# 1    A Comparison of Deep Learning Architectures for

# 2    Inferring Parameters of Diversification Models from

# 3            Extant Phylogenies

4

5    *Ismaël Lajaaiti[1,2,*], Sophia Lambert[3], Jakub Voznica[3], Hélène Morlon[3], Florian Hartig[1]*

6    *[1] Theoretical Ecology Lab, University of Regensburg, Regensburg, Germany*

7    *[2] Institut des Sciences de l'Évolution de Montpellier, Centre National de la Recherche*

8    *Scientifique, Montpellier, France*

9    *[3] Institut de Biologie de l'ENS (IBENS), École Normale Supérieure, CNRS, INSERM,*

10    *Université PSL, 75005 Paris, France*

11

12    *[*]Correspondence to be sent to: Ismaël Lajaaiti, Institut des Sciences de l'Évolution*

13    *de Montpellier (Bat. 22), Centre National de la Recherche Scientifique, Montpellier,*

14    *France; E-mail:* ismael.lajaaiti@gmail.com

15

16                 ABSTRACT

17    To infer the processes that gave rise to past speciation and extinction rates across

18    taxa, space and time, we often formulate hypotheses in the form of stochastic

19    diversification models and estimate their parameters from extant phylogenies using

20    Maximum Likelihood or Bayesian inference. Unfortunately, however, likelihoods can

21    easily become intractable, limiting our ability to consider more complicated

22    diversification processes. Recently, it has been proposed that deep learning (DL)

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

23    could be used in this case as a likelihood-free inference technique. Here, we explore

24    this idea in more detail, with a particular focus on understanding the ideal network

25    architecture and data representation for using DL in phylogenetic inference. We

26    evaluate the performance of different neural network architectures (DNN, CNN,

27    RNN, GNN) and phylogeny representations (summary statistics, Lineage Through

28    Time or LTT, phylogeny encoding and phylogeny graph) for inferring rates of the

29    Constant Rate Birth-Death (CRBD) and the Binary State Speciation and Extinction

30    (BISSE) models. We find that deep learning methods can reach similar or even

31    higher accuracy than Maximum Likelihood Estimation, provided that network

32    architectures and phylogeny representations are appropriately tuned to the

33    respective model. For example, for the CRBD model we find that CNNs and RNNs

34    fed with LTTs outperform other combinations of network architecture and phylogeny

35    representation, presumably because the LTT is a sufficient and therefore less

36    redundant statistic for homogenous BD models. For the more complex BiSSE model,

37    however, it was necessary to feed the network with both topology and tip states

38    information to reach acceptable performance. Overall, our results suggest that deep

39    learning provides a promising alternative for phylogenetic inference, but that data

40    representation and architecture have strong effects on the inferential performance.

41

42    *Keywords*: Cladogenesis, Machine Learning, Macroevolution, Stochastic Biodiversity

43    Model

44

45

46

47

## INTRODUCTION

48

49     Species richness varies greatly across taxonomic groups (G. E. Hutchison

50    1959), geological times (Barnosky et al. 2011) and geographical regions (Gaston

51    and Blackburn 2000). It is generally accepted that these patterns in species richness

52    emerge from variation in speciation and extinction rates in addition to dispersal and

53    migration processes. For instance, important ecological phenomena such as the

54    'Latitudinal Diversity Gradient' (Hillebrand 2004) may be partly explained by

55    variations in species net diversification rates defined as the balance of the speciation

56    and the extinction rate (Mittelbach et al. 2007; Rolland et al. 2014; Pontarp et al.

57    2019). Our general understanding of the processes that cause these rates changes,

58    however, is still very limited (Rabosky 2009a, 2009b; Condamine et al. 2013; Moen

59    and Morlon 2014).

60     To better understand the drivers of variation in diversification rates and their

61    consequences for biodiversity dynamics, it is crucial first to accurately estimate net

62    diversification rates (Pyron and Burbrink 2013), and then decompose them into

63    speciation and extinction rates (Stadler 2013). The latter is not trivial, as fossil data

64    are rarely available and reconstructed phylogenies do not include information on

65    extinct species. Indeed, without further constraints, the problem is ill-posed, meaning

66    that different diversification dynamics can lead to the same extant phylogeny (Louca

67    and Pennell 2020; Morlon et al. 2022). However, when making additional

68    assumptions about the functional form of extinction and speciation rates over time, it

69    is often possible to statistically estimate the parameters of these functions, and thus

70    infer the path that led to the currently observed species (Nee et al. 1994; Etienne

71    and Rosindell 2012; Morlon 2014).

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

72      Arguably one of the simplest diversification models is the Constant Rate Birth-

73      Death (hereafter 'CRBD') model, which assumes that the speciation and extinction

74      rates are both homogeneous across lineages and constant through time. Nee et al.

75      (1994) showed that extinction and speciation rates in this model can be estimated

76      based on reconstructed phylogenies with Maximum Likelihood using a distinct

77      pattern called the "pull of the present". In recent years, many studies have worked on

78      extending this model to account for various types of rate heterogeneity as well as

79      their potential drivers (Morlon et al. 2011; Etienne et al. 2012). These models

80      commonly allow for variations of speciation and extinction rates through time (time-

81      dependent models), across lineages (inhomogeneous models) as well as in

82      dependence of environmental (Condamine et al. 2013) or biotic factors (Etienne et

83      al. 2012). For instance, the time dependent Birth-Death model (Hallinan 2012) allows

84      to consider homogenous changes (*e.g.* an exponential decay) of speciation and

85      extinction rates over time (Rabosky and Lovette 2008; Morlon et al. 2011; Stadler

86      2011). Other diversification models allow lineage-specific shifts in diversification

87      rates that can be either discrete, as can be expected to occur with the appearance of

88      key innovations (Alfaro et al. 2009), or continuous, as can be expected to occur

89      given the gradual evolution of phenotypes (*e.g.* the 'ClaDS' model; Maliet et al. 2019)

90      . The specific innovations or phenotypes that drive these shifts can be either implicit,

91      as in the latter models, or explicitly modeled, as in State-dependent Speciation and

92      Extinction (SSE) models. A particular example of this is the Binary State Speciation

93      and Extinction model, hereafter 'BiSSE' (Maddison et al. 2007), which considers the

94      effect of a binary state on speciation and extinction rates.

95      A constant challenge when developing new diversification models is

96      establishing robust methods to fit them to data. For diversification models with low or

97  moderate complexity, the likelihood (*i.e.* the probability to observe a reconstructed

98  phylogeny given the model and its parameters) can be computed analytically or

99  approximated numerically (see Nee et al. 1994 for the CRBD model and Maddison et

100  al. 2007 for the BiSSE model). If the likelihood can be computed, model parameters

101  can be inferred using either Maximum Likelihood Estimation (MLE; see for instance

102  Ricklefs 2007) or Bayesian inference (*e.g.* Silvestro et al. 2011).

103  When the likelihood of a diversification model is analytically or numerically

104  intractable, simulation-based methods such as Approximate Bayesian Computation

105  (ABC) are typically being used. The ABC approach approximates the likelihood of a

106  model by comparing model predictions to data via summary statistics (Csilléry et al.

107  2010; see Saulnier et al. 2017 for an example). Although ABC can successfully infer

108  parameters from complex diversification models, it has two major drawbacks. First,

109  ABC typically requires a large number of simulations, which is often computationally

110  prohibitive in practical applications. Secondly, ABC suffers from the curse of

111  dimensionality, meaning that the computational increases sharply with the

112  dimensions of the summary statistics, which must be at least as large as the number

113  of model parameters to ensure sufficiency (Beaumont 2010; Csilléry et al. 2010;

114  Hartig et al. 2011). Together, both properties limit the applicability of ABC to

115  complicated models with many parameters.

116  Recent advances in the field of deep learning algorithms provide an alternative

117  for likelihood-free inference in phylogenetic inference. Unlike ABC, deep learning can

118  easily deal with high-dimensional data and even benefit from it (LeCun et al. 2015),

119  thus avoiding the necessity to find appropriate summary statistics. It was shown that

120  deep learning approaches based on Convolutional Neural Networks (CNNs) can

121  outperformed ABC for inference tasks in population genomics (Chan et al. 2018;

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

122  Schrider and Kern 2018). Moreover, as pointed out by (Schrider and Kern 2018),

123  neural networks are very flexible about their input data structure, possibly allowing a

124  far greater or more diverse set of inputs than in traditional statistical models.

