



HAL
open science

Deep learning from phylogenies to uncover the transmission dynamics of epidemics

J Voznica, Anna Zhukova, V Boskova, E Saulnier, Frédéric Lemoine, M Moslonka-Lefebvre, O Gascuel

► To cite this version:

J Voznica, Anna Zhukova, V Boskova, E Saulnier, Frédéric Lemoine, et al.. Deep learning from phylogenies to uncover the transmission dynamics of epidemics. 2021. hal-03428718

HAL Id: hal-03428718

<https://hal.science/hal-03428718>

Preprint submitted on 15 Nov 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

1 **DEEP LEARNING FROM PHYLOGENIES TO UNCOVER**
2 **THE TRANSMISSION DYNAMICS OF EPIDEMICS**

3
4 **AUTHORS**

5 Voznica J^{1,2,3,4*}, Zhukova A^{1,2,5*}, Boskova V⁶, Saulnier E^{1,2}, Lemoine F^{1,2,5}, Moslonka-Lefebvre M^{1,2}, Gascuel O^{1,2¶*}

6
7 **AFFILIATIONS**

8 ¹Unité de Bioinformatique Évolutive - Département Biologie computationnelle, Institut Pasteur, Paris, FRANCE

9 ²Unité de Biologie Computationnelle, USR 3756 CNRS, Paris, FRANCE

10 ³Université de Paris, Paris, FRANCE

11 ⁴Institut de Biologie de l'École Normale Supérieure, Ecole Normale Supérieure, CNRS, INSERM, Université Paris
12 Sciences et Lettres, 75005 Paris, FRANCE

13 ⁵Hub de Bioinformatique et Biostatistique - Département Biologie computationnelle, Institut Pasteur, Paris,
14 FRANCE

15 ⁶Center for Integrative Bioinformatics Vienna, Max Perutz Labs, University of Vienna and Medical University of
16 Vienna, Vienna, AUSTRIA

17 ¶ OG: current address

18
19 *** CO-CORRESPONDING AUTHORS**

20 jakub.voznica@pasteur.fr (JV), anna.zhukova@pasteur.fr (AZ), gascuelolivier@gmail.com (OG)

21

22 **ABSTRACT**

23 Widely applicable, accurate and fast inference methods in phylodynamics are needed to fully profit from the richness
24 of genetic data in uncovering the dynamics of epidemics. Standard methods, including maximum-likelihood and
25 Bayesian approaches, which are both model specific, often rely on complex mathematical formulae and
26 approximations, and do not scale well with dataset size. We develop a likelihood-free, simulation-based approach,
27 which combines deep learning with (1) a large set of summary statistics measured on phylogenies or (2) a complete
28 and compact vectorial representation of trees, which avoids potential limitations of summary statistics and applies to
29 any phylodynamic model. Our method enables both model selection and estimation of epidemiological parameters.
30 We demonstrate its speed and accuracy on simulated data, where it performs better than the state-of-the-art methods.
31 To illustrate its applicability, we assess the dynamics induced by superspreading individuals in an HIV dataset of men-
32 having-sex-with-men in Zurich.

33 **KEYWORDS**

34 deep learning, phylodynamics, molecular epidemiology, phylogenetic tree representation, summary statistics, HIV,
35 computer simulations, birth-death models, superspreading

36 INTRODUCTION

37 Pathogen phylodynamics is a field combining phylogenetics and epidemiology^[1]. Viral or bacterial samples from
38 patients are sequenced and used to infer a phylogeny, which describes the pathogen's spread among patients. The tips
39 of such phylogenies represent sampled pathogens, and the internal nodes transmission events. Moreover, transmission
40 events can be dated and thereby provide hints on transmission patterns. Such information is extracted by phylodynamic
41 methods to estimate epidemiological and population dynamic parameters^[2-4], assess the impact of population
42 structure^[2,5], and reveal the origins of epidemics^[6].

43 Birth-death models^[7] incorporate easily interpretable parameters common to standard infectious-disease
44 epidemiology, such as basic reproduction number R_0 , infectious period, *etc.* In contrast to the standard epidemiological
45 models, the birth-death models can be applied to estimate parameters from phylogenetic trees^[8]. In these models,
46 births represent transmission events, while deaths represent removal events for example due to treatment or recovery.
47 Upon a patient's removal their pathogen can be sampled, producing tips in the tree.

48 Here we focus on three specific, well-established birth-death models (**Fig. 1**): birth-death model (BD)^[8,9], birth-death
49 model with exposed and infectious classes (BDEI)^[5,10,11], and birth-death model with superspreading (BDSS)^[5,12].
50 These models were deployed using BEAST2^[12,13] to study the phylodynamics of such diverse pathogens as Ebola
51 virus^[10], Influenza virus^[12], Human Immunodeficiency Virus (HIV)^[5], Zika^[14] or Coronavirus^[15]. Using these models,
52 we will demonstrate the reliability of our deep learning-based approach.

53 While a great effort has been invested in the development of new epidemiological models in phylodynamics, the field
54 has been slowed down by the mathematical complexity inherent to these models. BD, the simplest model, has a closed
55 form solution for the likelihood formula of a tree for a given set of parameters^[8,10], but more complex models, *e.g.*
56 BDEI and BDSS, rely on a set of ordinary differential equations (ODEs) that cannot be solved analytically. To estimate
57 parameter values through maximum-likelihood and Bayesian approaches, these ODEs must be approximated
58 numerically for each tree node^[5,10-12]. These calculations become difficult as the tree size increases, resulting in
59 numerical instability and inaccuracy^[12], as we will see below.

60 Inference issues with complex models are typically overcome by approximate Bayesian computation (ABC)^[16,17].
61 ABC is a simulation-based technique relying on a rejection algorithm^[18], where from a set of simulated phylogenies

62 within a given prior (values assumed for parameter values), those closest to the analyzed phylogeny are retained and
63 give the posterior distribution of the parameters. This scheme relies on the definition of a set of summary statistics
64 aimed at representing a phylogeny and on a distance measure between trees. This approach is thus sensitive to the
65 choice of the summary statistics and distance metric, *e.g.* Euclidean distance. To address this issue *Saulnier et al.*^[19]
66 developed a large set of summary statistics. In addition, they used a regression step to select the most relevant statistics
67 and to correct for the discrepancy between the simulations and the analyzed phylogenies (see also *Blum et al.*^[20]).

68 We propose a likelihood-free, rejection-free, simulation-based method relying on deep learning from millions of trees
69 of varying size simulated within a broad range of parameter values. To describe these trees and use them as input for
70 the deep learner, we develop two tree representations: (1) a large set of summary statistics mostly based on *Saulnier*
71 *et al.*^[19], and (2) a complete and compact vectorial representation of phylogenies, including both the tree topology and
72 branch lengths. The summary statistics are derived from our understanding and knowledge of the epidemiological
73 processes. However, they can be incomplete and thus miss some important aspects of the studied phylogenies, which
74 can potentially result in low accuracy during inference. Moreover, it is expected that new phylodynamic models will
75 require design of new summary statistics, as confirmed by our results with BDSS. In contrast, our vectorial
76 representation is a raw data representation that preserves all information contained in the phylogeny and thus should
77 be accurate and deployable on any new model, provided the model parameters are identifiable. Our vectorial
78 representation naturally fits with deep learning methods, especially the convolutional architectures, which have
79 already proven their ability to extract relevant features from raw representations, for example in image analysis^[21,22]
80 or weather prediction^[23].

81 In the following, we introduce our vectorial tree representation and the new summary statistics designed for BDSS.
82 We then present the deep learning architectures trained on these representations and evaluate their accuracy on
83 simulated datasets in terms of both parameter estimation and model selection. The results are compared to those of
84 the gold standard method, BEAST2^[12,13]. Lastly, we showcase our methods on an HIV dataset^[24,25] from men-having-
85 sex-with-men (MSM) community from Zurich. All technical details are provided in Methods.

86 RESULTS

87 Neural networks are trained on numerical vectors from which they can learn regression and classification tasks. We
88 trained such networks on phylogenetic trees to estimate epidemiological parameters (regression) and select
89 phylodynamic models (classification). We undertook two strategies for representing phylogenetic trees as numerical
90 vectors, which we describe first, before showing the results with simulated and real data.

91 **Summary statistics (SS) representation.** We used a set of 83 SS developed by *Saulnier et al.*^[19]: 26 measures of
92 branch lengths, such as median of both internal and tip branch lengths; 8 measures of tree topology, such as tree
93 imbalance; 9 measures on the number of lineages through time, such as time and height of its maximum; and 40
94 coordinates representing the lineage-through-time plot. To capture more information on the phylogenies generated by
95 the BDSS model, we further enriched these SS with 14 new statistics on transmission chains describing the distribution
96 of the duration between consecutive transmissions (internal tree nodes). Our SS are diverse, complementary and
97 somewhat redundant. We used feed-forward neural networks (FFNN) with several hidden layers (**Fig. 2 b (i)**) that
98 select and combine relevant information from the input features. In addition to SS, we provide both the tree size, *i.e.*
99 the number of tips, and the sampling probability used to generate the tree, as input to our FFNN (**Fig. 2 a (vi)**). We
100 will refer to this method as FFNN-SS.

101 **Compact vectorial tree representation.** While converting raw information in a form of a phylogenetic tree into a set
102 of SS, information loss is unavoidable. This means not only that the tree cannot be fully reconstructed from its SS, but
103 also that depending on how much useful and relevant information is contained in the SS, the neural network may fail
104 to solve the problem at hand. As an alternative strategy to SS, and to prevent information loss in the tree representation,
105 we developed a representation called ‘Compact Bijective Ladderized Vector’ (CBLV).

106 Several vectorial representations of trees based either on polynomial^[26,27], Laplacian spectrum^[28] or F matrices^[29] have
107 been developed previously. However, they represent the tree shape but not the branch lengths^[26,27], may lose
108 information on trees^[28]. In addition, some of these representations^[27,29] require vectors or matrices of quadratic size
109 with respect to the number of tips.

110 Inspired by these approaches, we designed our concise, easily computable, compact, and bijective (*i.e.* 1-to-1) tree
111 representation that applies to trees of variable size and is appropriate as machine learning input (see **Methods** for

112 details). To obtain this representation, we first ladderize the tree, that is, for each internal node, the descending subtree
113 containing the most recently sampled tip is rotated to the left, **Fig. 2 a (ii)**. This ladderization step does not change the
114 tree but facilitates learning by standardizing the input data (**Supplementary Fig. 5**). Moreover, it is consistent with
115 trees observed in real epidemiological datasets, for example Influenza, where ladder-like trees reflect selection and
116 are observed for several pathogens^[1]. Then, we perform an inorder traversal^[30] of the ladderized tree, during which
117 we collect in a vector for each visited internal node its distance to the root and for each tip its distance to the previously
118 visited internal node. In particular, the first vector entry corresponds to the tree height. This transformation of a tree
119 into a vector is bijective, as one can unambiguously reconstruct any given tree from its vectorial representation
120 (**Supplementary Fig. 1**). The vector is as compact as possible, and its size grows linearly with the number of tips. We
121 complete this vector with zeros to reach the representation length of the largest tree contained in our simulation set,
122 and we add the known sampling probability (**Fig. 2 a (v), b (i)**).

123 Bijectivity combined with ladderization facilitates the training of neural networks, which do not need to learn that
124 different representations correspond to the same tree. However, unlike our SS, this full representation does not have
125 any high-level features. In CBLV identical subtrees will have the same representation in the vector whenever the roots
126 of these subtrees have the same height, while the vector representation of the tips in such subtrees will be the same no
127 matter the height of the subtree's root. Similar subtrees will thus result in repeated patterns along the representation
128 vector. We opted for Convolutional Neural Networks (CNN), which are designed to extract information on patterns
129 in raw data. Our CNN architecture (**Fig. 2 b (ii)**) includes several convolutional layers that perform feature extraction,
130 and maximum and average pooling that select relevant features and keep reasonable dimensions of feature maps. The
131 output of the CNN is then fed into an FFNN that combines the patterns found in the input to perform predictions. In
132 the rest of the manuscript, we refer to this method as CNN-CBLV.

