

# Applications of Machine Learning in Phylogenetics

Yu K. Mo<sup>1</sup>, Matthew W. Hahn<sup>1,2</sup>, and Megan L. Smith<sup>2,3,\*</sup>

<sup>1</sup>Department of Computer Science and <sup>2</sup>Department of Biology, Indiana University, Bloomington, IN 47405

<sup>3</sup>Department of Biological Sciences, Mississippi State University, Starkville, MS 39762

\*Corresponding Author: ms4438@msstate.edu

## Abstract

Machine learning has increasingly been applied to a wide range of questions in phylogenetic inference. Supervised machine learning approaches that rely on simulated training data have been used to infer tree topologies and branch lengths, to select substitution models, and to perform downstream inferences of introgression and diversification. Here, we review how researchers have used several promising machine learning approaches to make phylogenetic inferences. Despite the promise of these methods, several barriers prevent supervised machine learning from reaching its full potential in phylogenetics. We discuss these barriers and potential paths forward. In the future, we expect that the application of careful network designs and data encodings will allow supervised machine learning to accommodate the complex processes that continue to confound traditional phylogenetic methods.

## 1 Introduction

Phylogenetics aims to elucidate the evolutionary relationships among species. In recent decades, owing to rapid growth in the availability of genomic data, phylogenetic analysis has been able to use hundreds to thousands of loci (Delsuc et al., 2005). Using whole genomes, or even near-whole genomes, may allow for a more comprehensive view of the evolutionary events shaping species (Scornavacca et al., 2020). However, the accuracy of inference may be compromised when using such large datasets, as even small biases can be magnified many-fold. Biases in phylogenetics are often due to unmodeled heterogeneity in the evolutionary process, including heterogeneity across time, sites, genes, or lineages (Kapli et al., 2020). These processes may arise either individually or in combination, presenting challenges in subsequent analyses.

Recently, machine learning techniques have been used across fields, demonstrating impressive power in uncovering intricate relationships from data that contains extensive heterogeneity. Notable examples include successful applications in image classification (Krizhevsky et al., 2017), language models (Devlin et al., 2019), protein structure prediction (Jumper et al., 2021), and population

15 genetics (Schrider & Kern, 2018). Machine learning is comprised of two fundamental paradigms—  
16 supervised and unsupervised approaches. Supervised learning relies on the availability of labeled  
17 training data, where the true underlying state or value of the data is known. In phylogenetics and  
18 related fields, large amounts of labeled training data are generally unavailable, so simulations are  
19 often used to generate such data. The primary objective of supervised machine learning is to learn  
20 a function that can map input data to appropriate outputs. Within supervised learning, there  
21 are two primary tasks: classification and regression. While classification aims to predict discrete  
22 labels or categories, regression predicts continuous-valued outputs. In contrast, unsupervised  
23 learning operates without the need for labeled data, focusing instead on discerning underlying  
24 structures or patterns in the input data. Unsupervised approaches include tasks such as clustering  
25 and dimensionality reduction. Notably, deep learning is a specialized subset of machine learning  
26 that leverages neural networks (NNs) with many layers (hence "deep"). Some NN architectures  
27 are adept at automatically extracting hierarchical features from raw data, obviating the need for  
28 manual feature engineering—a significant advantage over traditional machine learning methods.

29 In the context of phylogenetics, machine learning algorithms are extremely flexible, both  
30 with regards to the structuring of input data, and the data used for training. Further, machine  
31 learning approaches can learn complex relationships from input implicitly, without the need for an  
32 explicit model. This facilitates the application of machine learning to complex models, especially  
33 scenarios in which standard likelihood and Bayesian inference may be intractable. Given the lack  
34 of analytical phylogenetic solutions that can be reasonably applied to large genomic datasets,  
35 machine learning offers the promise of moving beyond conventional methods.

36 Despite the promise that machine learning in general has for addressing many biological prob-  
37 lems, there is uncertainty about its superiority over conventional approaches in many applications  
38 to phylogenetics. While a growing number of papers have applied machine learning to multiple  
39 problems in the field, researchers have not yet seen a clear advantage to such approaches. Here,  
40 we review recent applications of machine learning to different tasks in phylogenetics, examining  
41 their limitations and strengths. We attempt to provide a general overview of the types of machine  
42 learning approaches that have been used—and those that could be used—in the hope that future  
43 work will bring the promise of machine learning to fruition.

## 44 **2 Tree Reconstruction**

45 Reconstructing evolutionary relationships among taxa is a central goal in evolutionary biology.  
46 A phylogenetic tree is composed of two primary components: a topology and a set of branch  
47 lengths. The topology serves as a representation of the hierarchical evolutionary relationships  
48 among species. The branch lengths represent evolutionary change, measured either in absolute  
49 time, in the number of nucleotide substitutions, or in other units. This section reviews machine  
50 learning approaches for inferring both components of phylogenetic trees.

### 51 **2.1 Topology inference**

52 Perhaps the most natural framing of the problem of topology inference is to use supervised  
53 machine learning approaches for classification, since the goal is to predict a discrete output  
54 (topology) from sequence data. Recall that supervised machine learning approaches require

55 labeled training data, which are generally unavailable in phylogenetics. Because of this, in  
56 most phylogenetic applications simulations are performed under each model of interest prior to  
57 inference, and these simulated data are used to train the machine learning network. When the  
58 goal is topology inference, the model space includes, at a minimum, the number of possible tree  
59 topologies. With as few as ten taxa, there are more than two million unrooted topologies, making  
60 it infeasible to use such approaches to infer tree topologies for even moderate numbers of taxa.  
61 The challenges associated with a large state-space of topologies are not unique to machine learning  
62 approaches: even conventional methods have difficulties in inferring trees for large numbers of  
63 species (Felsenstein, 1978b; Roch, 2006). To circumvent this problem, researchers have used three  
64 different types of approaches in order to apply machine learning to phylogenetic inference (Figure  
65 1). Here we review these approaches and the specific models that have been used.

