# Phylodynamic inference for emerging viruses using segregating sites

Authors: Yeongseon Park[1], Michael Martin[1], Katia Koelle[2,3,*]


[1] Graduate Program in Population Biology, Ecology, and Evolution, Emory University, Atlanta, GA 30322

[2] Department of Biology, Emory University, Atlanta, GA 30322

[3] Emory-UGA Center of Excellence for Influenza Research and Surveillance (CEIRS), Atlanta GA, USA


*katia.koelle@emory.edu

1    **Abstract**

2    Epidemiological models are commonly fit to case data to estimate model parameters and to infer

3    unobserved disease dynamics. More recently, epidemiological models have also been fit to viral

4    sequence data using phylodynamic inference approaches that generally rely on the

5    reconstruction of viral phylogenies. However, especially early on in an expanding viral population,

6    phylogenetic uncertainty can be substantial and methods that require integration over this

7    uncertainty can be computationally intensive. Here, we present an alternative approach to

8    phylodynamic inference that circumvents the need for phylogenetic tree reconstruction. Our

9    "tree-free" approach instead relies on quantifying the number of segregating sites observed in

10    sets of sequences over time and using this trajectory of segregating sites to infer epidemiological

11    parameters within a Sequential Monte Carlo (SMC) framework. Using forward simulations, we

12    first show that epidemiological parameters and processes leave characteristic signatures in

13    segregating site trajectories, demonstrating that these trajectories have the potential to be used

14    for phylodynamic inference. We then show using mock data that our proposed approach

15    accurately recovers key epidemiological quantities such as the basic reproduction number and

16    the timing of the index case. Finally, we apply our approach to SARS-CoV-2 sequence data from

17    France, estimating a reproductive number of approximately 2.2 and an introduction time of mid-

18    January 2021, consistent with estimates from epidemiological surveillance data. Our findings

19    indicate that "tree-free" phylodynamic inference approaches that rely on simple population

20    genetic summary statistics can play an important role in estimating epidemiological parameters

21    and reconstructing infectious disease dynamics, especially early on in an epidemic.

22

23

24

25

26

27

## Introduction

Phylodynamic inference methods use viral sequence data to estimate epidemiological quantities such as the basic reproduction number and to reconstruct epidemiological patterns of incidence and prevalence. These inference methods have been applied to sequence data across a broad range of RNA viruses, including HIV (Stadler and Bonhoeffer 2013; Popinga et al. 2014; Ratmann et al. 2017; Volz et al. 2017), ebola (Stadler et al. 2014; Vaughan et al. 2017; Volz and Siveroni 2018), dengue (Rasmussen et al. 2014), influenza (Rasmussen and Stadler), and most recently severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)(Danesh et al. 2020; Miller et al. 2020; Geidelberg et al. 2021). Most commonly, phylodynamic inference methods rely on underlying coalescent models or birth-death models. Coalescent-based approaches have been generalized to accommodate time-varying population sizes and parameter estimation for structured epidemiological models, for example, susceptible-exposed-infected-recovered (SEIR) models and models with spatial compartmentalization (Volz 2012; Volz and Siveroni 2018). Birth-death approaches (Stadler 2010; Stadler et al. 2012), where a birth in the context of infectious diseases corresponds to a new infection and death corresponds to a recovery from infection, instead carry other advantages, such as incorporating the role of demographic stochasticity in disease dynamics, which may be particularly important in emerging diseases that start with low infection numbers (Boskova et al. 2014). Both of these classes of phylodynamic inference approaches rely on time-resolved phylogenies and have been incorporated into the phylogenetics software package BEAST2 (Bouckaert et al. 2014: 2) to allow joint estimation of epidemiological parameters and dynamics while integrating over phylogenetic uncertainty (Stadler et al. 2013; Volz and Siveroni 2018). Integrating over phylogenetic uncertainty is crucial when applying these methods to viral sequence data that are sampled over a short period of time and contain only low levels of genetic diversity. However, integrating over phylogenetic uncertainty is computationally intensive. Moreover, phylodynamic approaches that use reconstructed trees for inference require estimation of parameters associated with models of sequence evolution, along with parameters that are of more immediate epidemiological interest.

55     Here, we present an alternative phylodynamic inference method that is particularly appropriate
56     to use when viral sequences are sampled over short time periods and when phylogenetic
57     uncertainty is considerable. This method does not rely on time-resolved phylogenies to infer
58     epidemiological parameters or to reconstruct patterns of viral spread. Instead, the "tree-free"
59     method we propose here fits epidemiological models to time series of the number of segregating
60     sites observed in a viral population that is sampled over time. Like existing coalescent-based
61     approaches, the approach we propose here allows for structured infectious disease models to be
62     considered in a straightforward "plug-and-play" manner. Like existing birth-death process
63     approaches, it incorporates the effect that demographic noise may have on epidemiological
64     dynamics. Below, we first describe how segregating site trajectories are calculated using
65     sequence data and how they are impacted by sampling effort, rates of viral spread, and
66     transmission heterogeneity. We then describe our proposed phylodynamic inference method
67     and apply it to simulated data to demonstrate the ability of this method to infer epidemiological
68     parameters and to reconstruct unobserved epidemiological dynamics. Finally, we apply our
69     segregating sites method to SARS-CoV-2 sequence data from France, arriving at quantitatively
70     similar parameter estimates to those arrived at using epidemiological data.

71     **New Approaches**

72     Mutations occur during viral replication within infected individuals and these have the potential
73     to be transmitted. During the epidemiological spread of an emerging virus, the virus population
74     (distributed across infected individuals) thus accrues mutations and diversifies genetically. This
75     joint process of viral spread and evolution can be simulated forward in time using compartmental
76     models, with patterns of epidemiological spread leaving signatures in the evolutionary trajectory
77     of the virus population. Parameters of these compartmental models that govern patterns of
78     epidemiological spread can thus be estimated using observed viral evolutionary trajectories.
79     Here, we develop a phylodynamic inference approach that fits compartmental epidemiological
80     models to times series of a low-dimensional evolutionary summary statistic. Specifically, we use
81     trajectories of the number of segregating sites from samples of the viral population taken over
82     time for phylodynamic inference. In Materials and Methods, we provide details on the simulation

83  of epidemiological models that incorporate viral evolution and thus can yield simulated time

84  series of the number of segregating sites. We further describe our phylodynamic inference

85  approach that relies on using particle filtering (otherwise known as Sequential Monte Carlo; SMC)

86  to infer parameters for these epidemiological models of arbitrary complexity and to reconstruct

87  unobserved disease dynamics.

## Results

89  **Segregating site trajectories are informative of epidemiological dynamics.**

90  Simulations of epidemiological models, as detailed in Materials and Methods, indicate that the

91  number of segregating sites that are observed over time in a viral population are sensitive to

92  sampling effort and are informative of epidemiological dynamics. To demonstrate this, we first

93  simulated a susceptible-exposed-infected-recovered (SEIR) model under an epidemic scenario

94  starting with a single infected individual (Figure 1A), further tracking the viral genotypes

95  according to the approach outlined in Materials and Methods. The effect of sampling effort is

96  shown in Figure 1B, which plots segregating site trajectories under dense sampling effort (40

97  sequences per 4-day time window) and under sparse sampling effort (20 sequences per 4-day

98  time window). At both of these sampling efforts, the number of segregating sites first increases

99  as the epidemic grows, as expected, with mutations accumulating in the virus population.

100  Following the peak of the epidemic, the number of segregating sites starts to decline as viral

101  lineages die out, reducing the amount of genetic variation present in the viral population.  At

102  lower sampling effort, less of the genetic variation present in the viral population over a given

103  time window is likely to be sampled, resulting in a lower number of observed segregating sites

104  during any time window.

105  To assess whether segregating site trajectories could be used for phylodynamic inference, we

106  first considered whether these trajectories differed between epidemics governed by different

107  basic reproduction numbers ($R_0$ values). Figure 1C shows simulations of the SEIR model under

108  two parameterizations of the basic reproduction number: an $R_0$ of 1.6, corresponding to the

109  simulation shown in Figure 1A, and a higher $R_0$ of 2.0. Differences in $R_0$ were implemented by

110  differences in the transmission rate. The epidemic with the higher $R_0$ grew more rapidly (Figure

111    1C) and, under the same sampling effort, resulted in a more rapid increase in the number of

112    segregating sites (Figure 1D). This indicates that segregating site trajectories can be informative

113    of $R_0$ early on in an epidemic.

114    We next considered the effect of transmission heterogeneity on segregating site trajectories.

115    Many viral pathogens are characterized by 'superspreading' dynamics, where a relatively small

116    proportion of infected individuals are responsible for a large proportion of secondary infections

117    (Lloyd-Smith et al. 2005). The extent of transmission heterogeneity is often gauged relative to

118    the 20/80 rule (the most infectious 20% of infected individuals are responsible for 80% of the

119    secondary cases (Woolhouse et al. 1997)), with some pathogens like SARS-CoV-2 exhibiting

120    extreme levels of superspreading, with as low as 6-15% of infected individuals responsible for 80%

121    of secondary cases (Althouse et al. 2020; Miller et al. 2020; Lemieux et al. 2021; Sun et al. 2021).

122    Because transmission heterogeneity is known to impact patterns of viral genetic diversity (Koelle

123    and Rasmussen 2012), we simulated the above SEIR model with transmission heterogeneity to

124    ascertain its effects on segregating site trajectories. Transmission heterogeneity was

125    implemented using a negative binomial distribution parameterized such that the most infectious

126    6% of infected individuals are responsible for 80% of the secondary cases (Materials and

127    Methods). Because transmission heterogeneity has a negligible impact on epidemiological

128    dynamics once the number of infected individuals is large (Keeling and Rohani 2008), these

129    simulated epidemiological dynamics should be quantitatively similar to one another, with

130    transmission heterogeneity simply expected to shorten the timing of epidemic onset in

131    simulations with successful invasion (Lloyd-Smith et al. 2005). Our simulations confirm this

132    pattern (Figure 1E). To compare segregating site trajectories between these simulations, we

133    therefore shifted the simulation with transmission heterogeneity later in time such that the two

134    simulated epidemics peaked at similar times (Figure 1E). Comparisons of segregating site

135    trajectories between these simulations indicated that transmission heterogeneity substantially

136    decreases the number of segregating sites during any time window (Figure 1F). These results

137    indicate that the number of segregating sites in principle could be informative of the extent of

138    transmission heterogeneity present in an unfolding epidemic. They also indicate that

139　transmission heterogeneity needs to be taken into consideration when estimating

140　epidemiological parameters using segregating site trajectories.

141　Finally, we wanted to assess whether changes in $R_0$ over the course of an epidemic would leave

142　signatures in segregating site trajectories. We considered this scenario because phylodynamic

143　inference has often been used to quantify the effect of public health interventions on $R_0$, most

144　recently in the context of SARS-CoV-2 (Danesh et al. 2020; Miller et al. 2020). We thus

145　implemented simulations with $R_0$ starting at 1.6 and then either remaining at 1.6 or reduced to

146　either 1.1 or 0.75 when the number of infected individuals reached 400 (Figure 1G). The

147　segregating site trajectories for these three simulations indicate that reductions in $R_0$ over the

148　course of an epidemic leave faint signatures in this low-dimensional summary statistic of viral

149　diversity, with the signature being more pronounced with a more precipitous drop in $R_0$ (Figure

150　1H).