133 **Simulated datasets**

134 For each phylodynamic model (BD, BDEI, BDSS), we simulated 4 million trees, covering a large range of values for
135 each parameter of epidemiological interest (R_0 , infectious period: $1/\gamma$, incubation period: $1/\varepsilon$, the fraction at
136 equilibrium of superspreading individuals: f_{SS} , and the superspreading transmission ratio: X_{SS}). Of the 4 million trees,
137 3.99 million were used as a training set, and 10,000 as a validation set for early stopping in the training phase^[31].
138 Additionally, we simulated another 10,000 trees, which we used as a testing set, out of which 100 were also evaluated

139 with the gold standard methods, BEAST2 and TreePar, which are more time consuming. Another 1 million trees were
140 used to define confidence intervals for estimated parameters. For BD and BDEI we considered two settings: one with
141 small trees (50 to 199 tips, in **Supplementary Fig. 2**) and a second one with large trees (200 to 500 tips, **Fig. 3**). For
142 BDSS, we considered only the setting with large trees, as the superspreading individuals are at low fraction and cannot
143 be detected in small trees (results not shown).

144 To increase the generality of our approach and avoid the arbitrary choice of the time scale (one unit can be a day, a
145 week, or a year), we rescaled all trees and corresponding epidemiological parameters, such that the average branch
146 length in a tree was equal to 1. After inference, we rescaled the estimated parameter values back to the original time
147 scale.

148 **Neural networks yield more accurate parameter estimates than gold standard methods**

149 We compared accuracy of parameter estimates yielded by our deep learning methods and those yielded by two state-
150 of-the-art phylodynamics inference tools, BEAST2^[12,13] and TreePar^[5]. The comparison shows that our deep learning
151 methods trained with SS and CBLV are either comparable (BD) or more accurate (BDEI and BDSS) than the state-
152 of-the-art inference methods (**Fig. 3**). The simple BD model has a closed form solution for the likelihood function,
153 and thus BEAST2 results are optimal in theory^[8,9]. Our results with BD are similar to those obtained with BEAST2,
154 and thus nearly optimal as well. For BDEI and BDSS our results are more accurate than BEAST2, which is likely
155 explained by numerical approximations of likelihood calculations in BEAST2^[5,10,11] for these models. These
156 approximations may lead to a lack of convergence (2% cases for BDEI and 15% cases for BDSS) or a convergence to
157 local minima. We suspect BEAST2 of converging to local optima when it converged to values with high relative error
158 (*i.e.* >1.0 ; 8% cases for BDEI and 11% cases for BDSS, **Fig. 3 b-c**). Furthermore, our deep learning approaches
159 showed a lower bias in parameter estimation than BEAST2 (**Supplementary Table 3**). As expected, both approaches,
160 FFNN-SS and CNN-CBLV, get more accurate with larger trees (**Supplementary Fig. 7**).

161 We tried to perform maximum likelihood estimation (MLE) implemented in the TreePar package^[5] on the same trees
162 as well. While MLE under BD model on simulations yielded as accurate results as BEAST2, for more complex models
163 it showed overflow and underflow issues (*i.e.* reaching infinite values of likelihood) and yielded inaccurate results,
164 such as more complex models (BDEI, BDSS) having lower likelihood than a simpler, nested one (BD) for a part of

165 simulations (results not shown). These issues were more prominent for larger trees. TreePar developers confirmed
166 these limitations and suggested using the last version of BEAST2 instead.

167 **CNN-CBLV has high potential for application to new models**

168 FFNN-SS and CNN-CBLV show similar accuracy across various settings (**Fig. 3, Supplementary Fig. 2,**
169 **Supplementary Table 2**). The advantage of the CBLV is its generality, meaning there is no loss of information
170 between the tree and its representation in CBLV regardless of which model the tree was generated under. This contrasts
171 with the pre-designed SS, which might need additional specific SS depending on the model at hand. This was
172 confirmed in our analyses of BDSS simulations. To estimate the parameters of this model, we added summary statistics
173 on transmission chains on top of the SS taken from *Saulnier et al.*^[19]. This improved the accuracy of superspreading
174 fraction estimates of the FFNN-SS (**Supplementary Fig. 6**), so that it was comparable to the CNN-CBLV.

175 **Neural networks are fast inference methods**

176 We compared the CPU time required by each of our inference approaches. For deep learning methods, while they
177 require longer computing times in the learning phase (*e.g.* in the order of a couple of BEAST2 inferences on large
178 trees under BDEI that we performed), inference is almost instantaneous and most of the time is spent on tree encoding
179 into SS or CBLV. Averaged over 10,000 trees generated by the three birth-death models, the parameter inference with
180 deep learning takes 6×10^{-5} CPU hours or ~ 0.2 seconds per tree.

181 For comparison, BEAST2 inference under the BD model with 5 million MCMC steps takes on average 0.18 CPU
182 hours per tree, and 0.009 CPU hours if only the time to convergence is considered. Inference under BDEI and BDSS
183 with 10 million MCMC steps takes 57 CPU hours (6 CPU hours to convergence) and 79 CPU hours (6 CPU hours to
184 convergence), respectively.

185 **Neural networks are accurate methods for model selection**

186 We trained CNN-CBLV and FFNN-SS on simulated trees to predict the birth-death model under which they were
187 simulated (BD or BDEI for small trees; BD, BDEI or BDSS for large trees). Note that for parameters shared between
188 multiple models, we used identical parameter value ranges across all these models (**Supplementary Table 1**). Then,
189 we assessed the accuracy of both of our approaches on 100 simulations obtained with each model and compared it

190 with the model selection under BEAST2 based on Akaike information criterion through Markov Chain Monte Carlo
191 (AICM)^[32,33]. The AICM, similar to deviance information criterion (DIC) by *Gelman et al.*^[32], does not add
192 computational load and is based on the average and variance of posterior log-likelihoods along the Markov Chain
193 Monte Carlo (MCMC).

194 FFNN-SS and CNN-CBLV have similar accuracy (**Supplementary Fig. 3**), namely 92% for large trees (BD vs BDEI
195 vs BDSS), and accuracy of 91% and 90%, respectively, for small trees (BD vs BDEI). BEAST2 yielded an accuracy
196 of 91% for large trees and 88% for small trees. The non-converging simulations were not considered for any of these
197 methods, *i.e.* 5% simulations for small trees and 24% for large trees.

198 The process of model selection with a neural network is as fast as the parameter inference, *i.e.* $6 \cdot 10^{-5}$ CPU hours (~ 0.2
199 seconds) per tree. This represents a practical, fast and accurate way to perform model selection in phylodynamics.

200 **Neural networks are well suited to learn complex models**

201 To assess the complexity of learned models, we explored other inference methods, namely: 1] linear regression as a
202 baseline model trained on summary statistics (LR-SS); 2] FFNN trained directly on CBLV (FFNN-CBLV); 3] CNN
203 trained on Compact Random Vector (CNN-CRV), for which the trees were randomly rotated, instead of being
204 ladderized as in **Fig. 2 (ii)**; and 4] two “null models”.

205 LR-SS yielded inaccurate results even for the BD model (**Supplementary Table 2**), which seems to contrast with
206 previous findings^[19], where LR approach combined with ABC performed only slightly worse than BEAST2. This can
207 be explained by the lack of rejection step in LR-SS, which enables to locally reduce the complexity of the relation
208 between the representation and the inferred values to a linear one^[18]. However, the rejection step requires a metric,
209 *e.g.* the Euclidean distance, which may or may not be appropriate depending on the model and the summary statistics.
210 Moreover, rejection has a computational cost with large simulation sets.

211 Neural networks circumvent these problems with rejection and allow for more complex than linear relation between
212 the tree representation and the inferred values to be captured. This also reflected in our results with FFNN-CBLV and
213 CNN-CRV, which both proved to be generally more accurate than LR-SS. However, FFNN-CBLV was substantially
214 less accurate than CNN-CBLV (**Supplementary Table 2, Supplementary Fig. 5**). This indicates the presence of
215 repeated patterns that may appear all along the vectorial representation of trees, such as subtrees of any size, which

216 are better extracted by CNN than by FFNN. In its turn, CNN-CRV required larger training sets to reach an accuracy
217 comparable to CNN-CBLV (**Supplementary Fig. 5**), showing that the ladderization and bijectivity of the CBLV
218 helped the training.

219 To assess how much information is actually learned, we also measured the accuracy of two “null models”: FFNN
220 trained to predict randomly permuted target values; and a random predictor, where parameter values were sampled
221 from prior distributions. Results show that the neural networks extract a considerable amount of information for most
222 of the estimated parameters (**Supplementary Table 2**). The most difficult parameter to estimate was the fraction of
223 superspreading individuals in BDSS model, with accuracy close to random predictions with small trees, but better
224 performance as the tree size increases (**Supplementary Fig. 7**).

225 **Showcase study of HIV in MSM subpopulation in Zurich**

226 The Swiss HIV Cohort is densely sampled, including more than 16,000 infected individuals^[24]. Datasets extracted
227 from this cohort have often been studied in phylodynamics^[8,25]. We studied a dataset of MSM subpopulation from
228 Zurich, which corresponds to a cluster of 200 sequences studied previously by *Rasmussen et al.*^[25], who focused on
229 the degree of connectivity and its impact on transmission between infected individuals. Using coalescent approaches,
230 they detected the presence of highly connected individuals at the beginning of the epidemic and estimated R_0 to be
231 between 1.0 and 2.5. We used their tree as input for neural networks and BEAST2.

232 To perform analyses, one needs an estimate of the sampling probability. We considered that: 1] the cohort is expected
233 to include around 45% of Swiss individuals infected with HIV^[24]; and 2] the sequences were collected from around
234 56% of individuals enrolled in this cohort^[34]. We used these percentages to obtain an approximation of sampling
235 probability of $0.45 \cdot 0.56 \sim 0.25$ and used this value to analyze the MSM cluster. To check the robustness of our
236 estimates, we also used sampling probabilities of 0.2 and 0.3 in our estimation procedures.

237 First, we performed a quick sanity check considering the resemblance of HIV phylogeny with simulations obtained
238 with each model. All three considered birth-death models passed this check (**Supplementary Fig. 8**). Then, we
239 performed model selection (BD vs BDEI vs BDSS) and parameter estimation using our two methods and BEAST2
240 (**Fig. 4 a-b**). Finally, we checked the model adequacy with a second, more specific, sanity check, derived from the
241 inferred values (**Fig. 4 c**).

242 Model selection with CNN-CBLV and FFNN-SS resulted in the acceptance of BDSS (probability of 1.00 versus 0.00
243 for BD and BDEI), and the same result was obtained with BEAST2 and AICM. These results are consistent with what
244 is known about HIV epidemiology, namely, the presence of superspreading individuals in the infected
245 subpopulation^[35] and the absence of incubation period without infectiousness such as is emulated in BDEI^[36].

246 We then inferred parameter values under the selected BDSS model (**Fig. 4 a-b**). The values obtained with FFNN-SS
247 and CNN-CBLV are close to each other and the 95% CI are nearly identical. We inferred an R_0 of 1.6 and 1.7, and an
248 infectious period of 10.2 and 9.8 years, with FFNN-SS and CNN-CBLV, respectively. Transmission by
249 superspreading individuals was estimated to be around 9 times higher than by normal spreaders and superspreading
250 individuals were estimated to account for around 7-8% of the population. Our R_0 estimates are consistent with the
251 results of a previous study^[8] performed on data from the Swiss cohort, and the results of *Rasmussen et al.*^[25] with this
252 MSM dataset^[25]. The infectious period we inferred is a bit longer than that reported by *Stadler et al.*, who estimated
253 it to be 7.74 [95% CI 4.39-10.99] years^[8]. The infectious period is a multifactorial parameter depending on treatment
254 efficacy and adherence, the times from infection to detection and to the start of treatment, *etc.* In contrast to the study
255 by *Stadler et al.*, whose data were sampled in the period between 1998 and 2008, our dataset covers also the period
256 between 2008 and 2014, during which life expectancy of patients with HIV was further extended^[37]. This may explain
257 why we find a longer infectious period (with compatible CIs). Lastly, our findings regarding superspreading are in
258 accordance with those of *Rasmussen et al.*^[25], and with a similar study in Latvia^[5] based on 40 MSM sequences
259 analyzed using a likelihood approach. Although the results of the latter study may not be very accurate due to the
260 small dataset size, they still agree with ours, giving an estimate of a superspreading transmission ratio of 9, and 5.6%
261 of superspreading individuals. Our estimates were quite robust to the choice of sampling probability, *e.g.* $R_0 = 1.54$,
262 1.60 and 1.66, with a sampling probability of 0.20, 0.25 and 0.30, respectively (**Fig. 4b**).

