# Identifying model violations under the multispecies coalescent model using P2C2M.SNAPP

**Drew J Duckett** [1] , **Tara A Pelletier** [2] , **Bryan C Carstens** [Corresp. 3]

[1] Department of Evolution, Ecology, & Organismal Biology, The Ohio State University, Columbus, OH, United States

[2] Biology Department, Radford University, Radford, VA, United States

[3] Department of Evolution, Ecology and Organismal Biology, The Ohio State University, Columbus, OH, United States

Corresponding Author: Bryan C Carstens
Email address: carstens.12@osu.edu

Phylogenetic estimation under the multispecies coalescent model (MSCM) assumes all incongruence among loci is caused by incomplete lineage sorting. Therefore, applying the MSCM to datasets that contain incongruence that is caused by other processes, such as gene flow, can lead to biased phylogeny estimates. To identify possible bias when using the MSCM, we present P2C2M.SNAPP. P2C2M.SNAPP is an R package that identifies model violations using posterior predictive simulation. P2C2M.SNAPP uses the posterior distribution of species trees output by the software package SNAPP to simulate posterior predictive datasets under the MSCM, and then uses summary statistics to compare either the empirical data or the posterior distribution to the posterior predictive distribution to identify model violations. In simulation testing, P2C2M.SNAPP correctly classified up to 83% of datasets (depending on the summary statistic used) as to whether or not they violated the MSCM model. P2C2M.SNAPP represents a user-friendly way for researchers to perform posterior predictive model checks when using the popular SNAPP phylogenetic estimation program. It is freely available as an R package, along with additional program details and tutorials.

1

# Identifying model violations under the Multispecies Coalescent Model using P2C2M.SNAPP

4

5

6

7  Drew J. Duckett[1], Tara A. Pelletier[2], Bryan C. Carstens[1*]
8  [1]Department of Evolution, Ecology, & Organismal Biology, The Ohio State University,
9  Columbus, OH, United States of America
10 [2]Biology Department, Radford University, Radford, VA, United States of America
11
12 Corresponding Author:
13 Bryan C. Carstens
14 Department of Evolution, Ecology, & Organismal Biology, Museum of Biological Diversity,
15 The Ohio State University, 1315 Kinnear Rd., Columbus, OH 43210, United States of America
16 Email address: carstens.12@osu.edu

17

## Abstract

19 Phylogenetic estimation under the multispecies coalescent model (MSCM) assumes all
20 incongruence among loci is caused by incomplete lineage sorting. Therefore, applying the
21 MSCM to datasets that contain incongruence that is caused by other processes, such as gene
22 flow, can lead to biased phylogeny estimates. To identify possible bias when using the MSCM,
23 we present P2C2M.SNAPP. P2C2M.SNAPP is an R package that identifies model violations
24 using posterior predictive simulation. P2C2M.SNAPP uses the posterior distribution of species
25 trees output by the software package SNAPP to simulate posterior predictive datasets under the
26 MSCM, and then uses summary statistics to compare either the empirical data or the posterior
27 distribution to the posterior predictive distribution to identify model violations. In simulation
28 testing, P2C2M.SNAPP correctly classified up to 83% of datasets (depending on the summary
29 statistic used) as to whether or not they violated the MSCM model. P2C2M.SNAPP represents a
30 user-friendly way for researchers to perform posterior predictive model checks when using the
31 popular SNAPP phylogenetic estimation program. It is freely available as an R package, along
32 with additional program details and tutorials.

33

34

35

36

37

38

## Introduction

40  Alleles that are shared across taxa present a formidable challenge to phylogenetic inference.
41  Species tree inference methods were introduced in an attempt to infer phylogeny without the
42  potentially confounding effects caused by ancestral alleles that were shared across OTUs
43  (*Maddison, 1997*; *Carstens & Knowles 2007*). Since the biological mechanisms that lead to this
44  process (i.e., incomplete lineage sorting) commonly occur at shallow levels of phylogenetic
45  divergence, species trees have largely (but not exclusively; e.g., *Prum et al., 2015*) been applied
46  near the species boundary, and often in clades where species limits are not entirely clear (*Satler,*
47  *Carstens & Hedin, 2013*). Such applications of the species tree model make the implicit
48  assumption that alleles shared across lineages result from incompletely sorted ancestral
49  polymorphism, even though gene flow is possible in closely related taxa. While gene flow was
50  once considered rare above the species level (at least in animals), recent investigations have
51  suggested that it is more common than previously recognized (e.g., snowshoe hares: *Melo-*
52  *Ferreira et al. (2014)*, chipmunks: *Sullivan et al. (2014)*, bears: *Kumar et al. (2017)*, and *Myotis*
53  bats: *Morales et al. (2017)*).
54
55  Given that gene flow has been shown to bias estimates of both topology and branch lengths when
56  it is not accounted for in a phylogenetic analysis (*Eckert & Carstens, 2008*; *Leache et al., 2013*),
57  evolutionary biologists should (at the least) consider the possibility that gene flow has interfered
58  with phylogeny estimation, particularly when inferring phylogeny from closely related species
59  where reproductive isolation may not be complete. One approach is to look for evidence of gene
60  flow in the data, for example by searching for alleles that are shared across non-sister taxa
61  because such alleles are more likely to result from gene flow than coalescent processes.
62  However, this is likely to be a laborious process, particularly in genomic datasets, and gene flow
63  can be easily missed in studies that do not analyze data from all possible
64  hybridization/introgression events. It is considerably more efficient to utilize statistical methods,
65  such as posterior predictive simulation, that seek to determine whether a given dataset violates
66  the model assumptions of the phylogenetic analysis (e.g., *Goldman, 1993*; *Reid et al., 2014*).
67
68  Posterior predictive approaches have been developed for several types of phylogenetic models,
69  including models of sequence evolution (*Huelsenbeck et al., 2001*; *Brown, 2014b*), species
70  delimitation (*Barley & Thomson, 2016*; *Barley, Brown & Thomson, 2018*), and species tree
71  estimation (*Reid et al., 2014*). The basic approach is to (i) draw parameter values from the
72  posterior distribution, (ii) simulate new datasets using these parameter values under the model
73  assumed by the analysis, (iii) analyze the simulated data to generate posterior predictive
74  distributions, and (iv) calculate and compare summary statistics from either the empirical data or
75  the posterior distribution to the posterior predictive distribution. Analytical models that represent
76  a good fit for the empirical data should produce summary statistics values that fall within the
77  distribution of values estimated under the correct model with posterior predictive datasets
78  (*Brown, 2014b*). Recently, posterior predictive checks have been incorporated into an R package

