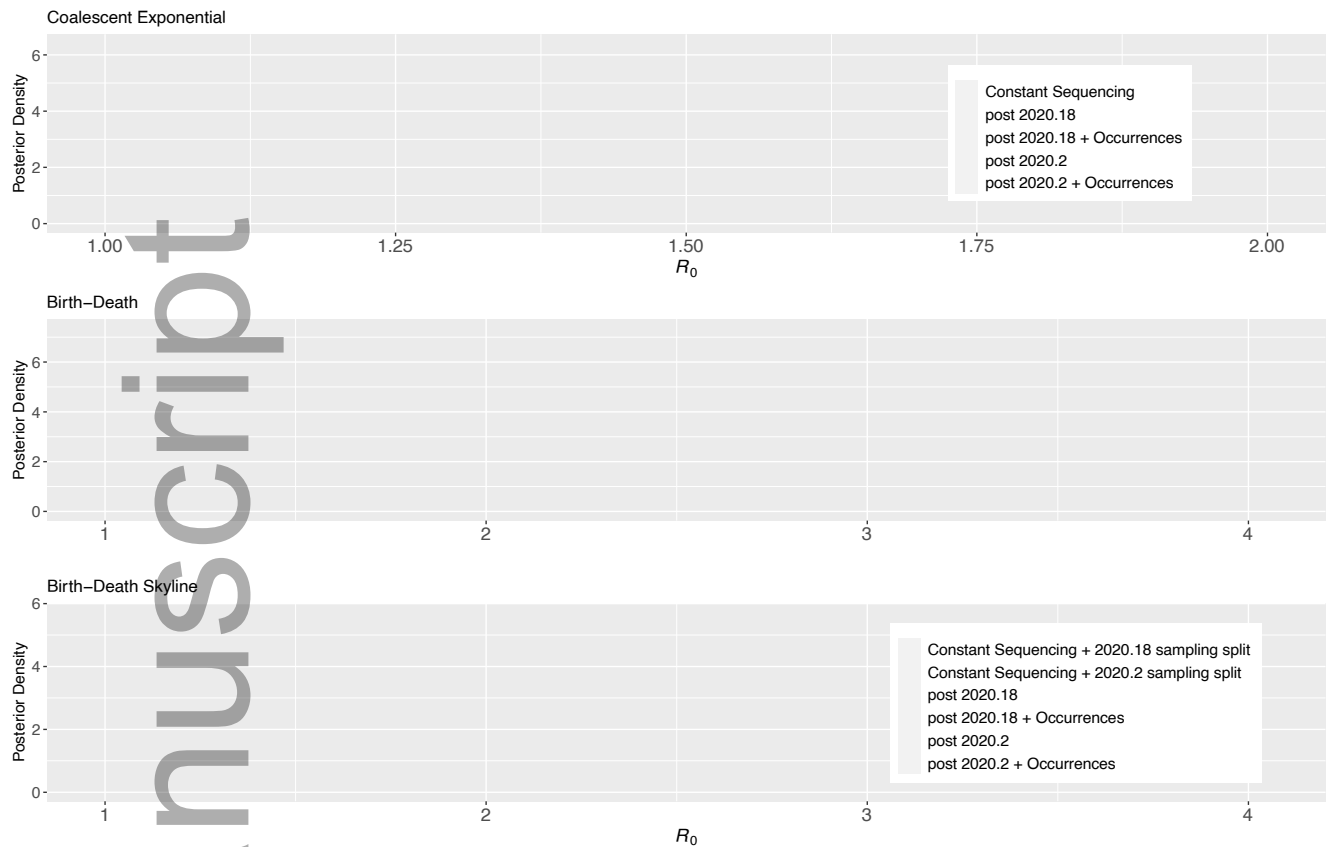
581 **Fig 1.** Example of a phylogenetic trees generated under a birth-death process with a basic
582 reproductive number (R_0) of 2, and a becoming uninfected rate (δ) of 100 for three analysis
583 scenarios. These trees are illustrative and simulated separately from those of our study data
584 sets. The solid line denotes the number of samples collected over time. In *constant* sampling
585 samples are collected and sequenced with probability $p =$ of 0.05 (i.e. $\psi=5$). In *late sampling*
586 samples are collected and sequenced after time T_s shown with the dashed line. In
587 *occurrence data* samples are collected constantly over time, but only sequenced after time
588 T_s , such that before T_s only occurrences (sampling times with no sequence data) are
589 included. Blue circles represent samples with sequence data, whereas those in orange
590 correspond with occurrences. In the *occurrence data* scenario, a Bayesian analysis would
591 integrate over their phylogenetic uncertainty. The solid line represents the number of
592 samples collected over time. In *late sampling* there are no samples collected before T_s , such
593 that assuming constant sampling can produce a bias in estimates of epidemiological
594 dynamics.



595

596 **Fig 2.** Posterior credible intervals of the basic reproductive number, R_0 , and the evolutionary
 597 rate for 100 simulations with true R_0 of 2 and an evolutionary rate of 0.01 subs/site/year. The
 598 bars represent the 95% highest posterior density (HPD) and the points are the mean.

