

**Evolution and Speciation of North American Toads**  
by

Kerry Allen Cobb

A dissertation submitted to the Graduate Faculty of  
Auburn University  
in partial fulfillment of the  
requirements for the Degree of  
Doctor of Philosophy

Auburn, Alabama  
December 9, 2023

Keywords: speciation, hybridization, hybrid zone, phylogenetics, introgression, admixture

Copyright 2023 by Kerry Allen Cobb

Approved by

Jamie R. Oaks, Chair, Associate Professor, Department of Biological Sciences  
Rita M. Graze, Associate Professor, Department of Biological Sciences  
Tonia S. Schwartz, Associate Professor, Department of Biological Sciences  
Laurie S. Stevison, Associate Professor, Department of Biological Sciences

# Abstract

1

...

2

# Acknowledgments

3

...

4

# Contents

5

<b>Abstract</b> . . . . .	<b>2</b>	6
<b>Acknowledgments</b> . . . . .	<b>3</b>	7
<b>Table of Contents</b> . . . . .	<b>6</b>	8
<b>List of Figures</b> . . . . .	<b>8</b>	9
<b>List of Tables</b> . . . . .	<b>9</b>	10
<b>1 Introduction</b> . . . . .	<b>10</b>	11
<b>2 Phylogeography</b> . . . . .	<b>12</b>	12
2.1 Introduction . . . . .	12	13
2.2 Methods . . . . .	15	14
2.2.1 Sampling and DNA Isolation . . . . .	15	15
2.2.2 RADseq Library Preparation . . . . .	16	16
2.2.3 Phylogenetic Data Processing . . . . .	18	17
2.2.4 Maximum Likelihood . . . . .	18	18
2.2.5 Multispecies Coalescent . . . . .	19	19
2.2.6 Test for Historic Admixture . . . . .	19	20
2.2.7 Population Structure Data Processing . . . . .	20	21
2.2.8 Population Structure . . . . .	20	22
2.2.9 Recent <i>A. fowleri</i> x <i>A. woodhousii</i> hybridization . . . . .	21	23
2.3 Results . . . . .	21	24

2.3.1	Assembly and alignment with <i>ipyrad</i>	21	25
2.3.2	Maximum Likelihood Phylogeny	22	26
2.3.3	Coalescent Phylogeny	22	27
2.3.4	Historic Introgression	23	28
2.3.5	Population Structure	24	29
2.3.6	Hybridization between <i>A. fowleri</i> and <i>A. woodhousii</i>	24	30
2.4	Discussion	24	31
2.4.1	Phylogenetic relationships	24	32
2.4.2	Divergence Time	27	33
2.4.3	Population Structure and Hybridization	27	34
2.4.4	Conclusion	28	35
References		28	36
2.5	Figures	33	37
2.6	Tables	44	38
<b>3</b>	<b>Hybrid Zone</b>	<b>53</b>	<b>39</b>
3.1	Introduction	53	40
3.2	Methods	57	41
3.2.1	Sampling and DNA Isolation	57	42
3.2.2	RADseq Library Preparation	58	43
3.2.3	Data Processing	60	44
3.2.4	Genetic Clustering & Ancestry Proportions	60	45
3.2.5	Genomic Cline Analysis	61	46
3.2.6	Genetic differentiation and Introgression	62	47
3.3	Results	63	48
3.3.1	Sampling and Data Processing	63	49
3.3.2	Genetic Clustering & Ancestry Proportions	63	50
3.3.3	Patterns of Introgression	64	51
3.3.4	Genomic Differentiation	65	52
3.4	Discussion	65	53

3.4.1	Evidence for ongoing hybridization . . . . .	65	54
3.4.2	Variability of introgression . . . . .	67	55
3.4.3	Relationship between introgression and differentiation . . . . .	69	56
3.4.4	Conclusion . . . . .	70	57
References	. . . . .	70	58
3.5 Figures	. . . . .	77	59
3.6 Tables	. . . . .	83	60
<b>4 Comparison of Linked versus Unlinked Character Models for Species Tree Inference</b>	. . . . .	<b>91</b>	61
4.1 Introduction	. . . . .	91	62
4.2 Methods	. . . . .	94	64
4.2.1 Simulations of error-free data sets	. . . . .	94	65
4.2.2 Introducing Site-pattern Errors	. . . . .	95	66
4.2.3 Assessing Sensitivity to Errors	. . . . .	95	67
4.2.4 Project repository	. . . . .	97	68
4.3 Results	. . . . .	97	69
4.3.1 Behavior of linked ( <i>StarBEAST2</i> ) versus unlinked ( <i>ecoevolity</i> ) character models	. . . . .	97	70
4.3.2 Analyzing all sites versus SNPs with <i>ecoevolity</i>	. . . . .	98	71
4.3.3 Coverage of credible intervals	. . . . .	98	72
4.3.4 MCMC convergence and mixing	. . . . .	99	73
4.4 Discussion	. . . . .	99	75
4.4.1 Robustness to character-pattern errors	. . . . .	100	76
4.4.2 Relevance to empirical data sets	. . . . .	102	77
4.4.3 Recommendations for using unlinked-character models	. . . . .	103	78
4.4.4 Other complexities of empirical data in need of exploration	. . . . .	104	79
References	. . . . .	105	80
4.5 Figures	. . . . .	109	81

# List of Figures

82

2.1.	Distribution of <i>americanus</i> group samples . . . . .	34	83
2.2.	Iqtree . . . . .	35	84
2.3.	Iqtree . . . . .	36	85
2.4.	Multispecies coalescent phylogeny . . . . .	37	86
2.5.	<i>f</i> -branch statistics . . . . .	38	87
2.6.	Americanus Population Structure . . . . .	39	88
2.7.	Fowleri Population Structure . . . . .	40	89
2.8.	Terrestris Population Structure . . . . .	41	90
2.9.	Woodhousii Population Structure . . . . .	42	91
2.10.	Fowleri Population Structure . . . . .	43	92
3.1.	Evanno method for optimal value of K in <i>STRUCTURE</i> . . . . .	77	93
3.2.	<i>STRUCTURE</i> iterations . . . . .	78	94
3.3.	Summarized <i>STRUCTURE</i> results for each value of K. . . . .	79	95
3.4.	Genetic evidence of hybridization between <i>A. americanus</i> and <i>A. terrestris</i> . . . . .	80	96
3.5.	Shape of genomic clines . . . . .	81	97
3.6.	Patterns of genomic divergence. . . . .	81	98
3.7.	Relationship between genetic divergence and introgression. . . . .	82	99
4.1.	Simulation model . . . . .	109	100
4.2.	Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 1000 base pair loci . . . . .	110	102
4.3.	Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 500 base pair loci . . . . .	111	104

4.4. Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 250 base pair loci . . . . .	112	105
4.5. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R \mu$ ) with 1000 base pair loci . . . . .	113	106
4.6. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R \mu$ ) with 500 base pair loci . . . . .	114	107
4.7. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R \mu$ ) with 250 base pair loci . . . . .	115	109
4.8. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D \mu$ ) with 1000 base pair loci . . . . .	116	112
4.9. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D \mu$ ) with 500 base pair loci . . . . .	117	113
4.10. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D \mu$ ) with 250 base pair loci . . . . .	118	117
4.11. Performance of <i>ecoevolity</i> with data sets simulated with unlinked characters . . . . .	119	121

# List of Tables

124

2.1	Samples used in this study . . . . .	45	125
3.1	Samples collected for this study . . . . .	83	126
3.2	Samples loaned from museums . . . . .	90	127

# Chapter 1

128

## Introduction

129

Studying speciation in evolutionary biology is of paramount importance as it provides a fundamental understanding of how biodiversity emerges and evolves over time. Speciation, the process by which new species arise from a common ancestor, holds the key to unraveling the intricate tapestry of life on Earth. It not only sheds light on the mechanisms driving genetic and phenotypic diversity but also offers crucial insights into the adaptation and survival of species in response to changing environments. Moreover, the study of speciation is essential for comprehending the origins of complex ecosystems and ecological interactions, as different species play distinct roles in shaping their environments. Ultimately, delving into the intricacies of speciation is not just an academic pursuit; it is a foundational pillar in our quest to understand life's past, present, and future, with implications for conservation, agriculture, and our own place in the natural world.

[KAC comment: Planning to move this to intro chapter.] It is clear that significant levels of admixture occur within some Bufonidae hybrid zones which could yield insights into the evolution of reproductive incompatibility. There are a few qualities that make these hybrid zones particularly attractive for further investigation. One of these qualities is the ease with which the primary behavioral isolating mechanisms, spawning period and advertisement call, can be measured and quantified in order to understand the strength of prezygotic mating barriers and possible patterns consistent with reinforce-

ment (Blair, 1974; Cocroft & Ryan, 1995; Kennedy, 1962). It has also been show that 149  
they can be readily bred in captivity (Blair, 1972). Many species produce thousands of 150  
offspring which are externally fertilized making a variety of embryological observations or 151  
manipulations possible (Blair, 1972). Breeding can be induced hormonally or performed 152  
in vitro, facilitating the planning and scheduling of experiments (Trudeau et al., 2010). 153  
Unlike many of the organisms which have undergone intensive study in the context of 154  
speciation such as *Drosophila*, *Mus*, and *Heliconius*, most Bufonidae have homomorphic 155  
sex chromosomes (Blair, 1972). This is an interesting contrast in light of the important 156  
roll of sex chromosomes have in the evolution of reproductive incompatibility and the 157  
roll of heterogamety in explaining evolutionary patterns such as Haldane's rule, faster 158  
male evolution, and faster-X evolution (Delph & Demuth, 2016). Furthermore, there is 159  
evidence of sex chromosome turnovers within Bufonidae which presents an opportunity to 160  
study differences in the evolution of reproductive incompatibility among closely related 161  
species with different sex determination systems (Dufresnes et al., 2020; Stöck et al., 162  
2011). All of these qualities along with the near global distribution of a large 642 species 163  
radiation present an excellent opportunity to further our understanding of the evolution 164  
of reproductive incompatibility (AmphibiaWeb, 2023). 165

Speciation is a major focus of evolutionary biology 166

# Chapter 2

167

## Phylogeography

168

### 2.1 Introduction

169

Many factors are understood to be important in driving and shaping the diversification and evolutionary history of organisms. Chief among them is the interplay between climatic conditions and geologic processes. Changes in these environmental variables can alter the distributions of organisms and result in changes in the connectivity of populations. Disconnected populations may undergo genetic divergence from one another due to adaptive evolution in response to changing abiotic or biotic conditions. Or they might simply diverge via neutral evolution driven by the effects of drift. Environmental changes can also reconnect previously isolated populations resulting in hybridization and gene flow, another very important process shaping patterns of diversity. Understanding the interplay of all of these factors is critical for understanding the evolutionary history of organisms. A critical step to understanding these processes is to obtain an accurate reconstruction of the evolutionary history of organisms.

The North American toads in the genus *Anaxyrus* are a group of organisms with a poorly understood evolutionary history. Although, not for lack of trying. Multiple studies of the evolutionary relationships among species in the genus have produced conflicting results (Fontenot et al., 2011; Graybeal, 1997; Masta et al., 2002; Portik et al., 2023; Pramuk et al., 2007; Pyron & Wiens, 2011). Particularly within the *americanus* group

composed of *A. americanus*, *A. baxteri*, *A. fowleri*, *A. hemiophrys*, *A. houstonensis*,  
187  
*A. microscaphus*, *A. terrestris*, and *A. woodhousii*. Two phylogenetic studies inferred  
188  
trees with *A. fowleri* forming a polytomy making them inconsistent with the current  
189  
taxonomy of *Anaxyrus* (Fontenot et al., 2011; Masta et al., 2002). These conflicting  
190  
results could be due to methodological differences such as the species included, the number  
191  
of individuals of each species sequenced, inference methods used, or the sequenced loci.  
192  
But the differences in inferred relationships could also result from real biological processes.  
193  
Incomplete lineage sorting is one potential source of discordance among datasets which  
194  
include different loci that arises from real biological processes and impacts phylogenetic  
195  
inference (Kubatko & Degnan, 2007). Incomplete lineage sorting could also produce the  
196  
polytypic relationship among *A. fowleri*.  
197

Gene flow is another potential source of discordance among genes which could drive  
198  
the differences in inferred relationships among studies using different loci and could also  
199  
produce the pattern seen in *A. fowleri* (Degnan & Rosenberg, 2009). While incomplete  
200  
lineage sorting is very likely to have impacted patterns of genetic variation in *Anaxyrus*,  
201  
gene flow due to hybridization is a distinct possibility as well. There are numerous  
202  
reports of natural hybridization between several different species of *Anaxyrus* (Green,  
203  
1996). A study of allozyme variation across a hybrid zone between *A. americanus* and  
204  
*A. hemiophrys* revealed introgression taking place across a more than 50km wide hybrid  
205  
zone. Hybrid zones are also suspected to exist between *A. americanus* and *A. terrestris*  
206  
and between *A. woodhousii* and *A. fowleri* though no study has yet been conducted to  
207  
characterize these putative hybrid zones (Green, 1996; Weatherby, 1982). Furthermore,  
208  
numerous laboratory crosses have been performed between pairs of *Anaxyrus* species with  
209  
currently overlapping distributions (Blair, 1963, 1972). Some of which produce viable and  
210  
fertile backcross progeny (Blair, 1963, 1972). These studies suggest that gene flow could  
211  
very easily have played a role in shaping patterns of diversity in *Anaxyrus*. However,  
212  
they provide only a snapshot in time with no indication of the long term consequences.  
213  
There are many potential consequences of hybridization such as adaptive introgression,  
214  
introgression of neutral genetic variation, reinforcement, lineage fusion, polyploidization,  
215

hybrid speciation, or transition to unisexual reproduction (Abbott et al., 2013). Inference  
of past introgression is an important starting point for exploring these outcomes yet it  
remains a challenging problem. The network structure of phylogenetic networks are far  
less tractable to infer than the more simple bifurcating phylogenies for which there has  
been extensive method development. There has been some recent work to overcome this  
challenge as well as increased feasibility of obtaining appropriate genome wide datasets  
to investigate past gene flow.

Apart from the significant evolutionary implications of hybridization which need to  
be understood, it also presents a valuable opportunity for investigating the mechanisms  
that drive divergence and the evolution of reproductive incompatibility (Rieseberg et al.,  
1999). Many generations of backcrossing within hybrid zones can produce a large number  
of highly recombinant genomes that allow for the observation of many possible hybrid  
genotypes under natural conditions in order to identify advantageous or disadvantageous  
hybrid genotypes. In most species it is not feasible to produce such a large number of  
highly recombinant offspring in order to make such observations. The evolutionary history  
of hybridizing species is important context to have when studying patterns of introgres-  
sion within hybrid zones. Context such as the phylogenetic relationship of hybridizing  
species relative to other closely related species, the amount of genetic differentiation, the  
time since divergence and by extension the biogeographic processes driving initial diver-  
gence and subsequent secondary contact in cases of allopatric divergence. This important  
context is currently missing for *Anaxyrus* which limits the inferences that can be made  
regarding hybrid zones in this genus.

