

1 **Species delimitation using machine learning recovers a phylogenomically consistent**
2 **classification for North American box turtles (*Terrapene* spp.)**

3

4 **Running title:** Machine learning species delimitation

5

6 **Authors:**

7 Bradley T. Martin^{1,5,*}, Tyler K. Chafin¹, Marlis R. Douglas¹, John S. Placyk Jr.^{2,6}, Roger
8 D. Birkhead³, Chris A. Phillips⁴, Michael E. Douglas¹

9

10 ¹Department of Biological Sciences, University of Arkansas, Fayetteville, Arkansas 72701, USA;

11 Email: (BTM) btm002@uark.edu (send reprint requests to this address); (TKC):

12 tkchafin@uark.edu; (MRD): mrd1@uark.edu; (MED): med1@uark.edu.

13 ²Department of Biology, University of Texas, Tyler, Texas, 75799, USA; Email:

14 japlacyk@gmail.com

15 ³Alabama Science in Motion, Auburn University, Auburn, AL 36849, USA; Email:

16 birkhrd@auburn.edu

17 ⁴Illinois Natural History Survey, Prairie Research Institute, University of Illinois, Champaign, IL

18 61820; Email: caphilli@illinois.edu

19

20 Present Address:

21 ⁵Global Campus, University of Arkansas, 2 E. Center St., Fayetteville, Arkansas 72701, USA

22 ⁶Science Division, Trinity Valley Community College, Athens, Texas 75751, USA

23

24 **Disclosure statement:** Authors have nothing to disclose

25 **ABSTRACT**

26 Model-based approaches to species delimitation are constrained by computational capacities as
27 well as frequently violated algorithmic assumptions applied to biologically complex systems. An
28 alternate approach employs machine learning to derive species limits without explicitly defining
29 an underlying species model. Herein, we demonstrate the capacity of these approaches to identify
30 phylogenomically relevant groups in North American box turtles (*Terrapene* spp.). We invoked
31 several machine learning-based species delimitation algorithms and a multispecies coalescent
32 approach to parse a large ddRAD sequencing SNP dataset. We highlight two major findings: 1)
33 Machine learning delimitations were variable among replicates, but heterogeneity only occurred
34 within major species tree clades; 2) in this sense unsupported splits echoed patterns of
35 phylogenetic discordance among several species-tree methods. Discordance, as corroborated by
36 previously observed patterns of differential introgression, may reflect biogeographic history, gene
37 flow, incomplete lineage sorting, or their combinations. Our study underscores machine learning
38 as a species delimitation method, and provides insight into how commonly observed patterns of
39 phylogenetic discordance may similarly affect machine learning classification.

40

41 **Keywords:** discordance, species tree, VAE, t-SNE, phylogenomics, ddRAD

42 1. INTRODUCTION

43 Delineating species is undeniably crucial for systematics, ecology, and the evolutionary
44 process. Species are the currency of biodiversity, as are inconsistencies in the application of what
45 constitutes a species ('multiplicity' of species definitions; Zachos 2018). This creates
46 downstream issues for conservation (Mace 2004), where spurious 'splitting' or 'lumping' of taxa
47 are impediments to equitable allocation of limited resources. For example, over-splitting may
48 redundantly allow threatened/endangered taxa to proliferate (Zachos *et al.* 2013; Sullivan *et al.*
49 2014), or conflate recovery goals more appropriately managed at separate scales along the
50 species-population continuum (Coates *et al.* 2018).

51 On the other hand, inappropriate lumping can mask potential extinctions and the
52 recognition of adaptive differentiation (Stanton *et al.* 2019). This can bias 'true' diversity, as
53 reflected by regional or clade-specific differences in taxonomic 'culture' (e.g. biases in trait-
54 delimitation or species-concepts), or 'inertia' (i.e. persistent knowledge gaps; Gippoliti *et al.*
55 2018). Both disproportionately promote 'species at peril' and subsequently drive inefficient
56 resource allocation (Morrison *et al.* 2009; Garnett & Christidis 2017), viewed divisively as
57 'taxonomic inflation' (Agapow *et al.* 2004; Isaac *et al.* 2004). Nevertheless, species
58 definitions/delineations are a critical dimension in conservation's 'agony of choice' regarding
59 resource allocation (Vane-Wright *et al.* 1991; Stanton *et al.* 2019). Delimiting species impacts
60 not only finite resource allocation across programs but also efforts to recover and protect
61 biodiversity.

62 Earlier work on species delimitation relied on few genes (or markers), resulting in limited
63 scope. Although genomic approaches have shown promise (Allendorf *et al.* 2010), conflicting

64 genome-wide signals from incomplete lineage sorting (ILS) and gene flow (Funk & Omland
65 2003) are still apparent. Contemporary species delimitation relies upon a probabilistic approach
66 to model gene tree conflicts (i.e. multispecies coalescent; MSC) (Yang & Rannala 2010).
67 However, some models assume all such conflicts stem from ILS, and thus ignore other sources
68 such as introgressive hybridization.

69 Two popular packages, BPP and BFD*/SNAPP (Yang & Rannala 2010; Leaché *et al.*
70 2014a), are not only intractable with large datasets, but also seemingly over-split in the presence
71 of high population structure (Sukumaran & Knowles 2017) or when broad, continuous
72 geographic distributions are involved (Chambers & Hillis 2019). Therein lies the difficulty when
73 species delimitation explicitly assumes an underlying process of speciation (i.e. not effectively
74 modeled as an aspect of high-dimensionality data; Chafin *et al.* 2019). Here, we advocate
75 recently developed machine learning algorithms as an alternative that does not rely upon *a priori*
76 assumptions regarding the speciation process, but instead evaluates the process in a relatively
77 unrestricted manner.

78 Machine learning is broadly divided into two components: supervised (SML) and
79 unsupervised (UML). The former requires a classification model be ‘trained’ with *a priori*
80 designations, from which a classification model is derived and optimized for assignment of
81 ‘unknown’ data. A popular SML approach invokes support vector machines (SVM) that
82 partitions groups using linear or non-linear vectors in multi-dimensional space. However, the
83 requirement of an *a priori* classification scheme from which to train the model limits its
84 applicability, particularly when the purpose is to define groups, as in species delimitation.
85 Additionally, SVM is often computationally demanding, and hence slow with respect to

86 alternatives (Suryachandra & Reddy 2016). UML, on the other hand, requires no *a priori*
87 classification, and relies instead upon inherent patterns in the data.

88 Several popular UML classifiers lend themselves to the species delimitation problem,
89 including: Random Forest (RF; Breiman 2001), t-distributed stochastic neighbor embedding (t-
90 SNE; Maaten & Hinton 2008), and variational autoencoders (VAE; Derkarabetian *et al.* 2019),
91 each with inherent strengths and weaknesses. For example, RF uses randomly replicated data
92 subsets (in the form of pairwise distances) as a mechanism to develop binary ‘decision trees’ for
93 a classification model. All randomly seeded decision trees are aggregated (=‘forest’), with
94 classification decisions parsed as a majority vote amongst all trees. The random sub-setting
95 approach is relatively robust to correlations among features (=summary statistics or principle
96 components used for prediction) and model overfitting (=over-training the model where it does
97 not generalize well with new data). One stipulation is that features must be of low occupancy and
98 without undue noise (Rodriguez-Galiano *et al.* 2012). By contrast, the goal of t-SNE is to create
99 diagnosable clusters in reduced-dimension space, typically a 2D plane extracted from a
100 distillation of multi-dimensional data. Thus, it conceptually resembles methods such as principle
101 components analysis [(PCA) (Maaten & Hinton 2008)].

102 Alternatively, VAE uses neural networks in an attempt to ‘learn’ or reconstruct
103 multidimensional data patterns from a compressed, low-dimensionality (=‘encoded’)
104 representation. Again, the approach conceptually resembles the dimensionality-reduction
105 employed by various ordination techniques, but without linear and orthogonal constraint being
106 imposed upon the informative components. This approach may also be more statistically
107 interpretable (Derkarabetian *et al.* 2019), and thus more appropriate for the capture of variability

108 within highly complex data. Yet, careful consideration must be paid to the derivation of
109 parameters (e.g. neural network ‘depth’) that controls the encoding process (Livingstone *et al.*
110 1997).

111 UML methods do not require *a priori* designations from which to train a classification
112 model yet may still be sensitive to priors and parameter settings. Thus, guidelines for appropriate
113 application must be clearly defined, particularly regarding complex, empirical datasets. Two
114 metrics that can influence the support of a given species delimitation hypothesis is concordance
115 among algorithms (Carstens *et al.* 2013), and the susceptibility of the underlying algorithms to
116 common sources of phylogenetic discordance. Some machine learning algorithms are robust to
117 processes such as gene flow (Derkarabetian *et al.* 2019; Newton *et al.* 2020; Smith & Carstens
118 2020), but more empirical tests in complex systems are warranted. For example, performance can
119 vary among datasets, with potential influences including data quality (e.g. missing data
120 proportions) and size (Newton *et al.* 2020), historical demography, evolutionary history, and
121 coalescent processes such as incomplete lineage sorting (Austerlitz *et al.* 2009). Thus, we
122 empirically apply some recently developed software packages (CLADES: Pei *et al.* 2018; RF, t-
123 SNE, VAE: Derkarabetian *et al.* 2019) and discuss their capacity for evaluating a group of
124 species historically recalcitrant to taxonomic resolution.

125

126 1.1. *The convoluted evolutionary history of Terrapene*

127 North American box turtles (Emydidae: *Terrapene*) are primarily terrestrial, with a
128 common name based on an anterior ventral hinge that allows the plastron (bottom part of shell) to

129 dorsally close against the carapace (Dodd 2001). There are five currently recognized species
130 (Minx 1996; Iverson *et al.* 2017): Eastern (*Terrapene carolina*), Ornate (*T. ornata*), Florida (*T.*
131 *bauri*), Coahuilan (*T. coahuila*), and Spotted (*T. nelsoni*), with a sixth (*T. mexicana*) proposed
132 (Martin *et al.* 2013, 2014). *Terrapene carolina* includes two subspecies (Woodland: *T. c.*
133 *carolina*; Gulf Coast: *T. c. major*) that inhabit the eastern U.S. from the Mississippi River to the
134 Atlantic Ocean, and south through the Gulf Coastal Plain (Fig. 1). The putative *T. mexicana*
135 contains three subspecies (Three-toed: *T. m. triunguis*; Mexican: *T. m. mexicana*; Yucatan: *T. m.*
136 *yucatanana*) ranging across the southeastern and midwestern United States, the Mexican state of
137 Tamaulipas, and the Yucatan Peninsula. The Ornate (*T. ornata ornata*) and Desert (*T. o. luteola*)
138 box turtles inhabit the Midwest and Southwest U.S. plus the Northwest corner of México, while
139 the Southern and Northern Spotted box turtles (*T. nelsoni nelsoni* and *T. n. klauberi*) occupy the
140 Sonoran Desert in western México. *Terrapene coahuila* is semi-aquatic and restricted to Cuatro
141 Ciénegas (Coahuila, México), and the Florida box turtles occur in Peninsular Florida.

142 Morphological analyses delineate *T. carolina/mexicana* as a single species, sister to *T.*
143 *coahuila* (Minx 1992, 1996), with anecdotal support from a subset of genetic studies (Feldman &
144 Parham 2002; Stephens & Wiens 2003). Alternatively, Martin *et al.* (2013) proposed the
145 elevation of *T. mexicana* as a separate species, with *T. coahuila* as a subgroup within *T. carolina*.
146 In this latter study, *T. c. carolina* was sister to *T. c. major/T. coahuila*, although potential gene
147 flow was suspected between *T. c. carolina* and *T. c. major* due to mito-nuclear discordance.
148 Accordingly, *T. c. major* was recently demoted to an intergrade population and its subspecific
149 status removed (Butler *et al.* 2011; Iverson *et al.* 2017), but Martin *et al.* (2013) disagreed and a
150 more recent study identified two potentially pure *T. c. major* populations in the Florida and

151 Mississippi panhandles (Martin *et al.* 2020). Likewise, *T. bauri* (formerly *T. carolina bauri*) was
152 recently elevated to a distinct species (Butler *et al.* 2011; Iverson *et al.* 2017). a possibility that
153 Martin *et al.* (2013) acknowledged, albeit cautiously as weak statistical support and inconsistent
154 phylogenetic placement were evident. For the sake of clarity, we herein follow the
155 recommendations of Martin *et al.* (2013, 2014), considering *T. c. major* a distinct entity and *bauri*
156 as a subspecies within *T. carolina*. The monophyly of *T. o. ornata/luteola* has also been
157 questioned; Herrmann and Rosen (2009) suggested distinct lineages using microsatellite analyses,
158 whereas Martin *et al.* (2013) suggested polyphyly and a lack of phylogenetic structure using
159 mitochondrial (mt)DNA and nuclear (n)DNA sequences.