### 66 2.1.1 Quartet-based methods

67 The first machine learning approaches in phylogenetics used quartet-based methods. In general,  
68 quartet-based methods involve extracting sets of four taxa from the full dataset, building trees  
69 for each set of four taxa, and then constructing a phylogeny from these quartet trees using one  
70 of several quartet amalgamation approaches, such as quartet puzzling (Bryant & Steel, 2001;  
71 Reaz et al., 2014; Snir & Satish, 2012). Because there are only three possible topologies for an  
72 unrooted quartet, such approaches are not plagued by the need to consider a very large state-space  
73 of topologies. Quartet-based methods therefore provide efficient inference algorithms that are  
74 scalable to very large datasets.

75 Several supervised learning approaches have been used to infer quartet trees. Zou et al. (2020)  
76 used a residual neural network, which takes as input one-hot encoded amino acid sequences.  
77 Machine learning algorithms generally require that input data are numerical, and one-hot encoding  
78 can be used to represent categorical variables. In this application, each site was represented by  
79 twenty channels, with each channel corresponding to an amino acid. For an individual site, the  
80 channel corresponding to the amino acid present in the position is set to one, while all other  
81 channels are set to zero. One-hot encoding may be more appropriate than integer-encoding (in  
82 which each amino acid would be treated as an integer between 1 and 20), since it avoids implicit  
83 ordered relationships among states. In Zou et al.’s approach, models were trained on amino acid  
84 sequences simulated on large, random trees, which were then pruned to subsets of four taxa.  
85 Both site and time heterogeneity were included in the simulations; additionally, the training data  
86 intentionally included a sizable proportion of trees susceptible to long-branch attraction, to ensure  
87 that a large number of difficult examples were included. When benchmarked against existing  
88 inference approaches, the residual network predictors consistently delivered better results with less  
89 computational time (not including training time), especially when dealing with several cases that  
90 confound existing methods—such as long branch attraction and heterotachy. By combining their  
91 approach with a quartet amalgamation approach, these authors were able to infer larger species  
92 trees with moderate accuracy.

93 In a similar approach, Suvorov et al. (2020) used a convolutional neural network (CNN) that  
94 takes integer-encoded nucleotide alignments as input. Simulations were carried out in a manner  
95 similar to Zou et al. (2020), under several substitution models and including site heterogeneity.  
96 In addition, these authors trained networks both with and without gaps. In the absence of gaps,

97 the CNN generally performed as well as or better than traditional approaches. On datasets that  
98 included gaps, the CNN substantially outperformed traditional approaches, likely because it better  
99 utilized this significant source of phylogenetic signal. The CNN initially exhibited reduced accuracy  
100 in some zones of branch length space (e.g., the Felsenstein zone; (Felsenstein, 1978a)). However,  
101 when more training data were included from these regions the CNN was able to outperform other  
102 approaches, highlighting the importance of carefully considering where to put effort in training  
103 such models.

104 Both of the methods described above treat alignments as images. While this approach to  
105 representing data has been found to be powerful in population genetics (Flagel et al., 2019), there  
106 are several limitations in the context of phylogenetics. For example, when inferring relationships  
107 among taxa, we would like the order in which sequences are included in the model to be irrelevant  
108 (a property referred to as “permutation equivariant”). However, most network architectures do not  
109 perform in this way. Zou et al. (2020) accommodated this behavior by including all permutations of  
110 the alignment when training, but such an approach increases the compute time and memory needed  
111 to train a neural network. Solís-Lemus et al. (2023) address this issue using a symmetry-preserving  
112 long short-term memory (LSTM) recurrent neural network (RNN). By avoiding the need to include  
113 permutations of the training alignments, they substantially improved compute times and memory  
114 usage compared to Zou et al. (2020).

115 Even though NNs can be very efficient for inferring quartet trees, considering larger trees  
116 remains prohibitive—the three approaches described above still must rely on quartet-amalgamation  
117 approaches to build larger trees. Additionally, as with all supervised machine learning, accuracy  
118 is likely limited in cases where the training data is not reflective of real data. Zaharias et al.  
119 (2022) explored these limitations by comparing the networks from Zou et al. (2020) to standard  
120 approaches on larger trees and on test datasets with higher rates of nucleotide evolution and/or  
121 shorter alignment lengths. They found that the neural networks only outperformed traditional  
122 approaches when the goal was to infer a quartet tree from relatively long amino acid sequences  
123 simulated under model conditions very similar to those used for training. Furthermore, when larger  
124 trees were considered, traditional approaches outperformed the combination of neural networks  
125 and quartet amalgamation. Machine learning approaches are therefore severely limited by their  
126 inability to directly infer trees from larger numbers of taxa, as well as by the specifics of the data  
127 used in training.

### 128 **2.1.2 Distance-based methods**

129 Rather than using machine learning to directly infer trees from sequence alignments, it is possible  
130 to instead infer evolutionary distances, which can then be used as input to standard distance-based  
131 approaches. Although often scoffed at by modern phylogeneticists, distance-based approaches  
132 such as neighbor joining (Saitou & Nei, 1987) are in fact guaranteed to infer the correct tree in  
133 most of parameter space, as long as distances are accurately inferred. In addition, they are much  
134 more accurate than maximum likelihood in the presence of high amounts of incomplete lineage  
135 sorting (Liu & Edwards, 2009; Mendes & Hahn, 2018). Therefore, it makes sense to apply machine  
136 learning to the task of accurately inferring distances.