599 Estimates are ordered from lowest to highest mean. We analysed the data by sequencing
 600 late in the outbreak only (i.e. after 0.75 of the tree height), with a constant sequencing effort
 601 (with all samples sequenced), and by including occurrences prior to sequences. The colours
 602 represent tree different tree priors; red for the coalescent exponential, blue for the birth-
 603 death skyline, and orange for the birth-death with constant sampling.



604

605 **Fig 3.** Posterior estimates of R_0 for SARS-CoV-2 genome data. Constant sampling refers to
 606 using all 821 genomes in the empirical dataset. Post 2020.18 refers to only including
 607 sequences from 2020-03-04 and afterwards. Post 2020.2 refers to the same from 2020-03-
 608 14 and afterwards. A) Posterior densities of the basic reproductive number, R_0 under the
 609 coalescent exponential. Late sequencing and corresponding late sequencing with
 610 occurrences posteriors are overlapping as expected. B) Posterior densities for estimates of
 611 the basic reproductive number, R_0 under the birth death. The same legend as in A applies.
 612 C) Birth-death skyline posteriors for R_0 .

613

614 **Tables**

615 **Table 1.** Priors used for analysis of simulated and empirical data. Priors used in empirical
 616 analyses were identical across treatments.

Simulated Data

Parameter	Birth-Death	Birth-Death Skyline	Coalescent Exponential
R_0	Lognormal($\mu = 0, \sigma = 1$)	Lognormal($\mu = 0, \sigma = 1$)	-
δ	$\Gamma(\text{mean} = 91, \text{sd} = 1)$	$\Gamma(\text{mean} = 91, \text{sd} = 1)$	-
p	0.05	0.05	-
origin	$U(0, \infty)$	$U(0, \infty)$	-
r	-	-	Laplace($\mu = 0.001, \text{scale} = 30.70$)
ϕ	-	-	$1/x$

Washington SARS-CoV-2 data

Parameter	Birth-Death	Birth-Death Skyline	Coalescent Exponential
R_0	Lognormal($\mu = 0, \sigma = 1$)	Lognormal($\mu = 0, \sigma = 1$)	-
δ	36.5	36.5	-
p	$\beta(2, 2)$	$\beta(2, 2)$	-
origin	$U(0, \infty)$	$U(0, \infty)$	-
r	-	-	Laplace($\mu = 0, \text{scale} = 100$)
ϕ	-	-	Exp(mean=100)

617

618

619 **Table 2.** Results of the simulation study with R_0 of 2 and evolutionary rate of 0.01
 620 subs/site/year. The rows correspond to the seven treatments. For R_0 and evolutionary rate
 621 (subs/site/year), columns denote the number of simulations (out of 100) where the value
 622 used to generate the data was contained within the 95% highest posterior density (HPD),
 623 also referred to as coverage and reflecting accuracy; average absolute error measures the
 624 average absolute difference between posterior mean R_0 and 2; and the average HPD width.
 625 BD stands for birth-death, CE for coalescent exponential, and BDSky to the birth-death
 626 skyline model with two sampling intervals.

Treatment	R_0 Within HPD	R_0 Mean Absolute Error	R_0 Mean HPD Width	Rate Within HPD	Rate Mean Absolute Error	Rate Mean HPD Width
CE + Constant Sequencing	89	0.115	0.481	92	0.00067	0.00298
CE + Late Sequencing	100	0.156	0.786	97	0.00138	0.00646
CE + Occurrences	96	0.163	0.748	91	0.00140	0.00605
BD + Constant Sequencing	94	0.072	0.364	94	0.00057	0.00285
BD + Late Sequencing	12	0.364	0.515	58	0.00233	0.00511
BD + Occurrences	94	0.076	0.384	95	0.00097	0.00450
BDSky + Constant Sequencing	95	0.095	0.468	92	0.00059	0.00290
BDSky + Late Sequencing	95	0.125	0.591	90	0.00122	0.00559
BDSky + Occurrences	97	0.109	0.529	93	0.00106	0.00495

627

628

629 **Table 3.** Posterior estimates of R_0 and p using the birth-death for the SARS-CoV-2 empirical
 630 dataset. Rows correspond to the 12 treatments.

Treatment	Mean R_0	95% HPD
BD Constant Sampling	1.96	(1.85, 2.07)
BD Post 2020.18	2.44	(2.31, 2.58)
BD Post 2020.18 + Occurrences	1.97	(1.87, 2.08)
BD Post 2020.2	3.53	(3.24, 3.82)
BD Post 2020.2 + Occurrences	1.96	(1.83, 2.09)
BDSky Constant Sequencing + 2020.18 sampling split	1.89	(1.75, 2.03)
BDSky Post 2020.18	1.57	(1.43, 1.71)
BDSky Post 2020.18 + Occurrences	1.93	(1.79, 2.07)
BDSky Constant Sequencing + 2020.2 sampling split	2.10	(1.95, 2.24)
BDSky Post 2020.2	1.48	(1.35, 1.63)
BDSky Post 2020.2 + Occurrences	2.16	(1.99, 2.34)
CE Constant Sampling	1.52	(1, 1.65)
CE Post 2020.18	1.51	(1.39, 1.64)
CE Post 2020.18 + Occurrences	1.50	(1.38, 1.62)
CE Post 2020.2	1.43	(1.3, 1.58)
CE Post 2020.2 + Occurrences	1.43	(1.29, 1.57)