160 One likely reason for the historically enigmatic classification of *T. carolina* and *T.*
161 *mexicana* includes contemporary hybridization and introgression occurring within a hybrid zone
162 in the southeastern U.S., with four taxa potentially involved (Auffenberg 1958, 1959; Milstead &
163 Tinkle 1967; Milstead 1969). Some researchers (Fritz & Havaš 2013, 2014) interpret
164 reproductive semi-permeability as evidence for lumping the southeastern taxa as a single species.
165 However, divergent selection reinforcing species boundaries in some southeastern *Terrapene* has
166 been suggested as a reason for re-examining their classificatory status, despite ongoing gene flow
167 (Martin *et al.* 2014, 2020). Alternatively, the close phylogenetic relationship between *T. c. major*
168 and *T. coahuila* is less well understood. This may result from ‘ghost’ admixture of *T. coahuila*
169 and/or *T. c. major* with the extinct *T. c. putnami* (Martin *et al.* 2013).

170 Herein, we evaluate the classification of *Terrapene* within the context of both UML and
171 coalescent model-based species delimitation approaches. In doing so, we empirically validate the
172 use of machine learning approaches with complex genetic datasets that, upon analysis, support a

173 well-characterized phylogenetic hypothesis. Of note, observed species delimitation classifications
174 are consistent with patterns of phylogenetic discordance, demonstrating an empirical application
175 where the sources for such discordance may similarly affect machine learning.

176

177 **2. MATERIALS AND METHODS**

178 *2.1. Sample collection, storage, and DNA extraction*

179 Tissue samples were obtained from various museums, organizations, agencies, and
180 volunteers (Table S1), then stored in 70%-95% ethanol or DMSO (di-methyl sulfoxide) buffer.
181 Non-invasive samples were also acquired from live specimens, with those more invasive (e.g.
182 toes, muscle) taken from road-kills. Upon receipt, samples were stored at -20°C. Genomic DNA
183 was extracted via the following spin-column kits: DNeasy Blood and Tissue Kits (QIAGEN),
184 QIAamp Fast DNA Tissue Kit (QIAGEN), and E.Z.N.A. Tissue DNA Kits (Omega Bio-tek).
185 Extracted DNA was quantified using Qubit (Thermo Fisher Scientific) broad-range dsDNA
186 fluorometry and tested for high-molecular weight DNA using gel electrophoresis.

187

188 *2.2. DNA library preparation*

189 We first estimated the expected number of loci recovered via ddRAD sequencing
190 (ddRADseq) through *in silico* digestion (Chafin *et al.* 2018) of the painted turtle (*Chrysemys*
191 *picta*) genome (Shaffer *et al.* 2013). This was done to optimize choice of base-cutters, size-
192 selection bounds, and multiplex-size, thus maximizing loci coverage while promoting high
193 sequencing depth. We also used the *in silico* digest to identify a candidate size-selection that

194 avoids restriction sites lying within repetitive genomic elements (Chafin *et al.* 2018). The
195 expected number of ddRADseq loci and depth of coverage were empirically verified by
196 performing a restriction enzyme digest on 1,000ng of DNA for a representative panel of 24
197 samples, followed by fragment analysis (Agilent 4200 TapeStation).

198 Samples with sufficient DNA quantity (≥ 50 ng/uL) were processed via ddRADseq
199 protocol (Peterson *et al.* 2012). Between 500-1,000ng of genomic DNA per sample was digested
200 using two restriction enzymes, *PstI* (5'-CTGCA|G-3') and *MspI* (5'-C|CGG-3'). Following a
201 digestion at 37°C for 24 hours, 5ul of each sample was visualized on a 2% agarose gel via
202 electrophoresis to verify DNA fragmentation. Samples were purified using an AMPure XP
203 (Beckman Coulter) solution at a concentration of 1.5X (relative to DNA volume), then
204 standardized at 100ng of DNA per sample. Unique barcoded adapters were ligated to each
205 individual before pooling 48 samples into a library. Taxa were spread across multiple libraries to
206 mitigate potential batch effects, and libraries were size-selected on a Pippin Prep (Sage Science)
207 using the *in silico* optimized range [378-433 base pairs (bp), excluding adapters]. Lastly, a
208 twelve-cycle polymerase chain reaction (PCR) was run with Phusion DNA Polymerase (New
209 England BioLabs), followed by 1x100 single-end sequencing on the Illumina Hi-Seq 4000,
210 pooling two indexed libraries (=96 individuals) per lane. Sequencing and additional quality
211 control (fragment visualization and qPCR) were performed at the Genomics and Cell
212 Characterization Core Facility, University of Oregon/Eugene.

213

214

215 2.3. Sequence quality control and assembly

216 FASTQC v. 0.11.5 was used to assess sequence quality (Andrews 2010), with IPYRAD
217 v0.7.28 employed to demultiplex the raw sequences and align reads (Eaton & Overcast 2020).
218 Demultiplexed reads were allowed a strict maximum of one barcode mismatch, given that
219 barcodes were designed with a minimum two-base distance. Reads with low PHRED quality
220 scores (<33) were excluded, with additional filtering to remove adapter sequences. We then
221 performed reference-guided assembly using the *Terrapene m. mexicana* reference genome
222 (GenBank Accession #: GCA_002925995.2) with a minimum identity threshold of 0.85.
223 Unmapped reads were removed, and retained loci exhibited $\geq 20X$ coverage depth to reduce
224 sequencing error bias (Nielsen *et al.* 2011) and maximize phylogenetically informative sites in
225 the alignment (Eaton *et al.* 2017). Loci were further excluded if they displayed <50% individual
226 occupancy, excessive heterozygosity ($\geq 75\%$ of individual SNPs), or more than two alleles per
227 sample (the latter two instances indicating over-merged paralogs).

228

229 2.4. Phylogenomic inference

230 To assess differences in phylogenetic inference, we generated species trees using three
231 contemporary algorithms. Admixture across *Terrapene* hybrid zones has been well-characterized
232 (Butler *et al.* 2011; Martin *et al.* 2013, 2020). Thus, to mitigate the impact of contemporary gene
233 flow on phylogenetic inference, we only utilized individuals confirmed to be parental types
234 (characterized in Martin *et al.* 2019), as modelled using NEWHYBRIDS (Anderson & Thompson

235 2002). In so doing, we partitioned *T. c. major* into two subsets comprising two putative parental
236 populations.

237 Maximum likelihood phylogenies have been commonly produced for decades, yet the
238 increased use of large-scale SNP datasets often inflates bootstrap support for concatenated
239 phylogenomic datasets (Salichos & Rokas 2013; Simmons & Goloboff 2014). Coalescent-based
240 approaches that account for independent gene tree histories are more applicable for SNP analysis,
241 and thus we employed SVDQUARTETS [(Chifman & Kubatko 2014), implemented in PAUP*
242 v4.0a164 (Swofford 2003)] to produce a species tree with individuals grouped into populations.
243 Unrooted four-taxon gene trees were generated to assess legitimate splits, then assembled to form
244 the full species tree. SVDQUARTETS performs better for concatenated SNP datasets than do
245 species tree methods utilizing summary statistics (Chou *et al.* 2015), and importantly works well
246 with the large amount of missing data typically produced by ddRADseq (Leaché *et al.* 2015).

247 To reduce linkage bias and because independent gene tree histories are assumed for each
248 site, only one SNP from each ddRADseq locus was included in the SVDQUARTETS alignment. To
249 assess sampling variance, we ran 100 bootstrap replicates and considered nodes resampled at
250 >70% as strongly supported. Taxon partitions were grouped at the lowest level of field
251 identification (i.e. subspecific designations, when available), and by U.S. and Mexican state
252 locality. Blanding's (*Emydoidea blandingii*) and spotted (*Clemmys guttata*) turtles were included
253 as outgroups. An exhaustive search of all possible quartets was performed, with the consensus
254 tree visualized in FIGTREE v1.4.2 (Rambaut 2014).

255 We also employed a polymorphism aware model (PoMo: Schrempf *et al.* 2016), as
256 implemented in IQ-TREE v1.6.9 (Nguyen *et al.* 2015), to generate a second species tree. We did

257 so because POMO allows within-population polymorphism to account for ILS. The full IPYRAD
258 alignment, including invariant sites, was input into POMO and executed with 1,000 ultrafast
259 bootstrap (UFB) replicates (Hoang *et al.* 2017) and a maximum virtual population size of 19. The
260 discrete gamma rate model was applied (N=4), and clades with bootstrap support $\geq 95\%$ were
261 considered strongly supported.

262 Finally, we generated a lineage-tree phylogeny (IQ-TREE v1.7.12; Nguyen *et al.* 2015) to
263 contrast with our species-trees. An edge-linked partition model with 1,000 UFB replicates was
264 run using MODELFINDER (Kalyaanamoorthy *et al.* 2017) to determine the optimal substitution
265 model for each separate ddRADseq locus. Given computational constraints, model selection was
266 restricted only the general time reversible (GTR) model. Following tree reconstruction, IQ-TREE
267 was used to calculate site-wise concordance factors (sCF; Minh *et al.* 2018) for each branch
268 because they are less susceptible than traditional bootstrapping to over-inflation (Philippe *et al.*
269 2011). The sCF were calculated from 100 quartets randomly sampled from internal branches of
270 the tree, as recommended by IQ-TREE for stable sCF values. UFB $\geq 95\%$ and sCF $\geq 50\%$ were
271 considered as strong support (per IQ-TREE documentation).

272 For statistical topology tests, we generated lineage trees with IQ-TREE under the
273 topological constraints supported by four species-tree hypotheses derived from: (a)
274 SVDQUARTETS and (b) POMO topologies, as generated herein; (c) Sanger sequencing with
275 mtDNA and nuclear introns (Martin *et al.* 2013); and (d) Morphological data (Minx 1996).
276 MODELFINDER was again employed to optimize substitution models for each locus, as partitioned
277 in a concatenated supermatrix, using a hierarchical clustering algorithm to minimize
278 computational burden in IQ-TREE (*-rcluster*). We also toggled the *-bnni* and *-opt-gamma-inv*

279 options to reduce the impact of severe model violation and more thoroughly explore gamma and
280 invariant site parameters. Nodal confidence of individual trees was assessed using 1,000 UFB.
281 We then compared support for the concatenated supermatrix among constraint trees using seven
282 topological tests and 10,000 re-samplings: (a) Raw log-likelihoods; (b) bootstrap proportion test
283 using the RELI approximation (Kishino *et al.* 1990); (c) Kishino-Hasegawa test (Kishino &
284 Hasegawa 1989); (d) Shimodaira-Hasegawa test (SH; Shimodaira & Hasegawa 1999); (e)
285 Approximately Unbiased test (Shimodaira 2002); and (f) Expected Likelihood Weights
286 (Strimmer & Rambaut 2002). To visualize support for each topology across the genome, site-
287 likelihood probabilities and pairwise site-likelihood score differences (ΔSLS) were calculated
288 between the best-supported *versus* remaining trees.

289

290 2.5. *Species delimitation*

291 We employed the multispecies coalescent Bayes Factor Delimitation approach [BFD*;
292 (Leaché *et al.* 2014a)] as a baseline to compare the machine learning-based methods. Because
293 BFD* is computationally intensive, taxa were subset to a maximum of five individuals that
294 contained the least amount of missing data (N=37, plus outgroups), with sampling locations
295 varied (excepting *T. c. bauri* and the extremely rare *T. m. mexicana* and *T. coahuila*, which occur
296 exclusively in Peninsular Florida and the Mexican states of Tamaulipas and Cuatro Ciénegas).
297 For consistency, the same subset of individuals was used across all approaches. Details for BFD*
298 prior selection and additional data filtering steps can be found in Supplemental Appendix 1.

