1  PipeMaster: inferring population divergence and demographic history with approximate
2  Bayesian computation and supervised machine-learning in R

3  Marcelo Gehara[1], Guilherme G. Mazzochinni[2], Frank Burbrink[1]

4  Corresponding author: Marcelo Gehara; email: marcelo.gehara@gmail.com

5  1 - American Museum of Natural History, Herpetology, Central Park West at 79[th] St, New York,
6  NY
7  2 - Programa de Pós-graduação em Biologia de Fungos, Algas e Plantas, Centro de Ciências
8  Biológicas, Universidade Federal de Santa Catarina, Florianópolis, Brazil

9  **Abstract**

10        Understanding population divergence involves testing diversification scenarios and

11  estimating historical parameters, such as divergence time, population size and migration rate.

12  There is, however, an immense space of possible highly parameterized scenarios that are

13  difficult or impossible to solve analytically. To overcome this problem researchers have used

14  alternative simulation-based approaches, such as approximate Bayesian computation (ABC)

15  and supervised machine learning (SML), to approximate posterior probabilities of hypotheses. In

16  this study we demonstrate the utility of our newly developed R-package to simulate summary

17  statistics to perform ABC and SML inferences. We compare the power of both ABC and SML

18  methods and the influence of the number of loci in the accuracy of inferences; and we show

19  three empirical examples: (i) the Muller's termite frog genomic data from Southamerica; (ii) the

20  cottonmouth and (iii) and the copperhead snakes sanger data from Northamerica. We found that

21  SML is more efficient than ABC. It is generally more accurate and needs fewer simulations to

22   perform an inference. We found support for a divergence model without migration, with a recent

23   bottleneck for one of the populations of the southamerican frog. For the cottonmouth we found

24   support for divergence with migration and recent expansion and for the copperhead we found

25   support for a model of divergence with migration and recent bottleneck. Interestingly, by using

26   an SML method it was possible to achieve high accuracy in model selection even when several

27   models were compared in a single inference. We also found a higher accuracy when inferring

28   parameters with SML.

30   ## Introduction

31      The process of population divergence and speciation is finally being realized across

32   many non-model organisms with the use of genetic data and advanced statistical models.

33   Understanding population divergence involves testing diversification scenarios and estimating

34   historical parameters, such as divergence time, historical demography and migration rate

35   (Nielsen & Beaumont, 2009). Under simple diversification scenarios it is possible to use the

36   coalescent model (Kingman, 1982) with the likelihood function and MCMC to infer model

37   probabilities and associated historical parameters (Beerli & Palczewski, 2010; Bouckaert et al.,

38   2014; Gronau, Hubisz, Gulko, Danko, & Siepel, 2011; Hey, 2010; Yang & Rannala, 2010). There

39   is, however, an immense space of possible diversification scenarios where several hypotheses

40   may translate into complex, highly parameterized models that are difficult or impossible to solve

41   analytically (Fagundes et al., 2007; Mayr, 1942).

42      To overcome these limitations, researchers have used alternative approaches to

43   approximate posterior probabilities or marginal likelihoods of population parameters by reducing

44     data to summary statistics (Beichman, Huerta-Sanchez, & Lohmueller, 2018). These summary

45     statistics can be used in approximate Bayesian computation (ABC) and Supervised Machine

46     Learning (SML) to test hypotheses in a flexible likelihood-free context. ABC uses simulations

47     generated from parameter values sampled from prior probabilities to infer posterior probabilities

48     by applying a rejection algorithm that discards all simulations where the distance to the

49     observed data falls above an arbitrary tolerance level (Beaumont, 2011; Csilléry, Blum,

50     Gaggiotti, & François, 2010). Alternatively, simulated summary statistics can be used in SML as

51     training data (Schrider & Kern, 2018). For the simulated data, the link between population

52     parameters and summary statistics is known, so the algorithm can learn this connection and

53     infer model probability and parameter values for observed summary statistics (Burbrink &

54     Gehara, 2018; Sheehan & Song, 2016). To perform ABC and SML, end-users need to create

55     custom scripts to sample parameters from prior distributions and pass them to a simulator. This

56     requires integration of many different packages in various languages and the user's ability to

57     control this workflow sets the limit on the testable diversity of scenarios and hypotheses.

58         ABC and SML algorithms were already implemented in different packages of the R

59     statistical platform (Csilléry, François, & Blum, 2012; Kuhn, 2008). However, there is currently no

60     R-package to generate simulations for simulation-based model inference. To fill that gap we

61     developed a new R-package, called *PipeMaster,* that can be used to build models, add prior

62     distribution to model parameters, and simulate coalescent data from these prior distributions.

63     *PipeMaster* can also calculate summary statistics for the empirical data to allow statistical

64     comparison between observed and simulated data.

65         Here we demonstrate the utility of our newly developed package for three empirical

66     examples and evaluate the power of ABC versus SML and the influence of the number of loci in

67     the accuracy of model inferences. In the first example, we tested 10 hypotheses of divergence

68  for the Muller's termite frog, *Dermatonotus muelleri,* using newly generated data of 2177 loci of

69  ultra-conserved elements (UCE). In the second and third examples, we tested six different

70  hypotheses for two species complexes of North American vipers, the cottonmouth and the

71  copperhead, using pre-existent multi-locus data (Burbrink & Guiher, 2015). We show that

72  *PipeMaster* can be used with other R-packages to perform model and parameter inference in a

73  single platform and to test complex diversification hypotheses to better understand the evolution

74  of organisms.


## 75  **Material and Methods**

76      The PipeMaster R-package is currently available on github

77  ([www.github.com/gehara/PipeMaster](www.github.com/gehara/PipeMaster)) and can be installed via the *install_github* function from

78  the devtools R-package. Below we describe the main features of the package and exemplify its

79  use for model and parameter inference using empirical data with Nexgen and Sanger

80  dimensions.


81  **The interactive menu**

82      *PipeMaster* has an interactive menu that allows the user to build models and set up

83  parameter priors. In addition, the *main.menu* function can take a *ms* simulator string (see

84  Hudson, 2002 for more information about ms) for model specification, which can be generated

85  interactively with the PopPlanner application (see Ewing, Reiff, & Jensen, 2015 for more

86  information about this application). Alternatively, the user can input a tree topology in newick

87  format as a backbone of a diversification model, thus generating a simple isolation model with

88  constant population size and divergence time parameters. This basic isolation model can be

89  modified by adding ancestral population size changes and migration parameters, or by removing

90 divergence parameters to simulate island models. The user can use the interactive menu to set

91 conditions for parameter sampling (e.g. Ne1 > Ne2: effective population size of population 1 is

92 larger than effective population size of population 2). In the current version, uniform and normal

93 prior distributions are allowed. When the user exits the menu, the model can be saved as an R

94 object. A previously generated model object can be used as a template for a different model

95 setup, eliminating the need to start from the beginning when generating a nested model.

96 Specific characteristics of the data regarding number of base pairs and samples per population

97 per gene can be obtained using the *get.data.structure* function. This function reads the

98 parameters of the observed data and replicates them in the model.


