

1  
2  
3  
4  
5  
6  
7  
8  
9  
10

# Accurate Detection of Incomplete Lineage Sorting via Supervised Machine Learning

Benjamin Rosenzweig<sup>1,\*</sup>, Andrew Kern<sup>2</sup>, and Matthew Hahn<sup>1,3</sup>

<sup>1</sup>Department of Computer Science, Indiana University, Bloomington, IN, 47408, U.S.A.

<sup>2</sup>Department of Biology, University of Oregon

<sup>3</sup>Department of Biology, Indiana University, Bloomington, IN, 47408, U.S.A.

\* [bkrosenz@iu.edu](mailto:bkrosenz@iu.edu)

## 11 Abstract

12 Gene tree discordance due to incomplete lineage sorting or introgression has been described in  
13 numerous genomic datasets. Among distantly related taxa, however, it is difficult to differentiate  
14 these biological sources of discordance from discordance due to errors in gene tree  
15 reconstruction, even when supervised machine learning techniques are used to infer individual  
16 gene trees. Here, rather than applying machine learning to the problem of inferring single tree  
17 topologies, we develop a model to infer important properties of a particular internal branch of the  
18 species tree via genome-scale summary statistics extracted from individual alignments and  
19 inferred gene trees. We show that our model can effectively predict the presence/absence of  
20 discordance, estimate the probability of discordance, and infer the correct species tree topology  
21 in the presence of multiple, common sources of error. While gene tree topology counts are the  
22 most salient predictors of discordance at short time scales, other genomic features become  
23 relevant for distantly related species. We validate our approach through simulation, and apply it  
24 to data from the deepest splits among metazoans. Our results suggest that the base of Metazoa  
25 experienced significant gene tree discordance, implying that discordant traits among current taxa  
26 can be explained without invoking homoplasy. In addition, we find support for Porifera as the  
27 sister clade to the rest of Metazoa. Overall, these results demonstrate how machine learning can  
28 be used to answer important phylogenetic questions, while marginalizing over individual gene  
29 tree—and even species tree—topologies.

30

31 keywords: *homoplasy, hemiplasy, gene tree discordance, deep learning, metazoa*

32

33

## 34 Introduction

35 Topological discordance among gene trees is found in nearly all phylogenomic datasets.  
36 Biological causes of discordance include both incomplete lineage sorting (ILS; Figure 1) and  
37 introgression (Maddison 1997), while technical sources of discordance include substitution rate  
38 heterogeneity and model misspecification. Rampant discordance has been described in recent  
39 radiations, such as cichlids (Ronco et al., 2021), horses (Jónsson et al. 2014), *Drosophila*  
40 (Pollard et al. 2006), and tomatoes (Pease et al. 2016), as well as in ancient radiations such as  
41 birds (Jarvis et al. 2014) and land plants (Wickett et al. 2014). When discordance is due to  
42 biological causes, phylogenetics provides powerful tools for understanding genotype-phenotype  
43 associations (Pease et al. 2016; Smith et al. 2020), as well as for distinguishing traits that have  
44 evolved more than once (homoplasy) from traits that have evolved once on discordant gene trees  
45 (hemiplasy; (Avice and Robinson 2008)). Therefore, determining whether gene tree discordance  
46 is due to biological factors is a key step in understanding trait evolution.

47  
48 For closely related species, individual gene trees can be recovered with high accuracy by  
49 maximum likelihood, Bayesian inference, or parsimony methods. For more distantly related taxa,  
50 gene tree reconstruction increasingly falls prey to errors due to homoplasy at individual sites (i.e.  
51 “multiple hits”). For truly ancient divergences, discordant topologies will be inferred at the same  
52 frequency as the species tree topology, regardless of the amount of discordance present due to  
53 biological factors (Figure 2). Evolutionary model misspecification, whether due to an incorrect  
54 substitution matrix or failure to account for site- and lineage-specific rate heterogeneities,  
55 compounds these errors. No method currently exists that can distinguish between true  
56 discordance—that due to biological processes—and the noise introduced by gene tree  
57 reconstruction error. Even when the species tree topology is known with high confidence, any  
58 assertion about the history of specific biological traits still requires accurately estimating the  
59 probability of hemiplasy, a procedure that requires knowledge of the level of true biological  
60 discordance among gene trees.

61  
62 A particularly contentious problem in phylogenetics involves relationships at the base of the  
63 Metazoa. Over the last decade, genomic data has led systematists to reconsider the relationship

64 between the ParaHoxozoa clade (comprising the phyla Bilateria, Cnidaria, and Placozoa) and the  
65 two other major clades of animals: the phyla Ctenophora and Porifera (Laumer et al. 2019;  
66 Simion et al. 2017; Whelan et al. 2015)). In the absence of biological discordance, the  
67 Ctenophora-sister hypothesis implies that major features such as nervous systems (Moroz et al.  
68 2014, 201), basement membranes, and the through-gut either evolved multiple times or arose  
69 once in the ancestor of all three groups and then were lost in the Porifera. While much effort has  
70 been devoted to investigating the impact of fast-evolving lineages and the choice of substitution  
71 model on species tree reconstruction at the base of Metazoa (Li et al. 2021), scant attention has  
72 been paid to the potential role of gene tree discordance in reconciling the evolution of complex  
73 traits during this period (King and Rokas 2017). In the presence of biological discordance, even  
74 complex traits can arise once and appear in multiple lineages without needing to invoke  
75 homoplasy (Hahn and Nakhleh 2016; Guerrero and Hahn 2018; Hibbins, Gibson, and Hahn  
76 2020).

77  
78 Supervised machine learning (SML) comprises a variety of parameter-rich regression algorithms  
79 that excel at learning nonlinear mappings from noisy, feature-rich data. SML methods have been  
80 successfully employed in a variety of contexts in phylogenetics, from inferring quartet topologies  
81 for single loci (Suvorov, Hochuli, and Schrider 2020; Zou et al. 2019), to enhancing nucleotide  
82 substitution model selection (Abadi et al. 2020), to tree-search proposal distributions (Azouri et  
83 al. 2021). Rather than attempting to use SML to improve the accuracy of individual gene trees,  
84 our goal here is to predict the fraction of gene trees in a dataset that are biologically discordant.  
85 Our model infers properties of an internal branch of the species tree given a collection of  
86 summary statistics from a set of gene trees. Given the noisy nature of gene trees inferred from  
87 deep divergences (Figure 2), SML offers a potentially powerful method for overcoming  
88 inference problems in such datasets.

89  
90 In this paper, we show that a variety of SML algorithms can effectively distinguish biological  
91 discordance from gene tree inference error across a wide range of parameter space. We  
92 demonstrate that simple feed-forward artificial neural network architectures are successful at (1)  
93 predicting the species tree topology, (2) predicting the fraction,  $p$ , of biological discordance in a  
94 set of gene trees, and (3) detecting the presence or absence of biological discordance in a given

95 dataset. We show that biological discordance can be identified under a wide range of biologically  
96 relevant scenarios, even when the assumptions of the training regime are violated.

97

## 98 Methods

### 99 Supervised Machine Learning Models

100 The regression and classification models were implemented in scikit-learn (Buitinck et al. 2013),  
101 with the exception of the deep neural network (DNN) models, which were implemented in  
102 pytorch (Mazza and Pagani 2021). Hyperparameters for the DNN architectures and optimization  
103 algorithm were selected via a Bayesian optimization search over 2000 candidate configurations.  
104 Three major categories of SML model were employed: Linear (regularized linear and logistic  
105 regression), Ensemble (random forest, adaptive boosting, and gradient boosting with decision  
106 trees), and a DNN with rectified linear activations. The DNN estimator of  $p$  (the expected  
107 frequency of concordant trees), DNN-Prob, and the species tree topology predictor, DNN-Top,  
108 were trained with mean squared logarithmic error and cross-entropy losses, respectively. The  
109 DNN classifier, DNN-Class, was obtained by thresholding the DNN-Prob regressor.

110

111

### 112 Simulation Conditions

113 All models were trained on the same training dataset, comprising 1.5 million synthetic sets of  
114 alignments, each of which included variable sequence lengths, site-specific evolutionary rates,  
115 and substitution models. Parameters for simulated training data are described in Table 1 and the  
116 overall process is shown in Figure 3.

117

118 Four-taxon gene trees were simulated with ms (Hudson 2002) and alignments with seq-gen  
119 (Rambaut and Grassly 1997). For the training set, all simulated trees initially had the same rate  
120 of sequence evolution and all loci were simulated without recombination; below we describe  
121 how heterotachy and recombination were added in the test sets. Training sequences were  
122 simulated using both LG and WAG substitution models with a discrete 5-category gamma model  
123 of site-specific rate heterogeneity. Individual gene trees were inferred with RaXML-NG  
124 (Kozlov et al. 2019), using a discrete gamma model of site-specific rate heterogeneity with 4

125 categories and empirical amino acid frequencies. Site concordance factors were computed with  
126 IQ-TREE (Minh, Hahn, and Lanfear 2020).

