

1

Version dated: January 23, 2023

2 RH: Deep Learning and Phylogeography

3 **Deep learning approaches to viral phylogeography are**
4 **fast and as robust as likelihood methods to model**
5 **misspecification**

6 AMMON THOMPSON^{1,*}, BENJAMIN LIEBESKIND¹, ERIK J. SCULLY¹, MICHAEL LANDIS^{2,*}

7 ¹*National Geospatial-Intelligence Agency, Springfield, VA, 22150, USA*

8 ²*Department of Biology, Washington University in St. Louis, Rebstock Hall, St. Louis, Missouri,*
9 *63130, USA*

10 ***Corresponding authors:** E-mail: Ammon.M.Thompson.ctr@nga.mil and

11 michael.landis@wustl.edu.

12 *Abstract.*— Analysis of phylogenetic trees has become an essential tool in epidemiology.
13 Likelihood-based methods fit models to phylogenies to draw inferences about the
14 phylodynamics and history of viral transmission. However, these methods are
15 computationally expensive, which limits the complexity and realism of phylodynamic
16 models and makes them ill-suited for informing policy decisions in real-time during rapidly
17 developing outbreaks. Likelihood-free methods using deep learning are pushing the
18 boundaries of inference beyond these constraints. In this paper, we extend, compare and
19 contrast a recently developed deep learning method for likelihood-free inference from trees.
20 We trained multiple deep neural networks using phylogenies from simulated outbreaks that
21 spread among five locations and found they achieve similar levels of accuracy to Bayesian
22 inference under the true simulation model. We compared robustness to model
23 misspecification of a trained neural network to that of a Bayesian method. We found that
24 both models had comparable performance, converging on similar biases. We also trained
25 and tested a neural network against phylogeographic data from a recent study of the
26 SARS-Cov-2 pandemic in Europe and obtained similar estimates of epidemiological
27 parameters and the location of the common ancestor in Europe. Along with being as
28 accurate and robust as likelihood-based methods, our trained neural networks are on
29 average over 3 orders of magnitude faster. Our results support the notion that neural
30 networks can be trained with simulated data to accurately mimic the good and bad
31 statistical properties of the likelihood functions of generative phylogenetic models.
32 (Keywords: phylogeography, SSE, phylodynamics, machine learning, deep learning,
33 epidemiology)

INTRODUCTION

34

35 Viral phylodynamic models use genomes sampled from infected individuals to trace the
36 evolutionary history of a pathogen and its spread through a population (Holmes and
37 Garnett 1994; Volz et al. 2013). By linking genetic information to epidemiological data,
38 such as the location and time of sampling, these generative models can provide valuable
39 insights into the transmission dynamics of infectious diseases, especially in the early stages
40 of cryptic disease spread when it is more difficult to detect and track (Holmes et al. 1995;
41 Rambaut et al. 2008; Lemey et al. 2009; Pybus et al. 2012; Worobey et al. 2016, 2020;
42 Lemey et al. 2021; Washington et al. 2021; Pekar et al. 2022). This information can be
43 used to inform public health interventions and improve our understanding of the evolution
44 and spread of pathogens. Many phylodynamic models are adapted from state-dependent
45 birth-death (SDBD) processes or, equivalently, state-dependent speciation-extinction (SSE)
46 models (Maddison et al. 2007; FitzJohn 2012; Kühnert et al. 2014; Beaulieu and O’Meara
47 2016). Here, we will refer to the state as location and the models as location-dependent
48 birth-death (LDBDS) models which include serial sampling (Kühnert et al. 2016).

49 Epidemiologists are increasingly using LDBDS models to estimate transmission
50 rates, migration rates between locations, and variation in these rates amongst populations
51 (Nadeau et al. 2021) and species (Lu et al. 2021). Analysts fit data to these models with
52 likelihood-based inference methods, such as maximum likelihood (Maddison et al. 2007;
53 Richter et al. 2020) or Bayesian Markov chain Monte Carlo (Kühnert et al. 2016; Scire
54 et al. 2020). Likelihood-based inference relies upon a likelihood function to evaluate the
55 relative probability (likelihood) that a given phylogenetic pattern (i.e., topology, branch
56 lengths, and tip locations) was generated by a phylodynamic process with particular model
57 parameter values. In this sense the likelihood of any possible phylodynamic data set is
58 mathematically encoded into the likelihood as a function of (unknown) data-generating
59 model parameters.

60 Computing the likelihood requires high-dimensional integration over a large and

61 complex space of evolutionary histories. Analytically integrated likelihood functions,
62 however, are not known for LDBDS models. Methods developers instead use ordinary
63 differential equation (ODE) solvers (Maddison et al. 2007; Kühnert et al. 2016) or data
64 augmentation (DA) methods (Maliot et al. 2019) to numerically approximate the
65 integrated likelihood. These clever approximations perform well statistically, but are too
66 computationally expensive to use with large epidemic-scale data sets. Thus, while
67 Nextstrain (Hadfield et al. 2018) and similar efforts provided useful visualizations to policy
68 makers during the COVID response, most phylogeographical methods are used forensically,
69 providing insight on the past, and are not used to provide parameter estimates in response
70 to emerging events to inform policy decisions in real-time due to the complexity and long
71 run-times of these models.

72 As phylodynamic models become more biologically realistic, they will necessarily
73 grow more mathematically complex, and therefore less able to yield likelihood functions
74 that can be approximated using ODE or DA methods. Because of this, phylodynamic
75 model developers tend to explore only models for which a likelihood-based inference
76 strategy is readily available. As a consequence, this impedes the design, study, and
77 application of richer phylodynamic models of disease transmission.

78 To avoid the computational limitations associated with likelihood-based methods,
79 deep learning inference methods that are likelihood-free have emerged as a complementary
80 framework for fitting a wide variety of evolutionary models (Bokma 2006). Deep learning
81 methods rely on training many-layered neural networks to extract information from data
82 patterns. These neural networks can be trained with simulated data as another way to
83 approximate the latent likelihood function (Cranmer et al. 2020). Once trained, neural
84 networks have the benefit of being fast, easy to use, and scalable. Recently, likelihood-free
85 deep learning neural network methods have successfully been applied to phylogenetics
86 (Suvorov et al. 2020; Suvorov and Schrider 2022b; Nesterenko et al. 2022; Solis-Lemus
87 et al. 2022; da Fonseca et al. 2020), and phylodynamic inference (Voznica et al. 2022;

88 Lambert et al. 2022).

89 Here we extend new methods of deep learning from phylogenetic trees (Voznica
90 et al. 2021; Lambert et al. 2022) to explore their potential when applied to phylogeographic
91 problems in geospatial epidemiology. Phylodynamics of birth-death-sampling processes
92 that include migration among locations have been under development for more than a
93 decade (Stadler 2010; Stadler et al. 2012; Kühnert et al. 2014, 2016; Scire et al. 2020; Gao
94 et al. 2021, 2022). Given the added complexity of location specific dynamics (e.g. location
95 specific birth rates) and recent successes in deep learning with phylogenetic time trees
96 (Voznica et al. 2022) under state-dependent diversification models (Lambert et al. 2022), we
97 sought to evaluate this approach when applied to viral phylodynamics and phylogeography
98 by including location data when training deep neural networks with phylogenetic trees.

99 One important limitation of likelihood-free approaches is that it is unknown how
100 brittle the inference machinery is when the assumptions used for simulation and training
101 are violated (Schmitt et al. 2022). For example, a brittle deep learning method would be
102 more easily misled by model misspecification when compared to a likelihood-based
103 method. Likelihood approaches may have some advantages because the simplifying
104 assumptions are explicit in the likelihood function while for trained neural networks it is
105 difficult to know how those assumptions are encoded for any given implementation.
106 However, with complex likelihood models, there may be unexpected interactions among
107 simplifying assumptions that result in large biases when applied to real-world data.
108 Characterizing the relative robustness and brittleness of these two inference paradigms is
109 essential for those who wish to confidently develop and deploy likelihood-free methods of
110 inference from real world data.

111 To explore relative robustness to model misspecification, we trained multiple deep
112 convolutional neural networks (CNNs) with transmission trees generated from epidemic
113 simulations. We show that simulation-trained CNNs are not only as accurate as
114 likelihood-based approaches but are no more sensitive to model violations than the

115 likelihood approach. Both methods consistently show similar biases induced by model
116 violations in test data sets. We find that for the models tested here, the migration rate
117 estimates are highly sensitive to misspecification of infection rate and sampling rates, but
118 that estimates of the infection and sampling rates are fairly robust to misspecification of
119 the migration models. We also show that the rate parameter estimates are fairly robust to
120 misspecification of both the number of locations in the model and phylogenetic error.
121 Finally, we compared a simulation-trained neural network to a recent phylodynamic study
122 of the first wave of the COVID pandemic in Europe (Nadeau et al. 2021) and obtain
123 similar inferences about the dynamics and history of SARS-CoV-2 in the European clade.

124 METHODS

125 LDBDS processes are stochastic branching processes that define location-dependent rates
126 for birth, death, migration, and sampling events to randomly generate time-scaled
127 phylogenies where taxa are associated with various locations. With serial sampling, many
128 chains of the transmission tree go undetected. Consequently, in phylogeography an absence
129 of evidence is not evidence of absence in time and space. This fact requires simulation of
130 not just the sampled/observed phylogenetic tree but the evolution of the underlying
131 population from which it is sampled. This underlying population is divided into
132 compartments of Susceptible individuals, infectious individuals, and
133 recovered/non-susceptible individuals. The dynamics of these compartments are describe
134 by the Susceptible-Infectious-Recovered (SIR) compartmental model.

135 First, we define the SIR model we assume here that is approximately equivalent to
136 the LDBDS model (Kühnert et al. 2016). Following that, is a description of the simulation
137 method to generate the training, validation, and test data sets of phylogenies under the
138 model. We next describe our implementation of simulation-trained deep learning inference
139 with convolutional neural networks (CNN) as well as a likelihood-based method using
140 Bayesian inference. We then describe our methods for measuring and comparing their

141 performance when tested against data sets generated by simulations under the inference
142 model as well as several data sets simulated under models that violate assumptions of the
143 inference model. Finally, we describe how we tested our simulation-trained CNN against a
144 real-world data set.

145 *Model definition*

146 We first define a general location-dependent SIR stochastic process used for simulations
147 and likelihood function derivation in the format of reaction equations we specified in
148 MASTER (Vaughan and Drummond 2013). Reaction equations 1 through 4 specify the
149 SIR compartment model with migration and serial sampling where S , I , and R denote the
150 number of individuals in each compartment. The S and I compartments are indexed by
151 geographic location using i and j . N_i is the total population size in location i and
152 $N_i = S_i + I_i + R_i$. The symbols for each rate parameter is placed above each reaction arrow.



153 We parameterize the model with the basic reproduction number in location i , R_{0_i} ,
154 which is related to β_i and δ_i by equation 5,

$$R_{0_i} = \frac{\beta_i}{\gamma + \delta_i}. \quad (5)$$

155 In particular, our study considers a location-independent SIR (LI-SIR) model with
156 sampling that assumes R_{0_i} was equal among all locations, and a location-dependent

157 (LD-SIR) model with sampling that assumes R_{0_i} varied among locations. During the
 158 exponential growth phase of an outbreak, the LI-SIR and LD-SIR models are equivalent to
 159 the location-independent birth-death-sampling (LIBDS) and location-dependent
 160 birth-death-sampling (LDBDS) models, respectively, that are often used in viral
 161 phylogeography (Kühnert et al. 2014, 2016; Douglas et al. 2021).

162 Each infectious individual transitions to recovered at rate γ . We assumed that
 163 sampling a virus in an individual occurs at rate δ_i in location i and immediately removes
 164 that individual from the infectious compartment and places them in the recovered
 165 compartment. Thus the effective recovery rate in location i is $\gamma + \delta_i$. The above reactions
 166 correspond to the following coupled ordinary differential equations.

$$\begin{aligned}
 \frac{dS_i}{dt} &= -\frac{\beta_i}{N_i} S_i I_i \\
 \frac{dI_i}{dt} &= \frac{\beta_i}{N_i} S_i I_i + \sum_{j \neq i}^n m_{ij} I_j - \sum_{j \neq i}^n m_{ji} I_i - (\gamma + \delta_i) I_i \\
 \frac{dR}{dt} &= \gamma \sum_{i=1}^n \delta_i I_i
 \end{aligned} \tag{6}$$

167 When the migration rate is constant among locations and the model is a
 168 location-independent SIR model, or equivalently, LIBDS, equation set 6 reduces to

$$\begin{aligned}
 \frac{dS_i}{dt} &= -\frac{\beta}{N_i} S_i I_i \\
 \frac{dI_i}{dt} &= \frac{\beta}{N_i} S_i I_i + m \left(\sum_{j \neq i}^n I_j - (n-1) I_i \right) - (\gamma + \delta) I_i \\
 \frac{dR}{dt} &= (\gamma + \delta) \sum_{i=1}^n I_i
 \end{aligned}$$

169 The number of infections and the migration of susceptible individuals is at
 170 negligible levels on the timescales investigated here. The infection rate is, therefore,

171 approximately constant and the migration of susceptible individuals can be safely ignored
172 requiring only migration of infectious individuals to be simulated.

173 At the beginning of an outbreak, it is often easier to know the recovery period from
174 clinical data than the sampling rate which requires knowing the prevalence of the disease.
175 Therefore, we treat the average recovery period as a known quantity and use it to make the
176 other two parameters (the sampling rate and the basic reproduction number R_0)
177 identifiable. This was done by fixing the corresponding rate parameter in the likelihood
178 function to the true simulated value for each tree, and by adding the true simulated value
179 to the training data for training the neural network.

180 *Simulated training and validation data sets*

181 Epidemic simulations of the SIR+migration model that approximates the LIBDS process
182 were performed using the MASTER package v. 6.1.2 (Vaughan et al. 2014) in BEAST 2 v.
183 2.6.6 (Bouckaert et al. 2019). MASTER allows users to simulate phylodynamic data sets
184 under user-specified epidemiological scenarios, for which MASTER simultaneously
185 simulates the evolution of compartment (population type) sizes and tracks the branching
186 lineages (transmission trees in the case of viruses) from which it samples over time. We
187 trained neural networks with these simulated data to learn about latent populations from
188 the shape of sampled and subsampled phylogenies. In addition to the serial sampling
189 process, at the end of the simulation 1% of infected lineages were sampled. In MASTER
190 this was approximated by setting a very high sampling rate and very short sampling time
191 such that the expected number sampled was approximately 1%. This final sampling event
192 was required to make a 1-to-1 comparison of the likelihood function used for this study (see
193 Likelihood method description below) which assumes at least one extant individual was
194 sampled to end the process. Coverage statistics from our MCMC samples closely match
195 expectations (see Likelihood method description below; SI Figure S2). Simulation

196 parameters under LIBDS and LDBDS models for training the neural network under the
197 phylogeography model were drawn from the following distributions:

$$\begin{aligned}R_0 &\sim \text{Uniform}(2, 8) \\ \delta &\sim \text{Unif}(0.0001, 0.005) \\ m &\sim \text{Uniform}(0.0001, 0.005) \\ \gamma &\sim \text{Unif}(0.01, 0.05)\end{aligned}\tag{7}$$

root location \sim Multinomial($k = 1, p_i = 1/n$), for n locations

198 All five locations had initial population sizes of 1,000,000 susceptible individuals
199 and one infected individual in one of the locations. Simulations were run for 100 time units
200 or until 50,000 individuals had been infected to restrict simulations to the approximate
201 exponential phase of the outbreak. For the experiments comparing the CNN to the
202 likelihood-based method under the LIBDS model, if this population threshold was reached
203 the simulation was rejected. This criterion was not enforced for simulations under the
204 LDBDS model. This ensured the LIBDS model used in the likelihood-based analyses are
205 equivalent to more complex density-dependent SIR models. After simulation, trees with
206 500 or more tips were uniformly and randomly downsampled to 499 tips and the sampling
207 proportion was recorded for training the neural networks and to adjust estimates of δ .

208 We simulated 410,000 outbreaks under these LIBDS settings to generate the
209 training, validation, and test sets for deep learning. Any simulation that generated a tree
210 with less than 20 tips was discarded, leaving a total of 111,157 simulated epidemiological
211 data sets. Of these, 104,157 data sets were used to train and 7,000 were used to validate
212 and test each CNN. A total of 193,110 LDBDS data sets were simulated, with 186,110 used
213 to train and 7,000 used to validate and test the LDBDS CNNs.

214 Training simulation parameters for the LDBDS process used to analyze the real

215 data set (Nadeau et al. 2021) were drawn from the same distributions as LIBDS except R_0
216 was unique for each location and was drawn from a hierarchical distribution to narrow the
217 magnitude of differences among locations within simulations to be within 8 of each other
218 but expand the magnitude of differences between simulations to range from 0.9 to 15:

$$\alpha \sim \text{Uniform}(4.9, 11)$$

$$R_{0_i} \sim \text{Uniform}(\alpha - 4, \alpha + 4)$$

219 For the empirical analysis, population sizes at each location were also set to 500,000
220 and instead of running the simulations for 100 time units, time was scaled by the recovery
221 period, $1/\gamma$, and was drawn from a uniform distribution:

$$\text{time} \sim \text{Uniform}(1, 20)$$

222 *Simulated test data sets with and without model misspecification*

223 We first simulated a test set of 138 trees under the training model to compare the
224 accuracy of the CNN and the likelihood-based estimates when the true model is specified.
225 These data sets were simulated by random draws of parameter values from the same
226 distributions described above for generating the training data set.

