1 **Towards the next generation of species delimitation methods: an**

2 **overview of Machine Learning applications**

3 Matheus M. A. Salles[a*], Fabricius M. C. B. Domingos[a]

4 [a]Departamento de Zoologia, Universidade Federal do Paraná, Curitiba 81531-

5 980, Brazil

6 [*]Corresponding author: matheus.salles@ufpr.br

7

8 ABSTRACT

9 Species delimitation is the process of distinguishing between populations of the

10 same species and distinct species of a particular group of organisms. Various

11 methods exist for inferring species limits, whether based on morphological,

12 molecular, or other types of data. In the case of methods based on DNA

13 sequences, most of them are rooted in the coalescent theory. However,

14 coalescence-based models have limitations, especially regarding complex

15 evolutionary scenarios, large datasets, and varying genetic data types. In this

16 context, machine learning (ML) can be considered as a promising analytical tool,

17 and provides an effective way to explore dataset structures when species-level

18 divergences are hypothesized. In this review, we examine the use of ML in

19 species delimitation and provide an overview and critical appraisal of existing

20 workflows. We also provide simple explanations on how the main types of ML

21 approaches operate, which should help uninitiated researchers and students

22 interested in the field. Our review suggests that while current ML methods

23 designed to infer species limits are analytically powerful, they also present

24 specific limitations and should not be considered as definitive alternatives to

25 coalescent methods for species delimitation. On the other hand, such variability

26  might also represent an advantage, highlighting the flexibility of ML algorithms.

27  Future enterprises should consider the constraints related to the use of simulated

28  data, as in other model-based methods relying on simulations. We also propose

29  best practices for the use of ML methods in species delimitation, offering insights

30  into potential future applications. We expect that the proposed guidelines will be

31  useful for enhancing the accessibility, effectiveness, and objectivity of ML in

32  species delimitation.

33  *Key words*: bioinformatics, molecular data, speciation, phylogenetics, artificial

34  intelligence, deep learning.

35

36  **1. Introduction**

37  *1.1. Inferring species limits*

38      Species represent fundamental entities across all biological disciplines.

39  Consequently, the review, categorization, and characterization of taxa within this

40  level constitute a pivotal aspect of biodiversity research (Bortolus, 2008; Vink et

41  al., 2012; Ely et al., 2017). The process of identifying, characterizing, and defining

42  a species is data-intensive and entails various practical dimensions. This

43  complexity arises from managing extensive biological data and dealing with a

44  range of theoretical elements, from the establishment of homologies, to taxon-

45  specific traits, and the very philosophical notion of species. Furthermore,

46  conceptual issues surrounding the definition of species concepts still attract

47  debates among taxonomists and evolutionary biologists (Pante et al., 2015;

48  Zachos, 2016). These discussions reach the realms of philosophy, because a

49  multitude of data and methodologies will probably not fully solve many

50  fundamental questions surrounding the nature of species (Zachos, 2016; Wilkins

51 et al., 2022), or the 'species ontology' (what a species really is or represents). A

52 complete resolution on this subject remains elusive, as it intertwines the empirical

53 evidence biologists are able to extract from nature with philosophical definitions

54 surrounding species concepts (Pigliucci, 2003).

55       One of the most popular modern definitions is the 'Biological Species

56 Concept' (de Queiroz, 2005a; Zachos, 2016), which defines species as

57 interbreeding populations reproductively isolated from others (Mayr, 1969; 1996;

58 2000). Yet, many challenges to this concept emerged throughout the years as

59 empirical data clearly shows that the history of life on Earth does not fit into a

60 bifurcating process (Edwards et al., 2016; Mallet et al., 2016), and a clear

61 delineation of reproductive barriers is hindered by instances of asexual

62 reproduction, natural hybridization and gene flow (Arnold, 1992; Shurtliff, 2013;

63 Gompert et al., 2017). Hence, taxonomists and evolutionary biologists must

64 recognize that multiple species definitions will coexist in the practice of species

65 delimitation, and these are usually chosen based on the biological context of the

66 organisms under study.

67       Another important concept, the General Lineage Concept (GLC), diverges

68 from many others by prioritizing the recognition of independently evolving

69 lineages over specific biological criteria such as reproduction or morphology (de

70 Queiroz, 1998; 1999; 2007). According to the GLC, a species is defined as an

71 independently evolving metapopulation lineage, emphasizing each species'

72 unique evolutionary identity across time and space (de Queiroz, 2007). While

73 unique morphological, ecological, or any other biological trait might be considered

74 relevant in supporting the investigation of the speciation process, they are not

75 mandatory criteria for species definition under the GLC perspective, but rather

additional evidence supporting lineage separation (de Queiroz, 2007). Thus, this concept accounts for the contingent nature of the speciation process, where different biological properties may support species limits in varying degrees. It also emphasizes the need for multiple lines of evidence to corroborate hypotheses of species divergence, aligning with Integrative Taxonomy approaches (Wiens & Penkrot, 2002; Dayrat, 2005; Padial et al., 2010; Fujita et al., 2012).

The GLC also provides a theoretical distinction between the 'species ontology problem' (what a species is) and the 'delimitation problem' (how to operationally distinguish among putative species) (de Queiroz, 2007). Interestingly, while a clear relationship exists between these components, namely the species concept and species delimitation, historically, a significant part of the scientific efforts has focused on the former (see Sites Jr and Marshall, 2004; Wiens, 2007; de Queiroz, 2011; Hausdorf, 2011). The development of theoretical considerations related to species delimitation, in particular that based on molecular data, occurred mainly in the last two decades, accompanied by the introduction of new criteria and statistical methods (Lukhtanov, 2019; Rannala and Yang, 2020). Historically, identifying species limits, and describing new species, have primarily relied on morphological data (Wiens, 2007; Rannala, 2015; Rannala and Yang, 2020). However, morphological traits can be influenced by environmental factors, leading to convergence or divergence without necessarily reflecting genetic or evolutionary relationships between lineages (Price et al., 2003; Wake et al., 2011; Jarvis et al., 2014). Thus, genomic data has emerged as a crucial tool for inferring species limits, offering a more objective

100   approach for species delimitation (Fujita et al., 2012), while complementing

101   traditional morphological methods (Jörger and Schrödl, 2013).

102   Modern species delimitation methods (SDMs) aiming at identifying

103   evolutionary units (Tautz et al., 2003; Vogler and Monaghan, 2007) mostly

104   operate with molecular data under the principles of Coalescent Theory, notably,

105   the multispecies coalescent (MSC; Rannala and Yang, 2003; Degnan and

106   Rosenberg, 2009). The MSC analytical framework addresses various

107   evolutionary assumptions while also managing different types of problems, such

108   as conflicts among different gene trees, **incomplete lineage sorting** (terms in

109   bold are defined in the Glossary, available in Appendix A), and errors in

110   phylogenetic inference (Knowles & Carstens, 2007; Carstens et al., 2013; Jacobs

111   et al., 2018). The use of modern SDMs has also grown due to advancements in

112   statistical frameworks for phylogenetic inference (Edwards, 2009; O'Meara,

113   2012), along with Molecular Biology tools (e.g., next-generation sequencing

114   (NGS; Slatko et al., 2018) and Bioinformatics (Searls, 2010).

115   Nonetheless, using SDMs with genetic data may fail to distinguish

116   population structure from species-level divergence (Sukumaran and Knowles,

117   2017), and may also be affected by other issues associated with the reliance on

118   the MSC model (Rannala and Yang, 2003; Degnan and Rosenberg, 2009;

119   Edwards, 2009; Fujita et al., 2012). Some methods also have their functionality

120   and performance compromised in scenarios when there is introgression between

121   groups that constitute potential species (Rannala and Yang, 2010; Leaché et al.,

122   2014; Jackson et al., 2017), and are more reliable in situations where gene flow

123   ceases immediately after population divergence (Fujita et al., 2012; Smith and

124   Carstens, 2020). Also, simulations have shown that ignoring gene flow leads the

125 MSC to overestimate **population sizes** and underestimate divergence times

126 (e.g., Leaché et al., 2014). Hence, the effectiveness of the MSC framework is

127 limited, to some extent, when additional processes influence divergence during

128 speciation (Smith and Carstens, 2020). In any case, different SDMs have varying

129 capabilities to address difficult evolutionary scenarios, and while such methods

130 may introduce biases in certain situations, they are not inherently useless.

131

132 *1.2. Machine learning, evolutionary biology, and species delimitation*

133     **Machine learning (ML)**, a branch of artificial intelligence (AI) known for its

134 computational efficiency and predictive accuracy, has recently gained popularity

135 in Evolutionary Biology mainly due to its ability to analyze and process large,

136 complex, and high-dimensional datasets (Chicco, 2017; Borowiec et al., 2022;

137 Fountain-Jones et al., 2021; Greener et al., 2021; Morimoto et al., 2021). In

138 general terms, ML can be defined as a group of computational programs that can

139 learn through experience (E) with respect to a class of tasks (T), and an

140 evaluation measure (P), if its performance on the tasks of T, evaluated by P,

141 increases with E (Mitchell, 1997). Many ML algorithms are known to be extremely

142 useful in various aspects of biology. This includes photo-based species

143 identification (Wäldchen and Mäder 2018), morphology-based species

144 delimitation and description (Domingos et al., 2014; Breitman et al., 2018),

145 biodiversity monitoring (McClure et al., 2020), behavioural studies (Valletta et al.,

146 2017; Wang, 2019), DNA sequencing (Libbrecht and Noble, 2015; Liu, 2019),

147 population genetics (Sheehan and Song 2016; Schrider and Kern, 2018; Fonseca

148 & Carstens, 2024), ecology (Christin et al., 2019; Scalon et al., 2020; Pichler et

149 al., 2020; Lürig et al., 2021; Silva et al., 2024), medicine (Sidey-Gibbons and

Sidey-Gibbons, 2019), microbiology (Qu et al., 2019), and more (see Borowiec et al., 2022; Fountain-Jones et al., 2021; Morimoto et al., 2021).

Therefore, its potential in evolutionary biology, and particularly in species delimitation, is evident (Karbstein et al., 2023). Specific examples can already be found in studies involving model selection in demography and phylogeography (Pudlo et al., 2016; Fonseca et al., 2021), speciation (Blischak et al., 2021), phylogenetics (Suvorov et al., 2020; Solis-Lemus et al., 2022 preprint; Smith & Hahn, 2023; Zaharias et al., 2022; Mo et al., 2024), and species delimitation (Pei et al., 2018; Derkarabetian et al., 2019; Smith & Carstens, 2020; Pyron et al., 2023), with the last one forming the primary focus of this review.

In the following sections, we provide an overview of ML applications in the context of species delimitation, with an emphasis on those that operate using molecular data.