125      Unlike for ABC, however, there is little experience on using deep learning

126  algorithms for parameter inference in diversification models, and there are a number

127  of options to implement DL for phylogenetic inference. One of the first questions to

128  solve is how to best encode phylogenetic data, which in turn affects the DL

129  architectures that could be considered for training. Recently (Voznica et al. 2022)

130  developed a full matrix representation for non-ultrametric phylogenies with the goal

131  of fitting epidemiological models, and compared the performance of feed-forward

132  neural networks (DNNs) and CNNs that were supplied either with this tree

133  representation or with conventional summary statistics to likelihood-based inference

134  techniques. (Lambert et al. 2022) adapted this approach with the goal of fitting birth-

135  death diversification models to reconstructed (ultrametric) phylogenies, extended the

136  matrix representation to the case of representing phylogenies with associated tip

137  state data, and similarly compared the performance of DNN, CNN and MLE. These

138  studies showed that the matrix encoding processed by CNNs performed very well,

139  leading to gains in predictive performance compared to likelihood-based methods

140  (Voznica et al. 2022). However, there are a number of further options regarding the

141  encoding and the network architectures, and we conjectured that the optimal

142  combinations of those would depend on the diversification model to be estimated.

143      When considering the choice of network architecture (Pichler and Hartig 2023),

144  the simplest option would be a standard fully connected Deep Neural Network

145  (DNN). The disadvantage of this option is that fully connected DNNs do not perform

146  favorably with high-dimensional structured input data (such as phylogenies).

147  Therefore, we can anticipate for successfully using DNNs, we would have to

148  summarize the phylogeny's shape (which by default is high dimensional with

149  neighborhood relationships being represented as a graph) using a limited number of

150  summary statistics (Saulnier et al. 2017), similar to the ABC approach. While

151  possible, this might result in a loss of information for the inference, depending on

152  whether those summary statistics are sufficient for the inference task.

153      A second option would be the use of Convolutional Neural Networks (CNNs).

154  CNNs apply one or several filters that slide incrementally across the input data to

155  detect spatial patterns. CNNs have shown great success for a variety of tasks,

156  including for fitting birth-death models to pathogen or species phylogenies encoded

157  in a full matrix representation (Voznica et al. 2022; Lambert et al. 2022). CNNs could

158  also be used on reduced representations of the phylogeny, such as the Lineage

159  Through Time (hereafter 'LTT') plot. This seems promising as it is known that the

160  slope of the LTT is a sufficient statistics for homogeneous models (Nee et al. 1994;

161  Ricklefs 2007).

162      When representing the phylogeny by its LTT, another option is the use of

163  Recurrent Neural Networks (RNNs). The neurons of RNNs include a temporal

164  feedback and, thanks to this feature, RNNs can process sequential data (*e.g.*, time

165  series) more efficiently. A common and efficient RNN neuron architecture is the Long

166  Short-Term Memory (LSTM) cell. RNNs based on LSTM cells can flexibly handle

167  long and short term dependencies of the input data (Hochreiter and Schmidhuber

168  1997; Yu et al. 2019) and have been applied with success in many fields, including

169  Natural Language Processing tasks (Huang et al. 2019), financial market forecasting

170  (Bukhari et al. 2020) or phoneme classification (Graves et al. 2005).

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

171    Finally, when considering the nature of a phylogeny, Graph Neural Networks

172    (GNNs) arise as a natural choice. GNNs generalize the idea of CNNs, which require

173    Euclidean neighborhoods, to graphs (Scarselli et al. 2009; Zhang et al. 2019; Wu et

174    al. 2021). GNNs are successfully used for many applications with a graph structure,

175    *e.g.* predicting molecules properties (Gilmer et al. 2017; Mansimov et al. 2019). The

176    potential advantage of using GNNs over CNNs is that GNNs can directly perform

177    convolutions along the phylogeny's topology, whereas the use of a CNN requires

178    transforming the phylogeny into a Euclidean structure, which potentially distorts

179    neighborhood relationships.

180    The hypotheses formulated in the previous paragraphs are based on the

181    theoretical knowledge about DL architectures, but so far, no systematic comparison

182    of the combination of phylogeny representation and neural network architecture has

183    been performed to confirm these conjectures. For example, while CNNs feeding on

184    encoded phylogenies were shown to infer rates with a good accuracy for different

185    birth-death models (Voznica et al. 2022; Lambert et al. 2022), we might expect that

186    this full encoding is redundant and thus suboptimal for homogeneous birth-death

187    models, because for these models, all the relevant information is in the LTT (Lambert

188    and Stadler 2013).

189    In this study, we systematically explore the interplay between model, data

190    representation and network architecture when inferring parameters of diversification

191    models using deep learning. As models, we considered the simple homogeneous

192    CRBD model and the more complex inhomogeneous BiSSE model, with the idea

193    that homogenous diversification models will require other data representations than

194    inhomogeneous models. We represented phylogenies with either: LTTs, sets of

195    summary statistics, phylogeny encodings from (Voznica et al. 2022 and Lambert et

196   al. 2022), or phylogeny graphs. Then, these representations were combined with the

197   suitable neural network architecture(s): summary statistics with DNNs, phylogeny

198   graphs with GNNs, encoded phylogenies with CNNs, and LTTs with CNNs or RNNs.

199   To evaluate the predictive performance of these deep learning inference methods,

200   we compared the prediction errors of each method and broke them down into three

201   terms: 1) variance, 2) uniform bias and 3) consistency bias (Smith and Rose 1995).

202   As a reference, we compared the deep learning results to the MLE for these models,

203   which is tractable in both cases.

204       Using this set-up, we ask three questions. First, can deep learning methods

205   accurately infer rates from diversification models, as suggested by (Lambert et al.

206   2022), and if so, can these methods outperform MLE regarding the prediction error?

207   We hypothesize that deep learning methods can theoretically outperform MLE as

208   they can trade off bias against variance to minimize the total error. Secondly, what is

209   the optimal representation of phylogeny data to infer rates for diversification models

210   of different complexity? We expect, for example, that the LTT is a sufficient statistic

211   for simple homogeneous models and is often preferable due to its simplicity. For

212   more complex models and specifically for inhomogeneous models, on the other

213   hand, we expect that simple representations such as the LTT will fail to provide the

214   neural network with the necessary information. More complex representations are

215   necessary for optimal inferential performance. Lastly, we ask how the choice of the

216   neural network architecture in combination with the chosen phylogeny representation

217   affects the inference. We hypothesize that the more complex the phylogeny

218   representation, the more important the neural network architecture becomes.

219                                   METHODS

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

220    *Diversification Models*

221    We chose two established diversification models as case studies to test the

222    performance of deep learning algorithms for phylogenetic inference: 1) the relatively

223    simple and homogeneous CRBD model, and 2) the more complex inhomogeneous

224    BiSSE model. The rationale for choosing these two models was to have two models

225    with tractable likelihoods but different complexity, and especially different in their

226    (in)homogeneity, as we conjectured that homogenous BD will profit less from

227    detailed representations of the phylogeny, given that it is known that all information

228    for their inference is contained in the LTT.

229    The Constant Rate Birth-Death (CRBD) model is one of the simplest

230    conceivable diversification models. It has two parameters, the speciation rate ($\lambda$) and

231    the extinction rate ($\mu$), that are homogenously constant over time. The likelihood of

232    this model is well-known (Nee et al., 1994) and depends only on the phylogeny's

233    branching times and not on its topology. To estimate diversification rates with the

234    MLE, e used the APE R package (Paradis et al. 2004).

235    In the Binary State Speciation and Extinction (BiSSE) model, lineages can

236    alternate between two states (0 and 1), where each state has its own speciation and

237    extinction rate. The switch through time between the two states translates character

238    transitions (*e.g.* sexual to asexual reproduction) that can have an impact on

239    diversification rates. The model has 6 parameters: 2 speciation rates ($\lambda_0$, $\lambda_1$), 2

240    extinction rates ($\mu_0$, $\mu_1$) and 2 two transition rates ($q_{01}$ for transition $0 \rightarrow 1$; $q_{10}$ for

241    transition $1 \rightarrow 0$). As for the CRBD model, the likelihood of the BiSSE model can be

242    computed (see Maddison et al., 2007). We imposed four constraints on the model to

243    simplify inference, thus decreasing the number of free parameters from six to two.