Ultimately, environmental change is what leads many populations divergence and and  
may lead to subsequent or concurrent gene flow. Therefore, environmental variables are  
also important context for making inferences from hybrid zones in addition to under-  
standing the process of diversification more generally. To date, there have not been any  
studies conducted to understand how the environment as driven diversification in North  
American toads. North America has had a very complex geologic and climatic history  
(Lyman & Edwards, 2022). The effects of which are often clade specific (Nuñez et al.,  
244

2023). But large scale environmental changes can impact multiple species simultaneously  
245  
(Oaks, 2019). There has been recent development in methods to infer these events (Oaks,  
246  
2019; Oaks et al., 2022). The identification of multiple pairs of lineages that underwent  
247  
divergence at the same time could provide evidence about environmental changes driving  
248  
diversification. Present day population structure could also provide further understanding  
249  
by revealing environmental factors that reduce gene flow assuming the biological limits  
250  
of present day species have not evolved dramatically from the ancestral condition.  
251

In this study, I investigate the evolutionary history of North American toads in the  
252  
genus *Anaxyrus* using genome wide sequence data. For this I obtained restriction enzyme-  
253  
associated DNA sequence (RADseq) data from 12 species of *Anaxyrus* including dense  
254  
sampling representing a large portion of the ranges of *A. americanus*, *A. fowleri*, *A.*  
255  
*terrestris*, and *A. woodhousii*. With these data I infer evolutionary relationships using  
256  
maximum likelihood analysis of a concatenated dataset of many broadly distributed sam-  
257  
ples in addition to using a multispecies coalescent analysis of a subset of the data. I also  
258  
test for the presence of shared divergence times which might suggest *Anaxyrus* diversifi-  
259  
cation has been driven by the same environmental changes and also estimate the absolute  
260  
timing of all divergences within the genus. With the robust estimate of phylogenetic re-  
261  
lationships among *Anaxyrus* species, I test for the presence of past introgression among  
262  
species. In order to identify the types of environmental factors that might have played  
263  
a role in isolating populations that would eventually diverge as species, I investigate  
264  
population structure within a subset of *Anaxyrus* species. Finally, I estimate proportions  
265  
admixture between *A. fowleri* and *A. woodhousii* to test the hypothesis that these species  
266  
form a hybrid zone in the central United States where their ranges meet.  
267

## 2.2 Methods 268

### 2.2.1 Sampling and DNA Isolation 269

I obtained tissue samples from museum tissue collections as well as from individuals I  
270  
collected from 2017 to 2020. I selected samples to represent as much of the range of each  
271

species of *Anaxyrus* as possible. I also included one *Rhinella marina* and one *Incilius nebulifer* for as outgroups for phylogenetic analyses. 272  
273

I isolated DNA from tissues by first lysing a piece of tissue approximately the size 274  
of a grain of rice in 300  $\mu\text{L}$  of a solution of 10mM Tris-HCL, 10mM EDTA, 1% SDS 275  
(w/v), and nuclease free water along with 6 mg Proteinase K that was incubated for 4-16 276  
hours at 55°C in a 1.5 mL microcentrifuge tube. To purify the DNA and separate it from 277  
the lysis product, I mixed the lysis product with a 2X volume of SPRI bead solution 278  
containing 1 mM EDTA, 10 mM Tris-HCl, 1 M NaCl, 0.275% Tween-20 (v/v), 18% PEG 279  
8000 (w/v), 2% Sera-Mag SpeedBeads (GE Healthcare PN 65152105050250) (v/v), and 280  
nuclease free water. I then incubated the samples at room temperature for 5 minutes, 281  
placed the beads on a magnetic rack, and discarded the supernatant once the beads had 282  
collected on the side of the tube. I then performed two ethanol washes by adding 1 mL of 283  
70% ETOH to the beads while still placed in the magnet stand and allowing it to stand 284  
for 5 minutes before discarding the ethanol. After removing all ethanol from the second 285  
wash, I removed the tube from the magnet stand and allowed the sample to dry for 1 286  
minute before mixing the beads with 100  $\mu\text{L}$  of TLE solution containing 10 mM Tris- 287  
HCL, 0.1 mm EDTA, and nuclease free water. After allowing the bead mixture to stand 288  
at room temperature for 5 minutes I returned the beads to the magnet stand, pipetted 289  
all of the TLE solution into another microcentrifuge tube, and discarded the beads. I 290  
quantified DNA with a Qubit fluorometer (Life Technologies, USA) and diluted samples 291  
with TLE solution to bring the concentration to 20 ng/ $\mu\text{L}$ . 292

## 2.2.2 RADseq Library Preparation 293

I prepared RADseq libraries using the 2RAD approach outlined by Bayona-Vásquez 294  
et al., 2019. On 96 well plates, I ligated 100 ng of sample DNA in 15  $\mu\text{L}$  of a solution 295  
with 1X CutSmart Buffer (New England Biolabs, USA; NEB), 10 units of XbaI, 10 units 296  
of EcoRI, 0.33  $\mu\text{M}$  XbaI compatible adapter, 0.33  $\mu\text{M}$  EcoRI compatible adapter, and 297  
nuclease free water with a 1 hour incubation at 37°C. I then immediately added 5  $\mu\text{L}$  of 298  
a solution with 1X Ligase Buffer (NEB), 0.75 mM ATP (NEB), 100 units DNA Ligase 299

(NEB), and nuclease free water and incubated at 22°C for 20 min and 37°C for 10 min for  
300 two cycles, followed by 80°C for 20 min to stop enzyme activity. For each 96 well plate,  
301 I pooled 10  $\mu$ L of each sample and split this pool equally between two microcentrifuge  
302 tubes. I purified each pool of libraries with a 1X volume of SpeedBead solution followed  
303 by two ethanol washes as described in the previous section except that the DNA was  
304 resuspended in 25  $\mu$ L of TLE solution.  
305

In order to be able to detect and remove PCR duplicates, I performed a single cycle  
306 of PCR with the iTru5-8N primer which adds a random 8 nucleotide barcode to each  
307 library construct. For each plate, I prepared four PCR reactions with a total volume of  
308 50  $\mu$ L containing 1X Kapa Hifi Buffer (Kapa Biosystems, USA; Kapa), 0.3  $\mu$ M iTru5-8N  
309 Primer, 0.3 mM dNTP, 1 unit Kapa HiFi DNA Polymerase, 10  $\mu$ L of purified ligation  
310 product, and nuclease free water. I ran reactions through a single cycle of PCR on a  
311 thermocycler at 98°C for 2 min, 60°C for 30 s, and 72°C for 5 min. I pooled all of the  
312 PCR products for a plate into a single tube and purified the libraries with a 2X volume  
313 of SpeedBead solution as described before and resuspended in 25  $\mu$ L TLE. I added the  
314 remaining adapter and index sequences unique to each plate with four PCR reactions with  
315 a total volume of 50  $\mu$ L containing 1X Kapa Hifi (Kapa), 0.3  $\mu$ M iTru7 Primer, 0.3  $\mu$ M  
316 P5 Primer, 0.3 mM dNTP, 1 unit of Kapa Hifi DNA Polymerase (Kapa), 10  $\mu$ L purified  
317 iTru5-8N PCR product, and nuclease free water. I ran reactions on a thermocycler with  
318 an initial denaturation at 98°C for 2 min, followed by 6 cycles of 98°C for 20 s, 60°C  
319 for 15 s, 72°C for 30 s and a final extension of 72°C for 5 min. I pooled all of the PCR  
320 products for a plate into a single tube and purified the product with a 2X volume of  
321 SpeedBead solution as described before and resuspended in 45  $\mu$ L TLE.  
322

I size selected the library DNA from each plate in the range of 450-650 base pairs using  
323 a BluePippin (Sage Science, USA) with a 1.5% dye free gel with internal R2 standards.  
324 To increase the final DNA concentrations I prepared four PCR reactions for each plate  
325 with 1X Kapa Hifi (Kapa), 0.3  $\mu$ M P5 Primer, 0.3  $\mu$ M P7 Primer, 0.3 mM dNTP, 1 unit  
326 of Kapa HiFi DNA Polymerase (Kapa), 10  $\mu$ L size selected DNA, and nuclease free water  
327 and used the same thermocycling conditions as the previous (P5-iTru7) amplification. I  
328

pooled all of the PCR products for a plate into a single tube and purified the product with  
329  
a 2X volume of SpeedBead solution as before and resuspended in 20  $\mu$ L TLE. I quantified  
330  
the DNA concentration for each plate with a Qubit fluorometer (Life Technologies, USA)  
331  
then pooled each plate in equimolar amounts relative to the number of samples on the  
332  
plate and diluted the pooled DNA to 5 nM with TLE solution. The pooled libraries were  
333  
pooled with other projects and sequenced on an Illumina HiSeqX by Novogene (China)  
334  
to obtain paired end, 150 base pair sequences.  
335

### 2.2.3 Phylogenetic Data Processing

336

To produce alignments for phylogenetic analysis, I first I demultiplexed the iTru7  
337  
indexes using the *process\_radtags* command from *Stacks* v2.6.4 (Rochette et al., 2019)  
338  
and allowed for two mismatches for rescuing reads. I removed PCR duplicates using the  
339  
the *clone\_filter* command from *Stacks*. To demultiplex individual samples I used *ipyrad*  
340  
v0.9.90 and allowed for one mismatch for rescuing reads. I assembled and aligned reads  
341  
with *ipyrad* using default parameters and a clustering threshold of 0.8. Using *ipyrad*, I  
342  
filtered loci not present in at least 75% of samples and filtered samples with fewer than  
343  
200 loci.  
344

### 2.2.4 Maximum Likelihood

345

Phylogenetic methods that do not account for incomplete lineage sorting do not per-  
346  
form well with data impacted by this process. However, methods that do account for  
347  
incomplete lineage sorting are far more computationally demanding. As a result, these  
348  
methods cannot be performed with a large number of samples. I therefore conducted con-  
349  
ducted maximum likelihood phylogenetic inference in order to infer a phylogeny with all  
350  
of the sequenced samples and to be able to identify samples that may be problematic for  
351  
other methods due to recent admixture or data quality. I conducted the maximum likeli-  
352  
hood phylogenetic inference with *IQ-TREE* v1.6.12 (Nguyen et al., 2015) with the *ipyrad*  
353  
alignment as input in order. I ran *IQ-TREE* with 1000 ultrafast bootstrap replicates  
354  
(Hoang et al., 2018) under the GTR substitution model.  
355

## 2.2.5 Multispecies Coalescent 356

In order to account for incomplete lineage sorting in the inference of phylogenetic relationships and to infer shared divergence times, I used the program *phycoeval*. I selected a subset of up to four samples from each species due to the infeasible run times for *phycoeval* with greater numbers of samples (see table 1). I excluded sample 006 from consideration due it having an anomalous position in the maximum likelihood tree. I used *ipyrad* to filter loci not present in at least 75% of samples. Using a custom script I filtered the phylip alignment file produced by *ipyrad* to exclude sites with more than two characters and output the filtered alignment to nexus format with a biallelic character encoding. I ran *phycoeval* with state frequencies fixed at 0.5. I set the mutation rate equal to one so that divergence times are in units of expected substitutions per site. I set the prior on the age of the root as an exponential distribution with a mean of 0.01. I ran *phycoeval* with the assumption of a single effective population size shared across all of the branches of the tree. The prior on the effective population size was a gamma prior with a shape of four and mean of 0.0005 I ran five independent MCMC chains for 10,000 generations, sampling every 10 generations. Each chain was started with a comb tree topology with all branches sharing the same divergence time. I summarized the posterior sample of tree topologies and parameters using *sumphycoeval*. To assess convergence and mixing, I used *sumphycoeval* to calculate the potential scale reduction factor (PSRF) and the effective sample size (ESS). I discarded the first 100 samples from each chain as burnin. I used *sumphycoeval* to rescale the branch lengths of the maximum a posteriori (MAP) tree produced by *sumphycoeval* so that the posterior mean root age was 16.5 million years ago based on the estimate of Feng et al., 2017.

## 2.2.6 Test for Historic Admixture 379

In order to test for a history of introgression between species of *Anaxyrus* I used the program *dsuite* v0.5r50 (Malinsky et al., 2021) to compute the *f*-branch statistic for each pair of *Anaxyrus* species for which the statistic can be calculated (Malinsky et al., 2018; Reich et al., 2009). I used *ipyrad* to filter all loci that were not found in at least 50% of

the samples that passed filtering and excluded one *A. fowleri* sample (sample 006 from ??) which falls outside of the *A. fowleri* clade inferred by *IQ-TREE*. For the input tree topology required to run *dsuite*, I used the topology inferred by *phycoeval* and I specified *Incilius nebulifer* as the outgroup species. I ran the *dsuite* Dtrios command to compute Patterson's the *f*4-ratio statistic for all possible trios with 20 block-jackknife replicates. I then ran the Fbranch command from *dsuite* to compute the *f*-branch statistics from the computed *f*4-ratio statistics. I plotted the *f*-branch statistics with *dtools* v0.1 which is packaged with the *dsuite* program (Malinsky et al., 2021).

Say something about how *f*-branch takes into account correlation among branches

### 2.2.7 Population Structure Data Processing

I processed reads differently for the analysis of population structure following PCR duplicate filtering. I demultiplexed individual samples, trimmed adapter sequence, and filtered reads with low quality scores as well as reads with any uncalled bases using the *process\_radtags* command and allowed for the rescue of restriction site sequence as well as barcodes with up to two mismatches. I allowed for 14 mismatches between alleles within, as well as between individuals (M and n parameters). This is equivalent to a sequence similarity threshold of 90% for the 140 bp length of reads post trimming. I also allowed for up to 7 gaps between alleles within and between individuals. I used the *populations* command from *Stacks* to filter loci missing in more than 5% of individuals, filter all sites with minor allele counts less than 3, filter any individuals with more than 90% missing loci, and randomly sample a single SNP from each locus.

### 2.2.8 Population Structure

To investigate population structure within *A. americanus*, *A. fowleri*, *A. terrestris*, and *A. woodhousii*, I used the demultiplexed and de-cloned reads used for the phylogenetic analyses for producing alignments. I assembled and aligned these reads using *Stacks* for each species separately. I allowed for 7 mismatches between alleles within, as well as between individuals (M and n parameters). This is equivalent to a sequence similarity

threshold of 95% for the 140 bp length of reads post trimming. I also allowed for up to  
411  
7 gaps between alleles within and between individuals. I used the *populations* command  
412  
from *Stacks* to filter loci missing from more than 5% of samples, filter all sites with minor  
413  
allele counts less than 3, filter any individuals with more than 90% missing loci and to  
414  
randomly sample a single site per locus.  
415

I ran the program *STRUCTURE* v2.3.4 (Pritchard et al., 2000) for each species  
416  
separately using the admixture model in order to cluster individuals and estimate ancestry  
417  
proportions for each individual. I ran *STRUCTURE* under four different models differing  
418  
in the number of populations assumed (K parameter), with the parameter ranging from 1-  
419  
4. I ran 10 iterations of *STRUCTURE* for each value of K for a total of 100,000 steps and  
420  
burnin of 50,000 for each iteration. I used the R package *POPHELP* v2.3.1 (Francis,  
421  
2017) to combine iterations for each value of K and to select the model producing the  
422  
largest  $\Delta K$  which is the the model that has the greatest increase in likelihood score  
423  
from the previous model having one fewer populations as described by (Evanno et al.,  
424  
2005). I also investigated population structure with a non-parametric approach, using  
425  
principle component analysis (PCA) implemented in the R package *adegenet* *adegenet*  
426  
v2.1.10 (Jombart, 2008).  
427

## 2.2.9 Recent *A. fowleri* x *A. woodhousii* hybridization

  
428

## 2.3 Results

  
429

### 2.3.1 Assembly and alignment with *ipyrad*

  
430

A total of 436,265,266 reads were obtained for all samples. After filtering low quality  
431  
reads and reads without restriction site sequence, 435,650,926 total reads remained for  
432  
assembly. The number of filtered reads per individual was highly variable with a mean of  
433  
4,538,030 ( $sd=3,619,076$ ). Prior to filtering there were 171,174 loci total loci which was  
434  
reduced to 659 after filtering loci not present in at least 75% of samples and filtering ??  
435  
samples which had fewer than 200 loci (Table 1). Mean sequence read coverage of the  
436

loci passing filter was 54x. The final alignment contained a total of 184,453 sites with 437  
20,361 SNPs with 14.96% of sites and 14.71% of SNPs missing. 438

### 2.3.2 Maximum Likelihood Phylogeny 439

The full majority rule consensus tree inferred by *IQ-TREE* is presented in 2.2. All 440  
species were inferred as a single monophyletic group with the exception of *A. fowleri*. A 441  
*A. fowleri* sample (sample 006) does not form a monophyletic group with other 442  
*A. fowleri* samples but is instead sister to the branch containing *A. woodhousii* and 443  
*A. fowleri* samples. A representation of the tree inferred by *IQ-TREE* with the tips 444  
within species specific clades collapsed is presented in 2.3. Each species specific clade for 445  
which there are at least two representatives samples all have ultrafast bootstrap support 446  
values of 100%. All branches below the level of the species specific clades have 447  
ultrafast bootstrap support values ranging from 70-100% with the majority being 100%. 448  
The most basal internal branch of the tree, marking the split between most of *Anaxyrus* 449  
and *A. punctatus* along with the outgroup *Incilius nebulifer* has an ultrafast bootstrap 450  
support value of 99%. The sister branch to *A. terrestris*, which contains the spurious 451  
*A. fowleri* sample (sample 006) and the clade containing *A. fowleri* and *A. woodhousii*, 452  
has an ultrafast bootstrap support value of 96%. The lowest ultrafast bootstrap support 453  
value is found on the branch sister to the *A. cognatus/A. speciosus* clade with a value of 454  
only 70%. 455