79 (Posterior Predictive Checks of Coalescent Models (P2C2M): *Gruenstaeudl et al., 2016*) for the
80 MSCM framework. P2C2M was designed to easily allow users to perform posterior predictive
81 analyses, but the program uses the species tree inference package *BEAST which is intended for
82 smaller, sub-genomic data sets (*Heled & Drummond, 2009*). Here, we expand P2C2M to the
83 genomic era so that it can be used to conduct posterior predictive checks using single nucleotide
84 polymorphisms (SNPs) in the SNAPP implementation of the MSCM (*Bryant et al., 2012*).
85

86 **Materials & Methods**
87 *Pipeline*
88 The posterior predictive simulation framework for SNAPP (P2C2M.SNAPP) has been
89 implemented as an R package (*R Core Team, 2018*), with detailed program settings described in
90 the package documentation and tutorial. P2C2M.SNAPP differs from the original P2C2M in the
91 input datatype (sequence data in the original versus SNP data in the SNAPP version) and
92 consequently the summary statistics used to compare empirical and posterior predictive datasets.
93 User input to P2C2M.SNAPP includes the SNAPP .xml formatted input file, the posterior
94 distribution of species trees and log file from a SNAPP analysis, and a metadata text file
95 containing the number of SNPs used, an estimated mutation rate, and the number of samples per
96 group. Importantly, P2C2M.SNAPP assumes users have properly performed SNAPP species tree
97 estimation analysis, including selecting the proper priors for their data and study system and
98 checking for Markov chain convergence. Because P2C2M.SNAPP relies on the posterior
99 distribution of species trees, users should retain at least 100 trees in the posterior distribution to
100 sample from. P2C2M.SNAPP proceeds as follows: (i) it samples, either uniformly or at random,
101 a user-specified number of species trees from the posterior distribution, (ii) extracts taxonomic
102 relationships and branch lengths from each tree, and (iii) for each tree sampled from the
103 posterior, it simulates a posterior predictive dataset under the MSCM using fastsimcoal2, a user-
104 specified number of simulations (*Excoffier et al., 2013*) and the parameters extracted from the
105 metadata text file (Figure 1). Posterior predictive datasets are converted to SNAPP .xml files,
106 and users conduct SNAPP analyses on each posterior predictive dataset using the .xml file output
107 by P2C2M.SNAPP. Prior distributions and Markov chain parameters for the posterior predictive
108 SNAPP analyses are recycled from those used in the original SNAPP analysis in order to
109 maintain consistency. Given the intense computational requirements of SNAPP, generation of
110 the posterior predictive species tree distributions is best conducted using parallel computation.
111 Example scripts for automating SNAPP analyses are included with the tutorial
112 (http://www.github.com/P2C2M/P2C2M_SNAPP). The results of SNAPP analyses on the
113 posterior predictive datasets (i.e., SNAPP .xml files, posterior species tree distributions, and log
114 files) are subsequently used as input for the second stage of the P2C2M.SNAPP analysis, where
115 summary statistics from the posterior and posterior predictive datasets are calculated and
116 compared to identify model violations.
117

118 *Summary Statistics*

119    Generally, summary statistics used in posterior predictive checks fall into two categories: data-
120    based, which compare the empirical and posterior predictive datasets themselves, and inference-
121    based, which compare the inferences produced by analyzing the empirical and posterior
122    predictive datasets (*Brown, 2014a; Barley and Thomson, 2016*). Inference-based statistics can
123    provide more insight as to whether a model violation affects the end result (e.g., the estimated
124    species tree), but can also be more computationally difficult because posterior predictive datasets
125    need to be analyzed with the same methods as the posterior (i.e., species trees need to be
126    estimated with SNAPP). In contrast, data-based statistics do not determine the effect a model
127    violation has on the inference, but are usually computationally efficient. Both data-based and
128    inference-based summary statistics were evaluated to determine which statistic identified model
129    violations to the MSCM with the highest accuracy. Data-based statistics included several based
130    on a fixation index ($F_{ST}$), and inference-based statistics included tree metrics based on Robinson-
131    Foulds or Kuhner-Felsenstein tree distances, and the mean and standard deviation of tree
132    likelihoods. $F_{ST}$ is a commonly used metric for measuring the amount of population structure,
133    and the value ranges from 0 to 1, with populations becoming more structured as $F_{ST}$ approaches
134    1 (*Wright, 1949*). Therefore, lineages exchanging genes should exhibit lower $F_{ST}$ values because
135    they will share alleles. Pairwise $F_{ST}$ was calculated across all loci in the KRIS package
136    (*Chaichoompu, 2018*). $F_{ST}$ summary statistics included mean $F_{ST}$, range of $F_{ST}$, and an $F_{ST}$ outlier
137    test. For the mean and range $F_{ST}$ statistics, the summaries are calculated for each posterior
138    predictive dataset and the empirical dataset. Similar to a two-tailed posterior predictive *p*-value
139    (*Brown, 2014a*; *Barley, Brown & Thomson, 2018*), a *p*-value is calculated by counting the
140    number of posterior predictive datasets with summary statistic values falling above and below
141    the empirical value, multiplying the lesser of these values by two (to emulate a two-tail test), and
142    then dividing by the total number of posterior predictive datasets. We consider *p*-values less than
143    $\alpha = 0.05$ to indicate a model violation. The $F_{ST}$ outlier test was conducted by calculating the
144    average difference between empirical and simulated values for each pairwise comparison, and
145    then conducting an outlier test using the function *boxplot.stats* in the grDevices package (*R core*
146    *team, 2018*). Since we consider any detected outlier to indicate a model violation, the pairwise
147    outliers identified by this approach can be used to identify lineages exchanging genes.
148
149    Two tree distance metrics were also examined, one that considers topology only and one that
150    considers topology and branch lengths. The Robinson-Foulds distance compares the topology
151    between two phylogenetic trees, with values ranging from 0 (no topology difference) to 1
152    (completely different topologies) (*Robinson & Foulds, 1981*). High rates of gene flow can
153    influence topology estimation and result in an errant clade consisting of two lineages that are not
154    closely related but ~~that~~ share alleles due to gene flow. However, it may be more likely that gene
155    flow may mislead the estimation of branch lengths even if the underlying topology is correct.
156    Therefore, a tree distance metric incorporating branch length differences as well as topology may
157    prove to be a useful summary statistic for comparing empirical and posterior predictive datasets.
158    One such metric is the Kuhner-Felsenstein distance, which also calculates values between 0 (no