299 For each BFD* model, SNAPP employed 48 path-sampling steps, 200,000 burn-in, plus
300 400,000 MCMC iterations, with sampling every 1,000 generations. The path-sampling steps were
301 conducted with 200,000 burn-in, 300,000 MCMC generations, $\alpha=0.3$, 10 cross-validation
302 replicates, and 100 repeats. Trace plots were visualized (TRACER v1.7.1) to confirm parameter
303 convergence and compute effective sample sizes (ESS; Rambaut *et al.* 2018). Bayes factors (BF)
304 were calculated as $[2 \times (\text{MLE}_1 - \text{MLE}_2)]$ from the normalized marginal likelihood estimates
305 (MLE). We considered the following scheme for BF model support: $0 < \text{BF} < 2$ =no model
306 differentiation; $2 < \text{BF} < 6$ =positive; $6 < \text{BF} < 10$ =strong; and $\text{BF} > 10$ =decisive support (Kass &
307 Raftery 1995).

308 The RF and t-SNE algorithms (Breiman 2001; Maaten & Hinton 2008) were run and
309 visualized using an R script developed by Derkarabetian *et al.* (2019). The data were represented
310 as scaled principle components (N=37 axes) generated in ADEGENET v2.1.1 (Jombart & Ahmed
311 2011) in R v3.5.1 (R Development Core Team 2018). We averaged 100,000 majority-vote
312 decision trees over 10,000 bootstrap replicates to generate RF predictions. Clustered RF output
313 was visualized using both classic and isotonic multidimensional scaling procedures (CMDS and
314 ISOMDS; Shepard *et al.* 1972; Kruskal & Wish 1978). We ran t-SNE for 10,000 iterations within
315 which equilibria of the clusters was visually confirmed. Perplexity, which limits the effective
316 number of t-SNE neighbors, was tested at values of five and ten.

317

318

319

320 2.6. Determining optimal K for random forests and t -SNE

321 Two common clustering algorithms, as implemented in the aforementioned R scripts
322 (Derkarabetian *et al.* 2019), were used to derive optimal K for both the RF and t -SNE analyses.
323 The first [Partitioning Around Medoids (PAM); Kaufman and Rousseeuw 1987] attempts to
324 minimize the distance between the center point *versus* all other points of K clusters. The program
325 requires K to be defined *a priori*, and thus $K=1-10$ were tested, with the gap statistic and highest
326 mean silhouette widths [(MSW) (Rousseeuw 1987; Tibshirani *et al.* 2001)] determining optimal
327 K . The second [Hierarchical Agglomerative Clustering (HAC); Fraley and Raftery 1998] merges
328 points with minimal dissimilarity metrics (based on pairwise distances) until all are clustered.

329

330 2.7. Variational autoencoders

331 The VAE UML approach (Derkarabetian *et al.* 2019) employs neural networks and deep learning
332 to infer the marginal likelihood distribution of sample means (μ) and standard deviations [(σ) (i.e.
333 ‘latent variables’)]. Clusters with non-overlapping σ are interpreted as distinct clusters, or
334 ‘species.’ Data were input as 80% training/20% validation, with model loss (~error) visualized to
335 determine the optimal number of ‘epochs’ (=cycles through the training dataset). VAE should
336 ideally be terminated when model loss converges on a minimum value between training and
337 validation datasets [(i.e. the ‘Goldilocks zone’; Fig. S1) (Al’Aref *et al.* 2019)]. An escalating
338 model loss in the validation dataset indicates overfitting, whereas a failure to acquire a minimum
339 value points to underfitting (i.e. inability to generalize across both training and unseen data).

340

341 2.8. Support vector machines

342 The CLADES software (Pei *et al.* 2018) derives six summary statistics for SVM: 1)
343 Proportion of private alleles; 2) a folded site-frequency spectrum (SFS); 3) pairwise F_{ST} values
344 within populations; 4) pairwise F_{ST} values among populations; 5) the pairwise difference ratio
345 ($d_{\text{between}}/d_{\text{within}}$); and 6) the longest shared tract (longest string shared by two sequences). More
346 extensive methodological descriptions of the UML and SML components of machine learning are
347 found in Supplemental Appendix 1.

348

349 3. RESULTS

350 3.1. Sampling and data processing

351 We sequenced 214 geographically widespread *Terrapene* (Fig. 1; Table S1) including all
352 recognized species and subspecies, save the exceptionally rare *T. nelsoni klauberi*. When
353 possible, we included a minimum of 10 individuals per taxon, though fewer were used per rare
354 clade (*T. m. yucatanana*, *T. m. mexicana*, *T. coahuila*, *T. n. nelsoni*, *T. o. luteola*, and *T. c. bauri*).
355 The IPYRAD pipeline recovered 134,607 variable sites across 13,353 loci that mapped to the *T. m.*
356 *mexicana* genome, with 90,777 being parsimoniously informative. The mean per-individual
357 coverage depth was 56.3X (Fig. S2).

358

359

360

361 3.2. Species tree inferences

362 The sCF tree contained N=214 tips (Fig. 2), whereas SVDQUARTETS and PoMo (Fig. 3)
363 grouped individuals into N=26 populations, again based on locality and subspecies (when
364 provided). The SVDQUARTETS alignment contained 10,299 unlinked SNPs, with 87,395,061
365 quartets employed to assemble the species tree (Fig. 3a). Concatenated ddRADseq loci were
366 included in the PoMo tree (Fig. 3b), to include both invariable and variable sites
367 ($N_{\text{sites}}=1,163,463$). All trees clearly delineated eastern *versus* western clades, with *T. mexicana*, *T.*
368 *carolina*, and *T. coahuila* composing the eastern clade and the west represented by the
369 monophyletic *T. ornata* and *T. nelsoni*. However, some differences among methodologies were
370 apparent within these clades.

371 All phylogenies clearly delineated the western *T. ornata* and *T. nelsoni*. However,
372 SVDQUARTETS paraphyletically nested *T. o. luteola* within *T. o. ornata*, whereas IQ-TREE and
373 PoMo represented them as distinct monophyletic clades. In the eastern clade, SVDQUARTETS
374 displayed two subdivisions: *Terrapene mexicana* (all subspecies) and *T. carolina* (all subspecies)
375 + *T. coahuila*. PoMo did likewise, but also placed *T. m. triunguis* as paraphyletic in *T. mexicana*.
376 Furthermore, SVDQUARTETS, PoMo, and IQ-TREE each differed regarding the placement of *T.*
377 *c. bauri*, *T. coahuila*, and two previously recognized clades within *T. c. major* (Martin *et al.*
378 2013, 2020). Specifically, SVDQUARTETS depicted *T. c. bauri* as ancestral in the
379 *bauri/major/coahuila/carolina* clade, whereas PoMo placed *T. c. major* from MS/*coahuila* as
380 ancestor to *T. c. major* (FL)/*bauri/carolina*. However, IQ-TREE placed 1) *T. c. bauri* sister to all
381 of *T. carolina/T. mexicana*, and 2) *T. coahuila/T. c. major* (MS) sister to *T. c. carolina/T. c.*

382 *major* (FL). IQ-TREE also placed one *T. c. major* individual within the *T. m. triunguis* clade, and
383 one *T. c. carolina* as ancestral to the Floridian *T. c. major* and remaining *T. c. carolina*.

384

385 3.3. *Species tree reconciliation*

386 Trees representing Sanger data and SVDQUARTETS were in agreement when we contrasted
387 our topology tests, whereas morphology-based and POMO trees were both significantly rejected
388 (Table 1). Although the SVDQUARTETS tree was ranked the highest, site-likelihood scores
389 indicated that each topology was determined by a small number of loci (Fig. S3), whereas the
390 remaining majority was relatively uninformative.

391

392 3.4. *Species delimitation methods compared*

393 BFD* supported two top models (Table 2): All taxa delimited ($K=9$), and all as distinct
394 save *T. o. ornata*/*T. o. luteola* ($K=8$; Fig. 4). BF did not distinguish between the top models (<2),
395 although both were decisively better than all others ($BF>10$). Convergence was confirmed for the
396 likelihood traces, and the mean per-model ESS were >300 (Table S2).

397 The majority of the RF and t-SNE runs (Fig. 4) also grouped *T. o. ornata* and *T. o.*
398 *luteola*. However, the remaining clusters were split conservatively relative to BFD*. All runs
399 clearly delineated *T. ornata*, *T. carolina* and *T. mexicana* ssp., with some also delimiting as
400 distinct entities *T. c. carolina*, two *T. c. bauri* clusters, and *T. m. mexicana*. Of note, the runs and
401 clustering algorithms exhibited high within- but not among-clade variability for *T. carolina* and
402 *T. mexicana*, excepting MSW using ISOMDS.

403 Each clustering algorithm and ordination technique displayed its own inherent
404 characteristics. Essentially, cMDS and the gap statistic were inclined to split subclades of *T.*
405 *carolina* and *T. mexicana*, ISOMDS and MSW were the most conservative, and t-SNE and HAC
406 were intermediate, though HAC oscillated in agreement with MSW and the gap statistic (Fig. 4).
407 RF, but not t-SNE, varied among the 100 replicates, which was most pronounced for cMDS.
408 Heightened cMDS run variation highlights its inherent sensitivity to low among-group variability
409 (Olteanu *et al.* 2013). Finally, t-SNE optimal K increased with perplexity.

410 VAE initially agreed with BFD* in recognizing $K=8$, clumping *T. o. ornata/T. o. luteola*
411 and splitting all other taxa (Fig. 5a). However, assessments of model-loss indicated overfitting in
412 the sense that given enough epochs, the predictive model can perfectly ‘learn’ the training
413 dataset, with predictive capacity rapidly decreasing for unseen test data. To mitigate, we
414 identified in the model loss plot the transition point, or ‘elbow’ (Fig. 5b), where predictive
415 accuracy falls off for the test data, such that test *versus* training sets diverge in accuracy. This
416 occurred at a much lower number of sampled epochs ($N=2,000$) and was subsequently re-
417 initiated at a new termination point. Once overfitting was eliminated, an optimal $K=3$ was derived
418 (Figs. 4, 5c, 5d), in agreement with other UML methods. The model was also tested with
419 $N=1,000$ epochs (not shown), for which $K=3$ clusters again persisted.

420

421 3.5. Supervised machine learning

422 CLADES yielded optimal $K=2$ ($P=1.44e^{-4}$; Fig. 4; Table S3), but with highly discordant
423 clusters compared with prior results and phylogenomic findings: *Terrapene c. carolina/T. c.*

424 *bauri* emerged as one species, and the remaining seven taxa (*T. ornata*, *T. mexicana*, and the
425 remaining *T. carolina*) as a consistently paraphyletic second species (Figs. 2-3). The possibility
426 of outliers misleading the delimitations was also explored by removing two *T. c. bauri* and North
427 Carolina *T. c. carolina* that, in a subset of UML runs either formed a potential second cluster or
428 clustered instead with *T. c. bauri*. However, CLADES provided similar output without
429 phylogenetic cohesiveness ($K=2$; $P=6.88e^{-6}$) with *T. c. bauri*/*T. c. major* (MS population) as one
430 species, and the remainder forming the second. In both cases, the estimated probability for
431 optimal K was quite low.

432

433 3.6. Relative performance among approaches

434 All UML species delimitation methods converged on $K=3$ if considering RF and t-SNE
435 classifications that did not inter-mix. Three *Terrapene* species (plus *T. nelsoni*) were corroborated
436 (Martin *et al.* 2013, 2020), whereas the clumping of *T. mexicana* and *T. carolina* (Minx 1996)
437 was rejected. Machine learning approaches were also markedly faster than BFD*. For example,
438 RF, t-SNE, and VAE required ~10-30 min run time on a Desktop PC utilizing one Intel i5-3570
439 CPU core and 16 GB RAM. Comparatively, the twenty BFD* runs required ~4,000 total wall-
440 time hours (~200 hours/model), parallelized across 24-48 threads and utilizing 200 GB
441 RAM/model.

442

443 4. DISCUSSION

444 We observed substantial heterogeneity among machine learning species delimitation
445 approaches in resolving the southeastern *Terrapene* taxa, echoing previous morphological and
446 single-gene results (Milstead 1967, 1969; Milstead & Tinkle 1967; Butler *et al.* 2011; Martin *et*
447 *al.* 2013). However, groups exhibiting such heterogeneity may indicate the involved taxa are one
448 species, whereas deficit groups may support distinctiveness. Additionally—as argued below—
449 these were interpreted as a more appropriate reflection of taxon-specific biological patterns. Our
450 results represent an empirical test for the *de novo* application of these software packages to other
451 taxonomically-complex systems.