631

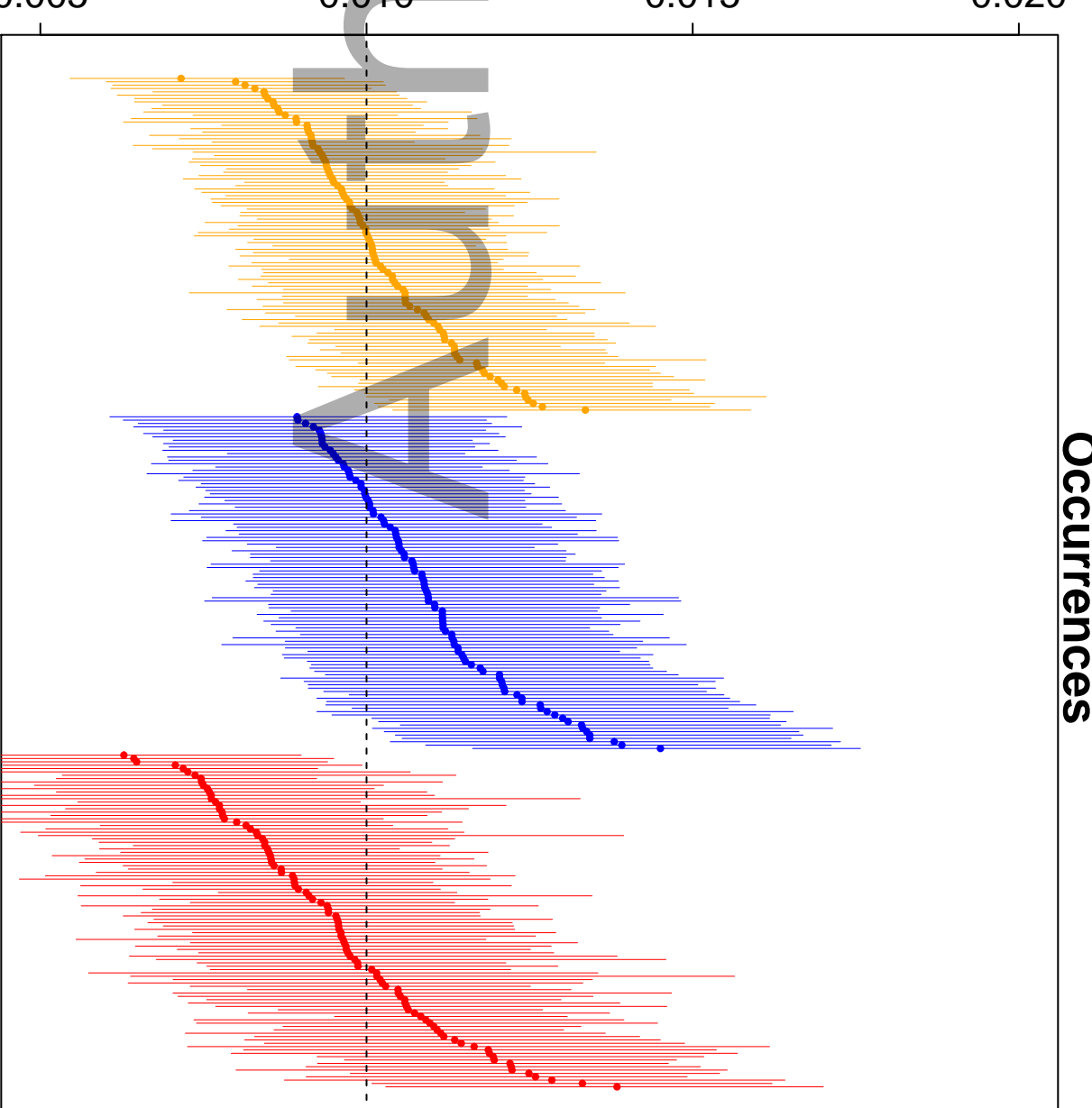