137 Nesterenko et al. (2022) used self-attention networks to infer evolutionary distances for up to  
138 100 species. Their model encapsulates alignment in a pairwise way, introducing a representation

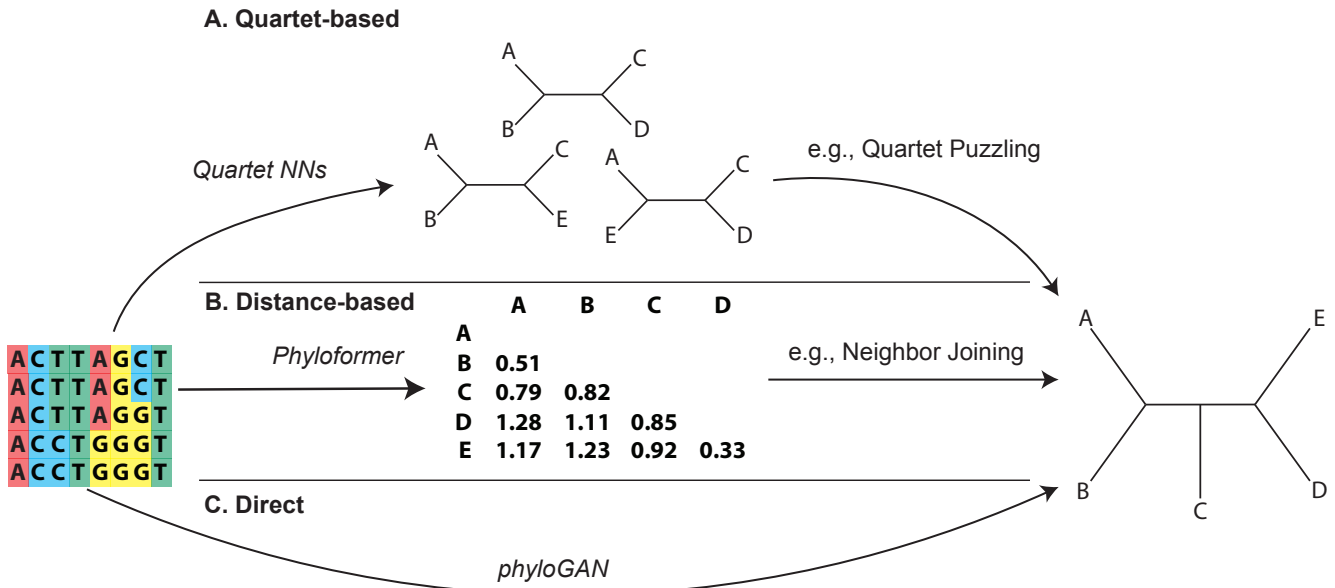


Figure 1. Methods for topology inference using machine learning. A. Quartet-based methods infer one of the three topologies possible with unrooted quartets. Trees from each quartet are inferred with NNs; a collection of such trees are then fed into existing quartet amalgamation algorithms (e.g. Quartet Puzzling) to infer a larger phylogeny. B. Distance-based methods estimate pairwise distances using NNs (e.g. Phyloformer). Distances are combined using standard methods (e.g. Neighbor Joining) to reconstruct trees. C. Direct methods infer a tree directly from an alignment using NNs (e.g. phyloGAN).

139 for each pair with the attention mechanism. The process entails an iterative sharing of information,  
140 first across sites within each pair (referred to as site-level attention) and subsequently across pairs  
141 within each site (termed pair-level attention). Such an approach is permutation-equivariant, and  
142 accommodates alignments of varying sizes. After inferring distances, these authors used neighbor  
143 joining for tree construction. Their approach outperformed traditional distance-based approaches,  
144 and was competitive with (and much faster than) maximum likelihood when training and testing  
145 data included similar numbers of species.

146 In a similar approach, Bhattacharjee and Bayzid (2020) used autoencoders to impute missing  
147 values in distance matrices. Alternatively, Jiang et al. (2023) use a CNN for phylogenetic  
148 placement—placing sequences from individual genes onto trees that may have been inferred using  
149 different genomic regions. In this case they inferred evolutionary distances for these new sequences,  
150 and then used a distance-based algorithm to place the new sequences on the tree (Balaban et al.,  
151 2022). Inferring evolutionary distances reframes phylogenetic inference as a regression problem,  
152 rather than as a classification problem. This reframing makes it possible to scale machine learning  
153 approaches to larger trees.

### 154 **2.1.3 Direct methods**

155 In maximum likelihood and Bayesian approaches to phylogenetic inference, the large number  
156 of possible topologies is accommodated by using heuristic searches to explore tree space; such  
157 approaches could also be used for direct inference of tree topologies from sequence data in machine  
158 learning contexts. Generative adversarial networks (GANs) consist of a generator, which aims to  
159 produce realistic data, and a discriminator, which aims to distinguish real and fake data (Goodfellow  
160 et al., 2020). Recently, Smith and Hahn (2023) proposed phyloGAN. phyloGAN consists of a  
161 generator, which generates topologies and branch lengths, and a CNN-based discriminator, which  
162 attempts to distinguish alignments simulated under these topologies and branch lengths from  
163 empirical (real) alignments. Ideally, at the end of training, it should be virtually impossible  
164 to distinguish simulated and empirical alignments. Once this level of accuracy is achieved, the  
165 topology that underpins the simulated data is considered to be the inferred topology. phyloGAN  
166 was tested on up to fifteen species, and a version incorporating gene tree heterogeneity was tested  
167 on six species. While phyloGAN worked well with small numbers of species (up to ten), it was  
168 computationally intensive, and several metrics indicated issues during training. Future work may  
169 explore alternative approaches for heuristically exploring model spaces using machine learning  
170 frameworks, including approaches covered in the next section here.

### 171 **2.1.4 Improving steps in topology inference**

172 Machine learning approaches have been used to assist standard phylogenetic approaches for  
173 topology inference. For example, machine learning approaches have been used to improve heuristic  
174 searches for tree topologies. Azouri et al. (2021) used a random forest (RF) regressor to predict  
175 likelihood scores for subtree-prune-regraft (SPR) moves, a standard and important step in heuristic  
176 tree searches. Given a starting topology, their network could accurately predict the change in  
177 likelihood associated with different SPR moves, which suggests that such an approach could be  
178 used to limit search space and therefore to reduce the computational requirements for heuristic  
179 searches. In a follow-up paper, Azouri et al. (2023) used reinforcement learning as an alternative



180 to traditional heuristic search algorithms. By allowing for suboptimal moves that, nonetheless,  
181 improved the final outcome of the search, this approach out-competed greedy search strategies.