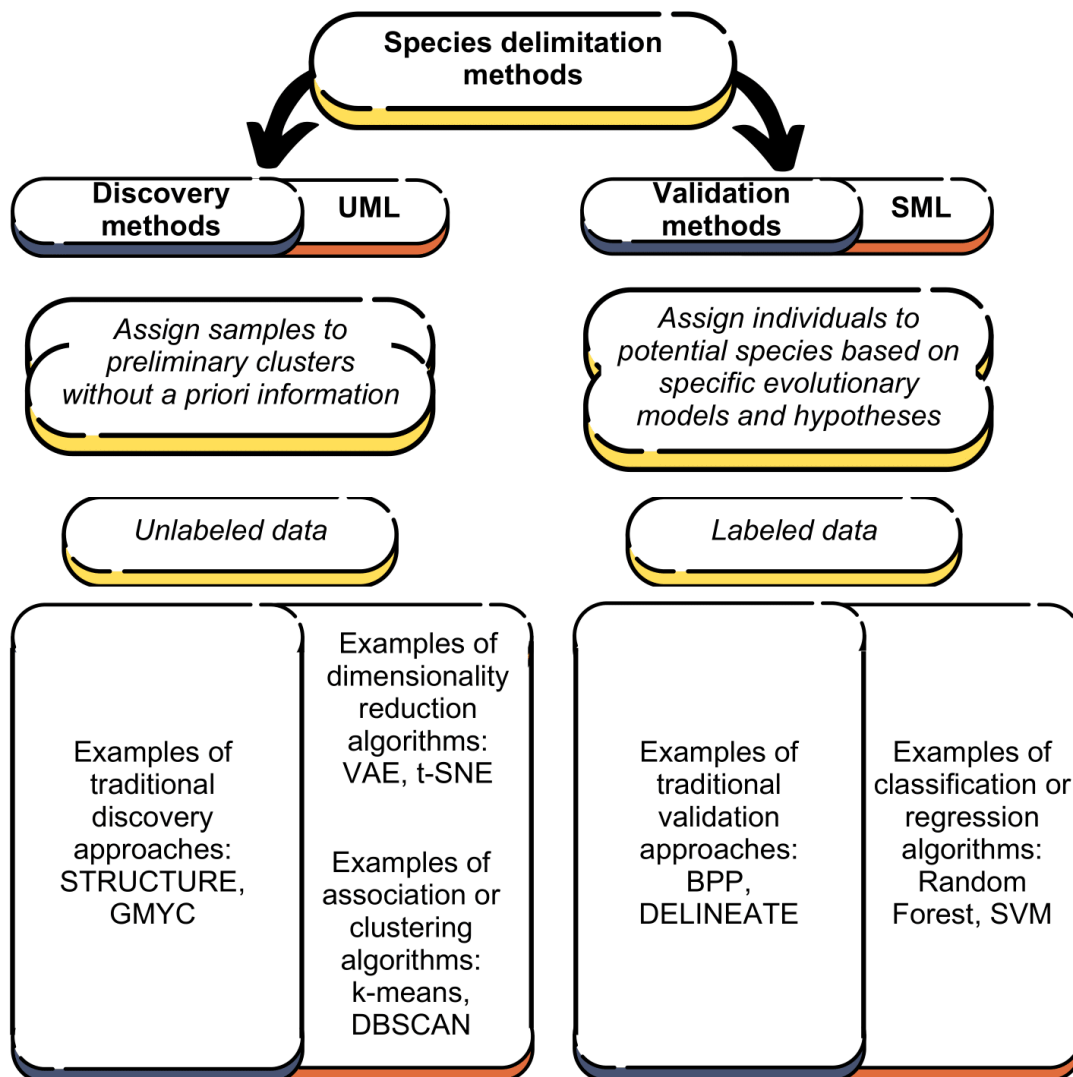## 2. Current ML applications for species delimitation

In the same way that there are two primary categories of ML, namely supervised and unsupervised learning (SML and UML, respectively), species delimitation methods can also be broadly categorized into two main groups: discovery and validation (see Carstens et al., 2013; Rannala, 2015). Discovery approaches involve grouping samples without prior information (Pons et al., 2006; O'Meara, 2010; Huelsenbeck et al., 2011), while validation approaches require researchers to first assign the samples to potential lineages (species hypotheses) before testing them (Flouri et al., 2018; Sukumaran et al., 2021). This draws a conceptual parallel between traditional discovery approaches and UML methods, and between validation methods and supervised algorithms (Fig.

175 1). Also, it is important to note that ML methods are likelihood-free species

176 delimitation approaches, offering several advantages over **likelihood-based**

177 **approaches**, such as eliminating the need for complex statistical calculations,

178 making them computationally efficient and suitable for analyzing large datasets

179 with many taxa.

180

181

**Fig. 1.** Comparative diagram categorizing species delimitation methods and machine learning algorithms, along with some of their key characteristics. Species delimitation methods can be broadly categorized as discovery and validation methods, akin to unsupervised and supervised machine learning algorithms, respectively.

182
183
184
185
186

187

188    Below, we present a comprehensive overview of recently applied ML

189    methods in the domain of molecular species delimitation, emphasizing their

190    computational attributes and underlying assumptions. Our selection process

191    involved a thorough search across scientific literature repositories, databases,

192    and online journals, with a specific emphasis on studies featuring ML methods

193    and workflows explicitly designed for species limits inference. We prioritized

194    research projects that either introduced novel methodologies (see Table 1) or

195    enhanced and tested existing techniques in this context (Table A.1 in Appendix

196    B). In our selection process, we focused exclusively on projects directly dedicated

197    to species delimitation, despite the abundant literature on ML within related fields

198    such as demography, population genetics, and phylogeography. Additionally, our

199    emphasis is on methods designed for analyzing DNA sequence data. The

200    categorized methods include SML, UML, and **deep learning**. While the backend

201    processes may differ among such ML categories, their main goal when it comes

202    to species delimitation usually remains the same: to analyze a given set of test

203    data and classify it into distinct outcomes that define the species represented

204    within the data.

205    Some studies applied ML techniques using other types of data rather than

206    molecular information, such as morphology or ecology, for species delimitation

207    and integrative taxonomy. A brief exploratory section regarding these particular

208    studies can be found in Appendix B.

209

210

211    **Table 1.** List of proposed ML applications specifically designed to work on inferences about species limits.

| Reference | Languages | Category | Algorithms | Simulator | Input | Data representation |
|---|---|---|---|---|---|---|
| CLADES: A Classification-based Machine Learning Method for Species Delimitation from Population Genetic Data (Pei et al., 2018)[1] | python | SML | Support vector machines | MCcoal | Multiple sequence alignment (MSA) or SNP matrix | Population genetics summary statistics |
| A demonstration of unsupervised machine learning in species delimitation (Derkarabetian et al., 2019)[2] | R/python | UML | Variational autoencoders and t-Distributed Stochastic Neighbor Embedding | NA | SNP data matrix | One-hot-encoding of the SNP data matrix and *axis* from a discriminant analysis of principal components |
| Process-based species delimitation leads to identification of more biologically relevant species (Smith & Carstens, 2020)[3] | python | SML | Random forest | fastsimcoal | SNP data matrix | Folded multi-dimensional SFS |
| Coalescent-based species delimitation meets deep learning: Insights from a highly fragmented cactus system (Perez et al., 2021)[4] | python | Deep learning | Convolutional neural networks | ms | SNP data matrix | Matrices (as images), with genotypes encoded as higher or lower frequency states |
| Speciation Hypotheses from Phylogeographic Delimitation Yield an Integrative Taxonomy for Seal Salamanders (*Desmognathus monticola*) (Pyron et al., 2023)[5] | R | UML | Self-organizing maps (SOMs) | NA | SNP data matrix | SNP matrix, in which the rows are individual specimens, the columns are the 2-4 possible states at each SNP locus, and the entries are the frequency of that state |

212    Online repositories where it is possible to find more information about the currently existing platforms. [1] https://github.com/pjweggy/CLADES;
213    [2] https://www.sciencedirect.com/science/article/abs/pii/S1055790319301721; [3] https://github.com/meganlsmith/delimitR; [4] https://github.com/manolofperez/CNN_spDelimitation_Piloso;
214    [5] https://github.com/kyleaoconnell22/Pyron_et_al_UML_sp_delim/tree/main

215    *2.1 Discovery and unsupervised methods*

216    Unsupervised machine learning (UML) relies solely on the inherent data structure to

217    find patterns within the data, whether by clustering similar data points together, reducing the

218    dimensionality of the data while retaining essential information, or by identifying unusual

219    patterns or outliers, which may indicate errors or novel phenomena (Hastie et al., 2009;

220    Libbrecht and Noble, 2015; Dike et al., 2018). Consequently, UML algorithms operate

221    without predefined assumptions about the dataset underlying structure, population

222    parameters, species numbers, or sample categorization, making them particularly suitable

223    for species delimitation where no prior hypotheses are put forward.

224    In terms of delimiting species, clustering or dimensionality reduction UML algorithms

225    are generally employed (Fig. 2). Clustering methods group input data into subsets, where

226    samples with high similarities are placed in the same cluster and exhibit less similarity with

227    samples in other clusters. Dimensionality reduction focuses on compressing data to identify

228    a smaller distinct set of variables that could capture essential features of the original data,

229    while minimizing information loss. Thus, UML dimensionality reduction may provide intuitive

230    data visualization and accommodate various data types (Libbrecht and Noble, 2015), being

231    particularly effective when coalescent methods tend to oversplit potential species

232    (Derkarabetian et al., 2019). In sum, UML algorithms enable the simultaneous use of diverse

233    data types, mainly by extracting and condensing the necessary information to identify limits

234    between biological groups (Pyron, 2023; Pyron et al., 2023).
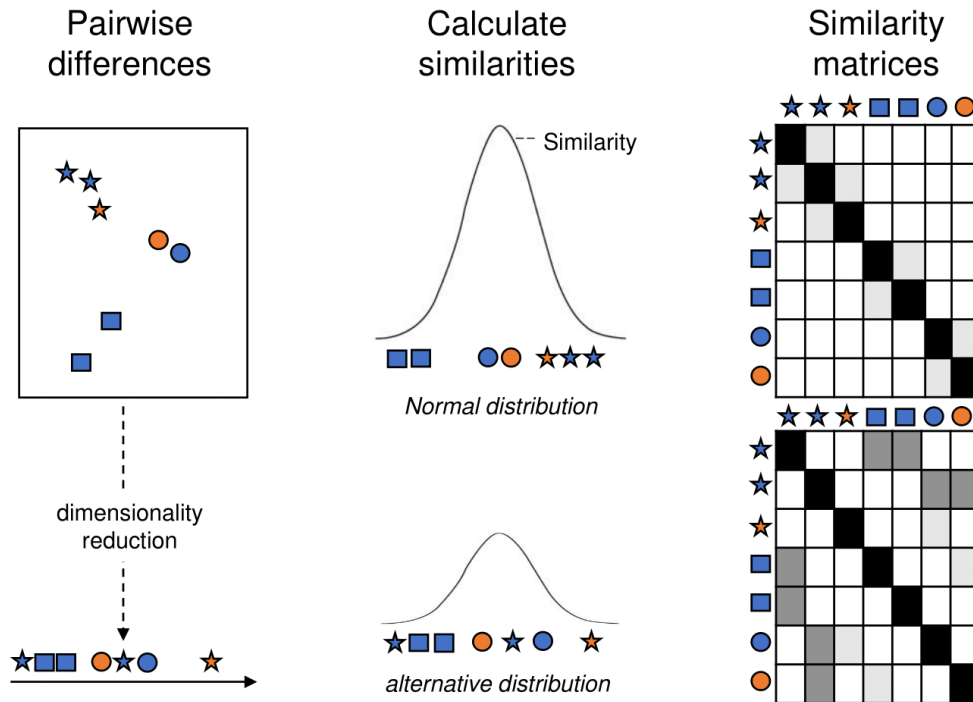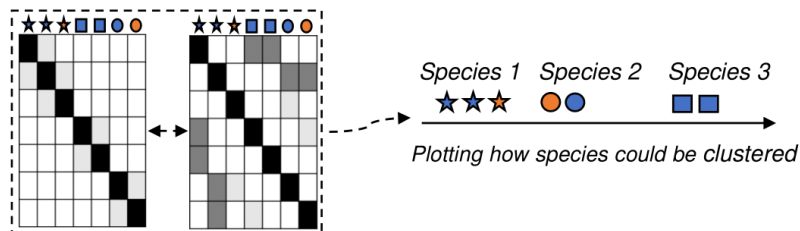
235

236

**a)** SNPs matrix (or transformations from it) representing the input data



**b)**

Pairwise differences

Calculate similarities

Similarity matrices

dimensionality reduction

Similarity

Normal distribution

alternative distribution

**c)** Minimize diferences, rearrange low-dimension matrix and iteratively compare it with the original one

Species 1   Species 2   Species 3

Plotting how species could be clustered

Fig. 2. Diagram outlining a potential UML workflow for species delimitation, utilizing the t-SNE algorithm (inspired by Derkarabetian et al., 2019). a) Data representation is the initial step, and it varies depending on the chosen ML tool, which may work with sequence data, SNP matrices, or population genetics metrics extracted from them. b) t-SNE, as a dimensionality reduction technique, iteratively finds a lower-dimensional representation of the original data. It identifies local similarity spaces between sample pairs by analyzing Gaussian and lower-dimensional distributions, such as the Cauchy or t-student with one degree of freedom. c) The algorithm's goal is to align the new similarity matrix with the original data by iteratively moving data points closer to their nearest neighbors in the higher-dimensional space and away from more distant ones. This process continues until the maximum number of iterations is reached or no further improvements can be made, resulting in the proper grouping of samples based on their similarities (e.g., individuals or populations assigned to a species based on the chosen data representation).