**Phylodynamic inference using segregating site trajectories**

152　To examine the extent to which phylodynamic inference based on segregating sites can be used

153　for parameter estimation, we generated a mock segregating site trajectory by forward simulating

154　an SEIR model with a $R_0$ of 1.6, sampling viral sequences from this simulation (Figure 2A), and

155　calculating a segregating site trajectory from these sampled sequences (Figure 2B). Because the

156　duration of the exposed period and the duration of the infectious period are generally known for

157　viruses undergoing phylodynamic analysis, we fixed these parameters at their true values and

158　first attempted to estimate only $R_0$ under the assumption that the timing of the index case $t_0$ is

159　known. We estimated an $R_0$ value of 1.59 (95% confidence interval of 1.49 to 1.64; Materials and

160　Methods; Figure 2C, 2C inset), demonstrating that phylodynamic inference using our segregating

161　sites approach applied to this simulated dataset is able to recover the true $R_0$ value of 1.6.

162　Because the timing of the index case is almost certainly not known for an emerging epidemic, we

163　further attempted to estimate both $R_0$ and $t_0$ using the segregating site trajectory shown in Figure

164　2B. To do this, we first considered the parameter space ranging from an $R_0$ of 1.2 to 2.5 and from

165　a $t_0$ of 60 days prior to the true start date of 0 to 56 days following this true start date. Considering

166　$R_0$ intervals of 0.02 and $t_0$ intervals of 2 days, we ran 10 SMC simulations for every parameter

167    combination. In Figure 3A, we plot the mean value of these 10 SMC log-likelihoods for every

168    parameter combination in the considered parameter space. Examination of this plot indicates

169    that there is a log-likelihood ridge that runs between early $t_0$/low $R_0$ parameter combinations and

170    late $t_0$/high $R_0$ parameter combinations. However, this ridge falls off on both edges, indicating

171    that the segregating sites approach can in principle estimate both $t_0$ and $R_0$. We therefore

172    calculated profile likelihoods for both $R_0$ and $t_0$ (Figures 3B, 3C; Materials and Methods), arriving

173    at an $R_0$ estimate of 1.50 (95% confidence = 1.34 to 1.67; Figure 3B) and a $t_0$ value of -13.8 (95%

174    confidence = -27.8 to 0.3; Figure 3C) for the simulated dataset. While the maximum likelihood

175    estimate for $R_0$ ran low and for $t_0$ ran early, the confidence intervals contained the true values of

176    $R_0 = 1.6$ and $t_0 = 0$, respectively. Our results indicate that joint estimation of these parameters is

177    thus possible. Using our estimates of $R_0$ and $t_0$, we reconstructed the dynamics of the segregating

178    sites (Figure 4A) and unobserved state variables: the number of susceptible, exposed, and

179    infected individuals over time (Figures 4B, C, D). These reconstructed state variables captured

180    the true epidemiological dynamics, demonstrating that our segregating sites phylodynamic

181    inference approach can be used to estimate epidemiological variables that generally go

182    unobserved.

**Phylodynamic inference for SARS-CoV-2 sequences from France**

184    We applied the segregating sites inference approach to a set of SARS-CoV-2 sequences sampled

185    from France between January 23, 2020 and March 17, 2020, when a country-wide lockdown was

186    implemented. We decided to apply our approach to this set of sequences for several reasons.

187    First, a large fraction of the 479 available full-genome sequences from France over this time

188    period appear to be genetically very similar to one another (Gámbaro et al. 2020), indicating that

189    one major lineage may have taken off in France (or at least, that most samples stemmed from

190    one major lineage). This lineage would be the focus of our analysis. Second, an in-depth analysis

191    previously inferred $R_0$ for France prior to the March 17 lock-down measures that were

192    implemented (Salje et al. 2020). This analysis fit a compartmental infectious disease model to

193    epidemiological data that included case, hospitalization, and death data. Because our

194    phylodynamic inference approach can accommodate epidemiological model structures of

195    arbitrary complexity, we can adopt the same model structure as in this previous analysis. We can

196    also set the epidemiological parameters that are assumed fixed in this previous analysis to their

197    same values. By controlling for model structure and the set of model parameters assumed as

198    given, we can ask to what extent sequence data corroborate the $R_0$ estimates arrived at from

199    detailed fits to epidemiological data.

200    To apply our segregating sites approach to the viral sequences from France, we first identified

201    the subset of the 479 sequences that constituted a single, large lineage. To keep with the "tree-

202    free" emphasis of our approach, we identified this subset of $n$ = 432 sequences without inferring

203    a phylogeny (Materials and Methods). Using phylogenetic inference, however, we confirmed that

204    our subset of sequences constituted a single evolutionary lineage (Figure S1). We calculated the

205    nucleotide distance from each sequence in this subset to Wuhan/Hu-1 (Wu et al. 2020)

206    (EPI_ISL_402125), a commonly used reference SARS-CoV-2 sequence that stemmed from a

207    sample collected in Wuhan, China in late December 2019. Using these nucleotide distances, we

208    estimated an evolutionary rate of 8.21 x $10^{-4}$ substitutions/site/yr (Figure 5A), consistent with the

209    range of inferred evolutionary rate estimates for SARS-CoV-2 (Duchene et al. 2020; Pekar et al.

210    2020). This provides another confirmation that this subset of sequences is a single evolutionary

211    lineage brought into France early on during the pandemic.

212    To generate a segregating site trajectory from these sequences, we established consecutive, non-

213    overlapping 4-day time windows such that the last time window ended on March 17, 2020. Figure

214    5B shows the number of sequences falling into each time window. Figure 5C shows the

215    segregating site trajectory calculated from these sequences. We jointly estimated $R_0$ and $t_0$ using

216    this segregating site trajectory, under the assumption that the most infectious 15% of SARS-CoV-

217    2 infected individuals are responsible for 80% of secondary infections, based on literature

218    estimates of the extent of SARS-CoV-2 transmission heterogeneity (Sun et al. 2021) (Materials

219    and Methods). We parameterized the model with a per genome, per transmission mutation rate

220    of $\mu$ = 0.33 using consensus sequence data from established SARS-CoV-2 transmission pairs that

221    were available in the literature (James et al. 2020; Popa et al. 2020; Braun et al. 2021; Lythgoe et

222    al. 2021) (Materials and Methods). Specifically, for each of the 87 transmission pairs we had

223    access to, we calculated the nucleotide distance between the consensus sequence of the donor

224    sample and that of the recipient sample and fit a Poisson distribution to these data (Figure 5D).

225    Using this approach, we estimated a $\mu$ value of 0.33 (95% confidence interval of 0.22 to 0.48),

226    corresponding approximately to one mutation occurring every 3 transmission events.

227    Similar to the approach we undertook with our simulated data to jointly estimate $R_0$ and $t_0$, we

228    first considered a broad parameter space over which to calculate log-likelihood values.

229    Specifically, we considered $R_0$ values between 1.2 and 3.4 (at intervals of 0.1) and $t_0$ values

230    between December 2, 2019 and February 16, 2020 (at intervals of 2 days). We ran 10 SMC

231    simulations and calculated the mean log-likelihood for each parameter combination (Figure 6A).

232    Similar to our findings on the simulated data set, we found evidence for a log-likelihood ridge

233    between early $t_0$/low $R_0$ and late $t_0$/high $R_0$ parameter combinations. Profile log-likelihoods for

234    $R_0$ and $t_0$ are shown in Figures 5B and 5C, respectively, yielding an estimate of $R_0$ =  2.22 (95%

235    confidence interval = 1.5 to 2.94) and an estimate of $t_0$ =  January 11 (95% confidence interval =

236    December 26, 2019 to January 28, 2020). Our maximum likelihood estimate of $R_0$ is somewhat

237    lower than the $R_0$ estimate arrived at through the epidemiological time series analysis that

238    presented the epidemiological model structure we adopted (Salje et al. 2020). That analysis

239    inferred an $R_0$ of 2.9 (95% confidence interval = 2.81 to 3.01) in France over this same time period.

240    However, the confidence intervals of our analyses are relatively broad for $R_0$, and their estimate

241    of $R_0$ = 2.9 falls within our 95% confidence interval. Our estimate is closer in line with estimates

242    of the reproduction number in Wuhan prior to travel restrictions being introduced ($R_0$ = 2.35,

243    with 95% CI of 1.15-4.77) (Kucharski et al. 2020) and with those estimated for Western European

244    countries using incidence data up through March 17, 2020 ($R_0$ = 2.2, with 95% CI of 1.9-2.6)

245    (Locatelli et al. 2021). Our estimate also aligns more closely with projections of $R_0$ made

246    specifically for France, using outbreak data from Wuhan (Hilton and Keeling 2020): $R_0$ = 2.2 and

247    $R_0$ = 2.7, under different assumptions related to age-dependent susceptibility and infectiousness.

248    Finally, our $R_0$ estimates can be juxtaposed against results from phylodynamic analyses that used

249    a birth-death model to infer $R_0$ during three distinct epochs in France using a similar set of

250    sequence data we analyze here (Danesh et al. 2020). Their second epoch spanned February 19

251    through March 7, and the $R_0$ inferred for this time period was 2.56 (95% credible interval = 1.66

252    to 4.74). Our maximum likelihood estimate of $t_0$ in the middle of January 2020 aligns well with

253    findings from Gámbaro et al. ( 2020) and is further consistent with the estimate from Salje et al.

254     (2020) that 58.65 (95% CI 37.85 – 88.37) individuals were present in the exposed ($E_1$ class) on

255     January 22, 2020 based on fitting the epidemiological model to epidemiological data.

256     As we had done in our analysis of the simulated data set, we reconstructed the unobserved state

257     variables using sampled particles from SMC simulations parameterized with $R_0$ and $t_0$ values that

258     were sampled from the parameter space shown in Figure 6, weighted according to the log-

259     likelihood values of the parameter combination.   Plotting of reconstructed segregating site

260     trajectories indicated a very good fit to the observed segregating site trajectory (Figure 7A). The

261     number of individuals in the $E_1$, $E_2$, and $I$ classes increased exponentially over the time period

262     considered (Figure 7B), as expected for an epidemic with an $R_0 > 1$. In Figure 7C, we plot the

263     reconstructed cumulative number of exposed individuals over time and the reconstructed

264     cumulative number of recovered individuals over time. These cumulative dynamics indicate that

265     by mid-March 0.004% to 0.069% of the population in France had become infected by this SARS-

266     CoV-2 lineage and that 0.001% to 0.017% of the population in France had recovered from

267     infection from this SARS-CoV-2 lineage. Depending on when seroconversion is assumed to occur,

268     these cumulative predictions can be compared against findings from a serological study that was

269     conducted over this time period in France (Le Vu et al. 2021). This study surveyed 3221 individuals,

270     finding that 0.41% of individuals (95% confidence interval = 0.05 to 0.88) had gotten infected

271     with SARS-CoV-2 by March 9 to 15, 2020. While these estimates fall slightly higher than our

272     predictions, we are considering only one SARS-CoV-2 lineage (albeit likely the dominant one

273     circulating during this time period), and would thus expect the cumulative positive proportion

274     we predict to be lower than overall (all lineage) serology estimates. Other reasons for possible

275     underestimation    involve    epidemiological    model    misspecification    and    inaccurate

276     parameterization, for example, of the extent of transmission heterogeneity $p_h$.