452

453 4.1. *Species Delimitation Approaches Reconciled in Terrapene*

454 Species trees provide a necessary phylogenetic context for species delimitation by
455 outlining hypothetical species compositions and identities. In our case, they underscored classic
456 discordance (Figs. 2-3), previously hypothesized via single-gene sequencing (Martin *et al.* 2013).
457 Differences were apparent in the ancestral progression of taxa, and in transitions between
458 monophyly *versus* paraphyly. Persistent uncertainties include: 1) Placement of *T. c. bauri*; 2)
459 monophyly of *T. mexicana* and 3) *T. o. luteola* subspecies status. Additionally, two individuals
460 were placed in unexpected clades, which was a far smaller proportion than previously seen in
461 single-gene datasets. We suggest the latter are examples of admixture, as both were collected
462 near a southeastern US hybrid zone (Martin *et al.* 2020), and suspect the other idiosyncrasies
463 represent either violations of the model or methodological artifacts.

464 Impacts of interspecific gene flow on species tree inference are well-characterized, yet
465 surprisingly, seldom modeled explicitly (Leaché *et al.* 2014b; Leaché & Oaks 2017). PoMo, for
466 example, constrains all nodes to the same N_e , a potentially poor assumption given contemporary
467 and possibly historical admixture (Martin *et al.* 2013, 2020). An examination of the species trees
468 alone reiterates previous single-gene taxonomic assessments. However, powerful species
469 delimitation assessments were utilized that provide a far more robust phylogenetic classification.

470 UML species delimitation inferences were consistent with the most recent phylogenetic
471 hypotheses (Martin *et al.* 2013; this study). *Terrapene o. ornata/luteola*, *T. c.*
472 *carolina/bauri/major/coahuila*, and *T. m. mexicana/triunguis* represent what we would consider
473 as species-level variants, within which each encompasses group assignment heterogeneity.
474 *Terrapene m. yucatana* falls within *T. mexicana*, and *T. nelsoni* as sister to *T. ornata*, although
475 their extreme rarity and concomitantly limited sampling (N=1) precluded them from species
476 delimitation analysis. Importantly, the variability in RF and t-SNE results primarily echoed
477 uncertainty found in the species tree analyses, including the distinctiveness of *T. o. luteola*
478 subspecific relationships within *T. carolina*. This variation also corresponded with the proclivities
479 of each algorithm. RF, for example, invokes a randomized process, with stochasticity perhaps
480 exacerbated by the phylogenetic discordance within *T. carolina*. t-SNE was influenced by its
481 perplexity parameter, with a second *T. c. major* group from Mississippi being a minor addition
482 for perplexity=5. This could underscore population structure among Mississippi and Florida *T. c.*
483 *major*. Finally, delimitations for both were strongly impacted by clustering algorithm, which
484 closely paralleled their own algorithmic tendencies. For example, the gap statistic often over-
485 estimates K (Dudoit & Fridlyand 2002; Yan & Ye 2007), MSW under-splits (Şenbabaoğlu *et al.*

486 2014), and outliers and noise particularly impact HAC (Kim *et al.* 2009; Şenbabaoğlu *et al.*
487 2014). This, in turn, may explain the varying extent of agreement between HAC and either MSW
488 or the gap statistic.

489 We suggest the variability observed among RF and t-SNE runs was due to a lack of
490 divergence within the more variable groups. Mixed classification was not observed among the *T.*
491 *mexicana*, *T. ornata*, and *T. carolina* groups, excepting RF ISOMDS based on MSW that only
492 differentiated *T. ornata* versus *T. carolina* (Fig. 4). The more conservative nature of ISOMDS was
493 reflected in several original empirical tests, which suggested a restriction to two clustering
494 dimensions may be more sensitive to higher genetic divergences (Derkarabetian *et al.* 2019), as
495 seems to be the case here. Otherwise, variability among taxa was constrained within respective
496 subspecific units.

497 VAE initially recovered results identical to BFD* ($K=8$), delimiting all taxa except *T. o.*
498 *ornata/T. o. luteola*. To ensure model training occurred appropriately, we more closely inspected
499 model loss and observed overfitting (Fig. 5b). The VAE script includes dropout regularization
500 methods, which randomly thin neural network nodes during model training to reduce overfitting
501 (Srivastava *et al.* 2014). However, regularization parameters can be sensitive to dataset properties
502 (e.g. large *versus* small/noisy *versus* tidy), and may not perform well for every dataset (Gal &
503 Ghahramani 2016; Derkarabetian *et al.* 2019). In model loss exploration, overfitting was
504 mitigated by early termination of model training when loss was at its minimum, though this could
505 also be accomplished by tuning dropout parameters. After correcting for overfitting, VAE also
506 delimited $K=3$ (i.e. *Terrapene mexicana*, *T. carolina*, and *T. ornata* ssp.), much like RF and t-
507 SNE if considering classification heterogeneity to indicate intra-specific relationships.

508

509 4.2. *Phylogenetic and biological support of species delimitations*

510 We suggest that identifying machine learning groups that consistently lack classification
511 overlap is one criterion to delimit species. In our case, these were corroborated as major species
512 tree clades, highlighting their complementary nature. In contrast, inconsistent species delimitation
513 assignments reflect many of the phylogenetic discordances observed in this and previous studies
514 (Butler *et al.* 2011; Martin *et al.* 2013). Potential underlying biological processes include
515 incomplete lineage sorting, ongoing primary divergence, hybridization, and/or complex
516 phylogeographic history [(e.g. isolation followed by secondary contact) (Mayr 1963; Barton &
517 Hewitt 1985; Rieseberg *et al.* 1999, 2007; Coyne & Orr 2004; Sousa & Hey 2013)]. Divergent
518 selection can counteract such processes and reinforce species boundaries (Feder *et al.* 2013). Our
519 species delimitation results are consistent with previously observed divergent selection at
520 candidate loci across *T. carolina* and *T. mexicana*, whereas it was absent for *T. c. carolina* and *T.*
521 *c. major* (Martin *et al.* 2020). Thus, *T. carolina* and *T. mexicana* may exhibit signatures of
522 secondary contact, whereas *T. c. major* and *T. c. carolina* may be earlier in the divergence
523 process. Alternatively, *T. c. major* could be an intergrade population between *T. c. carolina* and
524 *T. m. triunguis* (Butler *et al.* 2011), though the species trees disagree and two putative parental
525 populations persist (Martin *et al.* 2020). In this sense, *T. c. major* displays fairly disparate habitat
526 preferences, favoring salt marshes on the Gulf Coastal Plain, whereas *T. c. carolina* and *T. m.*
527 *triunguis* occupy mesic woodlands. The low differentiation between *T. c. major* and *T. c.*
528 *carolina* may result from *T. c. major* being restricted to the southeast. Here, *T. c. carolina*
529 possibly blocked northward expansion of *T. c. major*, with gene flow persisting across much of *T.*

530 *c. major*'s smaller range. Alternatively, it may have diverged more recently and now reflects
531 ongoing primary divergence.

532

533 4.3. Comparisons to other empirical studies

534 The capability of machine learning species delimitation algorithms to discount population
535 structure while isolating higher-level differentiation is corroborated by other recent studies
536 (Derkarabetian *et al.* 2019; Hedin *et al.* 2020). However, and Derkarabetian *et al.* (2019) and
537 Newton *et al.* (2020) emphasized the importance of integrative approaches, as they were able to
538 identify cryptic species by considering both VAE species delimitation and ecological niche
539 modeling. Given the increasing availability of geological resources, such integrative taxonomic
540 considerations may prove to be invaluable.

541 Excepting CLADES, the machine learning software used herein also seem robust to
542 hierarchical levels of genetic variation, having differentiated *T. carolina* versus *T. ornata* and the
543 less divergent *T. m. triunguis*. However, this hierarchical robustness may have limits, as one
544 recent geometric morphometric image-based deep learning study favored inter-generic over inter-
545 specific delimitations (Boer & Vos 2018). On the contrary, another recent study was more
546 accurate in recovering species-level delimitations rather than across genera, which they suggested
547 stemmed from less informed model training in low-diversity families with many unique species.
548 Recent and future work may also illuminate the impact of gene flow and population demography
549 on observed delimitations, processes that MSC approaches do not consider. For example,
550 DELIMITR incorporates models of secondary contact and divergence with gene flow into RF

551 classifiers for species delimitation (Smith & Carstens 2020). However, empirical tests of
552 DELIMITR tended to agree with the species delimited by BFD* and BP&P, whereas for *Terrapene*
553 RF, t-SNE, and VAE were more conservative than the MSC approaches. It may be that
554 *Terrapene* exhibits stronger population structure than the species included in the DELIMITR
555 applications, which can influence BFD* (Sukumaran & Knowles 2017). Finally, machine
556 learning frameworks may illuminate other potential sources of species tree discordance, with
557 recent applications predicting discordant species trees (Roettger *et al.* 2009), assessing historical
558 introgression despite ongoing gene flow (Burbrink & Gehara 2018), and identifying ILS
559 (Burbrink *et al.* 2020). Nevertheless, RF, t-SNE, and VAE are reported to at least be robust to
560 gene flow, with recent applications showing that they place admixed individuals between parental
561 clusters (Derkarabetian *et al.* 2019; Hedin *et al.* 2020; Newton *et al.* 2020). However, in our case
562 we avoided hybrids (as characterized by NewHybrids; Martin *et al.* 2020) due to frequent
563 introgression in the southeastern *Terrapene* hybrid zone.

564

565 4.4. *Conclusions*

566 UML approaches attempt to identify groups based on inherent structure in the data, and
567 accordingly are a natural extension to the species delimitation problem. In our case, a consensus
568 among UML approaches corroborated other axes of differentiation, whereas MSC-based
569 delimitations over-partitioned the data. Specifically, groups that were not supported by RF, t-
570 SNE, and VAE echoed classic patterns of phylogenetic uncertainty seen among our species trees,
571 which may be affected by previously observed genome-wide differential introgression.
572 Furthermore, in our case it seems likely that the phylogenetic signals affecting discordance are

573 similarly affecting the machine learning algorithms, which may include a combination of
574 historical biogeographic processes, gene flow, and incomplete lineage sorting. What is clear is
575 that delimiting almost every *Terrapene* taxon, as supported by BFD*, is probably not biologically
576 appropriate. Though MSC methods are undoubtedly still extremely useful, machine learning
577 provides a promising alternative for resolving long-standing biological problems. This may
578 particularly be the case for species that violate MSC model assumptions, as demonstrated by our
579 study system.

580

581 **ACKNOWLEDGEMENTS**

582 We would like to thank the numerous volunteers, organizations, and agencies that contributed
583 tissue samples (Table S1). We also thank both current students and alumni of the Douglas Lab, as
584 well as University of Arkansas faculty for support, advice, and guidance, to include: A. Alverson,
585 W. Anthonyamy, M. Bangs, J. Beaulieu, J. Koukl, S. Musmann, J. Pummill, and Z. Zbinden.
586 Sample collections were approved under three Animal Care and Use Committee (IACUC)
587 protocols: #113 (University of Texas at Tyler), and #16160 and #18000 (University of
588 Illinois/Champaign-Urbana). Funding sources included the Lucille F. Stickle Fund of the North
589 American Box Turtle Committee, the American Turtle Observatory (ATO), and two University
590 of Arkansas endowments [the Bruker Professorship in Life Sciences (MRD), and the 21st Century
591 Chair in Global Change Biology (MED)]. The Arkansas High Performance Computing Cluster
592 (AHPCC) and the Jetstream Cloud Service (NSF-XSEDE Research Allocation TG-BIO160065)
593 graciously supplied analytical resources.