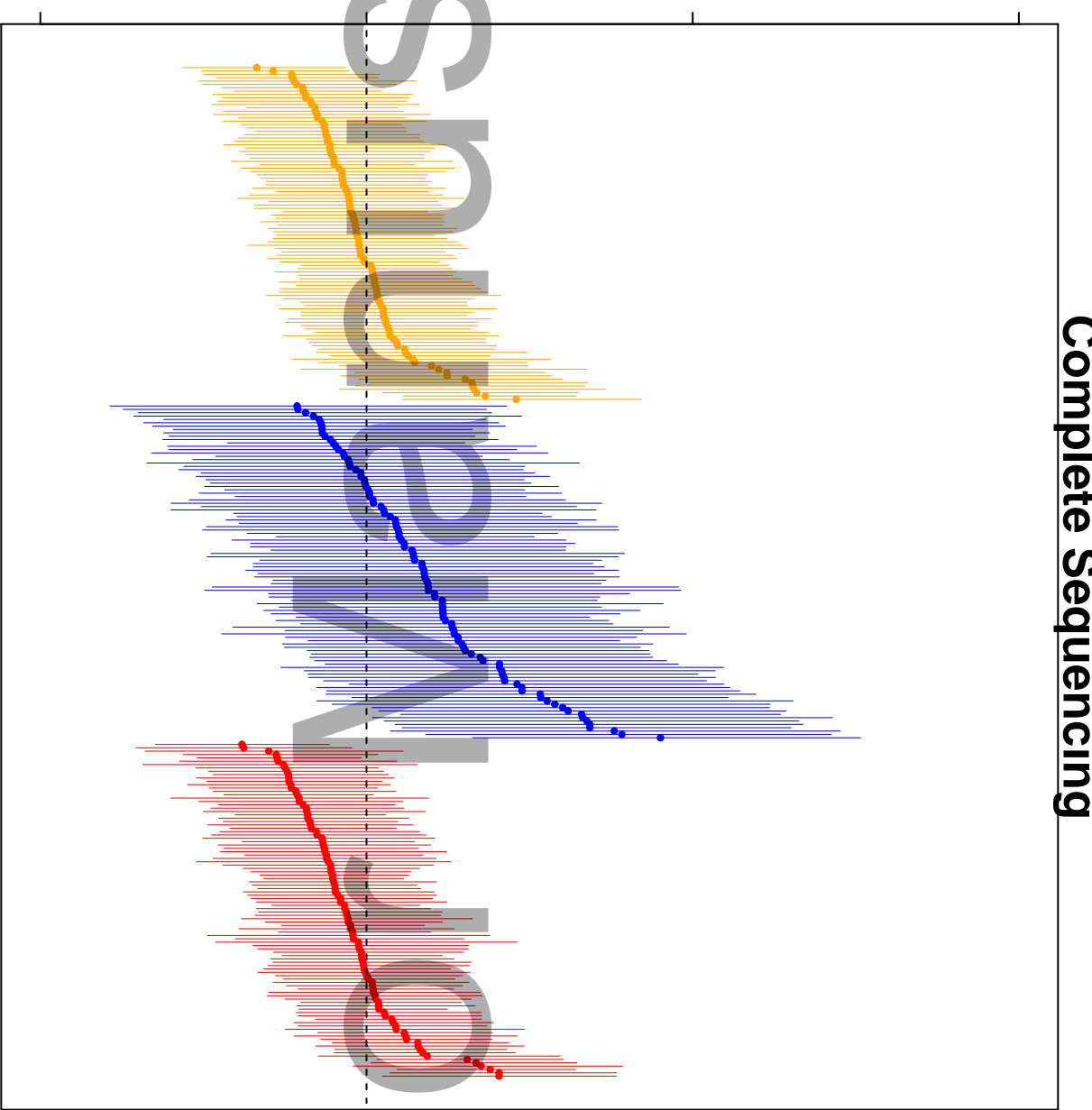
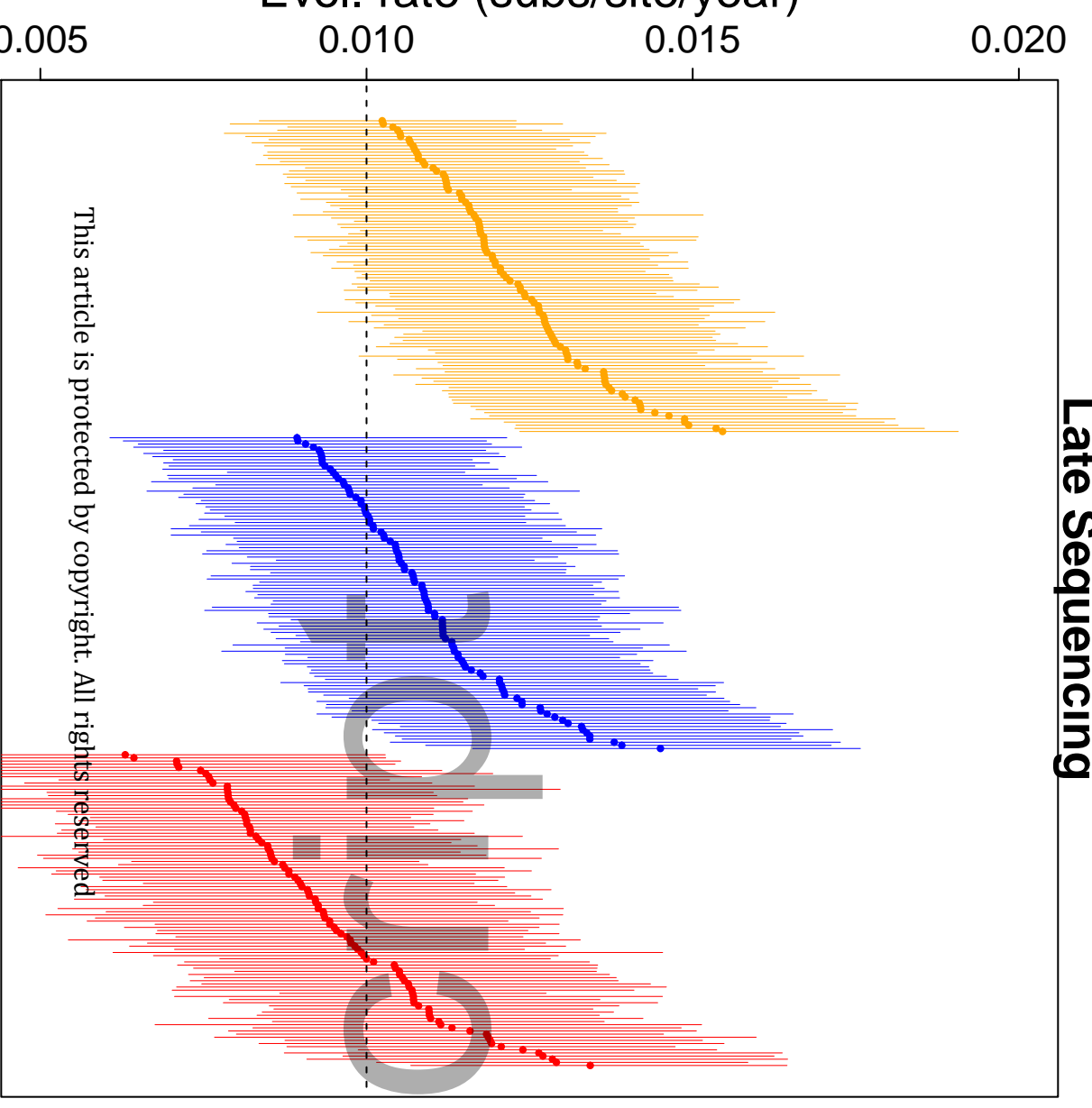
