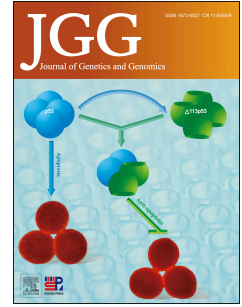# Journal Pre-proof

A critical evaluation of deep-learning based phylogenetic inference programs using simulated data sets

Yixiao Zhu, Yonglin Li, Chuhao Li, Xing-Xing Shen, Xiaofan Zhou

Please cite this article as: Zhu, Y., Li, Y., Li, C., Shen, X.-X., Zhou, X., A critical evaluation of deep-learning based phylogenetic inference programs using simulated data sets, *Journal of Genetics and Genomics*, https://doi.org/10.1016/j.jgg.2025.01.006.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1    A critical evaluation of deep-learning based phylogenetic inference

2                    programs using simulated data sets

3

4

5    **KEYWORDS**

6    molecular phylogenetics, deep learning, neural network, maximum likelihood

7

## Main text

Inferring phylogenetic trees from molecular sequences is a cornerstone of evolutionary

biology. Many standard phylogenetic methods (such as maximum-likelihood) rely on

explicit models of sequence evolution and thus often suffer from model

misspecification or inadequacy. The on-rising deep learning (DL) techniques offer a

powerful alternative. Deep learning employs multi-layered artificial neural networks to

progressively transform input data into more abstract and complex representations. DL

methods can autonomously uncover meaningful patterns from data, thereby bypassing

potential biases introduced by predefined features (Franklin, 2005; Murphy, 2012).

Recent efforts have aimed to apply deep neural networks (DNNs) to phylogenetics,

with a growing number of applications in tree reconstruction ( Suvorov et al., 2020;

Zou et al., 2020; Nesterenko et al., 2022; Smith and Hahn, 2023; Wang et al., 2023;),

substitution model selection (Abadi et al., 2020; Burgstaller-Muehlbacher et al., 2023)

and diversification rate inference (Voznica et al., 2022; Lajaaiti et al., 2023; Lambert et

al., 2023). In phylogenetic tree reconstruction, PhyDL (Zou et al., 2020) and

Tree_learning (Suvorov et al., 2020) are two notable DNN-based programs designed to

infer unrooted quartet trees directly from alignments of four amino acid (AA) and DNA

sequences, respectively. These two DNN programs offer pre-built models for

immediate analysis and the flexibility to train new models on user-defined data sets,

with benchmark tests showing performance comparable to or exceeding traditional

phylogenetic methods. However, DNNs encounter challenges as well. It is well known

that the effectiveness of a machine-learning algorithm heavily depends on the input-

30   data representation (Alzubaidi et al., 2021). Both PhyDL and Tree_learning are

31   supervised learning methods that need to be trained; however, in molecular

32   phylogenetics, simulation under explicit models of sequence evolution is the only

33   realistic source of training data. Therefore, while DNNs can outperform traditional

34   phylogenetic methods on benchmarks primarily consisting of simulated data

35   (Leuchtenberger et al., 2020), their performance might be compromised on biological

36   data, highlighting the need to understand the robustness of DL-based phylogenetic

37   methods when applied to out-of-distribution data. A recent study suggests that DNNs

38   struggle to match existing methods on datasets with branch-length and sequence-length

39   settings that differ significantly from those in the DNN training data (Zaharias et al.,

40   2022). In this study, we critically evaluated PhyDL and Tree_learning using simulated

41   data, highlighting critical constraints in current deep learning applications in molecular

42   phylogenetics and proposing suggestions to reduce the risk of inaccurate inferences in

43   practical use.

44

45   To investigate the strengths and weaknesses of PhyDL and Tree_learning, we first

46   designed a test to evaluate the performance of pre-built models provided by PhyDL and

47   Tree_learning, which are likely to be used out-of-the-box by the community (Fig. 1A).

48   Here, the test data sets were simulated under conditions deliberately selected to avoid

49   those well covered in the data used to train existing PhyDL and Tree_learning models.

50

51   PhyDL comes with three sets of pre-built DNN models, namely DNN1, DNN2, and

52      DNN3, differing in the simulation settings (e.g., heterogeneity level and branch length

53      distribution) of their training data. All these DNN models were trained with the long-

54      branch attraction (LBA) condition—also known as the Felsenstein zone—considered,

55      but relatively few long-branch repulsion (LBR) trees—those in the Farris zone—were

56      included in their training data (Table S1). These DNN models showed comparable or

57      superior performance than maximum-likelihood (ML) methods and other traditional

58      phylogenetic methods on data simulated from LBA-susceptible trees (Zou et al., 2020).