251     Derkarabetian et al. (2019) conducted a study to assess the performance of UML and

252     deep learning methods for species delimitation. Their research highlighted the effectiveness

253     of variational autoencoder (VAE) and t-Distributed Stochastic Neighbor Embedding (t-SNE)

254     algorithms for accurately identifying species clusters. In the case of VAE, single-nucleotide

255     polymorphism (SNP) matrices were converted via 'one-hot coding', where nucleotides were

256     transformed into binary variables (e.g., A = [1, 0, 0, 0]; C = [0, 1, 0, 0], and so on), including

257     ambiguous bases (e.g., M = [0.5, 0.5, 0.0, 0.0]). This VAE approach employed multiple

258     layers of encoding to compress high-dimensional input data, followed by the reconstruction

259     of data through successive decoding layers. The latent variables, represented as a normal

260     distribution with mean ($\mu$) and standard deviation ($\sigma$), offered a two-dimensional depiction of

261     the SNP matrix, facilitating a clear visualization that accounted for the uncertainty

262     surrounding groupings due to standard deviations among samples and groups. In the case

263     of t-SNE, data derived from a principal component analysis (PCA) was used as input

264     variables, followed by clustering tests using the output from the UML algorithms. Both

265     approaches yielded more readily interpretable outcomes compared to other methods

266     assessed by the authors, revealing distinct species groupings in a two-dimensional space

267     (Derkarabetian et al., 2019). Notably, the identified groups in this study corresponded to

268     those of an integrative taxonomy approach considered by the researchers in their

269     comparisons, suggesting that the limits identified by UML algorithms might correspond to

270     species-level divergence rather than population structure (Derkarabetian et al., 2019).

271     Pyron et al. (2023) introduced a novel UML approach designed for delineating

272     species limits from extensive genomic datasets, primarily grounded in **self-organizing**

273     **maps (SOMs)**. This approach produces discrete outcomes rather than continuous ones,

274     grouping genotypes based on similarity, and is posited as more advantageous than prior

275     workflows. Additionally, the authors propose determining the number of species by

276     analyzing the degree of grid occupancy in the SOM output. This quantification establishes

how many units, representing distinct clusters of genotypes, have been effectively mapped from the original SNP matrix. Subsequently, the method estimates the cumulative distances from each sample to its immediate neighbors. To effectively separate these candidate species, Pyron et al. (2023) recommend performing cluster analyses, such as k-means. The determination of the optimal number of **classes** or species in the dataset is achieved by selecting the value that maximizes the sequential reduction in the weighted sum of squares from k to k + 1. Also, we highlight that this technique is rooted in the assessment of similarity rather than dissimilarity. An extension of this method has been proposed in the form of a SuperSOM approach, incorporating the possibility of using several trait classes simultaneously, such as alleles, morphological and ecological variables (Pyron, 2023).
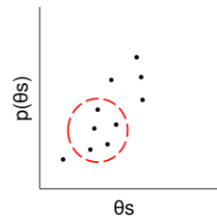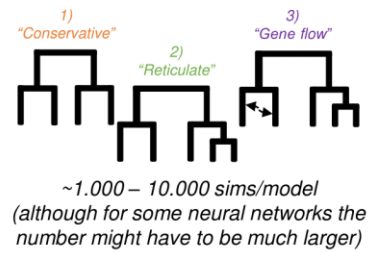
## 2.2. Validation and supervised methods

While UML approaches are powerful and widely applicable, there are situations where supervised machine learning (SML) will offer analytical advantages. Unlike UML, a workflow for applying any SML method to population genetic data generally include data simulation for various evolutionary scenarios, encoding both simulated and observed genetic data into **feature vectors**, **training** the algorithm, assessing its predictive performance through accuracy estimates, and applying it to new observed data points (Fig. 3). Thus, the use of simulated genetic data based on known evolutionary models is essential, given the scarcity of adequately sized datasets with high-confidence labels in Evolutionary Biology.
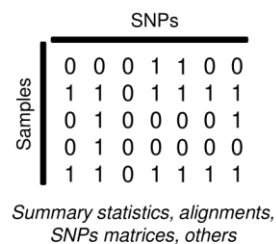
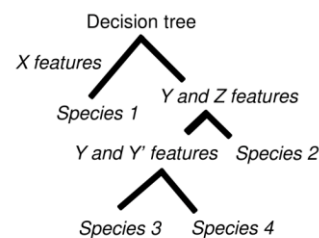**a)** Evolutionary models designing and prior distributions extraction

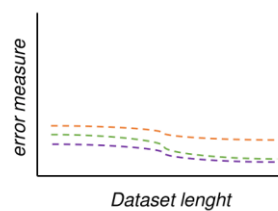**b)** Simulating data for each model and their respective prior distributions
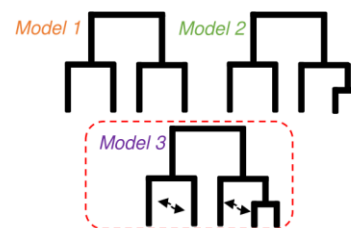


**c)** Choosing how to represent the biological data

**d)** Applying algorithm to the training set



**e)** Evaluating performance and optimizing parameters

**f)** Applying algorithm to the test set (empirical data), then choosing the best model



299
300 Fig. 3. Diagram illustrating a potential SML workflow for species delimitation (inspired by the work of Smith and
301 Carstens, 2020). a) The initial step involves designing priors for the evolutionary models considered in the
302 study. b) Simulated data is generated for each model, typically ranging from 1,000 to 10,000 simulations per
303 model, using relevant simulation software. c) The data is represented according to the requirements of the
304 chosen ML tool. d) Following data simulation and representation, ML model training begins, involving various
305 preliminary steps like data pre-processing, dataset division, feature selection, and algorithm choice. e) Model
306 performance (both in terms of biological accuracy and computationally) is assessed using statistical metrics,
307 allowing for retraining and adjustment based on the results. f) Once the model is adequately trained and
308 evaluated, it can be used to predict species categories for new data, which can be either newly simulated data
309 or empirical data consistent with the model's proposal, determining how many species exist in that particular
310 biological system.
311

312      The process of training and applying ML algorithms is influenced by the assumptions

313 of the underlying evolutionary processes, such as population size, selection strength, and

314 gene flow. Thus, the reliability of results obtained from SML methods rely on the

315 resemblance between the training data (typically simulated) and the actual biological data.

316 Anyhow, SML algorithms generally demand a significantly smaller amount of simulated data

compared to other methods based on simulations, such as **Approximate Bayesian Computation (ABC)**, resulting in reduced computational effort (e.g., a few thousand simulated datasets versus hundreds of thousands of simulations per scenario in most ABC approaches; Csilléry et al., 2010; Pudlo et al., 2016; Raynal et al., 2019).

CLADES (Pei et al., 2018), for example, is a SML approach designed for species delimitation, utilizing **classification models** trained and evaluated on *multilocus* sequence data. Notably, this study introduced the application of **support vector machines (SVM)** for species delimitation. For model training, datasets at the population level were simulated, with and without gene flow. Within this framework, species delimitation is framed as a classification task, where the goal is to classify pairs of populations as either belonging to the same or different species. Each training sample was represented as a list of summary statistics, and a SVM **regression** is calculated, through iterative training, to minimize the misclassification cost. Subsequently, the SVM classifier computed the probability of the training samples belonging to each potential grouping.

The training dataset was simulated based on a two-species model (A and B) where both species diverged at time $\tau$ with identical population size parameters ($\theta A = \theta B = \theta$). Each species further consisted of two populations that recently split at time $\tau_P$. **Migration** between species A and B was allowed at a rate of $M = Nm$ migrants per generation, with m representing the migration rate per generation. The MCcoal software (Rannala and Yang, 2003) was used to simulate multilocus sequence data of length L under various parameter combinations for training. For each possible parameter combination ($\theta$, $\tau$, $M$), sequences were simulated for 100 loci with a length of $L = 100Kbp$ for all populations. For each locus, 40 sequences were sampled, with 10 sequences per population. Additionally, symmetrical migration between species A and B was assumed before the populations of the species split at time $\tau_P$. All training samples were combined to train a global classifier, enabling it to adapt to various values of $\theta$ and $M$ instead of assuming fixed parameters. Longer loci improved

343 CLADES' efficiency, and it was robust to different modeling structures, accommodating

344 various demographic events and evolutionary parameters.

345 Smith and Carstens (2020) introduced delimitR, a SML approach designed to conduct

346 species delimitation in a model selection task; delimitR employs the multidimensional **site**

347 **frequency spectrum** (mSFS) with a **binning** strategy as a predictor variable for a **Random**

348 **Forest (RF)** classifier. Working with data summarized through the mSFS, delimitR facilitates

349 the evaluation of models that vary in terms of lineage numbers. In essence, this framework

350 aims to discriminate between various divergence models compatible with virtually any

351 species concept, as asserted by the authors. Given its supervised nature, delimitR demands

352 researchers to define reasonable priors, such as divergence times or migration rates, and

353 to make decisions about the inclusion of models within the set (Smith and Carstens, 2020).

354 For each model, Smith and Carstens (2020) simulated 10,000 mSFS. A RF classifier was

355 constructed using 1,000 **decision trees** to accommodate the extensive number of models.

356 delimitR's performance improved with larger SNP matrices and increasing divergence times.

357 Compared to ABC methods, delimitR showed lower error rates, even though the detection

358 of migration becomes challenging in cases of recent divergence between lineages (Smith

359 and Carstens, 2020). The authors acknowledge that further research is needed to elucidate

360 the association between the model space, number of parameters, and delimitation accuracy.

361

362 *2.3. Deep learning*

363 **Artificial neural networks (ANNs)** are increasingly employed in Evolutionary

364 Biology, often referred to as 'deep learning' (Sheehan and Song, 2016). Deep learning

365 techniques have found success in various fields in the Biological Sciences (Angermueller et

366 al., 2016; Sheehan and Song, 2016; Schrider and Kern, 2018). However, its adoption in

367 Evolutionary Biology is relatively recent (see Angermueller et al., 2016; Sheehan and Song,

368 2016; Blischak et al., 2021; Yelmen and Jay, 2023). The popularity of ANNs can be attributed

369  to their highly flexible data input and output structure, allowing networks trained for one task

370  to be repurposed for another by modifying their final **layers**, for instance, through **transfer**

371  **learning** approaches. This versatility enables the resolution of intricate tasks that might

372  prove challenging for **shallow learning** algorithms. Conversely, deep learning often

373  demands meticulous and more specific fine-tuning compared to shallow learning methods.

374      The fundamental stages involved in creating a supervised shallow learning

375  framework for species delimitation can be paralleled with the primary phases found in a deep

376  learning workflow. These encompass data simulation and representation, **model** training

377  and optimization, all the way to predicting the relevant categories from empirical data (Fig.

378  3). For a detailed description of how neural networks work, and their general structure, see

379  Sheehan and Song (2016), Borowiec et al. (2022), and Korfmann et al. (2023).

380

**a)** Simulate data under different evolutionary models

Pop. size

*Constant    Expansion   Bottleneck*

Topology

*Conservative        Splitter        Gene flow*

**b)** Convert the simulated data into image files

SNPs

Samples

0 0 0 1 0 0 0
1 1 0 1 0 1 1
0 1 0 0 0 0 1
0 1 0 0 1 0 0
1 1 0 1 0 1 1
1 1 0 1 1 0 1
0 0 0 0 1 1 0

*major freq.*

*minor freq.*

**c)** Train neural network with simulated data

SNPs

Samples

Convolutional layer

Pooling layer

Output models

ANN layers

Flatten layer

**d)** Predict the probability of each model from empirical data with the trained neural network
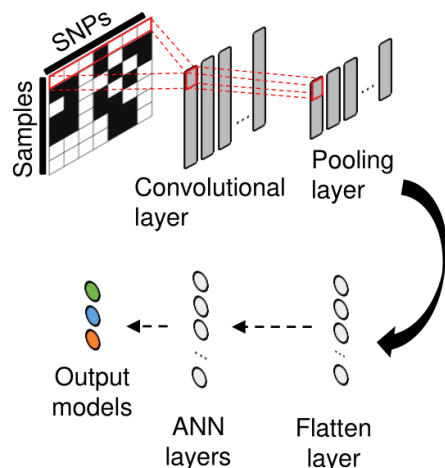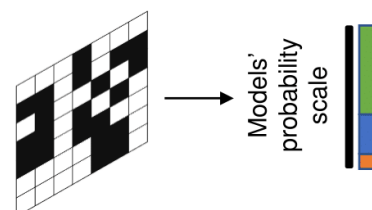
Models' probability scale

Fig. 4. Diagram illustrating a potential deep learning workflow applied in the context of species delimitation, using CNNs (inspired by Perez et al., 2021). a) The process typically begins with the simulation of biological data under various evolutionary models, considering factors like topology, population size, gene flow, and more, similar to SML. b) Next, data representation is crucial. For CNNs, SNP matrices are often converted into arrays or image files, where pixel contrast reflects differences in minor and major frequencies between samples. c) With the simulated and properly represented data, the network training phase can commence. The parameter configuration and network architecture may vary, depending on the specific study's requirements. d) Once each model is trained and its performance is rigorously evaluated, the final stage of the workflow involves predicting categories for new data. This can include using new simulated data with slight parametric modifications, still within the trained model's limits, as well as empirical data whose evolutionary history aligns with the proposed model. In both cases, the goal is to determine which delimitation model best applies to the biological system being investigated.