127  
128 For a given set of parameters, alignments were generated for a large number of loci (a minimum  
129 of 1,000 for each parameter set), and gene trees were inferred for each alignment. Raw features  
130 such as alignment lengths, site concordance factors, and patristic distances were computed from  
131 both the inferred trees and raw alignments. These feature values were then binned according to  
132 the rooted topology of the inferred tree. For every bin, summary statistics (central moments and  
133 order statistics) were computed across all values of each feature. This process resulted in 168  
134 total features per simulated dataset (Figure 3). Features such as concordance factors, patristic  
135 distances, and topology counts are not invariant to relabeling of taxa; training data was therefore  
136 augmented with label permutations corresponding to the three rooted species trees. At test time,  
137 predictions ( $\hat{p}$  for the regression algorithms, softmax outputs for DNN-Top) were generated  
138 separately for each of the three label permutations of the feature vector, and the resulting  
139 predicted values were averaged across all three permutations (the input data were also permuted  
140 in a similar manner). The average of the ensemble was considered the point estimate of our  
141 model.

142  
143 Additional conditions used for test datasets are described in Supplementary Table 1. Lineage-  
144 specific heterotachy was simulated by multiplying all ingroup external branches by independent  
145 Gamma(4, .25) multipliers. Recombination was simulated by randomly permuting subsequences  
146 within a synthetic dataset; this was done because simulations using the coalescent with  
147 recombination would have resulted in a training dataset that was too sparse. Both the number of  
148 recombination blocks per sequence and the fraction of loci experiencing recombination varied by  
149 dataset.

150

## 151 Metazoa Datasets

152 We used 13 previously published datasets containing 315 taxa from three Metazoan clades  
153 (Supplementary Table 2). Most of these studies utilized concatenated alignments with either the  
154 general time reversible (GTR) substitution model across all sites or a form of data partitioning  
155 (Whelan and Halanych 2017). Our analysis requires individual gene trees, but individual genes

156 do not provide enough data to reliably estimate GTR model parameters. We therefore used the  
157 LG substitution model (Le and Gascuel 2008) with Gamma-distributed site heterogeneity to infer  
158 gene trees for all data matrices. Our training set includes a mixture of substitution models (LG  
159 and WAG), as well as site heterogeneity. Although the training set did not include misspecified  
160 substitution models, the DNN methods are highly robust to substitution matrix misspecification  
161 (see next section). This justifies the use of a single substitution model across all gene trees for  
162 this analysis. Augmenting the training set with gene trees inferred from misspecified models  
163 could further improve predictive accuracy.

164

165 For each data matrix, we extracted alignments from all 4-taxon quartets comprising one member  
166 each from Ctenophora (Cte), Porifera (Por), ParaHoxozoa (PaH), and an outgroup from  
167 Choanoflagellata/Fungi/Ichthyospora/Filasterea (Out). For each gene, the corresponding subtree  
168 was rooted with the outgroup taxon. Quartets of taxa having fewer than 25 genes were excluded  
169 from the analysis. The distribution of inferred gene tree topologies from all quartets is displayed  
170 in Table 2.

171

## 172 Results

### 173 Predicting the probability of concordance

174 **Predicting  $p$ .** In our first set of experiments, we trained a variety of SML methods to predict  $p$ ,  
175 the probability that a 4-taxon gene tree matches the topology of its species tree. The value of  $p$   
176 ranges from  $1/3$  (all topologies are equiprobable) to 1 (all gene trees are concordant) and is a  
177 function of internal branch lengths (IBL; Hudson 1983). The quantity  $1-p$  therefore represents  
178 the true probability of discordance in a dataset.

179

180 Figure 4 shows the error,  $\hat{p} - p$ , of all six SML models across a range of external and internal  
181 branch lengths (EBL and IBL). SML models differ in their error rates (Kruskal-Wallis test,  $H =$   
182  $2747.88$ ,  $p < 10^{-300}$ ); one-sided Wilcoxon rank-sum tests were conducted between each pair  
183 of algorithms to explore differences in performance. Gradient Boosting (GradBoost) and the  
184 deep neural network (DNN-Pred) performed best, followed by Random Forest (RF), then  
185 Adaptive Boosting (AdaBoost), then the constant baseline predictor (Mean; equal to 0.999995,

186 the mean  $p$  across the training set), and finally the regularized linear regression (ElasticNet) (all  
187  $P < 10^{-80}$ , Holms-Bonferroni corrected  $\alpha = 3.34 \times 10^{-4}$ ). Although the overall error rate of  
188 DNN-Pred was not significantly better than GradBoost, its performance was more consistent: the  
189 50% interquartile ranges of error for all algorithms were: 0.4043 (ElasticNet), 0.399 (RF), 0.4040  
190 (AdaBoost), 0.0097 (GradBoost), 0.0006 (DNN-Pred).

191

192 In the remainder of the paper, we therefore focus on results from the deep neural networks.

193 Figure 5 shows the performance of DNN-Pred on simulated test data as a function of EBL and  
194 IBL.  $\hat{p}$  is most accurate when (a) both IBL and EBL are short or (b) IBL is long. For long EBL,  
195  $\hat{p}$  tends to be an overestimate of  $p$  for short IBL (where expected amino acid sequence  
196 divergence  $< 0.01$ ) and an underestimate for intermediate IBL (between .01 and 0.05). As  
197 expected, performance is best when IBL is extremely long ( $> 0.07$ , corresponding to  $p >$   
198  $0.99999$ ). In this region gene tree inference error is negligible even for long EBL.

199 Underestimating  $p$  for IBL in the 0.01 to 0.05 range is a feature shared with the gradient boosting  
200 algorithm, the only other method with reasonable performance across the regions of parameter  
201 space we explored.

202

## 203 Model Misspecification

204 **Recombination.** Long protein-coding genes have almost certainly experienced recombination,  
205 and therefore combine multiple topologies (Lanier and Knowles 2012). This increases the risk  
206 of hemiplasy (Mendes, Livera, and Hahn 2019), leading to inaccurate estimates of both  
207 divergence times and tree topologies. To investigate the impact of recombination on prediction  
208 accuracy we performed additional manipulations. To mimic the presence of intra-locus  
209 recombination, we varied the fraction of recombinant loci present in each dataset (Figure 6),  
210 drawing each recombinant locus uniformly from 2, 3, and 4 independently sampled blocks of  
211 sequence.

212

213 Although increased intragenic recombination has an effect on DNN-Pred's ability to predict  $p$ ,  
214 dataset size can effectively overcome this effect. For any single dataset size, we do see  
215 increasing error with an increasing fraction of recombinant sequences (Figure 6), but also  
216 increased error for smaller datasets even in the absence of recombination. To understand why

217 this is the case, note that DNN-Pred relies heavily on site concordance factors and gene tree  
218 topology counts for much of parameter space (see Section “Feature Importance”). These  
219 features are minimally impacted by intragenic recombination: site concordance factor-derived  
220 statistics are pooled across genes and, while recombination may cause branch lengths of  
221 individual gene trees to be overestimated (Mendes & Hahn, 2016), the topology is still more  
222 likely to match the “majority tree” of the constituent sequences, than either minor tree.

223  
224 **Evolutionary model misspecification.** We investigated the impact of several forms of  
225 evolutionary model misspecification: differing substitution matrices in training and test data, and  
226 hidden site- and lineage-specific rate heterogeneity. Each of these cases can be thought of an  
227 instance of what is called “data drift” in the machine learning world, where training data does not  
228 match empirical data which we wish to perform inference on. Misspecification of the  
229 substitution matrix has a negligible effect on predictive accuracy: prediction error on datasets  
230 simulated with the LG substitution matrix and inferred with WAG, or vice versa, do not differ  
231 significantly compared to datasets where trees were inferred with the correct substitution matrix  
232 (Supplementary Figure 1a). This suggests that sequence-based features (i.e. site concordance  
233 factors) provide enough information to estimate  $p$  despite errors in gene tree inference, or that  
234 substitution model misspecification does not systematically bias gene tree inference error in a  
235 way that obscures information regarding  $p$ . Further, to deal with rate heterogeneity we trained  
236 using simulations that included variation in rates among sites. This allowed for successful  
237 inference in the face of both site-specific heterogeneity and lineage-specific heterogeneity  
238 (heterotachy) in test data and did not significantly increase the error in estimating  $p$  (Kruskal-  
239 Wallis test;  $H = 0.686, P = 0.710$ ; Supplementary Figure 1b).

240

## 241 Binary Classification and Species Tree Prediction

242 In regions of parameter space where  $p$  cannot be inferred precisely, it is still useful to know  
243 whether *any* ILS exists in a particular dataset. Statistical power can be increased by introducing  
244 a binary classifier (DNN-Class) which predicts whether  $p$  is above or below a particular value.  
245 Note that distinguishing between the events  $p < 1$  (“potential ILS”) and  $p = 1$  (“no ILS  
246 whatsoever”) is not informative, as  $p = 1$  corresponds to an infinitely long internal branch. We  
247 instead chose a target value of 0.9: this is interpreted to mean that datasets with  $p < 0.9$  exhibit

248 biological discordance and those with  $p > 0.9$  do not. A more (or less) conservative target value  
249 than 0.9 can be chosen depending on the dataset. (For a dataset of  $N$  genes, for example,  
250 choosing a target value of  $1 - \frac{1}{N}$  would ensure that the expected number of discordant trees under  
251 the “no ILS” hypothesis will be less than 1). DNN-Class was constructed by applying a separate  
252 decision threshold to the output of DNN-Pred. The decision threshold of 0.9511 was chosen to  
253 optimize the geometric mean of the classifier’s precision and recall on the training set (Figure  
254 7a); a different target value will result in a different decision threshold.