594

595 5. REFERENCES

- 596 Agapow PM, Bininda-Emonds ORP, Crandall KA, Gittleman JL, Mace GM, Marshall JC, and
597 Purvis A (2004) The impact of species concept on biodiversity studies. *Quarterly Review of*
598 *Biology*, **79**, 161–179.
- 599 Al’Aref SJ, Anchouche K, Singh G, Slomka PJ, Kolli KK, Kumar A, Pandey M, Maliakal G,
600 Van Rosendael AR, and Beecy AN (2019) Clinical applications of machine learning in
601 cardiovascular disease and its relevance to cardiac imaging. *European Heart Journal*, **40**,
602 1975–1986.
- 603 Allendorf FW, Hohenlohe PA, and Luikart G (2010) Genomics and the future of conservation
604 genetics. *Nature Reviews Genetics*, **11**, 697–709.
- 605 Anderson EC and Thompson EA (2002) A model-based method for identifying species hybrids
606 using multilocus genetic data. *Genetics*, **160**, 1217–1229.
- 607 Andrews S (2010) FastQC: a quality control tool for high throughput sequence data.
608 <https://www.bibsonomy.org/bibtex/2b6052877491828ab53d3449be9b293b3/ozborn>.
- 609 Auffenberg W (1958) Fossil turtles of the genus *Terrapene* in Florida. *Bulletin of the Florida*
610 *State Museum*, **3**, 53–92.
- 611 Auffenberg W (1959) A Pleistocene *Terrapene* hibernaculum, with remarks on a second
612 complete box turtle skull from Florida. *Quarterly Journal of the Florida Academy of*
613 *Science*, **22**, 49–53.
- 614 Austerlitz F, David O, Schaeffer B, Bleakley K, Olteanu M, Leblois R, Veuille M, and Laredo C
615 (2009) DNA barcode analysis: a comparison of phylogenetic and statistical classification
616 methods. *BMC Bioinformatics*, **10**, S10.
- 617 Barton NH and Hewitt GM (1985) Analysis of hybrid zones. *Annual Review of Ecology and*
618 *Systematics*, **16**, 113–148.
- 619 Boer MJA and Vos RA (2018) Taxonomic classification of ants (Formicidae) from images using
620 deep learning. *BioRxiv*, 407452.
- 621 Breiman L (2001) Random Forests. *Machine Learning*, **45**, 5–32.
- 622 Burbrink FT and Gehara M (2018) The biogeography of deep time phylogenetic reticulation.
623 *Systematic Biology*, **67**, 743–755.
- 624 Burbrink FT, Grazziotin FG, Pyron RA, Cundall D, Donnellan S, Irish F, Keogh JS, Kraus F,
625 Murphy RW, and Noonan B (2020) Interrogating genomic-scale data for Squamata (lizards,
626 snakes, and amphisbaenians) shows no support for key traditional morphological
627 relationships. *Systematic Biology*, **69**, 502–520.

- 628 Butler JM, Dodd Jr. CK, Aresco M, and Austin JD (2011) Morphological and molecular evidence
629 indicates that the Gulf Coast box turtle (*Terrapene carolina major*) is not a distinct
630 evolutionary lineage in the Florida Panhandle. *Biological Journal of the Linnean Society*,
631 **102**, 889–901.
- 632 Carstens BC, Pelletier TA, Reid NM, and Satler JD (2013) How to fail at species delimitation.
633 *Molecular Ecology*, **22**, 4369–4383.
- 634 Chafin TK, Douglas MR, Martin BT, and Douglas ME (2019) Hybridization drives genetic
635 erosion in sympatric desert fishes of western North America. *Heredity*, **123**, 759–773.
- 636 Chafin TK, Martin BT, Musmann SM, Douglas MR, and Douglas ME (2018) FRAGMATIC: in
637 silico locus prediction and its utility in optimizing ddRADseq projects. *Conservation*
638 *Genetics Resources*, **10**, 325–328.
- 639 Chambers EA and Hillis DM (2019) The multispecies coalescent over-splits species in the case
640 of geographically widespread taxa. *Systematic Biology*, **69**, 184–193.
- 641 Chifman J and Kubatko L (2014) Quartet inference from SNP data under the coalescent model.
642 *Bioinformatics*, **30**, 3317–3324.
- 643 Chou J, Gupta A, Yaduvanshi S, Davidson R, Nute M, Mirarab S, and Warnow T (2015) A
644 comparative study of SVDquartets and other coalescent-based species tree estimation
645 methods. *BMC Genomics*, **16**, S2.
- 646 Coates DJ, Byrne M, and Moritz C (2018) Genetic Diversity and Conservation Units: Dealing
647 With the Species-Population Continuum in the Age of Genomics. *Frontiers in Ecology and*
648 *Evolution*, **6**, 165.
- 649 Coyne JA and Orr HA (2004) *Speciation*. Sinauer Associates, Sunderland, MA, USA.
- 650 Derkarabetian S, Castillo S, Koo PK, Ovchinnikov S, and Hedin M (2019) A demonstration of
651 unsupervised machine learning in species delimitation. *Molecular Phylogenetics and*
652 *Evolution*, **139**, 106562.
- 653 Dodd KC (2001) *North American Box Turtles, A Natural History*. University of Oklahoma Press,
654 Norman, OK, USA.
- 655 Dudoit S and Fridlyand J (2002) A prediction-based resampling method for estimating the
656 number of clusters in a dataset. *Genome Biology*, **3**, research0036.1.
- 657 Eaton DAR and Overcast I (2020) ipyrad: Interactive assembly and analysis of RADseq datasets.
658 *Bioinformatics*, **36**, 2592–2594.
- 659 Eaton DAR, Spriggs EL, Park B, and Donoghue MJ (2017) Misconceptions on missing data in
660 RAD-seq phylogenetics with a deep-scale example from flowering plants. *Systematic*
661 *Biology*, **66**, 399–412.
- 662 Feder JL, Flaxman SM, Egan SP, and Nosil P (2013) Hybridization and the build-up of genomic
663 divergence during speciation. *Journal of Evolutionary Biology*, **26**, 261–266.

- 664 Feldman CR and Parham JF (2002) Molecular phylogenetics of emydean turtles: Taxonomic
665 revision and the evolution of shell kinesis. *Molecular Phylogenetics and Evolution*, **22**, 388–
666 398.
- 667 Fraley C and Raftery AE (2002) *MCLUST: Software for model-based cluster and discriminant*
668 *analysis. Department of Statistics, University of Washington: Technical Report 415,*
669 *Department of Statistics, University of Washington, Seattle, WA, USA.*
670 <http://www.stat.washington.edu/raftery>.
- 671 Fritz U and Havaš P (2013) Order Testudines: 2013 update. In: Zhang, Z.-Q. (Ed.) *Animal*
672 *Biodiversity: An Outline of Higher-level Classification and Survey of Taxonomic Richness*
673 (Addenda 2013). *Zootaxa*, **3703**, 12–14.
- 674 Fritz U and Havaš P (2014) On the reclassification of Box Turtles (*Terrapene*): A response to
675 Martin et al. (2014). *Zootaxa*, **3835**, 295–298.
- 676 Funk DJ and Omland KE (2003) Species-level paraphyly and polyphyly: frequency, causes, and
677 consequences, with insights from animal mitochondrial DNA. *Annual Review of Ecology,*
678 *Evolution, and Systematics*, **34**, 397–423.
- 679 Gal Y and Ghahramani Z (2016) A theoretically grounded application of dropout in recurrent
680 neural networks. In: *Advances in Neural Information Processing Systems*, pp. 1019–1027.
- 681 Garnett ST and Christidis L (2017) Taxonomy anarchy hampers conservation. *Nature*, **546**, 25–
682 27.
- 683 Gippoliti S, Cotterill FPD, Zinner D, and Groves CP (2018) Impacts of taxonomic inertia for the
684 conservation of African ungulate diversity: an overview. *Biological Reviews*, **93**, 115–130.
- 685 Hedin M, Foldi S, and Rajah-Boyer B (2020) Evolutionary divergences mirror Pleistocene
686 paleodrainages in a rapidly-evolving complex of oasis-dwelling jumping spiders (Salticidae,
687 *Habronattus tarsalis*). *Molecular Phylogenetics and Evolution*, **144**, 106696.
- 688 Herrmann H and Rosen PC (2009) Conservation of aridlands turtles III: preliminary genetic
689 studies of the desert box turtle and yaqui slider. *Sonoran Herpetologist*, **22**, 38–43.
- 690 Hoang DT, Chernomor O, von Haeseler A, Minh BQ, and Vinh LS (2017) UFBoot2: improving
691 the ultrafast bootstrap approximation. *Molecular Biology and Evolution*, **35**, 518–522.
- 692 Isaac NJB, Mallet J, and Mace GM (2004) Taxonomic inflation: Its influence on macroecology
693 and conservation. *Trends in Ecology and Evolution*, **19**, 464–469.
- 694 Iverson JB, Meylan PA, and Seidel ME (2017) Testudines—Turtles. In: *Scientific and Standard*
695 *English Names of Amphibians and Reptiles of North America North of Mexico, with*
696 *Comments Regarding Confidence in Our Understanding* (ed Crother BI), pp. 82–91. SSAR
697 Herpetological Circular 43.
- 698 Jombart T and Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-wide SNP
699 data. *Bioinformatics*, **27**, 3070–3071.

- 700 Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, and Jermin LS (2017)
701 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods*, **14**,
702 587–589.
- 703 Kass RE and Raftery AE (1995) Bayes Factors. *Journal of the American Statistical Association*,
704 **90**, 773–795.
- 705 Kaufman L and Rousseeuw P (1987) Clustering by means of medoids. *Statistical Data Analysis*
706 *Based on the L1-Norm and Related Methods*, 405–416.
- 707 Kim E-Y, Kim S-Y, Ashlock D, and Nam D (2009) MULTI-K: accurate classification of
708 microarray subtypes using ensemble k-means clustering. *BMC Bioinformatics*, **10**, 260.
- 709 Kishino H and Hasegawa M (1989) Evaluation of the maximum likelihood estimate of the
710 evolutionary tree topologies from DNA sequence data, and the branching order in
711 hominoidea. *Journal of Molecular Evolution*, **29**, 170–179.
- 712 Kishino H, Miyata T, and Hasegawa M (1990) Maximum likelihood inference of protein
713 phylogeny and the origin of chloroplasts. *Journal of Molecular Evolution*, **31**, 151–160.
- 714 Kruskal JB and Wish M (1978) *Multidimensional Scaling*. Sage Publishing, Thousand Oaks, CA,
715 USA.
- 716 Leaché AD, Chavez AS, Jones LN, Grummer JA, Gottscho AD, and Linkem CW (2015)
717 Phylogenomics of phrynosomatid lizards: conflicting signals from sequence capture versus
718 restriction site associated DNA sequencing. *Genome Biology and Evolution*, **7**, 706–719.
- 719 Leaché AD, Fujita MK, Minin VN, and Bouckaert RR (2014a) Species delimitation using
720 genome-wide SNP data. *Systematic Biology*, **63**, 534–542.
- 721 Leaché AD, Harris RB, Rannala B, and Yang Z (2014b) The influence of gene flow on species
722 tree estimation: a simulation study. *Systematic Biology*, **63**, 17–30.
- 723 Leaché AD and Oaks JR (2017) The utility of single nucleotide polymorphism (SNP) data in
724 phylogenetics. *Annual Review of Ecology, Evolution, and Systematics*, **48**, 69–84.
- 725 Livingstone DJ, Manallack DT, and Tetko I V. (1997) Data modelling with neural networks:
726 Advantages and limitations. *Journal of Computer-Aided Molecular Design*, **11**, 135–142.
- 727 Maaten L van der and Hinton G (2008) Visualizing data using t-SNE. *Journal of Machine*
728 *Learning Research*, **9**, 2579–2605.
- 729 Mace GM (2004) The role of taxonomy in species conservation. *Philosophical Transactions of*
730 *the Royal Society B: Biological Sciences*, **359**, 711–719.
- 731 Martin BT, Bernstein NP, Birkhead RD, Koukl JF, Musmann SM, and Placyk JS (2013)
732 Sequence-based molecular phylogenetics and phylogeography of the American box turtles
733 (*Terrapene* spp.) with support from DNA barcoding. *Molecular Phylogenetics and*
734 *Evolution*, **68**, 119–134.
- 735 Martin BT, Bernstein NP, Birkhead RD, Koukl JF, Musmann SM, and Placyk Jr JS (2014) On