182 Machine learning approaches have also been used to guide researchers in their decisions about  
183 which standard approaches to use for topological inference. Leuchtenberger et al. (2020) developed  
184 a feed-forward neural network to classify alignments as belonging to the Farris (Siddall, 1998) or  
185 Felstenstein zone (Felsenstein, 1978a; Huelsenbeck & Hillis, 1993). They based their choice to  
186 use maximum parsimony (in the Farris Zone) or maximum likelihood (in the Felsenstein zone) on  
187 the predictions of this neural network. Using this approach resulted in higher overall accuracy  
188 compared to always using either maximum parsimony or maximum likelihood.

## 189 **2.2 Branch length inference**

190 In addition to a tree topology, most researchers are also interested in inferring the branch  
191 lengths of a tree. However, few studies have successfully inferred branch lengths using machine  
192 learning. While it may seem that this regression problem should be easier than the classification  
193 problem of inferring topologies, the size of the output vector depends on the number of edges in  
194 the tree—there are  $2n - 2$  branches in a rooted tree with  $n$  tips. The dependence on the number  
195 of tips complicates the use of machine learning approaches.

196 Suvorov and Schrider (2022) employed both a CNN and a multilayer perceptron (MLP)  
197 to infer branch lengths on fixed tree topologies with four or eight taxa. For the CNN-based  
198 approach, they adapted a previously proposed architecture (Suvorov et al., 2020). Instead of a  
199 classification task, the model was restructured for regression, aiming to predict all branch lengths  
200 simultaneously. Meanwhile, the MLP was fed with feature vectors derived from site pattern  
201 frequencies present within each alignment. Notably, the predictions generated by their models  
202 showed slightly superior accuracy compared to maximum likelihood estimates. Despite these  
203 promising results, there remains a degree of skepticism regarding the scalability of machine learning  
204 to infer branch lengths, especially when considering more species. Nevertheless, the flexibility  
205 of machine learning approaches with respect to the types of input data that can be considered  
206 offers many interesting possibilities. For instance, in the future such methods could facilitate the  
207 integration of heterogeneous fossil data in estimating time-calibrated trees.

## 208 **3 Other kinds of phylogenetic inferences**

209 In addition to phylogenetic tree inference, machine learning approaches have been applied  
210 to both upstream and downstream tasks in phylogenetics. Prior to tree inference using many  
211 approaches (e.g., Bayesian inference, maximum likelihood, neighbor joining) it is necessary to infer  
212 a sequence substitution model. After tree inference, researchers are often interested in detecting  
213 and quantifying discordance, testing for introgression, and inferring macroevolutionary parameters.  
214 Below, we review recent machine learning approaches to these upstream and downstream tasks.

### 215 **3.1 Substitution models**

216 It is crucial to select a suitable substitution model for accurate phylogenetic inference from  
217 sequence data, as it has long been known that misspecified models can lead to inaccurate estimates  
218 of trees (Buckley, 2002; Sanderson, 2002) and branch lengths (Abadi et al., 2019). Existing

219 methods for model selection infer the model that provides the best fit to the data, using one of  
220 several criteria. Popular criteria include likelihood ratio tests (LRTs), Akaike information criteria  
221 (AIC), corrected AIC (AICc), Bayesian information criteria (BIC), and decision theory (DT).  
222 However, these criteria rely on assumptions that are often not met in phylogenetics, and there  
223 is a lack of consensus regarding which criteria are the most appropriate (Abadi et al., 2019).  
224 Additionally, substitution model choice tends to impact branch length estimates more-so than  
225 topology inference (Abadi et al., 2019), but no criteria to-date have been designed to select the  
226 model best-suited for branch length inference. Finally, using these criteria to perform substitution  
227 model selection is computationally expensive, as it requires computation of the likelihood. Here  
228 we discuss two recent machine learning approaches that attempt to address these gaps.

229 ModelTeller (Abadi et al., 2020) is a machine learning approach that uses an RF regressor to  
230 rank 24 potential substitution models according to their accuracy in downstream branch length  
231 inference. Features fed into the model included over 50 summary statistics that can be broadly  
232 categorized into four primary groups: features inherent to the alignment, features drawn from  
233 an approximated tree inferred through a distance-based method, parameters inferred under a  
234 parameter-rich substitution model, and sequence similarity within certain subsets. ModelTeller’s  
235 primary distinction compared to traditional approaches lies in selecting a substitution model that  
236 improves accuracy in branch length inference. This leads to improved performance in terms of the  
237 accuracy of branch length estimates under the models selected using ModelTeller compared to  
238 models selected using more standard approaches, particularly on datasets simulated under realistic  
239 models. Additionally, ModelTeller was substantially faster than standard methods.

240 A later model, ModelRevelator (Burgstaller-Muehlbacher et al., 2023) aims to infer the  
241 correct generating model of nucleotide substitution using two neural networks. The first network,  
242 NNmodelfinder, takes as input a set of statistics calculated from pairwise alignments and predicts  
243 the best substitution model from a set of six possible models. The second network, NNalphafind,  
244 takes as input base composition profiles and predicts whether a site homogeneous model is  
245 appropriate or not. If a site homogeneous model is not appropriate, then NNalphafind estimates  
246 the  $\alpha$  parameter of a model with  $\Gamma$ -distributed rate heterogeneity among sites. Used together,  
247 these networks can predict the best substitution model for a given sequence alignment, whether  
248 rate heterogeneity should be included, and, when rate heterogeneity is included, the  $\alpha$  parameter  
249 to use in downstream inference. ModelRevelator performed comparably to maximum likelihood  
250 combined with substitution model selection under BIC as implemented in IQ-TREE (Minh et al.,  
251 2020), with substantially reduced computation times on large alignments.

252 Both ModelTeller and ModelRevelator are designed to select a substitution model that is  
253 suitable for inference; however, each uses different criteria for assessing suitability. ModelTeller is  
254 particularly focused on identifying a model that results in the most accurate estimates of branch  
255 lengths. The primary objective of ModelRevelator is to select the best substitution model and  
256 estimate the  $\alpha$  parameter when the best model includes rate heterogeneity. One can therefore use  
257 both methods together on a single dataset.

## 258 **3.2 Levels of discordance**

259 Gene tree topologies often differ from the species tree topology due to several biological factors,  
260 including incomplete lineage sorting, introgression, and gene duplication and loss (Maddison, 1997).