99 **The simulation functions work-flow**

100 *PipeMaster* uses *ms* (Hudson, 2002) as an internal R function, or *ms*ABC (Pavlidis &

101 Laurent, n.d.) as the essential source of simulation. The program *ms* simulates coalescent trees

102 under the Wright-Fisher model, and places segregating sites on these trees under the infinite

103 site model.

104 *PipeMaster* has three simulation functions for non-hierarchical models: i)

105 *sim.ms.sumstat*, used to simulate summary statistics optimized for Sanger-scale data; ii)

106 *sim.coaltrees*, to simulate coalescent trees; and iii) *sim.msABC.sumstat*, to simulate summary

107 statistics using the simulator *ms*ABC (Pavlidis, Laurent, & Stephan, 2010) as an external

108 program (**Figure 1a-c**). All functions take as input the model object generated by the *main.menu*

109 function. They have the same basic work-flow and are used to sample parameter values from

110 prior distributions, convert the values to coalescent scale, pass those values to a coalescent

111 simulator, and write the output in a text file. In the case of *sim.ms.sumstat,* the simulated data is

112 passed to PopGenome R-package (Pfeifer, Wittelsburger, Ramos-Onsins, & Lercher, 2014) for

113    summary statistics calculation and the entire simulation process is performed without calling any

114    external program.


**ABC and SML analyses**

116         We implemented two different simulation-based inference methods in this study,

117    approximate Bayesian computation (ABC) and supervised machine-learning (SML). In all

118    empirical examples, before proceeding with the inference, we evaluated model-fit by running a

119    PCA of simulated and observed data. For the ABC approach we used the *abc* R-package

120    (Csilléry et al., 2012). We performed an abc rejection using the *postpr* function to calculate

121    model probabilities by retaining 100 simulations with the closest distance from our observed

122    data. To calculate accuracy in model selection we used *cv4postpr* with 100 pseudoreplicates per

123    model and the same tolerance value. The final accuracy was calculated by dividing the total

124    number of correct classifications by the total number of pseudoreplicates.

125         For SML we used the simulated data to train a neural network with one hidden layer to

126    classify the data into different simulated scenarios using the *nnet* algorithm in *caret* R-package.

127    We preprocessed the summary statistics by centering and scaling the data. We used 75% of the

128    simulations as training data and the remaining 25% as testing data. To tune the parameters of

129    the neural network, such as number of nodes and decay value, we performed 10 bootstraps

130    with a maximum of 2,000 iterations in each learning replicate and retained the parameters

131    yielding the highest accuracy. After training and testing,  we used the neural network to classify

132    our observed summary stats.

133         To estimate parameters we used the *abc* function of the abc R-package with the

134    *neuralnet* regression method. Before proceeding with the estimation we simulated additional

135    data for the best selected model totalizing 1,000,000 data sets for the Sanger examples and

136    100,000 for the UCE example. The *abc* function first performs a rejection step, reducing the

137    dataset before neural network training. We evaluated tolerance and the accuracy of parameter

138    estimates using the *cv4abc* function with 100 replicates and two different tolerance values  for

139    the Sanger examples (0.01, 0.001) and three values for the UCE example(0.1, 0.01, 0.001) . We

140    then calculated the correlation (*r*) between true and estimated parameters for each tolerance

141    value. We selected the tolerance yielding the best correlations among parameters. All codes

142    used in the ABC and SML are available on github as part of a tutorial for the package

143    (github.com/gehara/PipeMaster).


144    **SML versus ABC and the influence of the number of loci in the accuracy of estimates**

145        To evaluate the influence of the number of loci and compare the performance of ABC

146    versus SML for estimating the true model, we ran a set of simulations experiments with four

147    treatments that varied in total number of loci (10, 100, 1000, 2177). We used the case study of

148    *Dermatonotus muelleri* below as an empirical basis for this experiment. Accordingly, simulation

149    parameters, models, priors and summary statistics were the same as simulated for *D. muelleri,*

150    while the different number of loci with their parameters (base pairs number of individuals per

151    population) were obtained by sub-sampling 10, 100, or 1000 loci from the total dataset of 2177

152    loci generated for *D. muelleri,* plus a fourth treatment that contained the entire dataset. In each

153    treatment we ran ABC and SML inferences for a group of pseudo-observed data (POD; i.e. test

154    data in machine-learning jargon). We repeated these calculations three times, varying the total

155    number of simulations per model (1,000, 10,000 and 100,000).

156        We also performed a simulation experiment based on the *D. muelleri* data to evaluate

157    the accuracy of parameter estimates under different number of loci. In this case we estimated

158    parameters for 100 POD under the IsBott2 model, which was the model with the highest

7

159 probability for *D. muelleri* (see details below). We simulated a total of 100,000 data sets to use

160 as reference data. To estimate parameters we used the *abc* function of the abc R-package with

161 the neural network regression. We retained 1000 simulations after the rejection step and used

162 these to train a neural network. We then calculated the correlation (*r*) between true value and

163 estimated value for each parameter. An *r* closer to 1 would indicate a lower error in parameter

164 estimates. We performed this calculation for the same treatments of 10, 100, 1000 and 2177 loci

165 and we tested different retention or tolerance values.


166 **Application with UCE data - testing diversification hypotheses for muller's termite frog**

167 As an empirical example we generated a dataset of 2177 loci of ultra-conserved

168 elements (UCE) for the neotropical frog, *Dermatonotus muelleri* (see details of molecular

169 protocol in the Supplementary methods)*.* This species is distributed along the dry diagonal of

170 open formations which separates the Amazon from the Atlantic Forest. It is an explosive

171 breeder, highly adapted to seasonal environments with pronounced periods of drought (Nomura,

172 Rossa-Feres, & Langeani, 2009). A previous study using three loci (Oliveira et al., 2018) found

173 that *D. muelleri* is composed by two deeply divergent populations, one distributed in the

174 Caatinga and north of Cerrado, and a second one distributed in the southwest part of Cerrado

175 (**Figure 2a**). Here we took a subsample of 88 individuals used in that study. After data assembly

176 population assignment tests (see Supplementary methods) confirmed the existence of  two

177 spatially structured clusters (Oliveira et al., 2018).

178 The geographic break separating these two populations falls in an area of high elevation,

179 which may have isolated the populations. Also, Pleistocene climatic cycles are expected to have

180 influenced the demographic history of at least the Northeast population (Gehara et al., 2017;

181 Oliveira et al., 2018). Oliveira et al. (2018) found support for a model of diversification without

migration and expansion only for the Northeastern population. To challenge these findings and

test alternative diversification hypotheses for *D. muelleri,* we tested 10 two-population models:

(i) a pure isolation scenario without migration and without demographic change (Is); (ii) an

isolation with migration scenario without demographic change (IM); (iii) an isolation with recent

expansion and no migration (IsExp); (iv) an isolation with migration and recent expansion

(IMExp); (v) isolation with recent bottleneck and expansion (IsBott); (vi) isolation with migration,

recent bottleneck and expansion (IMBott); (vii) isolation with recent expansion only for the

Northeastern population (IsExp2); (viii) isolation with migration with expansion only for the

Northeastern population (IMExp2); (ix) an isolation with bottleneck only for the Northeastern

population (IsBott2); (x) an isolation with migration scenario with a bottleneck only for the

Northeastern population (IMBott2) (**Figure 3**). Priors of population sizes and time of

demographic events were retrieved from Oliveira et al. (2018) and can be found in the

**Supplementary Table 1**.

We simulated 100,000 data sets of 38 summary statistics (see **Supplementary**

**Methods** and tutorial: github.com/gehara/PipeMaster) per model with *sim.msABC.sumstat*

function. We used two independent approaches for model inference, ABC and SML described

above.