263 Compared to BEAST2, the estimates of the infectious period and R_0 were similar for both approaches, but BEAST2
264 estimates were higher for the transmission ratio (14.5) and the superspreading fraction (10.6%). These values are in
265 accordance with the bias of BEAST2 estimates that we observed in our simulation study, that is, positive bias for both
266 the superspreading fraction and the transmission ratio, while our estimates were nearly unbiased (**Supplementary**
267 **Table 3**).

268 Finally, we checked the adequacy of BDSS model by resemblance of HIV phylogeny to simulations. Using inferred
269 95% CI, we simulated 10,000 trees and performed Principal Component Analysis on SS, to which we projected the
270 SS of our HIV phylogeny. This was close to simulations, specifically close to the densest swarm of simulations,
271 confirming adequacy of both the inferred values and the selected model (**Fig. 4 c**).

272 **PhyloDeep: Python package for parameter inference and model selection**

273 FFNN-SS and CNN-CBLV parameter inference, model selection, 95% CI computation and *a priori* check are
274 available via Python package PhyloDeep at <https://pypi.org>. They cover the parameter subspace as described in
275 **Supplementary Table 1**. Version 0.2.5 is used for all results in this article. The input is a dated phylogenetic tree with
276 50-500 tips and presumed sampling probability. The output is a PCA plot for *a priori* check, and a csv file with
277 probabilities for each model (for model selection) and point estimates and 95% CI values (for parameter inference).

278 **PERSPECTIVES**

279 In this manuscript, we presented new methods for parameter inference and model selection in phylodynamics based
280 on deep learning from phylogenies. Using extensive simulations, we established that these methods are at least as
281 accurate as the state-of-the-art methods. The main advantage of our approaches is the ease of deployment on new
282 models and computational speed once the networks are trained. We also applied our deep learning methods to the
283 Swiss HIV dataset from MSM and obtained results consistent with current knowledge of HIV epidemiology.

284 Using BEAST2, we obtained inaccurate results for some of the BDEI and BDSS simulations. While BEAST2 has
285 been successfully deployed on many models and tasks, it clearly suffers from approximations in likelihood
286 computation with these two models. However, these will likely improve in near future. In fact, we already witnessed
287 substantial improvements done by BEAST2 developers to the BDSS model, while carrying out this research.

288 Both of our neural network approaches circumvent likelihood computation and thereby represent a new way of using
289 molecular data in epidemiology, without the need to solve large systems of differential equations. This opens the door
290 to novel phylodynamics models, which would make it possible to answer questions previously too complex to ask.
291 This is especially true for CBLV representation, which does not require the design of new summary statistics, when
292 applied to trees generated by a new mathematical model.

293 A direction of further research is to extend our approach to the family of phylodynamic models based on structured
294 coalescent^[38,39], rather than birth-death models. Our methods could also be extended to the macroevolutionary birth-
295 death models, which are used to study the diversification dynamics of species, and are closely related to
296 epidemiological models^[40]. Other fields related to phylodynamics, such as population genetics, have been developing
297 likelihood-free methods^[41], for which our approach might serve as a source of inspiration, too.

298 Other advantages of the deep learning approaches are that they yield close to immediate estimates and apply to trees
299 of varying size. Collection of pathogen genetic data became standard in many countries, resulting in densely sampled
300 infected populations. Examples of such datasets include HIV in Switzerland and UK^[24,42], 2013 Ebola epidemics^[6],
301 several Influenza epidemics and the 2019 SARS-Cov-2 pandemic (www.gisaid.org)^[43]. For many such pathogens,
302 trees can be efficiently and accurately inferred^[44-46] and dated^[47-49] using standard approaches. When applied to such
303 dated trees, our methods can perform model selection and provide accurate phylodynamic parameter estimates within

304 a fraction of a second. Such properties are desirable for phylogeny-based real-time outbreak surveillance methods,
305 which must be able to cope with the daily influx of new samples, and thus increasing size of phylogenies, as the
306 epidemic unfolds, in order to study local outbreaks and clusters, and assess and compare the efficiency of healthcare
307 policies deployed in parallel.

308 **ACKNOWLEDGEMENT:** We would like to thank Dr Kary Ocaña and Tristan Dot for initiating experiments on
309 machine learning and phylogenetic trees in our laboratory. We would like to thank Quang Tru Huynh for
310 administrating GPU farm at Institut Pasteur and the INCEPTION program (Investissement d’Avenir grant ANR-16-
311 CONV-0005) that financed the GPU farm. We would like to thank Dr Christophe Zimmer from Institut Pasteur,
312 Sophia Lambert and Dr H el ene Morlon from Institut de Biologie de l’Ecole Normale Sup erieure IBENS and Dr Guy
313 Baele from Katholieke Universiteit KU Leuven for useful discussions and Dr Isaac Overcast from IBENS for critical
314 reading of the manuscript. We would like to thank Dr Tanja Stadler and J er emie Scir e for their help with BEAST2
315 and MLE approaches. JV is supported by Ecole Normale Sup erieure Paris-Saclay and by ED Fronti eres de l’Innovation
316 en Recherche et Education, Programme Bettencourt. VB would like to thank Swiss National Science Foundation for
317 funding (Early PostDoc mobility grant P2EZP3_184543). OG is supported by PRAIRIE (ANR-19-P3IA-0001).

318 **AUTHOR CONTRIBUTIONS:** JV, AZ and OG conceived and set up the methods; JV and VB performed the
319 experiments; JV and OG analyzed the results; JV and AZ wrote the Python package; JV and OG wrote the manuscript;
320 JV, AZ VB and OG edited the manuscript; OG initiated and supervised the project; all authors helped in this research,
321 discussed the results and read the final manuscript.

322 **COMPETING INTERESTS:** The authors have no competing interests to declare.

323 REFERENCES

- 324 1. Grenfell, B.T. et al. Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* **303**,
325 327-332 (2004).
- 326 2. Volz, E.M., Kosakovsky Pond, S.L., Ward, M.J., Leigh Brown, A.J., Frost, S.D. Phylodynamics of
327 infectious disease epidemics. *Genetics* **183**, 1421-30 (2009).
- 328 3. Drummond, A.J., Rambaut, A., Shapiro, B., Pybus, O.G. Bayesian Coalescent Inference of Past Population
329 Dynamics from Molecular Sequences. *Molecular Biology and Evolution*, **22**, 1185–1192 (2005).
- 330 4. Stadler, T. Birth–death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C
331 virus (HCV). *Proceedings of the National Academy of Sciences* **110**, 228-233 (2013)
- 332 5. Stadler, T., Bonhoeffer, S. Uncovering epidemiological dynamics in heterogeneous host populations using
333 phylogenetic methods. *Philosophical Transactions of the Royal Society B: Biological Sciences* **368**, (2013)
- 334 6. Gire, S.K. et al. Genomic surveillance elucidates Ebola virus origin and transmission during the 2014
335 outbreak. *Science* **345**, 1369-72 (2014).
- 336 7. Boskova, V., Bonhoeffer, S., Stadler, T. Inference of Epidemiological Dynamics Based on Simulated
337 Phylogenies Using Birth-Death and Coalescent Models. *PLOS Computational Biology* **10**, (2014).
- 338 8. Stadler, T. et al. Estimating the Basic Reproductive Number from Viral Sequence Data. *Mol. Biol. Evol.* **29**,
339 347–357 (2012).
- 340 9. Leventhal, G.E., Günthard, H.F., Bonhoeffer, S., Stadler, T. Using an Epidemiological Model for
341 Phylogenetic Inference Reveals Density Dependence in HIV Transmission. *Mol. Biol. Evol.* **31**, 6–17
342 (2014).
- 343 10. Stadler, T., Kuhnert, D., Rasmussen, D.A., du Plessis, L. Insights into the early epidemic spread of Ebola in
344 sierra leone provided by viral sequence data. *PLoS Curr.* **6**, (2014).
- 345 11. Kühnert, D., Stadler, T., Vaughan, T.G., Drummond, A.J. Phylodynamics with Migration: A
346 Computational Framework to Quantify Population Structure from Genomic Data. *Mol. Biol. Evol.* **33**,
347 2102-16 (2016).

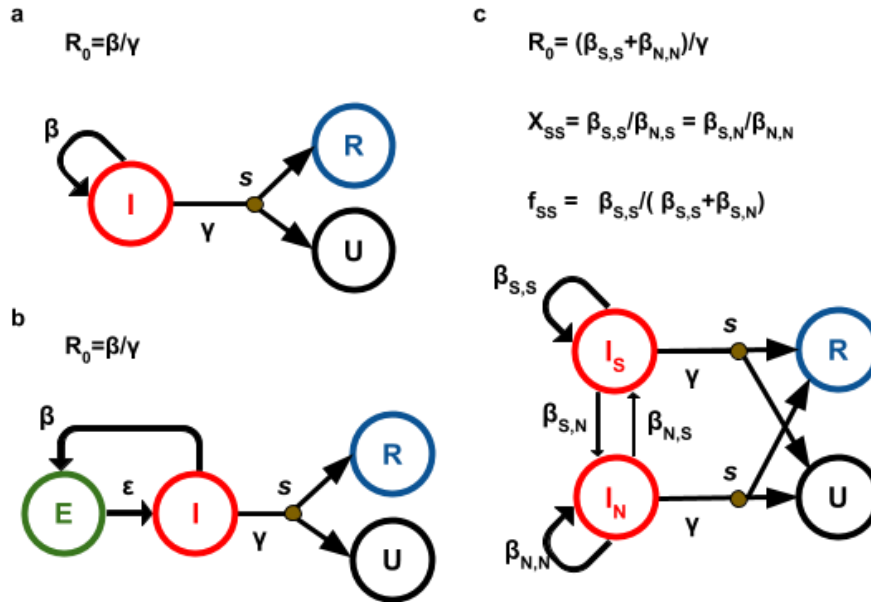
- 348 12. Sciré, J., Barido-Sottani, J., Kühnert, D., Vaughan, T.G., Stadler, T. Improved multi-type birth-death
349 phylodynamic inference in BEAST 2. Preprint at
350 <https://www.biorxiv.org/content/10.1101/2020.01.06.895532v1.full.pdf> (2020).
- 351 13. Bouckaert, R. et al. BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. *PLoS*
352 *Computational Biology* **10**, (2014).
- 353 14. Boskova, V., Stadler, T., Magnus, C. The influence of phylodynamic model specifications on parameter
354 estimates of the Zika virus epidemic. *Virus Evolution* **4**, January (2018).
- 355 15. Vaughan, T.G., Sciré, J., Nadeau, S.A., Stadler, T. Estimates of outbreak-specific SARS-CoV-2
356 epidemiological parameters from genomic data. Preprint at
357 <https://www.medrxiv.org/content/10.1101/2020.09.12.20193284v1.full.pdf> (2020).
- 358 16. Rubin, D.B. Bayesianly Justifiable and Relevant Frequency Calculations for the Applies Statistician. *The*
359 *Annals of Statistics* **12**, 1151-1172 (1984).
- 360 17. Beaumont, M.A., Zhang, W., Balding, D.J. Approximate Bayesian Computation in Population Genetics.
361 *Genetics* **164**, 2025-2035 (2002).
- 362 18. Csilléry, K., Blum, M.G.B., Gaggiotti, O.E., François, O. Approximate Bayesian Computation (ABC) in
363 practice. *Trends in Ecology & Evolution* **25**, 410-418 (2010).
- 364 19. Saulnier, E., Gascuel, O., Alizon, S. Inferring epidemiological parameters from phylogenies using
365 regression-ABC: A comparative study. *PLoS Comp. Biol.* **13**, (2017).
- 366 20. Blum, M.G.B. *Handbook Of Approximate Bayesian Computation Ch. Regression approaches for ABC*. 71–
367 85. (Chapman and Hall/CRC Press, Boca Raton, 2018).
- 368 21. LeCun, Y., Kavukcuoglu, K., Farabet, F. Convolutional networks and applications in vision. *Proc. IEEE*
369 *Int. Symp. Circuits Syst.* 253-256 (2010).
- 370 22. Krizhevsky, K., Sutskever, I., Hinton, G.E. ImageNet Classification with Deep Convolutional Neural
371 Networks. *Advances in neural information processing systems* 1097-1105 (2012).
- 372 23. Chattopadhyay, A., Hassanzadeh, P., Pasha, S. Predicting clustered weather patterns: A test case for
373 applications of convolutional neural networks to spatio-temporal climate data. *Sci. Rep.* **10**, 1317 (2020)
- 374 24. The Swiss HIV Cohort Study et al. Cohort Profile: The Swiss HIV Cohort Study. *International Journal of*
375 *Epidemiology* **39**, 1179–1189 (2010).