261 Two recent studies used deep learning to estimate the amount of discordance in phylogenetic  
262 datasets (Rosenzweig et al., 2022; Zhang et al., 2023). Rosenzweig et al. (2022) used several  
263 approaches, including a deep neural network (DNN), to estimate the amount of discordance in  
264 four-taxon datasets using a set of summary statistics calculated from alignments and inferred gene  
265 trees. Estimates from their DNN were more accurate than relying on inferred gene trees alone to  
266 estimate discordance, particularly when branch lengths were long. In addition to their network  
267 for estimating the amount of discordance, they introduced a network for inferring the quartet  
268 species tree topology from the same set of statistics. Similarly, Zhang et al. (2023) used CNNs to  
269 estimate the proportion of all different possible topologies for four and five-taxon datasets from  
270 multiple sequence alignments. Their CNN, called ERICA, was able to accurately infer topology  
271 proportions. The authors then used these inferred proportions to try to infer introgression and to  
272 identify potentially introgressed genomic windows. The ability of these approaches to estimate  
273 the proportions of quartet topologies more accurately than standard pipelines—which rely on  
274 inferred gene trees alone—offers promise for improving many quartet-based methods for species  
275 tree inference, as these generally assume that quartet frequencies are accurately estimated from  
276 input gene trees (Mirarab & Warnow, 2015).

### 277 3.3 Introgression

278 Most machine learning approaches for studying introgression have focused on population-scale  
279 data, rather than phylogenetic data. For example, Schrider et al. (2018) used ExtraTrees classifiers  
280 to detect introgressed regions between closely related species, while Ray et al. (2023) used a CNN  
281 and image segmentation for a similar task. Similarly, Gower et al. (2021) developed a CNN to  
282 detect adaptive introgression given data from three closely related populations or species. Several  
283 recent papers have also addressed introgression from a phylogenetic perspective using machine  
284 learning.

285 Two recent studies used supervised machine learning to determine whether there was evidence  
286 for reticulation in a dataset. Blischak et al. (2021) used a CNN to detect various types of  
287 reticulation in four-taxon trees, including hybrid speciation and introgression. Their CNN took  
288 as input mean and minimum values of  $d_{xy}$  (a measure of sequence divergence) between sets of  
289 populations. They compared HyDe-CNN to an RF classifier trained on several phylogenetic  
290 statistics for detecting introgression and found that HyDe-CNN had increased power. In a similar  
291 approach, Burbrink and Gehara (2018) trained a neural network to distinguish a bifurcating species  
292 tree from models including reticulation between two parent clades and one clade with a putative  
293 reticulate history. As input, their network takes pairwise distances between all sequences in the  
294 phylogeny (55 sequences from three clades). Their network had moderate power to distinguish  
295 among models with and without reticulations. When applied to their empirical data, the model  
296 supported a reticulate history for a clade in which reticulation was also inferred using SNaQ  
297 (Solís-Lemus & Ané, 2016). Most recently, Hibbins and Hahn (2022) used supervised machine  
298 learning to distinguish speciation and introgression histories. Under many regions of parameter  
299 space, gene trees and site patterns matching the introgression history can become more common  
300 than those matching the species tree, challenging many traditional approaches to species tree  
301 inference. By using several summary statistics calculated from gene trees, Hibbins and Hahn  
302 were able to accurately infer the speciation history for rooted three-taxon trees, even in regions

303 of parameter space where traditional approaches fail. While powerful, these approaches have  
304 primarily focused on four or fewer taxa. Future work may expand machine learning approaches to  
305 study introgression on larger trees.

### 306 **3.4 Diversification rates**

307 In addition to the kinds of inferences described above, recent studies have attempted to use  
308 inferred phylogenies for downstream inference of diversification rates. One challenge in any such  
309 analysis is determining the optimal way to encode phylogenetic trees. To address this issue,  
310 Voznica et al. (2022) introduced the compact bijective ladderized vector (CBLV), an encoding  
311 of phylogenetic trees that can be used as input into a CNN. They trained a CNN that took as  
312 input the CBLV to infer parameters of phylodynamic birth-death models and to perform model  
313 selection. They compared the performance of this CNN to a feed-forward neural network trained  
314 on summary statistics calculated from phylogenetic trees. Both networks were able to accurately  
315 infer parameters and distinguish among phylodynamic models. Lambert et al. (2023) used similar  
316 networks to infer speciation and turnover rates under a constant rate birth-death (CRBD) model  
317 and to infer the parameters of a binary state speciation and extinction (BiSSE) model. Lajaaiti  
318 et al. (2023) compared these networks to several other networks for inferring diversification  
319 parameters. They trained an additional CNN and RNN on lineage through time (LTT) plots.  
320 They also trained a graph neural network (GNN) that took phylogenies encoded as graphs directly  
321 as input. Under the CRBD model, the RNN and CNN trained on LTT plots outperformed the  
322 network trained on CBLV encodings. However, these same networks performed poorly under  
323 the BiSSE model, likely because the LTT plots did not include additional information about tip  
324 states, which was included in the other networks. Perhaps surprisingly, the GNN performed poorly  
325 across both models. These approaches highlight the importance of carefully choosing network  
326 architectures and data encodings for the task at hand.

## 327 **4 Discussion**

328 Recent progress has revealed the promise of machine learning in phylogenetics. However,  
329 inferences have often been limited to relatively small trees and relatively limited regions of parameter  
330 space. Moving forward, careful considerations of training datasets, network architectures, and  
331 data encodings will facilitate the use of machine learning to address fundamental challenges in  
332 phylogenetic inference.

333 Supervised machine learning requires a labeled training set. In the context of phylogenetics,  
334 however, we do not have labels for many real-world examples—we therefore have to simulate data.  
335 Despite attempts to simulate realistic data across a wide range of parameter space, it is inevitable  
336 that biases will creep in. For example, training data generated under one substitution model may  
337 not generalize to empirical datasets that were produced under a different model. In order to avoid  
338 such biases, we could take the computationally costly step of generating synthetic data across  
339 increasingly large sets of parameters. However, even when researchers attempt to consider a broad  
340 range of relevant parameters, there will inevitably be mismatches between training and empirical  
341 data, potentially leading to poor generalization to unseen data. To develop more robust networks,  
342 widely used techniques such as dropout, regularization, and ensemble methods can be employed.