59      We first followed the LBA benchmark design from Zou et al. (2020) to evaluate the

60      DNN models on data sets simulated under LBA/LBR conditions (Figs.S1–S5;

61      Supplementary File Text S1). To further examine the performance of DNN models, we

62      used datasets containing AA alignments simulated with progressively complex models

63      (LG+F+$\Gamma$, LG+C20+F+$\Gamma$, and LG+C60+F+$\Gamma$) (Wang et al., 2018) based on LBA and

64      LBR trees (Fig. 1B). We also analyzed these data sets using the ML phylogenetic

65      program IQ-TREE for comparison. For data simulated under LBA condition, none of

66      the three PhyDL models had an accuracy above 50%, while all ML phylogenetic

67      models performed substantially better than DNN models (Figs. 1C, S6). On LBR data

68      sets, the accuracies were 100% for DNN1 and DNN2 but nearly 0% for DNN3, while

69      the accuracies of ML models were between 65.00% and 99.97%. Additionally, we

70      investigated an unexpected performance of DNN3 regarding tree type, noting a high

71      frequency of "incorrect tree – other" on LBA data and "incorrect tree – LBR-I" on LBR

72      data (Figs. 1B, 1C, S7; Supplementary File Text S2). Furthermore, our investigation of

73      the performance of DNN models during their training processes revealed that DNN3 is

74  more vulnerable to model fluctuations during training compared to DNN1 and DNN2

75  (Fig. S8; Supplementary File Text S3). Overall, our results suggest that the DNN

76  models provided by PhyDL are less accurate than ML phylogenetic models on LBA

77  data.

78

79  We then employed the approach developed by Trost et al. (2023) to quantify the

80  disparity between our test data and the pre-built DNN training data. In brief, a Gradient

81  Boosted Trees (GBT) classifier was trained on one data set (e.g., the DNN1 training

82  data) and then applied on another (e.g., our LG+F+Γ LBA test data) to calculate a

83  balanced accuracy (BACC) (Brodersen et al., 2010) value (0 to 1.0, higher values

84  indicate greater differences) which reflects the difference between the two data sets

85  (Materials and Methods in Supplementary Text). As a result, the GBT analyses

86  accurately distinguished each of our test datasets from the training data of pre-built

87  DNN models (with BACC values above 0.99), indicating substantial differences

88  between our test data and the original training data (Table S2; Fig. S9).

89

90  In Suvorov et al. (2020), the Tree_learning CNN model trained on gapped data

91  performed much better than traditional phylogenetic methods on gapped alignments,

92  likely because it can extract additional phylogenetic signals from gaps (Suvorov et al.,

93  2020). Specifically, gaps in the training and test data were all simulated by INDELible,

94  and the phylogenetic signals carried by these indel gaps are expected to match the

95  underlying phylogenies. However, real data often contain random gaps (e.g., due to

96    incomplete genome assemblies, partial gene models, or errors in multiple sequence

97    alignments) that may add noise to phylogenetic analyses. To investigate whether the

98    inclusion of random gaps might impact the performance of pre-built CNN models, we

99    first simulated an ungapped data set (NOGAP.ori) and a gapped data set (INDEL.ori)

100   following the procedures of Suvorov et al., and then created two additional data sets,

101   NOGAP.extragaps and INDEL.extragaps, by introducing random gaps into the first two

102   data sets, respectively (Fig. 1D). We applied the CNN model trained on ungapped data

103   (referred to as "CNN.NOGAP.Ori") on NOGAP.ori, and the model trained on gapped

104   data (referred to as "CNN.INDEL.Ori") on the three data sets with gaps. For

105   comparison, we analyzed the data using IQ-TREE under two modes, including "IQ-

106   TREE.Standard", where gaps are treated as missing data with no information, and "IQ-

107   TREE.Recoded", where gaps are recognized as the fifth character in addition to A, T,

108   C, and G. Our evaluation of IQ-TREE and Tree-learning models on NOGAP.ori yielded

109   similar results to those reported by Suvorov et al. (Fig. S10; Supplementary File Text