- 376 25. Rasmussen, D.A., Kouyos, R., Günthard, H.F., Stadler, T. Phylodynamics on local sexual contact networks.
377 *PLOS Comp. Biol.* **13**, (2017).
- 378 26. Colijn, C. & Plazzotta, G. A metric on phylogenetic tree shapes. *Systematic Biology* **67**, 113–126 (2018).
- 379 27. Liu, P., Gould, M., Colijn, C. Polynomial Phylogenetic Analysis of Tree Shapes. Preprint at
380 <https://doi.org/10.1101/2020.02.10.942367> (2020).
- 381 28. Lewitus, E. & Morlon, H. Characterizing and Comparing Phylogenies from their Laplacian Spectrum.
382 *Systematic Biology* **65**, 495-507 (2016).
- 383 29. Kim, J., Rosenberg, N.A., Palacios, J.A. Distance metrics for ranked evolutionary trees. *Proceedings of the*
384 *National Academy of Sciences* **117**, 28876-28886 (2020).
- 385 30. Cormen, T.H., Leiserson, C.E., Rivest, R.L., Stein, C. *Introduction To Algorithms*. 286-307 (The MIT
386 Press, Cambridge, 2009).
- 387 31. Bengio, Y. *Neural Networks: Tricks Of The Trade, Ch. Practical Recommendations for Gradient-Based*
388 *Training of Deep Architectures*. (Springer, Berlin, Heidelberg 2002).
- 389 32. Gelman, A., Carlin, J.B., Stern, H.S., Rubin, D.B. *Bayesian Data Analysis: Second Edition*. (Chapman and
390 Hall/CRC Press, Boca Raton, 2004).
- 391 33. Baele, G. et al. Improving the accuracy of demographic and molecular clock model comparison while
392 accommodating phylogenetic uncertainty. *Mol. Biol. Evol.* **29**, 2157-67 (2012).
- 393 34. Kouyos, R.D. et al. Molecular epidemiology reveals long-term changes in HIV type 1 subtype B
394 transmission in Switzerland. *J. Infect. Dis.* **201**, 1488-97 (2010).
- 395 35. May, R.M. & Anderson, R.M. Transmission dynamics of HIV infection. *Nature* **326**, 137–142 (1987).
- 396 36. Brenner, B.G. et al. Quebec Primary HIV Infection Study Group. High rates of forward transmission events
397 after acute/early HIV-1 infection. *J. Infect. Dis.* **195**, 951-9 (2007).
- 398 37. Gueller, A. et al. Swiss National Cohort Life expectancy in HIV-positive persons in Switzerland. *AIDS* **31**,
399 427-436 (2017).
- 400 38. Rasmussen, D.A., Volz, E.M., Koelle, K. Phylodynamic Inference for Structured Epidemiological Models.
401 *PLoS Comput. Biol.* **10**, (2014).
- 402 39. Volz, E.M. & Siveroni, I. Bayesian phylodynamic inference with complex models. *PLoS Comput. Biol.* **14**,
403 (2018).

- 404 40. MacPherson, A., Louca, S., McLaughlin, A., Joy, J.B., Pennell, M.W. A General Birth-Death-Sampling
405 Model for Epidemiology and Macroevolution. Preprint at
406 <https://www.biorxiv.org/content/10.1101/2020.10.10.334383v2> (2020).
- 407 41. Sanchez, T., Cury, J., Charpiat, G., Jay, F. Deep learning for population size history inference: Design,
408 comparison and combination with approximate Bayesian computation. *Mol. Ecol. Resour.* **00**, 1-16. (2020).
- 409 42. Dunn, D. & Pillay, D. UK HIV drug resistance database: background and recent outputs. *J. HIV Ther.* **12**,
410 97–8 (2007).
- 411 43. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to reality.
412 *Euro Surveill.* **22**, 30494 (2017).
- 413 44. Minh, B.Q. et al. IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic
414 era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 415 45. Kozlov, A.M., Darriba, D., Flouri, T., Morel, B., Stamatakis, A. RAxML-NG: a fast, scalable and user-
416 friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453–4455 (2019).
- 417 46. Guindon, S. et al. New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing
418 the Performance of PhyML 3.0. *Systematic Biology* **59**, 307-21 (2010).
- 419 47. Sagulenko, P., Puller, V., Neher, R.A. TreeTime: Maximum-likelihood phylodynamic analysis. *Virus*
420 *Evol.* **4**, (2018).
- 421 48. To, T.H., Jung, M., Lycett, S., Gascuel, O. Fast Dating Using Least-Squares Criteria and Algorithms. *Syst*
422 *Biol.* **65**, 82-97 (2016).
- 423 49. Volz, E.M. & Frost, S.D.W. Scalable relaxed clock phylogenetic dating. *Virus Evol.* **3**, (2017).

424 **FIGURE LEGENDS**

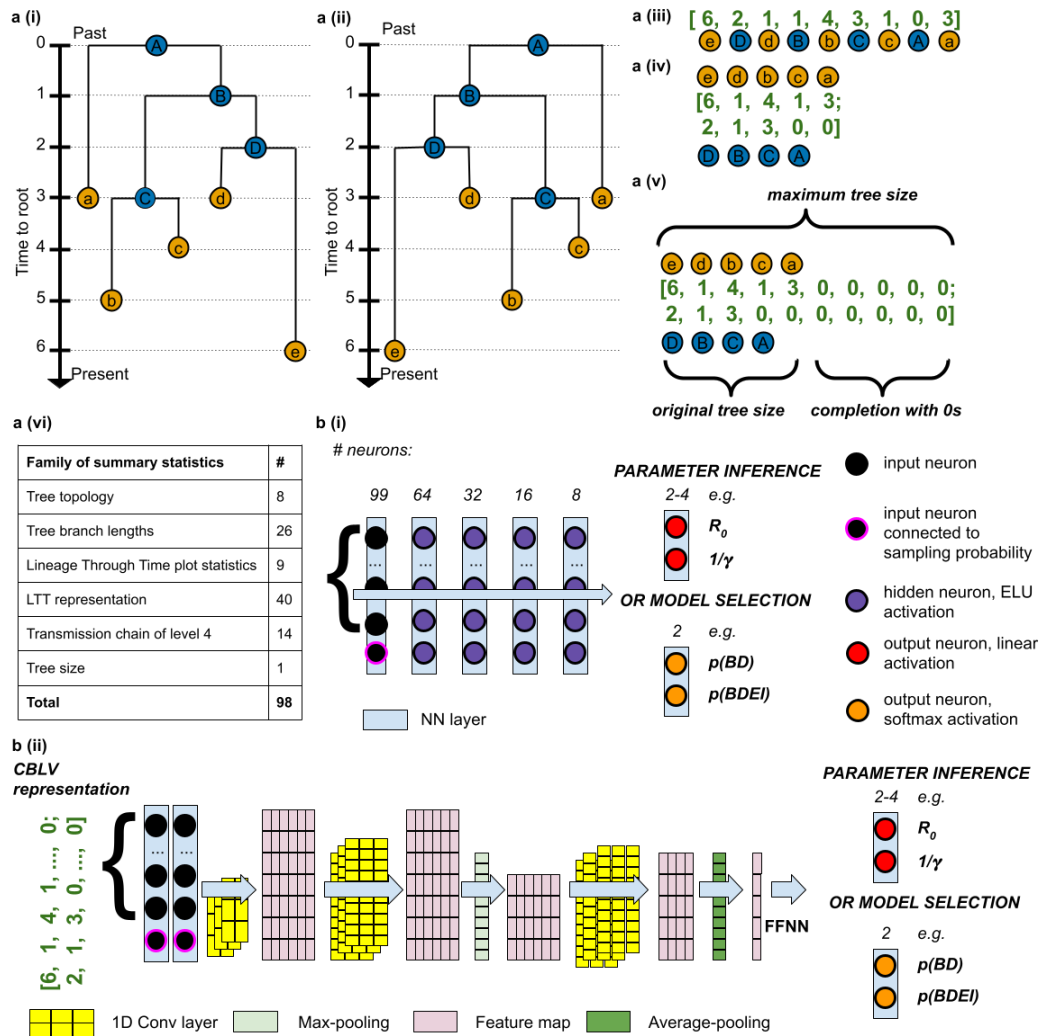
425 **Fig. 1: Birth-death models**



426

427 **a**, Birth-death model (BD)^[8,9], **b**, birth-death model with Exposed-Infectious individuals (BDEI)^[5,10,11] and **c**, birth-
 428 death model with SuperSpreading (BDSS)^[5,12]. BD is the simplest generative model, used to estimate R_0 and the
 429 infectious period ($1/\gamma$)^[8,9]. BDEI and BDSS are extended version of BD. BDEI enables to estimate latency period ($1/\epsilon$)
 430 during which individuals of exposed class E are infected, but not infectious^[5,10,11]. BDSS includes two populations
 431 with heterogeneous infectiousness: the so-called superspreading individuals (S) and normal spreaders (N).
 432 Superspreading individuals are present only at a low fraction in the population (f_{SS}) and may transmit the disease at a
 433 rate that is multiple times higher than that of normal spreaders (rate ratio = X_{SS})^[5,12]. Superspreading can have various
 434 complex causes, such as the heterogeneity of immune response, disease progression, co-infection with other diseases,
 435 social contact patterns or risk behavior, *etc.* Infectious individuals I (superspreading infectious individuals I_S and
 436 normal spreaders I_N for BDSS), transmit the disease at rate β ($\beta_{X,Y}$ for an individual of type X transmitting to an
 437 individual of type Y for BDSS), giving rise to a newly infected individual. The newly infected individual is either
 438 infectious right away in BD and BDSS or goes through an exposed state before becoming infectious at rate ϵ in BDEI.
 439 Infectious individuals are removed at rate γ . Upon removal, they can be sampled with probability s , becoming of
 440 removed sampled class R. If not sampled upon removal, they move to non-infectious unsampled class U.

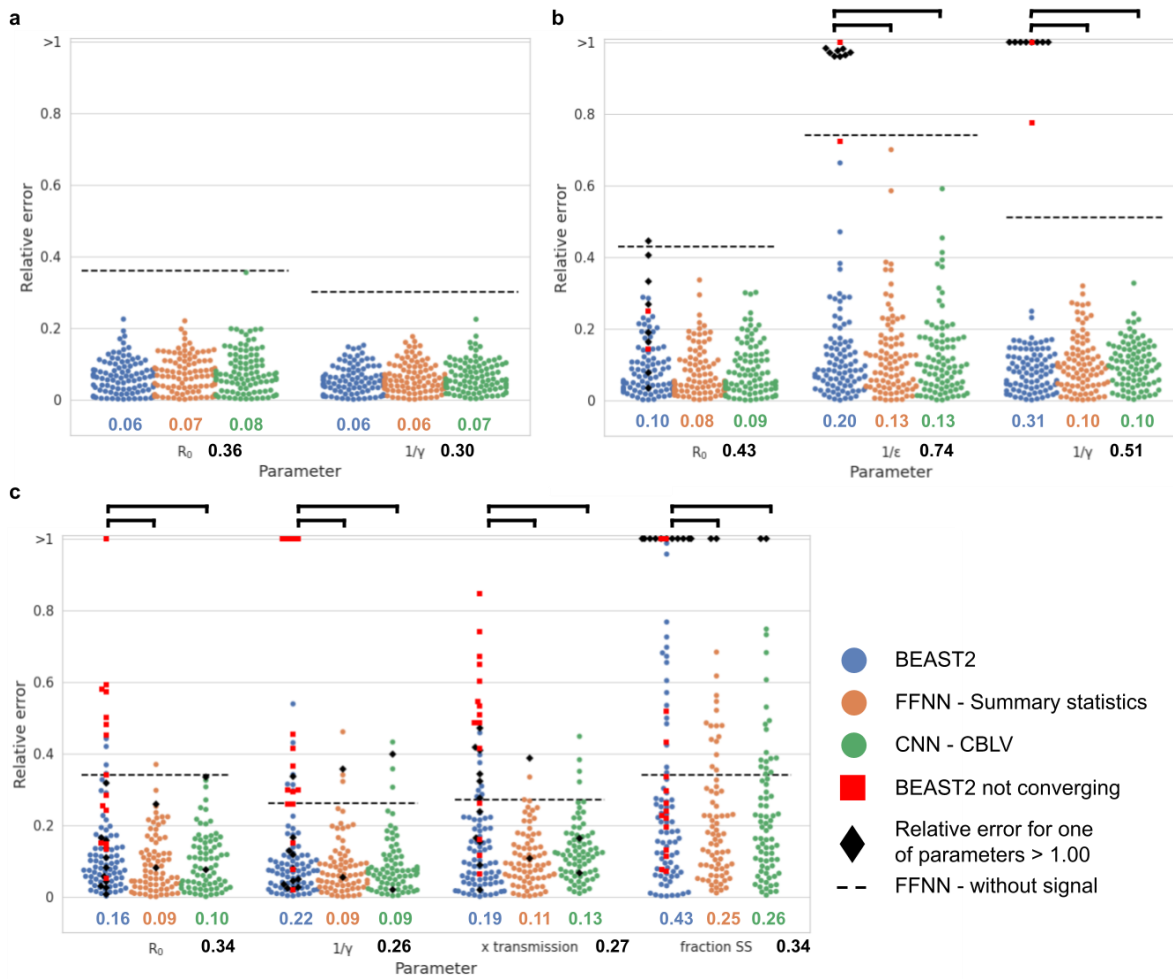
441 **Fig. 2: Pipeline for training neural networks on phylogenies**