255  
256 Although a task-specific DNN classifier can be trained independently (e.g. with a binary cross-  
257 entropy loss), we found that this did not improve over the procedure described above. At this  
258 threshold, the probability of incorrectly predicting no ILS (“false positives”) is negligible for all  
259 parameters tested. The probability of incorrectly predicting ILS (“false negatives”) is below 0.3  
260 in most regions of parameter space, with an elevated false negative rate for extreme values of  
261 EBL (Figure 7b). While DNN-Pred tends to overestimate  $p$  for  $IBL < .01$ , DNN-Class achieves  
262 zero error in this region. These results support the use of DNN-Class as a conservative estimator  
263 of concordance (i.e. “no ILS”).

264  
265 Similarly, we can construct a ternary classifier (DNN-Top) to predict which of the three possible  
266 topologies is most likely to be the underlying species tree. The accuracy of DNN-Top varies with  
267 EBL and IBL in a pattern similar to DNN-Pred and DNN-Class (Figure 8). When IBL is greater  
268 than 0.0015 (corresponding to  $p \approx 0.506$ ), DNN-Top can correctly infer the species tree  
269 topology with greater than chance accuracy for all values of EBL tested. For IBL greater than  
270 0.025 ( $p \approx 0.996$ ), accuracy is greater than 90%, even when a large fraction of gene trees are  
271 incorrectly inferred due to inference error.

272

## 273 Feature Importance

274 Tree-based classification and regression algorithms enable an intuitive notion of feature  
275 importance: the extent to which each covariate contributes to a prediction can be measured by  
276 the *Gini importance*: the mean decrease in impurity (for classification trees) or variance (for  
277 regression trees) at each node in which that covariate appears (Louppe et al. 2013).

278 Supplementary Figure 2a shows the features with the highest Gini importance for the gradient  
279 boosting regressor. Topology counts and sCF-derived features are the most informative features.  
280 This is expected, as topology counts alone are sufficient to recover  $p$  in the short external branch  
281 regime.

282  
283 Untangling the importance of features in a DNN is considerably more challenging, and this  
284 opacity has hampered the entry of DNN's into many fields of research. One method that has  
285 recently gained in popularity is the Shapley value: the average marginal contribution of each  
286 feature across all subsets of features. Exact computation of Shapley values is linear in the  
287 number of instances and exponential in the number of features, making it impractical for large,  
288 high-dimensional datasets. Numerous approximate methods have therefore been developed. We  
289 utilized Deep SHAP (Lundberg 2017), an approximate Shapley algorithm that leverages the  
290 compositional nature of neural networks. Like other measures of feature importance, the Shapley  
291 method assumes features are uncorrelated. This is not the case in our model, both because  
292 certain summary statistics for each feature are highly correlated (all measures of dispersion or of  
293 central tendency), and because the symmetry of the data makes certain combinations of features  
294 equivalent (permutations of taxa labels should not affect the output of DNN-Pred or DNN-  
295 Class). To overcome this, we computed Shapley values for the following groups of features:  
296 statistics derived from pairwise patristic distances  $d(i,j)$  (15 features each), statistics derived from  
297 counts of informative sites (15 features each), statistics derived from alignment lengths (15  
298 features each), statistics derived from site concordance factors (45 features) and counts of gene  
299 tree topologies (3 features). As with tree-based regressors, topology counts and site concordance  
300 factors are the most informative features at short distances (Supplementary Figure 2b). For  
301 deeper divergences, features derived from gene tree branch lengths become more important  
302 (Supplementary Figure 2c). Unlike the gradient boosting regressor, DNN-Pred relies heavily on  
303 sequence length and “number of informative sites” statistics. As these variables determine the  
304 quality of other features (such as the inferred gene tree topologies and patristic distances), they  
305 may be important in determining how the network weights such features for a given instance.

306

307 **Biological Discordance at the Base of the Metazoa**

308 In order to highlight the power of the SML approach described here, we applied our DNN  
309 models to a dataset in which the presence of discordance would have important biological  
310 implications. We used 13 datasets containing 315 taxa from three Metazoan clades to examine  
311 discordance at the base of animals. Because these datasets contain many more than four taxa, for  
312 each dataset we analyzed many possible sampled quartets: one species from each of Ctenophore,  
313 Porifera, ParaHoxozoa, and an Outgroup, using all genes shared by these four sampled species.  
314 Each quartet was analyzed using DNN-Pred (“how much discordance is there?”), DNN-Class  
315 (“is there *any* discordance?”), and DNN-Top (“what is the species tree?”).

316  
317 **Level of discordance.** Using DNN-Pred, estimates of  $p$  were highly variable across quartets both  
318 between datasets and within each dataset (Supplementary Figure 6). The distribution of  $\hat{p}$  values  
319 was bimodal, with peaks at 0.43 and 0.88); this variability may be due to underlying variability  
320 in the rate of evolution of individual genes or taxa sampled for each quartet. Taking 0.6 as the  
321 putative metazoan EBL (one half the average patristic distance between the two sister taxa across  
322 all studies), the prediction error of DNN-Pred on simulated data in this area of parameter space  
323 ranges from -0.022 to 0.17 as a function of the true (unknown) IBL (Figure 5); this corresponds  
324 to a relative error of -4% to +50%, at the base of metazoa. Given this level of prediction error  
325 and putative values of  $\hat{p}$ , we cannot be certain using DNN-Pred whether there is truly biological  
326 discordance in this dataset. In trees with more recent common ancestors,  $\hat{p}$  will be much more  
327 accurate (Figure 5).

328  
329 **Presence of discordance.** In contrast, at this same level of divergence, the maximum false  
330 negative rate and false positive rate of DNN-Class are 0.20 and 0.00, respectively, across all IBL.  
331 This gives us confidence that the presence or absence of ILS can still be detected even if the  
332 expected frequency of concordant gene trees cannot be precisely determined. In all but one  
333 dataset, DNN-Class infers  $p < 0.9$  in the overwhelmingly majority of quartets, providing strong  
334 evidence for ILS at the base of Metazoa (Figure 10a). The results from DNN-Class suggest that  
335 significant amounts of discordance are due to ILS, though we cannot say precisely how much.

336  
337 Looking within datasets, quartets in which DNN-Class predicted low  $p$  (high probability of ILS)  
338 have significantly more informative sites (two-sided  $t$ -tests,  $P < 10^{-8}$ ) and have higher site

339 concordance factors (sCFs) for alternate topologies (two-sided  $t$ -tests,  $P < 10^{-4}$ ), assuming that  
340 the species tree is Porifera-sister (see next section). Sequences with more informative sites carry  
341 more phylogenetic signal; feature values and predictions for these quartets should therefore be  
342 considered more reliable. The only dataset with a significant number of quartets supporting  $p >$   
343 0.9 is the Nosenko et al. (2013) ribosomal dataset. With the exception of some sCF statistics, all  
344 feature values for this dataset differed significantly from the other datasets examined (Mann-  
345 Whitney U tests with Holm-Sidak correction, all  $P < 0.001$ ). The features which varied most  
346 from other datasets were those derived from sequence length and pairwise distances involving  
347 the outgroup; the Nosenko et al. dataset has much shorter sequences than the other datasets  
348 examined, as well as longer outgroup branch length statistics (with sequence lengths at least 1.5  
349 standard deviations below and branch lengths at east 1.5 standard deviations above the global  
350 mean). This suggests that stochastic error (due to short sequences) and systematic error (due to  
351 long branch attraction) may be especially high in this dataset, making these predictions less  
352 reliable.

353  
354 **Species-tree prediction.** Across all datasets, there was consistent support for the Porifera-sister  
355 topology using predictions from DNN-Top (Figure 10b; Table 2). The output of DNN-Top is in  
356 the form of posterior probabilities for each topology, and therefore model confidence can be  
357 assessed by comparing such values. Site concordance factors are strongly associated with model  
358 confidence, measured as the absolute value of the log-likelihood ratio of the Porifera-sister and  
359 Ctenophora-sister hypotheses: quartets for which DNN-Top strongly favors Porifera-sister have  
360 high sCF values for the Porifera-sister topology; likewise quartets for which DNN-Top favors  
361 Ctenophora-sister have high sCF values for Ctenophore sister (Supplementary Figure 8).

362

## 363 Discussion

364 Topological discordance is a ubiquitous feature of phylogenomic datasets. When gene trees are  
365 accurately inferred, this discordance can be attributed to biological factors such as ILS or  
366 introgression. In this study we find that SML algorithms can provide useful information  
367 regarding ILS even in the presence of gene tree reconstruction error. Our estimators perform well  
368 even in worst-case scenarios, where test data contain multiple sources of noise not present in the

369 training data. Improvements in specific cases can be obtained by retraining models with  
370 simulated data that approximates the suspected sources of noise in a target dataset, for instance in  
371 the manner that we trained under site-specific rate heterogeneity.