**Application with Sanger data - testing diversification hypotheses for Copperhead and**

**Cottonmouth pit vipers**

We also performed a model selection for two species complexes of vipers widely

distributed in Eastern North America: the *Agkistrodon contortrix* complex (Copperheads), and

the *Agkistrodon piscivorus* complex (Cottonmouths). The dataset used contain one

mitochondrial and five nuclear loci.

205        The *A. contortrix* species complex comprises two species, *A. contortrix* and *A.*

206    *laticinctus,* which together cover a large portion of eastern and central United States*.*

207    *Agkistrodon contortrix* is associated with deciduous hardwoods and pine forests and has a wider

208    distribution in the Eastern and Midwestern US (**Figure 2b**). *Agkistrodon laticinctus* occurs in

209    drier grassland environments in the central US to the Trans-Pecos habitats of west Texas.

210    Diversification in this complex is likely ecological, since their contact zone falls in the transition

211    from forested habitats to grasslands (Burbrink & Guiher, 2015). Both species currently occur in

212    areas that were covered by ice sheet during the last glaciation and show genetic signs of

213    population expansion in the Pleistocene (Guiher & Burbrink, 2008).

214        The *A. piscivorus* is also composed of two species. One of them, *A. conanti,* is mainly

215    restricted to the Florida Peninsula. The other, *A. piscivorus,* is distributed north of the peninsula

216    up to southern Illinois and Indiana in the north, Eastern Texas in the west, and coastal North

217    Carolina in the east (**Figure 2c**). The contact zone of these two species in the Florida peninsula

218    represents a common phylogeographic break for several other organisms (Burbrink, Fontanella,

219    Alexander Pyron, Guiher, & Jimenez, 2008; Krysko, Nuñez, Lippi, Smith, & Granatosky, 2016;

220    Mckelvy & Burbrink, 2017; Soltis, Morris, McLachlan, Manos, & Soltis, 2006) and the

221    diversification of the complex was also likely influenced by the climatic cycles of the Quaternary

222    (Guiher & Burbrink, 2008).

223        Taking these aspects into account, we tested for both species complexes, six

224    diversification hypotheses (**Figure 3**).  We generated the six models (a subset of the models

225    simulated for the frog example above; see **Figure 3**) and simulated 100,000 datasets for each

226    model using the *sim.ms.sumstat* function of PipeMaster R-package. We used wide uniform prior

227    distributions according to Burbrink and Guiher (2015) (see parameter list and priors in the

228 supplementary material). We used a set of 17 summary statistics (see **Supplementary**

229 **Methods** and tutorial: github.com/gehara/PipeMaster)

230      For both species complexes and both methods used (ABC and SML) we compared the

231 models hierarchically. (i) first we compared all the Isolation models with each corresponding

232 version that included migration (e.g. IsD against IMD; IsBott against IMBott). (ii) Than we took

233 the best models resulting from the first comparisons and conducted a second comparison to find

234 the best model of all.


235 # Results

236 **SML versus ABC**

237      The simulation experiment shows a higher error in model selection when using ABC

238 relative to SML (**Figures 4 and 5**). The number of loci has a strong influence in the accuracy of

239 model inferences. The dataset with 2177 loci had highest accuracy while the 10 locus dataset

240 had the lowest. The number of simulations also influence accuracy with inferences performed

241 with a reference dataset of 1,000 simulations per model having the lower true model

242 probabilities (**Figure 4**), while the inferences performed with 100,000 simulations per model has

243 the highest, particularly for ABC. For the SML inference both reference datasets of 10,000 and

244 100,000 simulations per model yielded nearly identical accuracies. The number of loci also has

245 influence in parameter estimates. SML had higher precision when compared to ABC (**Figure 5**).

246 The number of retained simulations, the tolerance value, influences ABC and SML in different

247 ways. For ABC retaining a low number of simulations yielded higher $R$. For SML retaining more

248 simulations result in better algorithm training.


249 **Diversification of muller's termite frog**

250    Simulations presented a good fit to the data as shown by the PCA plots

251    (**Supplementary Figure 1**). The trained neural network is able to differentiate and classify the

252    10 models with an accuracy of 0.879 while the ABC had an accuracy of 0.83. Using the SML

253    approach the observed data was classified as the IsBott2 model with a probability higher than

254    0.99 (**Table 1**), where only the northeast population experienced a bottleneck with expansion.

255    The ABC inference suggest a different model, IMexp, where the two populations expand after

256    divergence. In this case, the probability of the model was considerably lower, 0.49. Because the

257    accuracy of the SML is higher, we consider the IsBott2 as the best model.

258    The divergence time can be estimated with high accuracy and suggest a split around 2.6

259    Ma between the two populations (**Table 2**). Estimated current sizes for population 1 suggest a

260    very large population after expansion but accuracy of this estimate is low**.** Estimates for

261    population 2 are more accurate (**Table 2**). The average estimated mutation rate was

262    2.2E-10/site/generation with an estimated standard deviation of 3.88E-10/site/generation.


263    **Diversification of Copperhead and Cottonmouth pit vipers**

264    For both species complexes, the simulated models had a good fit to the data, as

265    suggested by the  PCA (**Supplementary Figure 2**). In the first comparisons (1, 2 and 3; see

266    **Table 3**) for the *A. contortrix* complex, the accuracy varied from 0.79 – 0.85 for the SML  and

267    from 0.76 – 0.86 for the ABC. For comparison 1 (Is vs IM), ABC and SML showed conflicting

268    results, with the pure isolation model, Is, having the highest probability for the ABC and the

269    isolation with migration model, IM, having the highest probability in the SML. For comparisons 2

270    and 3, the two methods showed concordant results; models that included migration had higher

271    probabilities than the correspondent models without migration (**Table 3**). The final comparison

272    accuracies of ABC and SML were 0.78 and 0.79 respectively. Both methods converged in the

273    same best model for the diversification of *A. contortrix* complex, IMBott (**Table 3**). For all

274  comparisons, the SML showed higher probabilities for the selected model when compared to

275  ABC (**Table 3**).

276      In the first comparisons (1, 2 and 3; see **Table 3**) for the *A. piscivorus* complex, the

277  accuracy varied from 0.92 – 0.94 for the SML and from 0.89 – 0.93 for the ABC,  and the best

278  selected model were the same as the ones inferred for the *A. contortrix* complex (**Table 3**). In

279  the final comparison, the accuracy of the ABC was higher than the SML, 0.87 and 0.79

280  respectively. However, both methods suggest high probabilities for the same model, the IMexp,

281  which is an isolation with migration with expansion for both species (**Table 3**).