227 Sensitivity to model misspecification for each of the three rate parameters, R_0 , δ ,
228 and m , was tested. All sensitivity experiments used the same LIBDS model for inference
229 for both the CNN and the Likelihood-based methods. Sensitivity experiments we
230 conducted by simulating a test data set of trees that were generated by an epidemic
231 process that was more complex than or different from the LIBDS model.

232 The tree data set for the misspecified R_0 experiment consisted of simulating
233 outbreaks where each location had a unique R_0 drawn from the same distribution as above.
234 Likewise, the misspecified sampling model test set was generated by simulating outbreaks
235 where each location had a unique sampling rate, δ , drawn from the same distribution used
236 for the global sampling rate described above. For the misspecified migration model, a
237 random pair of coordinates, each drawn from a uniform(0,5) distribution in a plane, were
238 generated for the five locations, and a pairwise migration rate was computed such that
239 pairwise migration rates were symmetric and proportional to the inverse of their euclidean
240 distances and the average pairwise migration rate was equal to a random scalar which was
241 also drawn from a uniform distribution (see equations 7 above).

242 The tree set for the misspecified number of locations experiment was generated by
243 simulating outbreaks among ten locations instead of five. After simulations, six locations
244 were chosen at random and re-coded as being sampled from the same location.

245 To generate a test set where the time tree used for inference has incorrect topology
246 and branch lengths, we implemented a basic pipeline of tree inference from simulated
247 genetic data to mimic a worst case real world scenario. We simulated trees under the same
248 settings as before. Phylogenetic error was introduced in two ways: the amount of site data
249 (short sequences) and misspecification of the DNA sequence evolution inference model (*i.e.*
250 Using seq-gen V. 1.3.2 (Rambaut and Grassly 1997). We simulated the evolution of a 200
251 base-pair sequence under an HKY model with $\kappa = 2$, equal base frequencies and 4
252 discretized-gamma(2, 2) rate categories for among site rate variation. The simulated
253 alignment as well as the true tip dates (sampling times) was then used to infer test trees.
254 Test tree inference was done using iqtree v. 2.0.6 (Minh et al. 2020) assuming a
255 Jukes-Cantor model of evolution where all transition rates are equal. The inference model
256 also assumed no among-site rate variation. The number of shared branches between the
257 true transmission tree and the test tree inferred by IQ-Tree was measured using gotree v.
258 0.4.2 (Lemoine and Gascuel 2021). Polytomies were resolved using phytools (Revell 2012)

259 and a small, random number was added to each resolved branch. These trees were then
260 used for likelihood inference and CNN prediction.

261 *Deep learning inference method*

262 The resulting trees and location metadata generated by our pipeline were converted to a
263 modified cblv format (Voznica et al. 2022), which we refer to as the cblv+S (+State of
264 character, *e.g.* location) format. The cblv format uses an in-order tree traversal to
265 translate the topology and branch lengths of the tree into an $2 \times n$ matrix where n is the
266 number of tips in the tree. This representation gives each sampled tip a pair of coordinates
267 in ‘tree-traversal space’. Our cblv+S format associates geographic information
268 corresponding with each sampled taxon by appending each vector column with a one-hot
269 encoding vector of length g states to yield a $(2 + g) \times n$ cblv+S matrix. The cblv+S format
270 allows for multiple characters and/or states to be encoded, extending the single binary
271 character encoding format introduced by Lambert et al. (2022). Our study uses cblv+S to
272 encode a single character with $g = 5$ location-states. In addition to the the cblv+S data,
273 we also include a few tree summary statistics and known simulating parameters; the mean
274 branch length, the tree height and the recovery rate and the subsampling proportion. Trees
275 were rescaled such that their mean branch length was the default for phylodeep (Voznica
276 et al. 2021) before training and testing of the CNN. The mean pre-scaling branch length
277 and tree heights were also fed into the neural networks. Trees were not rescaled for the
278 likelihood-based analysis. Note that tree height did not vary for the LIBDS CNN training
279 set but did for the LDBDS training set.

280 Our CNNs were implemented in Python 3.8.10 using keras v. 2.6.0 and
281 tensorflow-gpu v. 2.6.0. (Chollet; Abadi et al. 2016). For each model, LIBDS and LDBDS,
282 we designed and trained two CNN architectures, one to predict epidemiological rate
283 parameters and the other to predict the outbreak location resulting in four total CNNs
284 trained by two training data sets (LIBDS and LDBDS). We used the mean-squared-error

285 for the regression neural loss function in the network trained to estimate epidemiological
286 rates, and the categorical cross-entropy loss function for the categorical network trained to
287 estimate outbreak location. We assessed the performance of the network by randomly
288 selecting 5,000 samples for validation before each round of training. We measured the
289 mean absolute error and accuracy using the validation sets. We used these measures to
290 compare architectures and determine early stopping times to avoid overfitting the model to
291 the training data. We also added more simulations to the training set until we could no
292 longer detect an improvement in error statistics. After comparing the performance of
293 several networks, we found that the CNN described in SI Figure S1 performed the best. In
294 brief, the networks have three parallel sets of sequential convolutional layers for the cblv+S
295 tensor and a parallel dense layer for the priors and tree statistics. The three sets of
296 convolution layers differed by dilation rate and stride lengths. These three segments and
297 the dense layer were concatenated and then fed into a segment consisting of a sequential
298 set of dense layers, each layer gradually narrowing to the output size to either three or five
299 for the rates and origin location networks, respectively, for the LIBDS model, and seven
300 and five for the seven rates and five locations, respectively, for the LDBDS model.

301 All layers of the CNN used rectified linear unit (ReLU) activation functions. We
302 used the Adam optimizer algorithm for batch stochastic gradient descent (Kingma and Ba
303 2017) with batch size of 128 and stopping after 15 epochs for the regression network and
304 ten epochs for the root location network. The output activation for the rates network was
305 linear with three nodes and for the outbreak location network was softmax with five nodes.
306 Otherwise the architecture was the same for all four networks. The LDBDS neural network
307 was trained with simulated trees where R_{0_i} varied among locations had output layer with
308 seven nodes; five for the each location's R_{0_i} and a node each for the sampling rate and the
309 migration rate. We tested networks with max-pooling layers between convolution layers as
310 well as dropout at several rates and found no improvement or a decrease in performance.

311

Likelihood-based method of inference

312 We compared the performance of our trained phylodynamic CNN to likelihood-based
313 Bayesian phylodynamic inferences. We specified LIBDS and LDBDS Bayesian models that
314 were identical to the LIBDS and LDBDS simulation models that we used to train our
315 CNNs. The most general phylodynamic model in the birth-death family applied to
316 epidemiological data is the state-dependent birth-death-sampling process (SDBDS;
317 (Kühnert et al. 2016; Scire et al. 2020)), where the state or type on which birth, death, and
318 sampling parameters are dependent is the location in this context. The basic model used
319 for experiments here is a phylogeographic model that is similar to the serially sampled
320 birth-death process (Stadler 2010) where rates do not depend on location, which we refer
321 to as the LDBDS model. The death rate, μ , is equivalent to the recovery rate, γ , in SIR
322 models. Standard phylogenetic birth-death models assume the birth and death rates, λ and
323 μ , are constant or time-homogeneous, while the SIR model's infection rate is proportional
324 to β and S and varies with time as S changes. However, when the number of infected is
325 small relative to susceptible people, as in the initial stages of an outbreak, the infection
326 rate, β , is approximately constant and approximately equal to the birth rate λ ;

$$\lambda = \frac{\beta S}{N} \approx \beta \quad (8)$$

327 The joint prior distribution was set to the same model parameter distributions that
328 were used to simulate the training and test sets of phylogenetic trees in the first section
329 with γ treated as known and the proportion of extant lineages sampled, ρ , set to 0.01 as in
330 the simulations. The likelihood was conditioned on the tree having extant samples (*i.e.* the
331 simulation ran for the allotted time without being rejected). All simulated trees in this
332 study had a stem branch and the outbreak origins were inferred for the parent node of the
333 stem branch.

334 We used Markov chain Monte Carlo (MCMC) to simulate random sampling from
335 the posterior distribution implemented in the TensorPhylo package
336 (<https://bitbucket.org/mrmay/tensorphylo/src/master/>) in RevBayes (Höhna et al. 2016).
337 After a burnin phase, a single chain was run for 7,500 cycles with 4 proposals per cycle and
338 at least 100 effective sample size (ESS) for all parameters. If the effective sample size (ESS)
339 was less than 100, the MCMC was rerun with a higher number of cycles. We also analyzed
340 the coverage of the 5, 10, 25, 50, 75, 90, and 95% highest posterior density (HPD) intervals
341 to verify that our simulation model and inference model are the same and that the MCMC
342 simulated draws from the true posterior distribution. Bayesian phylogeographic analysis
343 recovered the true simulating parameters (SI Figure S2) at the expected frequencies, thus
344 validating the simulations were working as expected and confirming that the MCMC was
345 accurately simulating draws from the true posterior distribution.

346 *Quantifying errors and error differences*

347 We measure the absolute percent error (APE) of the predictions from the CNN and the
348 mean posterior estimate (MPE) of the likelihood-based method. The formula for APE of a
349 prediction/estimate, y^{estimate} , of y^{truth} is

$$\text{APE} = \left| \frac{y^{\text{estimate}} - y^{\text{truth}}}{y^{\text{truth}}} \right| \times 100$$

350 The Bayesian alternative to significance testing is to analyze the posterior
351 distribution of parameter value differences between groups. In this framework, the
352 probability that a difference is greater than zero can be easily interpreted. We therefore
353 used Bayesian statistics to infer the median difference in error between the CNN and
354 likelihood-based methods and the increase in median error of each method when analyzing
355 misspecified data compared to when analyzing data simulated under the true inference
356 model.

357 We used Bayesian inference to quantify population error by performing three sets of
358 analyses: (1) inferred the population median APE under the true model (this will be the
359 reference group for analysis 3), (2) the effect of inference method — CNN or
360 likelihood-based (Bayesian) — on error by inferring the median difference between the
361 CNN estimate and the likelihood-based estimate, (3) the effect of misspecification on error
362 for each parameter by comparing the median error of estimates under misspecified
363 experiments and the reference group defined by analysis 1. See SI Figures S3 - S8 and SI
364 Table S1 for summaries and figures for all analyses for this section.

365 To infer these differences between groups we used the R package BEST (Meredith
366 and Kruschke). BEST assumes the data follow a t-distribution parameterized by a location
367 parameter, μ , a scale parameter, σ , and a shape parameter, ν , which they call the
368 "normality parameter" (*i.e.* if ν is large the distribution is more Normal). Because the
369 posterior distribution does not have a closed form, BEST uses Gibbs sampling to simulate
370 draws from the posterior distribution. 20,000 samples were drawn from the posterior
371 distribution for each BEST analysis. BEST uses automatic posterior predictive checks to
372 indicate that a model adequately describes the data distributions. Posterior predictive
373 checks indicate the BEST model adequately fits each data set analyzed below.

Inferring the median APE.— Before inferring differences between groups, we inferred the population median APE for predictions of R_0 , δ , and m from test data simulated under the inference model using the CNN and likelihood-based methods. Histograms of the sampled log-transformed APE appears to be symmetric with heavy tails so we fit the log APE to the BEST model. This implies that the sampled APE scores are drawn from a log-t distribution. The log-t distribution has a mean of ∞ and median of e^μ , we therefore focus our inference on estimating posterior intervals for the population median APE from the sampled APE values for each parameter estimated by the CNN method and likelihood-based method which we denote APE^{CNN} , and APE^{Like} respectively. The data

analyzed here and likelihood assumed by BEST is

$$y = \text{APE}^{\text{CNN}} \text{ or } \text{APE}^{\text{Like}}$$
$$\log y \mid \mu, \sigma, \nu \sim t_\nu(\mu, \sigma).$$

The priors were set to the vague priors that BEST provides by default,

$$\mu \sim \text{Normal}(\text{mean}(y), \text{sd}(y) \times 1000)$$
$$\sigma \sim \text{Uniform}(\text{sd}(y)/1000, \text{sd}(y) \times 1000)$$
$$\nu \sim \text{Exponential}(1/29) + 1.$$

95% highest posterior intervals (HPI) for the median APE, $\tilde{\mu}$, was estimated by the following transformation of simulated draws from the posterior distribution

$$\tilde{\mu} = e^\mu.$$

374 In summary, the results we present are 95% HPI from the posterior distributions of the
375 median error, $\tilde{\mu}$.

Inferring the relative accuracy of the CNN and likelihood-based method.— To quantify the difference in error between the CNN and the likelihood-based method, we fit the difference in sampled APE scores, ΔAPE , between the CNN method and the likelihood-based method to the BEST model. Histograms of ΔAPE appear symmetric with weak to strong outliers making the BEST model a good candidate for inference from this data. The data and likelihood are

$$\Delta y = \text{APE}^{\text{CNN}} - \text{APE}^{\text{Like}}$$
$$\Delta y \mid \mu, \sigma, \nu \sim t_\nu(\mu, \sigma)$$

376 We used the same default priors as above.

377 Because, Δy is not log-transformed, it is drawn from a t-distribution and the
378 marginal posterior of the parameter μ is an estimate of the population mean, μ^d . Because
379 the mean and the median are equivalent for a t-distribution, we again report the posterior
380 distribution of the median difference, $\tilde{\mu}^d$ to simplify the results.

381 In summary, the results we present are 95% HPI from the posterior distribution of
382 the median difference between the two methods, $\tilde{\mu}^d$.

383 When comparing CNN to the likelihood-based approach, positive values for $\tilde{\mu}^d$
384 indicate the CNN is less accurate, and negative indicate the likelihood-based estimates less
385 accurate. We emphasise that this quantity is the median difference in contrast to the
386 difference in medians, $\Delta\tilde{\mu}$, reported in the next section.

387 *Inferring sensitivity to model misspecification.*— Finally, to quantify the overall sensitivity
388 of each rate parameter to model misspecification under each inference method, we infer the
389 difference in median APE, $\tilde{\mu}$ of predictions under a misspecified model relative to
390 predictions under the true model. In other words we are inferring differences in medians
391 between experiments. For example, to infer the sensitivity of the CNN’s inference of the
392 sampling rate, δ , to phylogenetic error, we inferred the difference between the median APE
393 of the CNN’s predictions for misspecified trees and the median APE of CNN predictions
394 for true trees. The data is concatenated as below.

$$(y_1, y_2) = (\text{APE}^{\text{CNN}}, \text{APE}^{\text{CNN Ref}}) \text{ or}$$

$$(y_1, y_2) = (\text{APE}^{\text{Like}}, \text{APE}^{\text{Like Ref}})$$

395 We inferred the difference between group median APE scores, denoted $\Delta\tilde{\mu}$, by
396 assuming that the model parameters conditioned on the observed APE from the two
397 groups, y_1 and y_2 , follow a posterior distribution that is proportional to

$$P(y_1 | \mu_1, \sigma_1, \nu)P(y_2 | \mu_2, \sigma_2, \nu)P(\mu_1, \mu_2, \sigma_1, \sigma_2, \nu),$$

398 where $\log y_1$ and $\log y_2$ follow t distributions with means μ_1 and μ_2 and standard
399 deviations σ_1 and σ_2 , respectively while sharing a common normality parameter, ν .

400 The posterior sample of $\Delta\tilde{\mu}$ is obtained by transforming samples from the joint
401 marginal posterior distribution of μ_1 and μ_2 with the following equation,

$$\Delta\tilde{\mu} = e^{\mu_1} - e^{\mu_2}.$$

The two components of the likelihood are each t-distributed and share the ν parameter which means we assume both samples are drawn from a similarly shaped distribution (similarly heavy tails).

$$\log y_1 | \mu_1, \sigma_1, \nu \sim t_\nu(\mu_1, \sigma_1)$$

$$\log y_2 | \mu_2, \sigma_2, \nu \sim t_\nu(\mu_2, \sigma_2)$$

The prior distribution for the parameters of the model were set to the defaults for BEST,

$$\mu_1 \sim \text{Normal}(\text{mean}(\log y_1), \text{sd}(\log y_1) \times 1000)$$

$$\mu_2 \sim \text{Normal}(\text{mean}(\log y_2), \text{sd}(\log y_2) \times 1000)$$

$$\sigma_1 \sim \text{Uniform}(\text{sd}(\log y_1)/1000, \text{sd}(\log y_1) \times 1000)$$

$$\sigma_2 \sim \text{Uniform}(\text{sd}(\log y_2)/1000, \text{sd}(\log y_2) \times 1000)$$

$$\nu \sim \text{Exponential}(1/29) + 1$$

402 As before, interpretation of the posterior distribution of the difference in medians is
403 straightforward: the more positive the difference in median APE from the misspecified
404 model test set and the median APE from the true model test set, the more sensitive the
405 parameter is to model misspecification in the experiment.