110   S4). On INDEL.ori, which includes only indel gaps, CNN.INDEL.Ori and IQ-

111   TREE.Recoded achieved much higher accuracy compared to their performance on

112   NOGAP.ori, while the accuracy of IQ-TREE.Standard remained unchanged. However,

113   after random gaps were introduced into the test data, CNN.INDEL.Ori became

114   substantially less accurate on NOGAP.extragaps and INDEL.extragaps, while the two

115   IQ-TREE models had nearly the same accuracies (Fig. 1E). Additionally, we also tested

116   CNN.NOGAP.Ori, CNN.INDEL.Ori and IQ-TREE models across various branch-

117   length combinations (Fig. S11; Supplementary File Text S5). Our results indicated that

118    the inclusion of random noisy gaps in our test data impaired the performance of existing

119    Tree_learning models, rendering them less accurate than IQ-TREE. CNN models

120    trained on indel gaps likely misinterpreted random gaps as informative characters,

121    extracting misleading signals as a result.

122

123    In addition to offering pre-built models, both PhyDL and Tree_learning allow users to

124    train new models using custom data. Therefore, we tested if the performance of PhyDL

125    and Tree_learning on difficult data sets could be improved by targeted training using

126    data simulated under the same challenging conditions, either independently or in

127    conjunction with the original training data (Fig. 1F). Importantly, we examined the

128    performance of the new models under both target and non-target conditions to better

129    understand the outcome of this model optimization strategy.

130

131    We first examined if targeted training can produce PhyDL models with better

132    accuracies under LBA/LBR conditions. We simulated additional LBA and LBR data

133    sets under LG+C20+F+$\Gamma$. These data sets were used to train new DNN models,

134    including DNN_LBA10K (trained on 10,000 LBA alignments), DNN_LBR10K

135    (trained on 10,000 LBR alignments), and DNN_60K (training on 30,000 LBA and

136    30,000 LBR alignments). Additionally, we trained DNN_160K using the DNN_60K

137    data along with 100,000 alignments simulated similarly to the original DNN3 training

138    data. These new DNN models were applied on the same test data in our first test (Figs.

139    1G, S12). DNN_LBA10K demonstrated significantly improved performance on LBA

140    data (accuracy exceeding 95%), but showed notable bias when applied to LBR data

141    (Figs. 1G, S12). A similar trend was observed with DNN_LBR10K, which made

142    accurate inferences under LBR conditions, but its accuracy dropped on LBA data. We

143    also found that adding more simulated alignments from a denser sampling of branch

144    length combinations did not improve the performance of DNN_LBA10K and

145    DNN_LBR10K (Fig. S13). DNN_60K and DNN_160K demonstrated a more balanced

146    performance across LBA and LBR conditions, performing between DNN_LBA10K

147    and DNN_LBR10K on both types of test data (Figs. 1G, S12). Notably, DNN_160K

148    performed substantially better than DNN3 on our test data, and its accuracy on the

149    original DNN3 test data ("testing3_mixed") was still close to that of DNN3 itself (Table

150    S3). Unlike DNN3, errors made by all new DNN models were mostly of the expected

151    "incorrect tree – LBA" on LBA data sets, and distributed more evenly between two

152    types of incorrect trees on LBR data sets (Fig. 1G).

153

154    For Tree_learning, we trained two new CNN models, CNN.NOGAP.Extragaps and

155    CNN.INDEL.Extragaps, on data sets simulated under the NOGAP.extragaps and

156    INDEL.extragaps schemes, respectively, and tested their performance on NOGAP and

157    INDEL data sets with or without random gaps (Fig. 1H). Generally, the best-performing

158    model for each data set was the one whose training data were simulated in the same

159    way as the test data. CNN.INDEL.Extragaps had considerably higher accuracy than

160    CNN.INDEL.Ori    on    both    NOGAP.extragaps    (63.43%    vs.    38.57%)    and

161    INDEL.extragaps (84.54% vs. 70.16%) (Fig. 1H; Supplementary File Text S6). We

162  further enhanced the performance of CNN.INDEL.Ori on random gaps by conducting

163  additional training with alignments simulated under the INDEL.extragaps scheme. The

164  fine-tuned model (CNN.Fine-tuned) demonstrated significantly higher accuracy than

165  the original CNN.INDEL.Ori model on NOGAP.extragaps (68.65% vs. 38.57%) and

166  INDEL.extragaps (84.83% vs. 70.16%), while maintaining nearly identical

167  performance to CNN.INDEL.Ori on the ungapped dataset NOGAP.ori (69.42% vs.