159    difference between trees) and 1 (high difference between trees) (*Kuhner & Felsenstein, 1994*).

160    Both tree distance metrics were calculated using the ape package (*Paradis, Claude & Strimmer,*

161    *2018*). If posterior trees were estimated from a dataset that violates the MSCM model, we expect

162    that these trees will have large tree distances when compared to posterior predictive trees

163    simulated under the correct model (MSCM). Additionally, as all posterior trees reflect similar

164    processes in the empirical dataset, we expect that tree distances among trees in the posterior

165    under a model violation will be less than distances between the posterior and posterior predictive

166    trees. Therefore, for the tree distance metrics, 1000 comparisons were performed between

167    random trees sampled from the original SNAPP posterior distribution of species trees to create a

168    null distribution. Then 100 random trees from the posterior predictive distribution were

169    compared to the posterior tree they were simulated from, and this was repeated for each posterior

170    predictive dataset. A *p*-value was calculated by counting the number of posterior predictive to

171    posterior tree comparisons falling above the 95% null distribution (values below the 95% null

172    distribution represent high similarity between posterior and posterior predictive datasets, and

173    thus are not useful for detecting violations), and then dividing by the total number of

174    comparisons. We consider *p*-values greater than $\alpha = 0.05$ to indicate model violations. Finally,

175    because it is likely more difficult to estimate trees with high probability under an incorrect

176    model, we examined the mean and standard deviation of tree likelihoods as calculated from

177    SNAPP output. The evaluation of the likelihood statistics follows that of the ~~mean and range~~ $F_{ST}$

178    statistics, described above.

179

180    *Testing*

181    P2C2M.SNAPP was tested by simulating data under the MSCM and via a second simulation

182    under the MSCM with gene flow (i.e., MSCM+*m*) (Figure 2). One hundred replicates were

183    performed under each model. Note that the MSCM+*m* model is a clear violation of the

184    underlying coalescent model that is incorporated into SNAPP because an appreciable portion of

185    the shared polymorphism results from gene flow. All simulations were based on 2000 SNPs, 6

186    species with two individuals sampled per species, an effective population size ($N_e$) of 100K

187    individuals, and a symmetric topology with speciation event times of 5N, 10N, and 20N

188    generations. The number of SNPs simulated is lower than many empirical data sets, but it allows

189    SNAPP analyses to proceed in less time and should represent a conservative test of the ability of

190    P2C2M.SNAPP to detect model violations because the performance of SNAPP generally

191    improves with additional data (*Bryant et al. 2012*). The MSCM+*m* model was designed as a

192    secondary contact scenario, with gene flow between two lineages starting at 2.5N generations in

193    the past and continuing until the present. Both the species experiencing gene flow and the rate of

194    gene flow were selected at random, with the rate of gene flow having a uniform prior distribution

195    between 0.5 and 5 migrants per generation. Simulations were performed using fastsimcoal2

196    (*Excoffier et al., 2013*) and simulated datasets were converted to SNAPP .xml files using custom

197    Python scripts (http://www.github.com/P2C2M/P2C2M_SNAPP). SNAPP analyses were

198    conducted using the following parameters: a gamma prior on the rate of species divergence

199   (lambda) under the Yule speciation prior with α = 2 and β=200, a gamma rate prior on ancestral
200   effective population sizes (theta), mutation rate of μ=ν=1.0, and a Markov chain of 1M steps
201   with 100K burn-in steps and sampling every 1K steps. In order to evaluate the summary
202   statistics, the number of correct inferences, false positives, and false negatives were calculated
203   for each model (200 total) using the posterior and posterior predictive distributions from the
204   SNAPP analyses. False positives are defined as datasets simulated under the MSCM that were
205   indicated as model violations by P2C2M.SNAPP. Conversely, false negatives are defined as
206   datasets simulated under the MSCM+*m* model that were not detected as model violations by
207   P2C2M.SNAPP. Mathews Correlation Coefficient (MCC; *Mathews, 1975*) was also calculated
208   for each summary statistic with the R package *mltools* (*Gorman, 2018*). The MCC takes into
209   account false negatives and positives while measuring how well a binary classifier performs, in
210   this case whether a summary statistic correctly classifies a dataset. The coefficient ranges from -
211   1 to 1 with -1 indicating the classifier is completely wrong and 1 indicating it is completely
212   correct. Additionally, pairwise $F_{ST}$ outliers were compared to the MSCM+*m* simulation
213   parameters to assess if the statistic could identify the species exchanging genes to cause model
214   violations. Finally, *p*-values for each simulation were plotted against gene flow to identify any
215   trends between the level of gene flow and summary statistic performance.
216

## Results
218   P2C2M.SNAPP requires about five minutes on an average laptop (2.6GHz Intel Core i5, 8GB
219   RAM) to generate posterior predictive datasets at the beginning of the pipeline and to evaluate
220   summary statistics in order to identify model violations at the end of the pipeline. However, the
221   entire pipeline requires a considerable amount of time due to the demands of the SNAPP
222   program itself. For example, each replicate of our simulation testing required 300-450 CPU
223   hours on the Pitzer cluster (28 cores and 112GB RAM) at the Ohio Supercomputer Center (*Ohio*
224   *Supercomputer Center, 2018*). While this is clearly not an analysis that users would likely
225   conduct on a laptop computer, the time required for users to analyze their data using
226   P2C2M.SNAPP is still likely to be less than the time required to collect the samples, generate the
227   sequencing libraries, and conduct the bioinformatics.
228