442

443 **a. Tree representations.** **a (i)**, simulated binary trees. Under each model from Fig. 1, we simulate many trees of
 444 variable size (50 to 200 tips for ‘small trees’ and 200 to 500 tips for ‘large trees’). For illustration, we have here a tree
 445 with 5 tips. We encode the simulations into two representations, either **a (ii-v)**, in a complete and compact tree
 446 representation called ‘Compact Bijective Ladderized Vector’ abbreviated as CBLV or **a (vi)** with summary statistics
 447 (SS). CBLV is obtained through **a (ii)** ladderization or sorting of internal nodes so that the branch supporting the most
 448 recent leaf is always on the left and **a (iii)** an inorder tree traversal, during which we append to a real-valued vector
 449 for each visited internal node its distance to the root and for each visited tip its distance to the previously visited
 450 internal node. We reshape this representation into **a (iv)**, an input matrix in which the information on internal nodes
 451 and leaves is separated into two rows. Finally, **a (v)**, we complete this matrix with zeros so that the matrices for all
 452 simulations have the size of largest simulation matrices. For illustration purpose, we here consider that the maximal
 453 tree size covered by simulations is 10, and the representation is thus completed with 0s accordingly. SS consists of **a**
 454 **(vi)**, a set of 98 statistics: 83 published in *Saulnier et al.*^[19], 14 on transmission chains and 1 on tree size. The
 455 information on sampling probability is added to both representations. **b**, Neural networks are trained on these
 456 representations to estimate parameter values or to select the underlying model. For SS, we use, **b (i)**, a deep feed-
 457 forward neural network (FFNN) of funnel shape (we show the number of neurons above each layer). For the CBLV
 458 representation we train, **b (ii)**, Convolutional Neural Networks (CNN). The CNN is added on top of the FFNN. The
 459 CNN combines convolutional, maximum pooling and global average pooling layers, as described in detail in **Methods**.

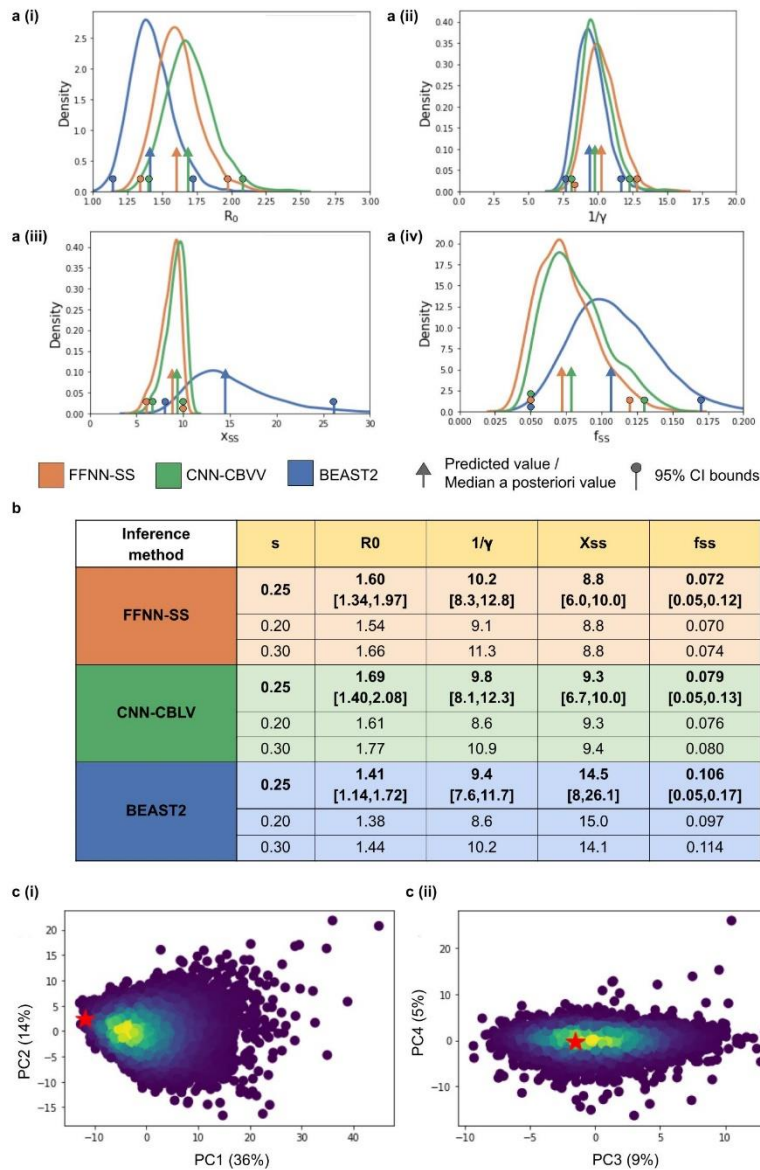
460 **Fig. 3: Assessment of deep learning accuracy**



461

462 **Comparison of inference accuracy** by BEAST2 (in blue), deep neural network trained on SS (in orange) and
 463 convolutional neural network trained on the CBLV representation (in green) on 100 test trees. The size of training and
 464 testing trees was uniformly sampled between 200 and 500 tips. We show the relative error for each test tree. The error
 465 is measured as the normalized distance between the median *a posteriori* estimate by BEAST2 or point estimates by
 466 neural networks and the target value for each parameter. We highlight simulations for which BEAST2 did not
 467 converge and whose values were thus set to median of the parameter subspace used for simulations by depicting them
 468 as red squares. We further highlight the analyses with a high relative error (>1.00) for one of the estimates as black
 469 diamonds. We compare the relative errors for **a**, BD-simulated, **b**, BDEI-simulated and **c**, BDSS-simulated trees.
 470 Average relative error is displayed for each parameter and method in corresponding color below each figure. The
 471 average error of a FFNN trained on summary statistics but with randomly permuted target is displayed as black dashed
 472 line and its value is shown in bold black below the x-axis. The accuracy of each method is compared by paired z-test;
 473 $P < 0.05$ is shown as thick full line; non-significant is not shown.

474 **Fig. 4: Parameter inference on HIV data sampled from MSM Zurich**



475

476 **Inference on Zurich HIV MSM dataset.** Using BDSS model with BEAST2 (in blue), FFNN-SS (in orange), and
 477 CNN-CBLV (in green) we infer, **a (i)**, basic reproduction number, **a (ii)**, infectious period (in years), **a (iii)**,
 478 superspreading transmission ratio and, **a (iv)**, superspreading fraction. For FFNN-SS and CNN-CBLV, we show the
 479 distributions and the 95% CIs obtained with a fast approximation of the parametric bootstrap (**Methods**). For
 480 BEAST2, the distributions and 95% CI were obtained considering all reported steps (9,000 in total) excluding the 10%
 481 burn-in. Arrows show the position of the original point estimates obtained with FFNN-SS and CNN-CBLV and the
 482 median *a posteriori* estimate obtained with BEAST2. Circles show lower and upper boundaries of 95% CI. **b**, these
 483 values are reported in a table, together with point estimates obtained while considering lower and higher sampling
 484 probability (0.20 and 0.30). **c**, 95% CI boundaries obtained with FFNN-SS are used to perform an *a posteriori* model
 485 adequacy check. We simulated 10,000 trees with BDSS while resampling each parameter from a uniform distribution,
 486 whose upper and lower bounds were defined by the 95% CI. We then encoded these trees into SS, performed PCA
 487 and projected SS obtained from the HIV MSM phylogeny (red stars) on these PCA plots. We show here the projection
 488 into **c (i)**, first two components of PCA, **c (ii)**, the 3rd and 4th components, together with the associated percentage of
 489 variance displayed in parentheses. Warm colors correspond to high density of simulations.

490	METHODS	
491	TABLE OF CONTENTS	
492		Page
493	TREE REPRESENTATION USING SUMMARY STATISTICS (SS)	26
494	<i>Saulnier et al.'s summary statistics</i>	26
495	<i>Additional summary statistics</i>	26
496	COMPLETE AND COMPACT TREE REPRESENTATION (CBLV)	27
497	<i>Tree ladderization</i>	27
498	<i>Tree traversal and encoding</i>	28
499	<i>Properties of CBLV</i>	28
500	<i>Alternative tree representations</i>	29
501	TREE RESCALING	29
502	REDUCTION AND CENTERING OF SUMMARY STATISTICS REPRESENTATION	30
503	PARAMETER INFERENCE USING NEURAL NETWORKS	30
504	<i>Deep feedforward neural network architecture for SS</i>	30
505	<i>Deep convolutional neural network for CBLV</i>	31
506	<i>Neural network setting and training</i>	31
507	<i>Preventing overfitting: Early stopping and Dropout</i>	31
508	<i>Neural networks for model selection</i>	31
509	CONFIDENCE INTERVALS (95% CI)	32
510	<i>Computation of 95% CI</i>	32
511	<i>Assessment of CI accuracy and width</i>	32
512	MODEL ADEQUACY	33
513	<i>A priori check</i>	33
514	<i>A posteriori check</i>	34

515	MODELS	34
516	<i>Constant rate birth-death model with incomplete sampling</i>	34
517	<i>Birth-death model with exposed-infectious classes</i>	34
518	<i>Birth-death model with superspreading</i>	35
519	SIMULATIONS	35
520	METHOD COMPARISON	36
521	<i>Parameter inference with BEAST2</i>	36
522	<i>Model selection with BEAST2</i>	38
523	<i>Linear regression</i>	38
524	<i>FFNN-CBLV</i>	39
525	<i>TreePar</i>	39
526	<i>Null models</i>	39
527	PERFORMANCE ASSESSMENT	40
528	<i>Mean relative error MRE</i>	40
529	<i>Mean relative bias MRB</i>	41
530	<i>Model selection accuracy</i>	41
531	<i>Comparison of time efficiency</i>	41
532	HIV DATASET	42
533	ADDITIONAL REFERENCES	42

534 **TREE REPRESENTATION USING SUMMARY STATISTICS (SS)**

535 We use a set of 98 summary statistics (SS), to which we add the sampling probability, summing to a vector of 99
536 values.

537 *Saulnier et al. summary statistics*

538 We use the 83 SS proposed by *Saulnier et al.*^[19]:

- 539 • 8 SS on tree topology
- 540 • 26 SS on branch lengths
- 541 • 9 SS on Lineage-Through-Time (LTT) plot
- 542 • 40 SS providing the coordinates of the LTT plot

543 The computing time of these statistics grows linearly with tree size. For details, see the original paper.

544 *Additional summary statistics*

545 In addition to *Saulnier et al.*^[19] statistics, we designed 14 SS on transmission chains. Moreover, we provide the number
546 of tips in the tree as input resulting in $83+14+1 = 98$ SS in total.

547 The statistics on transmission chains are designed to capture information on the superspreading population. A
548 superspreading individual transmits to more individuals within a given time period than a normal spreader. We thus
549 expect that with superspreading individuals we would have shorter transmission chains. To have a proxy for the
550 transmission chain length, we look at the sum of 4 subsequent shortest times of transmission for each internal node.
551 This gives us a distribution of time-durations of 4-transmission chains. We assume that information on the
552 transmission dynamics of superspreading individuals is retained in the lower, *i.e.* left, tail of 4-transmission-chain
553 lengths distribution which contains relatively many transmissions with short time to next transmission), while the
554 information on normal spreaders should be present in the rest of the distribution.