394    Perez et al. (2021) propose a species delimitation approach that accommodates the

395    integration of coalescence-based methods with model selection using **convolutional**

396    **neural networks (CNNs)**. Briefly, this approach can combine models from coalescent

397    analyses, such as using BPP (Flouri et al., 2018; 2020), allowing for the comparison of

398    different evolutionary scenarios. Thus, it allows for the test of species limits by integrating

399    data from various sources, including the possibility of incorporating knowledge from both

400    genetic analyses using coalescence-based methods and morphological hypotheses

401    reflecting diverse taxonomic arrangements. The initial steps involve simulating genetic data

402    for each delimitation hypothesis, with the study encompassing 10,000 simulations per

403    model, and transforming them into images. These images of simulated data are used to train

404    a neural network capable of recognizing simulations generated from each model. Then,

405    each species hypothesis probability can be predicted through CNNs using a **test set**. In the

406    same study, the authors conducted a comparison between their model selection approach

407    and ABC using empirical data. It is worth noting that while CNNs used 10,000 simulations

408    per model, ABC required 100,000 simulations per model. The CNNs consistently

409    demonstrated superior performance in distinguishing between the simulated demographic

410    scenarios, outperforming ABC in all cases, with fewer simulations and faster execution times

411    (Perez et al., 2021).

412

413    *2.4 How has machine learning changed our approach to delimit species so far?*

414    To date, relatively few studies (<20, also see Appendix B) have specifically explored

415    ML techniques for species delimitation, particularly when focusing on molecular data.

416    Among these, only five introduced novel ML approaches for species delimitation, providing

417    comprehensive details from initial simulations to statistical performance evaluations (Pei et

418    al., 2018; Derkarabetian et al., 2019; Smith and Carstens, 2020; Perez et al., 2021; Pyron

419    et al., 2023). These approaches, and also other ML frameworks applied in phylogeography

and demographic inferences, are often advocated by the researchers and developers themselves on the following arguments: i) challenges and limitations associated with the assumptions of coalescent methods (Derkarabetian et al., 2019; Smith and Carstens, 2020; Blischak et al., 2021; Martin et al., 2021; Derkarabetian et al., 2022); ii) ML computational efficiency and the capacity of handling complex evolutionary models (Pei et al., 2018; Martin et al., 2021; Perez et al., 2021; Derkarabetian et al., 2022; Pyron et al., 2023); and iii) ML acting as a likelihood-free approach, enabling the consideration of models where likelihood computation would be intractable (Smith and Carstens, 2020; Martin et al., 2021; Perez et al., 2021; Sanchez et al., 2020). While ML algorithms are often used similarly to simulation-based approaches like ABC, additional steps are generally incorporated, such as: i) selecting a more informative subset of summary statistics based on specific criteria (Smith and Carstens, 2020; Martin et al., 2021), and ii) handling larger or more complex genetic datasets compared to what Bayesian methods can do in a reasonable amount of time (Ghirotto et al., 2021; Smith and Carstens, 2020; Collin et al., 2021).

*2.5. What types of species ML methods might be detecting?*

A significant part of the studies we analyzed were philosophically based on species concepts grounded on evolutionary or genealogical independence criteria. This might stem from our focus on workflows using molecular data, which generally aims at identifying lineages and genetic clusters characterized by significant levels of genetic divergence and restricted amounts of gene flow. Also, some studies specifically model parameters like migration, which make them in line with concepts focused on reproductive criteria. While evolutionary and genealogical independence evidence (or reproductive criteria) may have their limitations in investigating species limits, results generated by ML methods in this context can still serve as hypotheses for further investigations (e.g., Fujita et al., 2012), aligning with the GLC perspective (de Queiroz, 1998; 1999; 2005b).

446       In this context, it is reasonable to assert that ML-based delimitation methods, just as

447    coalescence-based methods, might not always be identifying species *per se*, but rather: i)

448    incompletely separated (or incipient) species, which may eventually be classified as distinct

449    (Burbrink et al., 2021), or even as 'subspecies' (de Queiroz, 2020); or ii) population or

450    phylogeographic variation (Rosenblum et al., 2012; Sukumaran et al., 2021). Consequently,

451    while ML methods hold increasing promise for species boundaries inference, it is necessary

452    to evaluate the extent to which the ML methods could effectively discern evolutionary

453    independence among metapopulation lineages. So far, there are no definitive coalescent-

454    based solutions to differentiate between population structure and species (Sukumaran &

455    Knowles, 2017; Leaché et al., 2019). Thus, while model-based evolutionary lineage

456    structure detected through ML can be biologically relevant for species delimitation, additional

457    data and an evolutionary process-based perspective are crucial to discern the nature of the

458    inferred biological entities (Smith & Carstens, 2020; Sukumaran et al., 2021).

459       Inferring species limits from molecular data and integrating phenotypic data can be a

460    solution in some cases, but robust species delimitation still requires mechanistic hypotheses

461    about the speciation process itself (Padial & De la Riva, 2021; Pyron et al., 2023; Pyron et

462    al., 2024), because distinguishing between population structure and actively diverging or

463    collapsing species require explicit hypotheses and quantifiable tests (Sukumaran &

464    Knowles, 2017; Derkarabetian et al., 2019; Huang, 2020; Pyron et al., 2024). Just as

465    phenotypic, ecological, or other biological attributes are not mandatory criteria for

466    designating an evolutionary lineage as a species (de Queiroz, 2007; Pyron et al., 2023),

467    genetic or genealogical groupings identified using ML-based delimitation methods can be

468    similarly interpreted. Within this context, while the primary criterion for recognizing a species

469    can still be evolutionary independence, other characteristics may serve as secondary

470    evidence of divergence and could be also analyzed using ML frameworks.

471    Due to its great versatility in handling diverse data types, ML future applications to

472    infer species limits may also focus on evaluating which of the different biological properties

473    could be most effectively integrated into the species hypotheses testing process. They may

474    be useful in discerning between patterns of population structure and species-level

475    divergence, especially through the integration of distinct traits, such as genomic divergence,

476    gene flow, ecological adaptation, and phenotypic differentiation (Freedman et al., 2023;

477    Prates et al., 2023; Pyron et al., 2024). Again, this approach aligns with de Queiroz's GLC

478    (1998; 1999; 2005), providing a deeper understanding of the speciation processes through

479    multiple biological perspectives.

480    Only a few detailed ML pipelines have been proposed so far to explore the

481    relationships between evolutionary models and divergence scenarios in terms of distinct

482    characteristics, whether genetic, phenotypic, geographic or ecological. For example, Yang

483    et al. (2022) introduced a CNN method that successfully integrates morphological and

484    molecular data for species identification. Pyron (2023), on the other hand, implemented a

485    UML method using SOMs for learning high-dimensional associations between observations

486    (e.g., specimens) across a wide set of input features (e.g., genetics, geography,

487    environment, and phenotype). Future methodologies could further explore this integration of

488    multiple sources of information, both regarding species delimitation and integrative

489    taxonomy.

490

491    **3. Advantages, limitations and future perspectives**

492    *3.1. Strengths and benefits of using ML to delimit species*

493    In general, ML methods applied to infer species limits offer some advantages over

494    coalescent or traditional simulation-based methods. Despite particular constraints, ML

495    algorithms can perform as well as or even outperform (in terms of biological accuracy)

496    traditional model selection tools and likelihood-based species delimitation methods (Pei et

al., 2018; Smith and Carstens, 2020; Perez et al., 2021; Derkarabetian et al., 2022).

Moreover, being likelihood-free, they are computationally more efficient and generally can

be trained on models that are at times too intricate for formal statistical estimators (Pei et

al., 2018; Kuzenkov et al., 2020; Smith and Carstens, 2020; Suvorov et al., 2020; Martin et

al., 2021; Perez et al., 2021). Some of these algorithms have proven to be highly efficient in

complex evolutionary scenarios, including situations involving gene flow or population size

fluctuations (Pei et al., 2018; Perez et al., 2021). This efficiency does not compromise the

ability to distinguish between different models (Smith et al., 2017), and even simple SML

methods provide high selection accuracy when comparing multiple models in a single

analysis (Gehara et al., 2020 preprint).

Specifically, when it comes to deep learning, a major advantage is their capacity to

automatically extract information from alignments (commonly treated as images), as

opposed to relying on summary statistics typically required by other ML methods. This

facilitates accurate and efficient classification or regression tasks, as observed in studies by

Sanchez et al. (2020), Fonseca et al. (2021), Perez et al. (2021), and Borowiec et al. (2022),

thus holding promise in future species delimitation studies. Besides, especially in supervised

approaches, which often use explicit speciation models to validate species (e.g., Smith and

Carstens, 2020), ML enables a more in-depth exploration of the speciation and

phylogeographic processes that underlie the formation of independent evolutionary

lineages. Thus, given that properly sampled genomic datasets can offer sufficient data for

analyzing complex evolutionary models, ML might serve a dual role: providing primary

evidence for examining species limits patterns, and assisting in the investigation and

reconstruction of the evolutionary processes responsible for these patterns.

### 3.2. Constraints regarding ML and species delimitation

Certain algorithms, especially those in SML or deep learning, can be overly specialized. Modern ML methods are proficient at interpolating within the observed range of values in the training data, even in cases where specific values have not been encountered before, being adaptive and not solely reliant on memorizing specific training instances. Even so, because such algorithms are typically trained on simulated data with specific values of evolutionary parameters, such as θ and M, their performance might be compromised when applied far outside the training parameter space (Schrider and Kern, 2018; Borowiec et al., 2022). Besides, ML algorithms have some degree of **inductive bias** (Hüllermeier et al., 2013). Therefore, exploring in further details the association between training capacity and predictive power should be a priority for future studies.

Methods relying on a substantial volume of simulated data across diverse evolutionary scenarios need to consider the careful design of prior distributions to simulate models that closely resemble the real biological system under investigation. This challenge becomes more pronounced for non-model organisms, where data availability may severely limit the quality of parameter estimates (Tagu et al., 2014; Fonseca et al., 2016; Cerca et al., 2021; Jorna et al., 2021). Nonetheless, these simulation problems are not exclusive to ML-based workflows, as model selection frameworks such as ABC also employ simulated data (Beaumont et al., 2002; Bertorelle et al., 2010). All model-based methods depend on the specified models and its parameters, whether they are used for simulations or for direct likelihood estimation. Thus, traditional species delimitation methods that do not require simulations remain important alternatives for addressing delimitation challenges, in particular when there is no clear reference for simulations. Coalescence-based inferential methods are also limited in terms of their coverage of different evolutionary scenarios, but they possess optimality and iterability properties that span a reasonable portion of the parameter space, albeit at a considerable computational cost (e.g., Flouri et al., 2018;

549 Sukumaran et al., 2021). Nevertheless, methods not reliant on simulations can also be

550 sensitive to model misspecification, as the MSC deals with assumptions that may not be

551 appropriate for many biological systems.

552 Either wat, it may be unfeasible to simulate data or train an ML algorithm across an

553 entire parameter space, especially in complex evolutionary models (Rannala and Yang,

554 2020). Limited information is available regarding the asymptotic statistical performance of

555 most ML methods applied for species delimitation, and important phenomena may be

556 entirely missing from the simulations (e.g., background selection, Mo and Siepel (2023), or

557 missing data Arnab et al. (2023)). Thus, such models may never be comprehensive enough,

558 have limitations in representing real data, and demand substantial computational resources

559 (Arenas, 2012; Mangul et al., 2019a; Zaharias et al., 2022). This leads to an inherent

560 challenge in avoiding some degree of misspecification in the training data, even considering

561 the variety of powerful genetic data simulators currently available, such as SLiM (Messer,

562 2013), discoal (Kern and Schrider, 2016), msprime (Baumdicker et al., 2021), and

563 fastsimcoal2 (Excoffier et al., 2021).

564 Another crucial perspective to consider is that numerous studies, whether focusing

565 on species delimitation, demography, or population genetics, incorporate ML for inferences

566 based on summary statistics (Pei et al., 2018; Smith and Carstens, 2020; Collin et al., 2021;

567 Ghirotto et al., 2021). There are methodologies tailored for handling data derived from SNP

568 matrices (Derkarabetian et al., 2019; Sanchez et al., 2020; Smith and Carstens, 2020;

569 Blischak et al., 2021; Fonseca et al., 2021; Martin et al., 2021; Perez et al., 2021) or raw

570 sequence data (Pei et al., 2018; Ghirotto et al., 2021), and only a few pipelines offer

571 extensibility to various genetic markers (e.g., Collin et al., 2021). Notably, deep learning

572 techniques are valuable tools in this context, offering the capability to analyze both raw

573 genetic data and summary statistics (Korfmann et al., 2023).