372  
373 Previous applications of SML in phylogenetics have reported results from a single algorithm or  
374 neural network architecture, presumably the result of a laborious trial-and-error process that  
375 often goes unreported. Here we take a systematic approach to evaluating a range of SML  
376 algorithms, tuning their hyperparameters, and characterizing their performance and limitations.  
377 Given the well-justified opposition to adopting “black box” inference models in systematics (and  
378 biology more generally), this analysis provides a useful framework for future SML studies to  
379 emulate and improve upon. Recent studies have used convolutional neural networks to infer 4-  
380 taxa topologies for individual gene trees from multi-sequence alignments (Suvorov, Hochuli, and  
381 Schridder 2020; Zou et al. 2019; Solis-Lemus, Yang, and Zepeda-Nunez 2022). Although these  
382 approaches compare favorably with conventional approaches to topology inference, all methods  
383 fail on deep divergences, presumably due to the same recurrent mutation phenomena that  
384 confound likelihood, parsimony, and distance-based methods (Molloy and Warnow 2018).  
385 Furthermore, extending such methods to phylogenies with more than four taxa requires some  
386 form of tree reconciliation (Strimmer and von Haeseler 1996; Snir and Rao 2012), which has  
387 been shown to be less accurate than standard Neighbor-Joining or maximum likelihood inference  
388 (Zaharias, Grosshauser, and Warnow 2022). Here, we take a different approach: we develop  
389 SML methods to infer important properties of a particular internal branch in the species tree via  
390 genome-scale summary statistics extracted from individual alignments and inferred gene trees.  
391 This approach allows us to utilize state-of-the-art maximum likelihood methods to infer gene  
392 trees from large multisequence alignments, which can then be subsampled to obtain quartets that  
393 contain the branch of interest. This method provides enough signal to resolve a problematic  
394 internal branch and to accurately classify a dataset as containing biological discordance, even  
395 when inferred gene trees contain large amounts of noise.

396  
397 The approach taken here also demonstrates the advantage that SML can provide over traditional  
398 methods in answering highly tailored questions in phylogenetics. To this end, we emphasize that  
399  $p$  is a property of the dataset as a whole—our methods do not predict whether an individual gene

400 tree is concordant or discordant. Future research could likewise focus on training SML models to  
401 predict specific quantities of interest from raw data, rather than separately optimizing each step  
402 of an analysis pipeline or trying to solve every tree inference problem with a single model. SML  
403 with deep learning has been successfully applied in other phylogenomic studies to distinguish  
404 between reticulated and bifurcating phylogenies (Burbrink and Gehara 2018), between different  
405 hybridization scenarios (Blischak, Barker, and Gutenkunst 2021), and to identify selective  
406 sweeps in population genomic datasets (Isildak, Stella, and Fumagalli 2021; Kern and Schrider  
407 2018). A major disadvantage of all such approaches (including our own) is the large amount of  
408 training data required. When this data is produced via coalescent simulations as in the present  
409 study, this translates into months of compute time (~65 CPU-hours to generate 1000 gene trees,  
410 simulate their alignments, and infer maximum likelihood trees on a 2.5 GHz processor).  
411 Furthermore, our results show that a great deal of simulation and training resources are expended  
412 in exploring insoluble regions of parameter space: those areas where even the richest models  
413 with the largest training samples can do no better than a random classifier. Future SML  
414 applications might consider utilizing some form of active learning (Settles 2012) to simulate and  
415 train jointly, focusing effort on those regions where learning is tractable and ignoring regions  
416 where SML does not appear to provide any advantage.

417  
418 Summary methods for species tree inference such as ASTRAL (Zhang et al. 2018), MP-EST  
419 (Liu, Yu, and Edwards 2010), and ASTRID (Vachaspati and Warnow 2015) are statistically  
420 consistent when given accurately inferred input gene trees, even in the presence of ILS. These  
421 methods fail, however, when phylogenetic signal is obscured by high levels of gene tree  
422 reconstruction error (Molloy and Warnow 2018). The fact that DNN-Top and DNN-Class  
423 perform well across a similar range of parameter space suggests that that deep neural networks  
424 may be capable of inferring species trees along with their internal branch lengths in a manner  
425 similar to summary methods, and may even be able to improve upon them by being more robust  
426 to gene tree error. Future work could focus on modeling this gene tree error before it is input into  
427 summary methods.

428  
429 Applying our method to relationships at the base of the Metazoa, we find strong support for both  
430 the Porifera-sister hypothesis and for the presence of appreciable amounts of incomplete lineage

431 sorting. Branch-specific heterotachy has been cited as a key factor in the uncertainty regarding  
432 the base of Metazoa (Pisani et al. 2015; Kapli, Yang, and Telford 2020). In particular, the long  
433 branch leading to Ctenophora has prompted many researchers to explore amino acid recoding  
434 strategies to uncover hidden phylogenetic signal (Redmond and McLysaght 2021). While there  
435 are several important caveats to the inferences presented here, we have shown that our methods  
436 (i.e. DNN-Top) are robust to heterotachy, as well as a number of additional model violations. We  
437 therefore conclude that there is good reason to believe that we are accurately classifying the  
438 species topology.

439  
440 Even if a species tree topology is known with high confidence, any assertion about the history of  
441 specific biological traits still requires accurately estimating the probability of hemiplasy – a  
442 procedure that requires knowledge of the level of true biological discordance among gene trees.  
443 Much of the argument around the “true” set of relationships at the base of the Metazoa has  
444 presented alternative species topologies as uniquely tied to alternative trait histories. Despite our  
445 apparent inability to accurately estimate the exact level of biological discordance this deep in  
446 time (i.e. DNN-Pred), even the consideration of hemiplasy means that the species topology can  
447 be separated from trait histories. In this sense, the results from DNN-Class should play an  
448 important role in understanding the complex pattern of events that occurred during the early  
449 stages of animal evolution. Importantly, our results concluding that there was biological  
450 discordance during this period do not rely on any particular species tree topology, nor do they  
451 necessarily imply that discordance is due to ILS. Although it would be rash to make definitive  
452 conclusions at this time, the asymmetry in frequency of the two minor gene tree topologies  
453 (Table 2) and predicted species tree topologies (Figure 10b) could be due to introgression among  
454 early animals (cf. Huson and Bryant 2006). Given the ubiquity of introgression across the tree of  
455 life, there is no reason to think it was not also occurring early in animal evolution, possibly  
456 affecting key biological innovations.

457

458

459

## 460 5 Acknowledgements

461 This research was supported by a Department of Defense SMART fellowship (BKR) and  
462 National Science Foundation grant DEB-1936187 (MWH). ADK was supported under NIH  
463 award 1R01HG010774.

464

## 465 6 Supplementary Material

466 Code used in this paper is available at <https://github.com/bkrosenz/ml4ils>.

467

## 468 7 Bibliography

- 469 Abadi, Shiran, Oren Avram, Saharon Rosset, Tal Pupko, and Itay Mayrose. 2020. “Modelteller:  
470 Model Selection for Optimal Phylogenetic Reconstruction Using Machine Learning.”  
471 *Molecular Biology and Evolution* 37 (11): 3338–52.  
472 <https://doi.org/10.1093/molbev/msaa154>.
- 473 Avise, John C., and Terence J. Robinson. 2008. “Hemiplasy: A New Term in the Lexicon of  
474 Phylogenetics.” *Systematic Biology* 57 (3): 503–7.  
475 <https://doi.org/10.1080/10635150802164587>.
- 476 Azouri, Dana, Shiran Abadi, Yishay Mansour, Itay Mayrose, and Tal Pupko. 2021. “Harnessing  
477 Machine Learning to Guide Phylogenetic-Tree Search Algorithms.” *Nature*  
478 *Communications* 12 (1): 1–9. <https://doi.org/10.1038/s41467-021-22073-8>.
- 479 Blischak, Paul D, Michael S Barker, and Ryan N Gutenkunst. 2021. “Chromosome-Scale  
480 Inference of Hybrid Speciation and Admixture with Convolutional Neural Networks.” In  
481 *Molecular Ecology Resources*, 21:2676–88. <https://doi.org/10.1111/1755-0998.13355>.
- 482 Buitinck, Lars, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier  
483 Grisel, Vlad Niculae, et al. 2013. “API Design for Machine Learning Software:  
484 Experiences from the Scikit-Learn Project,” September. <http://arxiv.org/abs/1309.0238>.
- 485 Burbink, Frank T, and Marcelo Gehara. 2018. “The Biogeography of Deep Time Phylogenetic  
486 Reticulation.” *Systematic Biology* 67 (5): 743–55. <https://doi.org/10.1093/sysbio/syy019>.
- 487 Guerrero, Rafael F., and Matthew W. Hahn. 2018. “Quantifying the Risk of Hemiplasy in  
488 Phylogenetic Inference.” *BioRxiv* 115 (50): 391391. <https://doi.org/10.1101/391391>.
- 489 Hahn, Matthew W., and Luay Nakhleh. 2016. “Irrational Exuberance for Resolved Species  
490 Trees.” *Evolution* 70 (1): 7–17. <https://doi.org/10.1111/evo.12832>.
- 491 Hibbins, Mark S, Matthew J.S. Gibson, and Matthew W Hahn. 2020. “Determining the  
492 Probability of Hemiplasy in the Presence of Incomplete Lineage Sorting and  
493 Introgression.” *ELife* 9: 1–33. <https://doi.org/10.7554/ELIFE.63753>.
- 494 Hudson, Richard R. 2002. “Generating Samples under a Wright-Fisher Neutral Model of Genetic  
495 Variation.” *Bioinformatics* 18 (2): 337–38.  
496 <https://doi.org/10.1093/bioinformatics/18.2.337>.
- 497 Huson, Daniel H., and David Bryant. 2006. “Application of Phylogenetic Networks in  
498 Evolutionary Studies.” *Molecular Biology and Evolution* 23 (2): 254–67.  
499 <https://doi.org/10.1093/molbev/msj030>.