- 736 the reclassification of the *Terrapene* (Testudines: Emydidae): a response to Fritz & Havaš.
737 *Zootaxa*, **3835**, 292–294.
- 738 Martin BT, Douglas MR, Chafin TK, Placyk JS, Birkhead RD, Phillips CA, and Douglas ME
739 (2020) Differential introgression supports thermal adaptation and candidate genes shaping
740 species boundaries in North American box turtles (*Terrapene* spp.). *bioRxiv*, 752196.
- 741 Mayr E (1963) *Animal Species and Evolution*. Belknap Press at Harvard University Press,
742 Cambridge, MA.
- 743 Milstead WW (1967) Fossil box turtles (*Terrapene*) from central North America, and box turtles
744 of eastern Mexico. *Copeia*, **1967**, 168–179.
- 745 Milstead WW (1969) Studies on the evolution of the box turtles (genus *Terrapene*). *Bulletin of*
746 *the Florida State Museum, Biological Science Series*, **14**, 1–113.
- 747 Milstead WW and Tinkle DW (1967) *Terrapene* of Western Mexico, with comments on species
748 groups in the genus. *Copeia*, **1967**, 180–187.
- 749 Minh BQ, Hahn MW, and Lanfear R (2018) New methods to calculate concordance factors for
750 phylogenomic datasets. *bioRxiv*, 487801.
- 751 Minx P (1992) Variation in phalangeal formulas in the turtle genus *Terrapene*. *Journal of*
752 *Herpetology*, **26**, 234–238.
- 753 Minx P (1996) Phylogenetic relationships among the box turtles, Genus *Terrapene*.
754 *Herpetologica*, **52**, 584–597.
- 755 Morrison WR, Lohr JL, Duchen P, Wilches R, Trujillo D, Mair M, and Renner SS (2009) The
756 impact of taxonomic change on conservation: Does it kill, can it save, or is it just irrelevant?
757 *Biological Conservation*, **142**, 3201–3206.
- 758 Newton LG, Starrett J, Hendrixson BE, Derkarabetian S, and Bond JE (2020) Integrative species
759 delimitation reveals cryptic diversity in the southern Appalachian *Antrodiaetus unicolor*
760 (Araneae: Antrodiaetidae) species complex. *Molecular Ecology*, **29**, 2269–2287.
- 761 Nguyen L-T, Schmidt HA, von Haeseler A, and Minh BQ (2015) IQ-TREE: A fast and effective
762 stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology*
763 *and Evolution*, **32**, 268–274.
- 764 Nielsen R, Paul JS, Albrechtsen A, and Song YS (2011) Genotype and SNP calling from next-
765 generation sequencing data. *Nature Reviews Genetics*, **12**, 443.
- 766 Olteanu M, Nicolas V, Schaeffer B, Denys C, Missouf A, Kennis J, and Laredo C (2013)
767 Nonlinear projection methods for visualizing barcode data and application on two data sets.
768 *Molecular Ecology Resources*, **13**, 976–990.
- 769 Pei J, Chu C, Li X, Lu B, and Wu Y (2018) CLADES: A classification-based machine learning
770 method for species delimitation from population genetic data. *Molecular Ecology*
771 *Resources*, **18**, 1144–1156.

- 772 Peterson BK, Weber JN, Kay EH, Fisher HS, and Hoekstra HE (2012) Double digest RADseq: an
773 inexpensive method for de novo SNP discovery and genotyping in model and non-model
774 species. *PLOS ONE*, **7**, e37135.
- 775 Philippe H, Brinkmann H, Lavrov D V., Littlewood DTJ, Manuel M, Wörheide G, and Baurain D
776 (2011) Resolving difficult phylogenetic questions: why more sequences are not enough (D
777 Penny, Ed.). *PLOS Biology*, **9**, e1000602.
- 778 R Development Core Team (2018) R: A language and environment for statistical computing.
779 <https://cran.r-project.org/>.
- 780 Rambaut A (2014) FigTree v1.4.2. <http://tree.bio.ed.ac.uk/software/figtree/>.
- 781 Rambaut A, Drummond AJ, Xie D, Baele G, and Suchard MA (2018) Posterior summarization in
782 bayesian phylogenetics using Tracer 1.7 (E Susko, Ed.). *Systematic Biology*, **67**, 901–904.
- 783 Rieseberg LH, Kim S-C, Randell RA, Whitney KD, Gross BL, Lexer C, and Clay K (2007)
784 Hybridization and the colonization of novel habitats by annual sunflowers. *Genetica*, **129**,
785 149–165.
- 786 Rieseberg LH, Whitton J, and Gardner K (1999) Hybrid zones and the genetic architecture of a
787 barrier to gene flow between two sunflower species. *Genetics*, **152**, 713–727.
- 788 Rodriguez-Galiano VF, Ghimire B, Rogan J, Chica-Olmo M, and Rigol-Sanchez JP (2012) An
789 assessment of the effectiveness of a random forest classifier for land-cover classification.
790 *ISPRS Journal of Photogrammetry and Remote Sensing*, **67**, 93–104.
- 791 Roettger M, Martin W, and Dagan T (2009) A machine-learning approach reveals that alignment
792 properties alone can accurately predict inference of lateral gene transfer from discordant
793 phylogenies. *Molecular Biology and Evolution*, **26**, 1931–1939.
- 794 Rousseeuw PJ (1987) Silhouettes: A graphical aid to the interpretation and validation of cluster
795 analysis. *Journal of Computational and Applied Mathematics*, **20**, 53–65.
- 796 Salichos L and Rokas A (2013) Inferring ancient divergences requires genes with strong
797 phylogenetic signals. *Nature*, **497**, 327–333.
- 798 Schrempf D, Minh BQ, De Maio N, von Haeseler A, and Kosiol C (2016) Reversible
799 polymorphism-aware phylogenetic models and their application to tree inference. *Journal of*
800 *Theoretical Biology*, **407**, 362–370.
- 801 Şenbabaoğlu Y, Michailidis G, and Li JZ (2014) Critical limitations of consensus clustering in
802 class discovery. *Scientific Reports*, **4**, 1–13.
- 803 Shaffer HB, Minx P, Warren DE, Shedlock AM, Thomson RC, Valenzuela N, Abramyan J,
804 Amemiya CT, Badenhorst D, and Biggar KK (2013) The western painted turtle genome, a
805 model for the evolution of extreme physiological adaptations in a slowly evolving lineage.
806 *Genome Biology*, **14**, R28.
- 807 Shepard RN, Romney AK, and Nerlove SB (1972) *Multidimensional Scaling: Theory and*

- 808 *Applications in the Behavioral Sciences: I. Theory*. Seminar Press.
- 809 Shimodaira H (2002) An approximately unbiased test of phylogenetic tree selection. *Systematic*
810 *Biology*, **51**, 492–508.
- 811 Shimodaira H and Hasegawa M (1999) Multiple comparisons of log-likelihoods with applications
812 to phylogenetic inference. *Molecular Biology and Evolution*, **16**, 1114–1116.
- 813 Simmons MP and Goloboff PA (2014) Dubious resolution and support from published sparse
814 supermatrices: the importance of thorough tree searches. *Molecular Phylogenetics and*
815 *Evolution*, **78**, 334–348.
- 816 Smith ML and Carstens BC (2020) Process-based species delimitation leads to identification of
817 more biologically relevant species. *Evolution*, **74**, 216–229.
- 818 Sousa V and Hey J (2013) Understanding the origin of species with genome-scale data:
819 modelling gene flow. *Nature Reviews Genetics*, **14**, 404–414.
- 820 Srivastava N, Hinton G, Krizhevsky A, Sutskever I, and Salakhutdinov R (2014) Dropout: a
821 simple way to prevent neural networks from overfitting. *The Journal of Machine Learning*
822 *Research*, **15**, 1929–1958.
- 823 Stanton DWG, Frandsen P, Waples RK, Heller R, Russo IRM, Orozco-terWengel PA, Pedersen
824 CET, Siegismund HR, and Bruford MW (2019) More grist for the mill? Species delimitation
825 in the genomic era and its implications for conservation. *Conservation Genetics*, **20**, 101–
826 113.
- 827 Stephens PR and Wiens JJ (2003) Ecological diversification and phylogeny of emydid turtles.
828 *Biological Journal of the Linnean Society*, **79**, 577–610.
- 829 Strimmer K and Rambaut A (2002) Inferring confidence sets of possibly misspecified gene trees.
830 *Proceedings of the Royal Society of London. Series B: Biological Sciences*, **269**, 137–142.
- 831 Sukumaran J and Knowles LL (2017) Multispecies coalescent delimits structure, not species.
832 *Proceedings of the National Academy of Sciences of the United States of America*, **114**,
833 1607–1611.
- 834 Sullivan BK, Douglas MR, Walker JM, Cordes JE, Davis MA, Anthonysamy WJB, Sullivan KO,
835 and Douglas ME (2014) Conservation and management of polytypic species: The Little
836 Striped Whiptail Complex (*Aspidoscelis inornata*) as a case study. *Copeia*, **2014**, 519–529.
- 837 Suryachandra P and Reddy PVS (2016) Comparison of machine learning algorithms for breast
838 cancer. In: *Proceedings of the International Conference on Inventive Computation*
839 *Technologies, ICICT 2016*, pp. 1–6. Institute of Electrical and Electronics Engineers Inc.
- 840 Swofford DL (2003) *PAUP*. Phylogenetic Analysis Using Parsimony (*and Other Methods).*
841 *Version 4*. Sinauer Associates, Sunderland, Massachusetts.
- 842 Tibshirani R, Walther G, and Hastie T (2001) Estimating the number of clusters in a data set via
843 the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,

- 844 **63**, 411–423.
- 845 Vane-Wright RI, Humphries CJ, and Williams PH (1991) What to protect? Systemations and the
846 agony of choice. *Biological Conservation*, **55**, 235–254.
- 847 Yan M and Ye K (2007) Determining the number of clusters using the weighted gap statistic.
848 *Biometrics*, **63**, 1031–1037.
- 849 Yang Z and Rannala B (2010) Bayesian species delimitation using multilocus sequence data.
850 *Proceedings of the National Academy of Sciences*, **107**, 9264–9269.
- 851 Zachos FE (2018) Species concepts, species delimitation and the inherent limitations of
852 taxonomy. *Journal of Genetics*, **97**, 811–815.
- 853 Zachos FE, Apollonio M, Bärmann E V., Festa-Bianchet M, Göhlich U, Habel JC, Haring E,
854 Kruckenhauser L, Lovari S, McDevitt AD, Pertoldi C, Rössner GE, Sánchez-Villagra MR,
855 Scandura M, and Suchentrunk F (2013) Species inflation and taxonomic artefacts-A critical
856 comment on recent trends in mammalian classification. *Mammalian Biology*, **78**, 1–6.
- 857
- 858

859 **DATA ACCESSIBILITY**

860 The raw ddRADseq data is available on the GenBank Nucleotide Database at
861 <https://www.ncbi.nlm.nih.gov/bioproject/563121> (BioProject ID: 563121) [to be made public
862 upon publication]. Additional sequence alignments, Supplemental Appendix 1, and
863 supplementary materials will be available from a Dryad Digital Repository.

864

865 **AUTHOR CONTRIBUTIONS**

866 BTM and TKC designed the research, laboratory protocols, and scripts. BTM conducted the lab
867 work and bioinformatic analyses, analyzed the data, and wrote the manuscript. MRD and MED
868 were the study supervisors, guided the study design, and provided funding. JSP facilitated the
869 collection of thousands of *Terrapene* tissues and provided methodological expertise. RDB
870 collected hundreds of *Terrapene* tissues from southeastern North America and facilitated the
871 collection of many additional individuals. CAP provided many of the *T. ornata* tissues and
872 provided sampling expertise. All authors edited and revised the manuscript

873

874 **Table 1:** Topology tests for four hypothesized North American box turtle (*Terrapene*)
875 phylogenies. The morphology and Sanger sequencing trees are based on previously published
876 data (Minx 1996; Martin *et al.* 2013), whereas trees representing SVDQUARTETS and PoMo
877 (Polymorphism-Aware Model) were generated from ddRADseq data (Fig. 3, this study). Bolded
878 *P*-values with an asterisk (*) indicate supported trees ($P > 0.05$ or highly weighted).