## Discussion

278     Here, we developed a phylodynamic inference approach to estimate epidemiological parameters

279     from virus sequence data. Our inference approach is a "tree-free" approach in that it does not

280     rely on the reconstruction of viral phylogenies to estimate model parameters. One benefit of

281     using a "tree-free" approach for parameter estimation of emerging viral pathogens is that, early

282   on in an epidemic or pandemic, phylogenetic uncertainty is significant, and tree-based

283   phylodynamic inference approaches would need to integrate over this uncertainty, which is often

284   times computationally intensive. A second benefit of using a "tree-free" approach is that

285   parameters of the model of sequence evolution do not need to be estimated, reducing degrees

286   of freedom considerably. Instead of viral phylogenies being the data that statistically interface

287   with the epidemiological models, we use a low-dimensional summary statistic of the sequence

288   data, namely the number of segregating sites present in temporally-adjacent sets of viral

289   sequences. Beyond being a "tree-free" approach, our inference approach also benefits from

290   being "plug-and-play" in that it can easily accommodate any arbitrarily complex (or simple)

291   epidemiological model structure.

292   Based on fits to a simulated data set, we have shown that segregating site trajectories are highly

293   informative of epidemiological parameters such as $R_0$ and the timing of the index case $t_0$.  As far

294   as we are aware, only one other peer-reviewed tree-free phylodynamic inference method exists

295   (Kim et al. 2017), and future work should compare the approach developed here against this and

296   potentially other phylodynamic inferences approaches.

297   Although there are clear benefits of the phylodynamic inference approach detailed here, it still

298   relies on several assumptions that are also shared by other phylodynamic inference methods.

299   Most notably, it relies on an assumption of random sampling of individuals. However, in contrast

300   to coalescent-based models, the sampling rate does not have to be small relative to the number

301   of infected individuals. Phylodynamic inference based on birth-death-sampling models instead

302   requires the specification of a sampling process, such as a constant probability of an infected

303   individual being sampled upon recovery/death (Stadler 2010). Misspecification of the sampling

304   process can severely bias results, and much of the statistical power gained from these

305   approaches appears to arise from the sequence of sample times rather than genealogical

306   structure (Volz and Frost 2014). While our approach similarly requires an assumption of when

307   individuals are sampled, our approach provides considerable flexibility in what assumptions are

308   adopted, since the process model component of the state-space model can be easily

309   implemented under any number of assumptions of when individuals are available for sampling.

310   For example, in the compartmental model we used in the analysis of the France sequence data,

311  we could in principle assume that individuals could be sampled once they became infected during

312  a time window, rather than if they recovered during the time window.

313  The analysis we presented here focuses on phylodynamic inference using sequence data alone.

314  In recent years, there has also been a growing interest in combining multiple data sources – for

315  example, sequence data and epidemiological data or serological data - to more effectively

316  estimate model parameters. The few studies that have managed to incorporate additional data

317  while performing phylodynamic inference have shown the value in pursing this goal (Rasmussen

318  et al. 2011; Li et al. 2017). As a next step, we aim to extend the segregating sites approach

319  developed here to incorporate epidemiological data and/or serological data more explicitly.

320  Straightforward extension is possible due to the state-space model structure that is at the core

321  of the particle filtering routine we use. While the process model would stay the same, another

322  observation model can be added that relates the underlying state variables (e.g., $S$, $E$, $I$, $R$) to

323  observed case data for instance. This proposed approach mirrors a previously described

324  approach (Rasmussen et al. 2011), which showed that combining multiple data sources improved

325  parameter estimation.

326  Our analysis focused on phylodynamic inference based on sequence data belonging to a single

327  viral lineage, with a single index case. Our approach however can be expanded in a

328  straightforward manner to multiple viral lineages, each with their own index case. This is

329  especially useful in cases like SARS-CoV-2, where many regions have witnessed multiple clade

330  introductions in fueling the start of more local epidemics (Gonzalez-Reiche et al. 2020; Miller et

331  al. 2020). In this case, under the assumption that all lineages are phenotypically neutral and are

332  expanding in subpopulations experiencing the same epidemiological parameters (e.g., $R_0$), the

333  inference code can be expanded to estimate a single set of epidemiological parameters along

334  with multiple index case times, one corresponding to each viral lineage. When considering

335  multiple clades, a single segregating sites trajectory would be calculated for each clade, such that

336  multiple segregating site trajectories could be fit to at the same time.

337  Our approach can also be extended in a straightforward manner to consider multiple clades that

338  may be subject to different parameterizations for either intrinsic or extrinsic reasons. For

339    example, clades circulating in the same region may expand at different rates due to genetic

340    differences between the clades that confer a selective advantage of one clade over others. In this

341    case, multiple segregating site trajectories could again be calculated – one for each clade – and

342    phylodynamic inference would involve estimating epidemiological parameters, some of which

343    may be assumed to be similar across clades, while others such as $R_0$ may differ between clades.

344    As such, this inference method, which we initially developed for emerging pathogens with low

345    levels of genetic diversity, may continue to be useful for endemic pathogens when questions

346    involving emerging clades are a focus. Future work thus needs to determine when tree-free

347    phylodynamic inference provides advantages over tree-based phylodynamic inference, and

348    when tree-based methods provide better resolution into the dynamics of circulating virus

349    populations.

## Materials and Methods

351    **Epidemiological model simulations and calculation of segregating site trajectories.** We consider

352    epidemiological models of arbitrary complexity that incorporate demographic stochasticity using

353    Gillespie's $\tau$-leap algorithm. As a concrete example of such an epidemiological model, we here

354    use a susceptible-exposed-infected-recovered (SEIR) model whose dynamics are governed by the

355    following equations:

356    $S_{t+\Delta t} = S_t - N_{S \to E}$

357    $E_{t+\Delta t} = E_t + N_{S \to E} - N_{E \to I}$

358    $I_{t+\Delta t} = I_t + N_{E \to I} - N_{I \to R}$

359    $R_{t+\Delta t} = R_t + N_{I \to R}$

360    where:

361    $N_{S \to E} \sim Pois(\beta \frac{S_t}{N} I_t \Delta t)$

362    $N_{E \to I} \sim Pois(\gamma_E E_t \Delta t)$

363    $N_{I \to R} \sim Pois(\gamma_I I_t \Delta t)$

364   Here, $\beta$ is the transmission rate, $N$ is the host population size, $\gamma_E$ is the rate of transitioning from

365   the exposed to the infected class, $\gamma_I$ is the rate of recovering from infection, and $\Delta t$ is the $\tau$-leap

366   time step used. $R_0$ is given by $\beta / \gamma_I$. While the epidemiological dynamics of this model can be

367   simulated from the above equations alone, additional complexity is needed to incorporate virus

368   evolution throughout the time period of the simulation. To incorporate virus evolution, we

369   subcategorize both exposed individuals and infected individuals into genotype classes, with

370   genotype 1 being the reference genotype present at the start of the simulation. Mutations to the

371   virus occur at the time of transmission, with the number of mutations that occur in a single

372   transmission event given by a Poisson random variable with mean $\mu$, the per genome per

373   transmission event mutation rate. We assume infinite sites such that new mutations necessarily

374   result in new genotypes. New genotypes are numbered chronologically according to their

375   appearance. When new mutations are generated at a transmission event, the new genotype is

376   assumed to harbor the same mutation(s) as its infecting genotype plus any new mutations, which

377   are similarly numbered chronologically based on appearance. We use a sparse matrix approach

378   to store genotypes and their associated mutations to save on memory.

379   Given this model, during a time step $\Delta t$, $N_{E \to I}$ individuals are drawn at random from the set of

380   individuals who are currently exposed; these will be the individuals who will transition to the

381   infected class during this time step. Similarly, $N_{I \to R}$ individuals at drawn at random from the set

382   of individuals who are currently infected; these will be the individuals who will transition to the

383   recovered class during this time step. We further add $N_{S \to E}$ new individuals to the set of exposed

384   class during time step $\Delta t$. For each newly exposed individual, we randomly choose (with

385   replacement) a currently infected individual as its 'parent'. If no mutations occur during

386   transmission, then this new individual enters the same genotype class of its parent. If one or

387   more mutations occur during transmission, then this new individual enters a new genotype class,

388   and the sparse matrix is extended to document the new genotype and its associated mutations.

389   We start the simulation with one infected individual carrying a viral genotype that we consider

390   as the 'reference' genotype (genotype 1). To calculate a time series of segregating sites, we

391   define a time window length T ($T > \Delta t$) of a certain number of days and partition the simulation

392    time course into discrete, non-overlapping time windows. During simulation, we keep track of

393    the individuals that recover (transition from I to R) within a time window. For each time window

394    $i$, we then sample $n_i$ of these individuals at random, where $n_i$ is the number of sequences sampled

395    in a given time window based on the sampling scheme chosen. Because we have the genotypes

396    of the sampled individuals from the sparse matrix, we can calculate for any time window $i$, the

397    number of segregating sites $S_i$. $S_i$ is simply the number of polymorphic sites across the sampled

398    individuals in time window $i$.

399    **Phylodynamic inference using time series of segregating sites.** Our phylodynamic inference

400    approach relies on particle filtering, also known as Sequential Monte Carlo (SMC), to estimate

401    model parameters and reconstruct latent state variables. The underlying forward model we use

402    is formulated as a state-space model, with epidemiological variables (e.g., $S$, $E$, $I$, and $R$) being

403    latent/unobserved variables in the process model. The model is simulated using Gillespie's $\tau$-leap

404    algorithm, as described in the section above. The evolutionary component of the model also

405    contributes to the process model. For the observation model, we perform $k$ 'grabs' of sampled

406    individuals, with each 'grab' consisting of the following steps:

- pick (without replacement) $n_i$ individuals from the set of individuals who recovered during
407    
408        time window $i$, where $n_i$ is the number of samples observed in the empirical dataset in

409        window $i$.  We sample the same number of individuals as in the segregating sites dataset

410        that the model interfaces with, since sampling effort impacts the number of segregating

411        sites.

- calculate the simulated number of segregating sites $S_i^{sim}$, based on the genotypes of the
412    
413        sampled $n_i$ individuals (and their associated mutations).

414    Between 'grabs', replacement of previously sampled individuals occurs. We then calculate the

415    mean number of segregating sites for window $i$ by taking the average of all $k$ $S_i^{sim}$ values. Finally,

416    we calculate the probability of observing $S_i$ segregating sites in window $i$, given the model-

417    simulated mean number of segregating sites, using a Poisson probability density function

418    parameterized with the mean $S_i^{sim}$ value and evaluated at $S_i$. We use a Poisson probability density

419    function based on our observation that a Poisson distribution with the mean number of

420     segregating sites captures the distribution of $S_i^{sim}$ values from the 'grabs' effectively (Figure S2).