244   The constraints are: 1) $\lambda_1 = 2\lambda_0$; 2) $\mu_0 = 0$; 3) $\mu_1 = 0$; 4) $q_{01} = q_{10}$. Constraint 1)

245   ensures that states 0 and 1 have different speciation rates, 2) and 3) make the

246   model pure birth, and 4) is an assumption of symmetry that makes the probabilities

247   to switch from one state to another equal. Lambert et al. (2022) also used the CRBD

248   model, and a less constrained version of BiSSE. Here we simplified the model to

249   focus on the comparison of network architectures and data representation. To

250   estimate diversification rates with the MLE, we used the Diversitree R package

251   (FitzJohn 2012) and constrained the likelihood with the 4 constraints listed above.

252       To train the neural networks, we simulated 100,000 phylogenies for the CRBD

253   model and 1,000,000 phylogenies for the BiSSE model using the Diversitree library

254   (FitzJohn 2012) in R. We assumed complete sampling of the phylogenies. The latter

255   assumption (no missing species) could be relaxed (Lambert et al. 2022), but here

256   our focus is on the comparison between inference methods. For the CRBD model,

257   we draw the underlying parameters as follows: 1) we draw uniformly $\lambda$ in [0.1,1.0]; 2)

258   we draw uniformly the turnover rate $\varepsilon$ in [0, 0.9] from which we compute the

259   extinction rate $\mu = \varepsilon\lambda \in [0, 0.9\lambda]$. By doing so, we avoid the critical case where $\lambda \leq \mu$

260   (*i.e.,* speciation rate is inferior to or equivalent to the extinction rate). For the BISSE

261   model, we took $\lambda_0 \in [0.1,1.0]$ to stay in the same range as for the CRBD model and

262   $q_{01} \in [0.01,0.1]$ to ensure that one state is not overly represented compared to the

263   other one.

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

264 *Phylogeny Representation*

265    Most of the considered deep learning architectures cannot process phylogenies

266 as inputs, but require reformatting the phylogenies into a more regular data structure.

267 Here, we explain all options that we considered (see also Fig. 1).

268    *Summary statistics.*— Arguably the most basic option is to represent the

269 phylogeny by a number of summary statistics. We used a set of 84 summary

270 statistics inspired from the set of (Saulnier et al. 2017). Those summary statistics can

271 be split into three groups:

272    1) 8 statistics related to phylogeny topology (*e.g.* ratio of the width over the depth

273       of the phylogeny);

274    2) 25 branch lengths statistics (*e.g.* median of all branch lengths);

275    *3)* 51 LTT statistics (binned LTT coordinates and LTT slopes).

276 The list of the summary statistics and the changes compared to (Saulnier et al. 2017)

277 are detailed in the Table S1 of the Appendix.

278    For the BiSSE model we use the ratio of the number of tips in state 1 over the

279 total number of tips as an additional summary statistic.

280    *Lineage Through Time (LTT).*— LTTs illustrate the increase in lineages over

281 time. For a phylogeny of n tips, the LTT is composed of n-1 points where each point

282 is defined by two coordinates: time (t, abscissa) and the number of lineages (N,

283 ordinate). For a binary tree where each branching event results in two daughter

284 lineages, the LTT can be compressed to a 1-d array without loss of information, by

285 throwing away the number of lineages and keeping only the times of speciation

286 events.

287    *Phylogeny encoding (CBLV).*— Yet another alternative is to encode the

288    phylogeny in a real-values vector of length 2n-1 named the 'Compact Bijective

289    Ladderized Vector' ('CBLV', see Voznica et al., 2022). Each value of the vector

290    corresponds to one node (internal node or tip), thus there are n values for the n tips

291    and n-1 values for the internal nodes n-1, which result in a vector of 2n-1 values. The

292    encoding is done in two steps: 1) phylogeny ladderization, and 2) phylogeny

293    traversal. The principle of ladderization is to order each node's children, such that

294    the encoding is bijective (*i.e.*, one CBLV maps exactly to one phylogeny and

295    reversely). Here we ladderized the phylogeny such that for each node the left child is

296    the child which is further from the root. On the ladderized phylogeny, we perform an

297    inorder traversal using a classical recursive algorithm. If the visited node is an

298    internal node or the first tip is visited, its distance to the root is added to the vector.

299    Otherwise, its distance from its parent node is added to the vector. Moreover, for the

300    BiSSE model we add the n tip states to the vector. The tip states are ordered

301    according to the phylogeny traversal. Note that in the original method from Voznica

302    et al. (2022), the tip state information was added by adding one row to the distance

303    vector, thus resulting in a two-row matrix (first row for distances, second row for tip

304    states), while here we concatenate two 1-d vectors resulting in a 1-d vector. We

305    found that this modification did not affect the results.

306    *Phylogeny as a graph.*— The main motivation for the previous CBLV encoding

307    is to transform the phylogeny into a regular Euclidean format that can be easily

308    processed by CNNs. The disadvantage, however, is that neighborhood relationships

309    of the graph can be distorted in this process. Recent research in the field of machine

310    learning suggests that in such a case, it is often better to train neural networks

311    directly with the original graphs. This is done by so-called Graph Neural Networks or

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

312    GNNs, which extend the CNN idea to graphs. Here, we use the Pytorch Geometric

313    framework (Fey and Lenssen 2019) and the GraphNeuralNetworks Julia package

314    base on the Flux framework (Innes 2018) to represent the phylogeny as graph and

315    train GNNs. We provide the GNN with the phylogeny's topology and 4 attributes per

316    node: distance to the root and the lengths of the 3 edges linked to the node (1

317    incoming edge: parent $\rightarrow$ node; 2 outcoming edges: node $\rightarrow$child$_1$ and node $\rightarrow$

318    child$_2$). Moreover, for the BiSSE model we add as an attribute the tip states (0 or 1

319    for tips, -1 for nodes whose state is unknown).

320    *Neural Network Architectures*

321    For predicting model parameters from the formatted phylogenies, we

322    considered four different neural network architectures: DNN, CNN RNN and GNN

323    (see also Fig. 1). All architectures were built within the torch framework (Falbel and

324    Luraschi 2019) in R, except GNNs which were built with Torch in Python because

325    the PyTorch Geometric library (Fey and Lenssen, 2019) was needed for GNNs.