### 2.3.3 Coalescent Phylogeny 456

The maximum a posteriori (MAP) tree inferred under the multispecies coalescent 457  
model using *phycoeval* has a topology differs from the maximum likelihood topology 458  
inferred by *IQ-TREE* Fig. 2.4. The MAP tree produced by *phycoeval* does not have any 459  
shared divergence times among any of the 10 internal nodes of the tree. The frequency 460  
of topologies in the posterior sample that have 10 independent divergence times is 0.5. 461  
The next most frequent topology in the posterior are topologies with a single shared 462  
divergence time and nine independent divergences and occur with a frequency of 0.24. 463

One major difference between the maximum likelihood tree inferred by *IQ-TREE* and 464  
the MAP tree inferred by *phycoeval* is that the MAP tree has one multifurcation. This 465  
multifurcation happens at the ancestor of the *A. quercicus*, *A. speciosus*/*A. cognatus*, 466  
and *A. americanus* group lineages. However, this node has a low posterior probability of 467  
only 0.51. All other branches in the MAP tree have high posterior probabilities of 0.98 468  
or more. Most divergence events within *Anaxyrus* have occurred in the past 3.5 million 469  
years and most diversification within the *A. americanus* group is less than 2.5 million 470  
years old. 471

### 2.3.4 Historic Introgression

I used the program *dsuite* to compute the *f*-branchstatistic which is an estimate 472  
of excess allele sharing between species pairs that is not due to incomplete lineage 473  
sorting. I used the species tree topology produced by *phycoeval* for estimating the 474  
*f*-branchstatistics. The *f*-branchestimates for each species pair are presented with a 475  
heatmap in figure 2.5. Most *f*-branchestimates produced by *dsuite* were zero or very 476  
near zero. Only 24 out of 112 *f*-branchestimates were greater than 0 and 11 of those 477  
were greater than 0.05 Fig. 2.5. *A. americanus* and *A. woodhousii* had the largest number 478  
of estimates greater than zero associated with them with nearly every pairwise comparison 479  
greater than 0 Fig. 2.5. The highest *f*-branchstatistic values are between *A. americanus* 480  
and two other species: *A. hemiophrys* (0.24) and *A. baxteri* (0.22) Fig. 2.5. The values 481  
associated with *A. woodhousii* are appreciably lower with none exceeding 0.1 ???. The 482  
branch preceding *A. speciosus* and *A. cognatus* tested against *A. punctatus* along with 483  
the tests of *A. quercicus* with *A. cognatus* and *A. speciosus* all exceeded 0.1. 484

2.3.5 Population Structure	486
----------------------------	-----

2.3.6 Hybridization between <i>A. fowleri</i> and <i>A. woodhousii</i>	487
--	-----

2.4 Discussion	488
----------------	-----

2.4.1 Phylogenetic relationships	489
----------------------------------	-----

The maximum likelihood tree inferred by *IQ-TREE* Fig. 2.2 and ?? differs from trees inferred in previous studies of the relationships among *Anaxyrus* (Fontenot et al., 2011; Graybeal, 1997; Masta et al., 2002; Portik et al., 2023; Pramuk et al., 2007; Pyron & Wiens, 2011). Even among these previous studies there has been a great deal of inconsistency in the inferred relationships with the exception of a few taxa. As in all previous studies, the maximum likelihood tree inferred in this study places *A. punctatus* sister to all other *Anaxyrus*. I also found the *americanus* group to be monophyletic with *A. microscaphus* sister to all other *americanus* group species which is consistent with most previous studies. Two previous studies have inferred trees which do not place *A. fowleri* samples into a single monophyletic group (Fontenot et al., 2011; Masta et al., 2002). A single *A. fowleri* sample included in this study does not fall within a monophyletic group with the remaining *A. fowleri* samples but is instead sister to the clade containing all *A. fowleri* and *A. woodhousii* samples Fig. 2.3.

All of these studies have included different species, individuals, and loci, and also used different methods for alignment and phylogenetic inference. These differences in study design could result in the observed topology differences. The choice of locus in particular has a high likelihood of being the cause of these difference. Due to incomplete lineage sorting, the true histories of each gene may in fact differ (Kingman, 1982). The practice of concatenating multiple loci as all of the previous studies of *Anaxyrus* evolutionary relationships have done, can produce erroneous trees with high statistical support (Kubatko & Degnan, 2007).

To account for incomplete lineage sorting, I also inferred phylogenetic relationships among *Anaxyrus* species using the multispecies coalescent method *phycoeval* along with

a subset of individuals used for the maximum likelihood tree due to increased computational demands of multispecies coalescent methods. The topology of the *phycoeval* tree is substantially different from the maximum likelihood tree inferred in this study as well as trees from previous studies Fig. 2.4 (Fontenot et al., 2011; Graybeal, 1997; Masta et al., 2002; Portik et al., 2023; Pramuk et al., 2007; Pyron & Wiens, 2011). Unlike in any previous study or in the maximum likelihood tree, *A. americanus* and *A. terrestris* are placed sister to one another, whereas in all other trees it has had closer affinity to the *A. hemiophrys/A. baxteri* clade Fig. 2.3 (Portik et al., 2023; Pyron & Wiens, 2011). In the *phycoeval* tree, the *A. hemiophrys/A. baxteri* clade is instead sister to the *A. americanus/A. fowleri/A. terrestris/A. woodhousii* clade.

An unusual feature of *phycoeval* is that it can allow for multifurcations in inferred topologies (Oaks et al., 2022). This feature proved relevant for in this study as the inferred tree included one multifurcation at the ancestral node of *A. quercicus*, the *A. cognatus/A. speciosus* clade, and the *americanus* group.

Previous studies have produced trees with quite short internode branches at this part of the tree as did the *IQ-TREE* analysis in this study which is somewhat consistent with this. These methods can only produce bifurcations and thus would force any true multifurcation into bifurcations and estimate some branch length between them which would be expected to be short. In the *phycoeval* tree, the posterior probability of this split is low (0.51) so it may not be a perfect representation of the history of these lineages Fig. 2.4. More data may be necessary to have full resolution in this part of the tree. But it is clear that these three lineages diverged at least in rapid succession if not simultaneously. But I don't know of any significant implications these alternative scenarios would have for our understanding of *Anaxyrus* evolution.

Given the interest of hybridization in *Anaxyrus* and the promise of this group for understanding speciation it is important to consider the implications of these relationships with regards to hybridization.

*A. americanus* and *A. terrestris* are sister *A. fowleri* and *A. woodhousii* are sister This makes more sense with regards to hybridization between these two pairs of species and

better explains the sympatry of *A. americanus* and *A. terrestris* with *A. fowleri* and *A. woodhousii*. This does make the hybrid zone between *A. americanus* and *A. hemiophyrs* more surprising. This hybrid zone is quite narrow however. 542  
543  
544

Considering the relationships of the maximum likelihood tree with respect to hybridization, it is unsurprising to see *A. fowleri* (excluding sample 006) and *A. woodhousii* to be each other's closest relatives and separated by relatively short branch lengths. 545  
546  
547

Not clear if *A. fowleri* sample placement is due to incomplete lineage sorting or due to admixture. It does not fit the pattern of previous studies and is located near the contact zone of *A. woodhousii* and *A. fowleri*. Furthermore the *STRUCTURE* analysis suggests this sample has a lot of admixture which I discuss later. 548  
549  
550  
551

*A. americanus* and *A. terrestris* on the other hand is more surprising. *A. terrestris* is more closely related to *A. fowleri* than *A. americanus*, yet their ranges completely overlap and they appear to have strong reproductive isolation. *A. americanus* has a range that overlaps significantly with both *A. woodhousii* and *A. fowleri* but also appears to have strong reproductive isolation with it's sympatric congeners. 552  
553  
554  
555  
556

The relationship between *A. hemiophyrs* and *A. americanus* 557

*A. hemiophyrs* and *A. americanus* are separated but much longer branch lengths on the other but do not share any close relatives with which they. But The relationship between *A. americanus* and *A. hemiophyrs* is also unsurprising. Despite not being each others The *A. americanus* and *A. hemiophyrs* as well as *A. americanus* and *A. terrestris* pairs do not 558  
559  
560  
561  
562

If we consider these relationships with regards to hybrid zones Interestingly, only a single species pair (*A. fowleri* and *A. woodhousii*) for which there is evidence of a hybrid zone are each other's closest relatives. The branch lengths between *A. americanus* and *A. hemiophyrs* as well as between *A. americanus* and *A. terrestris* are long relative to the branches separating *A. woodhousii* and *A. fowleri*. 563  
564  
565  
566  
567

*A. terrestris* and *A. fowleri* are 568

Notably, only one species pair (*A. fowleri* and *A. woodhousii*) for which there is evidence a hybrid zone are sister to each other. *A. americanus* and *A. hemiophyrs* are 569  
570

close to sister as *A. hemiophrys* and *A. baxteri* have little genetic difference between them  
571  
Fig. 2.2. It is interesting that *A. americanus* and *A. terrestris* form a hybrid zone while  
572  
*A. fowleri* and *A. terrestris* which are more closely related in sympatry.  
573

## 2.4.2 Divergence Time 574

Only two studies have attempted to estimate the ages of splits within *Anaxyrus* Feng  
575  
et al included just two species *Anaxyrus canorus* (not included in this study) and *A. punctatus*. The estimated split between these two species was 12.3 million years ago.  
576  
The estimated split between *Incilius* and *Anaxyrus* was 16.5 million years ago which was  
577  
used to scale the branch lengths of the tree inferred with *phycoeval*. Portick et al. using a  
578  
concatenation method with a smaller number of loci and greater amount of missing data  
579  
estimated the split between *Incilius* and *Anaxyrus* to be 20.3 mya.  
580

split between *A. punctatus* group and rest of *Anaxyrus* was estimated to be 14.7  
582

*A. americanus* group diversification largely took place during the Pleistocene. A  
583  
period marked by repeated glaciations. The *phycoeval* analysis does not support shared  
584  
divergence events during this period suggesting that diversification in toads was driven  
585  
by different  
586

## 2.4.3 Population Structure and Hybridization 587

Given the appearance that there are many secondary contact zones. It seems probable  
588  
that toad species have undergone range expansions. Following these range expansions,  
589  
are there any barriers that are now reducing gene flow? We can test that by looking  
590  
for population structure within species that aligns with possible biogeographic barriers.  
591  
The maximum likelihood tree along with the *STRUCTURE* analyses do not support the  
592  
existence of any unrecognized species diversity or significant population structure as some  
593  
mitochondrial studies have suggested.  
594

Population structure in *A. woodhousii* with two overlapping mtDNA clades with  
595  
one more associated with the Southwest and one more associated with the great plains  
596  
(Masta et al., 2003)  
597

<b>2.4.4 Conclusion</b>	598
<b>References</b>	599
Abbott, R., Albach, D., Ansell, S., Arntzen, J. W., Baird, S. J. E., Bierne, N., Boughman, J., Brelandsford, A., Buerkle, C. A., Buggs, R., Butlin, R. K., Dieckmann, U., Eroukhmanoff, F., Grill, A., Cahan, S. H., Hermansen, J. S., Hewitt, G., Hudson, A. G., Jiggins, C., ... Zinner, D. (2013). Hybridization and speciation. <i>Journal of Evolutionary Biology</i> , 26(2), 229–246. <a href="https://doi.org/10.1111/j.1420-9101.2012.02599.x">https://doi.org/10.1111/j.1420-9101.2012.02599.x</a>	600 601 602 603 604 605
Bayona-Vásquez, N. J., Glenn, T. C., Kieran, T. J., Pierson, T. W., Hoffberg, S. L., Scott, P. A., Bentley, K. E., Finger, J. W., Louha, S., Troendle, N., Diaz-Jaimes, P., Mauricio, R., & Faircloth, B. C. (2019). Adapterama III: Quadruple-indexed, double/triple-enzyme RADseq libraries (2RAD/3RAD). <i>PeerJ</i> , 7, e7724. <a href="https://doi.org/10.7717/peerj.7724">https://doi.org/10.7717/peerj.7724</a>	606 607 608 609 610
Blair, W. F. (1963). Intragroup genetic compatibility in the <i>Bufo americanus</i> species group of toads. <i>The Texas Journal of Science</i> , 13, 15–34.	611 612
Blair, W. F. (1972). <i>Evolution in the genus Bufo</i> . University of Texas Press.	613
Degnan, J. H., & Rosenberg, N. A. (2009). Gene tree discordance, phylogenetic inference and the multispecies coalescent. <i>Trends in Ecology &amp; Evolution</i> , 24(6), 332–340. <a href="https://doi.org/10.1016/j.tree.2009.01.009">https://doi.org/10.1016/j.tree.2009.01.009</a>	614 615 616
Evanno, G., Regnaut, S., & Goudet, J. (2005). Detecting the number of clusters of individuals using the software structure: A simulation study. <i>Molecular Ecology</i> , 14(8), 2611–2620. <a href="https://doi.org/10.1111/j.1365-294X.2005.02553.x">https://doi.org/10.1111/j.1365-294X.2005.02553.x</a>	617 618 619
Feng, Y.-J., Blackburn, D. C., Liang, D., Hillis, D. M., Wake, D. B., Cannatella, D. C., & Zhang, P. (2017). Phylogenomics reveals rapid, simultaneous diversification of three major clades of Gondwanan frogs at the Cretaceous–Paleogene boundary. <i>Proceedings of the National Academy of Sciences</i> , 114(29). <a href="https://doi.org/10.1073/pnas.1704632114">https://doi.org/10.1073/pnas.1704632114</a>	620 621 622 623 624

Fontenot, B. E., Makowsky, R., & Chippindale, P. T. (2011). Nuclear–mitochondrial discordance and gene flow in a recent radiation of toads. <i>Molecular Phylogenetics and Evolution</i> , 59(1), 66–80. <a href="https://doi.org/10.1016/j.ympev.2010.12.018">https://doi.org/10.1016/j.ympev.2010.12.018</a>	625
Francis, R. M. (2017). POPHELPER: An R package and web app to analyse and visualize population structure. <i>Molecular Ecology Resources</i> , 17(1), 27–32. <a href="https://doi.org/10.1111/1755-0998.12509">https://doi.org/10.1111/1755-0998.12509</a>	628
Graybeal, A. (1997). Phylogenetic relationships of bufonid frogs and tests of alternate macroevolutionary hypotheses characterizing their radiation. <i>Zoological Journal of the Linnean Society</i> , 119(3), 297–338. <a href="https://doi.org/10.1111/j.1096-3642.1997.tb00139.x">https://doi.org/10.1111/j.1096-3642.1997.tb00139.x</a>	631
Green, D. M. (1996). The bounds of species: Hybridization in the <i>Bufo americanus</i> group of North American toads. <i>Israel Journal of Zoology</i> , 42, 95–109.	635
Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q., & Vinh, L. S. (2018). UFBoot2: Improving the Ultrafast Bootstrap Approximation. <i>Molecular Biology and Evolution</i> , 35(2), 518–522. <a href="https://doi.org/10.1093/molbev/msx281">https://doi.org/10.1093/molbev/msx281</a>	637
Huerta-Cepas, J., Serra, F., & Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. <i>Molecular Biology and Evolution</i> , 33(6), 1635–1638. <a href="https://doi.org/10.1093/molbev/msw046">https://doi.org/10.1093/molbev/msw046</a>	640
Jombart, T. (2008). Adegenet : A R package for the multivariate analysis of genetic markers. <i>Bioinformatics</i> , 24(11), 1403–1405. <a href="https://doi.org/10.1093/bioinformatics/btn129">https://doi.org/10.1093/bioinformatics/btn129</a>	643
Kingman, J. (1982). The coalescent. <i>Stochastic Processes and their Applications</i> , 13(3), 235–248. <a href="https://doi.org/10.1016/0304-4149(82)90011-4">https://doi.org/10.1016/0304-4149(82)90011-4</a>	646
Kubatko, L. S., & Degnan, J. H. (2007). Inconsistency of Phylogenetic Estimates from Concatenated Data under Coalescence (T. Collins, Ed.). <i>Systematic Biology</i> , 56(1), 17–24. <a href="https://doi.org/10.1080/10635150601146041">https://doi.org/10.1080/10635150601146041</a>	648
Lyman, R. A., & Edwards, C. E. (2022). Revisiting the comparative phylogeography of unglaciated eastern North America: 15 years of patterns and progress. <i>Ecology and Evolution</i> , 12(4), e8827. <a href="https://doi.org/10.1002/ece3.8827">https://doi.org/10.1002/ece3.8827</a>	651