421     These probabilities serve as the weights for the particles. Particle weights are calculated at the

422     end of each time window with $n_i > 0$. Particles are resampled at the end of each of these time

423     windows according to their assigned weights. Particles with stochastic extinction of the virus prior

424     to the end of the last time window with $n_i > 0$ have weights set to 0 in time window $i$. If the

425     number of sampled individuals $n_i$ in time window $i$ exceeds the total number of individuals who

426     recovered in time window $i$, the particle weight is similarly set to 0. We run 10 SMC simulations

427     for each parameter set considered, resulting in 10 log-likelihood values.

428     For maximum likelihood estimation, weighted quadratic fitting is used, which is adapted from

429     Ionides et al. (2017). First, we use local quadratic smoothing (*LOESS*) with a span of 0.75 to obtain

430     the peak of the log-likelihood surface. The weight of each data point is determined by the

431     distance between this peak, using the tri-cube weight function. After excluding data points with

432     smaller weights by filtering out the smallest $\lambda \times 100$ percent, a quadratic function is fitted to data

433     points based on weights. For Figure 2C, the $\lambda$ for the quadratic fit was set to 0.5. For Figure 3B,

434     the $\lambda$ was set to 0.75, and for Figure 3C, the $\lambda$ was set to 0.55. Latent state variables are

435     reconstructed by randomly sampling a particle's $x_{0:tend}$ at the end of an SMC simulation, where

436     $t_{end}$ is the date at which the last sampled time window ends. All of our SMC simulations were

437     performed with 200 particles. We used $k = 100$ 'grabs' for the simulated data and, in the interest

438     of time, $k = 50$ 'grabs' for the France data.

439     Note that the complexity of this phylodynamic method is largely independent of the number of

440     input sequences, in contrast to phylodynamic inference approaches that rely on integrating over

441     phylogenetic uncertainty with BEAST.

442     **Implementation of the transmission heterogeneity model.** We implement transmission

443     heterogeneity by subcompartmentalizing the infected classes into a high-transmission and a low-

444     transmission class, as has been done elsewhere (Volz and Siveroni 2018; Miller et al. 2020). For

445     an SEIR model, the model extended to incorporate transmission heterogeneity becomes:

446     $S_{t+\Delta t} = S_t - N_{S \to E}$

447 $\quad E_{t+\Delta t} = E_t + N_{S \to E} - N_{E \to I_h} - N_{E \to I_l}$

448 $\quad I_{h,t+\Delta t} = I_{h,t} + N_{E \to I_h} - N_{I_h \to R}$

449 $\quad I_{l,t+\Delta t} = I_{l,t} + N_{E \to I_l} - N_{I_l \to R}$

450 $\quad R_{t+\Delta t} = R_t + N_{I_h \to R} + N_{I_l \to R}$

451 $\quad$ where:

452 $\quad N_{S \to E} \sim Pois(\beta_h \frac{S_t}{N} I_{h,t} \Delta t) + Pois(\beta_l \frac{S_t}{N} I_{l,t} \Delta t)$

453 $\quad N_{E \to I} \sim Pois(\gamma_E E_t \Delta t)$

454 $\quad N_{E \to I_h} \sim Bin(N_{E \to I}, p_h)$

455 $\quad N_{E \to I_l} = N_{E \to I} - N_{E \to I_h}$

456 $\quad N_{I_h \to R} \sim Pois(\gamma_I I_{h,t} \Delta t)$

457 $\quad N_{I_l \to R} \sim Pois(\gamma_I I_{l,t} \Delta t)$

458 $\quad$ The parameter $p_h$ quantifies the proportion of exposed individuals who transition to the highly

459 $\quad$ infectious $I_h$ class. Parameters $\beta_h$ and $\beta_l$ quantify the transmission rates of the infectious classes

460 $\quad$ that have high and low transmissibility, respectively. We set the values of $\beta_h$ and $\beta_l$ based on a

461 $\quad$ given parameterization of overall $R_0$ and the parameter $p_h$. To do this, we first define, as in

462 $\quad$ previous work (Volz and Siveroni 2018; Miller et al. 2020), the relative transmissibility of infected

463 $\quad$ individuals in the $I_h$ and $I_l$ classes as $c = \frac{\beta_h}{\beta_l}$. We further define a parameter $P$ as the fraction of

464 $\quad$ secondary infections that resulted from a fraction $p_h$ of the most transmissible infected

465 $\quad$ individuals. Based on given values of $p_h$ and $P$, we set $c$, as in previous work (Miller et al. 2020),

466 $\quad$ to $\frac{\left[\frac{1-p_h}{p_h}\right]}{\left[\frac{1}{P}-1\right]}$. With $c$ defined in this way, $p_h$ is interpreted as the proportion of most infectious

467 $\quad$ individuals that result in $P = 80\%$ of secondary infections. Recognizing that $R_0 = \frac{p_h \beta_h + (1-p_h)\beta_l}{\gamma_I}$ in

468 $\quad$ this model, we can then solve for $\beta_l$: $\frac{R_0 \gamma_I}{p_h c + (1-p_h)}$, and set $\beta_h = c\beta_l$.

469 **Epidemiological model structure and parameterization used for the France analysis.**

470 The process model we use in our phylodynamic inference of the France sequence data is based

471 on a previously published epidemiological model for SARS-COV-2 in France (Salje et al. 2020). We

472 base our process model on this published model to allow for a direct comparison of inferred $R_0$

473 values between our sequence-based analysis and their analysis that focuses over a similar time

474 period. Their analysis was based on fitting an epidemiological model to a combination of case,

475 hospitalization, and death data. Their model structure, implemented using Gillespie's $\tau$-leap

476 algorithm, is given by:

477 $S_{t+\Delta t} = S_t - N_{S \to E1}$

478 $E_{1,t+\Delta t} = E_{1,t} + N_{S \to E1} - N_{E1 \to E2}$

479 $E_{2,t+\Delta t} = E_{2,t} + N_{E1 \to E2} - N_{E2 \to I}$

480 $I_{t+\Delta t} = I_t + N_{E2 \to I} - N_{I \to R}$

481 $R_{t+\Delta t} = R_t + N_{I \to R}$

482 where:

483 $N_{S \to E1} \sim Pois(\beta \frac{S_t}{N} I_t \Delta t) + Pois(\beta \frac{S_t}{N} E_{2,t} \Delta t)$

484 $N_{E1 \to E2} \sim Pois(\gamma_{E1} E_{1,t} \Delta t)$

485 $N_{E2 \to I} \sim Pois(\gamma_{E2} E_{2,t} \Delta t)$

486 $N_{I \to R} \sim Pois(\gamma_I I_t \Delta t)$

487 with $\beta$ being the transmission rate, the average duration of time spent in the $E_1$ class given by

488 $1/\gamma_{E1}$ = 4 days, the average duration of time spent in the $E_2$ class given by $1/\gamma_{E2}$ = 1 day, and the

489 average duration of time spent in the infected class given by $1/\gamma_I$ = 3 days. While exposed class

490 2 ($E_2$) and infected class $I$ both transmit as efficiently, their model contains this level of detail to

491 more effectively interface with the case data, where symptoms do not appear before an

492 individual is infected (in class $I$). We keep with this model, rather than reducing it to having only

493     a single exposed class and a single infectious class to keep the same distribution of infected times

494     as in their model.

495     Because SARS-CoV-2 dynamics are characterized by substantial levels of transmission

496     heterogeneity (Adam et al. 2020; Miller et al. 2020; Sun et al. 2021) and we have shown in Figure

497     1 that transmission heterogeneity impacts segregating site trajectories, we expanded the

498     compartmental epidemiological model described above to include transmission heterogeneity in

499     a manner similar to the one we used in Figures 1E, F. Based specifically on the analysis by Sun

500     and coauthors (Sun et al. 2021), we set $p_h$ to 0.15, such that 15% of infections are responsible for

501     80% of secondary infections.

502     **Estimation of the per genome, per transmission event mutation rate**

503     We set the per-genome, per-transmission mutation rate parameter $\mu$ to 0.33. This is based on

504     the fit of a Poisson distribution to the number of *de novo* substitutions between 87 transmission

505     pairs of SARS-CoV-2 from four studies (James et al. 2020; Popa et al. 2020; Braun et al. 2021;

506     Lythgoe et al. 2021). Accession numbers for 78/87 of these transmission pairs are available in

507     Table S1. Accession numbers for the remaining pairs were provided by the corresponding authors

508     of the relevant publication (Lythgoe et al. 2021) Sequence data were aligned to Wuhan/Hu-1

509     (MN908947.3) (Wu et al. 2020) using MAFFT v.7.464 (Katoh 2002). Insertions relative to

510     Wuhan/Hu-1 were removed and the first 55 and last 100 nucleotides of the genome were masked.

511     *De Novo* substitutions for each pair were identified in Python v.3.9.4 (http://www.python.org)

512     using NumPy v.1.19.4 (Harris et al. 2020). Ambiguous nucleotides were considered in the

513     identification of *de novo* substitutions (i.e. an R nucleotide was assumed to match both an A and

514     a G). The mean number of substitutions between transmission pairs is the Maximum Likelihood

515     Estimate for the λ parameter of the Poisson distribution. The 95% confidence intervals were

516     calculated using the exact method using SciPy v.1.5.4 (SciPy 1.0 Contributors et al. 2020) such

517     that the lower bound was $\frac{(X^2_{2Y,0.025})/2}{87}$ and the upper bound was $\frac{(X^2_{2(Y+1),0.975})/2}{87}$ where Y is the total

518     number of observed substitutions.

519   The value for $\mu$ = 0.33 is consistent with population-level substitution rate estimates for SARS-

520   CoV-2, which range from 7.9 x $10^{-4}$ to 1.1 x $10^{-3}$ substitutions per site per year (Duchene et al.

521   2020; Pekar et al. 2020). With a genome length of SARS-CoV-2 of approximately 30,000

522   nucleotides and a generation interval of approximately 4.5 days (Griffin et al. 2020), these

523   population-level substitution rates would correspond to per genome, per transmission mutation

524   rates of between 0.29 and 0.41, respectively.

525   **Estimation of segregating site trajectories for the France data.**

526   We downloaded all complete and high-coverage SARS-CoV-2 sequences with complete sampling

527   dates sampled through March 17[th], 2020 (https://www.france24.com/en/20200316-live-france-

528   s-macron-addresses-nation-amid-worsening-coronavirus-outbreak) in France and uploaded

529   through April 29[th], 2021 from GISAID (Shu and McCauley 2017). Sequences were aligned to

530   Wuhan/Hu-1 using MAFFT v.7.464 Insertions relative to Wuhan/Hu-1 were removed. Any

531   sequences with fewer than 28000 A, C, T, or G characters were removed. Following this filtering

532   protocol our dataset included 479 sequences. We masked the first 55 and last 100 nucleotides in

533   the genome as well as positions marked as "highly homoplasic" in early SARS-CoV-2 sequencing

534   data                (https://github.com/W-L/ProblematicSites_SARS-CoV2/blob/master/archived_vcf/

535   problematic_sites_sarsCov2.2020-05-27.vcf). Pairwise SNP distances were calculated in a

536   manner that accounted for IUPAC ambiguous nucleotides in Python using NumPy. To subset

537   these data to a single clade circulating within France, we identified the connected components

538   of this pairwise distance matrix with a cutoff of 1 SNP in Python using SciPy and identified the

539   shared SNPs relative to Wuhan/Hu-1 between all sequences in each connected component. The

540   largest connected component contained 308 sequences which shared the substitutions C241T,

541   C3037T, C14408T, and A23403G. Our final dataset included these 308 as well as 122 sequences

542   from connected components that shared these four substitutions relative to Wuhan/Hu-1. We

543   included connected components in which all sequences had an N at any of the four clade-defining

544   sites of the largest connected component. Two sequences were excluded as they differed from

545   all other sequences in the dataset by > 7 SNPs. This dataset is similar to the set of sequences

546   analyzed in Danesh et al. (2020). Sequences were binned into four-day windows, aligned such

547   that the last window ended on the latest sampling date, and the number of segregating sites in

548 each window calculated in Python using NumPy. Ambiguous nucleotides were considered in the

549 calculation of segregating sites.