282      The cross-validation for the parameter estimates suggest low correlation between

283  estimated and true values, particularly for *A. contortrix* (**Table 4**), suggesting high uncertainty in

284  estimates. In general, the parameters that can be estimated with higher confidence are the

285  current population sizes (**Table 4**).


286  ## Discussion

287      Our simulation experiment showed that supervised machine-learning outperforms

288  approximate Bayesian computation. This is particularly evident for datasets with genomic

289  dimentions, which is the current standard of molecular studies for non-model organisms. We

290  also show that much higher accuracies can be obtained with a SML as opposed to ABC, even

291  when using just 100 loci and a considerably low number of simulations per model (10,000).

292  Thus, because ABC requires a larger amount of simulations, it is more time consuming and less

293  efficient when compared to SML.

294      Our simulation experiment also show that the model parameters can be estimated with

295  higher accuracy with the increase in the number of loci. The SML approach also outperforms

296  ABC for parameter estimates (**Figure 4**). Some parameters, like current effective population

297    size and time of divergence, can be estimated with higher accuracy. However, ancestral

298    population sizes are harder to estimate (**Figure 5**; **Table 2**). Interestingly, posterior distributions

299    of the average and standard deviation of the mutation rate across all loci can be obtained with

300    high confidence, allowing a more relaxed assumption when compared to using a fixed mutation

301    rate for all loci.


302    **Diversification of muller's termite frog**

303          We found support for an isolation model with population contraction with expansion for

304    the northeast population. This partially agrees with Oliveira et al. (2018), who found support for

305    recent expansion without a contraction. Oliveira et al. (2018) analyzed only three loci while we

306    analyzed more than 2,000, thus our data certainly contains more information about historical

307    demography (e.g. Gill et al., 2013).

308          The inference of a population contraction in the northeast population reinforces the idea

309    of dynamic landscape changes in the northeast of Brazil along the Pleistocene. Currently, this

310    area is predominantly covered by the Caatinga semiarid environment, but many studies suggest

311    periods of increase in humidity in the last 1 Ma (Auler et al., 2004b; Cheng et al., 2013).

312    Travertine deposits suggest a long period of increase in humidity from approximately 460 to 330

313    K years (Auler et al., 2004a) which remarkably agrees with our estimated time for the reduction

314    in population size (mode: 337 Ky, CI: 195 – 437 Ky). These humid phases in the northeast of

315    Brazil may have allowed long distance dispersals between Amazon and Atlantic forest fauna

316    (Dal Vechio, Prates, Grazziotin, Zaher, & Rodrigues, 2018; Prates, Rivera, Rodrigues, &

317    Carnaval, 2016). The reduction in population size is followed by a population expansion starting

318    at around 230 K years (CI: 132 – 362 K years), in agreement of other studies that find

319    synchronous population expansion Caatinga's herpetofauna (Gehara et al., 2017).

320    The estimated divergence time at 2.6 Ma is considerably younger than previous

321    estimates (~4 Ma; see Oliveira et al., 2018). Our estimates places the divergence between the

322    northeast and southwest populations in the Pliocene-Pleistocene transition, after the mid

323    Pliocene warm period, when the average global temperature was 2 – 3º C higher than today.

324    This higher temperature may have allowed *D. muelleri,* a lowland species*,* to inhabit the

325    highlands of the Brazilian plateau. With the temperature cooling the highland climate may have

326    become unsuitable for the species and the Brazilian plateau became a vicariant barrier causing

327    diversification.


328    **Diversification of Copperhead and Cottonmouth pit vipers**

329    For both species complexes, we found support for demographic change and gene flow

330    between species pairs. For the *A. contortrix* complex, we found support for a reduction in

331    population size with subsequent expansion in the late Pleistocene. This species complex is

332    currently found in areas that were covered by ice sheets during glaciations. Accordingly, the

333    glaciation cycles would have restricted the distribution of the species to southern refugia,

334    causing a population contraction (Burbrink et al., 2016; Marshall, James, & Clarke, 2002). In

335    interglacial periods, the species would expand their range and their population sizes. It is also

336    possible that the climatic cycles influenced their divergence, driving speciation by the isolation of

337    populations in distinct refugia. Nevertheless, the presence of gene flow indicates that if isolation

338    happened during glaciations, they were likely followed by periods of contact. Gene flow may

339    also indicate the role of climatic gradients in diversification. *Agkistrodon contortrix* and *A.*

340    *latiscinctus* occur in distinctly different niches (Burbrink & Guiher, 2015; Gloyd & Conant, 1990)

341    and they likely present physiological adaptations to these different environments. Thus, hybrids

342    may have lower fitness when compared to non-hybrids (Gow, Peichel, & Taylor, 2007). Future

343  studies using thousands of loci will have the opportunity to test for selection across the climatic

344  gradients, and may shed more light on the evolution of the *A. contortrix* species complex.

345       For the *A. piscivorus* complex, we found no support for a bottleneck during the

346  Pleistocene. The most probable model suggests an isolation with gene flow and a recent

347  population expansion. Both *A. piscivorus* and *A. conanti* are mostly distributed in areas free from

348  broadscale effects of Pleistocene glaciation (Marshall et al., 2002). Accordingly, the supported

349  model suggests a relatively more stable population size, with recent population expansion for

350  both species. The contact zone between the species is in the northern area of the Florida

351  Peninsula. This region was isolated from the continent when sea levels were higher, so it is

352  likely that the diversification of the complex was influenced by sea level rise, which could have

353  isolated *A. conanti* in a continental island formed by part of the landmass that today represents

354  the Florida Peninsula (Hine, 2013; Krysko et al., 2016). In this scenario, gene flow between *A.*

355  *conanti* and *A. piscivorus* was favored during glacial periods when sea levels were low, while

356  isolation happened during interglacial periods while sea levels were high.


357  ## Conclusion

358       We demonstrated the use of coalescent simulations generated by our newly developed

359  R-package to infer the probability of complex diversification models in three different non-model

360  organisms. In the three cases, we were able to test relatively complex demographic models with

361  population size change, population structure and migration that are difficult, time consuming or

362  impossible to implement using a full Bayesian or likelihood approaches. Interestingly, by using a

363  SML method it was possible to achieve high accuracy in model selection even when several

364  models were compared in a single inference (**Table 1**).

365        Machine-learning algorithms are becoming increasingly available to the general scientific

366    community through R and Python applications, facilitating its use for an unprecedented number

367    of cases in evolutionary biology and ecology. Here we demonstrated its use comparing it with a

368    more traditional, ABC, for model inference in population genetics. Our results agree with the

369    recent literature (Schrider & Kern, 2018; Sheehan & Song, 2016) supporting the power of SML

370    in dealing with complex multi-dimensional problems such as the ones presented here.