326    Hyperparameter tuning for each architecture was performed by hand (see Appendix
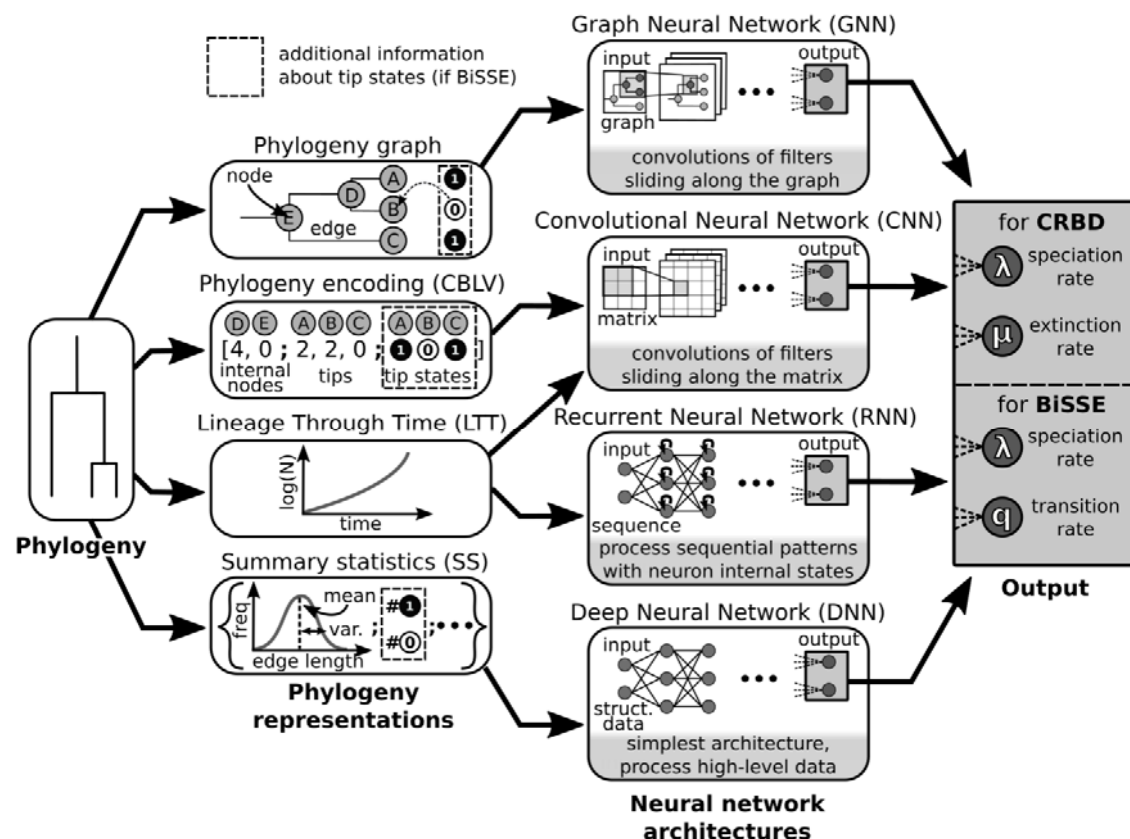
327    Tables S0-S5).

328

329

**Figure 1. Combining phylogeny representations with neural network architectures.** The five different combinations considered between phylogeny representations and neural network architectures are indicated by arrows from the first column to the second one. The rationales behind these five combinations are the following: 1) for phylogeny graph with GNN, GNN is the only architecture able to process graphs; 2) for CBLV with CNN, CNNs are able to detect patterns that are expected in the phylogeny encoding (Voznica et al. 2021); 3) for LTT with CNN, CNNs can detect patterns in the slope of the LTT; 4) for LTT with RNN, RNNs are designed to process time series; 5) SS with DNN, SS do not have spatial or temporal hierarchical order and therefore can be processed by a simple architecture. The third column describes the output of the inference methods for the two diversification models examined (CRBD and BiSSE).

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

342        *Training Neural Networks*

343        To train and validate the performance of the networks, we simulated 100,000

344    phylogenies under the CRBD model and 1,000,000 phylogenies under the BiSSE

345    model. For both models phylogenies were split randomly in three groups: 1) a

346    training set to train the neural networks (90% of the phylogenies); 2) a validation set

347    used to quantify performance of the neural network during the training (5% of the

348    phylogenies) and 3) a test set used to quantify the performance of the neural

349    network after the training (5% of the phylogenies). These split sets were the same for

350    all neural network architectures, so all of the architectures are trained and tested with

351    the same phylogenies. Moreover, the Maximum Likelihood Estimation ('MLE') was

352    tested on the phylogenies of the test set to fairly compare neural networks and MLE.

353        Training of the networks followed the standard practice of dividing the training

354    data into small batches (typically 64 phylogenies per batch, for more details see the

355    Appendix). During optimization, batches are successively provided to the network

356    and its weights are adjusted to reduce the prediction error (defined as the mean

357    squared error). Once the network has been fed with the whole training set, the

358    weights are frozen and the performance of the network is evaluated on the validation

359    set. This process of training and validation is called an 'epoch'. We repeated this

360    process until the prediction error on the validation set stopped to improve (early

361    stopping). This ensures that the network does not overfit the training data, which

362    tends to improve the performance of the model on new data (Bengio 2012). In

363    addition to early stopping, we also set at each new epoch a small proportion of the

364    network weights to zero (dropout, set to 1%), another standard method to avoid

365    overfitting. After the training of the networks, we evaluated their performances on the

366    test set.

367                                  *Interpreting Predictions Error*

368          Once the neural networks were trained, we use them to infer the parameter

369    values of the models from the test data (the phylogenies in the hold-out). We also

370    calculated the Maximum Likelihood Estimate for these phylogenies to establish a

371    baseline. Thus, for each method (the neural networks or the MLE) we have the

372    predicted model parameters ($y_{pred}$) *vs.* the true values ($y_{true}$).

373          To better understand the origin of deviations between the true and predicted

374    parameter values, we divide the prediction error into stochastic error and bias

375    according to Theil coefficients (Smith and Rose 1995). The idea of this approach is

376    that the sum squared differences $U_{tot} = \sum_{i=1}^{n}(y_{pred,i} - y_{true,i})^2$ can be expressed as

377    a sum    $U_{tot} = U_{uniformBias} + U_{consistencyBias} + U_{variance}$ where: 1) $U_{uniformBias} = $

378    $n\,\overline{(y_{pred} - y_{true})}^2$ corresponds the systematic error; 2) $U_{consistencyBias} = (b -$

379    $1)^2 \sum_{i=1}^{n}(y_{true,i} - \overline{y_{pred}})$ is the slope of the linear regression of $y_{pred}$ *vs.* $y_{true}$ and

380    $\overline{y_{pred}}$ the mean of the predicted values; 3) $U_{variance} = \sum_{i=1}^{n}(y_{pred,i} - y_{pred.lm,i})^2$ is

381    the stochastic error, measured by the variance of the $y_{pred}$ around the regression

382    slope $y_{pred.lm,i}$. Lastly, we normalize each term by the sum squared of the true

383    parameter values: $U_d = \sum_{i=1}^{n} y_{true,i}^2$. The error decomposition is illustrated in Figure

384    2.

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

385

**Figure 2. Decomposition of the prediction error.** The total error defined as the sum of squared differences is split in three terms also known as the Theil's coefficients: 1) variance that describe the width of the points distribution across the linear fit, 2) the uniform bias which quantifies the systematic error and 3) the consistency bias which expresses how far the slope of the linear fit is from 1 (Smith and Rose 1995). For the example we generate synthetic data such that: Predicted ~ 0.7·True + N($\mu$=0.05,$\sigma^2$=0.1) with 1,500 True values uniformly distributed in [0,1].This decomposition of the prediction error is summarized by a stack barplot (left, not to scale; for another example see Fig. 3). Note that the True axis has been mean centered such that the uniform bias can be measured at x=0.

396

To fairly compare neural networks and the Maximum Likelihood Estimator, we restricted the parameter space to test the inference models to an inner domain of the explored parameter space. Indeed, during the training phase the neural networks

400    learn to not predict values outside the parameter space they have seen, whereas the

401    Maximum Likelihood cannot do so. Thus, when estimating parameters close to the

402    boundaries of the explored parameter space, the variance of the neural network

403    predictions artificially decreases due to the specific set up of the phylogeny

404    simulations.