- 500 Isildak, Ulas, Alessandro Stella, and Matteo Fumagalli. 2021. “Distinguishing between Recent  
501 Balancing Selection and Incomplete Sweep Using Deep Neural Networks.” *Molecular*  
502 *Ecology Resources* 21 (8): 2706–18. <https://doi.org/10.1111/1755-0998.13379>.
- 503 Jarvis, Erich D., S. Mirarab, Andre J. Aberer, B. Li, P. Houde, Cai Li, S. Y. W. Ho, et al. 2014.  
504 *Whole-Genome Analyses Resolve Early Branches in the Tree of Life of Modern Birds*  
505 *(Supplement)*. Vol. 346. <http://www.sciencemag.org/cgi/doi/10.1126/science.1251385>.
- 506 Jónsson, Hákon, Mikkel Schubert, Andaine Seguin-Orlando, Aurélien Ginolhac, Lillian  
507 Petersen, Matteo Fumagalli, Anders Albrechtsen, et al. 2014. “Speciation with Gene  
508 Flow in Equids despite Extensive Chromosomal Plasticity.” *Proceedings of the National*  
509 *Academy of Sciences of the United States of America* 111 (52): 18655–60.  
510 <https://doi.org/10.1073/pnas.1412627111>.
- 511 Kapli, Paschalia, Ziheng Yang, and Maximilian J. Telford. 2020. “Phylogenetic Tree Building in  
512 the Genomic Age.” *Nature Reviews. Genetics* 21 (7): 428–44.  
513 <https://doi.org/10.1038/s41576-020-0233-0>.
- 514 Kern, Andrew D, and Daniel R Schrider. 2018. “DiploS/HIC: An Updated Approach to  
515 Classifying Selective Sweeps.” *G3 Genes|Genomes|Genetics* 8 (6): 1959–70.  
516 <https://doi.org/10.1534/g3.118.200262>.
- 517 King, Nicole, and Antonis Rokas. 2017. “Embracing Uncertainty in Reconstructing Early  
518 Animal Evolution.” *Current Biology* 27 (19): R1081–88.  
519 <https://doi.org/10.1016/j.cub.2017.08.054>.
- 520 Kozlov, Alexey M., Diego Darriba, Tomáš Flouri, Benoit Morel, and Alexandros Stamatakis.  
521 2019. “RAxML-NG: A Fast, Scalable and User-Friendly Tool for Maximum Likelihood  
522 Phylogenetic Inference.” *Bioinformatics* 35 (21): 4453–55.  
523 <https://doi.org/10.1093/bioinformatics/btz305>.
- 524 Lanier, Hayley C., and L. Lacey Knowles. 2012. “Is Recombination a Problem for Species-Tree  
525 Analyses?” *Systematic Biology* 61 (4): 691–701. <https://doi.org/10.1093/sysbio/syr128>.
- 526 Laumer, Christopher E., Rosa Fernández, Sarah Lemer, David Combosch, Kevin M. Kocot, Ana  
527 Riesgo, Sónia C.S. S. Andrade, et al. 2019. “Revisiting Metazoan Phylogeny with  
528 Genomic Sampling of All Phyla.” *Proceedings of the Royal Society B: Biological*  
529 *Sciences* 286 (1906). <https://doi.org/10.1098/rspb.2019.0831>.
- 530 Le, Si Quang, and Olivier Gascuel. 2008. “An Improved General Amino Acid Replacement  
531 Matrix.” *Molecular Biology and Evolution* 25 (7): 1307–20.  
532 <https://doi.org/10.1093/molbev/msn067>.
- 533 Li, Yuanning, Xing Xing Shen, Benjamin Evans, Casey W. Dunn, and Antonis Rokas. 2021.  
534 “Rooting the Animal Tree of Life.” *Molecular Biology and Evolution* 38 (10): 4322–33.  
535 <https://doi.org/10.1093/molbev/msab170>.
- 536 Liu, Liang, Lili Yu, and Scott V. Edwards. 2010. “A Maximum Pseudo-Likelihood Approach for  
537 Estimating Species Trees under the Coalescent Model.” *BMC Evolutionary Biology* 10  
538 (1): 302. <https://doi.org/10.1186/1471-2148-10-302>.
- 539 Louppe, Gilles, Louis Wehenkel, Antonio Suter, and Pierre Geurts. 2013. “Understanding  
540 Variable Importances in Forests of Randomized Trees.” *Advances in Neural Information*  
541 *Processing Systems*, 1–9.
- 542 Lundberg. 2017. “A Unified Approach to Interpreting Model Predictions Scott.” *NIPS* 32 (2):  
543 1208–17.
- 544 Maddison, Wayne P. 1997. “Gene Trees in Species Trees.” *Systematic Biology* 46 (3): 523–36.  
545 <https://doi.org/10.1093/sysbio/46.3.523>.

- 546 Mazza, Damiano, and Michele Pagani. 2021. "Automatic Differentiation in PCF." *Proceedings*  
547 *of the ACM on Programming Languages* 5 (POPL): 1–4.  
548 <https://doi.org/10.1145/3434309>.
- 549 Mendes, Fábio K., Andrew P. Livera, and Matthew W. Hahn. 2019. "The Perils of Intralocus  
550 Recombination for Inferences of Molecular Convergence." *Philosophical Transactions of*  
551 *the Royal Society B* 374 (1777). <https://doi.org/10.1098/RSTB.2018.0244>.
- 552 Minh, Bui Quang, Matthew W. Hahn, and Robert Lanfear. 2020. "New Methods to Calculate  
553 Concordance Factors for Phylogenomic Datasets." *Molecular Biology and Evolution* 37  
554 (9): 2727–33. <https://doi.org/10.1093/molbev/msaa106>.
- 555 Molloy, Erin K., and Tandy Warnow. 2018. "To Include or Not to Include: The Impact of Gene  
556 Filtering on Species Tree Estimation Methods." *Systematic Biology* 67 (2): 285–303.  
557 <https://doi.org/10.1093/sysbio/syx077>.
- 558 Moroz, Leonid L., Kevin M. Kocot, Mathew R. Citarella, Sohn Dosung, Tigran P. Norekian,  
559 Inna S. Povolotskaya, Anastasia P. Grigorenko, et al. 2014. "The Ctenophore Genome  
560 and the Evolutionary Origins of Neural Systems." *Nature* 510 (7503): 109–14.  
561 <https://doi.org/10.1038/nature13400>.
- 562 Pease, James B, David C Haak, Matthew W Hahn, and Leonie C Moyle. 2016. "Phylogenomics  
563 Reveals Three Sources of Adaptive Variation during a Rapid Radiation." *PLoS Biology*  
564 14 (2): 1–24. <https://doi.org/10.1371/journal.pbio.1002379>.
- 565 Pisani, Davide, Walker Pett, Martin Dohrmann, Roberto Feuda, Omar Rota-Stabelli, Hervé  
566 Philippe, Nicolas Lartillot, and Gert Wörheide. 2015. "Genomic Data Do Not Support  
567 Comb Jellies as the Sister Group to All Other Animals." *Proceedings of the National*  
568 *Academy of Sciences* 112 (50): 15402–7. <https://doi.org/10.1073/pnas.1518127112>.
- 569 Pollard, Daniel A., Venky N. Iyer, Alan M. Moses, and Michael B. Eisen. 2006. "Pollard et al. -  
570 2006 - Widespread Discordance of Gene Trees with Species Tree in Drosophila Evidence  
571 for Incomplete Lineage Sorting.Pdf." *PLoS Genetics* 2 (10): 1634–47.  
572 <https://doi.org/10.1371/journal.pgen.0020173>.
- 573 Rambaut, Andrew, and Nicholas C. Grassly. 1997. "Seq-Gen: An Application for the Monte  
574 Carlo Simulation of Dna Sequence Evolution along Phylogenetic Trees." *Bioinformatics*  
575 13 (3). <https://doi.org/10.1093/bioinformatics/13.3.235>.
- 576 Redmond, Anthony K., and Aoife McLysaght. 2021. "Evidence for Sponges as Sister to All  
577 Other Animals from Partitioned Phylogenomics with Mixture Models and Recoding."  
578 *Nature Communications* 12 (1). <https://doi.org/10.1038/s41467-021-22074-7>.
- 579 Settles, Burr. 2012. *Active Learning*. Morgan & Claypool.
- 580 Simion, Paul, Hervé Philippe, Denis Baurain, Muriel Jager, Daniel J. Richter, Arnaud Di Franco,  
581 Béatrice Roure, et al. 2017. "A Large and Consistent Phylogenomic Dataset Supports  
582 Sponges as the Sister Group to All Other Animals." *Current Biology* 27 (7): 958–67.  
583 <https://doi.org/10.1016/j.cub.2017.02.031>.
- 584 Smith, Stacey D., Matthew W. Pennell, Casey W. Dunn, and Scott V. Edwards. 2020.  
585 "Phylogenetics Is the New Genetics (for Most of Biodiversity)." *Trends in Ecology and*  
586 *Evolution* 35 (5): 415–25. <https://doi.org/10.1016/j.tree.2020.01.005>.
- 587 Snir, Sagi, and Satish Rao. 2012. "Quartet MaxCut: A Fast Algorithm for Amalgamating Quartet  
588 Trees." *Molecular Phylogenetics and Evolution* 62 (1): 1–8.  
589 <https://doi.org/10.1016/j.ympev.2011.06.021>.