555 The implementation of this 4-transmission-chain SS is the following. For each internal node, we sum the distances
556 from the internal node to its closest descendant nodes, descending exactly four times, that is, we take first the distance
557 from the given internal node to its closest child node (of level 1), then from the (level 1) child node, we take its distance
558 to its own closest child node (of level 2), *etc.* If one of the closest descendant nodes is a tip (except for the last one in

559 the chain), we do not retain any value for the given internal node. Other options, like the shortest 4-edge pathway,
560 could have been used as well and would likely give comparable results.

561 On the obtained distribution of 4-transmission-chain lengths, we compute 14 statistics:

- 562 • number of 4-transmission chains in the tree
- 563 • 9 deciles of 4-transmission-chain lengths distribution
- 564 • minimum and maximum values of 4-transmission-chain lengths distribution
- 565 • mean value of 4-transmission-chain lengths
- 566 • variance of 4-transmission-chain lengths

567 Adding the same summary statistics but on chains comprising 2, 3 and 5 consecutive transmissions had a negligible
568 impact on parameter inference accuracy (data not shown).

569 **COMPLETE AND COMPACT TREE REPRESENTATION (CBLV)**

570 Simulated dated trees are encoded in the form of real-valued vectors, which are then used as input for the neural
571 networks. The representation of a tree with n tips is a vector of length $2n-1$, where one single real-valued scalar
572 corresponds to one internal node or tip. This representation thus scales linearly with the tree size. The encoding is
573 achieved in two steps: tree ladderization and tree traversal.

574 *Tree ladderization*

575 The tree ladderization consists of ordering each node's children. Child nodes are sorted based on the sampling time of
576 the most recently sampled tip in their subtrees: for each node, the branch supporting the most recently sampled subtree
577 is rotated to the left, as in **Fig. 2 a (i-ii)**.

578 We considered several alternatives with different criteria for child (subtree) sorting instead of ladderization: sampling
579 time of the most anciently sampled tip, subtree length (*i.e.* sum of all branch lengths including the rooting branch),
580 diversification (*i.e.* number of tips), normalized branch lengths (*i.e.* subtree length divided by the number of tips), *etc.*
581 These did not yield better results than CBLV. We show in **Supplementary Fig. 5** the comparison of CBLV with
582 Compact Random Vector (CRV), for which internal nodes were sorted randomly before the tree traversal, showing
583 that CRV yields poorer results than CBLV, as expected.

584 ***Tree traversal and encoding***

585 Once the tree is sorted, we perform an inorder tree traversal, using a standard recursive algorithm from the depth first
586 family^[30]. When visiting a tip, we add its distance to the previously visited internal node or its distance to the root, for
587 the tip that is visited first (*i.e.* the tree height due to ladderization). When visiting an internal node, we add its distance
588 to the root. Examples of encoding are shown in **Fig. 2 a (ii-iii)**. This gives us the Compact Bijective Ladderized Vector
589 (CBLV). We then separate information relative to tips and to internal nodes into two rows (**Fig. 2 a (iv)**) and complete
590 the representation with zeros until reaching the size of the largest simulated tree for the given simulation set (**Fig. 2 a**
591 **(v)**).

592 ***Properties of CBLV***

593 CBLV has favorable features for deep learning. Ladderization does not actually change the input tree (phylogenies
594 are unordered trees), but by ordering the subtrees it standardizes the input data and facilitates the learning phase, as
595 observed with CRV (**Supplementary Fig. 5**). Then, the inorder tree traversal procedure is a bijective transformation,
596 as it transforms a tree into a vector, from which the (ordered) tree can be reconstructed unambiguously, using a simple
597 path-agglomeration algorithm shown in **Supplementary Fig. 1**. CBLV is “as concise as possible” being composed of
598 $2n-1$ real values **Fig. 2 a (iii)**, where n is the number of tips. A rooted tree has $2n-2$ branches, and thus $2n-2$ entries are
599 needed to represent the branch lengths. In our $2n-1$ vectorial encoding of trees, we not only represent the branch
600 lengths, but also the tree topology using only 1 additional entry.

601 The compactness and bijectivity of tree representation reduce the number of simulations required for training the
602 neural network (**Supplementary Fig. 5**). This is because the number of parameters to be trained remains reasonable
603 with compact representation. Moreover, the networks do not need to learn that several different inputs correspond to
604 the same tree.

605 Our neural networks are intended to apply to trees of variable sizes, *e.g.* trees of 200 to 500 tips in our experiments
606 with large trees. Thus, they are trained on representations of different lengths (*e.g.* a vector of length 399 for a tree of
607 200 tips), that we complete with zeroes to reach the length of the largest trees (*i.e.* 999 for 500 tips). We add an
608 additional zero to obtain a two-row matrix ($500*2$ for 500 tips).

609 *Alternative tree representations*

610 Our CBLV tree representation could likely be improved to ease the learning phase and obtain even better parameter
611 estimates. We tested several alternative representations, some inspired by the polynomial representation of small
612 subtrees^[26,27], the Laplacian spectrum^[28] and additive distance matrices that are equivalent to trees^[50]. None was by
613 far as convincing as CBLV, which is likely due to their large size (n^2 for distance matrices and polynomials) or
614 numerical instabilities and potential loss of information (for Laplacian spectrum).

615 Moreover, the margin for improvement of the accuracy of CNN-CBLV for the BD model, and likely for other models,
616 is low. This is due to the observation that the accuracy of CNN-CBLV is the similar to that of likelihood-based
617 approaches for the BD model, and the fact that we have an analytical likelihood formula for the BD model, making
618 the likelihood-based approach itself optimal^[8,9].

619 **TREE RESCALING**

620 Before encoding, the trees are rescaled so that the average branch length is 1, that is, each branch length is divided by
621 the average branch length of the given tree, called rescale factor. The values of the corresponding time-dependent
622 parameters, *i.e.* infectious period and incubation period, are divided by the rescale factor too. The NN is then trained
623 to predict these rescaled values. After parameter prediction, the predicted parameter values are multiplied by the
624 rescale factor and thus rescaled back to the original time scale.

625 This step enables us to overcome problems of arbitrary time scales of input trees and makes a pre-trained NN more
626 generally applicable. More specifically, an input tree with a time scale in days will be associated naturally with the
627 same output as the same tree with a time scale in years, since both these trees will be rescaled to the same intermediate
628 tree of average branch length of 1. Rescaling thus makes it possible to apply the same pre-trained NN to phylogenies
629 reconstructed from sequences of a pathogen associated with an infectious period on the scale of days (*e.g.* EboV) or
630 years (*e.g.* HIV), without the need to simulate new phylogenies and train a new NN.

631 **REDUCTION AND CENTERING OF SUMMARY STATISTICS REPRESENTATION**

632 Before training our NN and after having rescaled the trees to unit average branch length (see the sub-section above),
633 we reduce and center every summary statistic by subtracting the mean and scaling to unit variance. To achieve this,
634 we use the standard scaler from the scikit-learn package^[51], which is fitted to the training set.

635 **PARAMETER AND MODEL INFERENCE USING NEURAL NETWORKS**

636 We implemented deep learning methods in Python 3.6 using Tensorflow 1.5.0^[52], Keras 2.2.4^[53] and scikit-learn
637 0.19.1^[51] libraries. For each network, several variants in terms of number of layers and neurons, activation functions,
638 regularization, loss functions and optimizer, were tested. In the end, we decided for two specific architectures that best
639 fit our purpose: one deep FFNN trained on SS and one CNN trained on CBLV tree representation.

640 *Deep feedforward neural network architecture for SS*

641 The network consists of one input layer (of 99 input nodes both for trees with 50-199 and 200-500 tips), 4 sequential
642 hidden layers organized in a funnel shape with 64-32-16-8 neurons and 1 output layer of size 2-4 depending on the
643 number of parameters to be estimated. The neurons of the last hidden layer have linear activation, while others have
644 exponential linear activation^[54].

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 64)	6400
dense_2 (Dense)	(None, 32)	2080
dense_3 (Dense)	(None, 16)	528
dense_4 (Dense)	(None, 8)	136
dense_5 (Dense)	(None, 2)	18
Total params: 9,162		
Trainable params: 9,162		
Non-trainable params: 0		

645

646 **Architecture: Feedforward neural network architecture.** Example of FFNN trained on large trees to estimate the parameters
647 of the BD model (R_0 and infectious period $1/\gamma$). ‘Dense’ layer means that for each neuron, all the inputs are multiplied by
648 learned weights, summed together with the bias term. The activation function is then applied to the weighted sum before being
649 output to the next layer. Dense_1 to dense_4 are layers with neurons of exponential linear activation, while dense_5 is composed
650 either of softmax (in case of model selection) or linear neurons (in case of parameter estimation). The number of trainable

651 parameters in each layer is displayed (Param #): for example in the first layer, we have 99 input values and 1 bias for each of
652 the 64 neurons, giving us in total $(99+1)*64=6,400$ trainable parameters. Output by Keras^[53], the 'None' in the 'Output Shape'
653 means the network can input more than one training example at the time and that there is no constraint on the batch size (hence
654 'None').

655 *Deep convolutional neural network for CBLV*

656 The CNN consists of one input layer (of 400 and 1002 input nodes for trees with 50-199 and 200-500 tips,
657 respectively). This input is then reshaped into a matrix of size of $201*2$ and $501*2$, for small and large trees,
658 respectively, with entries corresponding to tips and internal nodes separated into separate rows (and one extra column
659 with one entry in each row corresponding to the sampling probability). Then, there are two 1D convolutional layers
660 of 50 kernels each, of size 3 and 10, respectively, followed by max pooling of size 10 and another 1D convolutional
661 layer of 80 kernels of size 10. After the last convolutional layer, there is a GlobalPoolingAverage1D layer and a FFNN
662 of funnel shape (64-32-16-8 neurons) with the same architecture and setting as the NN used with SS.

663 *Neural network setting and training*

664 For both NNs, we use the Adam optimisation algorithm^[55] as optimizer and the Mean Absolute Percentage Error
665 (MAPE) as loss function. The batch size is set to 8,000. To train the network, we split the simulated dataset into 2
666 groups: [1] proper training set (3,990,000 examples); [2] validation set (10,000).

667 *Preventing overfitting: Early stopping and Dropout*

668 To prevent overfitting during training, we use: [1] the early stopping algorithm evaluating MAPE on a validation set;
669 and [2] dropout that we set to 0.5 in the feed-forward part of both NNs^[56] (0.4, 0.45, 0.55 and 0.6 values were tried
670 for basic BD model without improving the accuracy).

671 *Neural networks for model selection*

672 For model selection, we use the same architecture for FFNN-SS and CNN-CBLV as those for parameter inference
673 described above. The only differences are: [1] the cost function: categorical cross entropy and [2] the activation
674 function used for the output layer, that is, softmax function (of size 2 for small trees, selecting between BD and BDEI
675 model, and of size 3 for large trees, selecting between BD, BDEI and BDSS). As we use the softmax function, the
676 outputs of prediction are the estimated probabilities of each model, summing to 1.

677 The FFNN-SS and CNN-CBLV are trained on 8×10^6 trees in the small tree setting (4×10^6 trees per model, BD and
678 BDEI). In the large tree setting, the FFNN-SS is trained on 12×10^6 trees (4×10^6 trees per model, BD, BDEI and BDSS)
679 and the CNN-CBLV is trained on 9×10^6 trees (3×10^6 trees per model, BD, BDEI and BDSS), instead of 12×10^6 for
680 GPU limitation purposes.

681 CONFIDENCE INTERVALS (95% CI)

682 *Computation of 95% CI*

683 We compute 95% CI using parametric bootstrap. To facilitate the deployment and speed-up the computation, we
684 perform an approximation using a separate set of 1,000,000 simulations for calculation of CI. For each simulation in
685 the CI set, we store the true parameter values (*i.e.* values with which we simulated the tree) and the parameter values
686 predicted with both of our methods. This large dataset of true/predicted values is used to avoid new simulations, as
687 required with the standard parametric bootstrap.