229   There was a dramatic difference across summary statistics in the ability of P2C2M.SNAPP to
230   identify model violations due to gene flow (Table 1). The mean and range of pairwise $F_{ST}$ values
231   correctly classified datasets in only 33% and 41% of simulations respectively (MCC equals -0.45
232   and -0.32 respectively; Figure 2). Each of these statistics exhibited a large number of false
233   positives in which a model violation was detected in a dataset that was simulated under the
234   assumptions of the MSCM. While the pairwise $F_{ST}$ outlier test classified 46% of datasets
235   correctly, the majority of misclassifications were false negatives (MCC equals -0.17).
236   Additionally, we examined the ability of the pairwise $F_{ST}$ outlier test to identify the OTUs
237   exchanging genes in the MSCM+*m* datasets. As the statistic only identified 3% of true model
238   violations, there were very few datasets to test. The pairwise $F_{ST}$ outlier test did not correctly

239     identify the OTUs exchanging genes in any of the datasets. Each tree statistic correctly classified
240     only around half of the datasets, with a high number of false negatives when using the Robinson-
241     Foulds distance and a similar number of false positives when using the Kuhner-Felsenstein
242     distance (MCC equals 0 for each). Similarly, evaluations by the mean tree likelihood statistic
243     were split evenly between correct inferences and false positives (MCC equals -0.29). Our results
244     identified one statistic that performed well. The standard deviation of tree likelihoods correctly
245     classified 83% of simulated datasets, with 14% false negatives and 3% false positives (MCC
246     equals 0.68). Only two summary statistics showed a trend between the rate of gene flow and the
247     $p$-value of posterior predictive checks (Figure 3). For the range of pairwise $F_{ST}$ and mean of tree
248     likelihoods, $p$-values decreased as the rate of gene flow increased.
249

## Discussion

251     While it has been known for some time that model violations can degrade the accuracy of
252     phylogenetic estimation (*Huelsenbeck et al., 2001*; *Eckert & Carstens 2008*, *Leache et al., 2013*;
253     *Brown, 2014a*; *Reid et al., 2014*; *Barley & Thomson, 2016*; *Barley, Brown & Thomson, 2018*),
254     few studies explore possible violations inherent in their datasets to the phylogenetic model used
255     in the analysis (e.g. *Morales et al., 2017; Diaz et al., 2018; Richards et al., 2018*). Apart from
256     the computational demands of the SNAPP analyses, P2C2M.SNAPP represents a user-friendly
257     and reasonably accurate method for identifying violations of the MSCM. The package and
258     tutorial, including examples for running analyses, are available on the P2C2M Github page
259     (https://github.com/P2C2M/P2C2M_SNAPP).
260

261     Our simulation testing indicates that the standard deviation of tree likelihoods is useful in
262     identifying datasets that contain SNP patterns resulting from gene flow between lineages, a clear
263     violation of SNAPP's analytical model. This statistic is likely useful because datasets that violate
264     the MSCM model will be more difficult to estimate and may exhibit posterior distributions with
265     poor convergence. Methods examining the variance within and between posterior and posterior
266     predictive datasets have previously proven useful for posterior predictive checks of Bayesian
267     phylogenetic models (*Gelfand and Ghosh, 1998*; *Lewis et al., 2013*). Users of P2C2M.SNAPP
268     should focus on the standard deviation of tree likelihoods when assessing their datasets.
269     Although higher rates of gene flow should result in a more egregious model violation, it does not
270     appear to be the case that model violations are easier to detect under scenarios with high rates of
271     gene flow. Two summary statistics (range of pairwise $F_{ST}$, mean of tree likelihoods) exhibit an
272     inverse correlation between gene flow and the resulting $p$-value, but both exhibited a high rate of
273     false positives which makes them a poor choice for use in posterior predictive checks. While this
274     relationship does not hold for the standard deviation of tree likelihoods, the statistic is able to
275     detect model violations equally well across a range of migration rates.
276

277     Several statistics were much less useful than we expected them to be. Although the tree distance
278     metrics are conceptually simple, their poor performance may be explained by the reliance of

279     posterior predictive simulation on the empirical phylogeny. Because the posterior predictive data
280     sets are simulated from the empirical phylogeny estimates, inaccuracies in topology and branch
281     lengths of the empirical phylogeny due to gene flow are translated into inaccurate topology and
282     divergence times in the posterior predictive simulations. The result is similar, but inaccurate
283     phylogeny estimates for each data type. $F_{ST}$ is a popular metric in population genetics for
284     examining population structure and gene flow, but may not be applicable to phylogenetic
285     analyses due to fixed differences among lineages. It is possible that including more samples per
286     lineage may increase the usefulness of $F_{ST}$ because more shared polymorphism may be evident,
287     but this may be unfeasible due to the computational requirements of SNAPP. Summary statistics
288     such as $F_{ST}$ are appealing because they can be computed from the posterior predictive datasets
289     without additional SNAPP runs, but many existing statistics were developed for population
290     genetic applications. Summary statistics such as the number of shared or private alleles may be
291     useful. Additionally, the calculation of effect sizes could be beneficial to users because it
292     provides information regarding the degree to which model violation has influenced their results
293     *Brown, 2014a*). Our simulation design investigated a relatively recent diversification scenario
294     because the presence of gene flow is likely to occur when lineages have not become completely
295     reproductively isolated. However, if gene flow occurs in older systems, it should presumably be
296     easier to differentiate from incompletely sorted ancestral polymorphism and thus more easily
297     recognized. Finally, other processes, such as natural selection, also violate the MSCM model,
298     and these additional model violations may also potentially be detectable using the posterior
299     predictive framework implemented in P2C2M.SNAPP, but further research is necessary to
300     identify summary statistics that can detect these violations.
301
302     While the detection of a model violation can have implications for the interpretation of a
303     phylogeny estimate, a model violation does not render the data useless. Minimally, researchers
304     should acknowledge the model violation and temper their interpretation of the patterns evident in
305     the phylogeny. Specifically, the possibility that a model violation may have confounded topology
306     estimates or, more likely, biased branch length/divergence time estimates should be addressed.
307     More preferably, researchers should conduct additional analyses to examine the cause of the
308     model violation, as such violations indicate interesting evolutionary processes not accounted for
309     by the MSCM model. In the case of gene flow, model violations can indicate unknown
310     hybridization among OTUs, and lead to the collection of population-level data that can be
311     analyzed using methods such as Migrate-*n* (*Beerli & Felsenstein, 2001*) or Bayesass (*Wilson &*
312     *Rannala, 2003*). Finally, many recently developed models attempt to infer gene flow and
313     phylogeny under the MSCM for small numbers of lineages (e.g. IMa3; *Hey et al., 2018*,
314     PhyloNet; *Wen et al., 2016*, SpeciesNetwork; *Zhang et al., 2018*). Model violations identified by
315     P2C2M.SNAPP are likely to point researchers to additional analyses that will enable them to
316     understand the history of their focal system.
317
318     **Conclusions**