406 *Real Data*

407 We compared the inferences of a LDBDS simulation trained neural network to that of a
408 phylodynamic study of the first COVID wave in Europe (Nadeau et al. 2021). These
409 authors analyzed a phylogenetic tree of viruses sampled in Europe and Hubei, China using
410 a location-dependent birth-death-sampling model in a Bayesian framework using priors
411 informed by myriad other sources of information. We simulated a new training set of trees
412 under an LDBDS model where R_{0_i} depends on the geographic location, and the sampling
413 process only consists of serial sampling and no sampling of extant infected individuals. We
414 then analyzed the whole tree from Fig. 1 in (Nadeau et al. 2021) as well as the European
415 clade which Nadeau et al. (2021) labeled as A2 in the same figure. We note that our
416 simulating model is not identical to the inference model used in (Nadeau et al. 2021). We
417 model migration with a single parameter with symmetrical migration rates among
418 locations and all locations having the same sample rate. Nadeau and colleagues
419 parameterize the migration process with asymmetric pairwise migration rates and assume
420 location-specific sampling rates. We also do not include the information the authors used
421 to inform their priors as that requires an extra level of simulation and training on top of
422 simulations done here, and is thus beyond the scope of this study.

423 The time tree from (Nadeau et al. 2021) was downloaded from GitHub
424 (<https://github.com/SarahNadeau/cov-europe-bdmm>). The recovery rate assumed in
425 (Nadeau et al. 2021) was 0.1 days^{-1} which was set to 0.05 to bring the recovery rate to
426 within the range of simulating values used to train the CNN. Consequently, the branch
427 lengths of the tree were then scaled by 2. The number of tips, tree height, and average

428 branch lengths were measured from the rescaled trees and fed into the network. The full
429 tree and A2 clade were then analyzed using the LDBD CNN and compared to the posterior
430 distributions from (Nadeau et al. 2021).

431 *Hardware used*

432 Simulations were run on a 16 core Intel(R) Xeon(R) Platinum 8175M CPU @ 2.50GHz.
433 For each simulation, an XML file with random parameter settings was generated using
434 custom scripts. These XML files were the inputs for MASTER which was run in the
435 BEAST2 platform. Neural network training and testing and predictions were conducted on
436 an 8 core Intel(R) Core(TM) i7-7820HQ CPU @ 2.90GHz laptop.

437 *Data and code availability.*— A repository containing data and code used in this study is
438 available here: Link to be provided soon.

439 RESULTS

440 *Comparing deep learning to likelihood*

441 Our first goal in this study was to train a CNN that produced phylodynamic parameter
442 point estimates that were as accurate as likelihood-based Bayesian posterior mean
443 estimates under the true model. This will serve as a reference for quantifying level of
444 sensitivity in our misspecification experiments. We focused on estimating important
445 epidemiological parameters – the reproduction number, R_0 , the sampling rate, δ , and the
446 migration rate, m – as well as the outbreak origin from viral phylogenies like those typically
447 estimated from serially sampled DNA sequences that were obtained as the virus spread.

448 Our CNN produced estimates that are as accurate as the mean posterior estimates
449 (MPE) under the true simulating model. We compared the absolute percent error (APE)

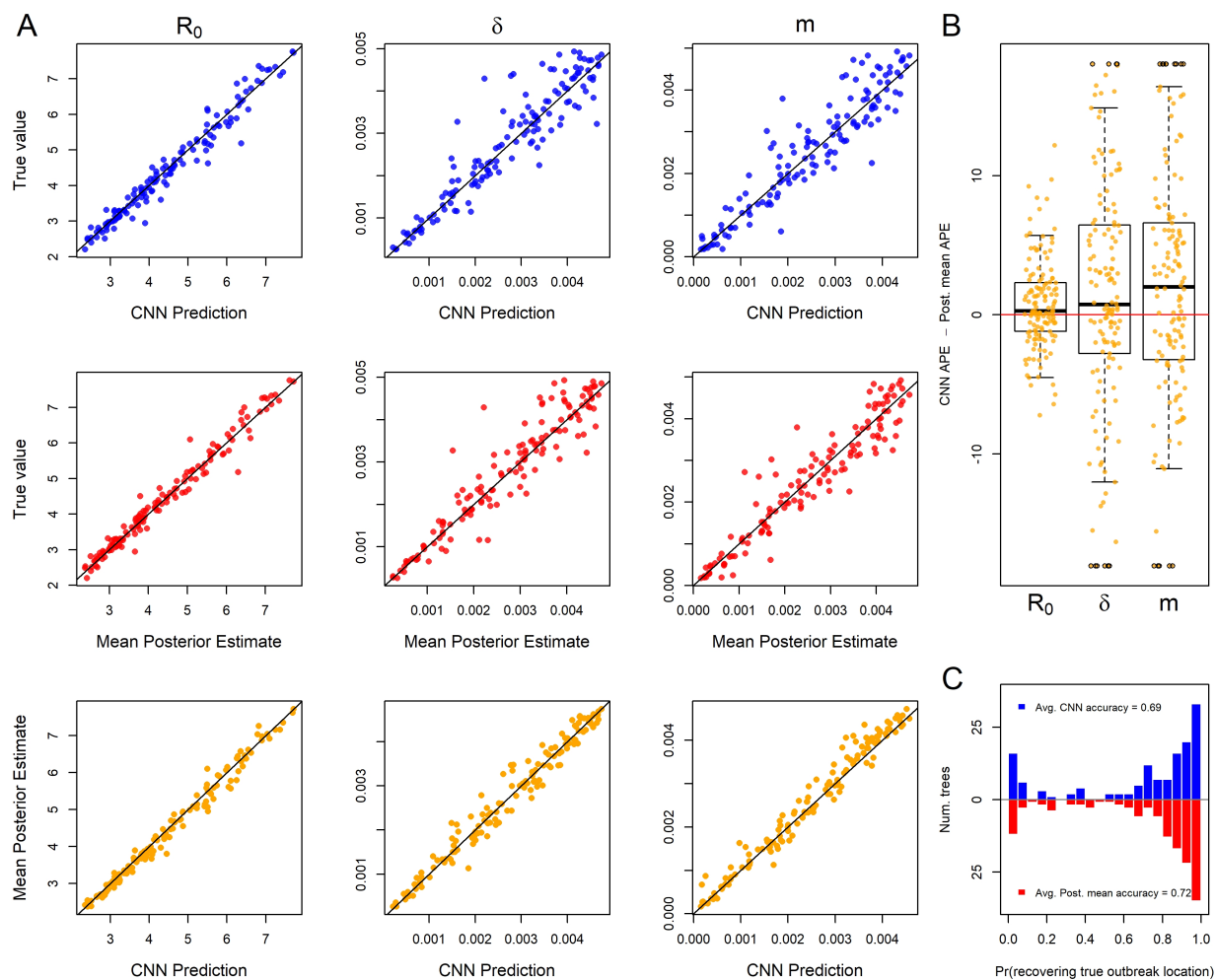


Figure 1: Inference under the true simulating model. (A) Scatterplot of CNN predictions and posterior mean estimates from Bayesian analyses against the true values (top two rows in blue and red respectively) of the basic reproduction number, R_0 , the sampling rate, δ , and the migration rate, m for 138 test trees. The bottom row in orange shows scatter plots of the CNN estimates against the posterior mean estimates for the same trees. (B) The difference in the absolute percent error (APE) of estimates for the two inference methods. Boxes show the inner 50% quantile of the data while whiskers extend 1.5 IQR. Dots with black circles show estimates that were truncated to the mean of the parameter with the most extreme outliers for visualization purposes. (C) Histograms of the probabilities of inferring the correct outbreak origin location for the same trees as in panels A and B.

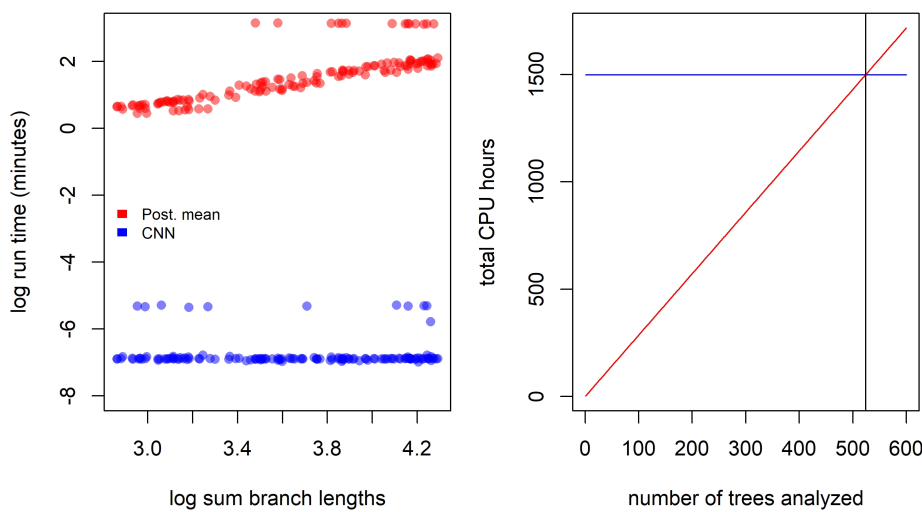


Figure 2: Left: Estimates of time to complete analysis of each of 138 trees relative to tree size. Right: The number of trees (524; gray vertical line) needed to analyze for total analysis time of Bayesian method (red line) to equal that of the entire simulation and CNN training and inference pipeline (blue line).

450 of the network predictions to the APE of the MPE of the Bayesian LIBDS model (Figure
451 1). The APE is straight-forward to interpret, e.g. an APE of < 10 means the estimate is
452 within 10 percentage points (ppts) of the true value. For the three epidemiological rate
453 parameters, R_0 , δ and m , both methods made very similar predictions for the 100 time tree
454 test set (Figure 1 panel A). The two methods appear to produce estimates that are more
455 similar to each other than to the ground truth labels (compare bottom row scatter plots in
456 orange to the blue and red scatter plots in panel A). Fig. 1 panel B shows that the inferred
457 median difference in APE, $\tilde{\mu}^d$, between the method's estimates for the three parameters is
458 close to zero ($|\tilde{\mu}^d|$ 95% highest posterior interval (HPI) is < 4 ppts; SI Table S1; SI Figure
459 S3). Fig. 1 Panel C shows that our predictions of the location of outbreak have similar
460 patterns of accuracy as those from the Bayesian method.

461 Our trained CNN provides nearly instantaneous estimates of model parameters.
462 While the run time of the likelihood approach employed in this study scales linearly with
463 the size of the tree, the neural network has virtually constant run times that are more than

464 three orders of magnitude faster. Because simulation-trained neural networks have a
465 one-time cost of simulating the training data set and then training the neural network,
466 these methods are often called amortized-approximators (Bürkner et al. 2022). This means
467 the time savings aren't recouped until a certain number of trees have been analyzed. For
468 example, here over 524 trees would need to be analyzed to realize the cost savings of
469 simulating data and training our neural network (Figure 2). This illustrates the importance
470 of simulation optimization and generality for likelihood-free approaches to inference.

471 *Comparing robustness to model misspecification*

472 To test the relative sensitivity of CNN estimates and the likelihood-based MPE to model
473 misspecification, we simulated several test data sets under different, more complex
474 epidemic scenarios and compared the decrease in accuracy (increase in APE).

475 Our first model misspecification experiment tested performance when assuming all
476 locations had the same R_0 when, in fact, each location had different R_{0_i} values. The
477 median APE for all three parameters increased to varying degrees (SI Fig. S4 Panel A)
478 compared to the median APE measured in Fig. S3. We found that both methods
479 converged on similar biased estimates for R_0 . In both the CNN and Bayesian method,
480 estimates of δ were relatively robust to misspecifying R_0 . In contrast, the migration rate
481 showed much more sensitivity to this model violation in both methods with both methods
482 also converging on similarly biased estimates (Figure 3 A). The median difference in error
483 between the two methods is close to zero for all rate parameters ($|\tilde{\mu}^d|$ 95% HPI < 6 ppts;
484 SI Table S1) (SI Figure S4 Panel B). The CNN appears to be slightly more sensitive than
485 the Bayesian approach when predicting the outbreak location. Nevertheless, their
486 distributions are quite similar (Fig. 3 Panel C).

487 Next, we measured method sensitivity when the sampling process of the test trees
488 violates assumptions in the inference model. In this set, each location had a unique and
489 independent sampling rate, δ , rather than a single δ shared among locations. The median

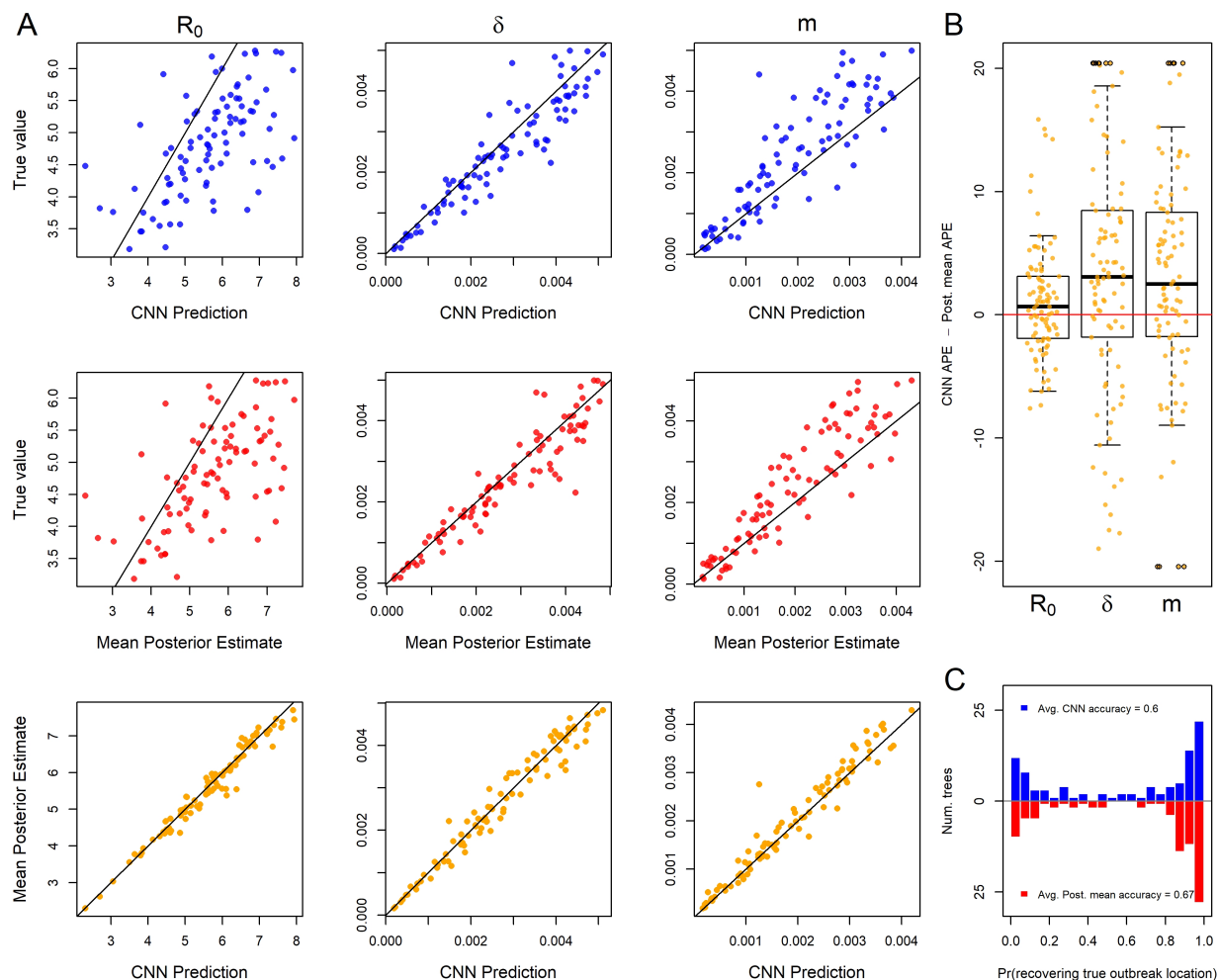


Figure 3: For 93 test trees where the R_0 parameter was misspecified: the simulating model for the test data specified 5 unique R_0 s among the five locations while the inference methods assumed one R_0 shared among locations. Because of this, the estimates for R_0 are plotted against mean of the five true R_0 values. See Figure 1 for general details about plots.

490 APE only increased for δ and m (SI Figure S5 Panel A). As expected, estimates of δ were
 491 highly biased for both methods (Figure 4 panel A). Panel A also shows that R_0 is virtually
 492 insensitive to sampling model misspecification, but that migration rate, again, is highly
 493 sensitive in both the CNN and likelihood method. The median difference in error between
 494 the two methods is close to zero for all the rate parameters ($|\tilde{\mu}^d|$ 95% HPI < 5 ppts; SI
 495 Table S1, SI Figure S5) (Figure 4 panel B). The location of outbreak prediction is also
 496 somewhat sensitive in both methods, with the CNN showing a slightly larger mean
 497 difference, but the overall distribution of accuracy of all the test trees again is similar

498 (Figure 4 panel C).