688 For a given simulated or empirical tree T , we obtain a set of predicted parameter values, $\{p\}$. The CI computation
689 procedure searches among stored data those that are closest to T in terms of tree size, sampling probability and
690 predicted values. We first subset:

- 691 • 10% of simulations within the CI set, which are closest to T in terms of size (number of tips), thus obtaining
692 100,000 CI sets of true/predicted parameter values.
- 693 • Amongst these, 10% of simulations that are closest to T in terms of sampling probability.

694 We thus obtain 10,000 CI sets of real/predicted parameter values, similar in size and sampling probability to T . For
695 each parameter value p predicted from T , we identify the 1,000 nearest neighboring values amongst the 10,000 true
696 values of the same parameter available in the CI sets, $R_{CI} = \{r_{i=1,1000}\}$, and keep the corresponding predicted values,
697 $P_{CI} = \{p_{i=1,1000}\}$. We then measure the errors for these nearest neighbors as $E_{CI} = \{e_i = p_i - r_i\}$. We center these
698 errors around p , using the median of errors, $median(E_{CI})$, which yields the distribution of errors for given prediction
699 $p: D = \{p + e_i - median(E_{CI})\}$, from which we extract the 95% CI around p . Individual points in the obtained
700 distribution that are outside of the parameter ranges covered through simulations are set to the closest boundary value
701 of the parameter range. For example, for f_{SS} , if for a point in the distribution we obtain a value lower than 0.05, we set
702 the value of that point to 0.05; and if we obtain a value larger than 0.20, we set it to 0.20. We undertake this procedure

703 for all parameters except for the time related ones, that is, infectious and incubation period as these depend on the time
704 scaling. The width of our 95% CIs is defined as the distance between the 2.5% and 97.5% percentile.

705 *Assessment of 95% CI coverage and width*

706 To assess this fast implementation of the parametric bootstrap, we used the test set of 10,000 simulations (and 100
707 simulations for comparison with BEAST2 95% CI). We measured the coverage being defined as the fraction of
708 simulations where the true/target parameter values are inside the obtained 95% CI:

$$709 \quad 95\% \text{ CI accuracy} = \frac{\# \text{ true values inside 95\% CI}}{\# \text{ simulations}}$$

710 We applied the same criteria for BEAST2. For comparison of all methods, we excluded BDEI and BDSS simulations
711 for which BEAST2 did not converge after 10 million steps. To draw BEAST2 CIs, we discarded the burn-in, *i.e.* the
712 first 10% of the MCMC, and calculated the CI on the remaining part of the chain. The CI width and coverage within
713 the CIs obtained either by NNs or BEAST2 are reported in **Supplementary Table 4**.

714 There exists a plethora of approaches for assessment of uncertainty and CI estimation. For example, [1] in a similar
715 ABC context, the use of neighboring trees (based on the Euclidean distance, not applicable to CBLV and questionable
716 with SS) combined with a regression-based correction similar to that explained above^[19,20]; [2] a (non-approximated)
717 parametric bootstrap^[56]; [3] predicting values from a distribution of trees reconstructed with Bayesian methods^[10]; *etc.*
718 We chose an approximation of parametric bootstrap for its easy deployability, speed, coverage and width of produced
719 CIs. The easy deployability comes from the fact that CIs are based on pre-calculated data stored in our CI set. The
720 speed of the method comes from it not requiring simulations of new trees, and thus producing CIs within 2-4 seconds.
721 The coverage and width are comparable to those of BEAST2 (**Sup. Table 4**), a Bayesian method, intended to estimate
722 the distribution of parameters and the uncertainty of inferences, with high computational cost.

723 **MODEL ADEQUACY**

724 *A priori check*

725 We performed a sanity check using the SS of the test set simulations and the SS measured on the empirical HIV
726 phylogeny. We reduced and centered the SS and performed a Principal Component Analysis (PCA) using the PCA
727 function from the scikit-learn^[51] package.

728 We highlighted the datapoint corresponding to the Zürich HIV MSM phylogeny in **Supplementary Fig. 8**, for each
729 model (BD, BDEI and BDSS). Dissemblance between the simulations and the HIV phylogeny would manifest by this
730 datapoint being outside the distribution corresponding to the simulations.

731 *A posteriori* check

732 We performed a test analogous to the *a priori* model adequacy check. This time, instead of using the test set as
733 representative of simulations, we simulated 10,000 additional simulations under the selected BDSS model. Parameter
734 values were resampled from uniform distribution with boundaries given by the 95% CIs, and sampling probability
735 fixed to presumed value of 0.25 (**Fig. 4**).

736 MODELS

737 The models we used for tree simulations are represented in the form of flow diagrams in **Fig. 1**. We simulated dated
738 binary trees for [1] the training of NNs and [2] accuracy assessment of parameter estimation and model selection. We
739 used the following three individual-based phylodynamic models:

740 *Constant rate birth-death model with incomplete sampling*

741 This model (BD^[8,9], **Fig. 1 a**) contains three parameters and three compartments: infectious (I), removed with sampling
742 (R) and removed unsampled (U) individuals. Infection takes place at rate β . Infectious individuals are removed with
743 rate γ . Upon removal, an individual is sampled with probability s .

744 For simulations, we re-parameterized the model in terms of: basic reproduction number, R_0 ; infectious period, $1/\gamma$;
745 sampling probability, s ; and tree size, t . We then sampled the values for each simulation uniformly at random in the
746 ranges given in **Supplementary Table 1**.

747 *Birth-death model with exposed-infectious classes*

748 This model (BDEI^[10-12], **Fig. 1 b**) is a BD model extended through the presence of an exposed class. More specifically,
749 this means that each infected individual starts as non-infectious (E) and becomes infectious (I) at incubation rate ϵ .
750 BDEI model thus has four parameters (β , γ , ϵ and s) and four compartments (E, I, R and U).

751 For simulations, we re-parameterized the model similarly as described for BD, set the ε value via $1/\gamma$ and incubation
752 ratio ($=\varepsilon/\gamma$). We sampled all parameters, including ε/γ , from a uniform distribution, just as with BD.

753 *Birth-death model with superspreading*

754 This model (BDSS^[5,10,11], **Fig. 1 c**) accounts for heterogeneous infectious classes. Infected individuals belong to one
755 of two infectious classes (I_S for superspreading and I_N for normal spreading) and can transmit the disease by giving
756 birth to individuals of either class, with rates $\beta_{S,S}$ and $\beta_{S,N}$ for I_S transmitting to I_S and to I_N , respectively, and $\beta_{N,S}$ and
757 $\beta_{N,N}$ for I_N transmitting to I_S and I_N , respectively. However, there is a restriction on parameter values: $\beta_{S,S} * \beta_{N,N} =$
758 $\beta_{S,N} * \beta_{N,S}$. There are thus superspreading transmission rates $\beta_{S_}$ and normal transmission rates $\beta_{N_}$ that are $X_{SS} (=$
759 $\frac{\beta_{S,S}}{\beta_{N,S}} = \frac{\beta_{S,N}}{\beta_{N,N}})$ times higher for superspreading. At transmission, the probability of the recipient to be superspreading is
760 $f_{ss} (= \frac{\beta_{S,S}}{\beta_{S,S} + \beta_{S,N}})$, the fraction of superspreading individuals at equilibrium. We consider that both I_S and I_N populations
761 are otherwise indistinguishable, that is, both populations share the same infectious period $(1/\gamma)^{[5,10,11]}$. The model thus
762 has six parameters, but only five need to be estimated to fully define the model^[5,10].

763 For simulations, we chose parameters of epidemiological interest for re-parameterization: basic reproduction number
764 $R_0 (= \frac{\beta_{S,S} + \beta_{N,N}}{\gamma})$, infectious period $1/\gamma$, f_{ss} , X_{ss} and sampling probability s . In our simulations, we used uniform
765 distributions for these 5 parameters, just as with BD and BDEI (**Supplementary Table 1**).

766 **SIMULATIONS**

767 For the parameters R_0 , $1/\gamma$, and s , that are common to all three birth-death models, the same value boundaries were
768 used across all models (**Supplementary Table 1**). We considered two spans of tree size: ‘small trees’ with 50 to 199
769 tips, and ‘large trees’ with 200 to 500 tips. We then sampled parameter values uniformly at random within these
770 parameter boundaries with standard Latin-hypercube sampling^[57] using PyDOE package. We created 3,990,000
771 parameter sets for training, 10,000 for validation and early stopping, another 10,000 for testing parameter inference
772 and model selection (comparison with BEAST2 used a subset of 100, for computing time reasons), and 1,000,000
773 parameter sets for fast computation of CIs.

774 With these parameter sets, we simulated trees under each birth-death model using our implementation in Python of
775 Gillespie algorithm^[58], based on a standard forward simulator. Comparable accuracies (as in **Fig. 3** and

776 **Supplementary Fig. 2**, both for BEAST2 and our methods) were reached on test simulations obtained with a well-
777 established, but slower, simulator TreeSim^[4,5,7] (data not shown).

778 Each simulation started with one infectious individual (the class was chosen randomly under the BDSS model) and
779 stopped when we obtained a tree with the given number of sampled individuals (tips). If the epidemic died away
780 stochastically, that is, there was no more infectious tips left due to stochastic death before reaching the given tree size,
781 we re-initialized the simulation up to 100 times. Only around 11% of simulations reached more than 2 iterations (20%
782 for BDSS), and less than 0.5% reached more than 50 iterations for all models. If still no tree of given size was obtained
783 after 100 iterations, we discarded the parameter set (less than 0.3% of all sets) and generated a new one to keep the
784 desired number of simulations. This enabled us to maintain a nearly uniform coverage of parameter space, within
785 selected parameter boundaries.

786 **METHOD COMPARISON**

787 *Parameter inference with BEAST2*

788 To assess the accuracy of our methods, we compared it with a well-established Bayesian method, as implemented in
789 BEAST2 (version 2.6.2). We used the BDSKY package^[4] (version 1.4.5) to estimate the parameter values of BD
790 simulations and the package *bdmm*^[12,13] (version 1.0) to infer the parameter values of BDEI and BDSS. Furthermore,
791 for the inference on BDSS simulations, instead of BEAST 2.6.2 we used the BEAST2 code up to the commit
792 nr2311ba7, which includes important fixes to operators critical for our analyses. We set the Markov Chain Monte
793 Carlo (MCMC) length to 5 million steps for the BD model, and to 10 million steps for the BDEI and BDSS models.
794 The xml files and command lines are available at:

795 https://github.com/evolbioinfo/phylodeep/tree/main/data_publication.

796 The sampling probability was fixed during the estimation. Since the BD, BDEI and BDSS models implemented in
797 BEAST2 do not use the same parametrizations as our methods, we needed to apply parameter conversions for setting
798 the priors for BEAST2 inference, and for translating the BEAST2 results back to parameterizations used in our
799 methods, in order to enable proper comparison of the results (**Table 1**). More specifically, the BEAST2 parameters
800 can be converted to those used in our methods, that is, instead of infectious period and incubation period, BEAST2
801 uses the inverse of these, namely the infectious rate and incubation rate, respectively; instead of superspreading

802 transmission ratio and superspreading fraction at equilibrium, it uses individual sub-component parameters $R_{0,SS}$,
803 $R_{0,SN}$, $R_{0,NS}$ and $R_{0,NN}$, which we will collectively refer to as “partial R_0 ”. For BDSS, the BEAST2 prior was thus not
804 the same as that of our simulations for BDSS (**Table 1** and **Supplementary Table 1**), since BEAST2 does not infer
805 the same parameters. We used the range of all parameter values used in our simulations to set the boundaries of
806 uniform prior distributions of parameters inferred by BEAST2. The initial values in the MCMC were set to the medians
807 observed in the training set. During the inference, the parameter values were constrained in the same way as in the
808 simulations, namely, we used the following constraint $R_{0,NN} * R_{0,SS} = R_{0,SN} * R_{0,NS}$ (equivalent to $\beta_{N,N} * \beta_{S,S} =$
809 $\beta_{S,N} * \beta_{N,S}$) in the BDSS model inference. Furthermore, the effective frequency of superspreading individuals
810 (parameter called “geo-frequencies” in *bdmm*) was constrained to be between 5% and 20%. Due to the parameter
811 conversions, and despite these constraints the inferred f_{ss} and X_{ss} can reach values outside the boundaries used for
812 simulations, in which case we set them to the closest boundary for fair comparison with deep learning methods in **Fig.**
813 **3** (e.g. if the median *a posteriori* f_{ss} was estimated to be larger than 0.20, it was set to 0.20 and if inferred f_{ss} was less
814 than 0.05, it was set to 0.05). The goal of this correction was to avoid penalizing BEAST2 when it converged to local
815 minima outside of the parameter boundaries used for simulations, which are implicitly known to NNs since they were
816 trained on simulations with parameters within these boundaries.