375    **Author Contributions**

376        MG and FB conceived the ideas; MG, GGM and FB designed methodology; FB and MG

377    collected the data; MG analyzed the data; MG and FB led the writing of the manuscript. All

378    authors contributed critically to the drafts and gave final approval for publication.

379    **Data Availability**

380        All codes used in the ABC and SML analyses are found in

381    github.com/gehara/PipeMaster. The assembled UCE data is available in the Dryad (upon

382    manuscript acceptance).

383    **References**

384

Auler, A. S., Wang, X., Edwards, R. L., Cheng, H., Cristalli, P. S., Smart, P. L., & Richards, D. A.
385        (2004a). Palaeoenvironments in semi-arid northeastern Brazil inferred from high precision

mass spectrometric speleothem and travertine ages and the dynamics of South American rainforests. *Speleogenesis and Evolution of Karst Aquifers*, *2*(2), 1–4.

Auler, A. S., Wang, X., Edwards, R. L., Cheng, H., Cristalli, P. S., Smart, P. L., & Richards, D. A. (2004b). Quaternary ecological and geomorphic changes associated with rainfall events in presently semi-arid northeastern Brazil. *Journal of Quaternary Science*, *19*(7), 693–701.

Beaumont, M. A. (2011). Approximate Bayesian Computation in Evolution and Ecology. *Annual Review of Ecology, Evolution, and Systematics*, *41*(1), 379–406.

Beerli, P., & Palczewski, M. (2010). Unified framework to evaluate panmixia and migration direction among multiple sampling locations. *Genetics*, *185*(1), 313–326.

Beichman, A. C., Huerta-Sanchez, E., & Lohmueller, K. E. (2018). Using Genomic Data to Infer Historic Population Dynamics of Nonmodel Organisms. *Annual Review of Ecology, Evolution, and Systematics*, *49*(1), annurev – ecolsys – 110617–062431.

Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.-H., Xie, D., … Drummond, A. J. (2014). BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Computational Biology*, *10*(4), e1003537.

Burbrink, F. T., Chan, Y. L., Myers, E. A., Ruane, S., Smith, B. T., Hickerson, M. J., & Sgro, C. (2016). Asynchronous demographic responses to Pleistocene climate change in Eastern Nearctic vertebrates. *Ecology Letters*, Vol. 19, pp. 1457–1467. doi: 10.1111/ele.12695

Burbrink, F. T., Fontanella, F., Alexander Pyron, R., Guiher, T. J., & Jimenez, C. (2008). Phylogeography across a continent: The evolutionary and demographic history of the North American racer (Serpentes: Colubridae: Coluber constrictor). *Molecular Phylogenetics and Evolution*, *47*(1), 274–288.

Burbrink, F. T., & Gehara, M. (2018). The Biogeography of deep time reticulation. *Systematic Biology*, *67*(5), 743–744.

Burbrink, F. T., & Guiher, T. J. (2015). Considering gene flow when using coalescent methods to delimit lineages of North American pitvipers of the genus Agkistrodon. *Zoological Journal of the Linnean Society*, *173*(2), 505–526.

Cheng, H., Sinha, A., Cruz, F. W., Wang, X., Edwards, R. L., D'Horta, F. M., … Auler, A. S. (2013). Climate change patterns in Amazonia and biodiversity. *Nature Communications*, *4*, 1411.

Csilléry, K., Blum, M. G. B., Gaggiotti, O. E., & François, O. (2010). Approximate Bayesian Computation (ABC) in practice. *Trends in Ecology & Evolution*, *25*(7), 410–418.

Csilléry, K., François, O., & Blum, M. G. B. (2012). Abc: An R package for approximate Bayesian computation (ABC). *Methods in Ecology and Evolution / British Ecological Society*, *3*(3), 475–479. Retrieved from arXiv.

Dal Vechio, F., Prates, I., Grazziotin, F. G., Zaher, H., & Rodrigues, M. T. (2018). Phylogeography and historical demography of the arboreal pit viper Bothrops bilineatus (Serpentes, Crotalinae) reveal multiple connections between Amazonian and Atlantic rain forests. *Journal of Biogeography*, *45*(10), 2415–2426.

Ewing, G. B., Reiff, P. A., & Jensen, J. D. (2015). *PopPlanner : visually constructing demographic models for simulation*. *6*(April), 1–4.

Fagundes, N. J. R., Ray, N., Beaumont, M., Neuenschwander, S., Salzano, F. M., Bonatto, S. L., & Excoffier, L. (2007). Statistical evaluation of alternative models of human evolution. *Proceedings of the National Academy of Sciences of the United States of America*, *104*(45), 17614–17619.

Gehara, M., Garda, A. A., Werneck, F. P., Oliveira, E. F., da Fonseca, E. M., Camurugi, F., … Burbrink, F. T. (2017). Estimating synchronous demographic changes across populations using hABC and its application for a herpetological community from northeastern Brazil.

434    *Molecular Ecology, 26*(May), 4756–4771.

435    Gill, M. S., Lemey, P., Faria, N. R., Rambaut, A., Shapiro, B., & Suchard, M. A. (2013).
436    Improving bayesian population dynamics inference: A coalescent-based model for multiple
437    loci. *Molecular Biology and Evolution, 30*(3), 713–724.

438    Gloyd, H. K., & Conant, R. (1990). Snakes of the Agkistrodon complex: a monographic review.
439    *Contributions to Herpetology*, 1–614. Oxford, Ohio: Society for the Study of Amphibians and
440    Reptiles.

441    Gow, J. L., Peichel, C. L., & Taylor, E. B. (2007). Ecological selection against hybrids in natural
442    populations of sympatric threespine sticklebacks. *Journal of Evolutionary Biology, 20*(6),
443    2173–2180.

444    Gronau, I., Hubisz, M. J., Gulko, B., Danko, C. G., & Siepel, A. (2011). Bayesian inference of
445    ancient human demography from individual genome sequences. *Nature Genetics, 43*(10),
446    1031–1034.

447    Guiher, T. J., & Burbrink, F. T. (2008). Demographic and phylogeographic histories of two
448    venomous North American snakes of the genus Agkistrodon. *Molecular Phylogenetics and
449    Evolution, 48*(2), 543–553.

450    Hey, J. (2010). Isolation with migration models for more than two populations. *Molecular Biology
451    and Evolution, 27*(4), 905–920.

452    Hine, A. C. (2013). *Geologic history of Florida: major events that formed the Sunshine State*.
453    Gainesville, USA: University Press of Florida.

454    Hudson, R. (2002). Ms a program for generating samples under neutral models. *Bioinformatics* ,
455    (2002), 337–338.

456    Kingman, J. F. C. (1982). The coalescent. *Stochastic Processes and Their Applications*, Vol. 13,
457    pp. 235–248. doi: 10.1016/0304-4149(82)90011-4

458    Krysko, K. L., Nuñez, L. P., Lippi, C. A., Smith, D. J., & Granatosky, M. C. (2016). Molecular
459    Phylogenetics and Evolution Pliocene – Pleistocene lineage diversifications in the Eastern
460    Indigo Snake ( Drymarchon couperi ) in the Southeastern United States q. *Molecular
461    Phylogenetics and Evolution, 98*, 111–122.