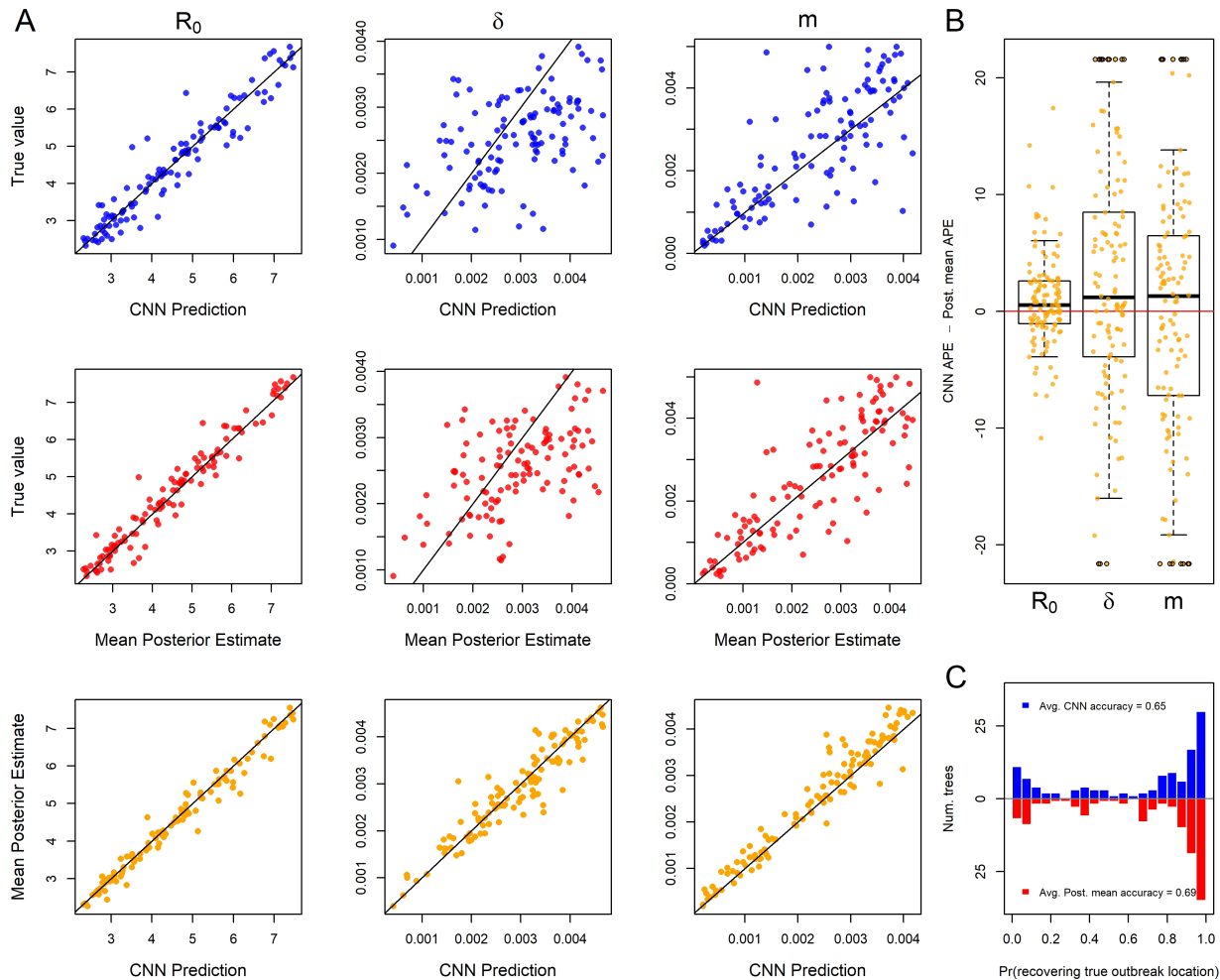


Figure 4: For 118 test trees where the sampling rate parameter was misspecified: the simulating model for the test data specified 5 unique sampling rates among the five locations while the inference methods assumed one sampling rate shared among locations. The estimates of δ are plotted against the mean true values of δ . See Figure 1 for general details about plots.

499 To explore sensitivity to migration model underspecification, we simulated a test set
 500 where the migration rates between locations is free to vary rather than being the same
 501 among locations as in the inference model. This implies $5!$ unique location-pairs and thus
 502 unique migration rates in the test data set. Results show that for both methods the
 503 parameters R_0 and δ are highly robust to this simplification (SI Fig. S6 Panel A). Though
 504 estimates of a single migration rate had a high degree of error compared to a single pair of

505 locations migration rates (Figure 5 panel A), the two methods still had similar estimates
 506 with the difference in APE centered near zero (Figure 5 panel B). The inferred median
 507 difference in APE was close to zero ($|\tilde{\mu}^d|$ 95% HPI < 3 ppts; SI Table S1; SI Figure S6
 508 Panel B). There was a slight but similar decrease in accuracy in predicting the outbreak
 509 location for both methods (Figure 5 panel C).

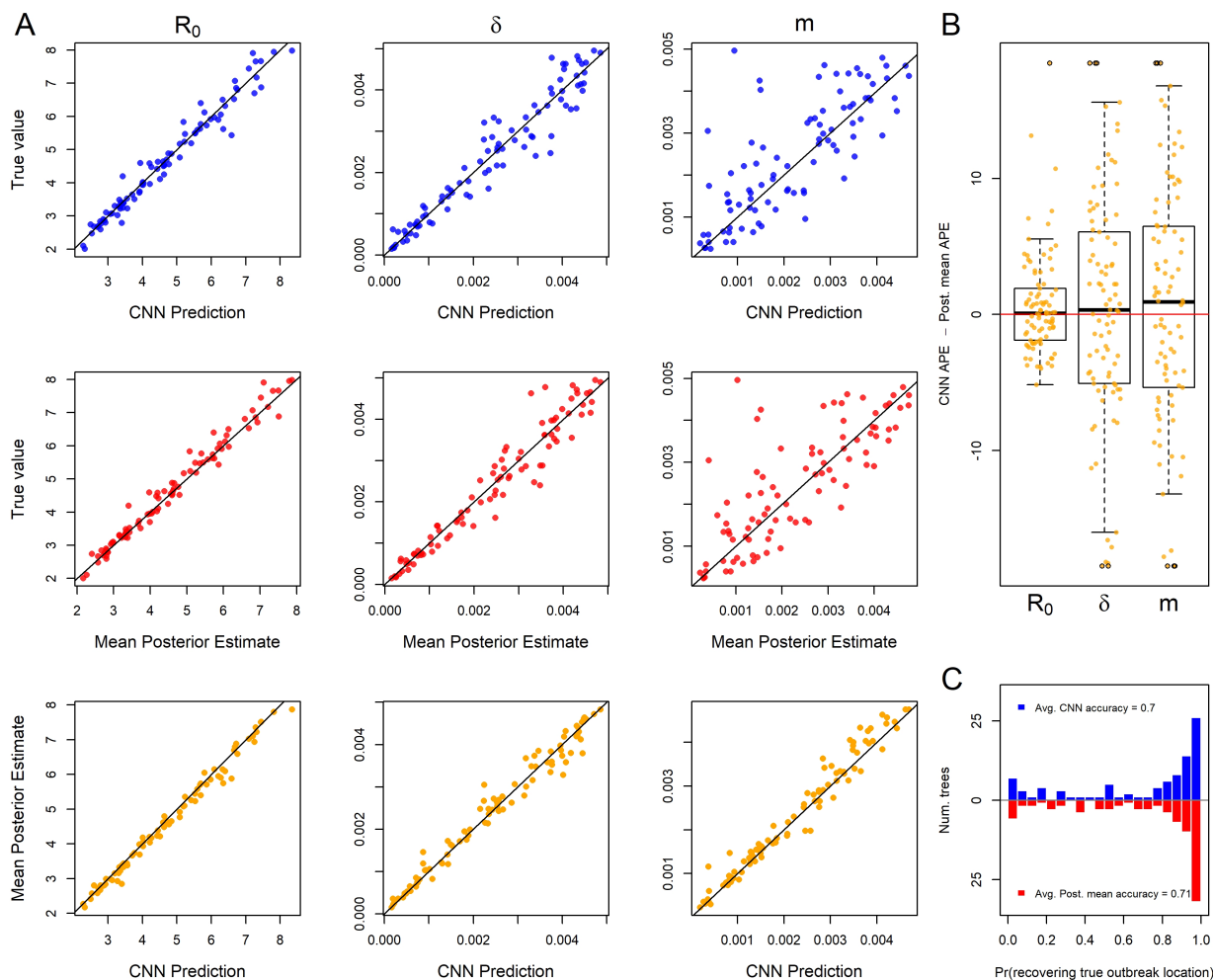


Figure 5: For 90 test trees where the migration rate parameter was misspecified: the simulating model for the test data specified 5! (120) unique migration rates among the unique pairs of the five locations while the inference methods assumed all migration rates were equal. The inferred migration rate is plotted against the mean pairwise migration rates of test data set. See Figure 1 for general details about plots.

510 When testing the sensitivity of the two methods to arbitrary groupings of locations,
 511 we found that both methods showed equal sensitivity to the same parameters (Fig. 6

512 Panels A and B). In particular, the migration rate showed a modest increase in median
513 APE and R_0 and sample rate showed virtually no sensitivity to arbitrary grouping of
514 locations (SI Figure S7 Panel A). The inferred median difference between method APE's
515 was again close to zero ($|\tilde{\mu}^d|$ 95% HPI < 4 ppts; SI Table S1; SI Figure S7 Panel B). This
516 suggests that for at least the exponential phase of outbreaks where rate parameters do not
517 vary among locations, these models have a fair amount of robustness to the decisions
518 leading to geographical division of continuous space into discrete space. The outbreak
519 location showed higher accuracy in both methods due to the fact that the test data was no
520 longer a flat distribution; the 6 combined locations should contain 60% of the outbreak
521 locations (Figure 6 panel C).

522 Finally, we explored the relative sensitivity of our CNN to amounts of phylogenetic
523 error that are present in typical phylogeographic analyses. Our simulated phylogenetic
524 error produced trees with average normalized Robinson-Foulds distances (Robinson and
525 Foulds 1981) between the inferred tree and the true tree of about 0.5 with 95% of
526 simulated trees having distances within 0.36 and 0.72. We again compared inferences
527 derived from the true tree and the tree with errors using the CNN and the Bayesian LIBDS
528 methods. Results show that migration rate was minimally affected but R_0 and δ were to a
529 some degree sensitive to phylogenetic error (Figure 7 panel A; SI Figure S8 Panel A), with
530 both methods again showing similar degrees of sensitivity (Figure 7 panel B). The inferred
531 median difference was, yet again, small ($|\tilde{\mu}^d|$ 95% HPI < 6 ppts. SI Table S1, SI Figure S8
532 Panel B). Inference of the origin location, were very similar for both methods (Fig. 7 Panel
533 C).

534 *Analysis of SARS CoV-2 tree*

535 We next compared our likelihood-free method to a recent study investigating the
536 phylodynamics of the first wave of the SARS CoV-2 pandemic in Europe (Nadeau et al.
537 2021). Despite simulating the migration and the sampling processes differently from

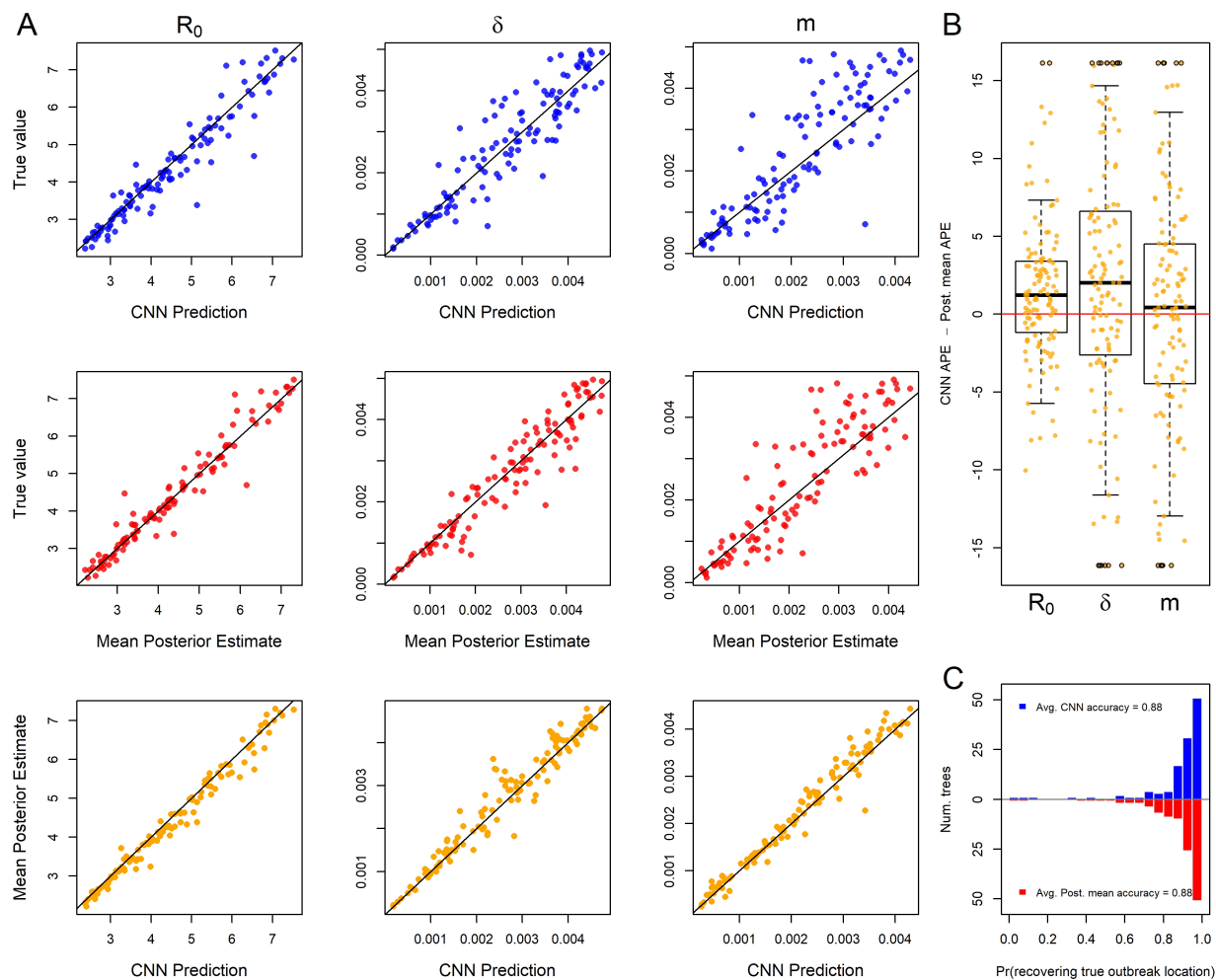


Figure 6: For 101 test trees where the number of locations was misspecified: the simulating model for the test data specified an outbreak among 10 locations with 6 locations subsequently combined into a single location while the inference methods assumed 5 locations with no arbitrary combining of locations. See Figure 1 for general details about plots.

538 Nadeau et al. (2021), our CNN produces similar estimates for the location-specific R_0 and
 539 the origin of the A2 clade (Figure 8). Whether the full tree or just the A2 clade is fed into
 540 the network, the predicted R_0 for each location was not far from the posterior estimates of
 541 Nadeau et al. (2021). The only significant discrepancy in the origin prediction is that their
 542 analysis suggests a much higher probability that the most recent common ancestor of the
 543 A2 clade was in Hubei than our CNN predicts. This is likely because our CNN only used
 544 the A2 clade to predict A2 origins which has no Hubei samples to infer the origin of the A2
 545 clade while Nadeau et al. (2021) used the whole tree. Notwithstanding this difference,

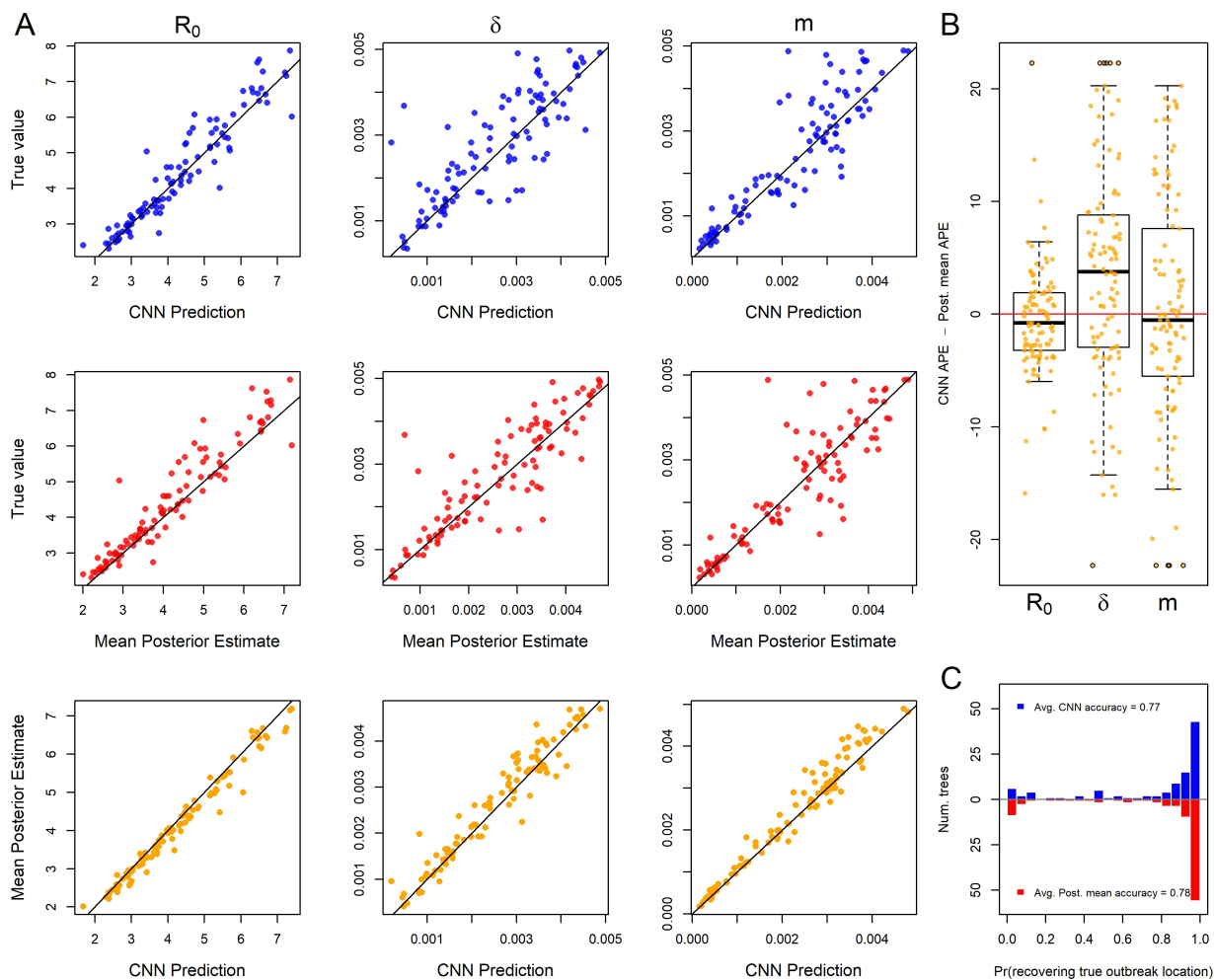


Figure 7: For 118 test trees where the time tree was misspecified: the true tree from the simulated test set was replaced with an inferred tree from simulated DNA alignments under the true tree. See Figure 1 for general details about plots.