- Malinsky, M., Matschiner, M., & Svardal, H. (2021). Dsuite - Fast D-statistics and related admixture evidence from VCF files. *Molecular Ecology Resources*, 21(2), 584–595. 654  
<https://doi.org/10.1111/1755-0998.13265> 655
- Malinsky, M., Svardal, H., Tyers, A. M., Miska, E. A., Genner, M. J., Turner, G. F., & Durbin, R. (2018). Whole-genome sequences of Malawi cichlids reveal multiple radiations interconnected by gene flow. *Nature Ecology & Evolution*, 2(12), 1940–1955. <https://doi.org/10.1038/s41559-018-0717-x> 656
- Masta, S. E., Laurent, N. M., & Routman, E. J. (2003). Population genetic structure of the toad *Bufo woodhousii* : An empirical assessment of the effects of haplotype extinction on nested cladistic analysis. *Molecular Ecology*, 12(6), 1541–1554. <https://doi.org/10.1046/j.1365-294X.2003.01829.x> 660
- Masta, S. E., Sullivan, B. K., Lamb, T., & Routman, E. J. (2002). Molecular systematics, hybridization, and phylogeography of the *Bufo americanus* complex in Eastern North America. *Molecular Phylogenetics and Evolution*, 24(2), 302–314. [https://doi.org/10.1016/S1055-7903\(02\)00216-6](https://doi.org/10.1016/S1055-7903(02)00216-6) 661
- Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2015). IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology and Evolution*, 32(1), 268–274. <https://doi.org/10.1093/molbev/msu300> 662
- Nuñez, L. P., Gray, L. N., Weisrock, D. W., & Burbrink, F. T. (2023). The phylogenomic and biogeographic history of the gartersnakes, watersnakes, and allies (Natricidae: Thamnophiini). *Molecular Phylogenetics and Evolution*, 186, 107844. <https://doi.org/10.1016/j.ympev.2023.107844> 663
- Oaks, J. R. (2019). Full Bayesian Comparative Phylogeography from Genomic Data (L. Kubatko, Ed.). *Systematic Biology*, 68(3), 371–395. <https://doi.org/10.1093/sysbio/syy063> 664
- Oaks, J. R., Wood, P. L., Siler, C. D., & Brown, R. M. (2022). Generalizing Bayesian phylogenetics to infer shared evolutionary events. *Proceedings of the National Academy of Sciences*, 119(29), e2121036119. <https://doi.org/10.1073/pnas.2121036119> 665

Portik, D. M., Streicher, J. W., & Wiens, J. J. (2023). Frog phylogeny: A time-calibrated, species-level tree based on hundreds of loci and 5,242 species. <i>Molecular Phylogenetics and Evolution</i> , 188, 107907. <a href="https://doi.org/10.1016/j.ympev.2023.107907">https://doi.org/10.1016/j.ympev.2023.107907</a>	683
Pramuk, J. B., Robertson, T., Sites, J. W., & Noonan, B. P. (2007). Around the world in 10 million years: Biogeography of the nearly cosmopolitan true toads (Anura: Bufonidae). <i>Global Ecology and Biogeography</i> , 0(0), 070817112457001–???. <a href="https://doi.org/10.1111/j.1466-8238.2007.00348.x">https://doi.org/10.1111/j.1466-8238.2007.00348.x</a>	684
Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of Population Structure Using Multilocus Genotype Data. <i>Genetics</i> , 155(2), 945–959. <a href="https://doi.org/10.1093/genetics/155.2.945">https://doi.org/10.1093/genetics/155.2.945</a>	685
Pyron, R. A., & Wiens, J. J. (2011). A large-scale phylogeny of Amphibia including over 2800 species, and a revised classification of extant frogs, salamanders, and caecilians. <i>Molecular Phylogenetics and Evolution</i> , 61(2), 543–583. <a href="https://doi.org/10.1016/j.ympev.2011.06.012">https://doi.org/10.1016/j.ympev.2011.06.012</a>	686
Reich, D., Thangaraj, K., Patterson, N., Price, A. L., & Singh, L. (2009). Reconstructing Indian population history. <i>Nature</i> , 461(7263), 489–494. <a href="https://doi.org/10.1038/nature08365">https://doi.org/10.1038/nature08365</a>	687
Rieseberg, L. H., Whitton, J., & Gardner, K. (1999). Hybrid zones and the genetic architecture of a barrier to gene flow between two sunflower species. <i>Genetics</i> , 152(2), 713–727.	688
Rochette, N. C., Rivera-Colón, A. G., & Catchen, J. M. (2019). Stacks 2: Analytical methods for paired-end sequencing improve RADseq-based population genomics. <i>Molecular Ecology</i> , 28(21), 4737–4754. <a href="https://doi.org/10.1111/mec.15253">https://doi.org/10.1111/mec.15253</a>	689
Wang, L.-G., Lam, T. T.-Y., Xu, S., Dai, Z., Zhou, L., Feng, T., Guo, P., Dunn, C. W., Jones, B. R., Bradley, T., Zhu, H., Guan, Y., Jiang, Y., & Yu, G. (2020). Treeio: An R Package for Phylogenetic Tree Input and Output with Richly Annotated and Associated Data. <i>Molecular Biology and Evolution</i> , 37(2), 599–603. <a href="https://doi.org/10.1093/molbev/msz240">https://doi.org/10.1093/molbev/msz240</a>	690

Weatherby, C. A. (1982). INTROGRESSION BETWEEN THE AMERICAN TOAD, BUFO AMERICANUS, AND THE SOUTHERN TOAD, B. TERRESTRIS, IN ALABAMA.	711
Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T.-Y. (2017). Ggtree: An r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. <i>Methods in Ecology and Evolution</i> , 8(1), 28–36. <a href="https://doi.org/10.1111/2041-210X.12628">https://doi.org/10.1111/2041-210X.12628</a>	712
	713
	714
	715
	716
	717

## 2.5 Figures

718

## Sampling Distribution

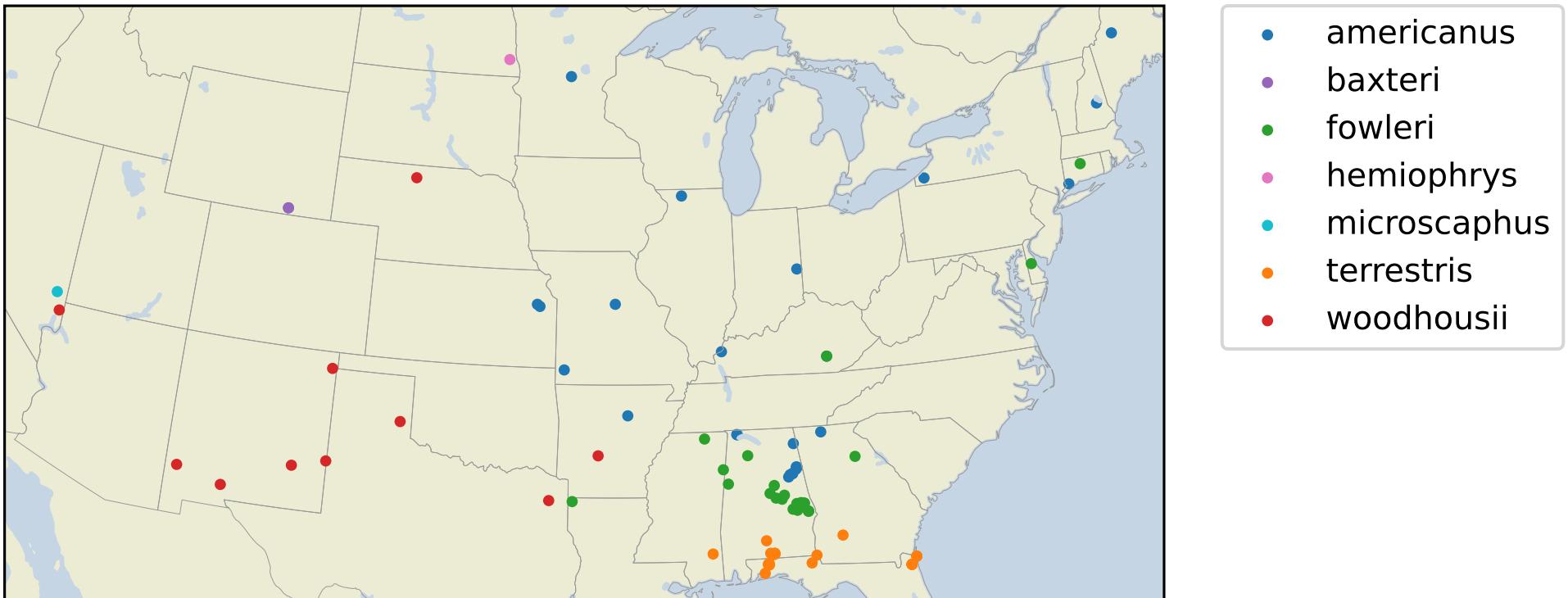


Figure 2.1. Distribution of *americanus* group samples

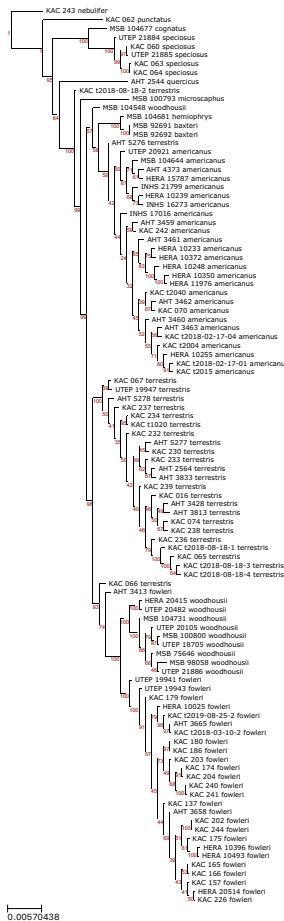


Figure 2.2. Maximum Likelihood Phylogeny Plotted using ETE 3.1.2 (Huerta-Cepas et al., 2016).

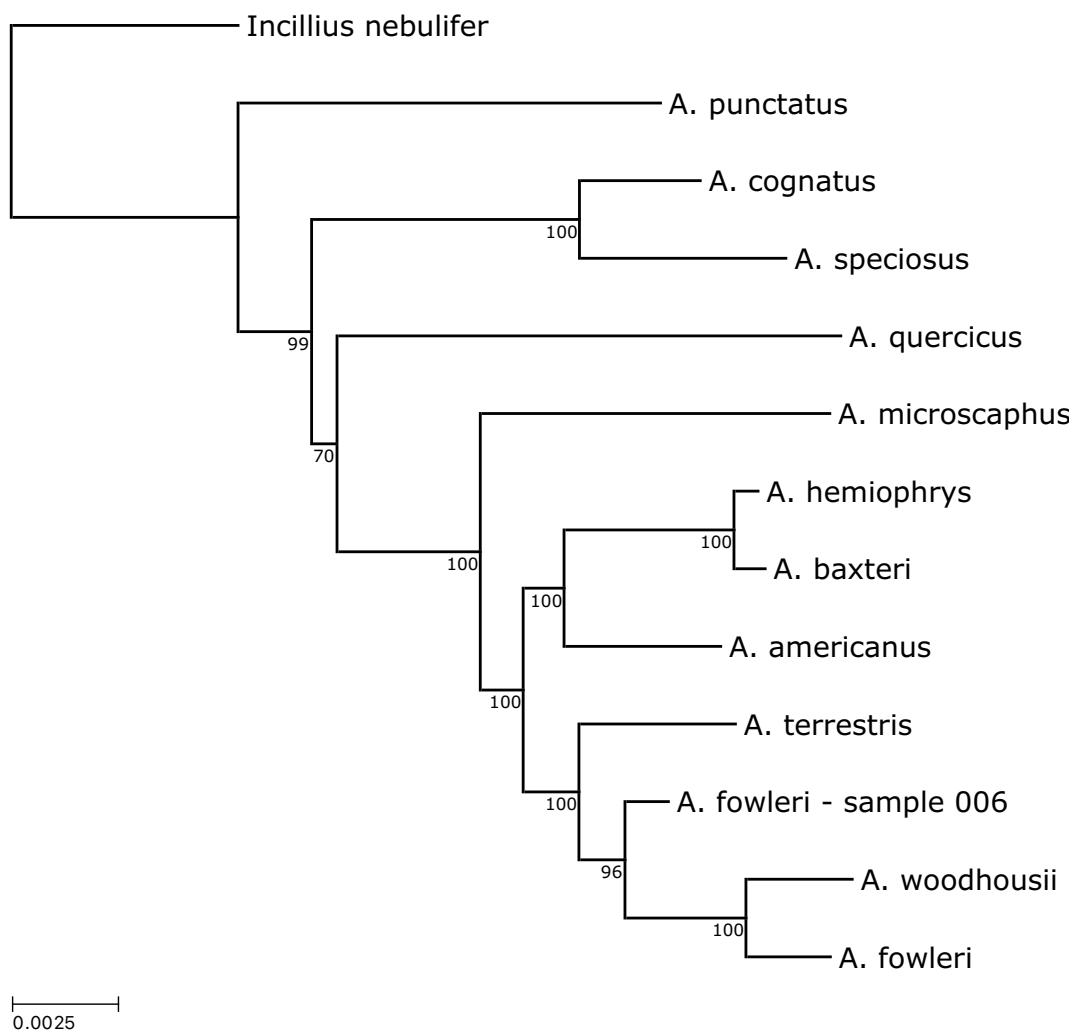
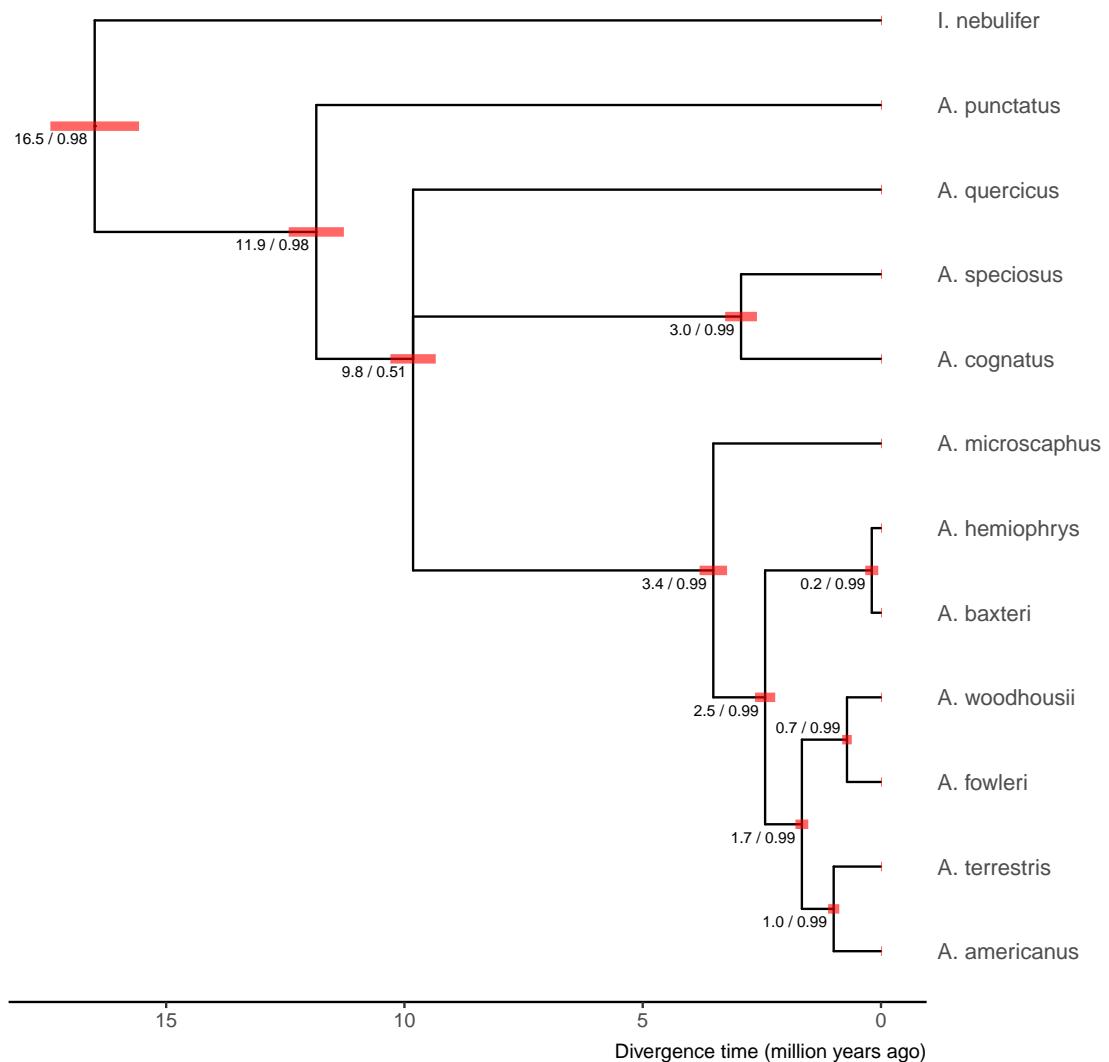


Figure 2.3. Maximum Likelihood Phylogeny with species clades collapsed. The lengths of tip branches are equal to the mean height of all collapsed tips from the base the collapsed clade. Plotted using ETE 3.1.2 (Huerta-Cepas et al., 2016).