- 590 Solis-Lemus, Claudia, Shengwen Yang, and Leonardo Zepeda-Nunez. 2022. “Accurate  
591 Phylogenetic Inference with a Symmetry-Preserving Neural Network Model.” arXiv.  
592 <https://doi.org/10.48550/arXiv.2201.04663>.
- 593 Strimmer, K, and A von Haeseler. 1996. “Quartet Puzzling: A Quartet Maximum-Likelihood  
594 Method for Reconstructing Tree Topologies.” *Molecular Biology and Evolution* 13 (7):  
595 964. <https://doi.org/10.1093/oxfordjournals.molbev.a025664>.
- 596 Suvorov, Anton, Joshua Hochuli, and Daniel R. Schrider. 2020. “Accurate Inference of Tree  
597 Topologies from Multiple Sequence Alignments Using Deep Learning.” *Systematic  
598 Biology* 69 (2): 221–33.
- 599 Vachaspati, Pranjal, and Tandy Warnow. 2015. “ASTRID: Accurate Species TRees from  
600 Internode Distances.” *BMC Genomics* 16 (10): S3. [https://doi.org/10.1186/1471-2164-16-  
601 S10-S3](https://doi.org/10.1186/1471-2164-16-S10-S3).
- 602 Whelan, Nathan V., and Kenneth M. Halaných. 2017. “Who Let the CAT out of the Bag?  
603 Accurately Dealing with Substitutional Heterogeneity in Phylogenomic Analyses.”  
604 *Systematic Biology* 66 (2): 232–55. <https://doi.org/10.1093/sysbio/syw084>.
- 605 Whelan, Nathan V., Kevin M. Kocot, Leonid L. Moroz, and Kenneth M. Halaných. 2015. “Error,  
606 Signal, and the Placement of Ctenophora Sister to All Other Animals.” *Proceedings of  
607 the National Academy of Sciences* 112 (18): 5773–78.  
608 <https://doi.org/10.1073/pnas.1503453112>.
- 609 Wickett, Norman J., Siavash Mirarab, Nam Nguyen, Tandy Warnow, Eric Carpenter, Naim  
610 Matasci, Saravanaraj Ayyampalayam, et al. 2014. “Phylotranscriptomic Analysis of the  
611 Origin and Early Diversification of Land Plants.” *Proceedings of the National Academy  
612 of Sciences of the United States of America* 111 (45): E4859–68.  
613 <https://doi.org/10.1073/pnas.1323926111>.
- 614 Zaharias, Paul, Martin Grosshauser, and Tandy Warnow. 2022. “Re-Evaluating Deep Neural  
615 Networks for Phylogeny Estimation: The Issue of Taxon Sampling.” *Journal of  
616 Computational Biology* 29 (1): 74–89. <https://doi.org/10.1089/cmb.2021.0383>.
- 617 Zhang, Chao, Maryam Rabiee, Erfan Sayyari, and Siavash Mirarab. 2018. “ASTRAL-III:  
618 Polynomial Time Species Tree Reconstruction from Partially Resolved Gene Trees.”  
619 *BMC Bioinformatics* 19 (S6): 153. <https://doi.org/10.1186/s12859-018-2129-y>.
- 620 Zou, Zhengting, Hongjiu Zhang, Yuanfang Guan, and Jianzhi Zhang. 2019. “Deep Residual  
621 Neural Networks Resolve Quartet Molecular Phylogenies.” *Molecular Biology and  
622 Evolution*, September, 787168. <https://doi.org/10.1101/787168>.

623

624

625

626

627

## 628 8 Tables and Figures

629

<b>Parameter</b>	<b>Values</b>
Alignment length (AA residues)	50, 100, 200, 500, 1000
Simulation sequence model	LG+G, WAG+G
Inference sequence model	LG+G, WAG+G
EBL (coalescent units)	20, 30, ..., 290, 300
IBL (coalescent units)	0.01, 0.025, 0.05, ..., 0.175, 0.2, 0.25, ..., 1.75, 2, 3, ..., 19, 20
Branch length to root (coalescent units)	20

630 Table 1: Simulated training conditions. At least 1,000 loci were simulated for each parameter  
631 combination.

632

633

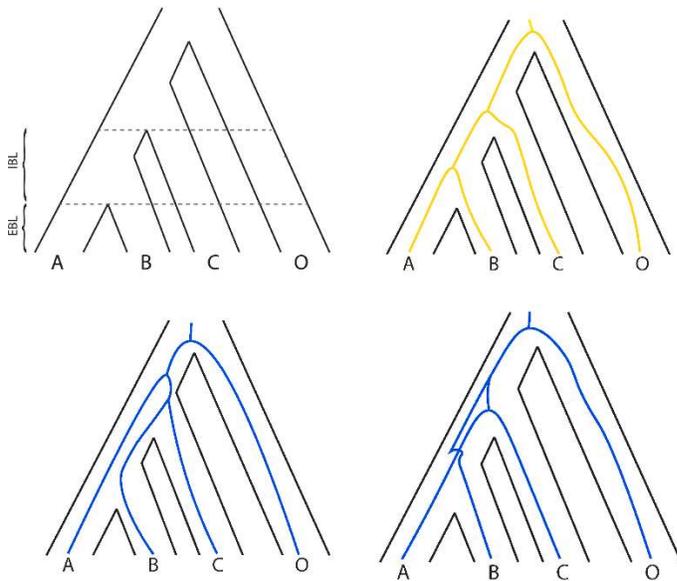
<b>Data Matrix</b>	<b>C</b>	<b>P</b>	<b>Pa</b>
Borowiec2015_Best108	6	<b>114</b>	60
Borowiec2015_Total1080	6	<b>122</b>	64
Nosenko2013_ribo_11057	239	<b>338</b>	242
Ryan2013_est	203	<b>854</b>	470
Ryan2013_est_only_choanozoa	87	<b>229</b>	143
Ryan2013_est_only_holozoa	134	<b>333</b>	217
Whelan2015_D10	15879	<b>36816</b>	23400
Whelan2015_D10_only_choanozoa	1916	<b>8296</b>	4647
Whelan2015_D1_only_choanozoa	1990	<b>8488</b>	4786
Whelan2015_D1_only_holozoa	6125	<b>12223</b>	7872
Whelan2015_D20_only_choanozoa	478	<b>978</b>	602
Whelan2017_full	250	<b>542</b>	342
Whelan2017_full_only_choanozoa	747	<b>1068</b>	747

634 Table 2: Number of quartets for which the most common gene tree topology supports  
635 Ctenophora-sister (C), Porifera-sister (P), or Parahoxozoa-sister (Pa). “Data Matrix” refers to the  
636 original papers from which datasets were obtained. For a full list and description of all datasets  
637 see Supplementary Table 2.

638

639

640



641

642 Figure 1: ILS among three taxa can produce one concordant (yellow) and two discordant (blue)

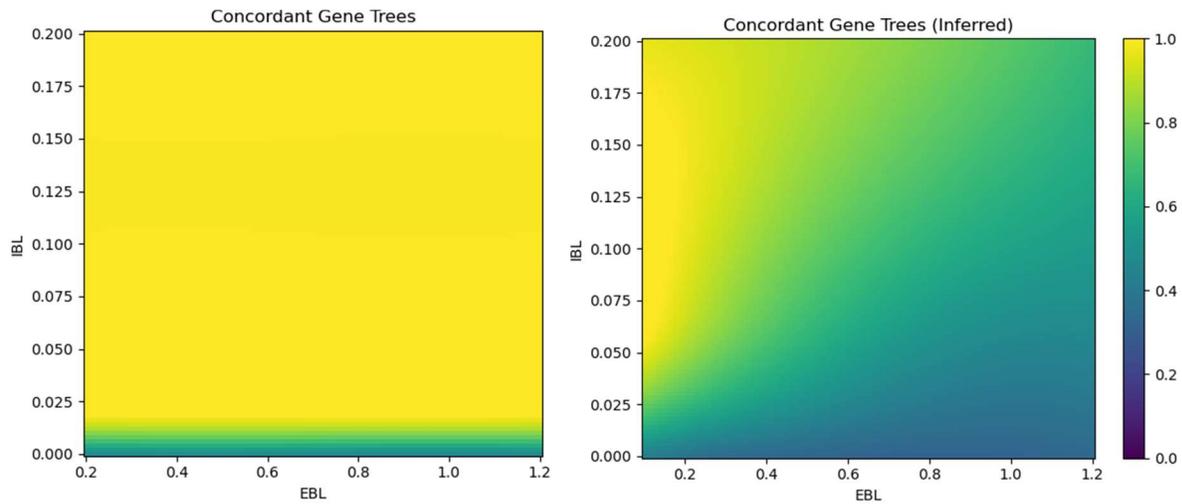
643 gene tree topologies. The probability of concordance depends only on the internal branch length

644 (IBL), while gene tree inference error depends on both IBL and external branch length (EBL).