546 among European locations, both methods predict Germany is the most likely location of
 547 the most recent common ancestor followed by Italy.

548 DISCUSSION AND CONCLUSIONS

549 Inference models are necessarily a simplified approximation of the real world. Both
 550 simulation-trained neural networks and likelihood-based inference approaches suffer from
 551 model under-specification and/or misspecification. When comparing inference methods it is
 552 important to assess the sensitivity of model inference to simplifying assumptions. In this

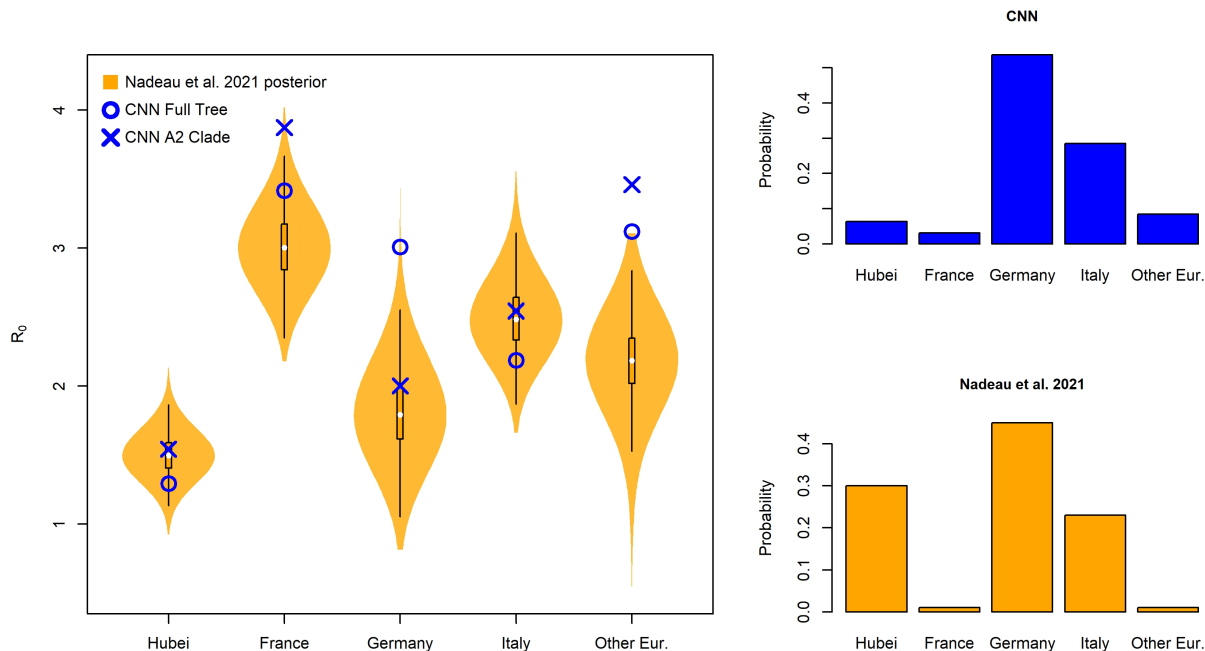


Figure 8: Location-dependent birth-death-sampling model (LDBDS) CNN comparison to (Nadeau et al. 2021) inference. Left violin plots show the posterior distributions of R_0 for each location in Europe as well as Hubei, China (orange). The blue X and O marks the MTBD CNN prediction from analyzing the full tree and the A2 (European) clade respectively. Right barplots show the LDBDS CNN prediction (blue) and posterior inference (orange) from (Nadeau et al. 2021) of the ancestral location of the A2 (European) clade (see Figure 1 (Nadeau et al. 2021)).

553 study we show that newer deep learning approaches and standard Bayesian approaches
554 behave and misbehave in similar ways under a panel of phylodynamic estimation tasks
555 where the inference model is correct as well as when it is misspecified.

556 By extending new approaches to encode phylogenetic trees in a compact data
557 structure (Voznica et al. 2021; Lambert et al. 2022), we have developed the first application
558 of phylodynamic deep learning applied to phylogeography with serial sampling. Our
559 approach is similar to that of (Lambert et al. 2022) in which they analyzed a binary SSE
560 model with exclusively extant sampling. By training a neural network on phylogenetic trees
561 generated by simulated epidemics, we were able to accurately estimate key epidemiological
562 parameters, such as the reproduction number and migration rate, in a fraction of the time

563 it would take with likelihood-based methods. Like Voznica et al. (2021) and Lambert et al.
564 (2022), we found that CNN estimators perform as well or nearly as well as likelihood-based
565 estimators under conditions where the inference model is correctly specified to match the
566 simulation model. The success of these separate applications of deep learning to different
567 phylodynamic problems is a testament to the versatility of the cblv encoding of trees.

568 We compared the sensitivity of deep learning and likelihood-based inference to
569 model misspecification. Because deep-learning methods of phylogenetic and phylodynamic
570 inference are new, few studies compare how simulation-trained deep learning methods fail
571 in comparison to likelihood methods in this way (Flagel et al. 2019). We assume that when
572 the inference model is correctly specified to match the simulation model, the trained CNN
573 will, at best, produce noisy approximations of likelihood-based parameter estimates. In
574 reality, issues related to training data set size, learning efficiency, and network overfitting
575 may cause our CNN-based estimates to contain excess variance or bias when compared to
576 Bayesian likelihood-based estimators. Our results from five model misspecification
577 experiments show that both methods of inference perform similarly when the simulating
578 model and the inference model assumptions do not perfectly match. These similarities
579 exist not only in aggregate, when comparing method performance across datasets, but also
580 when comparing performance for each individual dataset. This suggests that the CNN and
581 likelihood methods are truly estimating parameters using isomorphic criteria, despite the
582 fact that CNN heuristically learns these criteria through data patterns, while likelihood
583 precisely and mathematically defines these criteria through the model definition itself.

584 Results of comparative sensitivity experiments like this are important because if
585 likelihood-free methods using deep neural networks can easily be trained to yield estimates
586 that are as robust to model misspecification as likelihood-based methods, then analysis of a
587 large space of more complex outbreak scenarios for which tractable likelihood functions are
588 not available can be developed and applied to real world data. Additionally, sufficiently
589 realistic, pre-trained neural networks can yield nearly instantaneous inferences from data in

590 real time to inform analysts and policy makers.

591 We also tested location-dependent SIR simulation trained neural network against a
592 previous publication fitting a similar model – location-dependent birth-death-sampling
593 (LDBDS) model – on real-world data using a Bayesian method. Our CNN predicted
594 location-specific R_{0_i} and outbreak origin in Europe were similar to that inferred in (Nadeau
595 et al. 2021). This result and our model misspecification experiments suggest that
596 simulation-trained deep neural networks trained on phylogenetic trees can find patterns in
597 the training data that generalize well beyond the training data set.

598 Our study extends the results of Voznica et al. (2022) and Lambert et al. (2022) in
599 several important ways. This work showed that the new compact bijective ladderized
600 vector encoding of phylogenetic trees can easily be extended with one-hot encoding to
601 include metadata about viral samples. We extended it to include location data and were
602 able to train a neural network to not only predict important epidemiological parameters
603 such as R_{0_i} and the sampling rate, but also geographic parameters such as the migration
604 rate and the location of outbreak origination or spillover. We anticipate that much more
605 metadata can be added to train neural networks to bring more diverse and complex data to
606 make predictions about many important aspects of epidemiological spread such as the
607 relative roles of different demographic groups and the overlap of different species' ranges.

608 This approach can be readily applied to numerous compartment models used to
609 describe the spread of different pathogens among different species, locations, and
610 demographic groups, e.g. SEIR, SIRS, SIS, etc. (Ponciano and Capistrán 2011; Volz and
611 Siveroni 2018; Bjørnstad et al. 2020; Chang et al. 2020; O’Dea and Drake 2021) as well as
612 modeling super-spreader dynamics as in (Voznica et al. 2021). With fast, likelihood-free
613 inference afforded by deep learning, the technical challenges shift from exploring models for
614 which tractable likelihood functions can be derived towards models that produce realistic
615 empirical data patterns, have parameters that control variation of those patterns, and are
616 efficient enough to generate large training data sets. A growing number of advanced

617 simulators are rapidly expanding the possibilities for deep learning in phylogenetics. For
618 example, FAVITES (Moshiri et al. 2019) is a simulator of disease spread through large
619 contact networks that tracks transmission trees and simulates sequence evolution. Gen3sis,
620 MASTER, SLiM, and VGsim are flexible simulation engines for generating complex
621 ecological, evolutionary, and disease transmission simulations (Hagen et al. 2021; Vaughan
622 and Drummond 2013; Shchur et al. 2021; Haller and Messer 2019; Overcast et al. 2021).
623 Continued advances in epidemic simulation speed and flexibility will be essential for
624 likelihood-free methods to push the boundaries of epidemic modeling sophistication and
625 usefulness.

626 There are several avenues of development still needed to realize the potential of
627 likelihood-free inference in phylogeography using deep learning. The current setup is ideal
628 for simulation experiments, but it is more difficult to ensure that the optimal parameter
629 values for empirical data sets are within the range of training data parameters.
630 Standardizing input tree height, geographical distance, and other parameters help make
631 training data more universally applicable. Simulation-trained neural networks are often
632 called amortized methods (Bürkner et al. 2022; Schmitt et al. 2022) because the cost of
633 inference is front-loaded, *i.e.* it takes time to simulate a training set and train a neural
634 network. The total cost in time per phylogenetic tree amortizes as the number of trees
635 analyzed by the trained model increases. These methods are therefore important when a
636 model is intended to be widely deployed or be responsive to an emerging outbreak where
637 policy decisions must be formulated rapidly. Because amortized approximate methods
638 require multiple analyses to realize time savings, researchers need to generate training data
639 sets over a broad parameter and model space so that trained networks can be applied to
640 new and diverse data sets. Our research focuses on one phase of an outbreak (the
641 exponential phase), but there are many other scenarios to be investigated, such as when
642 the stage of an epidemic differs among locations (e.g. exponential, peaked, declining).

643 Quantifying uncertainty is also crucial to data analysis and decision making, and

644 Bayesian statistics provides a framework for doing so in a rigorous way. Quantifying
645 uncertainty in predictions from deep neural networks is a difficult problem, as these models
646 are trained to minimize prediction error, rather than to explicitly estimate uncertainty. In
647 typical machine learning, uncertainty is ignored or measured using ad hoc methods in
648 which interpretation requires care. Many of these approaches come with their own
649 challenges and limitations, and there is still much ongoing research in this area to address
650 the challenge of quantifying uncertainty in deep neural networks (Gal et al. 2022).

651 Another important challenge of inference with deep learning is the problem of
652 convergence to a location on the loss function surface that approximates the maximum
653 likelihood well. There are a number of basic heuristics that can help such as learning
654 curves but more rigorous methods of ascertaining convergence is the subject of active
655 research Bürkner et al. (2022); Schmitt et al. (2022).

656 With recent advances in deep learning in epidemiology, evolution, and ecology
657 (Battey et al. 2020; Schrider and Kern 2018; Voznica et al. 2022; Radev et al. 2021;
658 Lambert et al. 2022; Rosenzweig et al. 2022; Suvorov and Schrider 2022a) biologists can
659 now explore the behavior of entire classes of stochastic branching models that are
660 biologically interesting but mathematically or statistically prohibitive for use with
661 traditional likelihood-based inference techniques. Although we are cautiously optimistic
662 about the future of deep learning methods for phylogenetics, it will become increasingly
663 important for the field to diagnose the conditions where phylogenetic deep learning
664 underperforms relative to likelihood-based approaches, and to devise general solutions for
665 the field.

666

FUNDING

667 National Geospatial-Intelligence Agency. MJL was supported by the National Science
668 Foundation (DEB 2040347) and by an internal grant awarded by the Incubator for
669 Transdisciplinary Futures at Washington University.

670

ACKNOWLEDGEMENTS

671 We are grateful to Fábio Mendes, Sarah Swiston, Sean McHugh, Walker Sexton, and

672 Mariana Braga for helpful comments on the research.

*

673

674 References

675 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
676 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian
677 Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal
678 Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mane, Rajat
679 Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens,
680 Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay
681 Vasudevan, Fernanda Viegas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin
682 Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-Scale Machine Learning on
683 Heterogeneous Distributed Systems, March 2016.

684 CJ Battey, Peter L Ralph, and Andrew D Kern. Predicting geographic location from
685 genetic variation with deep neural networks. *eLife*, 9:e54507, June 2020. ISSN
686 2050-084X. doi: 10.7554/eLife.54507.

687 Jeremy M. Beaulieu and Brian C. O’Meara. Detecting Hidden Diversification Shifts in
688 Models of Trait-Dependent Speciation and Extinction. *Systematic Biology*, 65(4):
689 583–601, July 2016. ISSN 1063-5157, 1076-836X. doi: 10.1093/sysbio/syw022.

690 Ottar N. Bjørnstad, Katriona Shea, Martin Krzywinski, and Naomi Altman. The SEIRS
691 model for infectious disease dynamics. *Nature Methods*, 17(6):557–558, June 2020. ISSN
692 1548-7091, 1548-7105. doi: 10.1038/s41592-020-0856-2.

693 Folmer Bokma. Artificial neural networks can learn to estimate extinction rates from
694 molecular phylogenies. *Journal of theoretical biology*, 243(3):449–454, 2006.

695 Remco Bouckaert, Timothy G. Vaughan, Joëlle Barido-Sottani, Sebastián Duchêne,
696 Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, Graham Jones, Denise

697 Kühnert, Nicola De Maio, Michael Matschiner, Fábio K. Mendes, Nicola F. Müller,
698 Huw A. Ogilvie, Louis du Plessis, Alex Poppinga, Andrew Rambaut, David Rasmussen,
699 Igor Siveroni, Marc A. Suchard, Chieh-Hsi Wu, Dong Xie, Chi Zhang, Tanja Stadler, and
700 Alexei J. Drummond. BEAST 2.5: An advanced software platform for Bayesian
701 evolutionary analysis. *PLOS Computational Biology*, 15(4):e1006650, April 2019. ISSN
702 1553-7358. doi: 10.1371/journal.pcbi.1006650.

703 Paul-Christian Bürkner, Maximilian Scholz, and Stefan Radev. Some models are useful,
704 but how do we know which ones? Towards a unified Bayesian model taxonomy,
705 September 2022.

706 Sheryl L. Chang, Mahendra Piraveenan, Philippa Pattison, and Mikhail Prokopenko.
707 Game theoretic modelling of infectious disease dynamics and intervention methods: A
708 review. *Journal of Biological Dynamics*, 14(1):57–89, January 2020. ISSN 1751-3758,
709 1751-3766. doi: 10.1080/17513758.2020.1720322.

710 F. K. Chollet. Keras: The Python deep learning API. <https://keras.io/>.

711 Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based
712 inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062,
713 December 2020. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1912789117.

714 Emanuel Masiero da Fonseca, Guarino R. Colli, Fernanda P. Werneck, and Bryan C.
715 Carstens. Phylogeographic model selection using convolutional neural networks,
716 September 2020.