817 After we obtained the parameters of interest from the original parameters estimated by BEAST2, we evaluated the
818 Effective Sample Size (ESS) on all parameters. We reported the absolute percentage error of the median of *a posteriori*
819 values, corresponding to all reported steps (reported steps being spaced by 1,000 actual MCMC steps) past the 10%
820 burn-in. For simulations for which BEAST2 did not converge, we considered the median of the parameter distribution
821 used for simulations (**Fig. 3**, **Sup. Fig. 2**, **Sup. Tab. 2-3**) or excluded them from the comparison (**Sup. Tab. 2-3**,
822 values reported in brackets, **Sup. Tab. 4**).

823 For the HIV application, the prior of infectious period was set to [0.1, 30] years (uniform). All the other parameters
824 had the same prior distributions as used in simulations and shown in **Table 1**.

825

Table 1: BEAST2 priors and their relation to the parameters of interest.

Parameters	Name	BEAST2 parameter	Range	Initial value	Formula
γ	become uninfected rate	Yes	U(0.1,1.0)	0.55	
$1/\gamma$	infectious period	No	[1,10]	1.8	
s	sampling probability	Yes	fixed	true value	
R_0	basic reproduction number	Yes	U(1.0,5.0)	3.0	$= \beta/\gamma$ (BD & BDEI)
ϵ	incubation rate	Yes	U(0.02,5.0)	2.51	
$1/\epsilon$	incubation period	No	[0.2,50]	0.40	
$R_{0,SS}$	partial, within deme superspreading R_0	Yes	U(0.14,4.31)	1.25	
$R_{0,NN}$	partial, within deme normal spreading R_0	Yes	U(0.14,4.31)	1.44	
$R_{0,SN}$	partial, outside deme superspreading R_0	Yes	U(0.034,32.30)	9.0	
$R_{0,NS}$	partial, outside deme normal spreading R_0	Yes	U(0.034,32.30)	0.20	
R_0	basic reproduction number	No	[0.28,8.62]	2.69	$= R_{0,SS} + R_{0,NN}$ (BDSS)
X_{SS}	superspreading infectious ratio at equilibrium	No	$[4 \cdot 10^{-4}, 127]$	6.25	$= R_{0,SS}/R_{0,NS} = R_{0,SN}/R_{0,NN}$
f_{SS}	fraction of superspreading individuals at equilibrium	No	$[4 \cdot 10^{-3}, 0.99]$	0.12	$= R_{0,SS}/(R_{0,SS} + R_{0,SN})$ $= 1 - R_{0,NN}/(R_{0,NN} + R_{0,NS})$

826

827 This table shows parameters and their prior distributions used during inference with BEAST2. We display the
 828 parameters, their definitions and priors in BEAST2, that are common to all models (in yellow), common to BD
 829 and BDEI (in red), BDEI-specific (in purple) and BDSS-specific (in green). From these parameters, we deduce
 830 the values and distributions of parameters of interest as shown in the table. Note that the parameters of
 831 epidemiological interest are basic reproduction number and infectious period for BD, BDEI and BDSS, incubation
 832 period for BDEI, and superspreading infectious ratio and fraction of superspreading individuals for BDSS. We
 833 check convergence (ESS) on all parameters and extract median a posteriori and CI value exclusively for the
 834 parameters of epidemiological interest.

835

836 *Model selection with BEAST2*

837 We performed model selection under BEAST2 using Akaike's information criterion through MCMC (AICM)^[32,33].

838 The AICM is based on the following formula:

839

$$AICM = 2s_l^2 - 2l$$

840 where l and s_l^2 are the sample mean and variance of the posterior log-likelihoods. The AICM is an equivalent of AIC

841 and the model with lowest AICM value is selected.

842 For 100 simulations obtained with each model (BD, BDEI and BDSS for large trees, BD and BDEI for small trees),
843 we performed parameter estimation with BEAST2 under each model, computed AICM considering the whole MCMC,
844 but excluding 10% burn-in (*i.e.* 9,000 log-likelihood values for BDEI and BDSS considered in total, 4,500 for BD).
845 The results of model selection are shown in **Supplementary Fig. 3**. The BDEI and BDSS simulations for which
846 BEAST2 did not reach an ESS of 200 for all parameters were excluded from the computation of model selection
847 accuracy for all methods.

848 *Linear regression*

849 For each model, linear regression was trained using reduced and centered summary statistics (using scikit-learn
850 package, as with FFNN). Its bias and accuracy were assessed using the same criteria as for the NN approaches
851 (**Supplementary Tables 2-3, Supplementary Fig. 4**).

852 *FFNN-CBLV*

853 We trained an FFNN on CBLV representation. The FFNN architecture was close to the one described in *Architecture*
854 with one extra hidden layer, so 5 layers in total, organized in a funnel shape with 128-64-32-16-8 neurons and 1 output
855 layer of size 2-4 depending on the number of parameters to be estimated. The setting during the training and the sizes
856 of training, validation and testing sets were the same as for the CNN-CBLV. Its bias and accuracy were assessed using
857 the same criteria as for other NN approaches (**Supplementary Tables 2-3, Supplementary Fig. 5**).

858 *TreePar*

859 We used TreePar^[5] for MLE. With BD, we obtained results close to estimates under BEAST2, which is consistent
860 with former studies^[58]. TreePar^[5] uses an exact analytical formula of likelihood for BD and thus these (and BEAST2)
861 results are theoretically optimal.

862 We also performed several trials to do parameter inference for the more complex models, *i.e.* BDEI and BDSS, but in
863 a large number of cases we encountered numerical problems, *e.g.* underflow or overflow issues, which resulted in
864 infinite negative log-likelihood values, and eventually failed runs. When the calculations did not fail, we found that
865 many estimations under BDSS and BDEI had lower likelihood than estimations performed with (nested) BD on the

866 same input data. These numerical issues, without available solutions at the moment, were confirmed by the authors of
867 the TreePar package.

868 *Null models*

869 To assess how much information was learned on given problem, we compared FFNN-SS and CNN-CBLV to two null
870 models.

871 The first null model was the FFNN trained for each model on 4,000,000 simulations using SS, but with randomly
872 permuted target values (*i.e.* the initial correspondence between the SS and underlying parameter values was lost, while
873 the range of values was conserved). We then predicted parameters for 10,000 test simulations (100 for comparison
874 with BEAST2) and measured the mean absolute relative error (MRE, equivalent to MAPE; **Supplementary Table**
875 **2**). In such a case, the FFNN always predicted values close to the value with the lowest value of the cost function, *e.g.*
876 2.2 for parameter values uniformly sampled between 1 and 5 (and not random values inside the range). The MRE of
877 this approach represents then the lowest MRE that machine learning approaches can have in the absence of
878 information, but the knowledge of the parameter distribution. This can be used to get an idea of how well the trained
879 approaches perform and how much information regarding each parameter they can extract from the data.

880 The second null model was a set of random values sampled from the parameter ranges that were used for simulations
881 (**Supplementary Table 2**). In this model, as opposed to the previous null model, there is no training phase and we do
882 not learn the best compromise in the absence of information.

883 **PERFORMANCE ASSESSMENT**

884 *Mean relative error MRE*

885 To compare the accuracy of parameter estimation, we used 100 simulated trees per model. We computed the mean
886 absolute relative error (MRE, **Fig. 3, Supplementary Table 2, Supplementary Fig. 2**) between [1] the true (or target)
887 parameter values and the predicted values for machine learning approaches; and [2] the true (or target) parameter
888 values and the median *a posteriori* values obtained with BEAST2, which are more stable and accurate than maximum
889 *a posteriori* values:

$$890 \quad MRE = \frac{1}{100} \sum_{i=1}^{100} \frac{abs(predicted_i - target_i)}{target_i}.$$

891 We plotted individual absolute relative errors (RE) of predictions (**Fig. 3, Supplementary Fig. 2**) for each simulation
892 i , calculated as:

$$893 \quad RE_i = \frac{abs(predicted_i - target_i)}{target_i}$$

894 Not being limited by the computational cost for machine learning approaches, we computed the same metric but on
895 10,000 simulations (**Supplementary Figs. 4-6**; results from 1,000 simulations plotted in **Supplementary Fig. 7**).

896 We assessed the statistical significance of MRE differences using paired z-test. The two NN approaches were also
897 compared using the same test, but no significant differences were found.

898 *Mean relative bias MRB*

899 To compare the bias in parameter estimation, we used 100 simulated trees per model. We computed the mean relative
900 bias (MRB) between [1] the true (or target) parameter values and the predicted values for machine learning
901 approaches; and [2] the true (or target) parameter values and the median *a posteriori* values obtained with BEAST2
902 (**Supplementary Table 3**):

$$903 \quad MRB = \frac{1}{100} \sum_{i=1}^{100} \frac{(predicted_i - target_i)}{target_i}$$

904 *Model selection accuracy*

905 We performed model selection with CNN-CBLV, FFNN-SS and BEAST2 on 100 simulations obtained with each
906 model (10,000 for a sub-comparison of CNN-CBLV and FFNN-SS). Results are shown in **Supplementary Fig. 3** in
907 the form of confusion matrices, where the columns represent the true/target classes, and the rows are the predicted
908 classes.

909 We then computed the accuracy of each method:

$$910 \quad accuracy = \frac{\# \text{ true predictions}}{\# \text{ total predictions}}$$

911 For BEAST2 model selection and large trees, the chain did not converge (displayed as “ESS<200” in **Supplementary**
912 **Fig. 3**) for 24.3% simulations of large trees and 4.5% simulations of small trees. We did not consider these in accuracy
913 measurements, for all the methods.

914 *Comparison of time efficiency*

915 For FFNN-SS and CNN-CBLV, we reported the average CPU time of encoding a tree (average over 10,000 trees), as
916 reported by NextFlow workflow manager^[60], a pipeline software that we used. The inference time itself was negligible.

917 For BEAST2, we reported the CPU time averaged over 100 analyses with BEAST2 as reported by NextFlow. For the
918 analyses with BDEI and BDSS models, we report the CPU time to process 10 million MCMC steps, and for the
919 analyses with BD, we report the CPU time to process 5 million MCMC steps. To account for convergence, we re-
920 calculated the average CPU time considering only those analyses, for which the chain converged and the ESS of 200
921 was reached across all inferred parameters.

922 **HIV DATASET**

923 We used the original phylogenetic tree reconstructed by *Rasmussen et al.*^[25] from 200 sequences corresponding to the
924 largest cluster of HIV-infected men-having-sex-with-men (MSM) subpopulation in Zurich, collected as a part of the
925 Swiss Cohort Study^[24]. For details on tree reconstruction, please refer to their article.

926 **ADDITIONAL REFERENCES**

927 50. Zarestkii, K. Reconstructing a tree from the distances between its leaves. (in Russian) *Uspehi*
928 *Mathemicheskikh Nauk* **20**, 90-92 (1965).

929 51. Pedregosa, F. et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–
930 2830 (2011).

931 52. Abadi, M. et al. TensorFlow: Large-scale machine learning on heterogeneous systems. Preprint at
932 <https://arxiv.org/abs/1603.04467> (2015).

933 53. Chollet, F. K. <https://keras.io>. (2015).

934 54. Clevert, D.A., Unterthiner, T., Hochreiter, S. Fast and Accurate Deep Network Learning by Exponential Linear
935 Units (ELUs). *ICLR* (2016).

936 55. Kingma, D.P. & Ba, J. Adam: A Method for Stochastic Optimization. *ICLR* (2015).

937 56. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R. Dropout: A simple way to prevent
938 neural networks from overfitting. *Journal of Machine Learning Research* **15**,1929–1958 (2014).

939 57. Efron, B. *Breakthroughs In Statistics, Ch. Bootstrap Methods: Another Look at the Jackknife*. (Springer, New
940 York, 1999).

941 58. McKay, M., Beckman, R., Conover, W. A Comparison of Three Methods for Selecting Values of Input
942 Variables in the Analysis of Output from a Computer Code. *Technometrics* **21**, 239-245 (1979).

943 59. Gillespie, D.T. Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry* **81**,
944 2340-2361 (1977).

945 60. Di Tommaso, P. et al. Nextflow enables reproducible computational workflows. *Nat. Biotechnol.* **35**, 316–319
946 (2017)