645

646

647

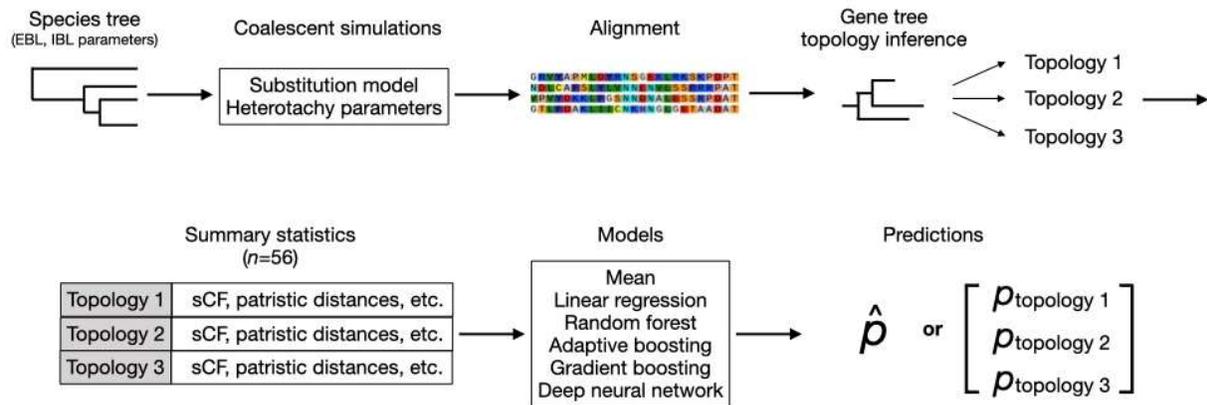


648

649 Figure 2: The fraction of true (left) versus inferred (right) gene tree topology frequencies that  
650 match the species tree topology ( $p$ ) as a function of IBL and EBL. X- and Y-axes are in units of  
651 expected (amino acid) sequence divergence.

652

653



654

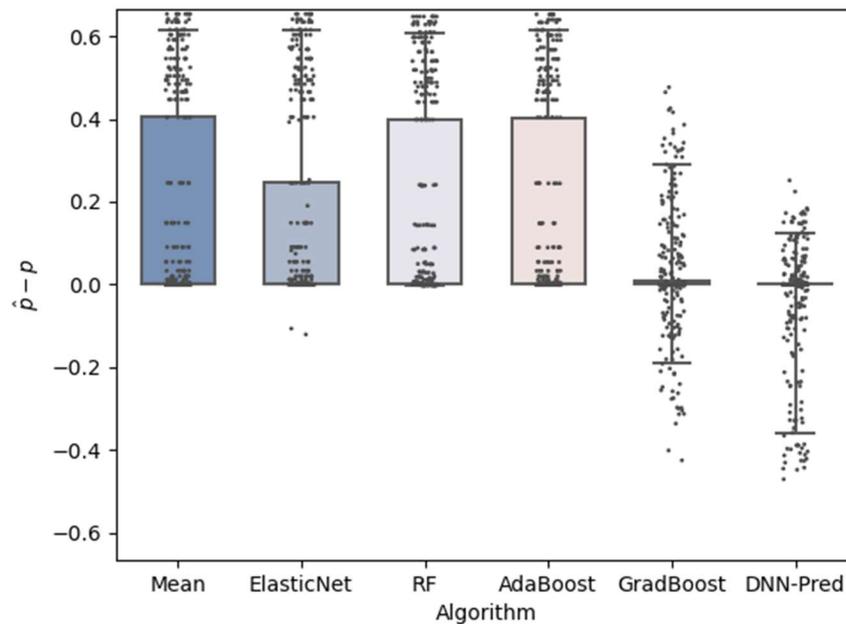
655 Figure 3: Flow chart of data processing pipeline. For a given set of species tree parameters,  
 656 coalescent simulations are carried out in conjunction with the simulation of an associated  
 657 alignment. For each alignment, gene trees are inferred, and sorted into bins by the inferred tree  
 658 topology. For each of the 3 topology bins, 56 summary statistics are computed, for a final  
 659 feature vector of length 168. A variety of supervised machine learning models are trained to  
 660 predict either  $\hat{p}$ , the probability of concordance (which can be used as a binary classifier with  
 661 suitable threshold), or the probability of each species tree topology.

662

663

664

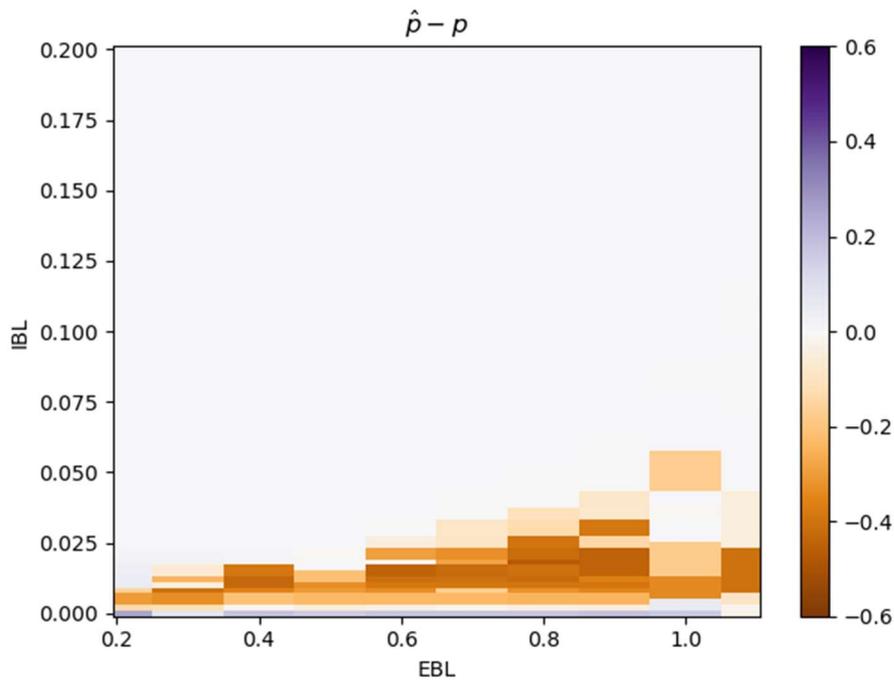
665



666

667 Figure 4: Error of all SML methods on a collection of 250-gene, 500-residue datasets. EBL  
668 ranges from 0.2 to 1.0, IBL from 0.001 to 0.18. The baseline Mean predictor is the mean  $p$  across  
669 all samples in the training sets. The other predictors from left to right are: regularized linear  
670 regression (ElasticNet), random forest (RF), adaptive boosting with decision trees (AdaBoost),  
671 Gradient Boosting with decision trees (GradBoost), and the deep neural network predictor  
672 (DNN-Pred). Boxes represent the 25% and 75% quantiles, whiskers the 5% and 95% quantiles.  
673

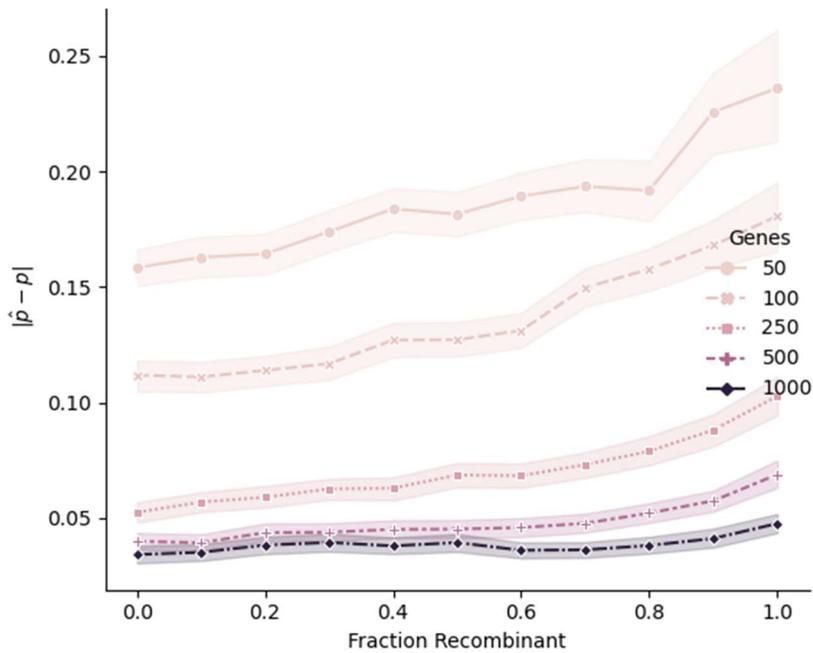
674



675

676 Figure 5: Accuracy of DNN-Pred. Average distance between true and predicted  $p$  for 500-gene,  
677 500-residue datasets. Branch length units on both axes are in expected amino acid sequence  
678 divergence.