717 Jordan Douglas, Fábio K Mendes, Remco Bouckaert, Dong Xie, Cinthy L Jiménez-Silva,
718 Christiaan Swanepoel, Joep de Ligt, Xiaoyun Ren, Matt Storey, James Hadfield, Colin R
719 Simpson, Jemma L Geoghegan, Alexei J Drummond, and David Welch. Phylodynamics
720 reveals the role of human travel and contact tracing in controlling the first wave of

- 721 COVID-19 in four island nations. *Virus Evolution*, 7(2), September 2021. ISSN
722 2057-1577. doi: 10.1093/ve/veab052.
- 723 Richard G. FitzJohn. Diversitree: Comparative phylogenetic analyses of diversification in
724 R. *Methods in Ecology and Evolution*, 3(6):1084–1092, 2012. ISSN 2041-210X. doi:
725 10.1111/j.2041-210X.2012.00234.x.
- 726 Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The Unreasonable Effectiveness of
727 Convolutional Neural Networks in Population Genetic Inference. *Molecular Biology and*
728 *Evolution*, 36(2):220–238, February 2019. ISSN 0737-4038, 1537-1719. doi:
729 10.1093/molbev/msy224.
- 730 Yarin Gal, Petros Koumoutsakos, Francois Lanusse, Gilles Louppe, and Costas
731 Papadimitriou. Bayesian uncertainty quantification for machine-learned models in
732 physics. *Nature Reviews Physics*, August 2022. ISSN 2522-5820. doi:
733 10.1038/s42254-022-00498-4.
- 734 Jiansi Gao, Michael R. May, Bruce Rannala, and Brian R. Moore. New Interval-Specific
735 Phylodynamic Models Improve Inference of the Geographic History of Disease
736 Outbreaks, December 2021.
- 737 Jiansi Gao, Michael R. May, Bruce Rannala, and Brian R. Moore. Model Misspecification
738 Misleads Inference of the Spatial Dynamics of Disease Outbreaks, August 2022.
- 739 James Hadfield, Colin Megill, Sidney M Bell, John Huddleston, Barney Potter, Charlton
740 Callender, Pavel Sagulenko, Trevor Bedford, and Richard A Neher. Nextstrain: real-time
741 tracking of pathogen evolution. *Bioinformatics*, 34(23):4121–4123, December 2018. ISSN
742 1367-4803. doi: 10.1093/bioinformatics/bty407. URL
743 <https://doi.org/10.1093/bioinformatics/bty407>.
- 744 Oskar Hagen, Benjamin Flück, Fabian Fopp, Juliano S. Cabral, Florian Hartig, Mikael
745 Pontarp, Thiago F. Rangel, and Loïc Pellissier. Gen3sis: A general engine for

- 746 eco-evolutionary simulations of the processes that shape Earth's biodiversity. *PLOS*
747 *Biology*, 19(7):e3001340, July 2021. ISSN 1545-7885. doi: 10.1371/journal.pbio.3001340.
- 748 Benjamin C Haller and Philipp W Messer. SLiM 3: Forward Genetic Simulations Beyond
749 the Wright–Fisher Model. *Molecular Biology and Evolution*, 36(3):632–637, March 2019.
750 ISSN 0737-4038. doi: 10.1093/molbev/msy228.
- 751 Sebastian Höhna, Michael J. Landis, Tracy A. Heath, Bastien Boussau, Nicolas Lartillot,
752 Brian R. Moore, John P. Huelsenbeck, and Fredrik Ronquist. RevBayes: Bayesian
753 Phylogenetic Inference Using Graphical Models and an Interactive Model-Specification
754 Language. *Systematic Biology*, 65(4):726–736, July 2016. ISSN 1063-5157, 1076-836X.
755 doi: 10.1093/sysbio/syw021.
- 756 Eddie C Holmes and Geoff P Garnett. Genes, trees and infections: molecular evidence in
757 epidemiology. *Trends in Ecology & Evolution*, 9(7):256–260, 1994.
- 758 Eddie C Holmes, Sean Nee, Andrew Rambaut, Geoff P Garnett, and Paul H Harvey.
759 Revealing the history of infectious disease epidemics through phylogenetic trees.
760 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*,
761 349(1327):33–40, 1995.
- 762 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization,
763 January 2017.
- 764 Denise Kühnert, Tanja Stadler, Timothy G. Vaughan, and Alexei J. Drummond.
765 Simultaneous reconstruction of evolutionary history and epidemiological dynamics from
766 viral sequences with the birth–death SIR model. *Journal of The Royal Society Interface*,
767 11(94):20131106, May 2014. ISSN 1742-5689, 1742-5662. doi: 10.1098/rsif.2013.1106.
- 768 Denise Kühnert, Tanja Stadler, Timothy G. Vaughan, and Alexei J. Drummond.
769 Phylodynamics with Migration: A Computational Framework to Quantify Population

- 770 Structure from Genomic Data. *Molecular Biology and Evolution*, 33(8):2102–2116,
771 August 2016. ISSN 0737-4038. doi: 10.1093/molbev/msw064.
- 772 Sophia Lambert, Jakub Voznica, and H el ene Morlon. Deep Learning from Phylogenies for
773 Diversification Analyses, September 2022.
- 774 Philippe Lemey, Andrew Rambaut, Alexei J. Drummond, and Marc A. Suchard. Bayesian
775 Phylogeography Finds Its Roots. *PLoS Computational Biology*, 5(9):e1000520,
776 September 2009. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1000520.
- 777 Philippe Lemey, Nick Ruktanonchai, Samuel L. Hong, Vittoria Colizza, Chiara Poletto,
778 Frederik Van den Broeck, Mandev S. Gill, Xiang Ji, Anthony Levasseur, Bas B.
779 Oude Munnink, Marion Koopmans, Adam Sadilek, Shengjie Lai, Andrew J. Tatem, Guy
780 Baele, Marc A. Suchard, and Simon Dellicour. Untangling introductions and persistence
781 in COVID-19 resurgence in Europe. *Nature*, June 2021. ISSN 0028-0836, 1476-4687. doi:
782 10.1038/s41586-021-03754-2.
- 783 Fr ed eric Lemoine and Olivier Gascuel. Gotree/Goalign: Toolkit and Go API to facilitate
784 the development of phylogenetic workflows. *NAR Genomics and Bioinformatics*, 3(3):
785 lqab075, September 2021. ISSN 2631-9268. doi: 10.1093/nargab/lqab075.
- 786 Lu Lu, Reina S. Sikkema, Francisca C. Velkers, David F. Nieuwenhuijse, Egil A. J. Fischer,
787 Paola A. Meijer, Noortje Bouwmeester-Vincken, Ariene Rietveld, Marjolijn C. A.
788 Wegdam-Blans, Paulien Tolsma, Marco Koppelman, Lidwien A. M. Smit, Renate W.
789 Hakze-van der Honing, Wim H. M. van der Poel, Arco N. van der Spek, Marcel A. H.
790 Spierenburg, Robert Jan Molenaar, Jan de Rond, Marieke Augustijn, Mark Woolhouse,
791 J. Arjan Stegeman, Samantha Lycett, Bas B. Oude Munnink, and Marion P. G.
792 Koopmans. Adaptation, spread and transmission of SARS-CoV-2 in farmed minks and
793 associated humans in the Netherlands. *Nature Communications*, 12(1):6802, December
794 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-27096-9.

- 795 Wayne P. Maddison, Peter E. Midford, and Sarah P. Otto. Estimating a Binary
796 Character's Effect on Speciation and Extinction. *Systematic Biology*, 56(5):701–710,
797 October 2007. ISSN 1076-836X, 1063-5157. doi: 10.1080/10635150701607033.
- 798 Odile Maliet, Florian Hartig, and H el ene Morlon. A model with many small shifts for
799 estimating species-specific diversification rates. *Nature Ecology & Evolution*, 3(7):
800 1086–1092, July 2019. ISSN 2397-334X. doi: 10.1038/s41559-019-0908-0.
- 801 Mike Meredith and John Kruschke. Bayesian Estimation Supersedes the t-Test. page 13.
- 802 Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D
803 Woodhams, Arndt von Haeseler, and Robert Lanfear. IQ-TREE 2: New Models and
804 Efficient Methods for Phylogenetic Inference in the Genomic Era. *Molecular Biology and*
805 *Evolution*, 37(5):1530–1534, May 2020. ISSN 0737-4038. doi: 10.1093/molbev/msaa015.
- 806 Niema Moshiri, Manon Ragonnet-Cronin, Joel O Wertheim, and Siavash Mirarab.
807 FAVITES: Simultaneous simulation of transmission networks, phylogenetic trees and
808 sequences. *Bioinformatics*, 35(11):1852–1861, June 2019. ISSN 1367-4803, 1460-2059.
809 doi: 10.1093/bioinformatics/bty921.
- 810 Sarah A. Nadeau, Timothy G. Vaughan, J er emie Scire, Jana S. Huisman, and Tanja
811 Stadler. The origin and early spread of SARS-CoV-2 in Europe. *Proceedings of the*
812 *National Academy of Sciences*, 118(9):e2012008118, March 2021. ISSN 0027-8424,
813 1091-6490. doi: 10.1073/pnas.2012008118.
- 814 Luca Nesterenko, Bastien Boussau, and Laurent Jacob. Phyloformer: Towards fast and
815 accurate phylogeny estimation with self-attention networks, June 2022.
- 816 Eamon B O'Dea and John M Drake. A semi-parametric, state-space compartmental model
817 with time-dependent parameters for forecasting COVID-19 cases, hospitalizations, and
818 deaths. page 32, 2021.

- 819 Isaac Overcast, Megan Ruffley, James Rosindell, Luke Harmon, Paulo AV Borges, Brent C
820 Emerson, Rampal S Etienne, Rosemary Gillespie, Henrik Krehenwinkel, D Luke Mahler,
821 et al. A unified model of species abundance, genetic diversity, and functional diversity
822 reveals the mechanisms structuring ecological communities. *Molecular Ecology*
823 *Resources*, 21(8):2782–2800, 2021.
- 824 Jonathan E. Pekar, Andrew Magee, Edyth Parker, Niema Moshiri, Katherine Izhikevich,
825 Jennifer L. Havens, Karthik Gangavarapu, Lorena Mariana Malpica Serrano, Alexander
826 Crits-Christoph, Nathaniel L. Matteson, Mark Zeller, Joshua I. Levy, Jade C. Wang,
827 Scott Hughes, Jungmin Lee, Heedo Park, Man-Seong Park, Katherine Zi Yan Ching,
828 Raymond Tzer Pin Lin, Mohd Noor Mat Isa, Yusuf Muhammad Noor, Tetyana I.
829 Vasylyeva, Robert F. Garry, Edward C. Holmes, Andrew Rambaut, Marc A. Suchard,
830 Kristian G. Andersen, Michael Worobey, and Joel O. Wertheim. The molecular
831 epidemiology of multiple zoonotic origins of SARS-CoV-2. *Science*, 0(0):eabp8337, July
832 2022. doi: 10.1126/science.abp8337.
- 833 José M. Ponciano and Marcos A. Capistrán. First Principles Modeling of Nonlinear
834 Incidence Rates in Seasonal Epidemics. *PLOS Computational Biology*, 7(2):e1001079,
835 February 2011. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1001079.
- 836 O. G. Pybus, M. A. Suchard, P. Lemey, F. J. Bernardin, A. Rambaut, F. W. Crawford,
837 R. R. Gray, N. Arinaminpathy, S. L. Stramer, M. P. Busch, and E. L. Delwart. Unifying
838 the spatial epidemiology and molecular evolution of emerging epidemics. *Proceedings of*
839 *the National Academy of Sciences*, 109(37):15066–15071, September 2012. ISSN
840 0027-8424, 1091-6490. doi: 10.1073/pnas.1206598109.
- 841 Stefan T. Radev, Frederik Graw, Simiao Chen, Nico T. Mutters, Vanessa M. Eichel, Till
842 Bärnighausen, and Ullrich Köthe. OutbreakFlow: Model-based Bayesian inference of
843 disease outbreak dynamics with invertible neural networks and its application to the

- 844 COVID-19 pandemics in Germany. *PLOS Computational Biology*, 17(10):e1009472,
845 October 2021. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1009472.
- 846 A. Rambaut and N. C. Grassly. Seq-Gen: an application for the Monte Carlo simulation of
847 DNA sequence evolution along phylogenetic trees. *Computer Applications in the*
848 *Biosciences*, 13:235–238, 1997.
- 849 Andrew Rambaut, Oliver G Pybus, Martha I Nelson, Cecile Viboud, Jeffery K
850 Taubenberger, and Edward C Holmes. The genomic and epidemiological dynamics of
851 human influenza a virus. *Nature*, 453(7195):615–619, 2008.
- 852 Liam J. Revell. Phytools: An R package for phylogenetic comparative biology (and other
853 things). *Methods in Ecology and Evolution*, 3(2):217–223, 2012. ISSN 2041-210X. doi:
854 10.1111/j.2041-210X.2011.00169.x.
- 855 Francisco Richter, Bart Haegeman, Rampal S. Etienne, and Ernst C. Wit. Introducing a
856 general class of species diversification models for phylogenetic trees. *Statistica*
857 *Neerlandica*, 74(3):261–274, 2020. ISSN 1467-9574. doi: 10.1111/stan.12205.
- 858 D. F. Robinson and L. R. Foulds. Comparison of phylogenetic trees. *Mathematical*
859 *Biosciences*, 53:131–147, 1981.
- 860 Benjamin K. Rosenzweig, Matthew W. Hahn, and Andrew Kern. Accurate Detection of
861 Incomplete Lineage Sorting via Supervised Machine Learning, November 2022.
- 862 Marvin Schmitt, Paul-Christian Bürkner, Ullrich Köthe, and Stefan T. Radev. Detecting
863 Model Misspecification in Amortized Bayesian Inference with Neural Networks, May
864 2022.
- 865 Daniel R. Schrider and Andrew D. Kern. Supervised Machine Learning for Population
866 Genetics: A New Paradigm. *Trends in Genetics*, 34(4):301–312, April 2018. ISSN
867 01689525. doi: 10.1016/j.tig.2017.12.005.

- 868 Jérémie Scire, Joëlle Barido-Sottani, Denise Kühnert, Timothy G. Vaughan, and Tanja
869 Stadler. Improved multi-type birth-death phylodynamic inference in BEAST 2. Preprint,
870 Evolutionary Biology, January 2020.
- 871 Vladimir Shchur, Vadim Spirin, Dmitry Sirotkin, Evgeni Burovski, Nicola De Maio, and
872 Russell Corbett-Detig. VGsim: Scalable viral genealogy simulator for global pandemic.
873 Preprint, Epidemiology, April 2021.
- 874 Claudia Solis-Lemus, Shengwen Yang, and Leonardo Zepeda-Nunez. Accurate Phylogenetic
875 Inference with a Symmetry-preserving Neural Network Model, January 2022.
- 876 Tanja Stadler. Sampling-through-time in birth–death trees. *Journal of Theoretical Biology*,
877 267(3):396–404, December 2010. ISSN 00225193. doi: 10.1016/j.jtbi.2010.09.010.
- 878 Tanja Stadler, Roger Kouyos, Viktor von Wyl, Sabine Yerly, Jürg Böni, Philippe Bürgisser,
879 Thomas Klimkait, Beda Joos, Philip Rieder, Dong Xie, Huldrych F. Günthard, Alexei J.
880 Drummond, Sebastian Bonhoeffer, and the Swiss HIV Cohort Study. Estimating the
881 Basic Reproductive Number from Viral Sequence Data. *Molecular Biology and Evolution*,
882 29(1):347–357, January 2012. ISSN 1537-1719, 0737-4038. doi: 10.1093/molbev/msr217.
- 883 Anton Suvorov and Daniel Schrider. Reliable estimation of tree branch lengths using deep
884 neural networks. *bioRxiv*, 2022a.
- 885 Anton Suvorov and Daniel R. Schrider. Reliable estimation of tree branch lengths using
886 deep neural networks, November 2022b.
- 887 Anton Suvorov, Joshua Hochuli, and Daniel R Schrider. Accurate Inference of Tree
888 Topologies from Multiple Sequence Alignments Using Deep Learning. *Systematic Biology*,
889 69(2):221–233, March 2020. ISSN 1063-5157, 1076-836X. doi: 10.1093/sysbio/syz060.
- 890 Timothy G. Vaughan and Alexei J. Drummond. A Stochastic Simulator of Birth–Death

891 Master Equations with Application to Phylodynamics. *Molecular Biology and Evolution*,
892 30(6):1480–1493, June 2013. ISSN 0737-4038. doi: 10.1093/molbev/mst057.

893 Timothy G. Vaughan, Denise Kühnert, Alex Poppinga, David Welch, and Alexei J.
894 Drummond. Efficient Bayesian inference under the structured coalescent.
895 *Bioinformatics*, 30(16):2272–2279, August 2014. ISSN 1367-4803, 1460-2059. doi:
896 10.1093/bioinformatics/btu201.

897 Erik M. Volz and Igor Siveroni. Bayesian phylodynamic inference with complex models.
898 *PLOS Computational Biology*, 14(11):e1006546, November 2018. ISSN 1553-7358. doi:
899 10.1371/journal.pcbi.1006546.

900 Erik M. Volz, Katia Koelle, and Trevor Bedford. Viral Phylodynamics. *PLOS*
901 *Computational Biology*, 9(3):e1002947, March 2013. ISSN 1553-7358. doi:
902 10.1371/journal.pcbi.1002947. URL [https:](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002947)
903 [//journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002947](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002947).
904 Publisher: Public Library of Science.

905 J Voznica, A Zhukova, V Boskova, E Saulnier, F Lemoine, M Moslonka-Lefebvre, and
906 O Gascuel. Deep learning from phylogenies to uncover the epidemiological dynamics of
907 outbreaks. preprint, *Bioinformatics*, March 2021. URL
908 <http://biorxiv.org/lookup/doi/10.1101/2021.03.11.435006>.