319 Here we present a new R package for assessing model violations in the species tree estimation
320 program SNAPP. The package uses posterior predictive simulations to identify model violations,
321 and is successful in testing with simulated datasets. P2C2M.SNAPP is the newest addition to a
322 small suite of user-friendly programs for conducting posterior predictive checks (*Gruenstaeudl et*
323 *al., 2016*). Due to the proven benefit of model checking for phylogenetic analyses, we
324 recommend researchers make posterior predictive checks a routine step in estimating
325 phylogenies.
326

## Acknowledgements

## References

334 Barley, A. J., & Thomson, R. C. (2016). Assessing the performance of DNA barcoding using
335 posterior predictive simulations. *Molecular Ecology*, *25*(9), 1944-1957.
336 Barley, A. J., Brown, J. M., & Thomson, R. C. (2018). Impact of model violations on the
337 inference of species boundaries under the multispecies coalescent. *Systematic Biology*, *67*(2),
338 269-284.
339 Beerli, P., & Felsenstein, J. (2001). Maximum likelihood estimation of a migration matrix and
340 effective population sizes in n subpopulations by using a coalescent approach. *Proceedings*
341 *of the National Academy of Sciences*, *98*(8), 4563-4568.
342 Brown, J. M. (2014a). Detection of implausible phylogenetic inferences using posterior
343 predictive assessment of model fit. *Systematic Biology*, *63*(3), 334-348.
344 Brown, J. M. (2014b). Predictive approaches to assessing the fit of evolutionary models.
345 *Systematic biology*, *63*(3), 289-292.
346 Bryant, D., Bouckaert, R., Felsenstein, J., Rosenberg, N. A., & RoyChoudhury, A. (2012).
347 Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a full
348 coalescent analysis. *Molecular Biology and Evolution*, *29*(8), 1917-1932.
349 Carstens B. C., & Knowles L. L. (2007). Estimating phylogeny from gene tree probabilities in
350 Melanoplus grasshoppers despite incomplete lineage sorting. Systematic Biology 56, 400-
351 411.
352 Chaichoompu, K., Abegaz, F., Tongsima, S., Shaw, P. J., Sakuntabhai, A., Pereira, L., & Van
353 Steen, K. (2018). KRIS: Keen and Reliable Interface Subroutines for Bioinformatic Analysis.
354 R package version 1.1.1. Available at: https://CRAN.R-project.org/package=KRIS.
355 Diaz, F., Lima, A. L. A., Nakamura, A. M., Fernandes, F., Sobrinho Jr, I., & De Brito, R. A.
356 (2018). Evidence for introgression among three species of the Anastrepha fraterculus group,
357 a radiating species complex of fruit flies. *Frontiers in genetics*, *9*, 359.
358 Eckert, A. J., & Carstens, B. C. (2008). Does gene flow destroy phylogenetic signal? The
359 performance of three methods for estimating species phylogenies in the presence of gene
360 flow. *Molecular Phylogenetics and Evolution*, *49*(3), 832-842.