**Figure 2.4.** The maximum a posteriori tree inferred under a multispecies coalescent model by *phycoeval*. Branch lengths are rescaled from expected substitutions per site to millions of years using secondary time calibrations (*Materials and Methods*). Numbers displayed at each node are the mean posterior node age followed by the approximate posterior probability of the node rounded down to the nearest hundredth. Red bars show the 95% HPDI for the scaled node age at each node. Created using ggplot2 (wickham2016), ggtree (Yu et al., 2017), and treeio (Wang et al., 2020)

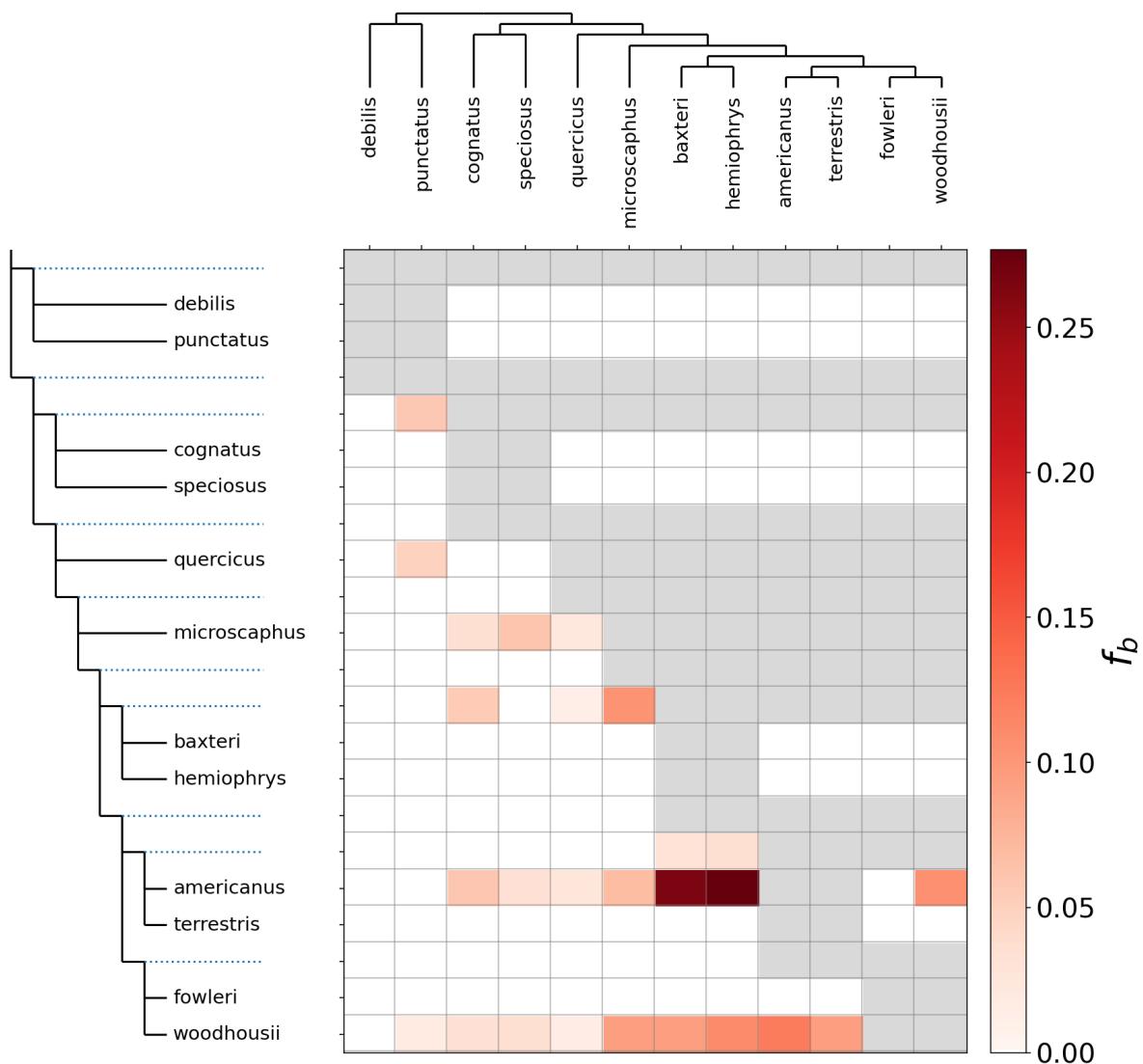


Figure 2.5. Heatmap showing the value of the  $f$ -branch statistic computed for all pairs of possible pairs of *Anaxyrus* species. The  $f$ -branch statistic indicates the proportion of excess allele sharing between a species on the x-axis and branch on the y-axis (relative to its sister branch). Excess allele sharing between species identifies possible gene flow between them. Grey boxes indicate that the given tips cannot be tested by Dsuite for the given tree topology.