343 Alternatively, noise can be added to training data to improve generalization (as is done with image  
344 augmentation). In the context of phylogenetics, adding noise could involve masking regions of the  
345 alignment during training. Alternatively, techniques from domain adaptation have emerged as  
346 promising solutions. Domain adaptation aims to develop networks that are robust to differences  
347 between the distribution of training data and the distribution of target or empirical data. Mo and  
348 Siepel (2023) used domain adaptation to make more accurate inferences of recombination rates  
349 and selection coefficients in the presence of domain differences. Their approach used adversarial  
350 domain-invariant feature extraction, which incorporates an additional layer to prevent the model  
351 from extracting features that differ between the training and target data. Such an approach  
352 promotes the extraction of domain-invariant features, and could be used to make robust inferences  
353 in phylogenetics.

354 A major intended advantage of machine learning is that, once trained, models can be applied to  
355 new datasets with minimal computational expenses. Even though a trained model makes inferences  
356 almost instantaneously, training remains computationally expensive. Ideally, trained networks  
357 would be applicable across a wide range of empirical datasets, but this is limited by the details  
358 of the training data used and the choice of network architectures. Specifically, many network  
359 architectures (e.g., most CNNs) are not invariant to dataset size. In other words, only datasets with  
360 the exact dimensions of the training data can be analyzed. However, in phylogenetics, datasets  
361 may vary in size due to different alignment lengths or different numbers of taxa. This challenge  
362 has been addressed in population genetics through padding (Flagel et al., 2019), and by designing  
363 appropriate network architectures that are size invariant (Sanchez et al., 2021). Some network  
364 architectures employed in phylogenetics have accommodated variable input sizes (Nesterenko et al.,  
365 2022), and moving forward this should be a central goal. To facilitate the reuse of networks in new  
366 empirical systems, techniques from transfer learning could also be used. Specifically, supervised  
367 transfer learning can be useful when limited training data are available from a new domain. For  
368 example, a network that has already been trained on data from one domain can be reused in a  
369 related, but distinct, domain. Supervised transfer learning and limited simulations in the new  
370 domain can be used to generate a robust network with reduced computational expenses compared  
371 to training the network from scratch. Combined, these approaches may facilitate more efficient  
372 uses of supervised machine learning in phylogenetic contexts.

373 Another major consideration is how to encode input data for neural networks. Most commonly,  
374 encoded alignments (Suvorov & Schrider, 2022; Suvorov et al., 2020; Zou et al., 2020), or summary  
375 statistics (Abadi et al., 2020; Burgstaller-Muehlbacher et al., 2023) have been used as input. When  
376 using encoded alignments, a primary challenge is scalability to longer alignments or more taxa.  
377 This is especially pertinent as available genomic data continues to grow. Encoded alignments can  
378 also pose challenges to network reusability, as discussed above. Alternatively, the input can be  
379 represented with summary statistics that are explanatory features drawn from alignments and trees  
380 for the task at hand. However, selecting a good set of features relies on prior knowledge, and the  
381 choice of statistics can heavily impact inference. Alternative strategies for representing alignments  
382 have been proposed, using attention mechanisms (Burgstaller-Muehlbacher et al., 2023; Nesterenko  
383 et al., 2022; Rao et al., 2021) or language models (Lupo et al., 2022). Such approaches can lead  
384 to networks that can accept variable input sizes, and are capable of incorporating relationships  
385 among sites and lineages simultaneously. It is also essential to develop a suitable representation for  
386 phylogenetic trees. Several efforts in this direction have been made, from explanatory summary

387 statistics (Voznica et al., 2022), to embeddings such as the CBLV (Voznica et al., 2022), to  
388 graphical representations in GNNs (Lajaaiti et al., 2023). While early uses are promising, these  
389 encodings have only been explored for a small set of inferential tasks, and it is unclear which  
390 encodings will prove most useful over a wider range of questions.

391 The promise of supervised machine learning is to efficiently consider a wide range of the complex  
392 processes that complicate phylogenetic inference. To date, most machine learning approaches for  
393 tree inference have largely ignored heterogeneity introduced by incomplete lineage sorting (ILS),  
394 gene duplication and loss, and introgression (though several exceptions have been described here).  
395 While standard phylogenetic approaches also have trouble modeling this heterogeneity, machine  
396 learning shows potential to include multiple of these processes at once. For example, if machine  
397 learning approaches can be used to more accurately infer quartet frequencies in the presence of  
398 these processes (as demonstrated in the case of ILS by (Rosenzweig et al., 2022; Zhang et al.,  
399 2023)) then the accuracy of phylogenetic trees could be improved. Moving forward, we expect that  
400 creative network architectures, data encodings, and task designs will facilitate the use of machine  
401 learning to improve phylogenetic inferences in the presence of complex processes that cannot be  
402 accommodated by standard approaches.