## RESULTS

406    Overall, we find that deep learning methods show good performances as they

407    predict rates with an error comparable to MLE in most cases (Fig. 3). Some deep

408    learning methods even outperformed MLE (CNN and RNN with LTTs for the CRBD

409    model, see Fig. 3a,b). From theory, one would expect that lower total error of DL

410    methods can be explained by the bias-variance trade-off, which allows DL methods

411    to trade off bias against a lower total error (*e.g.* Pichler & Hartig, 2023). Looking at

412    the decomposition of the error, however, we do not find that the deep-learning

413    methods show greater bias than the MLE (*e.g.* DNN with SS vs. MLE for the CRBD

414    model in Fig. 3a,b). The only notable exception to the good performance of the deep

415    learning models was the GNN as well as the RNN and CNNs fed with the LTT for the

416    BISSE model, which all performed notably poorer than the MLE.

417

418

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig
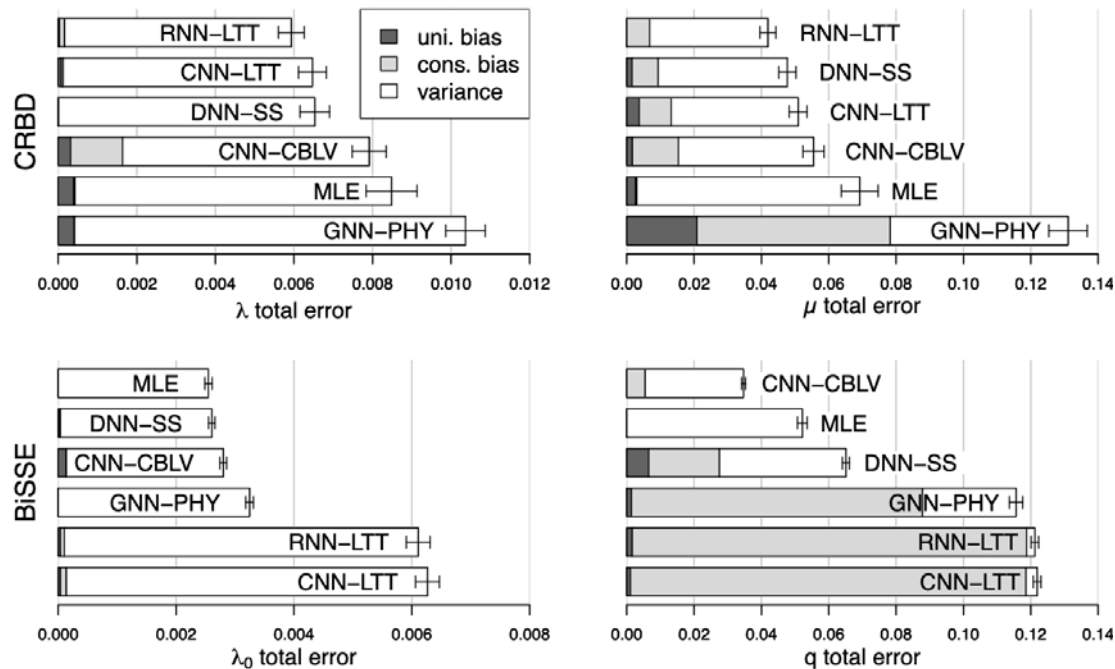
419

**Figure 3. Prediction error of deep learning methods and MLE for the CRBD (a-b) and the BiSSE (c-d) model.** From the CRBD we infer: a) the speciation rate ($\lambda$) and, b) the extinction rate ($\mu$). From the BiSSE model we infer: c) the speciation rate of state 0 ($\lambda_0$) and, d) the transition rate (q), the other parameters are constrained (see Methods for more details). The deep learning methods evaluated are the following: DNN with summary statistics (DNN-SS), CNN with LTTs (CNN-LTT), RNN with LTTs (RNN-LTT), CNN with encoded phylogeny (CNN-CBLV) and GNN with phylogeny graphs (GNN-PHY). The dataset contains 100,000 phylogenies for the CRBD model and 1,000,000 for the BiSSE model. The total summed-squared error is split in three terms: variance, uniform bias and consistency bias (for more details see Fig. 2 or Smith and Rose 1995). The error bars correspond to the 95% confidence interval.

432

Our second question was to understand how the optimal phylogeny representation depends on the complexity of the underlying diversification model.

435    Our findings show that for the CRBD model, the best models are based on the LTT

436    (CNN-LTT and RNN-LTT). Those models can even outperform the MLE (Fig. 3a,b).

437    In contrast, inference methods relying on more complex phylogeny representations

438    have a higher prediction error. This is especially clear for the two methods based on

439    the most complex representations: CNN with full phylogeny encodings and GNN with

440    phylogeny graphs. In sum, the more redundant information is added in the data, the

441    worse the deep learning methods perform.

442        For the BiSSE model, we found a different behavior. CNN with encoded

443    phylogenies and GNN with phylogeny graphs outperform simpler LTT-based models

444    (Fig. 3). To understand which information is missing in the LTT and thus is

445    responsible for the increased performance of the CNN with encoded phylogenies

446    and GNN with phylogeny graphs compared to LTT-based models, we provided

447    neural networks with data of more fine-grained increasing complexity. Specifically,

448    we compared a CNN fed with LTTs as a baseline (1) with a CNN and three

449    increasingly detailed encoded phylogenies (2-4). In the latter, we varied 2 features:

450    the number of tip states and the location of tip states in the phylogeny. Thus, the 3

451    options are: excluding all tip states information (2); including information about the tip

452    states, but not their location (3); including information about the tip states and their

453    location (4; see Fig. 4c). To sum up, we compare 4 options: (1) no information about

454    phylogeny's topology, tip state numbers and location, (2) information about topology

455    but not about tip states, (3) information about phylogeny topology, tip state numbers

456    but not about their location, and (4) information about phylogeny topology, tip state

457    numbers and their location.

458        This analysis confirms that the LTT alone is insufficient to infer rates accurately,

459    as can be seen by the CNN and RNN combined with LTTs being among the deep-

460  learning methods with the highest error (Fig. 3c,d). Moreover, the predictive

461  performance increases with each information that was added about tip states,

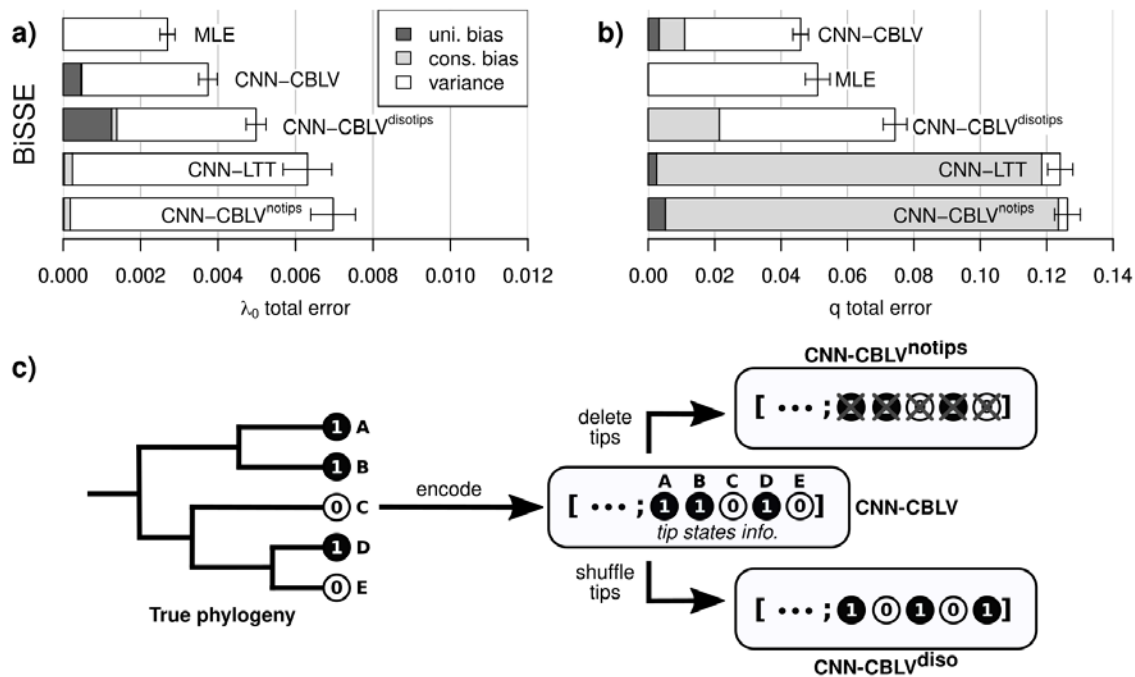462  suggesting that each of the intermediate steps considered loses information (Fig. 4).