550 **Phylogenetic analysis of SARS-CoV-2 sequences from France.**

551 To confirm that the subset of sequences from France obtained from finding connected

552 components formed an evolutionary lineage/clade, we first combined the 479 sequences

553 sampled from France with 100 randomly-selected complete, high-coverage, collected date

554 complete sequences sampled from outside France through March 17th, 2020 and uploaded to

555 GISAID through April 29th, 2021. These sequences were aligned to Wuhan/Hu-1 using MAFFT,

556 insertions were removed, and the sites described above were masked. This alignment was

557 concatenated with the aligned sequences from France. IQ-Tree v. 2.0.7 (Minh et al. 2020) was

558 used to construct a maximum likelihood phylogeny, and ModelFinder (Kalyaanamoorthy et al.

559 2017) was used to find the best fit nucleotide substitution model (GTR+F+I). Small branches were

560 collapsed. TreeTime v. 0.8.0 (Sagulenko et al. 2017) was used to remove any sequences with

561 more than four interquartile distances from the expected evolutionary rate, rooting at

562 Wuhan/Hu-1. Treetime was also used to generate a time-aligned phylogeny assuming a clock rate

563 of $1 \times 10^{-3}$ with a standard deviation of $5 \times 10^{-4}$, a skyline coalescent model, marginal time

564 reconstruction, accounting for covariation, and resolving polytomies.

565 Maximum likelihood phylogenies were visualized in Python using Matplotlib v. 3.3.3 (Hunter
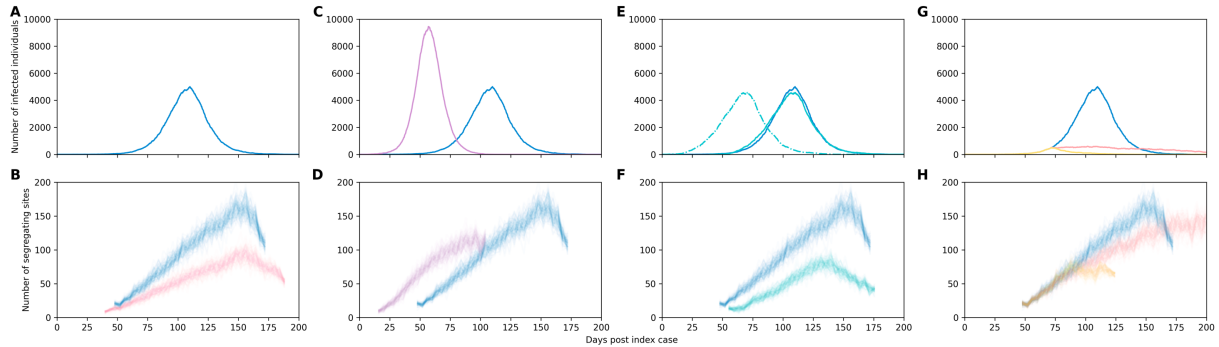
566 2007) and Baltic (https://github.com/evogytis/baltic).

567 **Availability of code.**

568 Python code used for generation of all figures is available on GitHub:

569 https://github.com/koellelab/segregating-sites

570

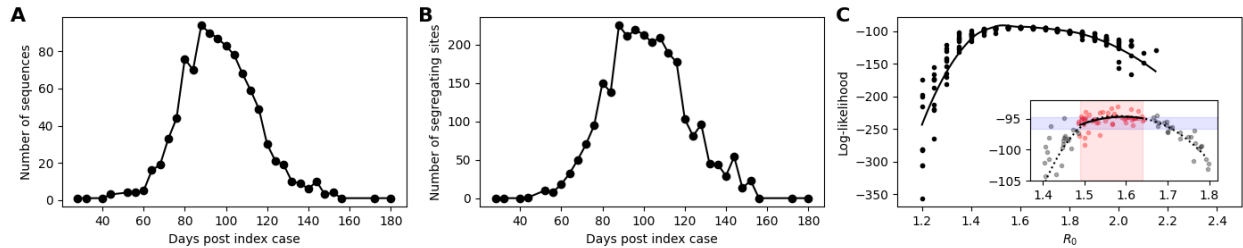571    **FIGURES**

572



573

574    **Figure 1. Segregating site trajectories under simulated epidemiological dynamics.** (A) Simulated

575    dynamics of infected individuals ($I$) under an SEIR model simulated with an $R_0$ of 1.6. (B) Segregating site

576    trajectories under dense and sparse sampling. Dense sampling (blue lines) corresponds to 40 sequences

577    sampled per time window. Sparse sampling (red lines) corresponds to 20 sequences sampled per time

578    window. (C) Simulated dynamics of infected individuals ($I$) under an SEIR model simulated with an $R_0$ of

579    2.0 (purple line) compared to those of the $R_0$ = 1.6 simulation (blue line). A higher transmission rate was

580    used to generate the higher $R_0$ value of 2.0. (D) Segregating site trajectories for the $R_0$ = 2.0 simulation

581    (purple lines) and the $R_0$ = 1.6 simulation (blue lines). Both simulations are densely sampled (40 sequences

582    sampled per time window). (E) Simulated dynamics of infected individuals ($I$) under an SEIR model with

583    an $R_0$ of 1.6 and incorporating transmission heterogeneity (teal, dashed line) compared to those of the

584    original $R_0$ = 1.6 simulation (blue line) without transmission heterogeneity. Transmission heterogeneity

585    was included by setting $p_h$ = 0.06, resulting in 6% of the most infectious individuals being responsible for

586    80% of secondary infections. For ease of comparing segregating site trajectories, the transmission

587    heterogeneity simulation was shifted later in time such that its epidemic peak aligned with the simulation

588    without transmission heterogeneity (teal, solid line). (F) Segregating site trajectories for the shifted

589    transmission heterogeneity simulation (teal lines) and the simulation without transmission heterogeneity

590    (blue line). Both simulations are densely sampled (40 sequences sampled per time window). (G) Simulated

591    dynamics of infected individuals ($I$) under an SEIR model simulated with changing $R_0$. Changes in $R_0$

592    occurred when the number of infected individuals reached 400. The simulation in red has $R_0$ decreasing

593    to 1.1. The simulation in yellow has $R_0$ decreasing to 0.75. The simulation in blue has $R_0$ remaining at 1.6.

594    (H) Segregating site trajectories for the three simulations shown in Figure 1G. All three simulations are

595     densely sampled (40 sequences sampled per time window). In all model simulations, $\gamma_E = 1/2$ days$^{-1}$,

596     $\gamma_I = 1/3$ days$^{-1}$, population size $N = 10^5$, and the per genome, per transmission mutation rate $\mu = 0.2$.

597     Initial conditions are $S(t_0)$ = N-1, $E(t_0) = 0$, $I(t_0) = 1$, and $R(t_0) = 0$. For the transmission heterogeneity

598     simulation (subplot E), initial conditions are $S(t_0)$ = N-1, $E(t_0) = 0$, $I_h(t_0) = 1$, $I_l(t_0) = 0$, and $R(t_0) = 0$. A time

599     step of $\tau = 0.1$ days was used in the Gillespie $\tau$-leap algorithm. Time windows of $T = 4$ days were used to

600     bin sequences for the segregating sites calculation. 100 different segregating site trajectories are shown

601     for each simulation.

602

603

**Figure 2. Phylodynamic inference on a simulated trajectory of segregating sites.** (A) The number of sampled sequences over time, by time window. Sampling was done in proportion to the number of individuals recovering in a time window. In all, 1000 sequences were sampled over the course of the simulated epidemic. The number of samples in a given time window was constrained to be ≤100. (B) Simulated segregating site trajectory from the sampled sequences. (C) Estimation of $R_0$ using SMC. Points show log-likelihood values from different SMC simulations across a range of $R_0$ values between 1.2 and 2.5, in 0.05 increments. Smoothed likelihood surface was obtained by *LOESS* smoothing with a span of 0.75. Inset: Maximum likelihood estimation of $R_0$ using quadratic fitting. Black points in inset show log-likelihood values from different SMC simulations across a range of $R_0$ values between 1.4 and 1.8. The vertical black dashed line shows the maximum likelihood estimate (MLE) of $R_0$ (1.59). The red band shows the 95% confidence interval of $R_0$ (1.49 − 1.64).  MLE and 95% CI were obtained from fitting a quadratic function to the log-likelihood values shown in the inset, using a similar approach to the one outlined in Ionides et al. (2017) with a $\lambda$ value of 0.5. 95% CI were set at the values of $R_0$ corresponding to the maximum likelihood value at the peak of the quadratic curve minus 1.92 log-likelihood units. Model parameters for the simulated data set are: $R_0$= 1.6,  $\gamma_E$= 1/2 days$^{-1}$, $\gamma_I$ = 1/3 days$^{-1}$, population size $N$ = $10^5$, $t_0$ = 0, and the per genome, per transmission mutation rate $\mu$ = 0.2. Initial conditions are $S(t_0)$ = N-1, $E(t_0)$ = 0, $I(t_0)$ = 1, and $R(t_0)$ = 0. A time step of $\tau$ = 0.1 days was used in the Gillespie $\tau$-leap algorithm. A time window of $T$ = 4 days was used to bin sequences for the segregating sites calculation.
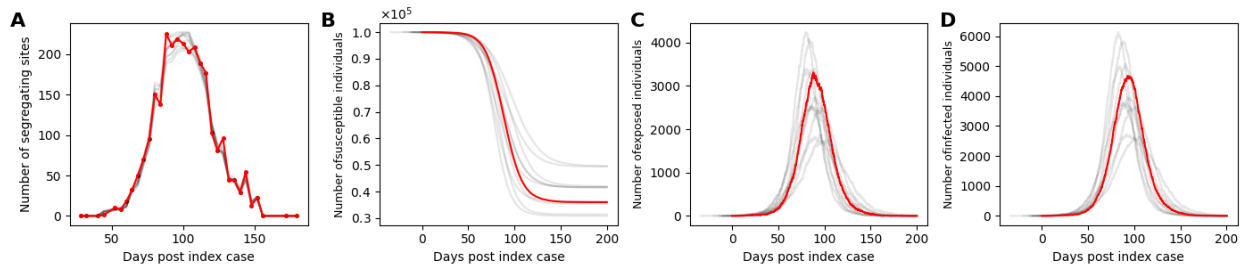
624

**Figure 3. Joint estimation of the basic reproduction number ($R_0$) and the timing of the index case ($t_0$) using simulated data.** (A) The likelihood surface based on the segregating site trajectory shown in Figure 2B is shown over a broad range of $R_0$ values (1.2 to 2.4, in 0.1 increments) and $t_0$ values (from 60 days prior to 56 days following the true $t_0$ of 0 in 2-day increments). Blank cells yielded log-likelihood values of <-1281. Log-likelihood values shown in each cell across this broad range of $R_0$ and $t_0$ are mean log-likelihood values calculated from 10 SMC simulations at each parameterization. (B-C) Profile likelihood for $R_0$ (B) and $t_0$ (C). Profile likelihoods were calculated using an approach similar to the one outlined in Ionides et al. (2017). The *LOESS* fit is shown with a dotted black line. The quadratic fit is shown with a solid black line. Points included in the quadratic fit are shown in red; points excluded from the quadratic fit are shown in gray. The shaded red area is the 95% confidence interval for the focal parameter. The shaded blue area shows the range of log-likelihood values that fall within 1.92 log-likelihood values of the quadratic fit's maximum value.