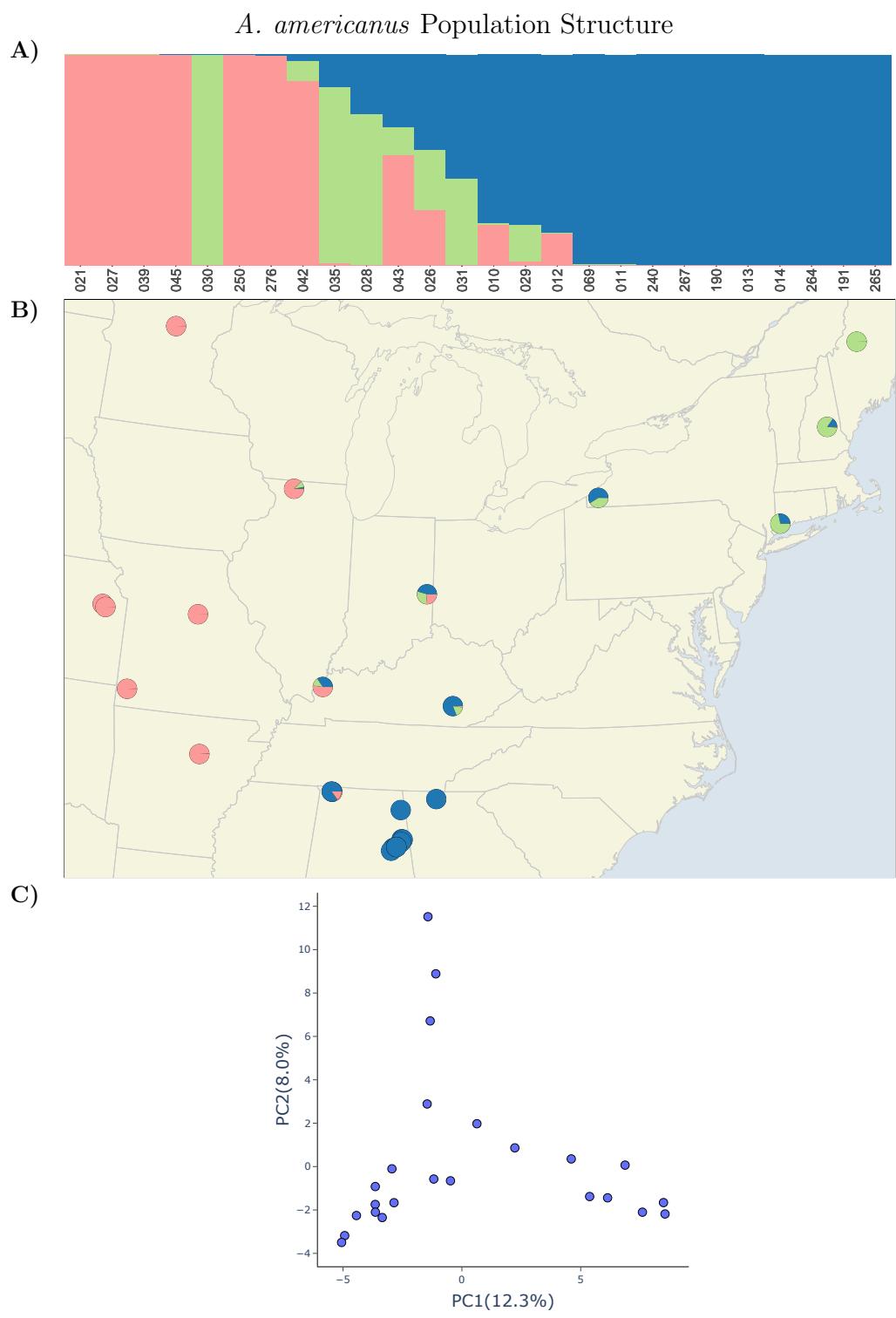


Figure 2.6. Population structure of *A. americanus* with PCA and map

*A. fowleri* Population Structure

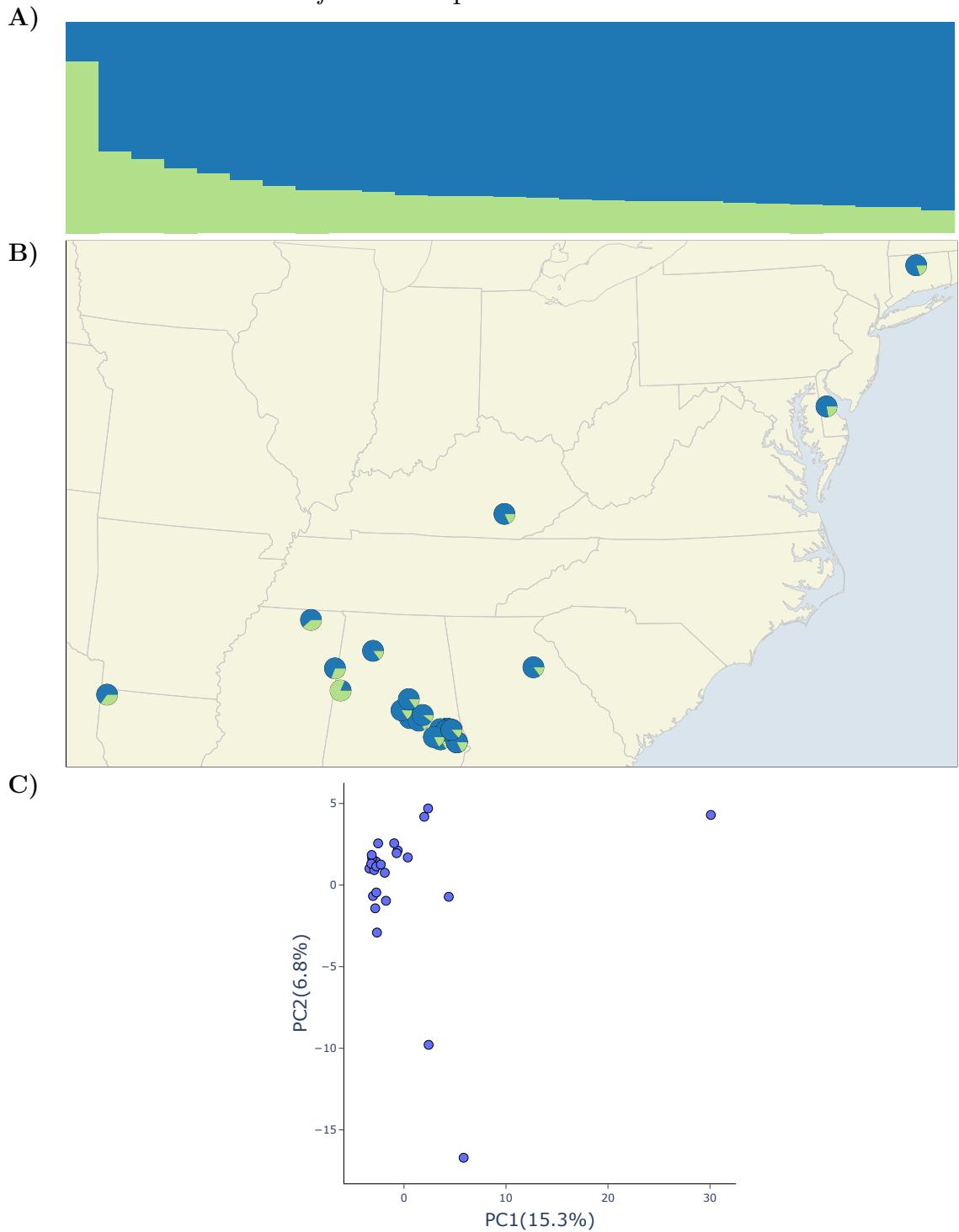


Figure 2.7. Population structure of *A. fowleri* with PCA and map

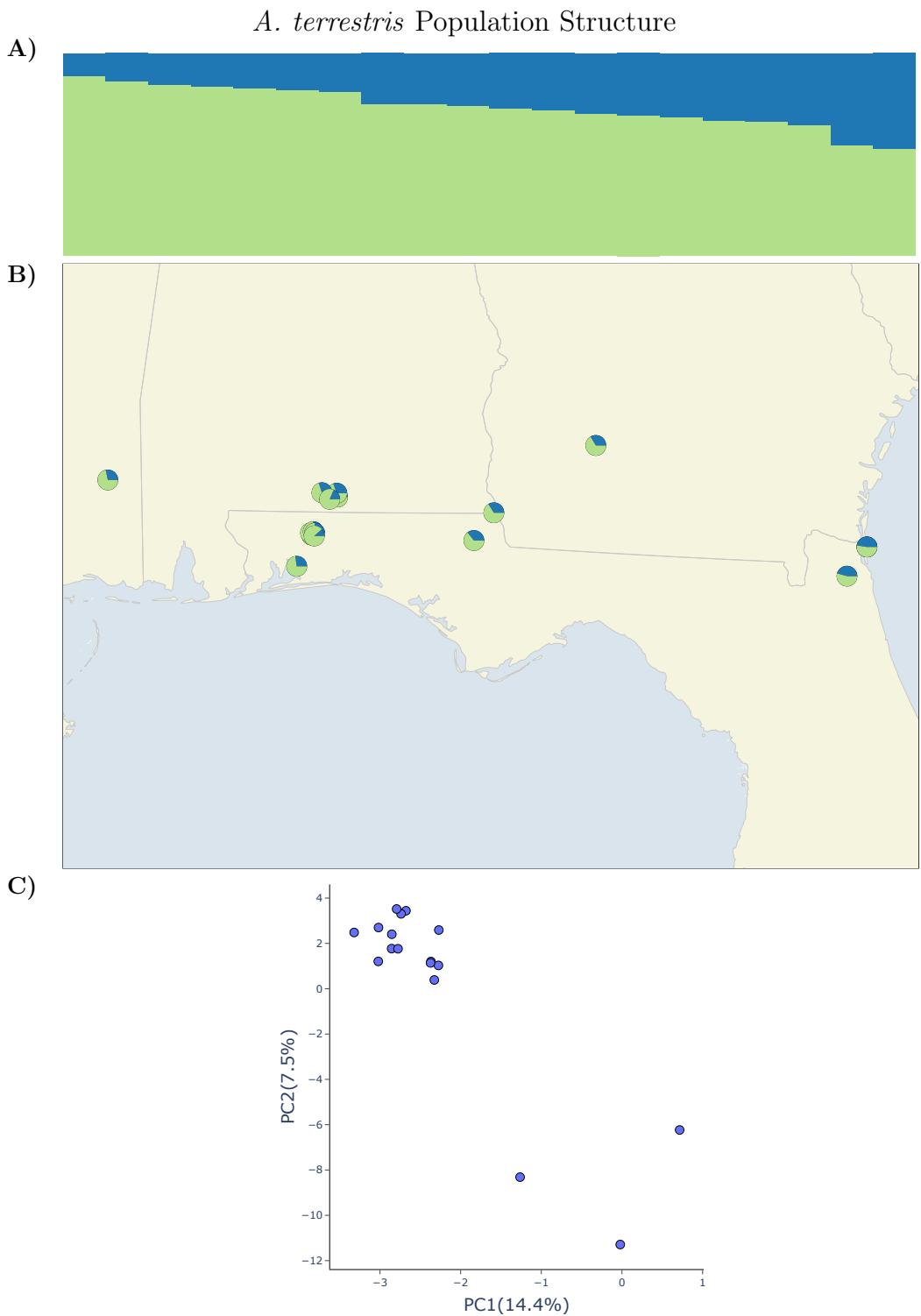


Figure 2.8. Population structure of *A. terrestris* with PCA and map

*A. woodhousii* Population Structure

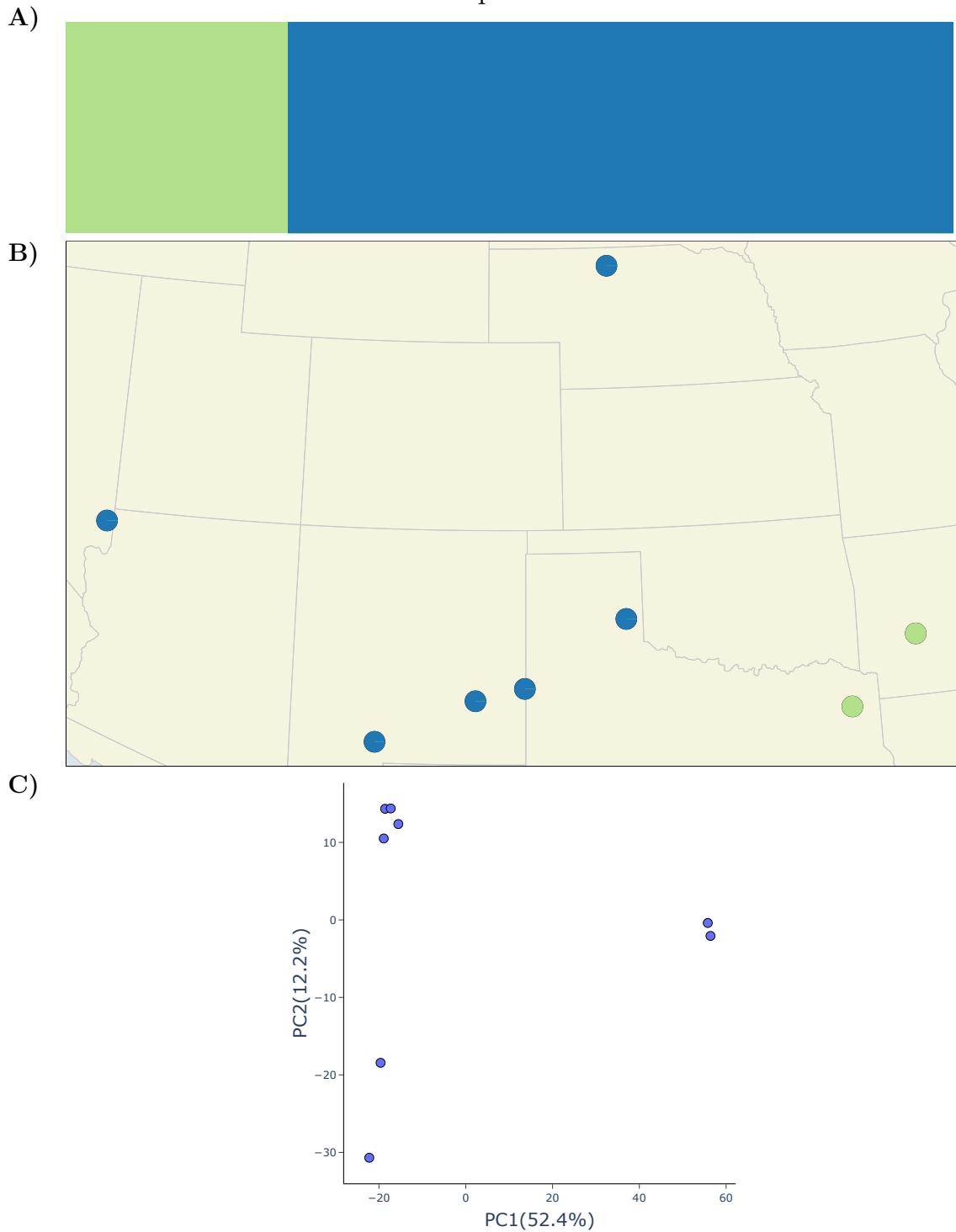


Figure 2.9. Population structure of *A. woodhousii* with PCA and map

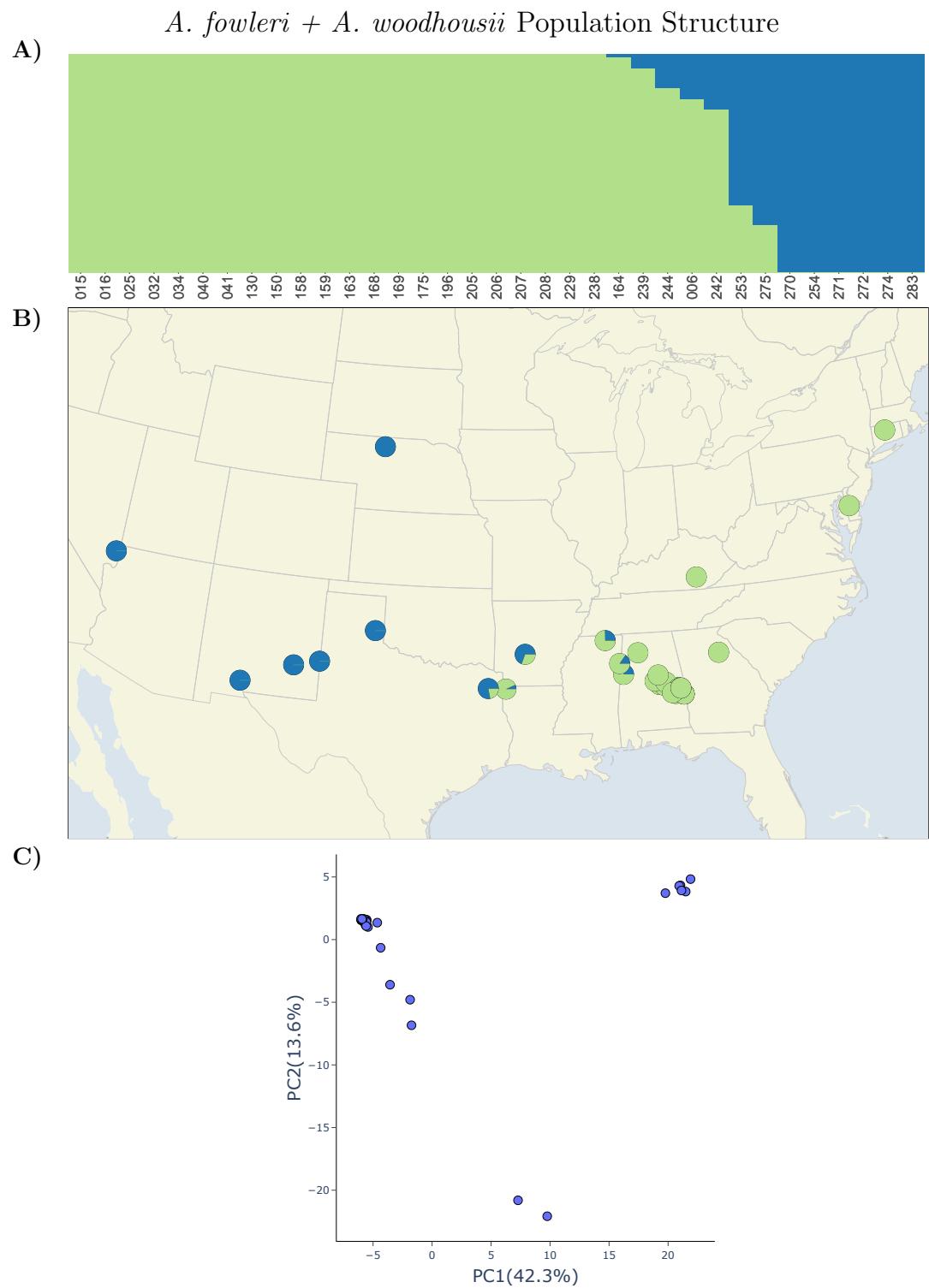


Figure 2.10. Population structure of *A. fowleri* with PCA and map

## 2.6 Tables

719

Table 2.1. Samples used in this study

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
003	AHT 2544	<i>quercicus</i>	30.99523	-86.23332	X	X	
004	AHT 2564	<i>terrestris</i>	31.55752	-84.04267	X	X	X
006	AHT 3413	<i>fowleri</i>	33.36940	-88.12941	X		X
009	AHT 3428	<i>terrestris</i>	31.12679	-86.54755	X		X
010	AHT 3459	<i>americanus</i>	34.88028	-87.71849	X		X
011	AHT 3460	<i>americanus</i>	33.78013	-85.58421	X		X
012	AHT 3461	<i>americanus</i>	34.88779	-87.74103	X		X
013	AHT 3462	<i>americanus</i>	33.77001	-85.55434	X		X
014	AHT 3463	<i>americanus</i>	33.71125	-85.59762	X		X
015	AHT 3658	<i>fowleri</i>	32.85842	-86.39697	X		X
016	AHT 3665	<i>fowleri</i>	32.81220	-86.17698	X		X
017	AHT 3813	<i>terrestris</i>	31.13854	-86.53906	X		
018	AHT 3833	<i>terrestris</i>	31.00422	-85.03427	X		X
021	AHT 4373	<i>americanus</i>	38.94913	-95.39818	X		X
022	AHT 5276	<i>terrestris</i>	31.55613	-86.82514			
023	AHT 5277	<i>terrestris</i>	31.15830	-86.55430	X		X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
024	AHT 5278	<i>terrestris</i>	31.16105	-86.69868	X		X
025	HERA 10025	<i>fowleri</i>	37.11151	-84.11812	X	X	X
026	HERA 10233	<i>americanus</i>	39.86453	-85.01037	X	X	X
027	HERA 10239	<i>americanus</i>	38.99151	-92.31078	X		X
028	HERA 10248	<i>americanus</i>	41.27319	-73.38974	X		X
029	HERA 10255	<i>americanus</i>	37.11151	-84.11812	X		X
030	HERA 10350	<i>americanus</i>	45.51396	-69.95928	X	X	X
031	HERA 10372	<i>americanus</i>	42.22795	-79.36759	X		X
032	HERA 10396	<i>fowleri</i>	41.80663	-72.73281	X	X	X
033	HERA 10484	<i>marina</i>	25.61296	-80.56606			
034	HERA 10493	<i>fowleri</i>	39.08588	-75.56844	X	X	X
035	HERA 11976	<i>americanus</i>	43.51819	-71.42336	X		X
036	HERA 13722	<i>fowleri</i>	36.55514	-89.18929			
037	HERA 14196	<i>retiformis</i>	33.34906	-112.49010			
038	HERA 14926	<i>microscaphus</i>	33.73033	-113.98078			
039	HERA 15787	<i>americanus</i>	38.88546	-95.29399	X	X	X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
040	HERA 20415	<i>woodhousii</i>	34.31743	-92.94602	X	X	X
041	HERA 20514	<i>fowleri</i>	33.95140	-83.36715	X		X
042	INHS 16273	<i>americanus</i>	42.30245	-89.55950	X		X
043	INHS 17016	<i>americanus</i>	37.46121	-88.18728	X		X
044	INHS 19127	<i>fowleri</i>	41.58247	-88.07273			
045	INHS 21799	<i>americanus</i>	46.01258	-94.26710	X		X
046	KAC 016	<i>terrestris</i>	30.54819	-86.93067	X		X
061	KAC 053	<i>fowleri</i>	32.78044	-86.73877			
062	KAC 060	<i>speciosus</i>	27.69185	-99.71955	X		
063	KAC 062	<i>punctatus</i>	29.43603	-103.50564	X		
064	KAC 063	<i>speciosus</i>	29.29522	-103.92916	X		
065	KAC 064	<i>speciosus</i>	29.29522	-103.92916	X		
066	KAC 065	<i>terrestris</i>	30.43282	-81.64088	X		
067	KAC 066	<i>terrestris</i>	30.43282	-81.64088			
068	KAC 067	<i>terrestris</i>	30.43282	-81.64088			
069	KAC 070	<i>americanus</i>	34.79963	-84.57678	X		X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
071	KAC 074	<i>terrestris</i>	30.77430	-85.22690	X		X
130	KAC 137	<i>fowleri</i>	33.01461	-86.60953	X		X
150	KAC 157	<i>fowleri</i>	32.43769	-85.63620	X		X
158	KAC 165	<i>fowleri</i>	32.66356	-85.48498	X		X
159	KAC 166	<i>fowleri</i>	32.66356	-85.48498	X		X
163	KAC 174	<i>fowleri</i>	32.62938	-85.63828	X		X
164	KAC 175	<i>fowleri</i>	32.64849	-85.64711	X		X
167	KAC 178	<i>fowleri</i>	32.38644	-85.23561			
168	KAC 179	<i>fowleri</i>	32.38644	-85.23561	X		X
169	KAC 180	<i>fowleri</i>	32.38644	-85.23561	X		X
175	KAC 186	<i>fowleri</i>	32.38579	-85.23565	X		X
190	KAC t2018-02-17-01	<i>americanus</i>	33.55274	-85.82913	X		X
191	KAC t2018-02-17-04	<i>americanus</i>	33.48548	-85.88857	X		X
196	KAC t2018-03-10-2	<i>fowleri</i>	32.93116	-86.08465	X		X
200	KAC t2018-08-18-1	<i>terrestris</i>	30.66902	-81.44013	X		X
201	KAC t2018-08-18-2	<i>terrestris</i>	30.66902	-81.44013			

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
202	KAC t2018-08-18-3	<i>terrestris</i>	30.43282	-81.64088	X	X	X
203	KAC t2018-08-18-4	<i>terrestris</i>	30.66902	-81.44013	X		X
205	KAC t2019-08-25-2	<i>fowleri</i>	34.21852	-87.36662	X		X
206	KAC 202	<i>fowleri</i>	33.25104	-86.43850	X		X
207	KAC 203	<i>fowleri</i>	32.62294	-85.49660	X		X
208	KAC 204	<i>fowleri</i>	32.62294	-85.49660	X		X
229	KAC 226	<i>fowleri</i>	32.48119	-85.79838	X		X
230	KAC 230	<i>terrestris</i>	30.80933	-86.77686	X		X
231	KAC 232	<i>terrestris</i>	30.80922	-86.78994	X		X
231	KAC 232	<i>terrestris</i>	30.80922	-86.78994	X		X
232	KAC 233	<i>terrestris</i>	30.80922	-86.78994	X		X
233	KAC 234	<i>terrestris</i>	30.80922	-86.78994	X		X
234	KAC 236	<i>terrestris</i>	30.82632	-86.80258	X		X
235	KAC 237	<i>terrestris</i>	30.83733	-86.77630	X		X
236	KAC 238	<i>terrestris</i>	30.82433	-86.76284	X		X
237	KAC 239	<i>terrestris</i>	30.80162	-86.76659	X		X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
238	KAC 240	<i>fowleri</i>	32.64328	-85.37114	X		X
239	KAC 241	<i>fowleri</i>	32.64328	-85.37114	X		X
240	KAC 242	<i>americanus</i>	34.50446	-85.63768	X		X
241	KAC 243	<i>nebulifer</i>	30.39140	-90.62049	X	X	
242	KAC 244	<i>fowleri</i>	32.89261	-93.88756	X		X
243	MSB 100793	<i>microscaphus</i>	37.27154	-114.46478	X	X	
244	MSB 100800	<i>woodhousii</i>	36.73612	-114.21972	X	X	X
245	MSB 100913	<i>microscaphus</i>	33.28038	-108.08868		X	
246	MSB 104548	<i>woodhousii</i>	36.49094	-103.20838			
247	MSB 104570	<i>fowleri</i>	34.00087	-95.38229			
248	MSB 104571	<i>americanus</i>	34.00917	-95.38058			
249	MSB 104608	<i>americanus</i>	34.00367	-94.82670			
250	MSB 104644	<i>americanus</i>	36.95124	-94.27782	X		X
251	MSB 104677	<i>cognatus</i>	46.39834	-97.20927	X	X	
252	MSB 104681	<i>hemiophrys</i>	46.47076	-97.04604	X		X
253	MSB 104731	<i>woodhousii</i>	42.61091	-100.65607	X	X	X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
254	MSB 75646	<i>woodhousii</i>	33.36365	-104.34282	X	X	X
255	MSB 92689	<i>baxteri</i>	41.21182	-105.82558			
256	MSB 92691	<i>baxteri</i>	41.21182	-105.82558	X	X	
257	MSB 92692	<i>baxteri</i>	41.21182	-105.82558	X	X	
258	MSB 96528	<i>debilis</i>	32.58239	-107.46348			
259	MSB 98058	<i>woodhousii</i>	32.83360	-108.60900			
260	MSB 98065	<i>cognatus</i>	32.63240	-108.73800		X	
261	KAC t1020	<i>terrestris</i>	31.10783	-86.62247	X		X
264	KAC t2004	<i>americanus</i>	33.58295	-85.73524	X		X
265	KAC t2015	<i>americanus</i>	33.58435	-85.74064	X		X
267	KAC t2040	<i>americanus</i>	33.58295	-85.73539	X		X
269	KAC t3040	<i>fowleri</i>	32.38644	-85.23561			
270	UTEP 18705	<i>woodhousii</i>	32.45198	-106.88317	X	X	X
271	UTEP 19941	<i>fowleri</i>	34.79137	-88.95715	X	X	X
272	UTEP 19943	<i>fowleri</i>	33.81998	-88.29533	X		X
273	UTEP 19947	<i>terrestris</i>	31.22432	-88.77548	X	X	X