361    Excoffier, L., Dupanloup, I., Huerta-Sanchez, E., Sousa, V. C., & Foll, M. (2013). Robust
362        demographic inference from genomic and SNP data. *PLoS Genetics*, *9*(10), e1003905.
363    Gelfand, A. E., & Ghosh, S. K. (1998). Model choice: a minimum posterior predictive loss
364        approach. *Biometrika*, *85*(1), 1-11.
365    Goldman, N. (1993). Statistical tests of models of DNA substitution. *Journal of Molecular
366        Evolution*, *36*(2), 182-198.
367    Gorman, B. (2018). mltools: Machine Learning Tools. R package version 0.3.5. https://CRAN.R-
368        project.org/package=mltools.
369    Gruenstaeudl, M., Reid, N. M., Wheeler, G. L., & Carstens, B. C. (2016). Posterior predictive
370        checks of coalescent models: P2C2M, an R package. *Molecular Ecology Resources*, *16*(1),
371        193-205.
372    Heled, J., & Drummond, A. J. (2009). Bayesian inference of species trees from multilocus
373        data. *Molecular biology and evolution*, *27*(3), 570-580.
374    Hey, J., Chung, Y., Sethuraman, A., Lachance, J., Tishkoff, S., Sousa, V. C., & Wang, Y. (2018).
375        Phylogeny Estimation by Integration over Isolation with Migration Models. *Molecular
376        biology and evolution*, *35*(11), 2805-2818.
377    Huelsenbeck, J. P., Ronquist, F., Nielsen, R., & Bollback, J. P. (2001). Bayesian inference of
378        phylogeny and its impact on evolutionary biology. *Science*, *294*(5550), 2310-2314.
379    Leaché, A. D., Harris, R. B., Rannala, B., & Yang, Z. (2013). The influence of gene flow on
380        species tree estimation: a simulation study. *Systematic Biology*, *63*(1), 17-30.
381    Lewis, P. O., Xie, W., Chen, M. H., Fan, Y., & Kuo, L. (2013). Posterior predictive Bayesian
382        phylogenetic model selection. *Systematic biology*, *63*(3), 309-321.
383    Kuhner, M. K., & Felsenstein, J. (1994). A simulation comparison of phylogeny algorithms
384        under equal and unequal evolutionary rates. *Molecular biology and evolution*, *11*(3), 459-
385        468.
386    Kumar, V., Lammers, F., Bidon, T., Pfenninger, M., Kolter, L., Nilsson, M. A., & Janke, A.
387        (2017). The evolutionary history of bears is characterized by gene flow across
388        species. *Scientific reports*, *7*, 46487.
389    Maddison, W. P. (1997). Gene trees in species trees. *Systematic biology*, *46*(3), 523-536.
390    Matthews, B. W. (1975). Comparison of the predicted and observed secondary structure of T4
391        phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, *405*(2), 442-451.
392    Melo-Ferreira, J., Seixas, F. A., Cheng, E., Mills, L. S., & Alves, P. C. (2014). The hidden
393        history of the snowshoe hare, Lepus americanus: extensive mitochondrial DNA introgression
394        inferred from multilocus genetic variation. *Molecular ecology*, *23*(18), 4617-4630.
395    Morales, A. E., Jackson, N. D., Dewey, T. A., O'Meara, B. C., & Carstens, B. C. (2017).
396        Speciation with gene flow in North American Myotis bats. *Systematic Biology*, *66*(3), 440-
397        452.
398    Ohio Supercomputer Center. (2018). Pitzer Supercomputer. Columbus, OH: Ohio
399        Supercomputer Center.
400    Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and evolution in
401        R language. *Bioinformatics*, *20*(2), 289-290.

402 Prum, R. O., Berv, J. S., Dornburg, A., Field, D. J., Townsend, J. P., Lemmon, E. M., &
403     Lemmon, A. R. (2015). A comprehensive phylogeny of birds (Aves) using targeted next-
404     generation DNA sequencing. *Nature*, *526*(7574), 569.
405 R Core Team. (2018). R: A Language and Environment for Statistical Computing. (R Foundation
406     for Statistical Computing). Vienna, Austria. Availabe at http://www.r-project.org.
407 Reid, N. M., Hird, S. M., Brown, J. M., Pelletier, T. A., McVay, J. D., Satler, J. D., & Carstens,
408     B. C. (2013). Poor fit to the multispecies coalescent is widely detectable in empirical
409     data. *Systematic Biology*, *63*(3), 322-333.
410 Richards, E. J., Brown, J. M., Barley, A. J., Chong, R. A., & Thomson, R. C. (2018). Variation
411     across mitochondrial gene trees provides evidence for systematic error: How much gene tree
412     variation is biological?. *Systematic biology*, *67*(5), 847-860.
413 Robinson, D. F., & Foulds, L. R. (1981). Comparison of phylogenetic trees. *Mathematical*
414     *biosciences*, *53*(1-2), 131-147.
415 Satler J. D., Carstens B. C., Hedin M. (2013). Multilocus species delimitation in a complex of
416     morphologically conserved trapdoor spiders (Mygalomorphae, Antrodiaetidae, Aliatypus).
417     Systematic Biology 62, 805-823.
418 Sullivan, J., Demboski, J. R., Bell, K. C., Hird, S., Sarver, B., Reid, N., & Good, J. M. (2014).
419     Divergence with gene flow within the recent chipmunk radiation (Tamias). *Heredity*, *113*(3),
420     185.
421 Wen, D., Yu, Y., Zhu, J., & Nakhleh, L. (2018). Inferring phylogenetic networks using
422     PhyloNet. Systematic biology, 67(4), 735-740.
423 Wilson, G. A., & Rannala, B. (2003). Bayesian inference of recent migration rates using
424     multilocus genotypes. *Genetics*, *163*(3), 1177-1191.
425 Wright, S. (1949). The genetical structure of populations. *Annals of eugenics*, *15*(1), 323-354.
426 Zhang, C., Ogilvie, H. A., Drummond, A. J., & Stadler, T. (2017). Bayesian inference of species
427     networks from multilocus sequence data. *Molecular biology and evolution*, *35*(2), 504-517.

# Figure 1

Workflow of the P2C2M.SNAPP pipeline.