574     While summary statistics can also be derived from the original genetic data and are

575     valuable for distinguishing between simulated models, not all of them may be suitable for

576     making inferences about species limits. The practical implementation of such statistics on

577     the detection of specific evolutionary processes often encounters confounding factors that

578     can mimic similar effects on gene histories (Flagel et al., 2019). For example, Tajima's D is

579     a statistic sensitive to both positive selection and changes in population size (Simonsen et

580     al., 1995). Moreover, since different studies often employ their specific set of summary

581     statistics, comparing the results of ML applications is not always straightforward, or feasible,

582     without acknowledging the significant nuances tied to the biological context considered in

583     each approach. Thus, the tendency of some ML algorithms to rely on specific

584     representations of data rather than the original dataset can be seen as a drawback in certain

585     scenarios. Unless we precisely know which type of data is truly sufficient to represent the

586     target data, an approach solely based on a particular set of summary statistics can inevitably

587     result in a degree of information loss (Rannala and Yang, 2020).

588     An alternative to learning from summary statistics is to consider the alignment itself

589     as input, as demonstrated in the CNNs approach introduced by Perez et al. (2021). Along

590     with other deep learning techniques, CNNs implicitly enable dimensionality reduction while

591     capturing structures within the input data. Thus, comparing different ML approaches might

592     be misleading due to the variability in the biological foundations employed in each workflow.

593     In other words, it is not always reasonable to strictly compare results produced by different

594     ML approaches, as they are generally trained on specific parameters and data

595     representation.

596

597     *3.3. Possible avenues and prospects for future studies*

598     Regarding ML, one approach to mitigate the effects of misspecification during

599     simulation involves designing or using a simulator that enforces greater compatibility

600 between simulated and actual data. Generative adversarial networks (GANs), a type of deep

601 learning algorithm commonly used for creating synthetic images and voices (Chadha et al.,

602 2021), have shown promise in this regard (see Callier, 2022; Wang et al., 2021). GANs

603 operate with two networks, the generator and the discriminator, trained together (Goodfellow

604 et al., 2014). While the generator generates simulated data, the discriminator distinguishes

605 between real and fake data. Over the course of training, the generator network becomes

606 more powerful at producing realistic **examples**, and the discriminator network becomes

607 more skilled at distinguishing between real and synthetic data. Once training is complete,

608 the generator network can be utilized to generate new examples that are indistinguishable

609 from real data, providing a reliable way to work with **labeled data**. Researchers have already

610 assessed the utility of GANs in various fields, including genomics, phylogenetics, and

611 population genetics (Booker et al., 2023; Nesterenko et al., 2022 preprint; Yelmen and Jay,

612 2023). Smith and Hahn (2023), for instance, introduced phyloGAN, a workflow that takes a

613 concatenated alignment (or a set of alignments) as input and infers a phylogenetic tree,

614 potentially accounting for gene tree heterogeneity.

615 While such approaches perform effectively in relatively straightforward scenarios,

616 challenges still emerge as the complexity of evolutionary model space increases. This

617 complexity might stem from more variables in evolutionary models or larger trees and

618 alignments, resulting in potential issues related to accuracy and execution time (Nesterenko

619 et al., 2022 preprint; Smith and Hahn, 2023). Consequently, it is important to recognize that

620 applications of GANs in the field of Evolutionary Biology are still in the early stages of

621 development. To fully harness the potential of this tool in species delimitation, further efforts

622 are required to refine estimates of genetic population parameters (e.g., Wang et al., 2021).

623 Future advancements in GANs within the realm of evolutionary biology should focus, for

624 instance, on enhancing the efficiency of exploring parameter spaces, reducing

625 computational training times, and accommodating more complex models (Smith and Hahn,

626 2023).

627 Besides, issues related to potential errors in data simulation can be likened to a

628 "domain adaptation" problem, where a model trained on one data distribution is applied to a

629 dataset originating from a different distribution (Farahani et al., 2021; Mo and Siepel, 2023).

630 A classic illustration of domain adaptation is found in image classification. Consider a

631 situation in which a recognition model needs to identify different dog breeds from

632 photographs ("target domain"), but there is an abundance of labeled training data available

633 only in cartoon drawings of dogs ("source domain"). In such cases, a ML model must be

634 trained on one dataset with the expectation of performing well on another, even in the

635 presence of systematic differences between the two distributions.

636 Recent approaches typically involve learning a "domain-invariant" data

637 representation through a feature extractor neural network. This is accomplished by

638 minimizing domain disparities (Rozantsev et al., 2018), utilizing adversarial networks (Ganin

639 and Lempitsky, 2015; Liu and Tuzel, 2016; Bousmalis et al., 2017), or employing auxiliary

640 reconstruction tasks (Ghifary et al., 2016). Domain adaptation techniques have found

641 applications in fields such as genomics (Cochran et al., 2022) and population genetics (Mo

642 and Siepel, 2023), particularly as an unsupervised domain adaptation problem. Through

643 extensive simulation studies, Mo and Siepel (2023) convincingly demonstrated that their

644 domain-adapted models significantly outperformed standard networks across various

645 simulation misspecification scenarios. This outcome underscores the potential of domain

646 adaptation techniques as a promising avenue for developing more robust deep learning

647 models in the realm of population genetic inference (Mo and Siepel, 2023), potentially

648 including species delimitation.

649 In addition to the limitations regarding simulations and training models in specific

650 parameter spaces, there is the issue associated with the manipulation of data attributes and

different types of input data. This becomes even more relevant considering that ML

techniques are lauded for their adaptability, especially considering transfer learning

frameworks. A neural network initially trained for a specific task can be repurposed for

different learning contexts with the simple modification of some of its layers, even though

reusing trained models can be very challenging due to differences in data dimensionality

(Sanchez et al., 2020). As an example, a deep learning **architecture** originally trained for

inferring historical population sizes can be repurposed for classifying demographic

scenarios (Pan and Yang, 2010). Also, deep learning methods used for phylogeographic

model selection (Fonseca et al., 2021) could be easily applied to species limits issues with

minimal adaptations.

## 4. Optimizing the use of ML in the context of species delimitation

*4.1. Enhancing Species Delimitation through accessible and purpose-built ML*

The introduction of new ML approaches will increasingly enhance researchers' ability

to make biologically precise decisions, especially when these methods are purpose-built,

from conception to implementation, for the specific task of delimiting evolutionary lineages.

A critical step in any species delimitation study is to select the appropriate methods to be

employed, considering the available data and putative evolutionary scenarios. With a

multitude of possibilities in the modern Evolutionary Biology toolkit, the ideal choice should

not only consider an appropriate fit with the biological problem under investigation, but also

a statistical evaluation and performance optimization (Greener et al., 2021; Morimoto et al.,

2021), under various diversification scenarios, while estimating historical parameters like

divergence time, population size, and migration rate. It is important to assess in which

specific evolutionary scenarios coalescent methods might exhibit strong limitations, and

whether a new ML workflow might outperform others in terms of performance. Thus, a

676    comprehensive analysis of the methods characteristics, advantages, disadvantages, and

677    overall performance compared to existing SDMs is desired.

678         Such evaluation should also encompass both the algorithm's biological predictions

679    and computational performance. Comparisons should be performed considering the

680    inherent properties of the used ML algorithms, such as how the workflows manipulate the

681    data attributes, and the different types of input and output data. In nearly all studies using

682    ML methods to infer species limits, at least a minimal approach to estimating error or noise

683    is typically employed (Pei et al., 2018; Smith & Carstens, 2020; Martin et al., 2021;

684    Derkarabetian et al., 2022). For example, it is common for researchers to evaluate the ML

685    model's performance using genetic datasets of varying sizes, or alignments of different

686    dimensions. Then, the quantity and quality of data clearly influence the effectiveness of ML

687    applications, as analyses conducted on larger, well-filtered datasets consistently yield better

688    results (Pei et al., 2018; Smith & Carstens, 2020; Martin et al., 2021; Derkarebetian, et al.,

689    2022). This effect is pronounced in UML approaches, as they tend to be more susceptible

690    to data-related issues (Martin et al., 2021).

691         From a practical perspective, evaluating the suitability of an ML tool for species

692    delimitation also involves assessing its accessibility, particularly when compared to

693    traditional SDMs. To promote the widespread adoption of ML tools in species delimitation,

694    it is crucial to ensure that analyses are accessible and reproducible. For example, a

695    thorough description of the ML method, but without a detailed reference to the dataset, can

696    lead to significant issues within the workflow (Chicco, 2017; Greener et al., 2021). The same

697    rationale extends to the availability of the trained models. For example, Derkarebetian et al.

698    (2022) assessed a ML approach's capability to delimit cryptic species, and constructed a

699    "customized" training dataset from a well-studied lineage with biological characteristics akin

700    to their focal taxon. In cases like these, where a specific ML classifier has been designed

701    and trained with a particular dataset based on a specific evolutionary model's parameters,

702    it is also important to ensure both the dataset and the classifier are meticulously described

703    and made accessible to the public. Such precautions minimize the need to construct entirely

704    new workflows for each study, involving tasks such as data simulation, model training, and

705    the selection of evaluation metrics, enabling researchers to evaluate and enhance the

706    method without needing to start from scratch (Greener et al., 2021; Heil et al., 2021).

707    Additionally, ML's ability to efficiently compare a wide range of models using large

708    datasets in less computational time provides an important advantage over traditional model

709    comparison approaches. Nonetheless, access to adequate computing resources remains a

710    challenge for many researchers in species delimitation and various scientific disciplines

711    (Veretnik et al., 2008; Truong et al., 2012; Helmy et al., 2016; Mangul et al., 2019b). Then,

712    efforts to provide resources like graphics processing units, cloud storage, and computational

713    clusters are all crucial steps toward making ML more accessible and inclusive for scientists

714    across diverse domains of knowledge, including species delimitation.

715

716    *4.2. Combining analytical frameworks to investigate complex delimitation models*

717    All models, while inherently limited in representing the underlying nature of species

718    diversification and, hence, of the current species limits among the tested entities, will be

719    more or less useful depending on their effectiveness in extracting relevant evolutionary

720    information from the available data. Accordingly, in some systems, certain methods should

721    be prioritized based on the processes driving divergence, and using multiple methods with

722    similar biases might not always enhance biological interpretability. For instance, Smith and

723    Cartens (2020) argue that traditional methods like BPP can accurately infer the number of

724    species but may overlook significant processes, such as secondary contact, something that

725    ML workflows like delimitR could be more efficient in dealing with. In this context, the choice

726    on which species delimitation method to use should be done before hypothesis-testing,

727 considering the nature of the available data, and possibly prior relevant biological

728 information regarding the evolution of the organisms.

729       One approach that would greatly benefit from the combination of coalescence-based

730 methods and ML algorithms, and that could shape the future direction of genetic-based

731 species delimitation, involves the empirical validation of speciation-based models, which can

732 provide a nuanced understanding of the speciation process. Different speciation-based

733 delimitation models, whether relying on ML, coalescence, or a combination of both, can be

734 employed to capture different facets of the evolutionary divergence process, and to test

735 different increasingly complex scenarios, with model validation serving as the means to

736 articulate expert knowledge and the available statistical tools for hypothesis testing. In sum,

737 while currently no universally superior species delimitation method exists, ML algorithms

738 offer promising prospects for their integration into systematic protocols tailored for species

739 delimitation.

740

741 **5. Conclusions**

742       Relatively few studies have explored ML techniques for species delimitation using

743 molecular data so far. They are generally employed due to coalescence-based methods'

744 specific assumptions and limitations. Besides, they are computationally efficient, can be

745 easily integrated with traditional methods, and clearly provides a concrete and robust way

746 to explore dataset structures when species-level divergences are hypothesized. The

747 flexibility of ML-based methods allows them to accommodate complex evolutionary

748 scenarios. Furthermore, likelihood-free approaches such as ML can provide more accurate

749 estimates of species limits and population parameters, particularly in cases where traditional

750 methods may struggle to converge or produce biased results.