637

638

639



**Figure 4. Trajectories of reconstructed unobserved state variables for the simulated dataset.** (A) Simulated trajectory of the number of segregating sites (red), alongside reconstructed trajectory of the number of segregating sites from 10 sampled SMC particles (gray). For each SMC particle, a combination of $t_0$ and $R_0$ values of 10 SMC iterations were 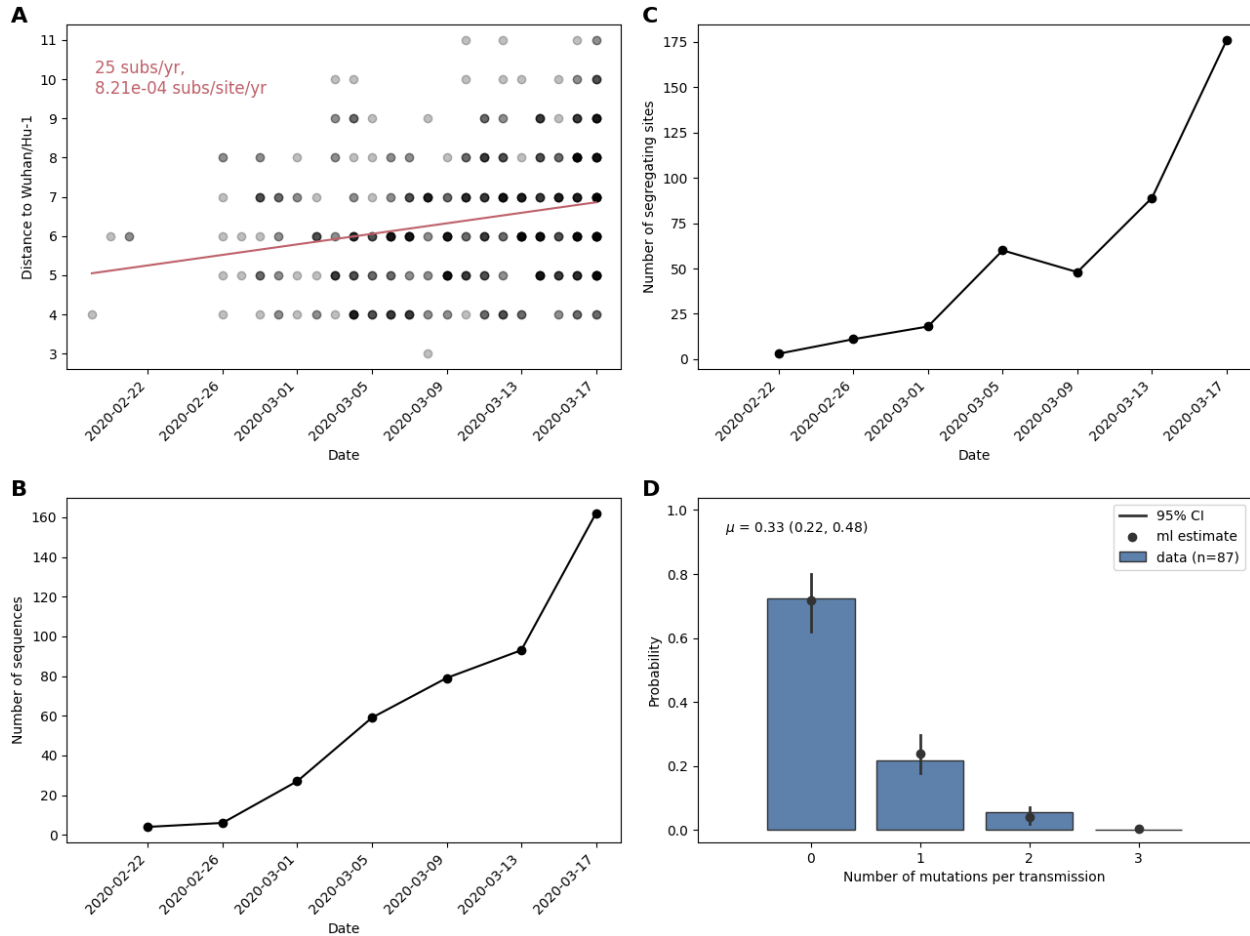randomly chosen based on their log-likelihood values. (B) Simulated dynamics of susceptible individuals (red), alongside reconstructed dynamics of susceptible individuals from these SMC simulations (gray). (E) Simulated dynamics of exposed individuals (red), alongside reconstructed dynamics of exposed individuals (gray). (F) Simulated dynamics of infected individuals (red), alongside reconstructed dynamics of infected individuals (gray).

648

649

**Figure 5. Sequences and parameters used in the estimation of $R_0$ and $t_0$ for the France data.** (A)
Sequences used in the phylodynamic analysis, plotted by their collection date and their nucleotide
divergence from the Wuhan/Hu-1 reference sequence. (B) The number of sampled sequences over time,
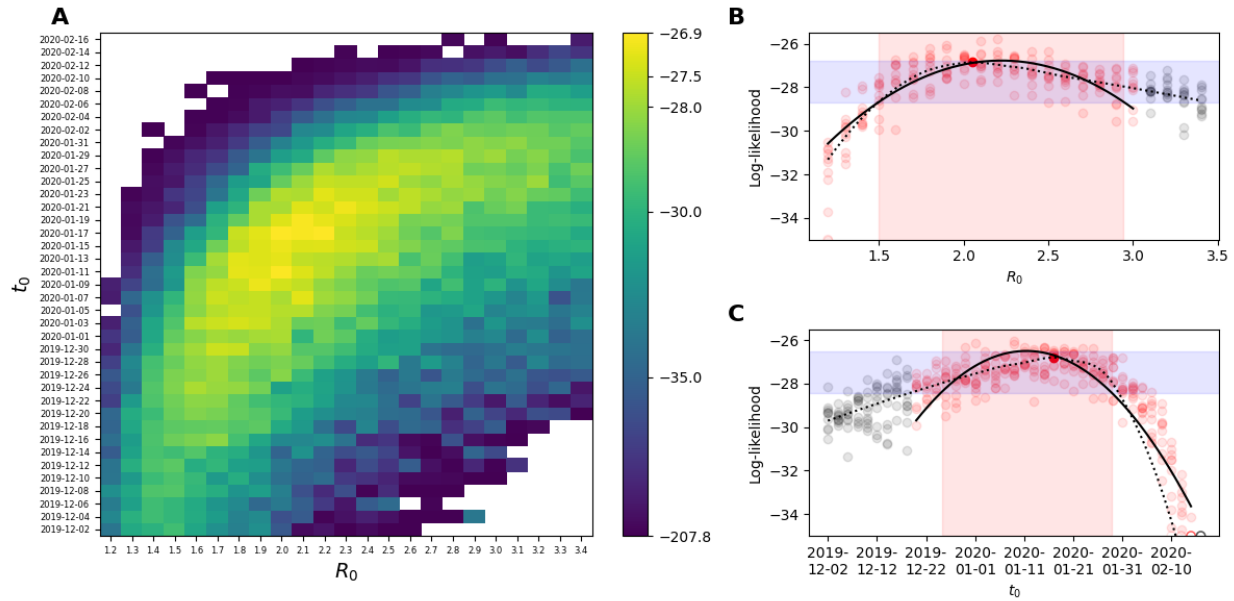calculated using a $T$ = 4 day time window. (C) The segregating site trajectory calculated from the sampled
sequences, using the same $T$ = 4 day time window shown in (B). (D) Estimation of the per genome, per
transmission mutation rate $\mu$. Blue histogram plots the fraction of transmission pairs with consensus
sequences that differed from one another by the number of mutations shown on the x-axis. The Poisson
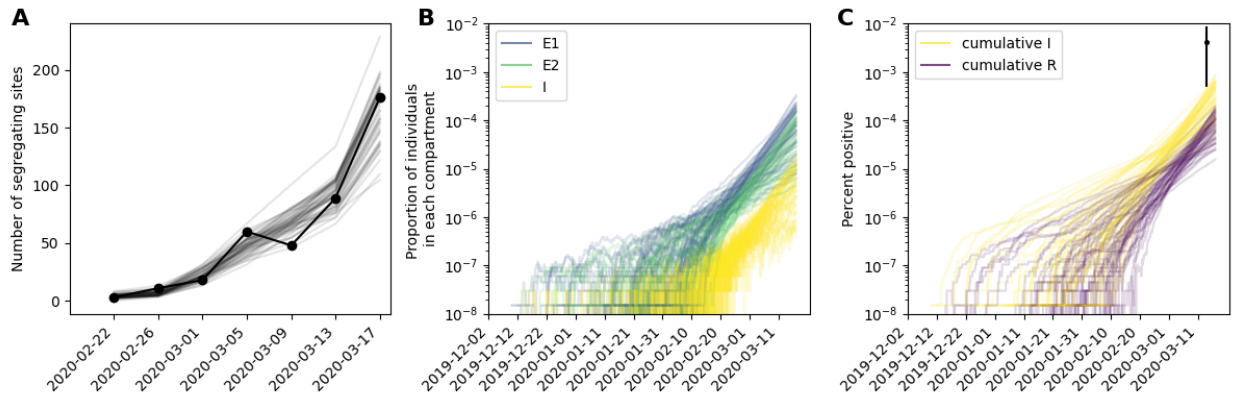estimate from these data, shown in black, was $\mu$ = 0.33 (95% CI = 0.22-0.48).