#### 403 **4.1 Acknowledgements**

404 This work was supported by a National Science Foundation (NSF) grant to M.W.H. (DEB-  
405 1936187).

#### 406 **References**

- 407 Abadi, S., Avram, O., Rosset, S., Pupko, T., & Mayrose, I. (2020). ModelTeller: Model selection  
408 for optimal phylogenetic reconstruction using machine learning. *Molecular Biology and*  
409 *Evolution*, *37*(11), 3338–3352.
- 410 Abadi, S., Azouri, D., Pupko, T., & Mayrose, I. (2019). Model selection may not be a mandatory  
411 step for phylogeny reconstruction. *Nature Communications*, *10*(1), 934.
- 412 Azouri, D., Abadi, S., Mansour, Y., Mayrose, I., & Pupko, T. (2021). Harnessing machine learning  
413 to guide phylogenetic-tree search algorithms. *Nature Communications*, *12*(1), 1983.
- 414 Azouri, D., Granit, O., Albuquerque, M., Mansour, Y., Pupko, T., & Mayrose, I. (2023). The  
415 tree reconstruction game: Phylogenetic reconstruction using reinforcement learning. *arXiv*.  
416 <https://doi.org/10.48550/arXiv.2303.06695>
- 417 Balaban, M., Jiang, Y., Roush, D., Zhu, Q., & Mirarab, S. (2022). Fast and accurate distance-based  
418 phylogenetic placement using divide and conquer. *Molecular Ecology Resources*, *22*(3),  
419 1213–1227.
- 420 Bhattacharjee, A., & Bayzid, M. S. (2020). Machine learning based imputation techniques for  
421 estimating phylogenetic trees from incomplete distance matrices. *BMC Genomics*, *21*(1),  
422 497.
- 423 Blischak, P. D., Barker, M. S., & Gutenkunst, R. N. (2021). Chromosome-scale inference of hybrid  
424 speciation and admixture with convolutional neural networks. *Molecular Ecology Resources*,  
425 *21*(8), 2676–2688.

- 426 Bryant, D., & Steel, M. (2001). Constructing optimal trees from quartets. *Journal of Algorithms*,  
427 38(1), 237–259.
- 428 Buckley, T. R. (2002). Model misspecification and probabilistic tests of topology: Evidence from  
429 empirical data sets. *Systematic Biology*, 51(3), 509–523.
- 430 Burbrink, F. T., & Gehara, M. (2018). The biogeography of deep time phylogenetic reticulation.  
431 *Systematic Biology*, 67(5), 743–755.
- 432 Burgstaller-Muehlbacher, S., Crotty, S. M., Schmidt, H. A., Reden, F., Drucks, T., & von Haeseler,  
433 A. (2023). ModelRevelator: Fast phylogenetic model estimation via deep learning. *Molecular*  
434 *Phylogenetics and Evolution*, 188, 107905.
- 435 Delsuc, F., Brinkmann, H., & Philippe, H. (2005). Phylogenomics and the reconstruction of the  
436 tree of life. *Nature Reviews Genetics*, 6(5), 361–375.
- 437 Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional  
438 transformers for language understanding. *Proceedings of the 2019 Conference of the North*  
439 *American Chapter of the Association for Computational Linguistics: Human Language*  
440 *Technologies, Volume 1 (Long and Short Papers)*, 4171–4186.
- 441 Felsenstein, J. (1978a). Cases in which parsimony or compatibility methods will be positively  
442 misleading. *Systematic Zoology*, 27(4), 401–410.
- 443 Felsenstein, J. (1978b). The number of evolutionary trees. *Systematic Zoology*, 27(1), 27–33.
- 444 Flagel, L., Brandvain, Y., & Schrider, D. R. (2019). The unreasonable effectiveness of convolutional  
445 neural networks in population genetic inference. *Molecular Biology and Evolution*, 36(2),  
446 220–238.
- 447 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., &  
448 Bengio, Y. (2020). Generative adversarial networks. *Communications of the ACM*, 63(11),  
449 139–144.
- 450 Gower, G., Picazo, P. I., Fumagalli, M., & Racimo, F. (2021). Detecting adaptive introgression in  
451 human evolution using convolutional neural networks. *eLife*, 10, e64669.
- 452 Hibbins, M. S., & Hahn, M. W. (2022). Distinguishing between histories of speciation and  
453 introgression using genomic data. *bioRxiv*. <https://doi.org/10.1101/2022.09.07.506990>
- 454 Huelsenbeck, J., & Hillis, D. (1993). Success of phylogenetic methods in the four-taxon case.  
455 *Systematic Biology*, 42(3), 247–264.
- 456 Jiang, Y., Blaban, M., Zhu, Q., & Mirarab, S. (2023). DEPP: Deep learning enables extending  
457 species trees using single genes. *Systematic Biology*, 72(1), 17–34.
- 458 Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool,  
459 K., Bates, R., Žídek, A., Potapenko, A., Bridgland, A., Meyer, C., Kohl, S., Ballard, A.,  
460 Cowie, A., Romera-Paredes, B., Nikolov, S., Jain, R., Adler, J., ... Hassabis, D. (2021).  
461 Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873), 583–589.
- 462 Kapli, P., Yang, Z., & Telford, M. J. (2020). Phylogenetic tree building in the genomic age. *Nature*  
463 *Reviews Genetics*, 21(7), 428–444.
- 464 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2017). ImageNet classification with deep convolu-  
465 tional neural networks. *Communications of the ACM*, 60(6), 84–90.
- 466 Lajaaiti, I., Lambert, S., Voznica, J., Morlon, H., & Hartig, F. (2023). A comparison of deep learning  
467 architectures for inferring parameters of diversification models from extant phylogenies.  
468 *bioRxiv*. <https://doi.org/10.1101/2023.03.03.530992>