462    Kuhn, M. (2008). Building Predictive Models in R Using the caret Package. *Journal of Statistical
463    Software, 28*(5), 1–26. Retrieved from arXiv.

464    Marshall, S. J., James, T. S., & Clarke, G. K. C. (2002). North American Ice Sheet
465    reconstructions at the last glacial maximum. *Quaternary Science Reviews, 21*(1-3),
466    175–192.

467    Mayr, E. (1942). *Systematics and the Origin of Species, from the Viewpoint of a Zoologist*.
468    Harvard University Press.

469    Mckelvy, A. D., & Burbrink, F. T. (2017). Molecular Phylogenetics and Evolution Ecological
470    divergence in the yellow-bellied kingsnake ( Lampropeltis calligaster ) at two North
471    American biodiversity hotspots. *Molecular Phylogenetics and Evolution, 106*, 61–72.

472    Nielsen, R., & Beaumont, M. A. (2009). Statistical inferences in phylogeography. *Molecular
473    Ecology*, Vol. 18, pp. 1034–1047. doi: 10.1111/j.1365-294x.2008.04059.x

474    Nomura, F., Rossa-Feres, D. C., & Langeani, F. (2009). Burrowing behavior of Dermatonotus
475    muelleri (Anura, Microhylidae) with reference to the origin of the burrowing behavior of
476    Anura. *Journal of Ethology, 27*(1), 195–201.

477    Oliveira, E. F., Gehara, M., São-Pedro, V. A., Costa, G. C., Burbrink, F. T., Colli, G. R., … Garda,
478    A. A. (2018). Phylogeography of Muller's termite frog suggests the vicariant role of the
479    Central Brazilian Plateau. *Journal of Biogeography, 45*(11), 2508–2519.

480    Pavlidis, P., & Laurent, S. (n.d.). *msABC : A modification of Hudson ' s ms to facilitate
481    multi-locus ABC analysis User ' s Guide & Manual Table of Contents*. doi:

482   10.1111/j.1755-0998.2010.02832.x

483 Pavlidis, P., Laurent, S., & Stephan, W. (2010). MsABC: A modification of Hudson's ms to
484   facilitate multi-locus ABC analysis. *Molecular Ecology Resources*, *10*(4), 723–727.

485 Pfeifer, B., Wittelsburger, U., Ramos-Onsins, S. E., & Lercher, M. J. (2014). PopGenome: An
486   efficient swiss army knife for population genomic analyses in R. *Molecular Biology and*
487   *Evolution*, *31*(7), 1929–1936.

488 Prates, I., Rivera, D., Rodrigues, M. T., & Carnaval, A. C. (2016). A mid-Pleistocene rainforest
489   corridor enabled synchronous invasions of the Atlantic Forest by Amazonian anole lizards.
490   *Molecular Ecology*, *25*(20), 5174–5186.

491 Schrider, D. R., & Kern, A. D. (2018). Supervised Machine Learning for Population Genetics: A
492   New Paradigm. *Trends in Genetics: TIG*, *xx*, 1–12.

493 Sheehan, S., & Song, Y. S. (2016). Deep Learning for Population Genetic Inference. *PLoS*
494   *Computational Biology*, *12*(3). doi: 10.1371/journal.pcbi.1004845

495 Soltis, D. E., Morris, A. B., McLachlan, J. S., Manos, P. S., & Soltis, P. S. (2006). Comparative
496   phylogeography of unglaciated eastern North America. *Molecular Ecology*, *15*, 4261–4293.

497 Yang, Z., & Rannala, B. (2010). Bayesian species delimitation using multilocus sequence data.
498   *Proceedings of the National Academy of Sciences of the United States of America*, *107*(20),
499   9264–9269.

## Tables

**Table 1:** Model probabilities and accuracies calculated with ABC and SML for the comparison of 10 simulated models for the frog *Dermatonotus muelleri* 2177 UCE data. (see **Figure 3** for a schematic representation of the models).

| Model | Probability SML | Probability ABC |
|---|---|---|
| IM | 0 | 0 |
| IMBott | 0 | 0.17 |
| IMBott2 | 0 | 0 |
| IMExp | 0 | **0.49** |
| IMExp2 | 0 | 0 |
| Is | 0 | 0 |
| IsBott | 0.0037 | 0.17 |
| **IsBott2** | **0.9963** | 0.13 |
| IsExp | 0 | 0 |
| IsExp2 | 0 | 0.04 |
| *Accuracy* | *0.8798* | *0.826* |

**Table 2:** Parameter priors, posterior estimates and correlation (*r*) result calculated with the cross-validation experiment for the frog species (UCE data). See **Supplementary Table 1** for a complete list of priors and parameters.

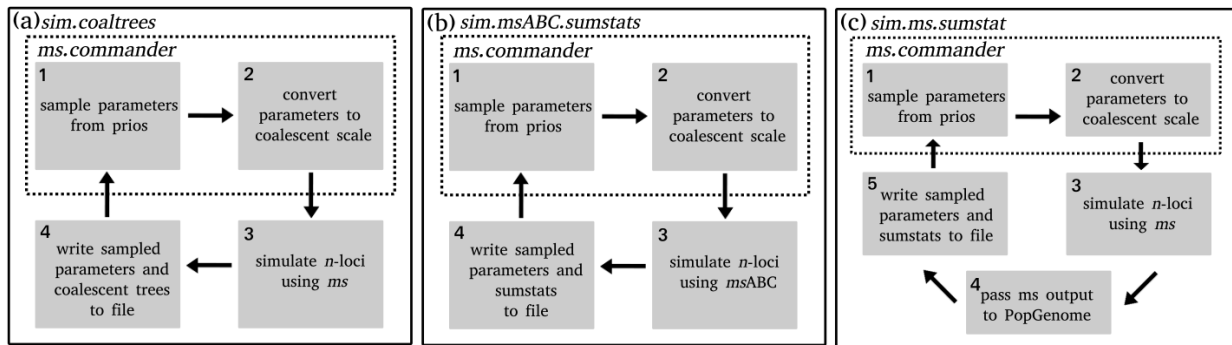| Paramater | Prior (min – max) | 2.50% | Median | Mean | Mode | 97.50% | *r* |
|---|---|---|---|---|---|---|---|
| Ne0.pop1 | 100,000 – 5,000,000 | 2,450,945 | 3,756,958 | 3,845,324 | 3,543,630 | 5,538,957 | 0.58 |
| Ne0.pop2 | 100,000 – 5,000,000 | -347,089 | 1,012,398 | 1,007,289 | 993,517 | 2,270,467 | 0.82 |
| Ne1.pop1 | 1,000 – 50,000 | 5,275 | 24,930 | 24,360 | 31,128 | 39,783 | 0.55 |
| Ne2.pop1 | 50,000 – 5,000,000 | 1,784,442 | 3,633,865 | 3,728,091 | 3,342,328 | 5,908,761 | 0.2 |
| Ne1.pop2 | 50,000 – 5,000,000 | -727,952 | 123,083 | 133,108 | 132,045 | 1,166,070 | 0.82 |
| join1 | 500,000 – 8,000,000 | 1,139,788 | 2,615,899 | 2,600,426 | 2,620,037 | 4,173,138 | 0.84 |
| t.Ne1.pop1 | 20,000 – 500,000 | 132,078 | 238,909 | 241,542 | 233,411 | 362,085 | 0.64 |
| t.Ne2.pop1 | 20,000 – 500,000 | 195,084 | 324,834 | 321,999 | 336,868 | 427,928 | 0.51 |
| t.Ne1.pop2 | 500,000 – 8,000,000 | 1,002,084 | 2,470,257 | 2,453,955 | 2,464,068 | 3,996,727 | 0.84 |
| mean.rate | 1E-11 – 1E-9 | 6.50E-11 | 2.28E-10 | 2.38E-10 | 2.27E-10 | 5.04E-10 | 0.77 |
| sd.rate | 1E-11 – 1E-9 | 2.71E-10 | 3.99E-10 | 4.12E-10 | 3.88E-10 | 6.50E-10 | 0.79 |