679

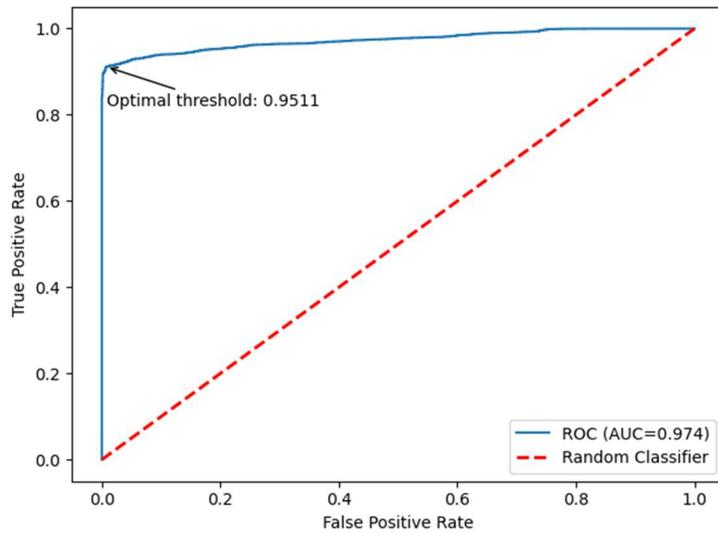


680

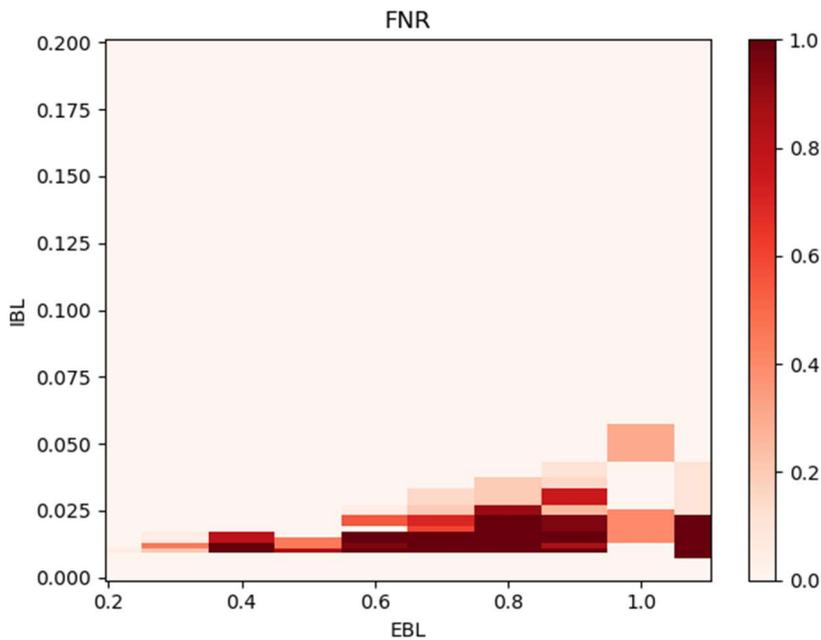
681 Figure 6: Intragenic recombination reduces the accuracy of  $p$ . Average absolute error between  
682 true and predicted  $p$  for datasets of length 500 residues is shown across a range of parameter  
683 space. Recombinant genes are drawn uniformly from 2-, 3-, and 4-block conditions, and both  
684 the number of genes and fraction of recombinant genes per test set are varied.

685

686



687



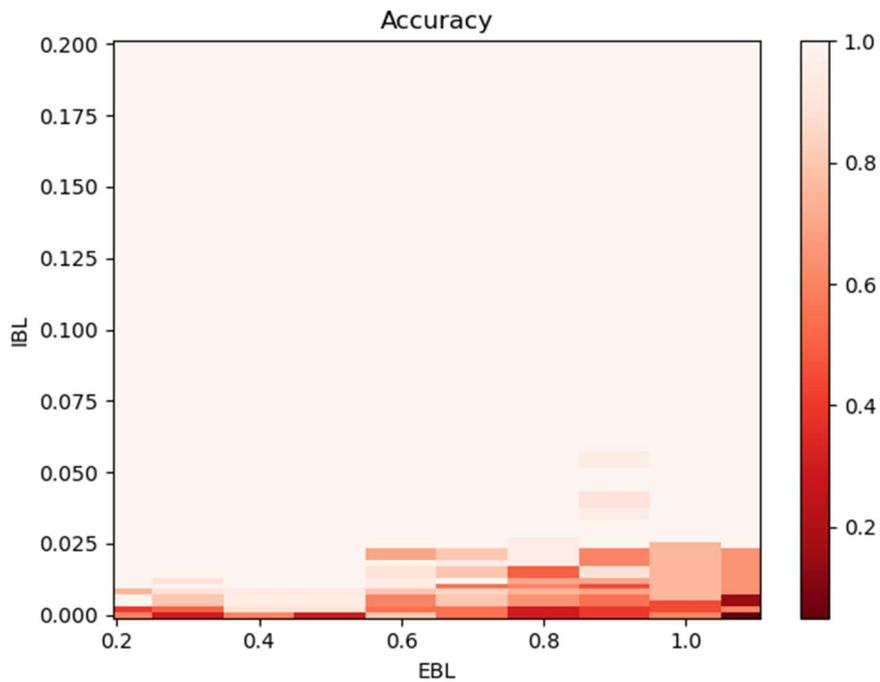
688

689 Figure 7: Performance of DNN-Class on 500-gene, 500-residue datasets. a) Receiver Operating  
690 Characteristic (ROC) curve. The area under the curve (AUC) is also reported. b) False negative  
691 rate (FNR) at the optimal threshold in panel a) as a function of internal and external branch  
692 length. Branch length units on both axes are in expected amino acid sequence divergence.

693

694

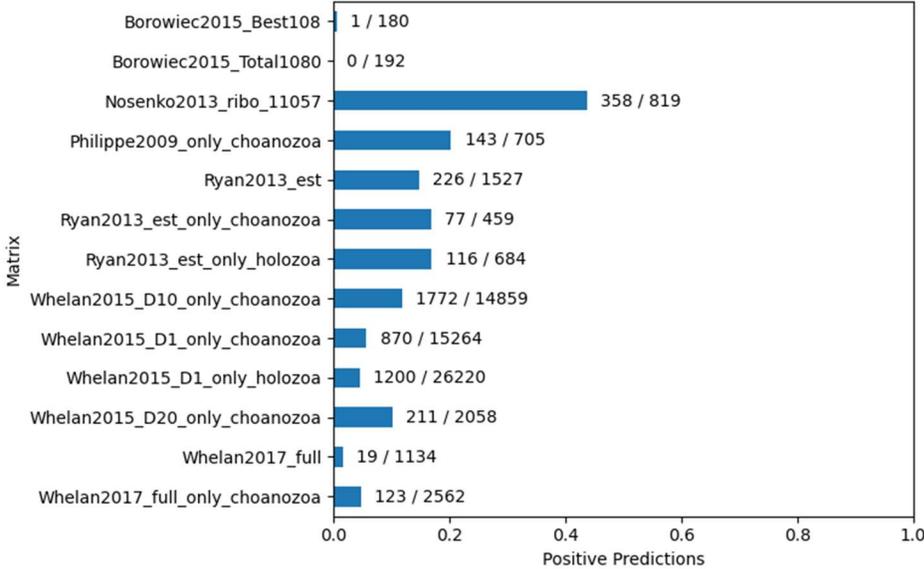
695  
696  
697  
698



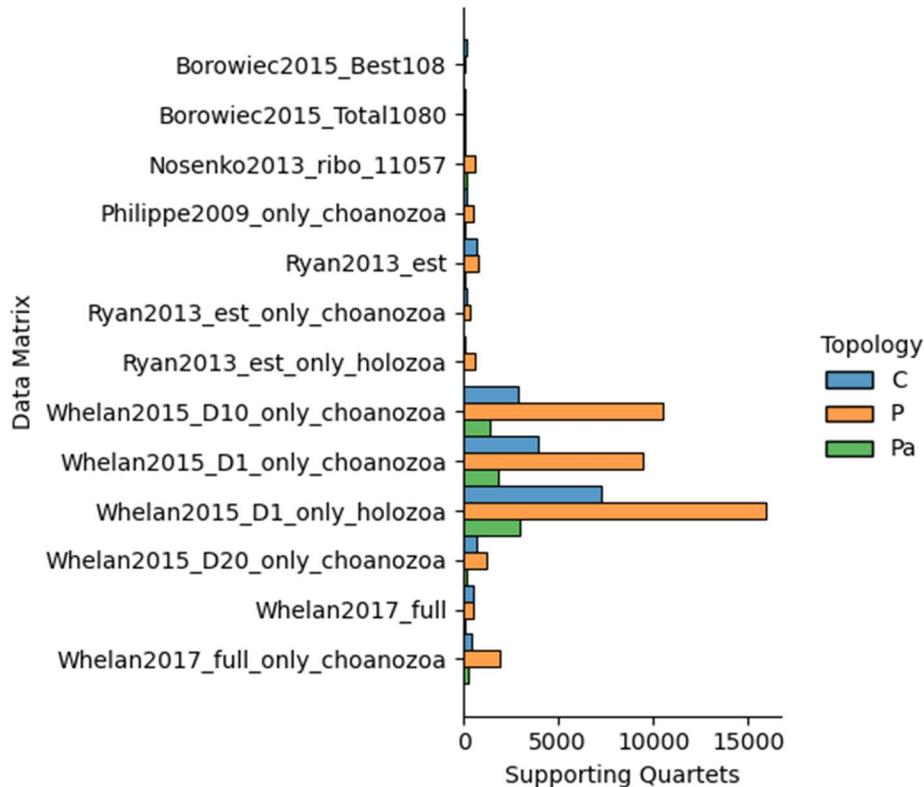
699  
700  
701  
702  
703  
704

Figure 8: Accuracy of DNN-Top on 500-gene, 500-residue datasets. Accuracy is measured as the fraction of topology predictions that match the true species tree topology. Branch length units on both axes are in expected amino acid sequence divergence.

705



706



707

708 Figure 10: Discordance in the Metazoa. a) The fraction of cases with no discordance (DNN-  
 709 Class predicts  $p > 0.9$ ) across all Cte/Por/Par/Out quartets for each of 13 Metazoan data  
 710 matrices. The number of such cases and total number of quartets in each dataset is also shown.  
 711 b) Species tree topologies predicted by DNN-Top for the Metazoa datasets. Across all data  
 712 matrices, the majority of quartets support the Porifera-sister species topology.

714