463



464

465  **Figure 4. Prediction error of deep learning methods using different levels of**

466  **tip states information for the BiSSE model.** We infer: a) the speciation rate of

467  state 0 ($\lambda_0$) and, b) the transition rate (q), the other parameters are constrained (see

468  Methods for more details). c) The deep-learning methods evaluated are CNN with

469  encoded phylogenies: excluding tip states information (CNN-CLBV$^{notips}$), containing

470  tip states but randomly shuffled (CNN-CBLV$^{disotips}$), and containing tip states in order

471  (CNN-CBLV) and are presented in panel c). Additionally, we consider CNN with LTT

472  to have a reference of another deep learning method that does not use tip states

473  information, and MLE which is our baseline. The dataset for training and testing

474  neural networks contains 100,000 phylogenies. The total summed-squared error is

475    split in three terms: variance, uniform bias and consistency bias (for more details see

476    Fig. 2). The error bars correspond to the 95% confidence interval.

477

478    When looking at the influence of the neural network architecture, we see that

479    this choice interacts with the data representation: for simple data representations of

480    the phylogeny, CNN and RNN combined with LTTs achieved a similar prediction

481    accuracy in all cases (Fig. 3), and they were overall the best option for the CRBD

482    model, whereas for the more complicated BiSSE model, the even simpler DNN

483    architecture with summary statistics outperformed the LTT-based architectures,

484    whereas the CNN with encoded phylogenies performed best overall, including the

485    MLE (Figs. 3, 4). The information provided to this architecture is equivalent, although

486    differently structured, than the information provided to the GNN. The latter, however,

487    performed notably poorer.

## DISCUSSION

489    The goal of this study was to compare the ability of different deep learning

490    architectures to infer parameters of diversification models from reconstructed

491    phylogenies. Our main findings are that deep learning methods are surprisingly

492    accurate for this task and can even outperform the MLE in some cases. Looking in

493    more detail at the deep learning method, we find that the optimal phylogeny

494    representation and network architecture depend on the diversification model. For the

495    simple homogeneous CRBD model, we found that feeding the networks with the LTT

496    was the main factor that improved inference, presumably because the LTT is a

497    sufficient statistic and thus provides less redundant information to the networks. For

498    the inhomogeneous BISSE model, it was optimal to provide the data in its most

499    complex form, and the choice of the network architecture played a far greater role.

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

500    We speculate that these results will generalize: if simple sufficient statistics exist for

501    a model, these should be used and the choice of the network architecture is likely

502    secondary. If no simple sufficient data representation exists for a given model, more

503    complicated network architectures have to be used, and the choice of the network

504    architecture becomes more critical.

505        Note that, in line with these thoughts, the performance of the DNN with SS for

506    the BISSE model might be improved by adding more informative summary statistics

507    that encode tip state information. This idea is supported by our observation that

508    predictions of CNNs with CBLV outperform predictions of DNN with SS for the

509    transition rate by a great margin, suggesting that tip state information that are

510    contained in the CBLV encoding is critical for the parameter inference.

511                                    *Total Error and Bias*

512        We ranked the different inference methods primarily by their total error, which

513    sums up bias and variance (Smith and Rose 1995). This is in line with general

514    practice in machine learning. To achieve a lower error, machine learning methods

515    often trade off bias against variance (Pichler and Hartig 2023), whereas in statistics,

516    bias is often seen as more crucial than variance, because an unbiased estimator

517    allows a field to accumulate evidence over time (Shmueli 2010). Yet, given that there

518    are usually no independent replicates of a phylogeny, we find it defensible to use the

519    total error as the primary performance metric. Moreover, as we discuss below, DL

520    algorithms did not generally exhibit larger bias than the MLE, suggesting that the

521    problem of bias may be less severe than one might have anticipated.

522        Our results regarding total error on inferred parameters support earlier findings

523    of (Voznica et al. 2022), who also reported that deep learning methods can

524    outperform state-of-the-art methods when inferring model parameters from

525     phylogenies (BEAST 2 in his study, MLE in our study). Contrary to the general idea

526     of the bias-variance trade-off expectations, we did not find that deep learning

527     methods generate this performance advantage over the MLE by leveraging the bias-

528     variance trade-off. Indeed, several deep learning methods often had a lower total

529     error and a lower bias than MLE (where we expected that they traded off variance

530     against a greater bias to achieve lower error). To explain these findings, we

531     speculate that the MLE apparently also has some small-sample bias for the CRBD

532     model, especially for low net diversification rates ($i.e.$ $\mu \lesssim \lambda$).

533     *Explaining the poor performance of the GNNs*

534     Against our expectations, GNNs did not outperform other DL architectures and

535     specifically the CNN with encoded phylogenies. This is surprising given the generally

536     positive results in the machine learning literature about GNNs (Zhang et al. 2019),

537     and the fact that GNNs are fed with the most accurate and natural representation of

538     the data, which is the phylogeny itself. We speculate that our results can be

539     explained by two well-known limitations of GNNs: 'hop neighborhood' (Nikolentzos et

540     al. 2020) and 'over-smoothing' (Oono and Suzuki 2019).

541     First, 'hop neighborhood' refers to the fact that, in a GNN, at each new layer,

542     node attributes are updated by aggregating their neighbor attributes (Scarselli et al.

543     2009). Thus, after k layers, each node aggregates only nodes that can be reached in

544     k or less 'hops'. Thus, node attributes only contain information about local graph

545     structures but do not encode macroscopic graph information (Kriege et al. 2018). We

546     speculate that for inhomogeneous birth-death models, the most important info about

547     parameters related to the inhomogeneity is contained in subclades that have

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

548    diversified long ago, and thus typically have large graph distances. In such cases,

549    the "short-sightedness" of GNNs might be a major limitation.

550    We tested to increase the number of GNN convolutional layers from 2 to 5 (for

551    the exact architecture see Appendix Table S5-6) but it did not lead to a performance

552    improvement. Indeed, the naïve way to counteract this limitation is to increase the

553    number of GNN layers as the size of the 'k-hop neighborhood' increases with k the

554    number of layers. However, by doing so one will encounter the second limitation:

555    'over-smoothing' (Oono and Suzuki 2019): when stacking many layers in a GNN,

556    node attributes become indistinguishable from one another due to the cumulative

557    aggregations occurring at each layer. This phenomenon is referred to as 'over-

558    smoothing' (Li et al. 2018). Thus, adding layers in GNNs often results in a

559    deterioration of predictive performances (Li et al., 2018).

560    The combination of the two latter limitations might explain why 'classical' GNNs

561    failed to infer rates accurately from reconstructed phylogenies. We speculate that our

562    results regarding GNNs could be improved by considering specific GNN variants

563    designed to overcome these limitations (for over-smoothing see Li et al. 2019; and

564    for hop neighborhood see Nikolentzos et al. 2020), and encourage further research

565    in this direction.

566    *Designing Efficient Deep-Learning Methods for Rate Inference*

567    Our results reinforce previous analyses (Voznica et al. 2022; Lambert et al.

568    2022) suggesting that deep learning methods are a viable alternative to infer birth-

569    death rates from phylogenetic data. This is particularly interesting for cases where

570    likelihoods are not tractable or hard to compute. We refine these insights by

571    demonstrating that the performance of the deep learning methods may strongly

572    depend on data representation and network architecture. Regarding the complexity

573 of the phylogeny representation, we find that avoiding to provide redundant data

574 seems to help the algorithms. For instance, for homogeneous models, the topology

575 of the phylogeny does not influence the likelihood (Lambert and Stadler 2013). Thus,

576 the LTT provides a less redundant representation of the phylogeny while still

577 containing all information useful for the inference (sufficient statistic). For more

578 complex diversification models, in particular for models with intractable likelihoods, it

579 will likely be harder to find appropriate sufficient statistics, and deep learning

580 methods that use the whole phylogeny might be needed. For the BiSSE model, we

581 found that the fully encoded phylogeny (Voznica et al. 2022) containing phylogeny

582 topology and tip states information was the best choice among the representations

583 that we considered. In short, choosing a good representation of the data is important:

584 over-complicated representation either costs unnecessary computational power or

585 deteriorates the accuracy of the predictions for simple diversification models, while

586 choosing an over-simplified representation for a complex diversification model may

587 lead to poor predictions as the representation does not contain enough information

588 to predict parameters properly.