**Table 3:** Model probabilities estimated with ABC and SML with respective accuracies of estimates for the two snake species complex. Models were compared hierarchically, first comparisons 1, 2 and 3 were carried out independently. The final comparison included the best models of comparison 1, 2 and 3. Bold probabilities indicate the selected model for each comparison (see Figure 3 for a schematic representation of the models).

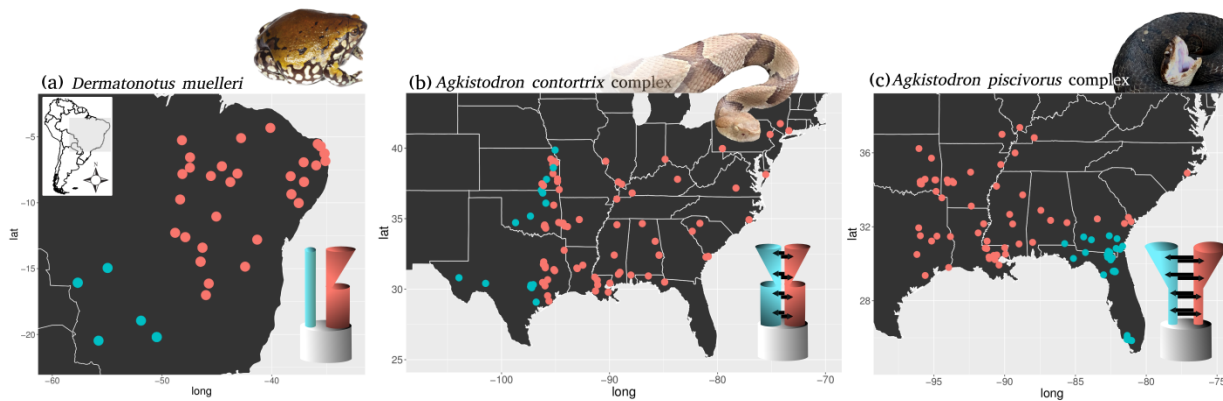| | *Agkistrodon piscivorus* | | | |
| --- | --- | --- | --- | --- |
| | **1** | **2** | **3** | **Final** |
| | **Is vs IM (Accuracy)** | **IsExp vs IMExp (Accuracy)** | **IsBott vs IMBott (Accuracy)** | **Best 1 vs Best 2 vs Best 3 (Accuracy)** |
| SML | 0.11 / **0.89** (0.94) | 0.01 / **0.99** (0.94) | 0.04 / **0.96** (0.92) | 0.01 / **0.85** / 0.14 (0.79) |
| ABC | **0.78** / 0.23 (0.93) | 0.19 / **0.82** (0.91) | 0.49 / **0.51** (0.89) | 0.13 / **0.61** / 0.26 (0.87) |
| | | | | |
| | *Agkistrodon contortrix* | | | |
| | **1** | **2** | **3** | **Final** |
| | **Is vs IM (Accuracy)** | **IsExp vs IMExp (Accuracy)** | **IsBott vs IMBott (Accuracy)** | **Best 1 vs Best 2 vs Best 3 (Accuracy)** |
| SML | 0.03 / **0.97** (0.79) | 0.00 / **1.00** (0.85) | 0.00 / **1.00** (0.79) | 0.12 / 0.01 / **0.87** (0.79) |
| ABC | **0.59** / 0.42 (0.76) | 0.08 / **0.92** (0.86) | 0.18 / **0.82** (0.77) | 0.33 / 0.01 / **0.66** (0.78) |

**Table 4:** Parameter priors, posterior estimates and R result calculated with the cross-validation experiment for the two snake species complexes (Sanger data). See **Supplementary Table 1** for a complete list of priors and parameters.

| | *Agkistrodon contortrix* | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| **Paramater** | **Prior (min – max)** | **2.50%** | **Median** | **Mean** | **Mode** | **97.50%** | *r* |
| Ne0.pop1 | 20,000 – 1,000,000 | 26,112 | 123,258 | 140,965 | 77,843 | 325,493 | 0.45 |
| Ne0.pop2 | 20,000 – 1,000,000 | 92,136 | 491,649 | 506,789 | 202,242 | 988,100 | 0.19 |
| Ne1.pop1 | 1,000 – 10,000 | 1,360 | 7,030 | 6,637 | 8,902 | 9,924 | 0.06 |
| Ne2.pop1 | 20,000 – 1,000,000 | 5,895 | 556,886 | 535,324 | 883,327 | 967,817 | 0.01 |
| Ne1.pop2 | 1,000 – 10,000 | 1,809 | 6,932 | 6,612 | 9,094 | 9,910 | 0.12 |
| Ne2.pop2 | 20,000 – 1,000,000 | 303,312 | 710,668 | 691,052 | 850,608 | 983,683 | 0.34 |
| join1 | 60,000 – 3,000,000 | 285,633 | 1,625,645 | 1,611,305 | 2,039,677 | 2,874,342 | 0.06 |
| t.Ne1.pop1 | 9,000 – 300,000 | 31,942 | 127,478 | 139,292 | 99,724 | 289,616 | 0.46 |
| t.Ne2.pop1 | 9,000 – 300,000 | 73,265 | 191,863 | 198,353 | 148,644 | 346,239 | 0.38 |