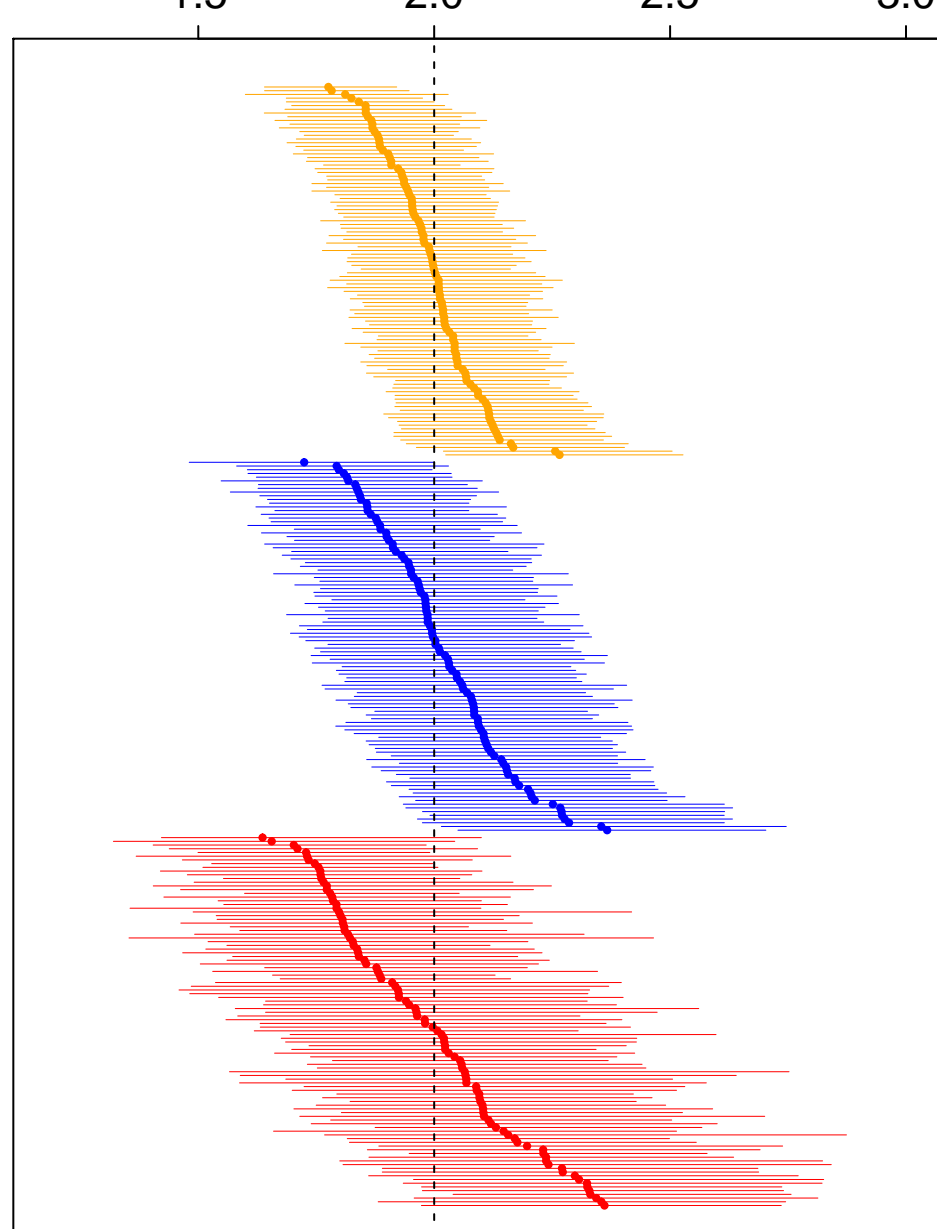