Continued on next page

Table 2.1 – continued from previous page

ID	Sample ID	Species	Latitude	Longitude	Passed Filtering	Phycoeval	Structure
274	UTEP 20105	<i>woodhousii</i>	33.62853	-103.08198	X		X
275	UTEP 20482	<i>woodhousii</i>	32.90708	-94.74945	X		X
276	UTEP 20921	<i>americanus</i>	35.55405	-91.83443	X		X
277	UTEP 21284	<i>debilis</i>	31.25968	-105.33402		X	
278	UTEP 21286	<i>speciosus</i>	31.70140	-105.47958			
279	UTEP 21724	<i>speciosus</i>	31.26087	-104.60168			
280	UTEP 21881	<i>cognatus</i>	35.53600	-100.44035		X	
281	UTEP 21884	<i>speciosus</i>	32.75472	-101.43208	X		
282	UTEP 21885	<i>speciosus</i>	32.20195	-100.34345	X		X
283	UTEP 21886	<i>woodhousii</i>	35.07800	-100.43392	X		X

# Chapter 3

720

## Hybrid Zone

721

### 3.1 Introduction

722

Speciation is the process by which genetic divergence leads to reproductive isolation between divergent lineages. It is a continuous process during which there may be ongoing gene flow or introgression via hybridization following a period of isolation and subsequent secondary contact (Mallet, 2008; Wu, 2001). Introgression is possible because genetic barriers to introgression that accumulate within the genome are a property of genomic regions rather than a property of the entirety of the genome (Gompert, Parchman, et al., 2012; Wu, 2001). Natural hybridization between divergent lineages has become increasingly appreciated as a widespread phenomenon in recent years (Mallet, 2005; Moran et al., 2021). It is a phenomenon that can have important evolutionary consequences. Hybridization can be a source of adaptive variation (Hedrick, 2013). It can also introduce deleterious genetic load which persists long term within a population (Moran et al., 2021). Hybridization can create conditions where selection favors the evolution of traits that enhance assortative mating and reduce the production of unfit hybrid offspring which drives further genetic divergence and reinforcement of reproductive barriers between lineages (Servedio & Noor, 2003). If hybrids do not suffer any negative fitness effects, hybridization could lead to the erosion of differences between divergent populations (Taylor et al., 2006). Potentially resulting in populations that are genetically distinct

from either parent species which can themselves eventually evolve reproductive isolation  
740  
from the parent species (Moran et al., 2021).  
741

Aside from having important evolutionary consequences which need to be understood,  
742  
hybridization is also an excellent opportunity to investigate the processes that result in  
743  
the evolution of reproductive incompatibility and divergence between evolutionary lin-  
744  
eages. Hybrid zones are particularly suitable due to the production of a large numbers  
745  
of recombinant genomes carrying many possible combinations of genomic elements from  
746  
parent species resulting from many generations of backcrossing (Rieseberg et al., 1999).  
747  
The many generations of backcrossing and recombination make it possible to distinguish  
748  
between the effects of closely linked genes (Rieseberg et al., 1999). The many generations  
749  
and large number of individuals producing these genetic combinations are not feasible  
750  
to produce experimentally in the vast majority of species (Rieseberg et al., 1999). Fur-  
751  
thermore, the combination of genes produced are exposed to selection under natural  
752  
conditions. This is important as the effect of hybrid incompatibilities can be dependent  
753  
on on environmental conditions and can only truly understood in this context (Miller &  
754  
Matute, 2016).  
755

Despite being a fundamental evolutionary process, our understanding of speciation  
756  
is far from complete (Butlin et al., 2011). Only a few loci, in a few species, have been  
757  
pinpointed as the direct cause of reproductive incompatibility between species (Black-  
758  
man, 2016; Nosil & Schlüter, 2011). Consequently, our understanding of the processes  
759  
that drive the evolution of loci resulting in reproductive incompatibility is limited (Butlin  
760  
et al., 2011). Studies of introgression within hybrid zones have identified highly variable  
761  
rates of introgression among loci (Barton & Hewitt, 1985; Gompert et al., 2017). This  
762  
heterogeneity can arise via genetic drift occurring within hybrid zones, but will also be  
763  
caused by differences among loci in the strength of selection against them in a hybrid  
764  
genomic background (Barton & Hewitt, 1985; Gompert et al., 2017). It has also been ob-  
765  
served that the levels of genetic divergence between species are highly variable across the  
766  
genome (Nosil et al., 2009). Much of this heterogeneity is the result of divergent selection  
767  
acting on each species independently (Nosil et al., 2009). Regions with particularly high  
768

levels of divergence between closely related species have been coined "genomic islands of divergence" (Wolf & Ellegren, 2017). It is assumed, particularly in the case of speciation with gene flow, that these genomic island harbor genes that reduce interbreeding between species. When speciation occurs with gene flow, divergent selection can cause adaptive divergence in habitat use, phenology, or mating signals, and reduce the frequency or success of interspecific matings. When species diverge in geographic isolation, divergent selection and reproductive isolation could be decoupled and reproductive isolation could just be a result of genetic drift. Whether loci under divergent selection between two species also contribute to reproductive isolation has not been widely explored. A handful of studies have found evidence for a modest relationship between genetic divergence and selection against introgression (Gompert, Lucas, et al., 2012; Larson et al., 2013; Nikolakis et al., 2022; Parchman et al., 2013). How consistent and widespread this pattern is remains to be seen. At least one study has found no association (Jahner et al., 2021).

In this study I investigate hybridization between the American toad (*Anaxyrus americanus*) and Southern toad (*Anaxyrus terrestris*) at a suspected hybrid zone in the Southern United States to assess the extent of introgression between them and test for a relationship between introgression and genetic divergence. This suspected hybrid zone has not been investigated with genetic data previously but it bears many hallmarks of a tension zone (Barton & Hewitt, 1985). Under the tension zone model of hybridization, species boundaries are maintained by a balance between dispersal and selection against individuals carrying incompatible hybrid genotypes (Barton & Hewitt, 1985). The ranges of *A. americanus* and *A. terrestris* abut with an abrupt transition and no apparent overlap along a long contact zone which from Louisiana to Virginia. This contact zone closely corresponds with a prominent physiographic feature known as the "fall line" (Mount, 1975; Powell et al., 2016). The Fall line is the boundary between the Southern coastal plain to the South and the Appalachian Highlands to the North (Shankman & Hart, 2007). These regions differ in their underlying geology, topography, and elevation (Shankman & Hart, 2007). The distribution of *A. terrestris* is restricted to the coastal plain extending from the Mississippi River in the West to Virginia in the East (Fig. 3.4).

The distribution of the American Toad encompasses nearly all of the Eastern North American with the exception of the Southern coastal plain (Fig. 3.4). Tension zones are expected to correspond with natural features that reduce dispersal or abundance (Barton, 1979). Such a sudden transition is difficult to explain if not the result of the processes characteristic of tension zones. For there to be no mutually hospitable areas permitting some range overlap is implausible without there being an extreme level of competition or extreme degree of adaptation by each species to their respective environments. The two species have only slight differences in male advertisement call and in morphological appearance (Cocroft & Ryan, 1995; Weatherby, 1982). They only differ slightly in the timing of their spawn (Mount, 1975) However, there is some overlap in the spawning period and male Bufonidae are famously indiscriminate in their choice of mates (Đorđević & Simović, 2014; Weatherby, 1982). They have also been shown to have a degree of reproductive compatibility through laboratory crossing experiments which produced viable  $F_2$  offspring (Blair, 1963). Analysis of morphological variation in central Alabama by Weatherby, 1982 suggests there has been introgression between them.

The "true toads" in the family Bufonidae, to which *A. americanus* and *A. terrestris* belong, have been a prominent group of organisms in the literature on hybridization. W.F. Blair and colleagues performed a remarkable 1,934 separate experimental crosses to quantify the degree of reproductive incompatibility between species pairs within this family (Blair, 1972; Malone & Fontenot, 2008). These experiments demonstrated a high degree of compatibility between some closely related species pairs in which hybrids were capable of producing viable backcross or  $F_2$  hybrid offspring (Blair, 1963). Furthermore, numerous cases of natural hybridization among toad species have been reported with several apparent or clear hybrid zones (Colliard et al., 2010; Green, 1996; Van Riemsdijk et al., 2023; Weatherby, 1982). Despite the interest in and appreciation for hybridization in Bufonidae, only a small amount of work has been done to understand patterns of introgression within hybrid zones. A clinal pattern of admixture at 26 allozyme loci has been shown within the *Anaxyrus americanus* X *Anaxyrus hemiophrys* hybrid zone in Ontario, Canada(Green, 1983). Almost no admixture was detected at 7 microsatellite loci

within the suspected *Bufo siculus* X *Bufo balearicus* hybrid zone in Sicily, Italy (Colliard et al., 2010). The most comprehensive study of introgression within a Bufonidae hybrid zone found significant levels of genome wide admixture, fitting a clinal pattern, at two separate transects at either end of the *Bufo bufo* x *Bufo spinosus* hybrid zone in Southern France (Van Riemsdijk et al., 2023).

The suspected *A. americanus*, *A. terrestris* hybrid zone has great potential to expand our understanding of speciation. This will be dependent on the degree of ongoing introgression, if any, between these species. In this study I use genome-wide sequence data to characterize patterns of introgression within the hybrid zone using model based inference of admixture proportions, bayesian genomic cline analysis, and estimates of parental population differentiation. With these approaches I specifically address the following questions: 1) Is there evidence of ongoing hybridization and admixture between the two species, 2) Do any loci have outstanding patterns of introgression consistent with them being linked to reproductive incompatibility, and 3) Is there any relationship between patterns of introgression and levels of genetic differentiation between parental lineages?

## 3.2 Methods

### 3.2.1 Sampling and DNA Isolation

I collected genetic samples from *A. americanus* and *A. terrestris* by driving roads during rainy nights between 2017 and 2020 in a region of central Alabama where hybridization has previously been inferred from the presence of morphological intermediates (Weatherby, 1982). I euthanized individuals with immersion in buffered MS-222. I removed liver and/or toes and preserved them in 100% ethanol and fixed specimens with 10% Formalin solution. Genetic samples and formalin fixed specimens were deposited at the Auburn Museum of Natural History. Additional samples were also provided by museums (see ??).

I isolated DNA by lysing a small piece of liver or toe approximately the size of a grain

of rice in 300  $\mu$ L of a solution of 10mM Tris-HCL, 10mM EDTA, 1% SDS (w/v), and  
nuclease free water along with 6 mg Proteinase K and incubating for 4-16 hours at 55°C.  
To purify the DNA and separate it from the lysis product, I mixed the lysis product  
with a 2X volume of SPRI bead solution containing 1 mM EDTA, 10 mM Tris-HCl,  
1 M NaCl, 0.275% Tween-20 (v/v), 18% PEG 8000 (w/v), 2% Sera-Mag SpeedBeads  
(GE Healthcare PN 65152105050250) (v/v), and nuclease free water. I then incubated  
the samples at room temperature for 5 minutes, placed the beads on a magnetic rack,  
and discarded the supernatant once the beads had collected on the side of the tube. I  
then performed two ethanol washes by adding 1 mL of 70% ETOH to the beads while  
still placed in the magnet stand and allowing it to stand for 5 minutes before removing  
and discarding the ethanol. After removing all ethanol from the second wash, I removed  
the tube from the magnet stand and allowed the sample to dry for 1 minute before  
thoroughly mixing the beads with 100  $\mu$ L of TLE solution containing 10 mM Tris-HCL,  
0.1 mM EDTA, and nuclease free water. After allowing the bead mixture to stand at  
room temperature for 5 minutes I returned the beads to the magnet stand, collected the  
TLE solution, and discarded the beads. I quantified DNA in the TLE solution with a  
Qubit fluorometer (Life Technologies, USA) and diluted samples with additional TLE  
solution to bring the concentration to 20 ng/ $\mu$ L.

### 3.2.2 RADseq Library Preparation

I prepared RADseq libraries using the 2RAD approach developed by Bayona-Vásquez  
et al., 2019. On 96 well plates, I ligated 100 ng of sample DNA in 15  $\mu$ L of a solution  
with 1X CutSmart Buffer (New England Biolabs, USA; NEB), 10 units of XbaI, 10 units  
of EcoRI, 0.33  $\mu$ M XbaI compatible adapter, 0.33  $\mu$ M EcoRI compatible adapter, and  
nuclease free water with a 1 hour incubation at 37°C. I then immediately added 5  $\mu$ L of  
a solution with 1X Ligase Buffer (NEB), 0.75 mM ATP (NEB), 100 units DNA Ligase  
(NEB), and nuclease free water and incubated at 22°C for 20 min and 37°C for 10 min for  
two cycles, followed by 80°C for 20 min to stop enzyme activity. For each 96 well plate,  
I pooled 10  $\mu$ L of each sample and split this pool equally between two microcentrifuge

tubes. I purified each pool of libraries with a 1X volume of SpeedBead solution followed  
882 by two ethanol washes as described in the previous section except that the DNA was  
883 resuspended in 25  $\mu$ L of TLE solution and combined the two pools of cleaned ligation  
884 product.  
885

In order to be able to detect and remove PCR duplicates, I performed a single cycle  
886 of PCR with the iTru5-8N primer which adds a random 8 nucleotide barcode to each  
887 library construct. For each plate, I prepared four PCR reactions with a total volume of  
888 50  $\mu$ L containing 1X Kapa Hifi Buffer (Kapa Biosystems, USA; Kapa), 0.3  $\mu$ M iTru5-8N  
889 Primer, 0.3 mM dNTP, 1 unit Kapa HiFi DNA Polymerase, 10  $\mu$ L of purified ligation  
890 product, and nuclease free water. I ran reactions through a single cycle of PCR on a  
891 thermocycler at 98°C for 2 min, 60°C for 30 s, and 72°C for 5 min. I pooled all of the  
892 PCR products for a plate into a single tube and purified the libraries with a 2X volume  
893 of SpeedBead solution as described before and resuspended in 25  $\mu$ L TLE. I added the  
894 remaining adapter and index sequences which were unique to each plate with four PCR  
895 reactions with a total volume of 50  $\mu$ L containing 1X Kapa Hifi (Kapa), 0.3  $\mu$ M iTru7  
896 Primer, 0.3  $\mu$ M P5 Primer, 0.3 mM dNTP, 1 unit of Kapa Hifi DNA Polymerase (Kapa),  
897 10  $\mu$ L purified iTru5-8N PCR product, and nuclease free water. I ran reactions on a  
898 thermocycler with an initial denaturation at 98°C for 2 min, followed by 6 cycles of 98°C  
899 for 20 s, 60°C for 15 s, 72°C for 30 s and a final extension of 72°C for 5 min. I pooled all  
900 of the PCR products for a plate into a single tube and purified the product with a 2X  
901 volume of SpeedBead solution as described before and resuspended in 45  $\mu$ L TLE.  
902