Guide Tree	Log-likelihood	ΔLL	bp-RELL	p-KH	p-SH	c-ELW	p-AU
Morphology	-2639307.9	601.5	0.00	0.01	0.02	0.00	0.01
PoMo	-2639200.2	493.8	0.01	0.03	0.06*	0.01	0.03
Sanger	-2638898.4	192.0	0.23*	0.24*	0.41*	0.23*	0.26*
SVDquartets	-2638706.4	0.0	0.75*	0.76*	1.00*	0.75*	0.81*

879 ΔLL =change in log-likelihood

880 bp-RELL=Bootstrap proportions using REll method (weights sum to 1)

881 p-KH=Kishino-Hasegawa test

882 p-SH=Shimodaira-Hasegawa test

883 c-ELW=Expected likelihood weight (sum to 1)

884 p-AU=Approximately unbiased test

885

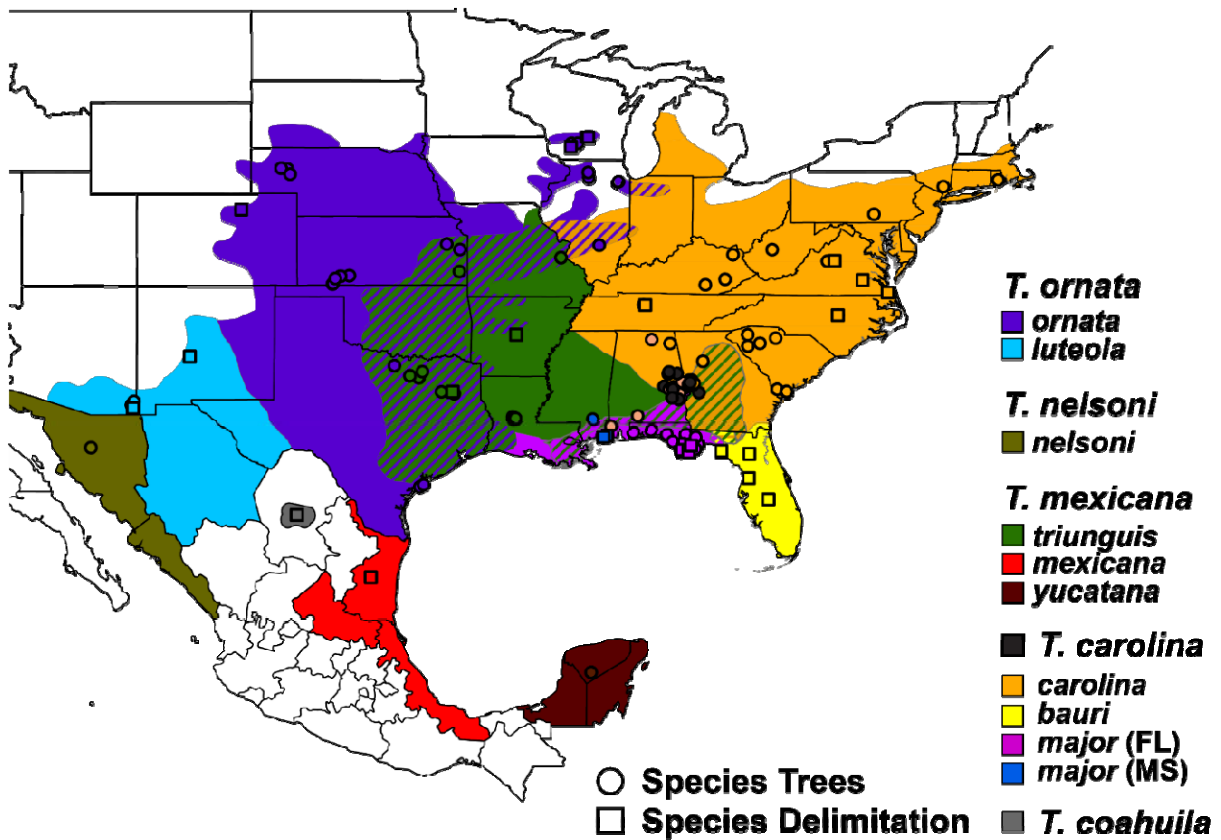
886

887 **Table 2:** Species delimitation results from BFD* (Bayes Factor Delimitation, *with genomic
 888 data) in 37 North American box turtle (*Terrapene* spp.) individuals. BFD* was run with 179
 889 unlinked, bi-allelic single nucleotide polymorphisms (SNPs) generated using ddRADseq. Bayes
 890 factors (BF) were used to identify support among models and were calculated as $2 \times (\text{MLE}_1 - \text{MLE}_2)$. An * indicates the best supported models; “+” shows taxa that were grouped together; “/”
 891 delineates multiple groupings. DS=Desert (*T. o. luteola*), ON=Ornate (*T. o. ornata*),
 892 EA=Woodland (*T. c. carolina*), GUFL=Gulf Coast (*T. c. major*) from Florida, GUMS=Gulf
 893 Coast (*T. c. major*) from Mississippi, CH=Coahuilan (*T. coahuila*), FL=Florida (*T. c. bauri*),
 894 TT=Three-toed (*T. m. triunguis*), and MX=Mexican (*T. m. mexicana*) box turtles. East=all *T.*
 895 *carolina* and *T. mexicana*, West=all *T. ornata*. The outgroup (not shown) included the spotted
 896 turtles (*Clemmys guttata*).
 897

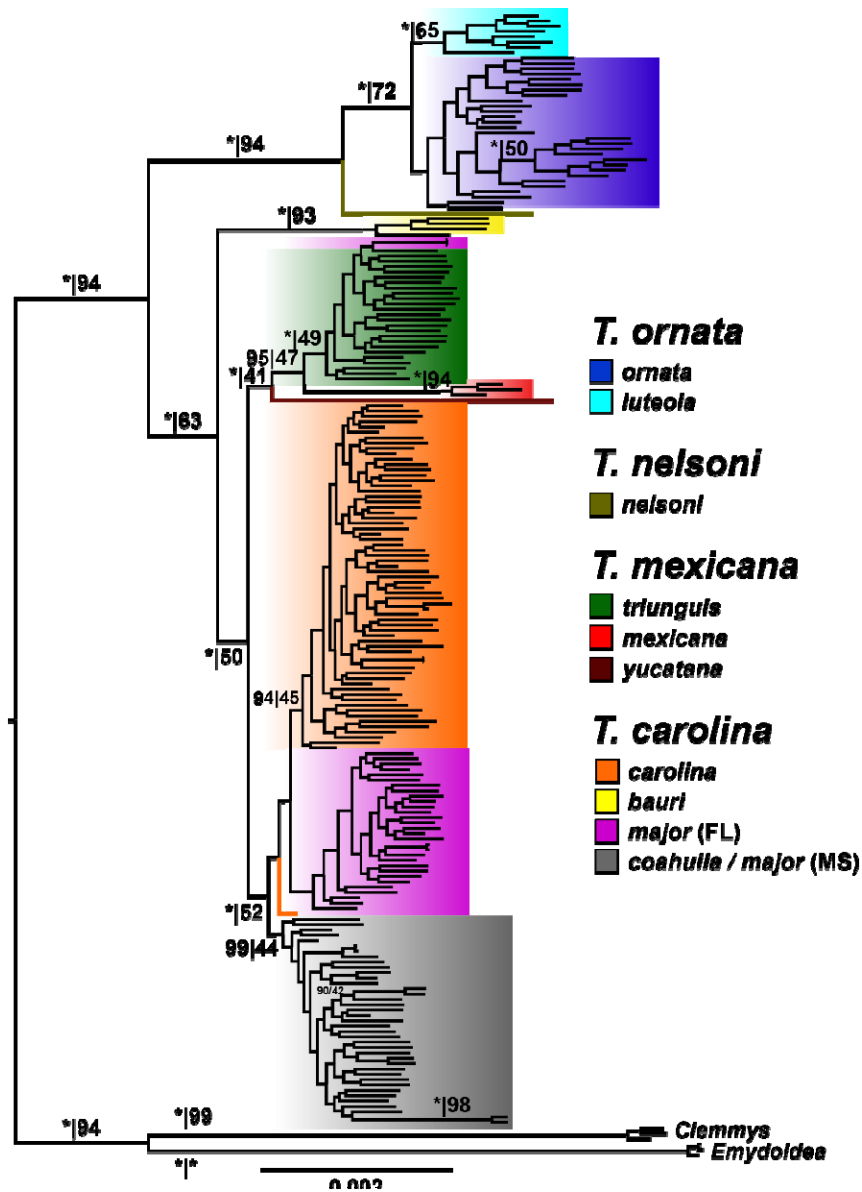
BFD* Model	MLE†	K‡	Rank§	BF¶
All Separate*	-2403.39	10	1	-
DS+ON*	-2404.34	9	2	1.90
EA+GUFL	-2417.84	9	3	28.91
GUMS+GUFL	-2427.58	9	4	48.39
GUMS+CH	-2448.61	9	5	90.44
GUMS+CH/GUFL+EA	-2461.28	8	6	115.79
GUMS+GUFL+CH	-2489.62	8	7	172.45
EA+FL	-2511.83	9	8	216.89
GUMS+GUFL+CH+EA	-2514.86	7	9	222.94
EA+FL+GUFL	-2552.22	8	10	297.66
EA+FL/CH+GUMS	-2555.16	8	11	303.53
EA+FL+GUFL/CH+GUMS	-2594.91	7	12	383.04
EA+CH+GUMS+GUFL+TT	-2607.72	6	13	408.66
EA+CH+GUMS+GUFL+MX	-2657.48	6	14	508.19
EA+FL+CH+GUMS+GUFL	-2693.37	6	15	579.96
EA+CH+GUMS+GUFL+TT+MX	-2719.02	5	16	631.27
ON+DS/EA+TT+MX+CH+GUMS+GUFL/FL	-2720.23	4	17	633.69
EA+FL+CH+GUMS+GUFL+TT	-2800.56	5	18	794.35
EA+FL+CH+GUMS+GUFL+TT+MX	-2926.20	4	19	1045.62
East/West	-2926.56	3	20	1046.35

898 †MLE=Marginal likelihood estimates
 899 ‡K=# tips
 900 §Rank=model ranking based on MLE (lower=better)
 901 ¶BF=Bayes factors
 902
 903

904



905
906 **Figure 1:** Range map for North American box turtles, *Terrapene*, depicting sample localities.
907 Cross-hatched areas indicate known hybrid zones. Circles represent samples used for the species trees
908 trees, whereas squares were also used only in species delimitation analyses. Headings and
909 subheadings correspond to species and subspecies, respectively, following (Martin *et al.* 2013),
910 and include the Ornate (*T. ornata ornata*), Desert (*T. o. luteola*), Spotted (*T. nelsoni*), Three-toed
911 (*T. mexicana triunguis*), Mexican (*T. m. mexicana*), Yucatan (*T. m. yucatana*), Woodland (*T.*
912 *carolina carolina*), Florida (*T. c. bauri*), Gulf Coast from distinct Florida and Mississippi
913 populations (*T. c. major*), and Coahuilan (*T. coahuila*) box turtles. Localities with black circles
914 indicate *T. carolina* samples lacking subspecific field identifications.
915



916
 917 **Figure 2:** Maximum likelihood phylogeny (IQ-TREE) reflecting relationships among 214
 918 *Terrapene* samples. The tree was generated from 11,962 unlinked ddRADseq loci with 1,000
 919 ultrafast bootstrap (UFB) replicates. Site concordance-factors (SCF) were calculated from 100
 920 quartets randomly sampled from internal branches. Branch support values represent UFB
 921 replicates and SCF on the left and right of each vertical line, respectively. UFBs $\geq 95\%$ and
 922 SCF $\geq 50\%$ were considered strong support. UFBs < 90 and SCF < 40 were omitted for visual clarity,
 923 with the latter rounded to the nearest integer. Asterisks (*) indicate 100% support. Legend
 924 headers and subheaders depict species and subspecies from Martin et al. (2013): Ornate
 925 (*Terrapene ornata ornata*), Desert (*T. o. luteola*), Spotted (*T. nelsoni*), Three-toed (*T. mexicana*
 926 *triunguis*), Mexican (*T. m. mexicana*), Yucatan (*T. m. yucatanae*), Florida (*T. carolina bauri*), Gulf
 927 Coast (*T. c. major*; two populations from Florida and Mississippi), Woodland (*T. c. carolina*),
 928 and Coahuilan (*T. coahuila*) box turtles. The Spotted (*Clemmys guttata*) and Blanding's
 929 (*Emydoidea blandingii*) turtles were used as outgroups.