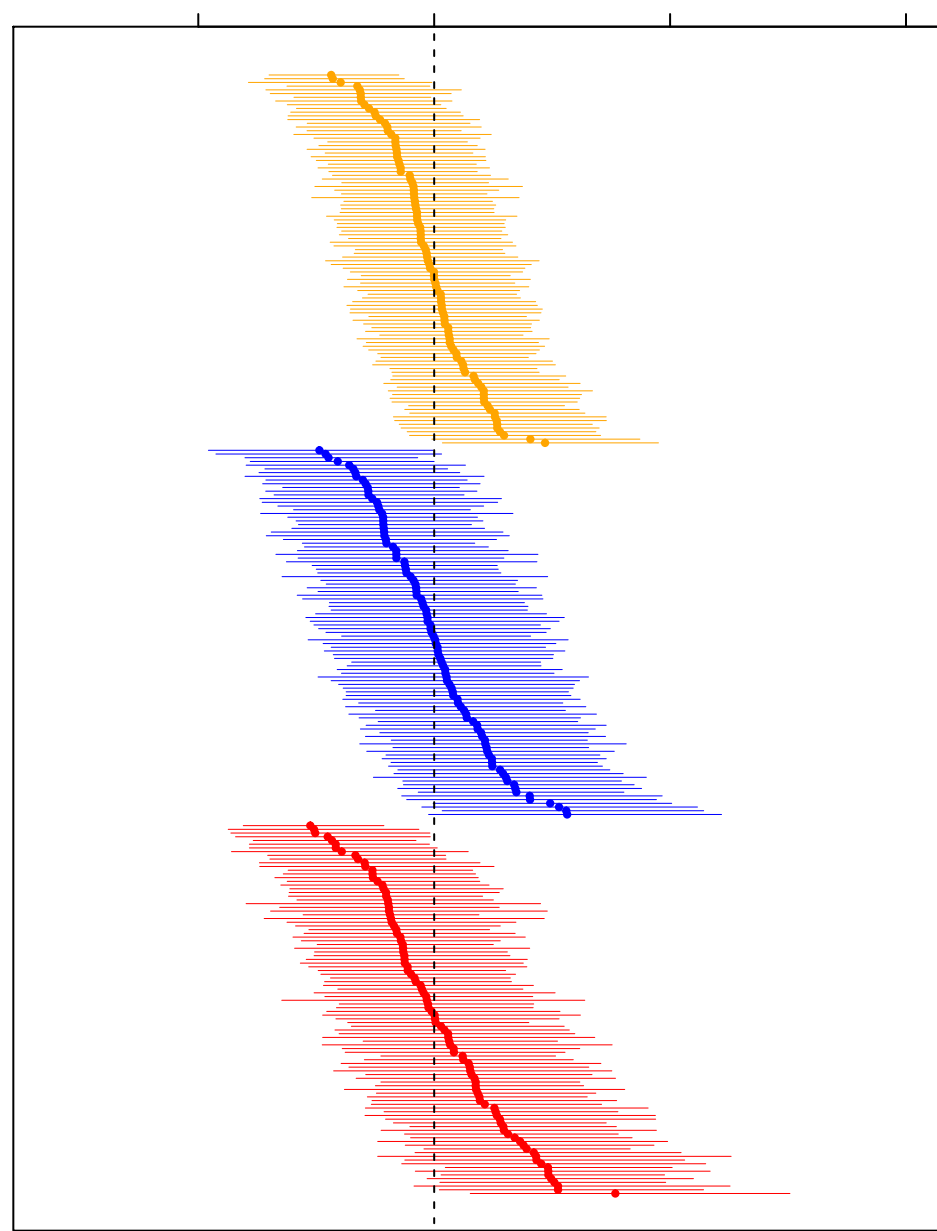