751       Both ML approaches and coalescence-based methods provide a wide array of

752 choices, necessitating careful selection considering multiple factors. Particularly, ML

753 algorithms offer promising prospects but require thorough evaluation, comparison, and

754 adaptation to specific biological problems. Besides, selecting an appropriate ML method for

755 species delimitation should prioritize suitability for specific data and research questions over

756 popularity. This assessment includes biological predictions, computational performance,

757 and comparisons to existing methods, even considering that comparing existing methods

758 can be challenging.

759 Some specific challenges can be highlighted regarding the utilization of ML

760 frameworks to infer species limits. For example, overly specialized algorithms might perform

761 well within observed ranges of evolutionary parameters but can struggle outside the training

762 space. This gains importance as ML applications in Evolutionary Biology rely heavily on

763 simulated data. Besides, model specialization for simulated data can hinder generalizability

764 and transferability across different studies or data types. To address this issue, there are

765 some potential solutions and emerging approaches. For example, GANs enable the creation

766 of more realistic simulated data, and domain adaptation techniques to transfer knowledge

767 across datasets with systematic differences. Another challenge relies on handling data

768 derived from distinct genetic markers, hindering the comparison of different ML approaches.

769 Just as some coalescence-based methods, ML-based delimitation methods may not

770 always discern species, but might identify incompletely separated species or ephemeral

771 population variations. Therefore, ML should be progressively developed and used alongside

772 traditional methods to enhance objectivity and robustness in species delimitation processes,

773 combining the strengths of distinct analytical structures for hypothesis testing. This approach

774 may allow for the accurate estimation of the speciation process, facilitating a clearer

775 differentiation between population structure and evolutionary independence. Also, future

776 applications of ML methods in species delimitation may focus on integrating various

777 biological properties into species hypothesis testing. Finally, there is potential in utilizing ML

778 methods in Integrative Taxonomy approaches, as combining morphological, ecological, and

779 molecular data, is crucial for robust species delimitation and may benefit from the flexibility

780 of these AI-based approaches. As these conditions are increasingly met, ML is poised to

781 become an integral part of the toolkit used by scientists not only in the field of species

782 delimitation, but for various Evolutionary Biology applications worldwide.

783

784 **Declaration of Competing Interest**

785 The authors declare that they have no known competing financial interests or personal

786 relationships that could have appeared to influence the work reported in this paper.

787

794

795 **References**

796 References identified with an asterisk (*) are cited only within the Appendices.

797 Angermueller, C., Pärnamaa, T., Parts, L., & Stegle, O. 2016. Deep learning for computational

798    biology. Molecular Systems Biology 12.

799 Arenas, M. 2012. Simulation of molecular data under diverse evolutionary scenarios. PLoS computational

800    biology 8.

801 Arnab, S.P., Amin, M. R., & DeGiorgio, M. 2023. Uncovering footprints of natural selection through spectral

802    analysis of genomic summary statistics. Molecular Biology and Evolution, 40.

803 Arnold, M. L. 1992. Natural hybridization as an evolutionary process. Annual review of Ecology and

804    Systematics, *23*, 237–261.

805 Baumdicker, F., et al. 2021. Efficient ancestry and mutation simulation with msprime 1.0.  Genetics 220.

806    doi:10.1093/ genetics/iyab229

807    Beaumont, M.A., Zhang, W., Balding, D.J. 2002. Approximate Bayesian computation in population genetics.

808         Genetics 162, 2025–2035. doi:10.1093/genetics/162.4.2025

809    Bertorelle, G., Benazzo, A., Mona, S. 2010. ABC as a flexible framework to estimate demography over space

810         and time: some cons, many pros. Mol Ecol. 19, 2609–2625. doi:10.1111/j.1365-294X.2010.04690.x

811    Blischak, P.D., Barker, M.S., & Gutenkunst, R. N. 2021. Chromosome-scale inference of hybrid speciation

812         and admixture with convolutional neural networks. Molecular Ecology Resources 21, 2676–2688.

813         https://doi.org/10.1111/1755-0998.13355

814    Booker, W.W., Ray, D.D., & Schrider, D.R. 2023. This population does not exist: learning the distribution of

815         evolutionary histories with generative adversarial networks. Genetics, 224(2), iyad063.

816    Borowiec, M.L., Dikow, R.B., Frandsen, P.B., McKeeken, A., Valentini, G., & White, A.E. 2022. Deep learning

817         as a tool for ecology and evolution. Methods in Ecology and Evolution 13, 1640–1660.

818    Bortolus, A. 2008. Error cascades in the biological sciences: the unwanted consequences of using bad

819         taxonomy in ecology. AMBIO: A journal of the human environment 37, 114–118.

820    Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., & Krishnan, D. 2017. Unsupervised pixel-level domain

821         adaptation with generative adversarial networks. Proceedings of the IEEE conference on computer

822         vision and pattern recognition, 3722–3731.

823    Breitman, M.F., Domingos, F.M., Bagley, J.C., Wiederhecker, H.C., Ferrari, T.B., Cavalcante, V.H., ... & Colli,

824         G.R. 2018. A new species of *Enyalius* (Squamata, Leiosauridae) endemic to the Brazilian Cerrado.

825         Herpetologica 74, 355–369.

826    Burbrink, F.T., & Ruane, S. 2021. Contemporary philosophy and methods for studying speciation and

827         delimiting species. Ichthyology & Herpetology 109, 874–894.

828    Callier, V. 2022. Machine learning in evolutionary studies comes of age. Proceedings of the National

829         Academy of Sciences 119.

830    Carstens, B. C., Pelletier, T. A., Reid, N. M., & Satler, J. D. 2013. How to fail at species

831         delimitation. Molecular Ecology 22, 4369–4383.

832    Cerca, J., Maurstad, M.F., Rochette, N. C., Rivera-Colón, A.G., Rayamajhi, N., Catchen, J.M., & Struck, T H.

833         2021. Removing the bad apples: A simple bioinformatic method to improve loci-recovery in de novo

834         RADseq data for non-model organisms. Methods in Ecology and Evolution 12, 805–817.

835    Chadha, A., Kumar, V., Kashyap, S., & Gupta, M. 2021. Deepfake: An overview. In Proceedings of Second

836         International Conference on Computing, Communications, and Cyber-Security, pp. 557-566.

837         Springer, Singapore.

838 Chicco, D. 2017. Ten quick tips for machine learning in computational biology. BioData Mining 10, 1–17.

839 https://doi.org/10.1186/s13040-017-0155-3

840 Christin, S., Hervet, É., & Lecomte, N. 2019. Applications for deep learning in ecology. Methods in Ecology

841 and Evolution 10, 1632–1644.

842 Cochran, K., Srivastava, D., Shrikumar, A., Balsubramani, A., Hardison, R.C., Kundaje, A., Mahony, S. 2022.

843 Domain adaptive neural networks improve cross-species prediction of transcription factor binding.

844 Genome Res. 32, 512–523.

845 Collin, F.D., Durif, G., Raynal, L., Lombaert, E., Gautier, M., Vitalis, R., Marin, J.M., & Estoup, A. 2021.

846 Extending approximate Bayesian computation with supervised machine learning to infer

847 demographic history from genetic polymorphisms using DIYABC Random Forest. Molecular Ecology

848 Resources 21, 2598–2613. https://doi.org/10.1111/1755-0998.13413.

849 Csilléry, K., Blum, M.G., Gaggiotti, O.E., & François, O. 2010. Approximate Bayesian computation (ABC) in

850 practice. Trends in Ecology & Evolution 25, 410–418.

851 Dayrat, B. 2005. Towards integrative taxonomy. Biological journal of the Linnean society, 85, 407–417.

852 Degnan, J. H. & Rosenberg, N.A. 2009. Gene tree discordance, phylogenetic inference and the multispecies

853 coalescent. Trends Ecol. Evol. 24, 332–340.

854 Derkarabetian, S., Castillo, S., Koo, P.K., Ovchinnikov, S., & Hedin M. 2019. A demonstration of

855 unsupervised machine learning in species delimitation. Molecular Phylogenetics and Evolution 139.

856 https://doi.org/10.1016/j.ympev.2019.106562

857 Derkarabetian, S., Starrett, J., & Hedin, M. 2022. Using natural history to guide supervised machine learning

858 for cryptic species delimitation with genetic data. Frontiers in Zoology 19, 1–15.

859 Dike, H.U., Zhou, Y., Deveerasetty, K.K., & Wu, Q. 2018. Unsupervised learning based on artificial neural

860 network: A review. In 2018 IEEE International Conference on Cyborg and Bionic Systems (CBS), pp.

861 322-327.

862 Domingos, F.M., Bosque, R.J., Cassimiro, J., Colli, G.R., Rodrigues, M.T., Santos, M.G., & Beheregaray, L.

863 B. 2014. Out of the deep: cryptic speciation in a Neotropical gecko (Squamata, Phyllodactylidae)

864 revealed by species delimitation methods. Molecular Phylogenetics and Evolution 80, 113–124.

865 *Duan, L., Han, L.N., Liu, B., Leostrin, A., Harris, A.J., Wang, L., Arslan, E., Ertuğrul, K., Knyazev, M.,

866 Hantemirova, E., Wen, J., & Chen, H.F. 2023. Species delimitation of the liquorice tribe

867 (Leguminosae: Glycyrrhizeae) based on phylogenomic and machine learning analyses. Journal of

868 Systematics and Evolution 61, 22–41. https://doi.org/10.1111/jse.12902.

869    Edwards, S. V. 2009. Is a new and general theory of molecular systematics emerging? Evolution 63, 1–19.

870    Edwards, S.V., Xi, Z., Janke, A., Faircloth, B.C., McCormack, J.E., Glenn, T.C., ... & Davis, C.C. 2016.

871        Implementing and testing the multispecies coalescent model: a valuable paradigm for

872        phylogenomics. Molecular Phylogenetics and Evolution 94, 447–462.

873    Edwards, S. V., Potter, S., Schmitt, C. J., Bragg, J. G., & Moritz, C. 2016. Reticulation, divergence, and the

874        phylogeography–phylogenetics continuum. Proceedings of the National Academy of Sciences, 113,

875        8025–8032.

876    Ely, C.V., de Loreto Bordignon, S.A., Trevisan, R., & Boldrini, I.I. 2017. Implications of poor taxonomy in

877        conservation. Journal for Nature Conservation 36, 10–13.

878    Excoffier, L. et al. 2021. fastsimcoal2: demographic inference under complex evolutionary   scenarios.

879        Bioinformatics 37, 4882–4885. doi:10.1093/bioinformatics/btab468.

880    *Fan, X.K., Wu, J., Comes, H.P., Feng, Y., Wang, T., Yang, S.Z., Iwasaki, T., Zhu, H., Jiang, Y., Lee, J., & Li,

881        P. 2023. Phylogenomic, morphological, and niche differentiation analyses unveil species delimitation

882        and evolutionary history of endangered maples in Acer series Campestria (Sapindaceae). Journal of

883        Systematics and Evolution 61, 284–298. https://doi.org/10.1111/jse.12919.

884    Farahani, A., Voghoei, S., Rasheed, K., & Arabnia, H.R. 2021. A brief review of domain adaptation.

885        Advances in data science and information engineering: proceedings from ICDATA 2020 and IKE

886        2020, 877–894.

887    Flagel, L., Brandvain, Y., & Schrider, D.R. 2019. The unreasonable effectiveness of convolutional neural

888        networks in population genetic inference. Molecular Biology and Evolution 36, 220–238.

889    Flouri, T., Jiao, X., Rannala, B., Yang, Z. 2018. Species Tree Inference with BPP using Genomic Sequences

890        and the Multispecies Coalescent. Molecular Biology and Evolution 35, 2585–2593.

891        doi:10.1093/molbev/msy147.

892        2020. A Bayesian implementation of the multispecies coalescent model with introgression for

893        phylogenomic analysis. Molecular Biology and Evolution 37, 1211–1223.

894    Fonseca, R.R. et al. 2016. Next-generation biology: sequencing and data analysis approaches

895        for non-model organisms. Marine genomics 30, 3–13.

896    Fonseca, E. M., Colli, G. R., Werneck, F. P., & Carstens, B. C. 2021. Phylogeographic model selection using

897        convolutional neural networks. Molecular Ecology Resources 21, 2661–2675.

898        https://doi.org/10.1111/1755-0998.13427.

899    Fonseca, E. M., & Carstens, B. C. (2024). Artificial intelligence enables unified analysis of historical and

900          landscape influences on genetic diversity. *Molecular Phylogenetics and Evolution*, 108116.