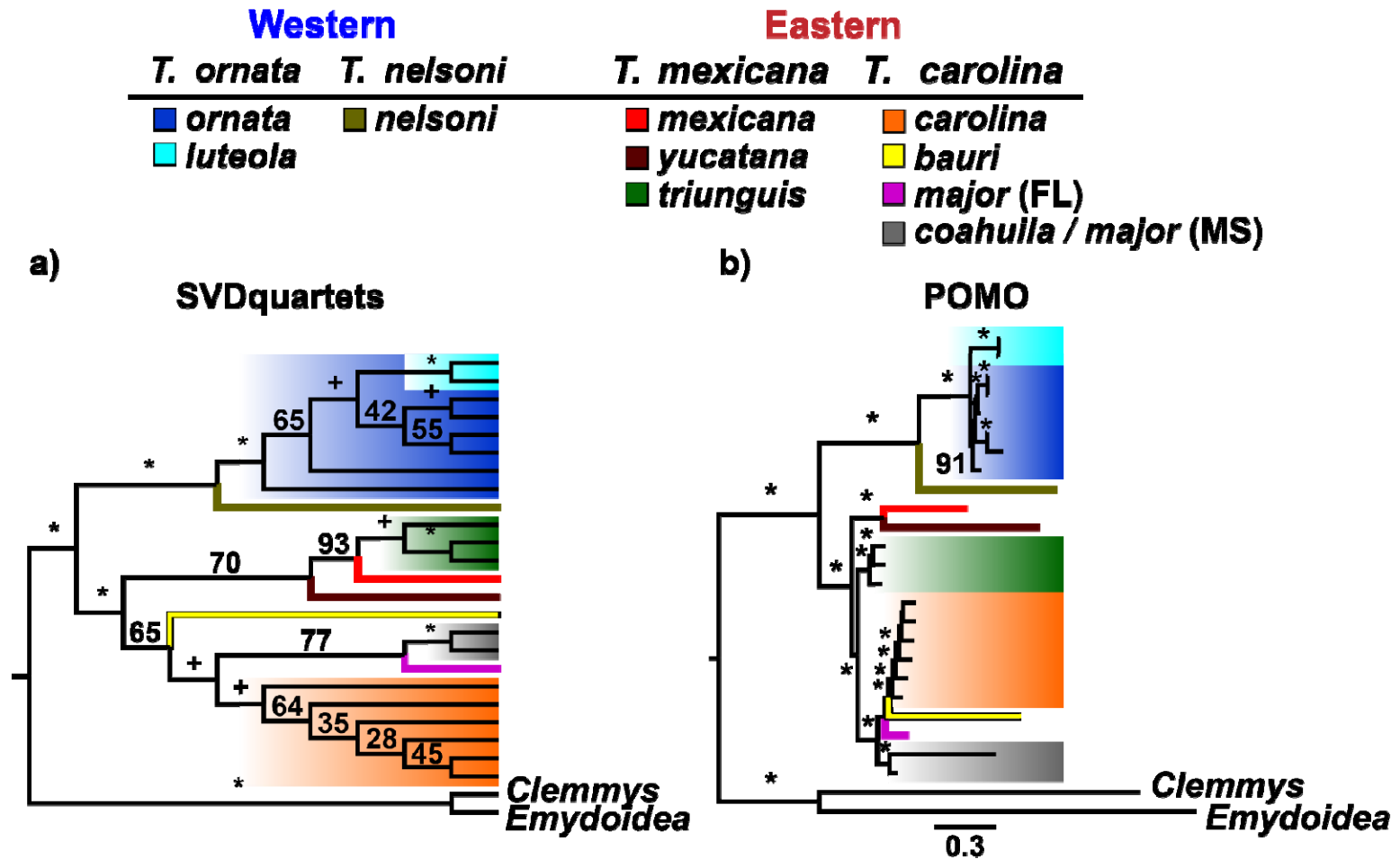


Figure 3: Species trees based on 10,299 unlinked ddRADseq loci depicting phylogenetic relationships of 214 North American box turtle (*Terrapene* spp.) samples. Species trees were generated using a) SVDQUARTETS and b) POMO (Polymorphism Aware Model). Each tree contained 26 populations grouped by specific or subspecific designations (if available), and U.S. or Mexican State locality. Spotted (*Clemmys guttata*) and Blanding's (*Emydoidea blandingii*) turtles were used as outgroups to root the trees. * and + above branches represent nodes with 100% and $\geq 95\%$ bootstrap support.

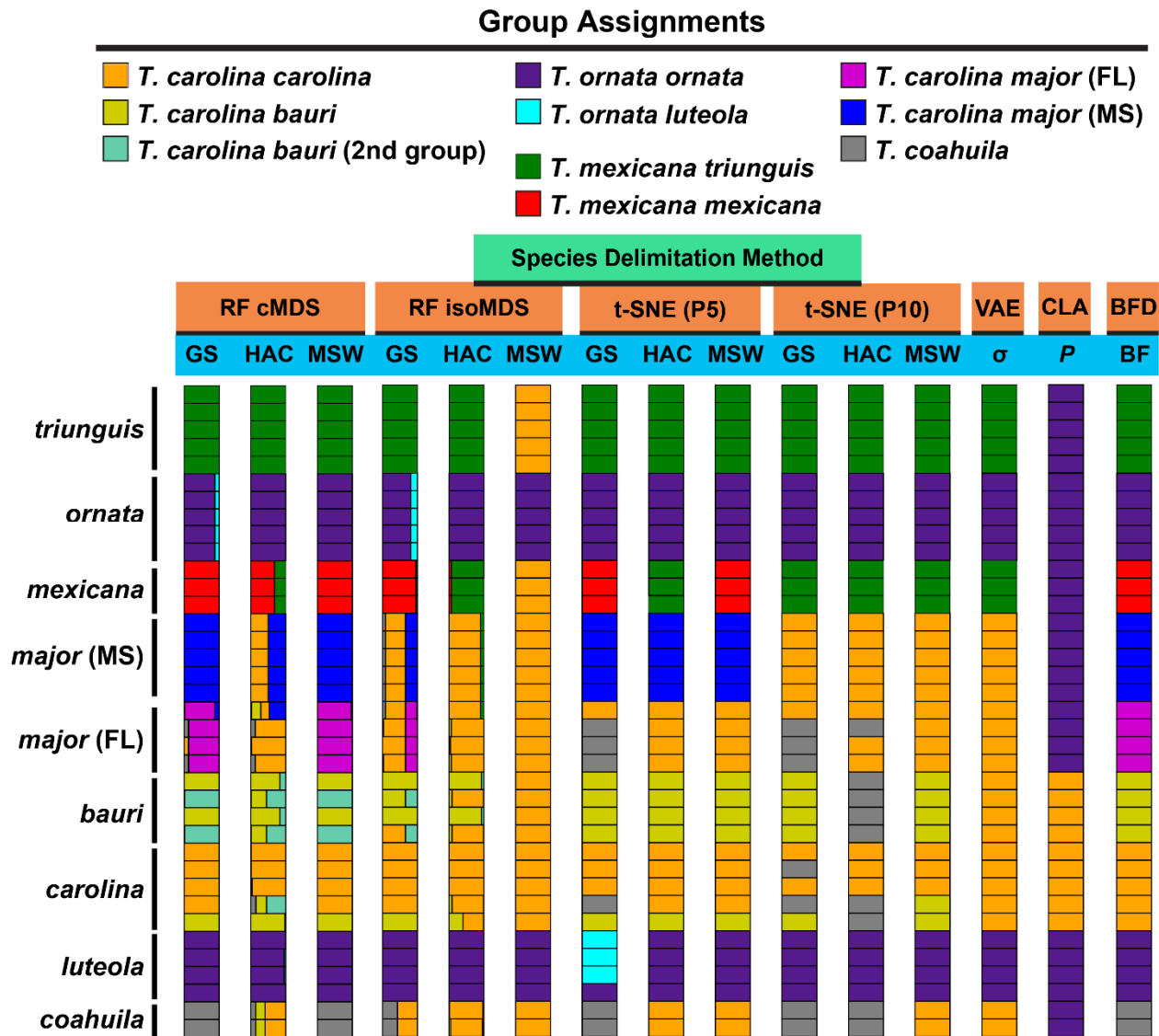


Figure 4: Species delimitations among 37 *Terrapene* samples using 7,395 unlinked ddRADseq single nucleotide polymorphisms (SNPs). Species delimitation groups were derived using unsupervised (UML) and supervised (SML) machine learning algorithms, plus a multispecies coalescent (MSC) approach. UML algorithms include RF=random forest, visualized using cMDS and isoMDS (classic and isotonic multidimensional scaling) ordination, t-SNE=t-distributed stochastic neighbor embedding, and VAE=variational autoencoders. Each cMDS, isoMDS, and t-SNE column represents a summarization of 100 independent runs, with colors indicating percent group assignments per method. Mixed colors show clustering variation among runs. t-SNE was run with perplexity settings of five and ten (P5 and P10). RF and t-SNE optimal K 's were assessed using hierarchical agglomerative clustering (HAC), partition around medoids using mean silhouette widths (MSW) and the gap statistic (GS), whereas standard deviation (σ) overlap was used for VAE. Optimal K 's for CLA=CLADES (SML) and BFD=Bayes Factor Delimitation were determined using probabilities (P) and Bayes Factors (BF).

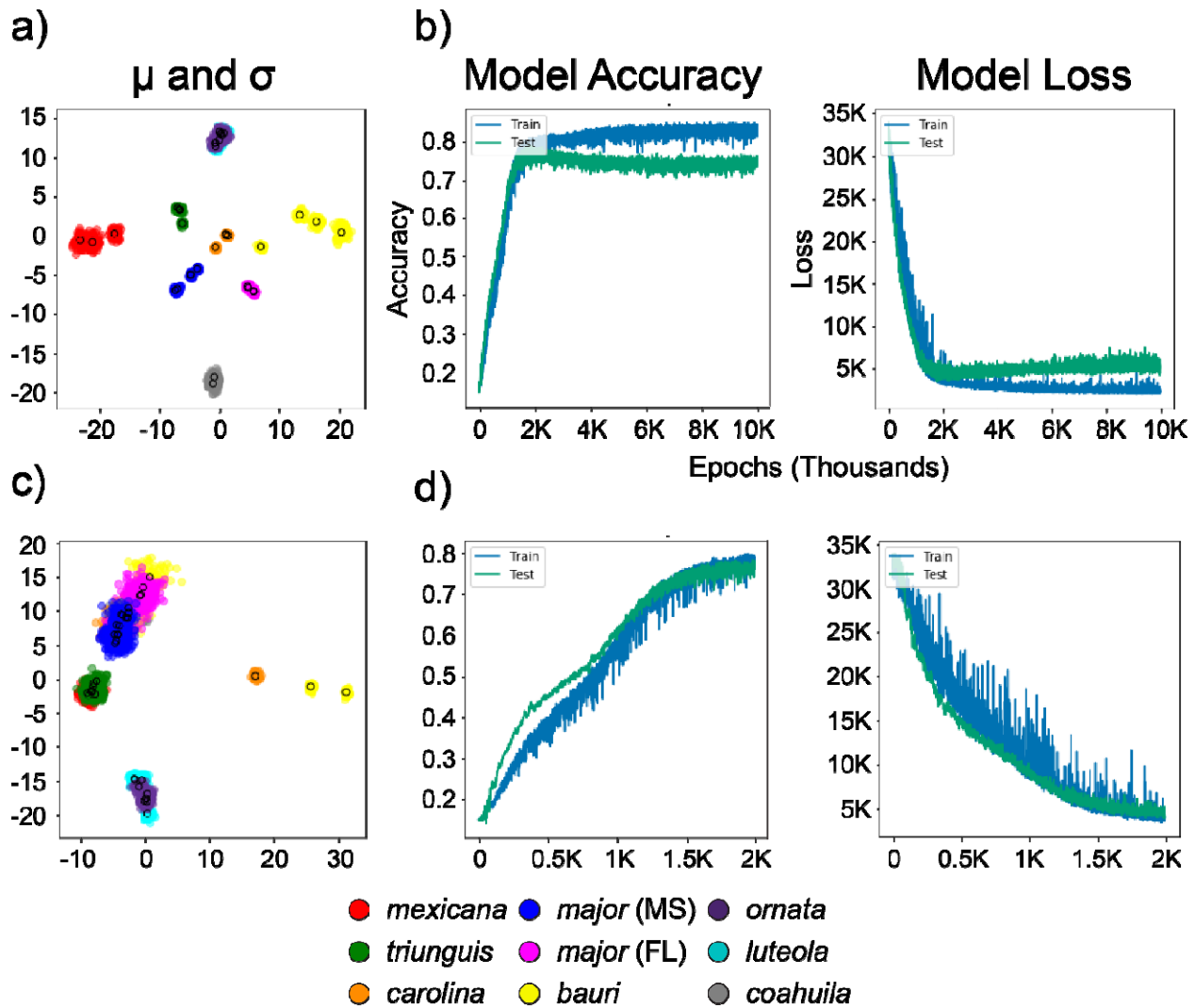


Figure 5: Results from variational autoencoder (VAE) machine learning species delimitation among 37 *Terrapene* samples and 7,395 unlinked ddRADseq single nucleotide polymorphisms (SNPs). Each circle represents the mean (μ) of one individual in the reconstructed parameter space, and the surrounding amorphous area are the standard deviations (σ) across a) 10,000 and c) 2,000 epochs. The model accuracy and loss traces depict the fit of the model to test (green) and training (blue) data across b) 10,000 and d) 2,000 epochs. The colors depict eight subspecies across *T. mexicana*, *T. carolina*, and *T. ornata*, plus monotypic *T. coahuila*.