658

659

**Figure 6. Joint estimation of the basic reproduction number $R_0$ and the timing of the index case $t_0$ using the France data.** (A) The joint log-likelihood surface based on the estimated segregating site trajectory for the France data. Each cell is colored according to the mean of log-likelihood for a $t_0$, $R_0$ combination obtained from 10 SMC simulations. (B-C) Profile likelihood for $R_0$ (B) and $t_0$ (C). Profile likelihoods were calculated using an approach similar to the one outlined in Ionides et al. (2017). The *LOESS* fit is shown with a dotted black line. The quadratic fit is shown with a solid black line. Points included in the quadratic fit are shown in red; points excluded from the quadratic fit are shown in black. The shaded red and blue areas are, as in Figures 3B and 3C, the 95% confidence interval for the focal parameter and the range of log-likelihood values that fall within 1.92 log-likelihood values of the quadratic fit's maximum value.
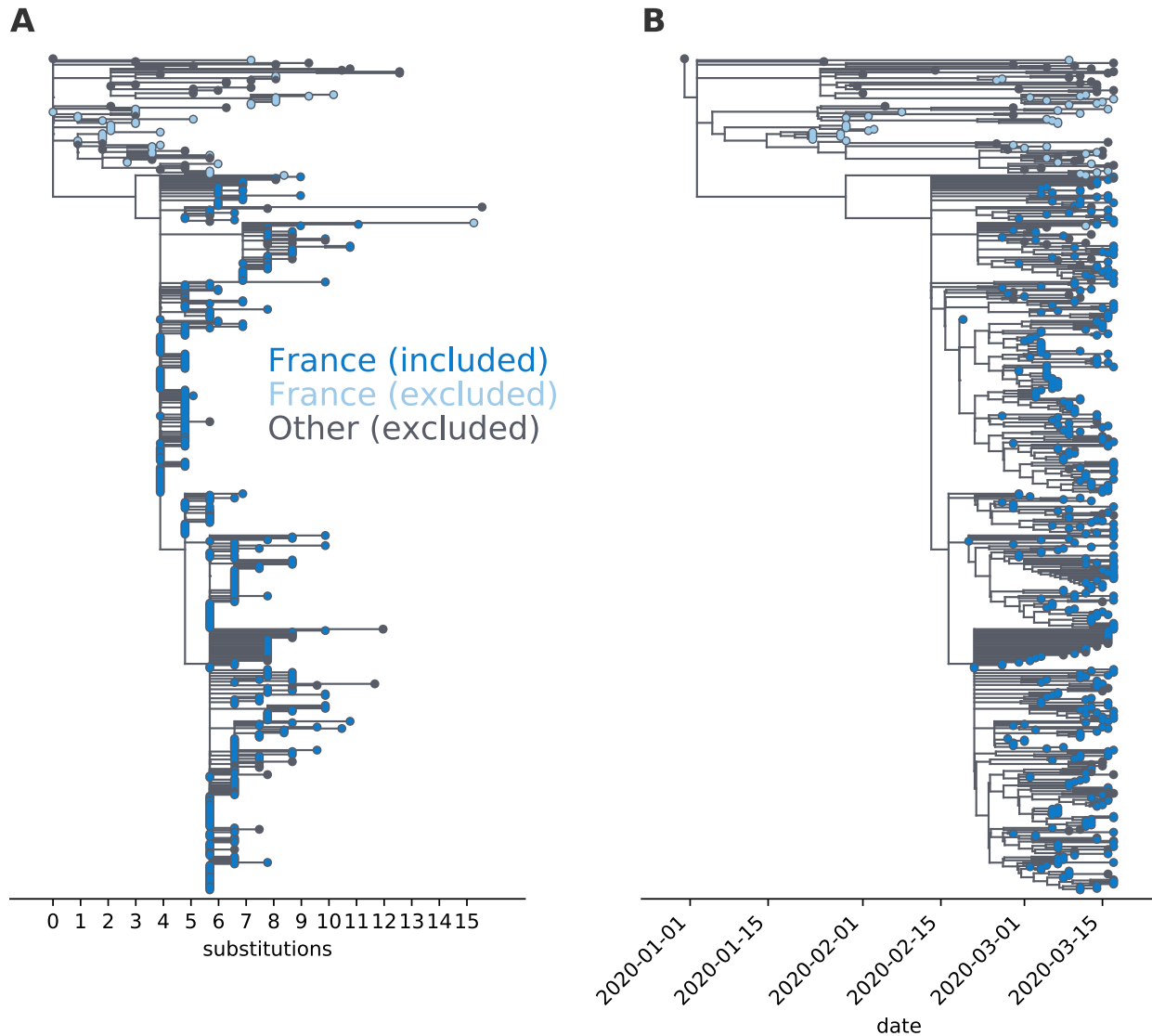
669

670

671



**Figure 7. Trajectories of reconstructed unobserved state variable for the France data.** For the reconstruction of all state variables shown, a combination of two parameters, $R_0$ and $t_0$, are sampled based on their log-likelihood values from 10 SMC simulations. (A) Segregating site trajectory for the France data, alongside segregating site trajectories from 10 sampled SMC particles. (B) Reconstructed dynamics of the number of individuals in the first exposed class ($E_1$), the second exposed class ($E_2$), and the infected class ($I$). (C) Cumulative number of exposed individuals (yellow) and cumulative number of recovered individuals (purple) over time. The maximum likelihood estimate of the fraction of the population that had been infected with SARS-CoV-2 by mid-March, and the 95% confidence interval of this estimate, are shown in black. Estimates are from a serological study conducted during the time window March 9-15, 2020 (Le Vu et al. 2021).

682

683

**Figure S1.** Inferred phylogenies for the sequences sampled from France, January 23-March 17, 2020. (A) Divergence tree, showing the number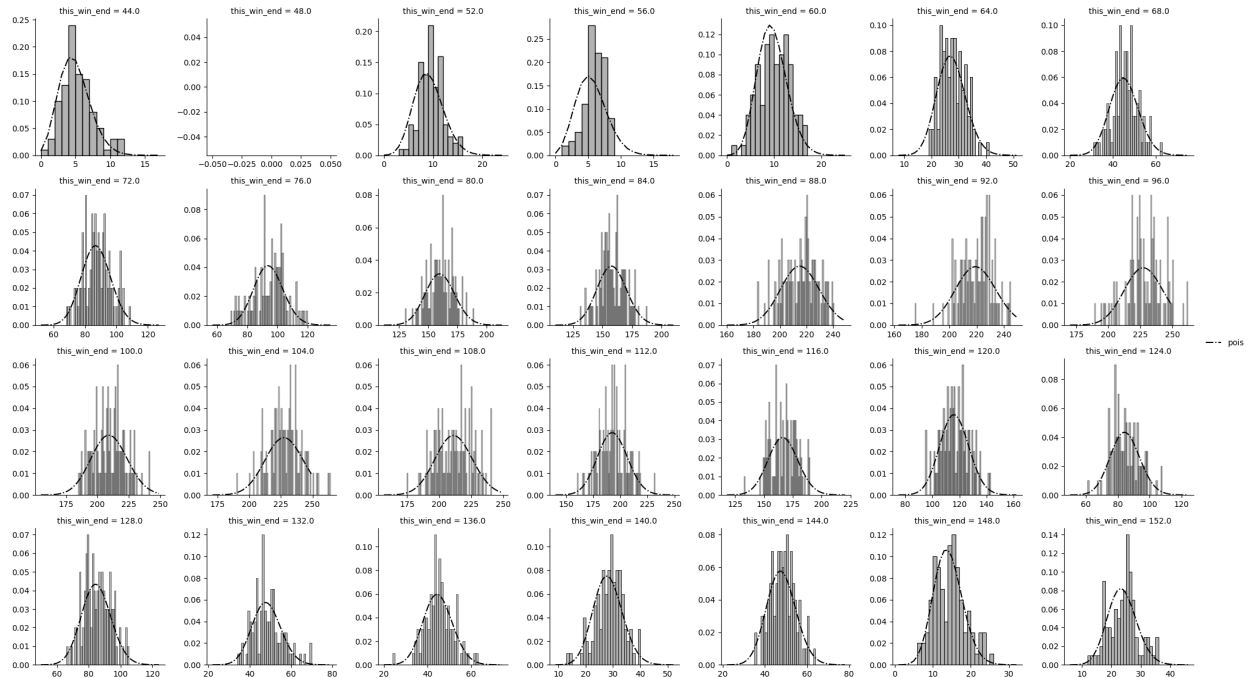 of nucleotide substitutions from Wuhan/Hu-1. Sequences from France are colored in blue, with dark blue coloring indicating sequences that were included in our single-lineage analysis and light blue coloring indicating sequences that were excluded from our analysis. Tips colored in gray denote genetically similar sequences sampled from outside of France during this time period. (B) Time-aligned maximum likelihood phylogeny, with coloring of sequences as in (A).

690

691

**Figure S2. Appropriateness of the Poisson distribution in the observation model.** Each subplot shows a

time window *i*, with the blue vertical line indicating the observed value in that time window, $S_i$. Each time

window further shows a histogram of $S_i^{sim}$ values from 100 'grabs' from a single randomly sampled particle.

The dash-dotted black curves show Poisson probability mass functions, parameterized with the average

of the $S_i^{sim}$ values.

697

**Table S1.** Transmission pairs used to estimate the per genome, per transmission event mutation rate $\mu$.

Accession numbers of the consensus sequences from the donor and the recipient of the transmission pair

are provided.

701

702

# References

Adam DC, Wu P, Wong JY, Lau EHY, Tsang TK, Cauchemez S, Leung GM, Cowling BJ. 2020. Clustering and superspreading potential of SARS-CoV-2 infections in Hong Kong. *Nat Med* 26:1714–1719.

Althouse BM, Wenger EA, Miller JC, Scarpino SV, Allard A, Hébert-Dufresne L, Hu H. 2020. Superspreading events in the transmission dynamics of SARS-CoV-2: Opportunities for interventions and control. *PLoS Biol* 18:e3000897.

709    Boskova V, Bonhoeffer S, Stadler T. 2014. Inference of epidemiological dynamics based on simulated
710         phylogenies using birth-death and coalescent models.Koelle K, editor. *PLoS Computational*
711         *Biology* 10:e1003913.

712    Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, Suchard MA, Rambaut A, Drummond AJ.
713         2014. BEAST 2: a software platform for Bayesian evolutionary analysis.Prlic A, editor. *PLoS*
714         *Computational Biology* 10:e1003537.

715    Braun K, Moreno G, Wagner C, Accola MA, Rehrauer WM, Baker D, Koelle K, O'Connor DH, Bedford T,
716         Friedrich TC, et al. 2021. Limited within-host diversity and tight transmission bottlenecks limit
717         SARS-CoV-2 evolution in acutely infected individuals. Evolutionary Biology Available from:
718         http://biorxiv.org/lookup/doi/10.1101/2021.04.30.440988

719    Danesh G, Elie B, Michalakis Y, Sofonea MT, Bal A, Behillil S, Destras G, Boutolleau D, Burrel S, Marcelin
720         A-G, et al. 2020. Early phylodynamics analysis of the COVID-19 epidemic in France. Epidemiology
721         Available from: http://medrxiv.org/lookup/doi/10.1101/2020.06.03.20119925

722    Duchene S, Featherstone L, Haritopoulou-Sinanidou M, Rambaut A, Lemey P, Baele G. 2020. Temporal
723         signal and the phylodynamic threshold of SARS-CoV-2. *Virus Evolution* 6:veaa061.