- 469 Lambert, S., Voznica, J., & Morlon, H. (2023). Deep learning from phylogenies for diversification  
470 analyses. *Systematic Biology*, syad044.
- 471 Leuchtenberger, A. F., Crotty, S. M., Drucks, T., Schmidt, H. A., Burgstaller-Muehlbacher, S.,  
472 & von Haeseler, A. (2020). Distinguishing felsenstein zone from farris zone using neural  
473 networks. *Molecular Biology and Evolution*, 37(12), 3632–3641.
- 474 Liu, L., & Edwards, S. V. (2009). Phylogenetic analysis in the anomaly zone. *Systematic Biology*,  
475 58(4), 452–460.
- 476 Lupo, U., Sgarbossa, D., & Bitbol, A.-F. (2022). Protein language models trained on multiple  
477 sequence alignments learn phylogenetic relationships. *Nature Communications*, 13(1), 6298.
- 478 Maddison, W. (1997). Gene trees in species trees. *Systematic Biology*, 46(3), 523–536.
- 479 Mendes, F. K., & Hahn, M. W. (2018). Why concatenation fails near the anomaly zone. *Systematic*  
480 *Biology*, 67(1), 158–169.
- 481 Minh, B., Schmidt, H., Chernomor, O., Schrempf, D., Woodhams, M., von Haeseler, A., & Lanfear,  
482 R. (2020). IQ-TREE 2: New models and efficient methods for phylogenetic inference in the  
483 genomic era. *Molecular Biology and Evolution*, 37(5), 1530–1534.
- 484 Mirarab, S., & Warnow, T. (2015). ASTRAL-II: Coalescent-based species tree estimation with  
485 many hundreds of taxa and thousands of genes. *Bioinformatics*, 31(12), i44–i52.
- 486 Mo, Z., & Siepel, A. (2023). Domain-adaptive neural networks improve supervised machine learning  
487 based on simulated population genetic data. *bioRxiv*. [https://doi.org/10.1101/2023.03.01.  
488 529396](https://doi.org/10.1101/2023.03.01.529396)
- 489 Nesterenko, L., Boussau, B., & Jacob, L. (2022). Phyloformer: Towards fast and accurate phylogeny  
490 estimation with self-attention networks. *bioRxiv*. <https://doi.org/10.1101/2022.06.24.496975>
- 491 Rao, R. M., Liu, J., Verkuil, R., Meier, J., Canny, J., Abbeel, P., Sercu, T., & Rives, A. (2021).  
492 MSA transformer. *International Conference on Machine Learning*, 8844–8856.
- 493 Ray, D. D., Flagel, L., & Schrider, D. R. (2023). IntroUNET: Identifying introgressed alleles via  
494 semantic segmentation. *bioRxiv*. <https://doi.org/10.1101/2023.02.07.527435>
- 495 Reaz, R., Bayzid, M. S., & Rahman, M. S. (2014). Accurate phylogenetic tree reconstruction from  
496 quartets: A heuristic approach. *PLOS ONE*, 9(8), e104008.
- 497 Roch, S. (2006). A short proof that phylogenetic tree reconstruction by maximum likelihood is  
498 hard. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 3(1), 92–94.
- 499 Rosenzweig, B. K., Kern, A. D., & Hahn, M. W. (2022). Accurate detection of incomplete lineage  
500 sorting via supervised machine learning. *bioRxiv*. <https://doi.org/10.1101/2022.11.09.515828>
- 501 Saitou, N., & Nei, M. (1987). The neighbor-joining method: A new method for reconstructing  
502 phylogenetic trees. *Molecular Biology and Evolution*, 4(4), 406–425.
- 503 Sanchez, T., Cury, J., Charpiat, G., & Jay, F. (2021). Deep learning for population size history  
504 inference: Design, comparison and combination with approximate bayesian computation.  
505 *Molecular Ecology Resources*, 21(8), 2645–2660.
- 506 Sanderson, M. J. (2002). Estimating absolute rates of molecular evolution and divergence times: A  
507 penalized likelihood approach. *Molecular Biology and Evolution*, 19(1), 101–109.
- 508 Schrider, D. R., Ayroles, J., Matute, D. R., & Kern, A. D. (2018). Supervised machine learning  
509 reveals introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLoS*  
510 *Genetics*, 14(4), e1007341.
- 511 Schrider, D. R., & Kern, A. D. (2018). Supervised machine learning for population genetics: A  
512 new paradigm. *Trends in Genetics*, 34(4), 301–312.



- 513 Scornavacca, C., Delsuc, F., & Galtier, N. (2020). *Phylogenomics in the genomic era*. Open access  
514 book. <https://hal.inria.fr/PGE>
- 515 Siddall, M. (1998). Success of parsimony in the four-taxon case: Long-branch repulsion by likelihood  
516 in the Farris zone. *Cladistics*, *14*(3), 209–220.
- 517 Smith, M. L., & Hahn, M. W. (2023). Phylogenetic inference using generative adversarial networks.  
518 *Bioinformatics*, *39*(9), btad543.
- 519 Snir, S., & Satish, R. (2012). Quartet MaxCut: A fast algorithm for amalgamating quartet trees.  
520 *Molecular Phylogenetics and Evolution*, *62*(1), 1–8.
- 521 Solís-Lemus, C., & Ané, C. (2016). Inferring phylogenetic networks with maximum pseudolikelihood  
522 under incomplete lineage sorting. *PLoS Genetics*, *12*(3), e1005896.
- 523 Solís-Lemus, C., Yang, S., & Leonardo, Z.-N. (2023). Accurate phylogenetic inference with a  
524 symmetry-preserving neural network model. *arXiv*. [https://doi.org/10.48550/arXiv.2201.](https://doi.org/10.48550/arXiv.2201.04663)  
525 [04663](https://doi.org/10.48550/arXiv.2201.04663)
- 526 Suvorov, A., Hochuli, J., & Schrider, D. R. (2020). Accurate inference of tree topologies from  
527 multiple sequence alignments using deep learning. *Systematic Biology*, *69*(2), 221–233.
- 528 Suvorov, A., & Schrider, D. R. (2022). Reliable estimation of tree branch lengths using deep neural  
529 networks. *bioRxiv*. <https://doi.org/10.1101/2022.11.07.515518>
- 530 Voznica, J., Zhukova, A., Boskova, V., Saulnier, E., Lemoine, F., Moslonka-Lefebvre, M., & Gascuel,  
531 O. (2022). Deep learning from phylogenies to uncover the epidemiological dynamics of  
532 outbreaks. *Nature Communications*, *13*(1), 3896.
- 533 Zaharias, P., Grosshauser, M., & Warnow, T. (2022). Re-evaluating deep neural networks for  
534 phylogeny estimation: The issue of taxon sampling. *Journal of Computational Biology*,  
535 *29*(1), 74–89.
- 536 Zhang, Y., Zhu, Q., Shao, Y., Jiang, Y., Ouyang, Y., Zhang, L., & Zhang, W. (2023). Inferring  
537 historical introgression with deep learning. *Systematic Biology*, syad033.
- 538 Zou, Z., Zhang, H., Guan, Y., & Zhang, J. (2020). Deep residual neural networks resolve quartet  
539 molecular phylogenies. *Molecular Biology and Evolution*, *37*(5), 1495–1507.