168  69.51%) and exhibiting slightly reduced accuracy on INDEL.ori (85.89% vs. 88.17%)

169  (Fig. 1H). Additionally, we tested if targeted training can produce Tree-learning models

170  with better performance under LBA/LBR conditions (TableS4; Supplementary File

171  Text S7). Our results indicate that our targeted optimization effort has successfully

172  enhanced the model's capability to handle random gaps, albeit with a slight compromise

173  on its performance on phylogenetically informative indels.

174

175  In conclusion, our critical evaluation of PhyDL and Tree_learning provides practical

176  evidence that ML methods generally outperformed DNN programs, especially when

177  data properties were unfamiliar to the pre-built DNN models. While DNN performance

178  can be enhanced by training new models tailored to these specific conditions, this

179  comes at the cost of reduced generalizability. Additionally, several challenges must be

180  addressed before DL-based phylogenetic methods can compete with traditional

181  approaches: first, existing DL methods like PhyDL and Tree_learning can only infer

182  quartet trees instead of full phylogenies (in cases of more than four sequences); second,

183  DL methods need to demonstrate their ability to learn patterns from empirical MSAs;

184 third, few DL methods can successfully infer branch lengths (Supplementary File Text

185 S8).

186

187 Based on our results, we recommend assessing the differences between training and test

188 data prior to conducting tree inference to avoid potential pitfalls in phylogenetic

189 reconstruction with DNN programs (Fig. 1I). Our examination of the difference

190 between the pre-built DNN training data and our test data with GBT classifier may

191 serve as an example (Table S2). Overall, our evaluation provides valuable insights for

192 the future development of DNN-based phylogenetic methods and offers practical

193 guidance for their application.

194

195 **Data availability**

196 All gene alignments and gene trees are available on the figshare repository

197 (https://doi.org/10.6084/m9.figshare.23617767 – please note that this link will become

198 active upon publication). For access during the peer review process, please use the

199 private link (https://figshare.com/s/2c806c4ab1ebc43472c6).

200

201 **Conflict of interest**

202 The authors declare no competing financial interests.

203

## Acknowledgments

## References

Abadi, S., Avram, O., Rosset, S., Pupko, T., Mayrose, I., 2020. ModelTeller: model selection for optimal phylogenetic reconstruction using machine learning. Mol. Biol. Evol. 37, 3338–3352.

Alzubaidi, L., Zhang, J., Humaidi, A.J., Al-Dujaili, A., Duan, Y., Al-Shamma, O., Santamaría, J., Fadhel, M.A., Al-Amidie, M., Farhan, L., 2021. Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. J. Big Data 8, 53.

Brodersen, K.H., Ong, C.S., Stephan, K.E., Buhmann, J.M., 2010. The balanced accuracy and its posterior distribution. Presented at the 2010 20th international conference on pattern recognition, IEEE, pp. 3121–3124.

Burgstaller-Muehlbacher, S., Crotty, S.M., Schmidt, H.A., Reden, F., Drucks, T., von Haeseler, A., 2023. ModelRevelator: Fast phylogenetic model estimation via deep learning. Mol. Phylogenet. Evol. 188, 107905.

Franklin, J., 2005. The elements of statistical learning: data mining, inference and prediction. Math. Intell. 27, 83–85.

Lajaaiti, I., Lambert, S., Voznica, J., Morlon, H., Hartig, F., 2023. A comparison of deep learning architectures for inferring parameters of diversification models from extant phylogenies. bioRxiv 2023.03.03.530992.

Lambert, S., Voznica, J., Morlon, H., 2023. Deep learning from phylogenies for diversification analyses. Syst. Biol. 72, 1262–1279.

Leuchtenberger, A.F., Crotty, S.M., Drucks, T., Schmidt, H.A., Burgstaller-Muehlbacher, S., 2020. Distinguishing felsenstein zone from farris zone using neural networks. Mol. Biol. Evol. 37, 3632–3641.

236 Murphy, K.P., 2012. Machine learning: a probabilistic perspective. MIT press.

237 Nesterenko, L., Boussau, B., Jacob, L., 2022. Phyloformer: towards fast and accurate phylogeny
238 estimation with self-attention networks. bioRxiv 2022.06.24.496975 .