901    Fountain-Jones, N.M., Smith, M.L., & Austerlitz, F. 2021. Machine learning in molecular ecology. Molecular

902          Ecology Resources 21, 2589–2597. https://doi.org/10.1111/1755-0998.13532.

903    Fujita, M.K., Leaché, A.D., Burbrink, F.T., McGuire, J.A., & Moritz, C. 2012. Coalescent-based species

904          delimitation in an integrative taxonomy. Trends in Ecology & Evolution 27, 480–488.

905    Ganin, Y., & Lempitsky, V. 2015. Unsupervised domain adaptation by backpropagation. In International

906          conference on machine learning, 1180–1189.

907    Gehara, M., Mazzochinni, G.G., & Burbrink, F. 2020. PipeMaster: inferring population divergence and

908          demographic history with approximate Bayesian computation and supervised machine-learning in

909          R. BioRxiv, 2020-12. https://doi.org/10.1101/2020.12.04.410670

910    Ghifary, M, Kleijn, W.B., Zhang, M., Balduzzi, D., Li, W. 2016. Deep Reconstruction Classification Networks

911          for Unsupervised Domain Adaptation. In: Leibe B, Matas J, Sebe N, Welling M, editors. Computer

912          Vision ECCV 2016. Lecture Notes in Computer Science. Cham: Springer International Publishing. p.

913          597

914    Ghirotto, S., Vizzari, M.T., Tassi, F., Barbujani, G. & Benazzo, A. 2021. Distinguishing among complex

915          evolutionary models using unphased whole-genome data through random forest approximate

916          Bayesian computation. Molecular Ecology Resources 21, 2614–2628. https://doi.org/10.1111/1755-

917          0998.13263.

918    Gompert, Z., Mandeville, E. G., & Buerkle, C. A. 2017. Analysis of population genomic data from hybrid

919          zones. Annual Review of Ecology, Evol., and Syst., 48, 207–229.

920    Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio Y.

921          2014. Generative adversarial nets. Advances in Neural Information Processing Systems, 2672–

922          2680.

923    Greener, J.G., Kandathil, S.M., Moffat, L., & Jones, D.T. 2021. A guide to machine learning for biologists.

924          Molecular Cell Biology 23, 40–55. https://doi.org/10.1038/s41580-021-00407-0.

925    Hastie, T., Tibshirani, R., & Friedman, J. 2009. Unsupervised learning. In The elements of statistical

926          learning (pp. 485-585). Springer, New York, NY.

927    Hausdorf, B. 2011. Progress toward a general species concept. Evolution 65, 923–931.

928    Heil, B.J., Hoffman, M. M., Markowetz, F., Lee, S.I., Greene, C.S. & Hicks, S.C. 2021. Reproducibility

929          standards for machine learning in the life sciences. Nature Methods 18, 1132–1135.

930 Helmy, M., Awad, M., & Mosa, K.A. 2016. Limited resources of genome sequencing in developing countries:

931    challenges and solutions. Applied & translational genomics 9, 15–19.

932 *Hodel, R.G., Winslow, S.K., Liu, B.B., Johnson, G., Trizna, M., White, A.E., ... & Wen, J. 2023. A

933    phylogenomic approach, combined with morphological characters gleaned via machine learning,

934    uncovers the hybrid origin and biogeographic diversification of the plum genus. bioRxiv, 2023-09.

935    https://doi.org/10.1101/2023.09.13.557598

936 Hüllermeier, E., Fober, T. & Mernberger, M. 2013. Inductive bias. Encyclopedia of systems biology, 1018–

937    1019.

938 Huang, J. P. 2020. Is population subdivision different from speciation? From phylogeography to species

939    delimitation. Ecology and Evolution 10, 6890–6896.

940 Huelsenbeck, J.P., Andolfatto, P., Huelsenbeck, E.T. 2011. Structurama: Bayesian inference of population

941    structure. Evolutionary Bioinformatics 7, 55–59.

942 Jackson, N.D., Carstens, B.C., Morales, A.E. & O'Meara B.C. 2017. Species delimitation with gene

943    flow. Systematic Biology 66, 799–812.

944 Jackson, N.D., Morales, A.E., Carstens, B.C. & O'Meara B.C. 2017. PHRAPL: phylogeographic inference

945    using approximate likelihoods. Systematic Biology 66, 1045–1053.

946 Jacobs, S. J., Kristofferson, C., Uribe-Convers, S., Latvis, M., & Tank, D. C. (2018). Incongruence in

947    molecular species delimitation schemes: What to do when adding more data is difficult. Molecular

948    Ecology 27, 2397–2413.

949 *Jamdade, R., Al-Shaer, K., Al-Sallani, M., Al-Harthi, E., Mahmoud, T., Gairola, S., & Shabana, H.A. 2022.

950    Multilocus marker-based delimitation of *Salicornia persica* and its population discrimination assisted

951    by supervised machine learning approach. PLoS ONE 17.

952    https://doi.org/10.1371/journal.pone.0270463.

953 Jarvis, E.D., Mirarab, S., Aberer, A.J., Li, B., Houde, P., Li, C., ... & Zhang, G. 2014. Whole-genome

954    analyses resolve early branches in the tree of life of modern birds. Science 346, 1320–1331.

955 Jörger, K.M., & Schrödl, M. 2013. How to describe a cryptic species? Practical challenges of molecular

956    taxonomy. Frontiers in Zoology 10, 1–27.

957 Jorna, J. et al. 2021. Species boundaries in the messy middle—A genome-scale validation of species

958    delimitation in a recently diverged lineage of coastal fog desert lichen fungi. Ecology and Evolution

959    11, 18615-18632.

960     Karbstein, K., Kösters, L., Hodač, L., Hofmann, M., Hörandl, E., Tomasello, S., ... & Wäldchen, J. (2023).

961        Species delimitation 4.0: integrative taxonomy meets artificial intelligence. *Trends in Ecology &*

962        *Evolution*.

963     *Khalighifar, A., Brown, R.M., Goyes Vallejos, J., & Peterson, A.T. 2021. Deep learning improves acoustic

964        biodiversity monitoring and new candidate forest frog species identification (genus *Platymantis*) in

965        the Philippines. Biodiversity and Conservation, 30, 643-657.

966     Knowles, L. L., & Carstens, B. C. 2007. Delimiting species without monophyletic gene trees. Syst. Biol., 56,

967        887–895.

968     Korfmann, K., Gaggiotti, O.E. & Fumagalli, M. 2023. Deep learning in population genetics. Genome Biology

969        and Evolution. https://doi.org/10.1093/gbe/evad008.

970     Kuzenkov, O., Morozov, A., & Kuzenkova, G. 2020. Exploring evolutionary fitness in biological systems using

971        machine learning methods. Entropy 23, 1–35.

972     Leaché, A.D., Harris, R.B., Rannala, B. & Yang, Z. 2014. The influence of gene flow on species tree

973        estimation: a simulation study. Systematic Biology 63, 17–30.

974     Leaché, A.D., Zhu, T., Rannala, B., & Yang, Z. 2019. The spectre of too many species. Systematic Biology

975        68, 168–181.

976     Libbrecht, M.W. & Noble, W.S. 2015. Machine learning applications in genetics and genomics. Nature

977        Reviews Genetics 16, 32–332.

978     *Lima, A.P. et al. 2020a. Not as widespread as thought: Integrative taxonomy reveals cryptic diversity in the

979        Amazonian nurse frog *Allobates tinae* Melo-Sampaio, Oliveira and Prates, 2018 and description of a

980        new species. Journal of Zoological Systematics and Evolutionary Research, 58(4), 1173–1194.

981     *Lima, L.R. et al. 2020b. Below the waterline: cryptic diversity of aquatic pipid frogs (*Pipa carvalhoi*) unveiled

982        through an integrative taxonomy approach. Systematics and Biodiversity, 18(8), 771–783.

983     Liu, B. 2019. BioSeq-Analysis: a platform for DNA, RNA and protein sequence analysis based on machine

984        learning approaches. Briefings in bioinformatics 20, 1280–1294.

985     Liu, M.Y. & Tuzel, O. 2016. Coupled Generative Adversarial Networks. In: Advances in Neural Information

986        Processing Systems 29. Curran Associates, Inc.

987        https://papers.nips.cc/paper/2016/hash/502e4a16930e414107ee22b6198c578f-Abstract.html.

988     Lukhtanov, V.A. 2019. Species Delimitation and Analysis of Cryptic Species Diversity in the XXI Century.

989        Entmol. Rev. 99, 463–472. https://doi.org/10.1134/S0013873819040055.

990 Lürig, M.D., Donoughe, S., Svensson, E.I., Porto, A. & Tsuboi, M. 2021. Computer vision, machine learning,

991     and the promise of phenomics in ecology and evolutionary biology. Frontiers in Ecology and

992     Evolution 9.

993 *Magalhães, F.D.M., Lyra, M.L., De Carvalho, T.R., Baldo, D., Brusquetti, F., Burella, P., ... & Garda, A.A.

994     2020. Taxonomic review of South American Butter Frogs: Phylogeny, geographic patterns, and

995     species delimitation in the *Leptodactylus latrans* species group (Anura:

996     Leptodactylidae). Herpetological Monographs, 34(1), 131–177.

997 Mallet, J., Besansky, N., & Hahn, M. W. 2016. How reticulated are species? BioEssays, 38, 140–149.

998 Mangul, S. et al. 2019a. Systematic benchmarking of omics computational tools. Nature communications 10.

999         2019b. How bioinformatics and open data can boost basic science in countries

1000        and universities with limited resources. Nature biotechnology 37, 324–326.

1001 Martin, B.T., Chafin, T.K., Douglas, M.R., Placyk, Jr J.S., Birkhead, R.D., Phillips, C.A., & Douglas, M.E.

1002        2021. The choices we make and the impacts they have: Machine learning and species delimitation in

1003        North American box turtles (*Terrapene* spp.). Molecular Ecology Resources 21, 2801–2817.

1004 Mayr, E. M. 1969. The biological meaning of species. Biological Journal of the Linnean

1005        society, 1, 311–320.

1006            1996. What is a species, and what is not? Philosophy of science, 63, 262–277.

1007            2000. The biological species concept. Species concepts and phylogenetic theory: a debate,

1008            17–29.

1009 McClure, E.C., Sievers, M., Brown, C.J. Buelow, C.A., Ditria, E.M., Hayes, M.A., ... & Connolly R.M. 2020.

1010        Artificial intelligence meets citizen science to supercharge ecological monitoring. Patterns 1.

1011 Messer, P. W. 2013. SLiM: simulating evolution with selection and link-age. Genetics 194, 1037–1039.

1012        doi:10.1534/genetics.113. 152181.

1013 Mitchell, T.M. 1997. Machine Learning. McGraw-Hill, New York.

1014 Mo, Z., & Siepel, A. 2023. Domain-adaptive neural networks improve supervised machine learning based on

1015        simulated population genetic data. PLOS Genetics, 19.

1016 Mo, Y. K., Hahn, M. W., & Smith, M. L. 2024. Applications of machine learning in phylogenetics. *Molecular

1017        Phylogenetics and Evolution*, *196*, 108066.

1018 Morimoto, J., Ponchon, A., Sofronov, G., & Travis, J. 2021. Editorial: Applications of Machine Learning to

1019        Evolutionary Ecology Data. Frontiers in Ecology and Evolution.

1020  Nesterenko, L., Boussau, B., & Jacob, L. 2022. Phyloformer: towards fast and accurate phylogeny estimation

1021      with self-attention networks. bioRxiv, 2022-06. https://doi.org/10.1101/2022.06.24.496975

1022  *Newton, L.G., Starrett, J., Hendrixson, B.E., Derkarabetian, S., & Bond, J.E. (2020).

1023      Integrative species delimitation reveals cryptic diversity in the southern Appalachian *Antrodiaetus*

1024      *unicolor* (Araneae: Antrodiaetidae) species complex. Molecular Ecology 29, 2269–2287.

1025  O'Meara B. C. 2010. New heuristic methods for joint species delimitation and species tree inference.

1026      Systematic Biology 59, 59–73.

1027      2012. Evolutionary inferences from phylogenies: a review of methods. Annual

1028      Review of Ecology, Evolution, and Systematics 43, 267–285.

1029  Padial, J. M., Miralles, A., De la Riva, I., & Vences, M. 2010. The integrative future of

1030      taxonomy. Frontiers in zoology, 7, 1–14.

1031  Pan, S. J. & Yang, Q. 2010. A survey on transfer learning. IEEE Transactions on Knowledge and Data

1032      Engineering 22, 1345–1359.

1033  Pante, E., Puillandre, N., Viricel, A., Arnaud-Haond, S., Aurelle, D., Castelin M., ... & Samadi, S. 2015.

1034      Species are hypotheses: avoid connectivity assessments based on pillars of sand. Molecular

1035      Ecology 24, 525–544.

1036  Pei, J., Chu, C., Li, X., Lu, B., & Wu, Y. 2018. CLADES: A classification-based machine learning method for

1037      species delimitation from population genetic data. Molecular Ecology Resources 18, 1144–1156.