22

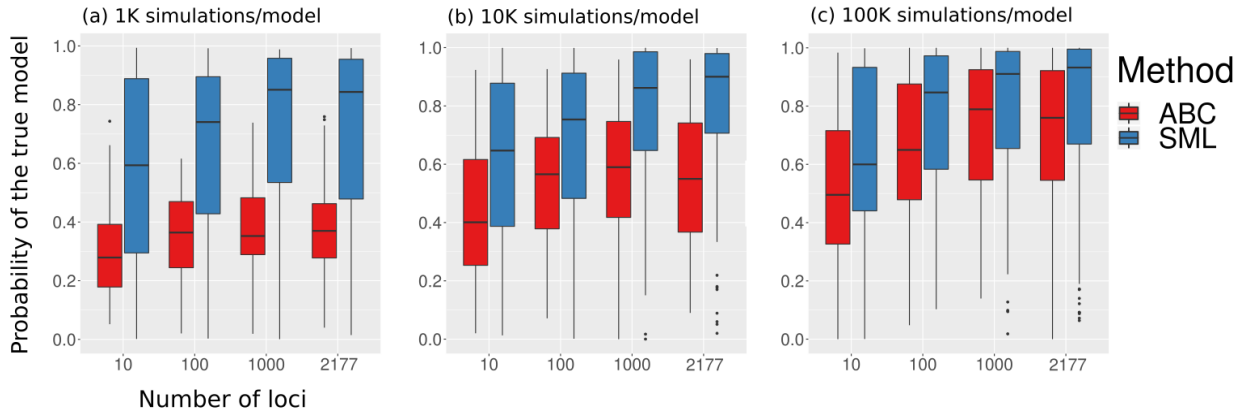| Paramater | Prior (min – max) | 2.50% | Median | Mean | Mode | 97.50% | r |
|---|---|---|---|---|---|---|---|
| t.Ne1.pop2 | 9,000 – 300,000 | 36,938 | 132,674 | 140,750 | 106,220 | 271,105 | 0.28 |
| t.Ne2.pop2 | 9,000 – 300,000 | 74,227 | 198,593 | 197,694 | 202,263 | 316,021 | 0.06 |
| mig0.1_2 | 0 – 2 | 0.76 | 1.5 | 1.47 | 1.75 | 2.06 | 0.18 |
| mig0.2_1 | 0 – 2 | 0.07 | 1.05 | 1.04 | 1.75 | 1.95 | 0.06 |
| | | | | | | | |
| | | | *Agkistrodon piscivorus* | | | | |
| **Paramater** | **Prior (min – max)** | **2.50%** | **Median** | **Mean** | **Mode** | **97.50%** | *r* |
| Ne0.pop1 | 10,000 – 500,000 | 104,779 | 258,656 | 269,769 | 198,668 | 475,458 | 0.4 |
| Ne0.pop2 | 10,000 – 500,000 | 59,975 | 216,254 | 236,734 | 130,139 | 475,953 | 0.47 |
| Ne1.pop1 | 1,000 – 10,000 | 3,534 | 41,156 | 42,827 | 32,328 | 89,432 | 0.42 |
| Ne1.pop2 | 1,000 – 10,000 | 8,510 | 50,965 | 50,489 | 54,790 | 91,127 | 0.49 |
| Ne2.pop2 | 10,000 – 1,000,000 | 163,818 | 552,602 | 556,057 | 564,831 | 956,823 | 0.05 |
| join1 | 9,9000 – 9,900,000 | 276,722 | 4,627,486 | 4,741,277 | 1,325,740 | 9,626,987 | 0.55 |
| t.Ne1.pop1 | 9,000 – 210,000 | 20,944 | 104,030 | 109,209 | 53,650 | 210,963 | 0.35 |
| t.Ne1.pop2 | 9,000 – 210,000 | -414 | 79,993 | 86,828 | 27,249 | 195,189 | 0.53 |
| t.Ne2.pop2 | 9,900 – 9,900,000 | 246,730 | 4,629,289 | 4,743,890 | 1,309,786 | 9,653,230 | 0.56 |
| mig0.1_2 | 0 – 2 | 0.19 | 0.87 | 0.93 | 0.72 | 1.91 | 0.28 |
| mig0.2_1 | 0 – 2 | 0.03 | 0.74 | 0.8 | 0.51 | 1.75 | 0.29 |

516 **Figure 1**: work-flow of the main simulation functions of PipeMaster and schematic
517 representation of the simulated models in the toy example. (a) work-flow of the *sim.ms.sumstat*
518 function; (b) schematic representation of the diversification models simulated in the toy example;
519 (c) work-flow of the *sim.coaltrees* function; (d) work-flow of the *sim.msABC.sumstat* function.
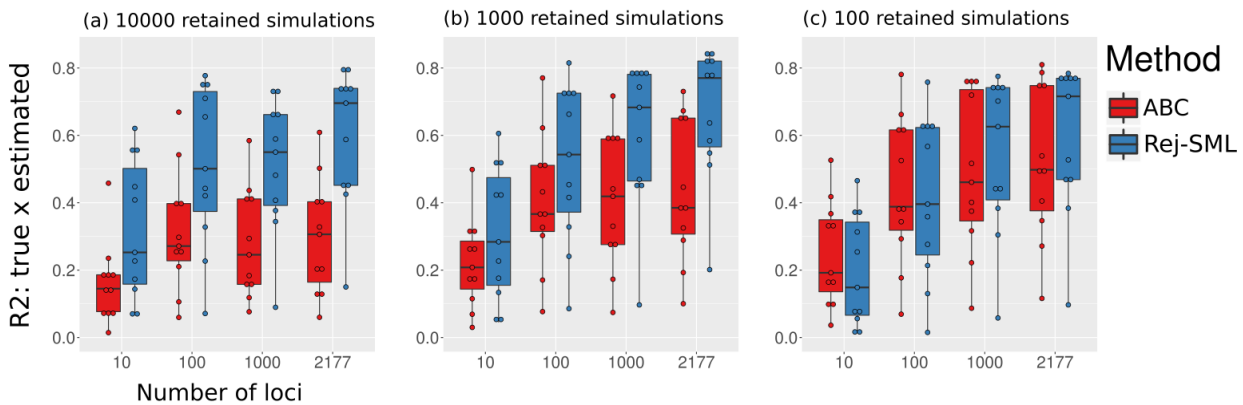


520 **Figure 2:** Distribution maps and best model for each data set analyzed in this study.

**Figure 3:** Schematic representation of the diversification models tested in the two *Agkistrodon* species complexes and *Dermatonotus muelleri*. Dotted line indicate the six models tested for *A. contortrix* and *A. piscivorus* complex. For *D. muelleri* we tested all 10 models. See **Supplementary Table 1** for a complete list of priors and parameters.

**Figure 4:** Results of the simulation experiment to compare the accuracy of ABC and SML for model inference in different conditions. The y-axis represents the probability of the true model, the x-axis represent different data dimensions. Each box plot represent probabilities of the true model for 100 pseudo observed data, 10 per model. For the ABC analysis, 100 simulations are retained in the rejection step, for the SML all simulations are used for algorithm training. (a) estimates performed with 1K simulations per model totalizing 10,000 simulations in the reference table. (b) estimates performed with 10K simulations per model totalizing 100,000 simulations in the reference table. (c) estimates performed with 100K simulations per model totalizing 1,000,000 simulations in the reference table.



**Figure 5:** Results of the simulation experiment to evaluate the influence of number of loci and tolerance values on parameter estimates of ABC and rejection with SML. The y-axis represents the correlation between estimated and true values for 100 pesudo-observed data for the 11 parameters of the model. (a) estimates are performed by retaining 10,000 closest simulations. (b) estimates are performed by retaining 1,000 closest simulations. (c) estimates are performed by retaining 100 closest simulations.