239 Smith, M.L., Hahn, M.W., 2023. Phylogenetic inference using generative adversarial networks.
240 Bioinformatics 39, btad543.

241 Suvorov, A., Hochuli, J., Schrider, D.R., 2020. Accurate inference of tree topologies from multiple
242 sequence alignments using deep learning. Syst. Biol. 69, 221–233.

243 Trost, J., Haag, J., Hohler, D., Nesterenko, L., Jacob, L., Stamatakis, A., Boussau, B., n.d. Simulations
244 of sequence evolution: how (un)realistic they really are and why. bioRxiv 2023.07.11.548509.

245 Voznica, J., Zhukova, A., Boskova, V., Saulnier, E., Lemoine, F., Moslonka-Lefebvre, M., Gascuel, O.,
246 2022. Deep learning from phylogenies to uncover the epidemiological dynamics of outbreaks. Nat.
247 Commun. 13, 3896.

248 Wang, H.-C., Minh, B.Q., Susko, E., Roger, A.J., 2018. Modeling site heterogeneity with posterior
249 mean site frequency profiles accelerates accurate phylogenomic estimation. Syst. Biol. 67, 216–
250 235.

251 Wang, Z., Sun, J., Gao, Y., Xue, Y., Zhang, Y., Li, K., Zhang, W., Zhang, C., Zu, J., Zhang, L., 2023. Fusang:
252 a framework for phylogenetic tree inference via deep learning. Nucleic Acids Res. 51, 10909–10923.

253 Zaharias, P., Grosshauser, M., Warnow, T., 2022. Re-evaluating deep neural networks for phylogeny
254 estimation: the issue of taxon sampling. J. Comput. Biol. 29, 74–89.

255 Zou, Z., Zhang, H., Guan, Y., Zhang, J., 2020. Deep residual neural networks resolve quartet
256 molecular phylogenies. Mol. Biol. Evol. 37, 1495–1507.

257

258

259 Yixiao Zhu

260 College of Agriculture and Biotechnology and Centre for Evolutionary & Organismal

261 Biology, Zhejiang University, Hangzhou 310058, China

262 Yonglin Li, Chuhao Li

263 Guangdong Laboratory for Lingnan Modern Agriculture, Guangdong Province Key

264 Laboratory of Microbial Signals and Disease Control, Integrative Microbiology

265 Research Centre, South China Agricultural University, Guangzhou 510642, China

266 Xing-Xing Shen[*]

267 College of Agriculture and Biotechnology and Centre for Evolutionary & Organismal

268 Biology, Zhejiang University, Hangzhou 310058, China

269                                      Xiaofan Zhou[*]

270    Guangdong Laboratory for Lingnan Modern Agriculture, Guangdong Province Key

271    Laboratory of Microbial Signals and Disease Control, Integrative Microbiology

272    Research Centre, South China Agricultural University, Guangzhou 510642, China

273

274    [*] Corresponding authors

275    *E-mail addresses:* xingxingshen@zju.edu.cn (X.-X. Shen);

276    xiaofan_zhou@scau.edu.cn (X. Zhou)

277

278

## Figure Legends

**Fig. 1.** Evaluation of deep learning-based phylogenetic inference programs on simulated datasets. **A**: Schematics of performance evaluations for pre-built models conducted in this study. **B**: Illustrations of the three possible inference outcomes for a four-sequence AA alignment under LBA or LBR conditions, as inferred by IQ-TREE and PhyDL models. **C**: Proportions of different types of trees inferred by IQ-TREE and PhyDL models from test data sets simulated under LBA or LBR conditions. **D**: Schematics of the procedures for simulating the four distinct DNA test datasets used for tree inference with various IQ-TREE and Tree_learning models. **E**: Proportions of correctly inferred trees for various IQ-TREE and Tree_learning models on four simulated test datasets. **F**: Schematics of the performance evaluations for custom-trained models conducted in this study. **G**: Performance of optimized PhyDL models on simulated protein sequence alignments across various branch length combinations. **H**: Performance of new Tree_learning models optimized for the presence of random gaps on simulated DNA sequence alignments. **I**: Schematics of a potential solution to mitigate risks arising from differences between training and testing data.

**A**



**B**

Possible inference outcomes for data simulated under LBA conditions

Possible inference outcomes for data simulated under LBR conditions



**C**



**D**



**E**



**F**



**G**



**H**



**I**