909 J. Voznica, A. Zhukova, V. Boskova, E. Saulnier, F. Lemoine, M. Moslonka-Lefebvre, and
910 O. Gascuel. Deep learning from phylogenies to uncover the epidemiological dynamics of
911 outbreaks. *Nature Communications*, 13(1):3896, July 2022. ISSN 2041-1723. doi:
912 10.1038/s41467-022-31511-0.

913 Nicole L. Washington, Karthik Gangavarapu, Mark Zeller, Alexandre Bolze, Elizabeth T.
914 Cirulli, Kelly M. Schiabor Barrett, Brendan B. Larsen, Catelyn Anderson, Simon White,
915 Tyler Cassens, Sharoni Jacobs, Geraint Levan, Jason Nguyen, Jimmy M. Ramirez,

916 Charlotte Rivera-Garcia, Efren Sandoval, Xueqing Wang, David Wong, Emily Spencer,
917 Refugio Robles-Sikisaka, Ezra Kurzban, Laura D. Hughes, Xianding Deng, Candace
918 Wang, Venice Servellita, Holly Valentine, Peter De Hoff, Phoebe Seaver, Shashank Sathe,
919 Kimberly Gietzen, Brad Sickler, Jay Antico, Kelly Hoon, Jingtao Liu, Aaron Harding,
920 Omid Bakhtar, Tracy Basler, Brett Austin, Duncan MacCannell, Magnus Isaksson,
921 Phillip G. Febbo, David Becker, Marc Laurent, Eric McDonald, Gene W. Yeo, Rob
922 Knight, Louise C. Laurent, Eileen de Feo, Michael Worobey, Charles Y. Chiu, Marc A.
923 Suchard, James T. Lu, William Lee, and Kristian G. Andersen. Emergence and rapid
924 transmission of SARS-CoV-2 B.1.1.7 in the United States. *Cell*, 184(10):2587–2594.e7,
925 May 2021. ISSN 00928674. doi: 10.1016/j.cell.2021.03.052.

926 Michael Worobey, Thomas D Watts, Richard A McKay, Marc A Suchard, Timothy
927 Granade, Dirk E Teuwen, Beryl A Koblin, Walid Heneine, Philippe Lemey, and
928 Harold W Jaffe. 1970s and ‘patient 0’ hiv-1 genomes illuminate early hiv/aids history in
929 north america. *Nature*, 539(7627):98–101, 2016.

930 Michael Worobey, Jonathan Pekar, Brendan B. Larsen, Martha I. Nelson, Verity Hill,
931 Jeffrey B. Joy, Andrew Rambaut, Marc A. Suchard, Joel O. Wertheim, and Philippe
932 Lemey. The emergence of SARS-CoV-2 in Europe and North America. *Science*, 370
933 (6516):564–570, October 2020. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abc8169.

934

SUPPLEMENTAL TABLES AND FIGURES

Table S1: BEST comparisons between CNN and Bayesian absolute percent errors (APEs) for model parameters across all experiments.

95% HPD intervals of average relative error from BEST analysis			
True inference model (Reference for misspecification experiments)	Median CNN APE	Median Like.-based APE	median(CNN APE - Like.-based APE)
R_0	2.4, 3.5	2.1, 3.1	0.1, 1.2
δ	7.0, 10.5	5.7, 8.9	0.2, 3.0
m	9.5, 14.1	8.4, 12.1	0.4, 3.2
Misspecified R_0 experiment			
	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	11.8, 17.8	11.0, 16.9	-0.1, 1.6
δ	0.8, 7.6	-0.6, 5.3	1.3, 5.8
m	8.2, 17.9	6.5, 15.9	1.3, 4.7
Misspecified sample rate experiment			
	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.3, 1.7	0.03, 1.7	0.1, 1.3
δ	12.0, 21.2	12.6, 21.4	0.1, 4.0
m	3.3, 12.0	5.6, 14.4	-1.2, 2.7
Misspecified migration rate experiment			
	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.9, 0.8	-0.6, 1.0	-0.5, 0.8
δ	-2.3, 3.3	0.1, 5.8	-1.4, 2.3
m	4.0, 15.2	5.0, 16.2	-1.3, 2.6
Misspecified number of locations experiment			
	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.3, 1.5	-0.7, 0.8	0.5, 1.9
δ	-0.3, 4.9	-0.5, 4.2	0.4, 3.5
m	3.4, 11.1	5.8, 13.5	-0.9, 1.6
Phylogenetic error experiment			
	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	0.7, 3.0	1.7, 4.4	-1.4, 0.1
δ	2.3, 9.6	1.5, 7.2	1.4, 5.3
m	-1.2, 6.0	-1.8, 5.4	-1.7, 2.4

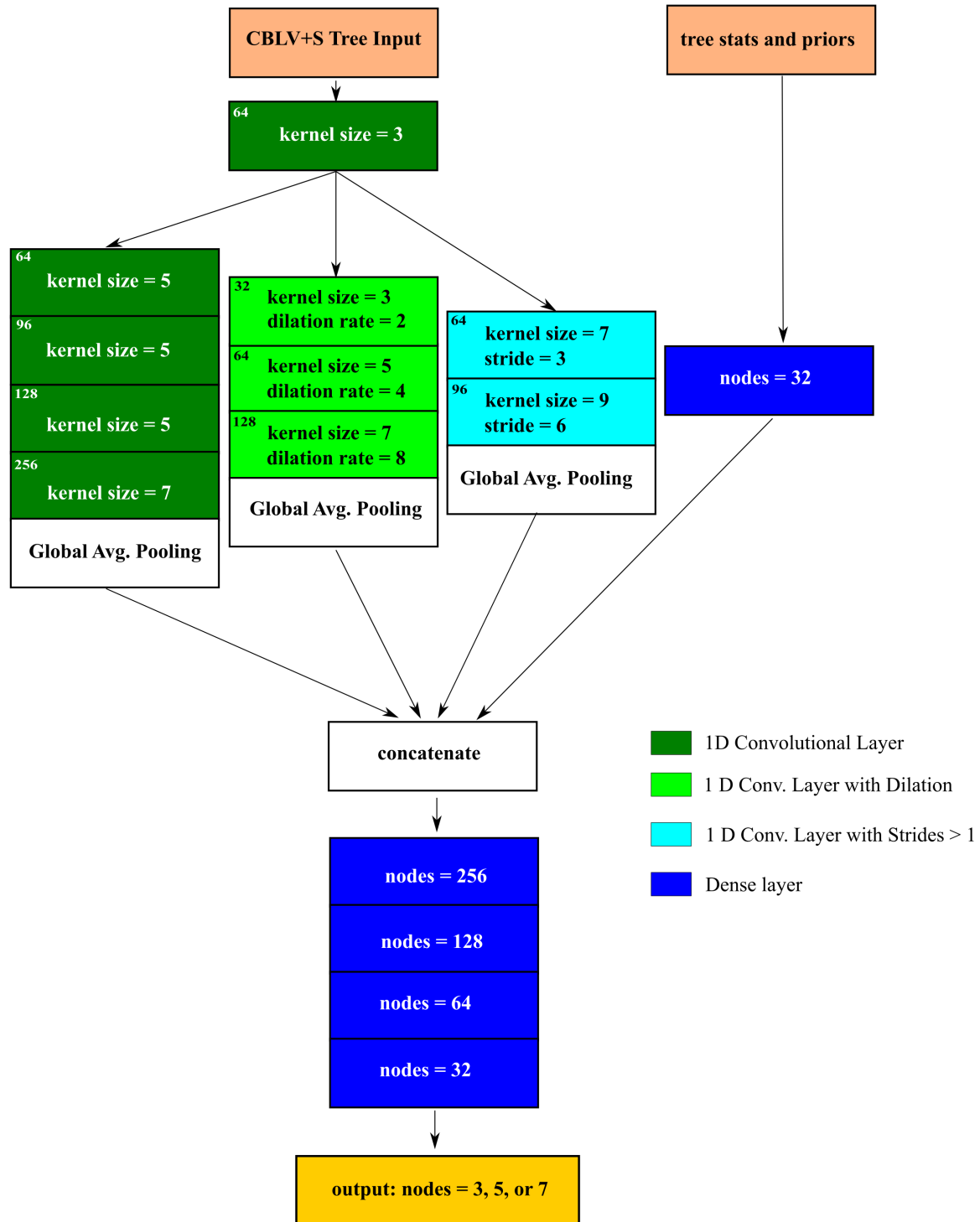


Figure S1: Diagram of deep neural network trained to make 2 kinds of predictions (rates and origin location) under two models (MTBD and SDMTBD).

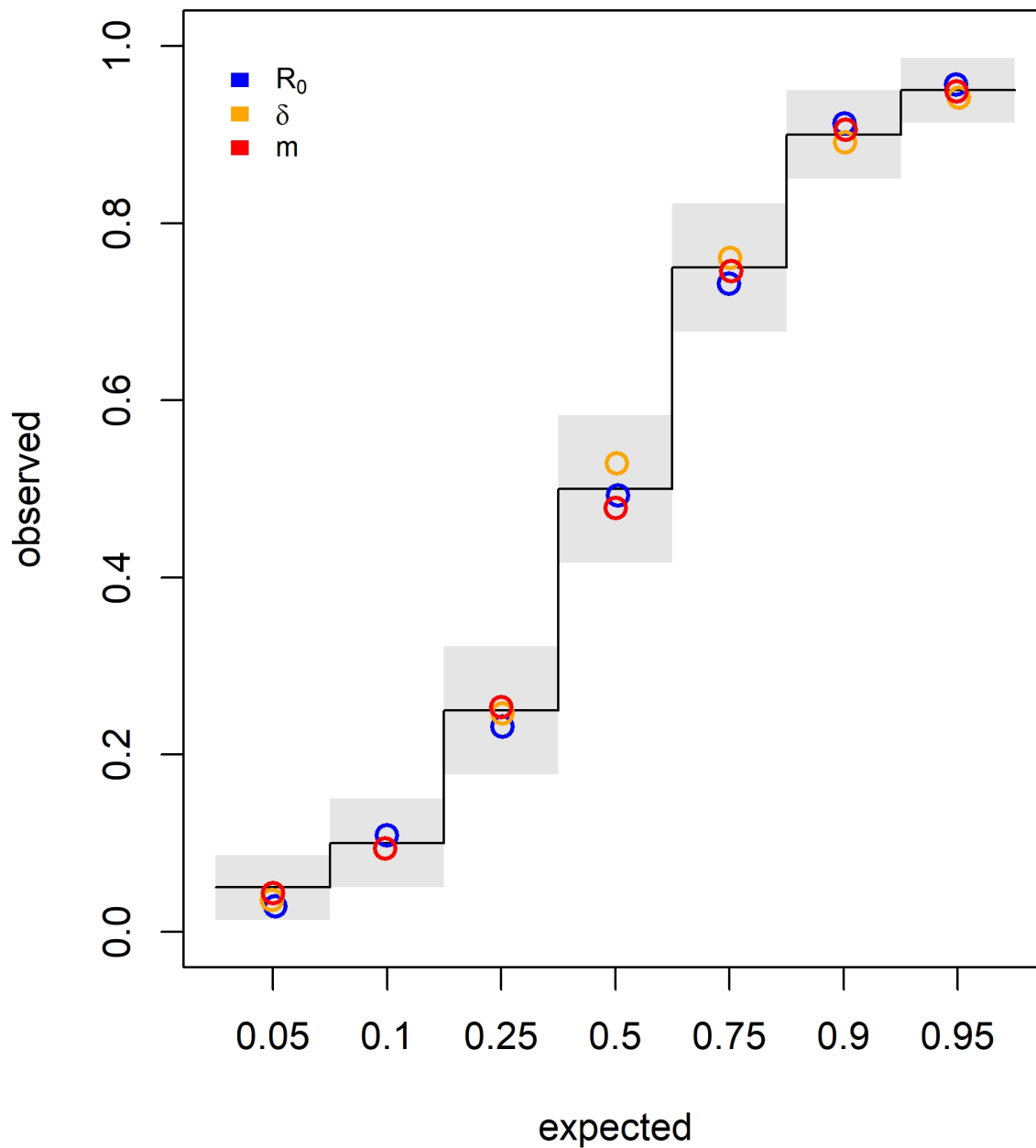


Figure S2: Coverage of posterior distributions simulated with TensorPhylo. Seven different HPD intervals were measured for coverage (labeled horizontal). The expected frequency of coverage for each of the categories is shown at the black steps. Gray bands indicate the 95% confidence intervals for estimates of the binomial proportion at each of the expected values. The colored circles indicate the observed coverage of the three rate parameters at each of the expected values.

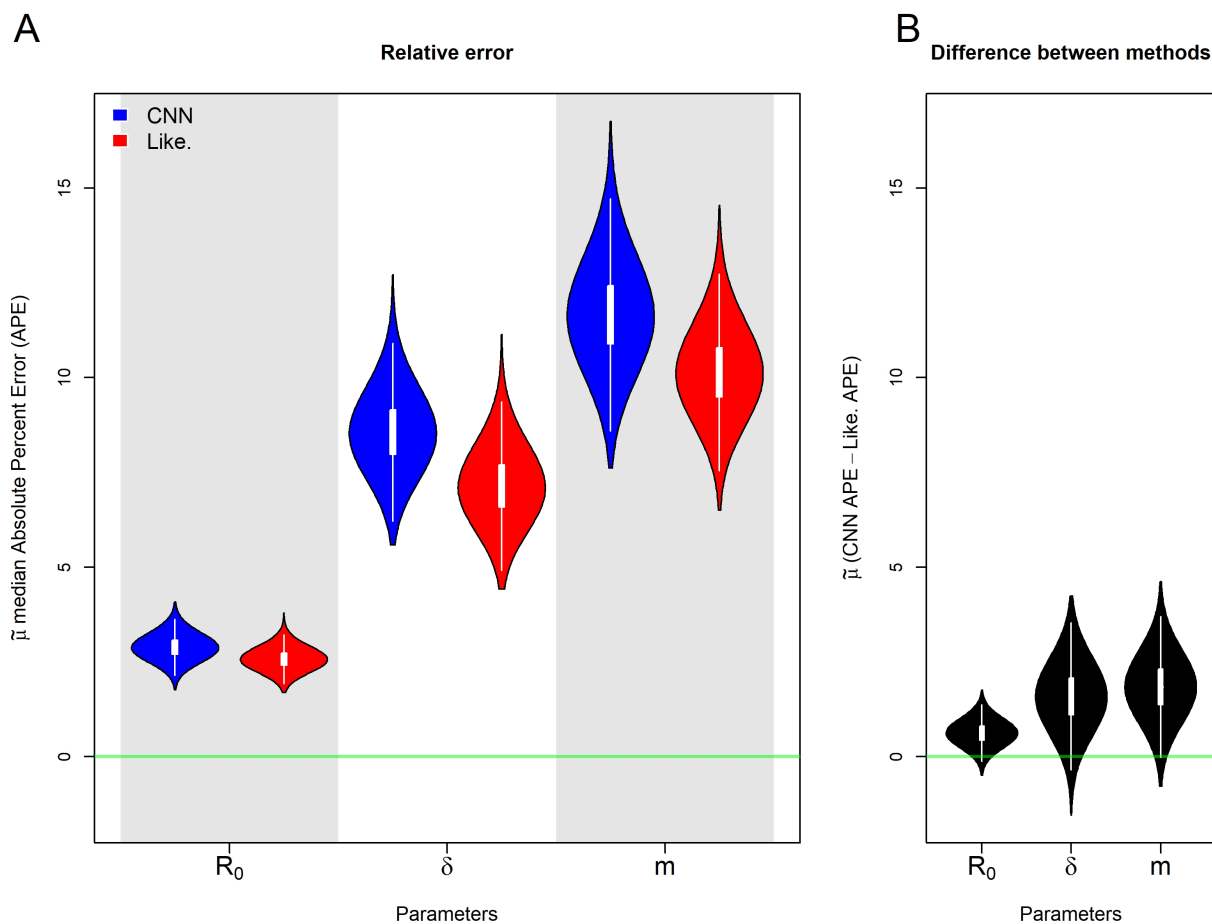


Figure S3: Posterior distributions of the population median, $\tilde{\mu}$, APE estimates of the rate parameters R_0 , δ , and m under the true model. A) shows posterior distribution of the median APE for each of the 3 rate parameters estimated by the CNN (blue) and the likelihood-based method (red). The green line indicates no error. B) shows the posterior distribution for the median difference between the CNN estimate's APE and the likelihood-based estimate's APE. The green line indicates the median APE difference is zero.