589   Once a suitable phylogeny representation has been found, the appropriate

590 neural network architecture has to be selected. First, the choice of the architecture is

591 restricted by the data type *e.g.*, a time series can go either with a CNN or an RNN

592 but do not fit with a DNN which is not designed to detect patterns or to process

593 temporal data. Secondly, the importance of the choice of the network architecture

594 among the remaining possibilities is dependent on the complexity of the phylogeny

595 representation: the more complex the representation, the more important the choice

596 of the architecture is. In other words, when data become more difficult to process,

597 more attention should be paid to the choice of the neural network architecture that

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

598    will process the data. In theory, network architectures could be further fine-tuned

599    regarding parameters such as the width and depth of the network layers, activation

600    functions, training parameters and regularization settings. Due to computational

601    reasons, we did not perform such a hyperparameter tuning, but rather set the

602    different classes of network architectures to sensible default settings.

603                                *Outlook and Further Research*

604            While our results are encouraging regarding the use of deep learning models

605    for parameter inference, they are at this stage only a proof of concept, because both

606    diversification models considered here are still relatively simple. Moreover, to keep

607    the analysis and training times tractable, we simplified the inhomogeneous BiSSE

608    model by assuming four constraints on the rate parameters. We hold that these

609    choices were appropriate for the purpose of this study, which was to establish

610    general principles of the problem, but it is clear that the practical benefit of this

611    methods would be to perform inference on more complicated diversification models

612    with intractable likelihoods and probably a higher number of parameters (such as,

613    *e.g.*, Hagen et al. 2021). How well the deep learning methods perform in this

614    scenario remains to be explored by further research. Results in (Voznica et al. 2022)

615    and Lambert et al. (2022) suggest that the deep learning methods still perform well

616    when increasing model complexity.

617            Moreover, inferring parameters with deep learning methods has a number of

618    drawbacks compared to traditional MLE approaches. Most importantly, training

619    networks requires generating a large set of simulated phylogenies under the

620    considered diversification model. To avoid that each user has to undergo this

621    process, it seems obvious that already trained neural networks must be provided, but

622    those would be likely still much larger and possibly less portable than current

623    statistical procedures.

## CONCLUSIONS

625    In conclusion, our study supports the idea that deep learning methods can be

626    used to infer diversification rates from reconstructed phylogenies. In our results, they

627    often showed similar and sometimes even better performance than the MLE. Thus,

628    deep learning methods offer a promising approach to likelihood-based statistical

629    inference, in particular for models whose likelihoods are intractable or hard to

630    compute. However, our results also provide a warning that deep learning methods

631    are very diverse and the choice of the method (*i.e.* the combination of a phylogeny

632    representation and a neural network architecture) must be adjusted to the model

633    considered. If that can be achieved, however, they may be instrumental in expanding

634    our options to perform statistical inference for complicated macroevoluationary

635    models.

## CONFLICT OF INTEREST

637    The authors declare no conflict of interest.

## SUPPLEMENTARY MATERIALS

639    The code to reproduce the results in this study is available at

640    https://github.com/ilajaait/phylo-inference-ml. Simulations data and the online

641    appendix are available at https://doi.org/10.5061/dryad.7h44j0zzq.

## ACKNOWLEDGEMENTS

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

645     MSc Thesis by Felix Gottschlich at the University of Regensburg that examined the

646     possibility to use GNNs for phylogenetic inference.

647

648

649

650                                    REFERENCES

651    Alfaro M.E., Santini F., Brock C., Alamillo H., Dornburg A., Rabosky D.L., Carnevale
652            G., Harmon L.J. 2009. Nine exceptional radiations plus high turnover explain
653            species diversity in jawed vertebrates. Proc. Natl. Acad. Sci. 106:13410–
654            13414.

655    Barnosky A.D., Matzke N., Tomiya S., Wogan G.O.U., Swartz B., Quental T.B.,
656            Marshall C., McGuire J.L., Lindsey E.L., Maguire K.C., Mersey B., Ferrer E.A.
657            2011. Has the Earth's sixth mass extinction already arrived? Nature. 471:51–
658            57.

659    Beaumont M.A. 2010. Approximate Bayesian Computation in Evolution and Ecology.
660            Annu. Rev. Ecol. Evol. Syst. 41:379–406.

661    Bengio Y. 2012. Neural Networks: Tricks of the Trade. Springer Berlin, Heidelberg.

662    Bukhari A.H., Raja M.A.Z., Sulaiman M., Islam S., Shoaib M., Kumam P. 2020.
663            Fractional Neuro-Sequential ARFIMA-LSTM for Financial Market Forecasting.
664            IEEE Access. 8:71326–71338.

665    Chan J., Perrone V., Spence J.P., Jenkins P.A., Mathieson S., Song Y.S. 2018. A
666            Likelihood-Free Inference Framework for Population Genetic Data using
667            Exchangeable Neural Networks. Adv. Neural Inf. Process. Syst. 31:8594–
668            8605.

669    Condamine F.L., Rolland J., Morlon H. 2013. Macroevolutionary perspectives to
670            environmental change. Ecol. Lett. 16:72–85.

671    Csilléry K., Blum M.G.B., Gaggiotti O.E., François O. 2010. Approximate Bayesian
672            Computation (ABC) in practice. Trends Ecol. Evol. 25:410–418.

673    Etienne R.S., Haegeman B., Stadler T., Aze T., Pearson P.N., Purvis A., Phillimore
674            A.B. 2012. Diversity-dependence brings molecular phylogenies closer to
675            agreement with the fossil record. Proc. R. Soc. B Biol. Sci. 279:1300–1309.

676    Etienne R.S., Rosindell J. 2012. Prolonging the Past Counteracts the Pull of the
677            Present: Protracted Speciation Can Explain Observed Slowdowns in
678            Diversification. Syst. Biol. 61:204.

679  Falbel D., Luraschi. 2019. torch: Tensors and Neural Networks with "GPU"
680       Acceleration. Available from https://torch.mlverse.org/docs/index.html.

681  Fey M., Lenssen J.E. 2019. Fast Graph Representation Learning with PyTorch
682       Geometric. ArXiv190302428 Cs Stat.

683  FitzJohn R.G. 2012. Diversitree: comparative phylogenetic analyses of diversification
684       in R. Methods Ecol. Evol. 3:1084–1092.

685  G. E. Hutchison. 1959. Homage to Santa Rosalia or Why Are There So Many Kinds
686       of Animals? | The American Naturalist: Vol 93, No 870. Available from
687       https://www.journals.uchicago.edu/doi/abs/10.1086/282070.

688  Gaston K.J., Blackburn T.M. 2000. Pattern and Process in Macroecology. John
689       Wiley & Sons, Ltd.

690  Gilmer J., Schoenholz S.S., Riley P.F., Vinyals O., Dahl G.E. 2017. Neural Message
691       Passing for Quantum Chemistry. Proc. 34th Int. Conf. Mach. Learn.:1263–
692       1272.

693  Graves A., Fernández S., Schmidhuber J. 2005. Bidirectional LSTM Networks for
694       Improved Phoneme Classification and Recognition. Artif. Neural Netw. Form.
695       Models Their Appl. – ICANN 2005.:799–804.

696  Hagen O., Flück B., Fopp F., Cabral J.S., Hartig F., Pontarp M., Rangel T.F.,
697       Pellissier L. 2021. gen3sis: A general engine for eco-evolutionary simulations
698       of the processes that shape Earth's biodiversity. PLOS Biol. 19:e3001340.

699  Hallinan N. 2012. The generalized time variable reconstructed birth–death process.
700       J. Theor. Biol. 300:265–276.

701  Hamilton W.L., Ying R., Leskovec J. 2018. Inductive Representation Learning on
702       Large Graphs. ArXiv170602216 Cs Stat.

703  Hartig F., Calabrese J.M., Reineking B., Wiegand T., Huth A. 2011. Statistical
704       inference for stochastic simulation models – theory and application. Ecol. Lett.
705       14:816–827.

706  Hillebrand H. 2004. On the Generality of the Latitudinal Diversity Gradient. Am. Nat.
707       163:192–211.

708  Hochreiter S., Schmidhuber J. 1997. Long Short-Term Memory. Neural Comput.
709       9:1735–1780.

710  Huang L., Ma D., Li S., Zhang X., Wang H. 2019. Text Level Graph Neural Network
711       for Text Classification. .

712  Innes M. 2018. Flux: Elegant machine learning with Julia. J. Open Source Softw.
713       3:602.

714  Kipf T.N., Welling M. 2017. Semi-Supervised Classification with Graph Convolutional
715       Networks. ArXiv160902907 Cs Stat.