Blue arrows represent the path of the data. Steps outlined in blue are those performed by the user and steps outlined in red are performed by P2C2M.SNAPP. Workflow proceeds from the top of the figure.
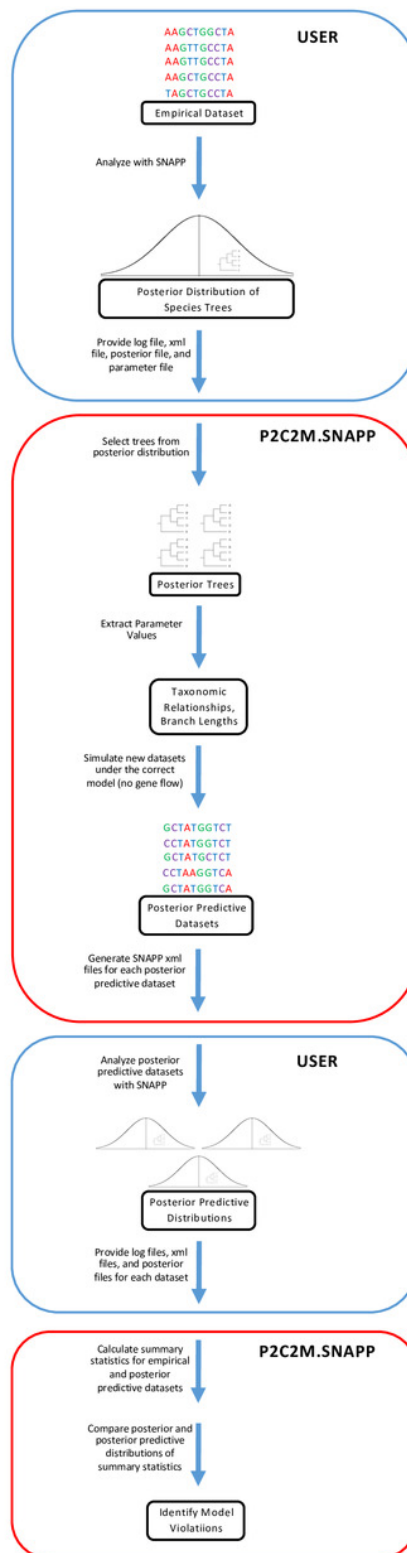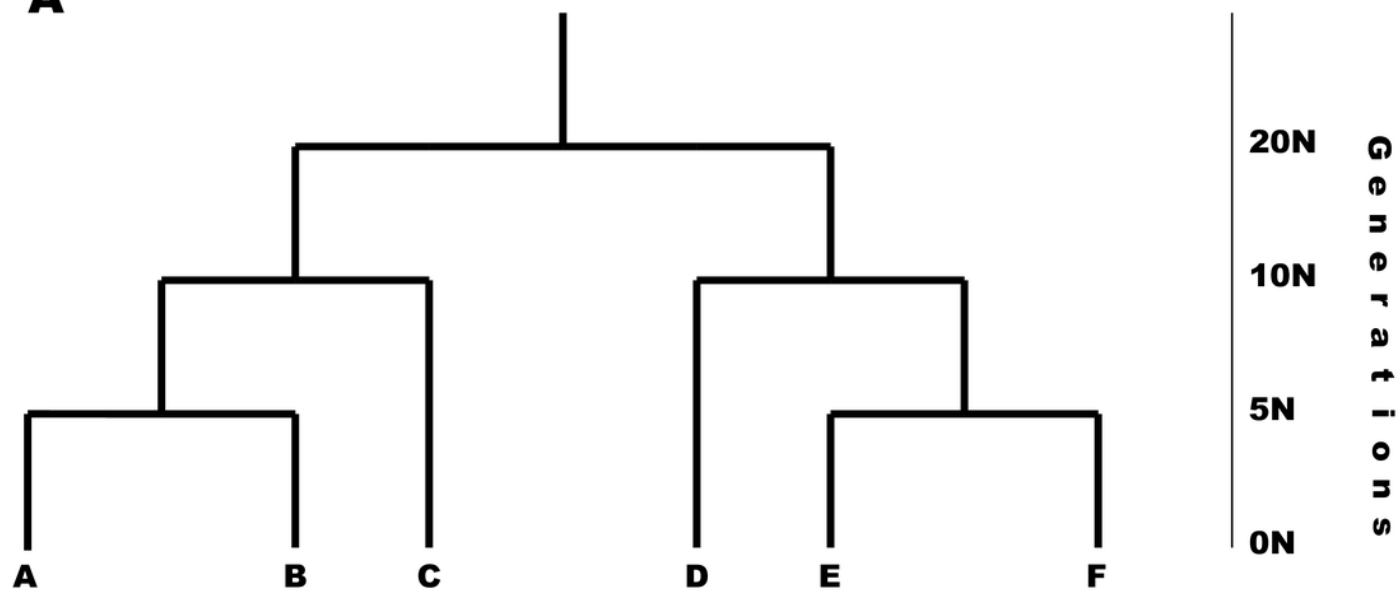
# Figure 2

Models used in simulation testing.

A) MSCM model used for simulation testing. B) Example of the MSCM+$m$ model that includes gene flow violating the MSCM model implemented in SNAPP. The amount of gene flow and taxa exchanging genes were randomly selected for each simulation replicate.
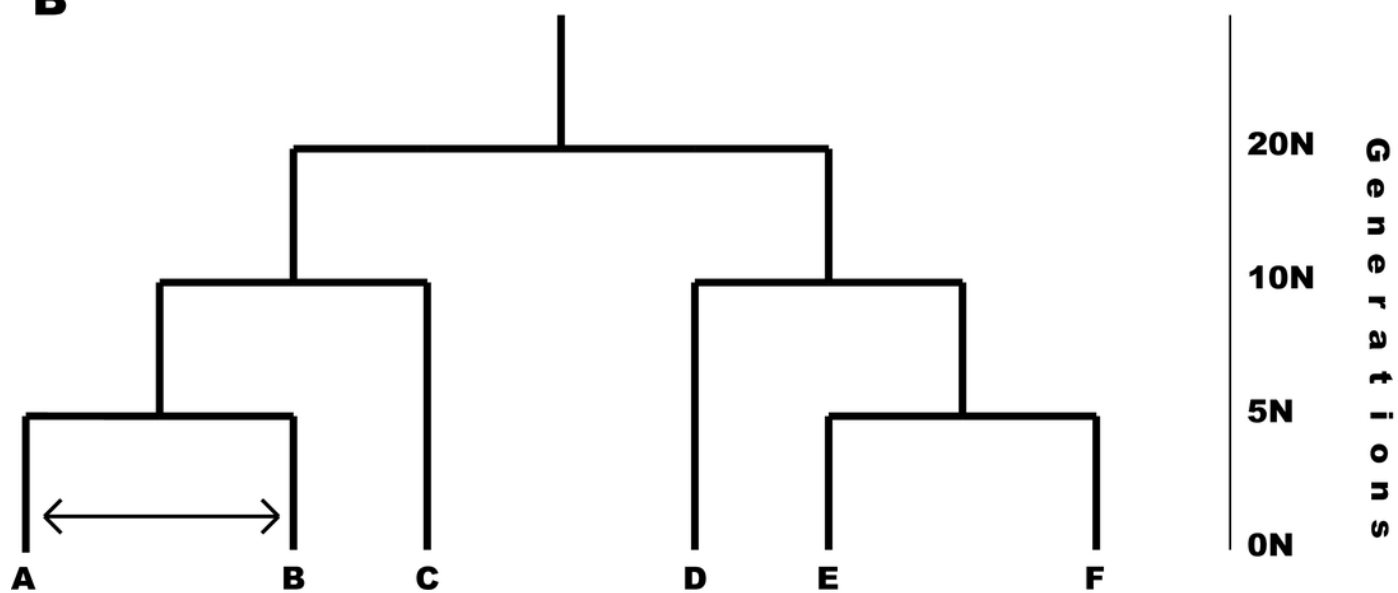
# Figure 3

Correlations between the level of gene flow and the ability of each summary statistic to identify model violations.