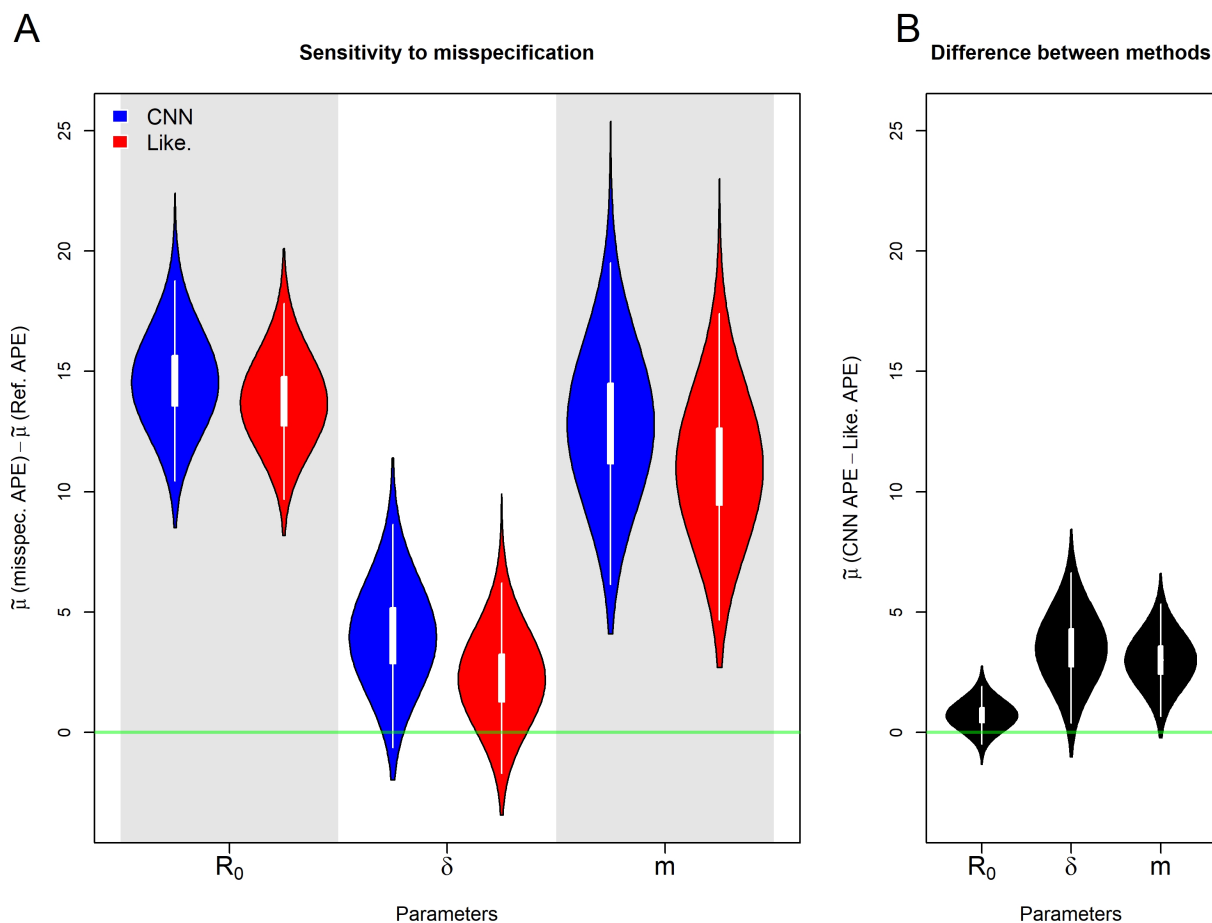


Figure S4: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified R_0 experiment. A) shows posterior distribution of the difference between the median error under the misspecified model and the the median error under the true, reference model. B) shows the posterior distribution for the population median difference between the CNN estimate's APE and the likelihood-based estimate's APE.

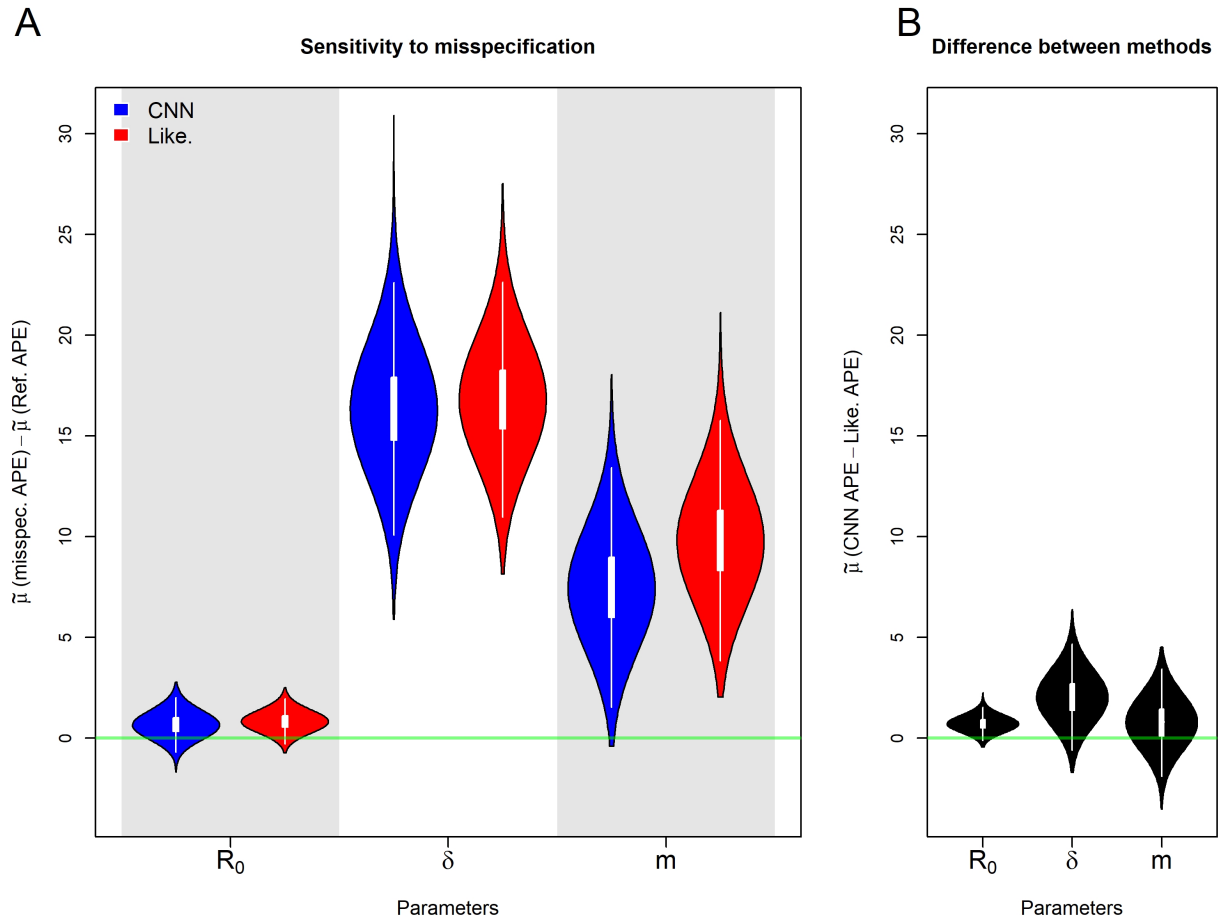


Figure S5: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified sampling rate, δ , experiment. Details are the same as in S4

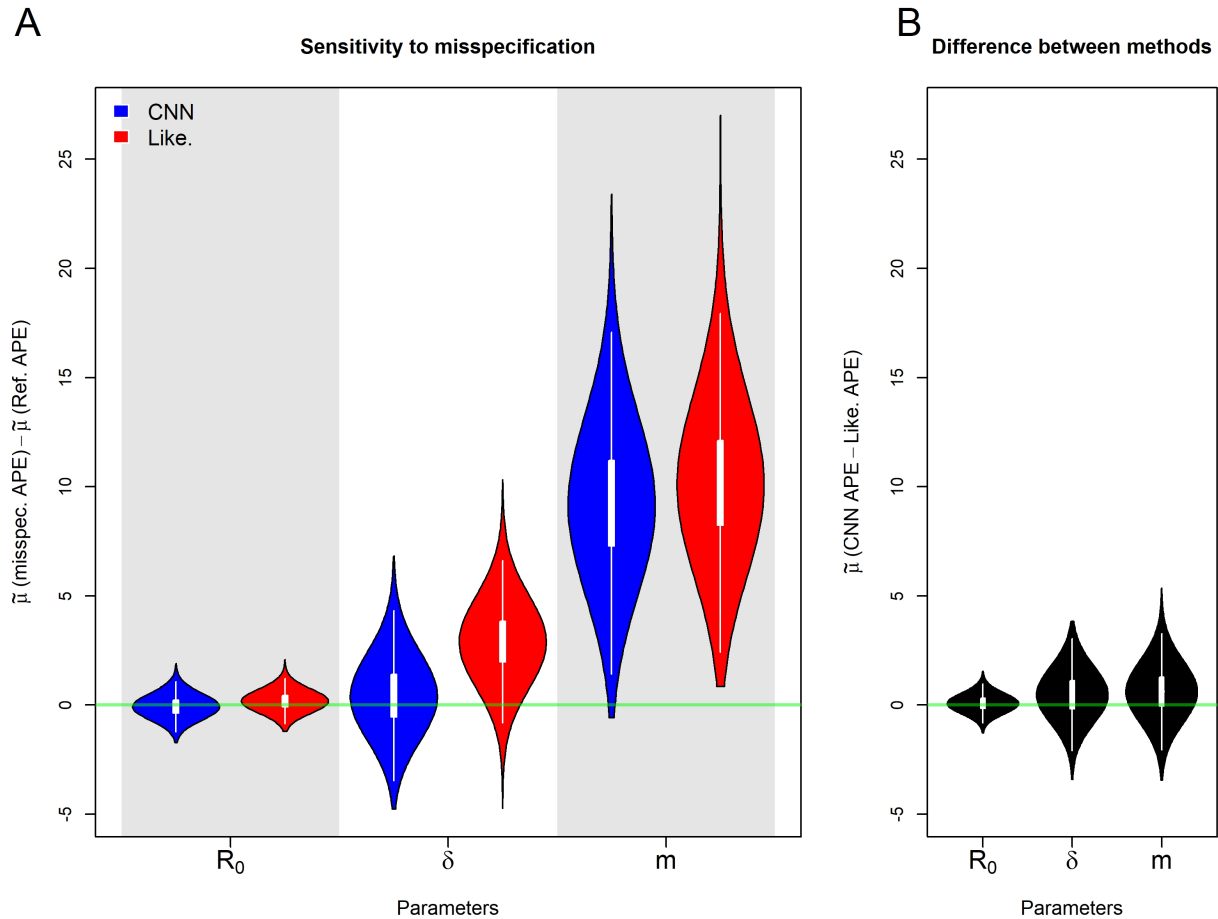


Figure S6: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified migration rate, m , experiment. Details are the same as in S4

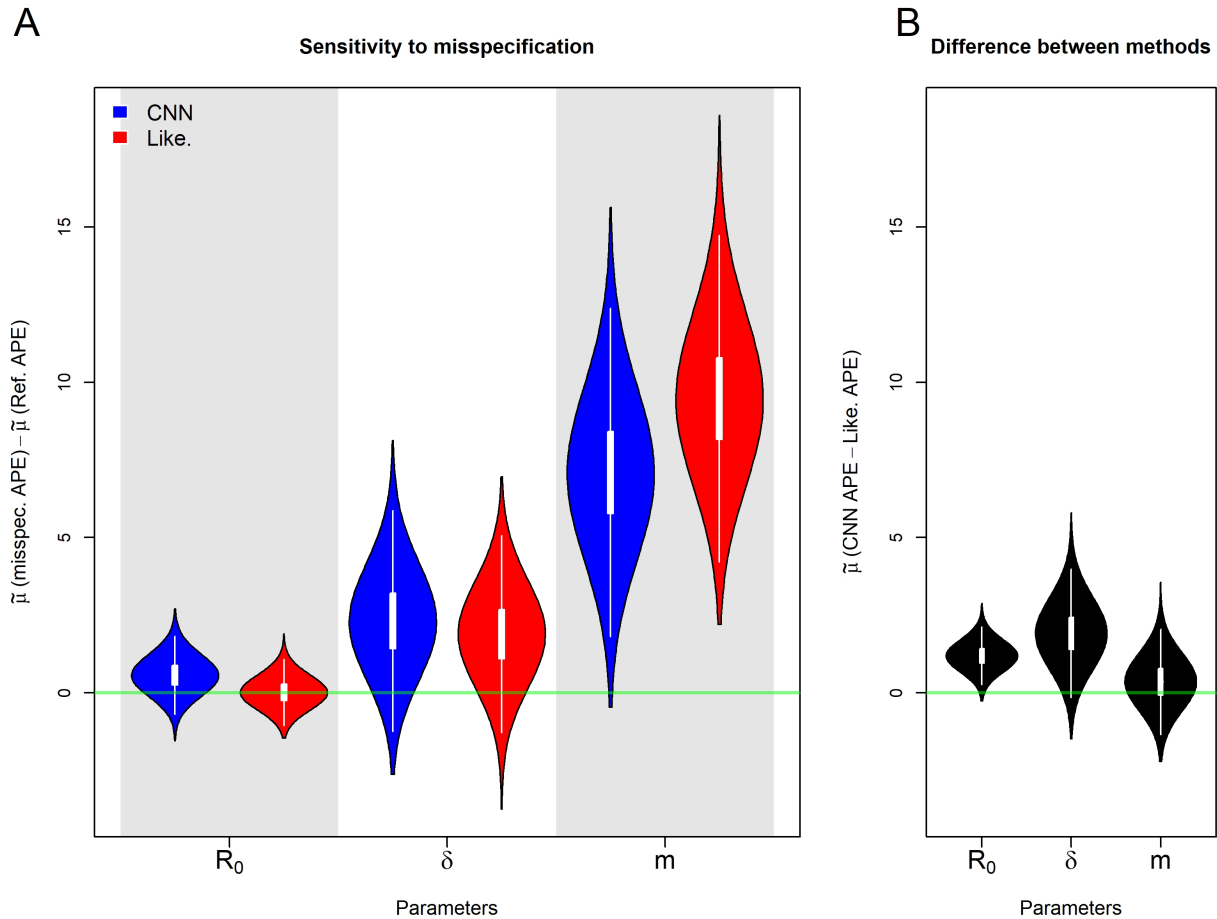


Figure S7: Posterior distributions of the median APE when the model is misspecified for the number of locations. Details are the same as in S4

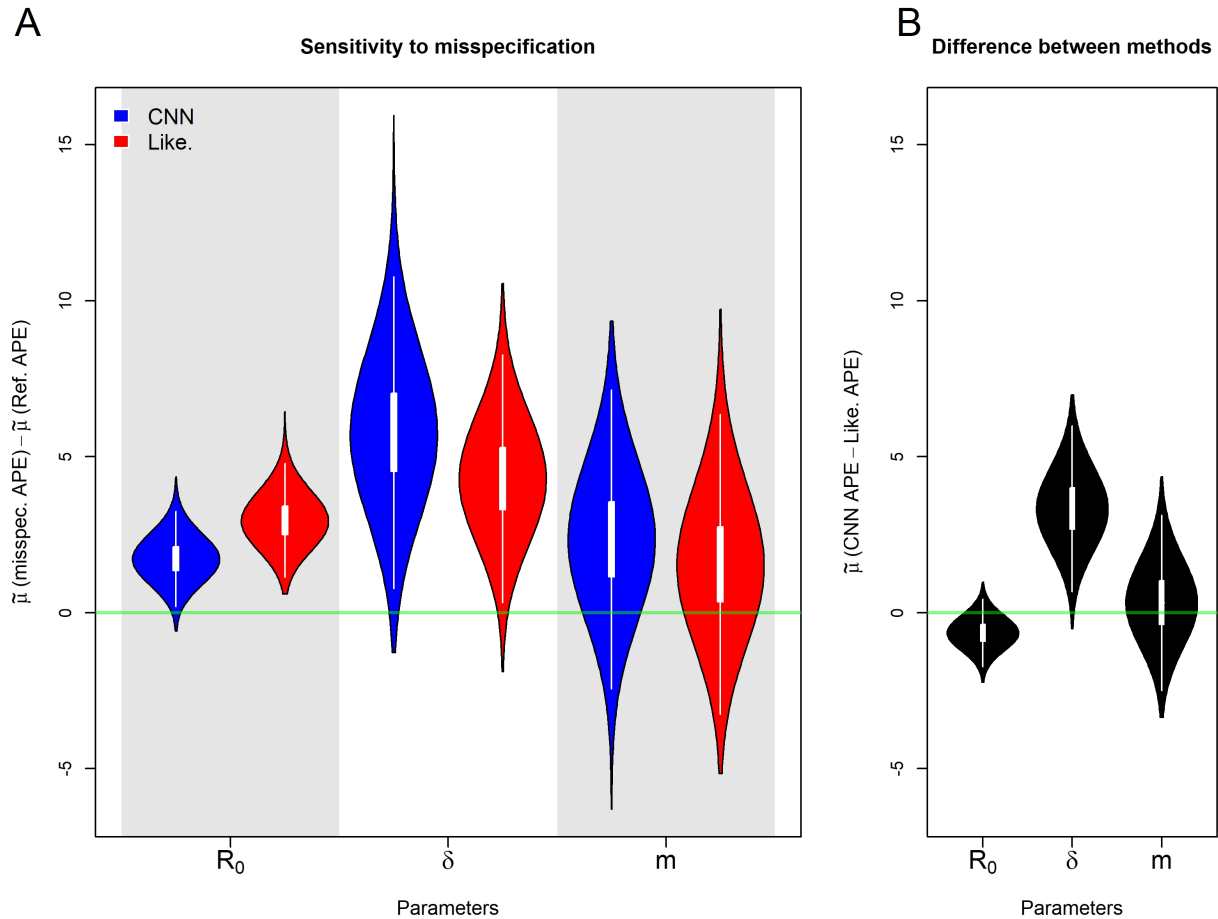


Figure S8: Posterior distributions of the median APE when the phylogenetic tree is incorrect. Details are the same as in S4