716    Kriege N.M., Morris C., Rey A., Sohler C. 2018. A Property Testing Framework for
717         the Theoretical Expressivity of Graph Kernels. Proc. Twenty-Seventh Int. Jt.
718         Conf. Artif. Intell.:2348–2354.

719    Lambert A., Stadler T. 2013. Birth–death models and coalescent point processes:
720         The shape and probability of reconstructed phylogenies. Theor. Popul. Biol.
721         90:113–128.

722    Lambert S., Voznica J., Morlon H. 2022. Deep Learning from Phylogenies for
723         Diversification Analyses. :2022.09.27.509667.

724    LeCun Y., Bengio Y., Hinton G. 2015. Deep learning. Nature. 521:436–444.

725    Li G., Muller M., Thabet A., Ghanem B. 2019. DeepGCNs: Can GCNs Go As Deep
726         As CNNs? :9267–9276.

727    Li Q., Han Z., Wu X.-M. 2018. Deeper Insights into Graph Convolutional Networks
728         for Semi-Supervised Learning. .

729    Louca S., Pennell M.W. 2020. Extant timetrees are consistent with a myriad of
730         diversification histories. Nature. 580:502–505.

731    Maddison W.P., Midford P.E., Otto S.P. 2007. Estimating a Binary Character's Effect
732         on Speciation and Extinction. Syst. Biol. 56:701–710.

733    Maliet O., Hartig F., Morlon H. 2019. A model with many small shifts for estimating
734         species-specific diversification rates. Nat. Ecol. Evol. 3:1086–1092.

735    Mansimov E., Mahmood O., Kang S., Cho K. 2019. Molecular Geometry Prediction
736         using a Deep Generative Graph Neural Network. Sci. Rep. 9:20381.

737    Mittelbach G.G., Schemske D.W., Cornell H.V., Allen A.P., Brown J.M., Bush M.B.,
738         Harrison S.P., Hurlbert A.H., Knowlton N., Lessios H.A., McCain C.M.,
739         McCune A.R., McDade L.A., McPeek M.A., Near T.J., Price T.D., Ricklefs
740         R.E., Roy K., Sax D.F., Schluter D., Sobel J.M., Turelli M. 2007. Evolution and
741         the latitudinal diversity gradient: speciation, extinction and biogeography.
742         Ecol. Lett. 10:315–331.

743    Moen D., Morlon H. 2014. Why does diversification slow down? Trends Ecol. Evol.
744         29:190–197.

745    Morlon H. 2014. Phylogenetic approaches for studying diversification. Ecol. Lett.
746         17:508–525.

747    Morlon H., Parsons T.L., Plotkin J.B. 2011. Reconciling molecular phylogenies with
748         the fossil record. Proc. Natl. Acad. Sci. 108:16327–16332.

749    Morlon H., Robin S., Hartig F. 2022. Studying speciation and extinction dynamics
750         from phylogenies: addressing identifiability issues. Trends Ecol. Evol.

751  Nee S., Holmes E.C., May R.M., Harvey P.H., Lawton J.H., May R.M. 1994.
752      Extinction rates can be estimated from molecular phylogenies. Philos. Trans.
753      R. Soc. Lond. B. Biol. Sci. 344:77–82.

754  Nikolentzos G., Dasoulas G., Vazirgiannis M. 2020. k-hop graph neural networks.
755      Neural Netw. 130:195–205.

756  Oono K., Suzuki T. 2019. Graph Neural Networks Exponentially Lose Expressive
757      Power for Node Classification. .

758  Paradis E., Claude J., Strimmer K. 2004. APE: Analyses of Phylogenetics and
759      Evolution in R language. Bioinformatics. 20:289–290.

760  Pichler M., Hartig F. 2023. Machine learning and deep learning—A review for
761      ecologists. Methods Ecol. Evol. n/a.

762  Pontarp M., Bunnefeld L., Cabral J.S., Etienne R.S., Fritz S.A., Gillespie R., Graham
763      C.H., Hagen O., Hartig F., Huang S., Jansson R., Maliet O., Münkemüller T.,
764      Pellissier L., Rangel T.F., Storch D., Wiegand T., Hurlbert A.H. 2019. The
765      Latitudinal Diversity Gradient: Novel Understanding through Mechanistic Eco-
766      evolutionary Models. Trends Ecol. Evol. 34:211–223.

767  Pyron R.A., Burbrink F.T. 2013. Phylogenetic estimates of speciation and extinction
768      rates for testing ecological and evolutionary hypotheses. Trends Ecol. Evol.
769      28:729–736.

770  Rabosky D.L. 2009a. Ecological limits and diversification rate: alternative paradigms
771      to explain the variation in species richness among clades and regions. Ecol.
772      Lett. 12:735–743.

773  Rabosky D.L. 2009b. Heritability of Extinction Rates Links Diversification Patterns in
774      Molecular Phylogenies and Fossils. Syst. Biol. 58:629–640.

775  Rabosky D.L., Lovette I.J. 2008. Explosive Evolutionary Radiations: Decreasing
776      Speciation or Increasing Extinction Through Time? Evolution. 62:1866–1875.

777  Ricklefs R.E. 2007. Estimating diversification rates from phylogenetic information.
778      Trends Ecol. Evol. 22:601–610.

779  Rolland J., Condamine F.L., Jiguet F., Morlon H. 2014. Faster Speciation and
780      Reduced Extinction in the Tropics Contribute to the Mammalian Latitudinal
781      Diversity Gradient. PLOS Biol. 12:e1001775.

782  Saulnier E., Gascuel O., Alizon S. 2017. Inferring epidemiological parameters from
783      phylogenies using regression-ABC: A comparative study. PLOS Comput. Biol.
784      13:e1005416.

785  Scarselli F., Gori M., Tsoi A.C., Hagenbuchner M., Monfardini G. 2009. The Graph
786      Neural Network Model. IEEE Trans. Neural Netw. 20:61–80.

787  Schrider D.R., Kern A.D. 2018. Supervised Machine Learning for Population
788      Genetics: A New Paradigm. Trends Genet. 34:301–312.

Ismaël Lajaaiti, Sophia Lambert, Jakub Voznica, Hélène Morlon, Florian Hartig

789    Shmueli G. 2010. To Explain or to Predict? Stat. Sci. 25:289–310.

790    Silvestro D., Schnitzler J., Zizka G. 2011. A Bayesian framework to estimate
791        diversification rates and their variation through time and space. BMC Evol.
792        Biol. 11:311.

793    Smith E.P., Rose K.A. 1995. Model goodness-of-fit analysis using regression and
794        related techniques. Ecol. Model. 77:49–64.

795    Stadler T. 2011. Mammalian phylogeny reveals recent diversification rate shifts.
796        Proc. Natl. Acad. Sci. 108:6187–6192.

797    Stadler T. 2013. Recovering speciation and extinction dynamics based on
798        phylogenies. J. Evol. Biol. 26:1203–1219.

799    Voznica J., Zhukova A., Boskova V., Saulnier E., Lemoine F., Moslonka-Lefebvre M.,
800        Gascuel O. 2021. Deep learning from phylogenies to uncover the
801        transmission dynamics of epidemics. :2021.03.11.435006.

802    Voznica J., Zhukova A., Boskova V., Saulnier E., Lemoine F., Moslonka-Lefebvre M.,
803        Gascuel O. 2022. Deep learning from phylogenies to uncover the
804        epidemiological dynamics of outbreaks. Nat. Commun. 13:3896.

805    Wu Z., Pan S., Chen F., Long G., Zhang C., Yu P.S. 2021. A Comprehensive Survey
806        on Graph Neural Networks. IEEE Trans. Neural Netw. Learn. Syst. 32:4–24.

807    Yu Y., Si X., Hu C., Zhang J. 2019. A Review of Recurrent Neural Networks: LSTM
808        Cells and Network Architectures. Neural Comput. 31:1235–1270.

809    Zhang S., Tong H., Xu J., Maciejewski R. 2019. Graph convolutional networks: a
810        comprehensive review. Comput. Soc. Netw. 6:11.

811

812

813