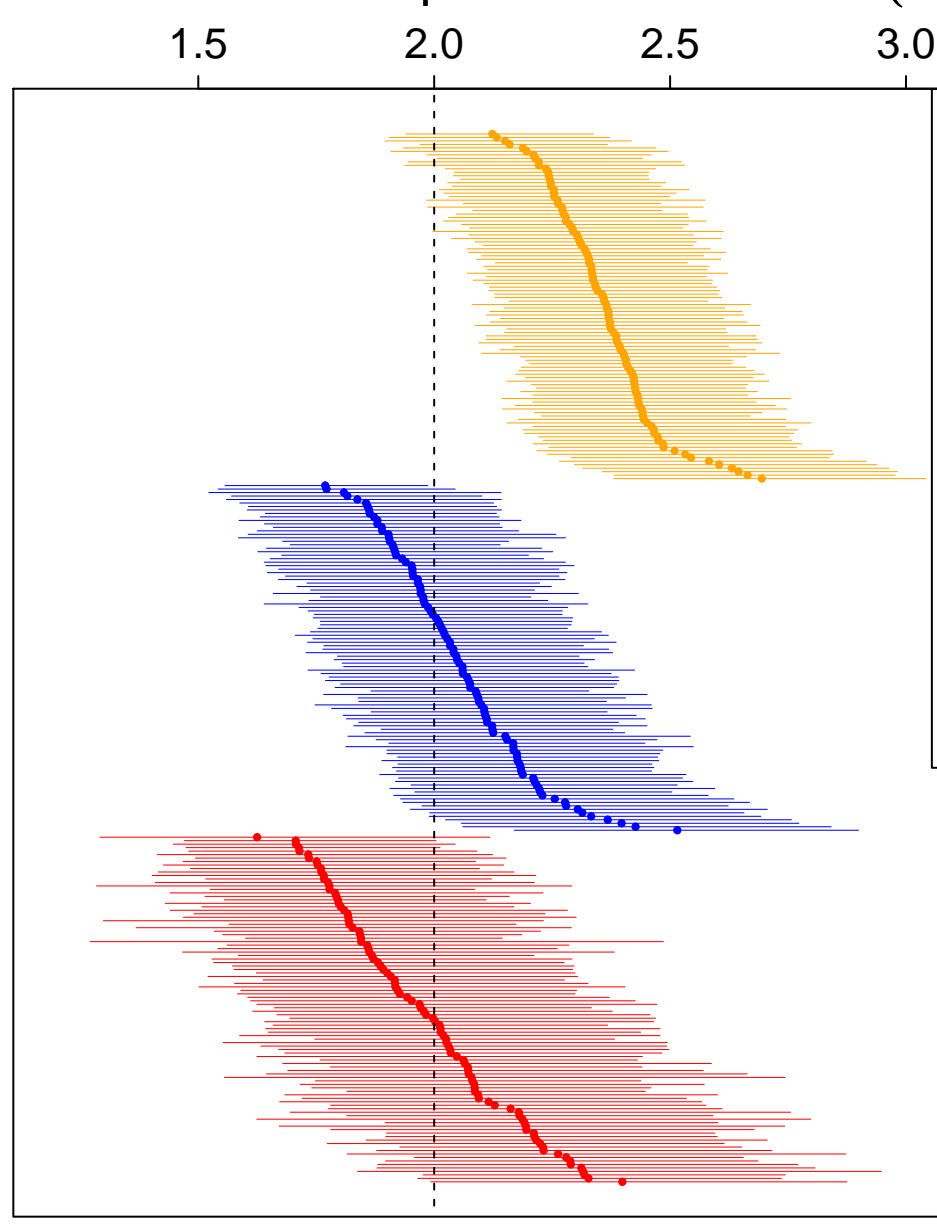
Author Manuscript