724    Gámbaro F, Behillil S, Baidaliuk A, Donati F, Albert M, Alexandru A, Vanpeene M, Bizard M, Brisebarre A,
725         Barbet M, et al. 2020. Introductions and early spread of SARS-CoV-2 in France, 24 January to 23
726         March 2020. *Eurosurveillance* [Internet] 25. Available from:
727         https://www.eurosurveillance.org/content/10.2807/1560-7917.ES.2020.25.26.2001200

728    Geidelberg L, Boyd O, Jorgensen D, Siveroni I, Nascimento FF, Johnson R, Ragonnet-Cronin M, Fu H,
729         Wang H, Xi X, et al. 2021. Genomic epidemiology of a densely sampled COVID-19 outbreak in
730         China. *Virus Evolution* 7:veaa102.

731    Gonzalez-Reiche AS, Hernandez MM, Sullivan MJ, Ciferri B, Alshammary H, Obla A, Fabre S, Kleiner G,
732         Polanco J, Khan Z, et al. 2020. Introductions and early spread of SARS-CoV-2 in the New York City
733         area. *Science*:eabc1917.

734    Griffin JM, Collins AB, Hunt K, McEvoy D, Casey M, Byrne AW, McAloon CG, Barber A, Lane EA, More SJ.
735         2020. A rapid review of available evidence on the serial interval and generation time of COVID-
736         19. Epidemiology Available from: http://medrxiv.org/lookup/doi/10.1101/2020.05.08.20095075

737    Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S,
738         Smith NJ, et al. 2020. Array programming with NumPy. *Nature* 585:357–362.

739    Hilton J, Keeling MJ. 2020. Estimation of country-level basic reproductive ratios for novel Coronavirus
740         (SARS-CoV-2/COVID-19) using synthetic contact matrices. *PLoS Comput Biol* 16:e1008031.

741    Hunter JD. 2007. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* 9:90–95.

742    Ionides EL, Breto C, Park J, Smith RA, King AA. 2017. Monte Carlo profile confidence intervals for
743         dynamic systems. *Journal of The Royal Society Interface* 14:20170126.

744 James SE, Ngcapu S, Kanzi AM, Tegally H, Fonseca V, Giandhari J, Wilkinson E, Chimukangara B, Pillay S,
745      Singh L, et al. 2020. High Resolution analysis of Transmission Dynamics of Sars-Cov-2 in Two
746      Major Hospital Outbreaks in South Africa Leveraging Intrahost Diversity. Infectious Diseases
747      (except HIV/AIDS) Available from:
748      http://medrxiv.org/lookup/doi/10.1101/2020.11.15.20231993

749 Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. 2017. ModelFinder: fast model
750      selection for accurate phylogenetic estimates. *Nat Methods* 14:587–589.

751 Katoh K. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier
752      transform. *Nucleic Acids Research* 30:3059–3066.

753 Keeling MJ, Rohani P. 2008. Modeling Infectious Diseases in Humans and Animals. Princeton University
754      Press

755 Kim K, Omori R, Ito K. 2017. Inferring epidemiological dynamics of infectious diseases using Tajima's D
756      statistic on nucleotide sequences of pathogens. *Epidemics* 21:21–29.

757 Koelle K, Rasmussen DA. 2012. Rates of coalescence for common epidemiological models at equilibrium.
758      *J. R. Soc. Interface* 9:997–1007.

759 Kucharski AJ, Russell TW, Diamond C, Liu Y, Edmunds J, Funk S, Eggo RM, Sun F, Jit M, Munday JD, et al.
760      2020. Early dynamics of transmission and control of COVID-19: a mathematical modelling study.
761      *The Lancet Infectious Diseases* 20:553–558.

762 Le Vu S, Jones G, Anna F, Rose T, Richard J-B, Bernard-Stoecklin S, Goyard S, Demeret C, Helynck O,
763      Escriou N, et al. 2021. Prevalence of SARS-CoV-2 antibodies in France: results from nationwide
764      serological surveillance. *Nat Commun* 12:3025.

765 Lemieux JE, Siddle KJ, Shaw BM, Loreth C, Schaffner SF, Gladden-Young A, Adams G, Fink T, Tomkins-
766      Tinch CH, Krasilnikova LA, et al. 2021. Phylogenetic analysis of SARS-CoV-2 in Boston highlights
767      the impact of superspreading events. *Science* 371.

768 Li LM, Grassly NC, Fraser C. 2017. Quantifying Transmission Heterogeneity Using Both Pathogen
769      Phylogenies and Incidence Time Series. *Molecular Biology and Evolution* 34:2982–2995.

770 Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM. 2005. Superspreading and the effect of individual
771      variation on disease emergence. *Nature* 438:355–359.

772 Locatelli I, Trächsel B, Rousson V. 2021. Estimating the basic reproduction number for COVID-19 in
773      Western Europe.Khudyakov YE, editor. *PLoS ONE* 16:e0248731.

774 Lythgoe KA, Hall M, Ferretti L, de Cesare M, MacIntyre-Cockett G, Trebes A, Andersson M, Otecko N,
775      Wise EL, Moore N, et al. 2021. SARS-CoV-2 within-host diversity and transmission. *Science*
776      372:eabg0821.

777 Miller D, Martin MA, Harel N, Tirosh O, Kustin T, Meir M, Sorek N, Gefen-Halevi S, Amit S, Vorontsov O,
778      et al. 2020. Full genome viral sequences inform patterns of SARS-CoV-2 spread into and within
779      Israel. *Nat Commun* 11:5518.

780    Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R. 2020. IQ-
781        TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic
782        Era.Teeling E, editor. *Molecular Biology and Evolution* 37:1530–1534.

783    Pekar J, Worobey M, Moshiri N, Scheffler K, Wertheim JO. 2020. Timing the SARS-CoV-2 Index Case in
784        Hubei Province. Evolutionary Biology Available from:
785        http://biorxiv.org/lookup/doi/10.1101/2020.11.20.392126

786    Popa A, Genger J-W, Nicholson MD, Penz T, Schmid D, Aberle SW, Agerer B, Lercher A, Endler L, Colaço
787        H, et al. 2020. Genomic epidemiology of superspreading events in Austria reveals mutational
788        dynamics and transmission properties of SARS-CoV-2. *Sci. Transl. Med.* 12:eabe2555.

789    Popinga A, Vaughan T, Stadler T, Drummond AJ. 2014. Inferring epidemiological dynamics with Bayesian
790        coalescent inference: the merits of deterministic and stochastic models. *Genetics*.

791    Rasmussen DA, Boni MF, Koelle K. 2014. Reconciling phylodynamics with epidemiology: the case of
792        dengue virus in southern Vietnam. *Molecular Biology and Evolution* 31:258–271.

793    Rasmussen DA, Ratmann O, Koelle K. 2011. Inference for nonlinear epidemiological models using
794        genealogies and time series. *PLoS Comput Biol* 7:e1002136.

795    Rasmussen DA, Stadler T. Coupling adaptive molecular evolution to phylodynamics using fitness-
796        dependent birth-death models. *Evolutionary Biology*:24.

797    Ratmann O, Hodcroft EB, Pickles M, Cori A, Hall M, Lycett S, Colijn C, Dearlove B, Didelot X, Frost S, et al.
798        2017. Phylogenetic Tools for Generalized HIV-1 Epidemics: Findings from the PANGEA-HIV
799        Methods Comparison. *Molecular Biology and Evolution* 34:185–203.

800    Sagulenko P, Puller V, Neher R. 2017. TreeTime: maximum likelihood phylodynamic analysis.

801    Salje H, Tran Kiem C, Lefrancq N, Courtejoie N, Bosetti P, Paireau J, Andronico A, Hozé N, Richet J,
802        Dubost C-L, et al. 2020. Estimating the burden of SARS-CoV-2 in France. *Science* 369:208–211.

803    SciPy 1.0 Contributors, Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D,
804        Burovski E, Peterson P, Weckesser W, et al. 2020. SciPy 1.0: fundamental algorithms for
805        scientific computing in Python. *Nat Methods* 17:261–272.

806    Shu Y, McCauley J. 2017. GISAID: Global initiative on sharing all influenza data – from vision to reality.
807        *Eurosurveillance* [Internet] 22. Available from:
808        https://www.eurosurveillance.org/content/10.2807/1560-7917.ES.2017.22.13.30494

809    Stadler T. 2010. Sampling-through-time in birth–death trees. *Journal of Theoretical Biology* 267:396–
810        404.

811    Stadler T, Bonhoeffer S. 2013. Uncovering epidemiological dynamics in heterogeneous host populations
812        using phylogenetic methods. *Phil. Trans. R. Soc. B* [Internet] 368. Available from:
813        http://rstb.royalsocietypublishing.org/content/368/1614/20120198

814   Stadler T, Kouyos R, Wyl V von, Yerly S, Böni J, Bürgisser P, Klimkait T, Joos B, Rieder P, Xie D, et al. 2012.
815         Estimating the basic reproductive number from viral sequence data. *Mol Biol Evol* 29:347–357.

816   Stadler T, Kuhnert D, Bonhoeffer S, Drummond AJ. 2013. Birth-death skyline plot reveals temporal
817         changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proceedings of the National*
818         *Academy of Sciences* 110:228–233.

819   Stadler T, Kühnert D, Rasmussen DA, du Plessis L. 2014. Insights into the Early Epidemic Spread of Ebola
820         in Sierra Leone Provided by Viral Sequence Data. *PLoS Curr* [Internet]. Available from:
821         https://currents.plos.org/outbreaks/article/insights-into-the-early-epidemic-spread-of-ebola-in-
822         sierra-leone-provided-by-viral-sequence-data/

823   Sun K, Wang W, Gao L, Wang Y, Luo K, Ren L, Zhan Z, Chen X, Zhao S, Huang Y, et al. 2021. Transmission
824         heterogeneities, kinetics, and controllability of SARS-CoV-2. *Science* 371:eabe2424.

825   Vaughan TG, Leventhal GE, Rasmussen DA, Drummond AJ, Welch D, Stadler T. 2017. Directly Estimating
826         Epidemic Curves From Genomic Data. Available from:
827         http://biorxiv.org/lookup/doi/10.1101/142570

828   Volz EM. 2012. Complex population dynamics and the coalescent under neutrality. *Genetics* 190:187–
829         201.

830   Volz EM, Frost SDW. 2014. Sampling through time and phylodynamic inference with coalescent and
831         birth-death models. *Journal of The Royal Society Interface* 11:20140945–20140945.

832   Volz EM, Ndembi N, Nowak R, Kijak GH, Idoko J, Dakum P, Royal W, Baral S, Dybul M, Blattner WA, et al.
833         2017. Phylodynamic analysis to inform prevention efforts in mixed HIV epidemics. *Virus*
834         *Evolution* [Internet] 3. Available from: https://academic.oup.com/ve/article-
835         lookup/doi/10.1093/ve/vex014

836   Volz EM, Siveroni I. 2018. Bayesian phylodynamic inference with complex models.Darling AE, editor.
837         *PLoS Comput Biol* 14:e1006546.

838   Woolhouse MEJ, Dye C, Etard J-F, Smith T, Charlwood JD, Garnett GP, Hagan P, Hii JLK, Ndhlovu PD,
839         Quinnell RJ, et al. 1997. Heterogeneities in the transmission of infectious agents: Implications for
840         the design of control programs. *Proceedings of the National Academy of Sciences* 94:338–342.

841   Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian J-H, Pei Y-Y, et al. 2020. A new
842         coronavirus associated with human respiratory disease in China. *Nature* 579:265–269.

843