The $p$-value for each MSCM+$m$ simulation is plotted against the amount of gene flow simulated with that dataset. FSTA: average pairwise $F_{ST}$, FSTR: range of pairwise $F_{ST}$, KF: Kuhner-Felsenstein distance, MLM: Mean of the maximum likelihood of posterior trees, MLSD: standard deviation of the maximum likelihood of posterior trees, RF: Robinson-Foulds distance.
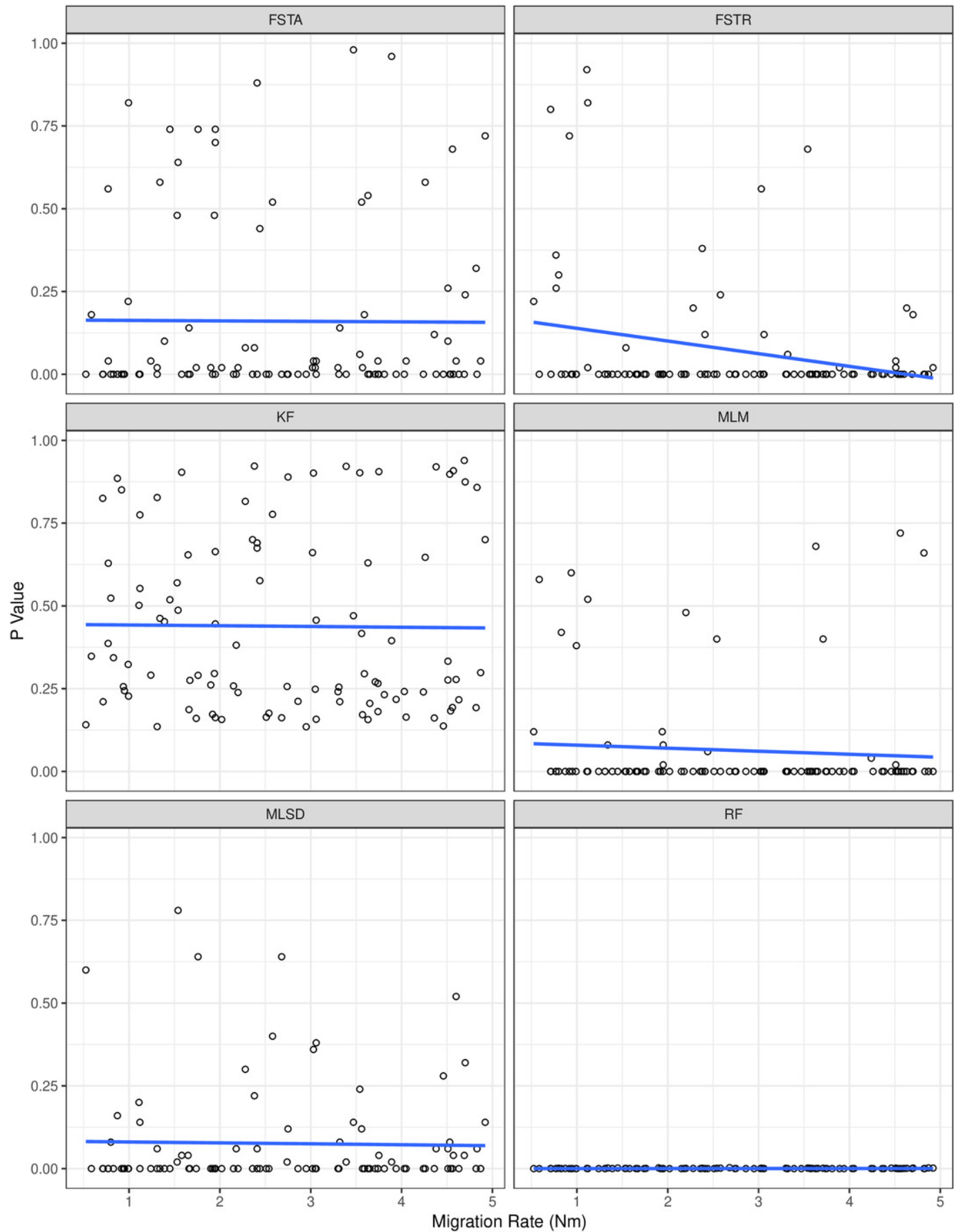
**Table 1**(on next page)

Results of Simulation Testing.

Results include all simulations with both the MSCM and MSCM+$m$ models. False positives are datasets simulated under the MSCM model which P2C2M.SNAPP classified as a model violation. False negatives are datatsets simulated under the MSCM+$m$ model that P2C2M.SNAPP classified as not violating the model implemented in SNAPP.

| Statistic | True Positives | True Negatives | False Positives | False Negatives | Matthew's Correlation Coefficient (MCC) |
|---|---|---|---|---|---|
| Average Pairwise $F_{ST}$ (FSTA) | 66 | 0 | 100 | 34 | -0.45 |
| Range of Pairwise $F_{ST}$ (FSTR) | 81 | 0 | 100 | 19 | -0.32 |
| $F_{ST}$ Outlier Test (PFST) | 3 | 88 | 12 | 97 | -0.17 |
| Kuhner-Felsenstein Distance (KF) | 100 | 0 | 100 | 0 | 0.00 |
| Robinson-Foulds Distance (RF) | 0 | 100 | 0 | 100 | 0.00 |
| Mean of Maximum Likelihood (MLM) | 84 | 0 | 100 | 16 | -0.29 |
| Standard Deviation of Maximum Likelihood (MLSD) | 71 | 95 | 5 | 29 | 0.68 |

1