Late Sequencing

Complete Sequencing

Occurrences

● Birth-Death Constant
● Birth-Death Skyline
● Coalescent Exponential



Late Sequencing

Complete Sequencing

Occurrences

This article is protected by copyright. All rights reserved

Evol. rate (subs/site/year)

Basic Reproductive Number (R_0)

0.005

0.010

0.015

0.020

1.5

2.0

2.5

3.0

3.5

0.005

0.010

0.015

0.020

1.5

2.0

2.5

3.0

3.5

0.005

0.010

0.015

0.020

1.5

2.0

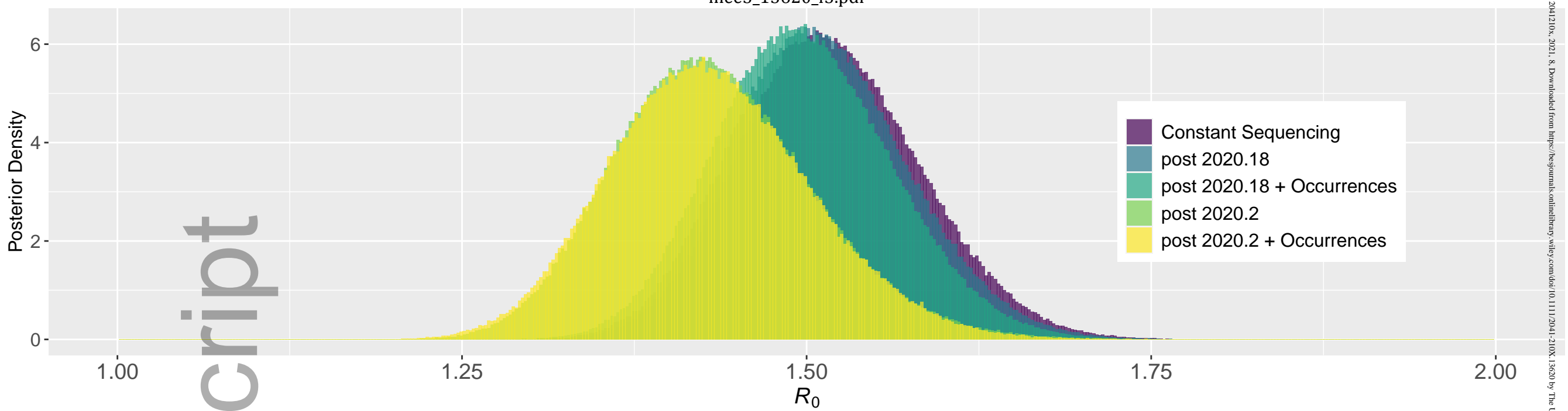
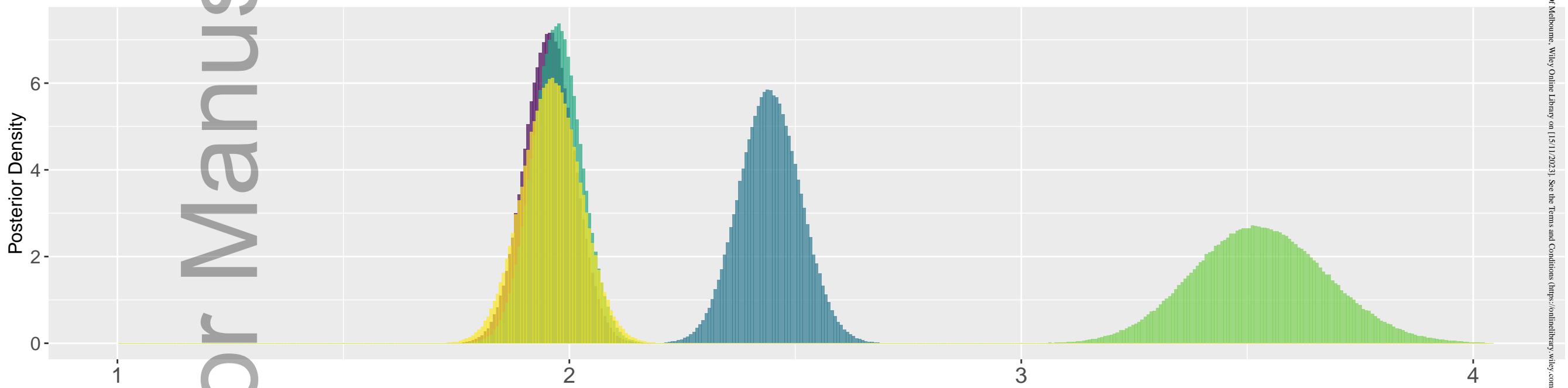
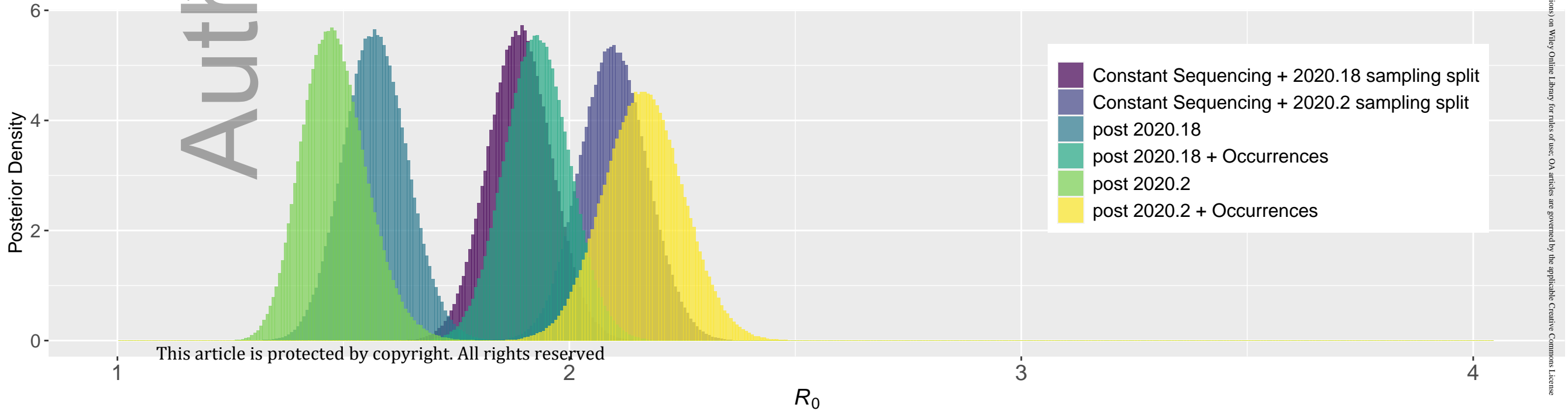
2.5

3.0

3.5

A Coalescent Exponential

mee3_13620_f3.pdf

**B** Birth–Death**C** Birth–Death Skyline

Author Manuscript

This article is protected by copyright. All rights reserved

2041210x, 2021, 8, Downloaded from https://onlinelibrary.wiley.com/doi/10.1111/2041-210X.13620 by The University Of Melbourne, Wiley Online Library on [15/11/2023]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License