1038      https://doi.org/10.1111/1755-0998.12887.

1039  Perez, M.F., Bonatelli, I.A.S., Romeiro-Brito, M., Franco, F.F., Taylor, N.P., Zappi, D.C. et al. 2021.

1040      Coalescent-based species delimitation meets deep learning: Insights from a highly fragmented

1041      cactus system. Molecular Ecology Resources.

1042  Pichler, M., Boreux, V., Klein, A. M., Schleuning, M. & Hartig F. 2020. Machine learning algorithms to infer

1043      trait-matching and predict species interactions in ecological networks. Methods in Ecology and

1044      Evolution 11, 281–293.

1045  Pigliucci, M. 2003. Species as family resemblance concepts: the (dis-)solution of the species problem?

1046      BioEssays, 25, 596–602.

1047  Pons, J., Barraclough, T.G., Gomez-Zurita, J. et al. 2006. Sequence-based species delimitation for the DNA

1048      taxonomy of unde-scribed insects. Systematic Biology 55, 595–609.

1049     Price, T.D., Qvarnström, A., & Irwin, D.E. 2003. The role of phenotypic plasticity in driving genetic

1050          evolution. Proceedings of the Royal Society of London. Series B: Biological Sciences 270, 1433–

1051          1440.

1052     *Pritchard, J.K., Stephens, M., Donnelly, P. 2000. Inference of population structure using multilocus

1053          genotype data. Genetics 155, 945–959.

1054     Pudlo, P., Marin, J.M., Estoup, A., Cornuet, J.M., Gautier, M., & Robert, C.P. 2016. Reliable ABC model

1055          choice via random forests. Bioinformatics 32, 859–866. https://doi.org/10.1093/bioinformatics/btv684.

1056     Pyron, R.A. 2023. Unsupervised Machine Learning for Species Delimitation, Integrative Taxonomy, and

1057          Biodiversity Conservation. Molecular Phylogenetics and Evolution, 189.

1058     Pyron, R.A., O'Connell, K.A., Duncan, S.C., Burbrink, F.T., & Beamer, D.A. 2023. Speciation hypotheses

1059          from phylogeographic delimitation yield an integrative taxonomy for Seal Salamanders

1060          (*Desmognathus monticola*). Systematic Biology, 72, 179–197.

1061     Pyron, R. A., Kakkera, A., Beamer, D. A., & O'Connell, K. A. 2024. Discerning structure versus speciation in

1062          phylogeographic analysis of Seepage Salamanders (*Desmognathus aeneus*) using demography,

1063          environment, geography, and phenotype. Molecular Ecology, 33, e17219.

1064     de Queiroz, K. 1998. The general lineage concept of species, species criteria, and the process

1065          of speciation. Endless forms: species and speciation.

1066               1999. The General Lineage Concept of Species and the Defining Properties of

1067                  the Species Category. In book: Species: New Interdisciplinary Essays, Chapter: 3,

1068                  Publisher: MIT Press, Editors: Robert A. Wilson.

1069               2005a. Ernst Mayr and the modern concept of species. Proceedings of the National

1070                  Academy of Sciences, 102, 6600–6607.

1071               2005b. Different species problems and their resolution. BioEssays 27,

1072                  1263–1269.

1073               2007. Species concepts and species delimitation. Syst. Biol. 56, 879–886.

1074               2011. Branches in the lines of descent: Charles Darwin and the evolution of the species

1075                  concept. Biol. J. Linn. Soc. 103, 19–35.

1076               2020. An updated concept of subspecies resolves a dispute about the taxonomy of

1077                  incompletely separated lineages. Herpetological Review.

1078     Qu, K., Guo, F., Liu, X., Lin, Y., & Zou, Q. 2019. Application of machine learning in microbiology. Frontiers in

1079          Microbiology 10.

1080    Rannala, B. 2015. The art and science of species delimitation. Current Zoology 61, 846–853.

1081    Rannala, B., & Yang, Z. 2003. Bayes estimation of species divergence times and ancestral population sizes

1082        using DNA sequences from multiple loci. Genetics 164, 1645–1656.

1083            2010. Bayesian species delimitation using multilocus sequence data. Proceedings of the

1084            National Academy of Sciences 107, 9264–9269.

1085            2020. Species Delimitation. In: Phylogenetics in the genomic era.

1086    Rannala, B., Edwards, S.V., Leaché, A., & Yang, Z. 2020. The Multi-species Coalescent Model and Species

1087        Tree Inference. Scornavacca, Celine; Delsuc, Frédéric; Galtier, Nicolas. Phylogenetics in the

1088        Genomic Era, No commercial publisher | Authors open access book.

1089    Raynal, L., Marin, J.M., Pudlo, P., Ribatet, M., Robert, C.P., & Estoup, A. 2019. ABC random forests for

1090        Bayesian parameter inference. Bioinformatics 35, 1720–1728.

1091    Rozantsev, A., Salzmann, M. & Fua, P. 2018. Beyond sharing weights for deep domain adaptation. IEEE

1092        transactions on pattern analysis and machine intelligence 41, 801–814.

1093    *Saryan, P., Gupta, S. & Gowda, V. 2020. Species complex delimitations in the genus Hedychium: A

1094        machine learning approach for cluster discovery. Applications in Plant Sciences 8.

1095        https://doi.org/10.1002/aps3.11377.

1096    Sanchez, T., Cury, J., Charpiat, G. & Jay, F. 2020. Deep learning for population size history inference:

1097        Design, comparison and combination with approximate Bayesian computation. Molecular Ecology

1098        Resources 21, 2645–2660.

1099    Scalon, M.C., Domingos, F. M.C.B., Cruz, W.J.A., Marimon-Júnior, B. H., Marimon, B.S., & Oliveras, I. 2020.

1100        Diversity of functional trade-offs enhances survival after fire in Neotropical savanna species. Journal

1101        of Vegetation Science, 31, 139-150.

1102    Schrider, D.R. & Kern, A.D. 2016. Discoal: flexible coalescent simulations with selection. Bioinformatics 32,

1103        3839–3841.  doi:10.1093/ bioinformatics/btw556.

1104            2018. Supervised Machine Learning for Population Genetics: A New Paradigm. Trends in

1105            Genetics 34, 301–312. https://doi.org/10.1016/j.tig.2017.12.005

1106    Searls, D.B. 2010. The Roots of Bioinformatics. PLoS Comput Biol 6.

1107        https://doi.org/10.1371/journal.pcbi.1000809.

1108    Sheehan, S., & Song, Y.S. 2016. Deep learning for population genetic inference. PLoS computational

1109        biology 12.

1110    Shurtliff, Q. R. 2013. Mammalian hybrid zones: a review. Mammal Review, 43, 1–21.

Sidey-Gibbons, J.A., & Sidey-Gibbons, C.J. 2019. Machine learning in medicine: a practical introduction. BMC medical research methodology 19, 1–18.

Silva, D.C., Oliveira, H.F.M., & Domingos, F.M.C.B. 2024. Cerrado bat community assembly is determined by both present-day and historical factors. Journal of Biogeography.

Simonsen, K.L., Churchill, G.A., Aquadro, C.F. 1995. Properties of statistical tests of neutrality for DNA polymorphism data. Genetics 1411, 413–429.

Sites, Jr J.W. & Marshall, J.C. 2004. Operational criteria for delimiting species. Annual Review of Ecology, Evolution, and Systematics, 199-227.

Slatko, B.E., Gardner, A.F. & Ausubel, F.M. 2018. Overview of next-generation sequencing technologies. Current protocols in molecular biology 122.

Smith, M.L., Ruffley, M., Espíndola, A., Tank, D.C., Sullivan, J. & Carstens, B.C. 2017. Demographic Model Selection using Random Forests and the Site Frequency Spectrum. Molecular Ecology.

Smith, M.L. & Carstens B.C. 2020. Process-based species delimitation leads to identification of more biologically relevant species. Evolution 74, 216–229. https://doi.org/10.1111/evo.13878.

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

Solis-Lemus, C., Yang, S., & Zepeda-Nunez, L. 2022. Accurate phylogenetic inference with a symmetry-preserving neural network model. arXiv preprint arXiv:2201.04663.

Sukumaran, J. & Knowles, L.L. 2017. Multispecies coalescent delimits structure, not species. Proceedings of the National Academy of Sciences 114, 1607–1612.

Sukumaran, J., Holder, M.T., & Knowles, L.L. 2021. Incorporating the speciation process into species delimitation. PLoS Computational Biology 17.

Suvorov, A., Hochuli, J. & Schrider, D.R. 2020. Accurate inference of tree topologies from multiple sequence alignments using deep learning. Systematic biology 69, 221–233.

Tagu, D., Colbourne, J.K. & Nègre, N. 2014. Genomic data integration for ecological and evolutionary traits in non-model organisms. BMC genomics 15, 1–16.

Tautz, D., Arctander, P., Minelli, A., Thomas, R.H., Vogler, A.P. 2003. A plea for DNA taxonomy. Trends Ecol. Evol. 18, 70–74.

Truong, H.L., Pham T.V., Thoai N., & Dustdar S. 2012. Cloud computing for education and research in developing countries. In Cloud computing for teaching and learning: strategies for design and implementation, pp. 64–80. IGI Global.

1142    Valletta, J.J., Torney, C., Kings, M., Thornton, A. & Madden J. 2017. Applications of machine learning in

1143           animal behaviour studies. Animal Behaviour 124, 203–220.

1144    Veretnik, S., Fink, J.L. & Bourne, P.E. 2008. Computational biology resources lack persistence and usability.

1145           PLoS computational biology 4.

1146    Vink, C.J., Paquin, P., & Cruickshank, R.H. 2012. Taxonomy and irreproducible biological science.

1147           BioScience 62, 451–452.

1148    Vogler, A.P., Monaghan, M.T. 2007. Recent advances in DNA taxonomy. J. Zool. Syst. Evol. Res. 45, 1–10.

1149    Wake, D.B., Wake, M.H. & Specht C.D. 2011. Homoplasy: from detecting pattern to determining process and

1150           mechanism of evolution. Science 331, 1032–1035.

1151    Wäldchen, J. & Mäder, P. 2018. Machine learning for image-based species identification. Methods in

1152           Ecology and Evolution 9, 2216–2225.

1153    Wang, G. 2019. Machine learning for inferring animal behavior from location and movement data. Ecological

1154           informatics 49, 69–76.

1155    Wang, Z., Wang, J., Kourakos, M., Hoang, N., Lee, H.H., Mathieson, I., & Mathieson, S. 2021. Automatic

1156           inference of demographic parameters using generative adversarial networks. Molecular ecology

1157           resources 21, 2689–2705.

1158    Wiens, J. J., & Penkrot, T. A. 2002. Delimiting species using DNA and morphological variation and

1159           discordant species limits in spiny lizards (*Sceloporus*). Syst. Biol., 51, 69–91.

1160    Wiens, J. J. 2007. Species delimitation: new approaches for discovering diversity. Syst. Biol. 56, 875–8.

1161    Wilkins, J. S., Zachos, F. E., & Pavlinov, I. Y. (Eds.). 2022. Species Problems and Beyond: Contemporary

1162           Issues in Philosophy and Practice. CRC Press.

1163    Yang, B., Zhang, Z., Yang, C.Q., Wang, Y., Orr, M.C., Wang, H., & Zhang, A.B. 2022. Identification of

1164           species by combining molecular and morphological data using convolutional neural networks.

1165           Systematic Biology, 71, 690–705.

1166    Yelmen, B. & Jay, F. 2023. An Overview of Deep Generative Models in Functional and Evolutionary

1167           Genomics. Annual Reviews of Biomedical Data Science. https://doi.org/10.1146/annurev-biodatasci-

1168           020722.

1169    Zachos, F. E. 2016. Species concepts in biology (Vol. 801). Cham: Springer.

1170           2018. (New) Species concepts, species delimitation and the inherent limitations of

1171           taxonomy. Journal of genetics, 97, 811–815.

1172    Zaharias, P., Grosshauser, M. & Warnow, T. 2022. Re-evaluating Deep Neural Networks for Phylogeny

1173        Estimation: The Issue of Taxon Sampling. Journal of Computational Biology 29, 74–89.

1174        https://doi.org/10.1089/cmb.2021.0383.