I size selected the library DNA from each plate in the range of 450-650 base pairs using  
903 a BluePippin (Sage Science, USA) with a 1.5% dye free gel with internal R2 standards.  
904 To increase the final DNA concentrations I prepared four PCR reactions for each plate  
905 with 1X Kapa Hifi (Kapa), 0.3  $\mu$ M P5 Primer, 0.3  $\mu$ M P7 Primer, 0.3 mM dNTP, 1 unit  
906 of Kapa HiFi DNA Polymerase (Kapa), 10  $\mu$ L size selected DNA, and nuclease free water  
907 and used the same thermocycling conditions as the previous (P5-iTru7) amplification. I  
908 pooled all of the PCR products for a plate into a single tube and purified the product with  
909 a 2X volume of SpeedBead solution as before and resuspended in 20  $\mu$ L TLE. I quantified  
910

the DNA concentration for each plate with a Qubit fluorometer (Life Technologies, USA) 911  
then pooled each plate in equimolar amounts relative to the number of samples on the 912  
plate and diluted the pooled DNA to 5 nM with TLE solution. The pooled libraries were 913  
pooled with other projects and sequenced on an Illumina HiSeqX by Novogene (China) 914  
to obtain paired end, 150 base pair sequences. 915

### 3.2.3 Data Processing 916

I demultiplexed the iTru7 indexes using the *process\_radtags* command from *Stacks* 917  
v2.6.4 (Rochette et al., 2019) and allowed for two mismatches for rescuing reads. To 918  
remove PCR duplicates, I used the *clone\_filter* command from *Stacks*. I demultiplexed 919  
inline sample barcodes, trimmed adapter sequence, and filtered reads with low quality 920  
scores as well as reads with any uncalled bases using the *process\_radtags* command again 921  
and allowed for the rescue of restriction site sequence as well as barcodes with up to two 922  
mismatches. I built alignments from the processed reads using the *Stacks* pipeline. I 923  
allowed for 14 mismatches between alleles within, as well as between individuals (M and 924  
n parameters). This is equivalent to a sequence similarity threshold of 90% for the 140 bp 925  
length of reads post trimming. I also allowed for up to 7 gaps between alleles within and 926  
between individuals. I used the *populations* command from *Stacks* to filter loci missing 927  
in more than 5% of individuals, filter all sites with minor allele counts less than 3, filter 928  
any individuals with more than 90% missing loci, and randomly sample a single SNP 929  
from each locus. 930

### 3.2.4 Genetic Clustering & Ancestry Proportions 931

To cluster individuals and characterize patterns of genetic differentiation and ad- 932  
mixture between clusters, I used the Bayesian inference program *STRUCTURE* v2.3.4 933  
(Pritchard et al., 2000) with *STRUCTURE*'s admixture model which returns an estimate 934  
of ancestry proportions for each sample. To evaluate the assumption that samples are 935  
best modeled as inheriting their genetic variation just two groups corresponding to the 936  
species identification made in the field, I ran *STRUCTURE* under four different models, 937

each with a different number of assumed clusters of individuals (K parameter) ranging 938  
from 1 to 4. For each value of K, I ran 20 iterations for 100,000 total steps with the first 939  
50,000 as burnin. I used the R package *POPHelper* v2.3.1 (Francis, 2017) to combine 940  
iterations for each value of K and to select the model producing the largest  $\Delta K$  which 941  
is the the model that has the greatest increase in likelihood score from the model with 942  
one fewer populations as described by (Evanno et al., 2005). I also examined genetic 943  
clustering and evidence of admixture using a non-parametric approach with a principal 944  
component analysis (PCA) implemented in the R package *adegenet* v2.1.10 (Jombart, 945  
2008). I visualized the relationship between the first principal component axis and the 946  
estimated admixture proportion for each individual to check for agreement between the 947  
parametric *STRUCTURE* analysis and the non-parametric PCA analysis. 948

### 3.2.5 Genomic Cline Analysis 949

To investigate patterns of introgression across the hybrid zone I used the bayesian 950  
genomic cline inference tool *BGC* v1.03 (Gompert & Buerkle, 2012) to infer parameters 951  
under a genomic cline model. Explain a bit how *BGC* works??? ... I classified a sample as 952  
being admixed if it had an inferred admixture proportion of <95% for one species under 953  
the model with a K of two in the *STRUCTURE* analysis. I used *VCFtools* vXX.XX.XX to 954  
filter all non-biallelic sites from the the VCF file produced by the *populations* command in 955  
*Stacks*. I converted the VCF formatted data into the *BGC* format using *bgc\_utils* v0.1.0, 956  
a *Python* package that I developed for this project [github.com/kerrycobb/bgc\\_utils](https://github.com/kerrycobb/bgc_utils). I 957  
ran *BGC* with 5 independent chains, each for 1,000,000 steps and sampling every 1000. I 958  
visualized MCMC output, discarded samples from the posterior as burnin, combined the 959  
independent chains, summarized the posterior samples, and identified outlier loci with 960  
*buc\_utils*. A primary goal of *BGC* analysis is to identify loci which have exceptional 961  
patterns of introgression. These loci, or loci in close linkage to them, are expected to 962  
be enriched for genetic regions affected by selection due to reproductive incompatibility 963  
between the two species. I identified loci with exceptional patterns of introgression using 964  
two approaches described by Gompert and Buerkle, 2011. (1) If locus specific introgres- 965

sion differed from the genome-wide average which I will refer to as "excess ancestry" 966  
following Gompert and Buerkle, 2011. More specifically, I classified a locus as having 967  
excess ancestry if the 90% highest posterior density interval (HPDI) for the alpha or beta 968  
parameter did not cover zero. (2) If locus specific introgression is statistically unlikely 969  
relative to the genome-wide distribution of locus specific introgression which I will refer 970  
to as "outliers" following Gompert and Buerkle, 2011. I classified a locus as an outlier if 971  
the median of the posterior sample for the  $\alpha$  or  $\beta$  parameters for a locus were not con- 972  
tained the interval from 0.05 to 0.95 of the probability density functions  $Normal(0, \tau_\alpha)$  973  
or  $Normal(0, \tau_\beta)$  respectively, where  $\tau_\alpha$  and  $\tau_\beta$  are the median values from the posterior 974  
sample for the conditional random effect priors on  $\tau_\alpha$  and  $\tau_\beta$ . These conditional priors 975  
describe the genome-wide variation of locus specific  $\alpha$  and  $\beta$ . I further classified outlier 976  
 $\alpha$  parameter estimates for a locus based on whether the median of the posterior sample 977  
was positive or negative. Positive estimates of  $\alpha$  mean there is a greater probability of 978  
*A. americanus* ancestry in individuals at the locus relative to their hybrid index whereas 979  
negative estimates of  $\alpha$  mean there is a greater probability of *A. terrestris* ancestry. 980

### 3.2.6 Genetic differentiation and Introgression 981

To test for a relationship between patterns of introgression and genetic divergence I 982  
used *VCFtools* to calculate the Weir and Cockerham, 1984  $F_{ST}$  between each species using 983  
only the samples inferred through the *STRUCTURE* analysis to have >95% ancestry for 984  
one species under the model with a K of two (Danecek et al., 2011). The Weir and 985  
Cockerham  $F_{ST}$  is calculated per site and I calculated the per site  $F_{ST}$  for the same 986  
sites as those used in the *BGC* analysis. To determine if patterns of introgression are 987  
correlated with population differentiation a performed a Pearson Correlation to test if 988  
 $F_{ST}$  correlates with either the  $\alpha$  or  $\beta$  parameters. I ran the correlation test with the 989  
absolute value of the median of the posterior sample for the  $\alpha$  parameter and the median 990  
of the posterior sample for the  $\beta$  parameter. I binned loci based on their status as  $\alpha$  991  
parameter outliers. I categorized loci as being  $\alpha$  outliers greater expected *A. americanus* 992  
ancestry,  $\alpha$  outliers with greater expected *A. terrestris* ancestry, and estimates of  $\alpha$  993

that are not outliers. I performed a Kruskal-Wallis test using *SciPy* v1.10.1 to test  
994 whether there were significant differences in values of  $F_{ST}$  at each locus between these  
995 groups (Virtanen et al., 2020). I then performed Mann-Whitney tests between all pairs  
996 of groups using *scikit-posthocs* to test which groups differ significantly from each other  
997 [github.com/maximtrp/scikit-posthocs](https://github.com/maximtrp/scikit-posthocs).  
998

### 3.3 Results

999

#### 3.3.1 Sampling and Data Processing

1000

I prepared reduced-representation sequencing libraries from 173 samples collected for  
1001 this study (Table 3.1) and 19 samples available from existing collections (Table 3.2)).  
1002 The *Stacks* pipeline assembled reads into 432,336 loci with a mean length of 253.31 bp.  
1003 Prior to filtering the mean coverage per sample was 32X. After filtering loci missing  
1004 from greater than 5% of samples, filtering sites with minor allele counts less than 3,  
1005 filtering individuals with greater than 90% missing loci, and randomly sampling a single  
1006 SNP from each locus, 1194 sites remained and 43 samples were excluded from further  
1007 analyses leaving a total of 149. For the included samples, 56 had been identified as most  
1008 closely resembling *A. americanus* and 93 had been identified as most closely resembling  
1009 *A. terrestris*.  
1010

#### 3.3.2 Genetic Clustering & Ancestry Proportions

1011

A visual inspection of the *STRUCTURE* results shows that each iteration with same  
1012 value for K converged on very similar results (Fig. 3.2). The *STRUCTURE* model with  
1013 the largest  $\Delta K$  was the model with a K of two (Fig. 3.1). Furthermore, individuals  
1014 are inferred as having ancestry derived largely from only two ancestral groups even for  
1015 K values of three and four. For these values of K, only a small amount of ancestry is  
1016 attributed to the third or fourth ancestral groups for any individual sample (Fig. 3.3)  
1017 Using a 95% estimated ancestry proportion as a cutoff for considering individuals to have  
1018 pure ancestry, 36 samples were classified as pure *A. americanus*, 75 as pure *A. terrestris*,  
1019

and 38 as being admixed. The proportions of admixture among the samples shows a clear  
1020 gradient between 0 and 1 which is consistent with many individuals being the product  
1021 of advanced generation hybrids beyond the  $F_1$  generation. The transition of admixture  
1022 proportions from one species to the other increase with distance from the locations of  
1023 pure individuals with proportions closest to 0.5 being found in the center of this transition  
1024 (Fig. 3.4).  
1025

### 3.3.3 Patterns of Introgression

  
1026

Visualization of the MCMC output with trace plots and histograms of each parameter  
1027 indicated that each of the five chains run in *BGC* converged on the same parameter space  
1028 and that each chain quickly reached stationarity. I conservatively discarded the first 10%  
1029 of samples as burnin. The median of the posterior sample for  $\alpha$  ranged from -0.525-  
1030 0.494. The  $\beta$  parameter was less variable and ranged from -0.158-0.220. I identified 16  
1031 loci with excess ancestry for the  $\alpha$  parameter relative to the genome wide average; i.e.,  
1032 the 90% HDPI does not cover 0. Of these, the median of the posterior sample for 5 of  
1033 these loci was negative and for 11 loci was positive. Negative values represent a greater  
1034 probability of *A. americanus* ancestry at a locus relative to the hybrid index whereas  
1035 positive values represent a greater probability of *A. terrestris* ancestry. I did not identify  
1036 any loci for which the estimates of  $\beta$  were outliers relative to the genome-wide average. I  
1037 identified 116 loci as outliers for the  $\alpha$  parameter relative to the genome-wide distribution  
1038 of locus specific introgression. Of these, the median of the posterior sample for 24 of these  
1039 loci was negative and for 92 loci was positive (Fig. 3.5). I did not identify any loci for  
1040 which the estimates of  $\beta$  were outliers relative to the genome-wide distribution of locus  
1041 specific introgression. All 16 of the loci identified as having excess ancestry for the  $\alpha$   
1042 parameter relative to the genome-wide average were also identified as outliers relative to  
1043 the genome-wide distribution of locus specific introgression.  
1044

### 3.3.4 Genomic Differentiation

1045

Genetic differentiation between *A. americanus* and *A. terrestris* was highly variable 1046  
(Fig. 3.6). Locus-specific  $F_{ST}$  between non-admixed *A. americanus* and *A. terrestris* had 1047  
a mean of 0.07.  $F_{ST}$  values for 249 loci were 0. Only a single locus had fixed differences 1048  
between species with an  $F_{ST}$  of 1.0. There is little apparent relationship between  $\alpha$  or 1049  
 $\beta$  and  $F_{ST}$  except at that the highest  $\alpha$  and  $\beta$  estimates have non-zero  $F_{ST}$  estimates 1050  
(Fig. 3.7). The Pearson correlation test estimates a weak correlation between  $\alpha$  and  $F_{ST}$  1051  
( $r=0.29$ ,  $p=1.62e-23$ ) a week correlation between  $\beta$  and  $F_{ST}$  ( $r=0.32$ ,  $p = 8.28 \times 10^{-30}$ ). 1052  
The result of the Kruskal-Wallis test are consistent with there being significant differences 1053  
between the  $F_{ST}$  values of loci with outlier  $\alpha$  estimates and non-outlier  $\alpha$  estimates on 1054  
average ( $p = 1.32 \times 10^{-40}$ ) (Fig. 3.6). The results of the post hoc pairwise Mann-Whitney 1055  
tests are consistent with both categories of loci with outlier  $\alpha$  estimates having greater 1056  
 $F_{ST}$  values on average than the non-outlier estimates of  $\alpha$ . The difference between non- 1057  
outlier loci and loci with greater probability of *A. americanus* ancestry was slightly higher 1058  
( $p = 2.72 \times 10^{-38}$ ) than the difference between non-outlier loci and loci with greater *A.* 1059  
*terrestris* ancestry ( $p = 8.16 \times 10^{-6}$ ). 1060

## 3.4 Discussion

1061

### 3.4.1 Evidence for ongoing hybridization

1062

With the genome-wide sequence data obtained in this study, I find evidence of sub- 1063  
stantial gene flow across the hybrid zone of these two species. The *STRUCTURE* analysis 1064  
inferred 38 out of 149 samples as having a proportion of ancestry of at least 5% of sites 1065  
attributable to admixture (Fig. 3.4). The admixture proportions inferred in the *STRU- 1066  
TURE* analysis range from 0.05%-0.5% which is consistent with hybrids being viable, 1067  
fertile, and capable of backcrossing over multiple generations [CITATION NEEDED!] 1068  
(Fig. 3.4). When backcrossing occurs over multiple generations in combination with mi- 1069  
gration of hybrid progeny and selection against introgressing alleles, a cline will form 1070  
across the hybrid zone with introgressing alleles becoming more uncommon with distance 1071

from the cline center (Barton & Hewitt, 1985). The results of the *STRUCTURE* analysis are largely consistent with this. Inferred admixture coefficients are highest at the center of the hybrid zone and decrease and approach zero with distance from the center (Fig. 3.4).

Admixed samples were located quite far from the center of the hybrid zone. In fact samples with greater than 5% admixture proportions are located all the way at the North-eastern and Southwestern edges of the sampling area. The width of a hybrid zone is a product of the strength of selection for or against introgression and the average dispersal distance of individuals within their reproductive lifespan (Barton & Hewitt, 1985). Breeden, 1987 estimated that 27% of individual *A. fowleri* breed at non-natal breeding ponds with some individuals dispersing at least as much as 2 km. Female *A. americanus* can migrate at least 1 km between breeding sites and post-breeding locations (Forester et al., 2006) Invasive cane toads (*Rhinella marina*) in Australia are estimated to have expanded their range at a rate of 10-15 km per year shortly after their introduction (Urban et al., 2008). The presence of samples with little to no admixture in close proximity to toads with high proportions of admixture shows that dispersal has an important roll in shaping the patterns of this hybrid zone. Individuals would be expected to appear more like their neighbors if dispersal rates and distances were very low. It is also likely that this hybrid zone may be more appropriately described as a mosaic hybrid zone rather than a more simple tension zone (Harrison, 1986). However, the sampling for this study or too sparse and irregular to definitively test this. Another possibility is that some of this inferred admixture is the result of a statistical artifact or due to error. Some reassurance is provided by the result of the PCA which is largely consistent with the *STRUCTURE* results although it is possible that they could be affected by the same bias or error introduced at by data collection and processing (Fig. 3.4) [CITATION NEEDED!].

The tension zone model of hybrid zones predicts that location of hybrid zones centers will be dependent on the effects of selection along with population density and natural dispersal barriers (Barton, 1979). The *STRUCTURE* results show that in two areas, there is a clear transition from samples with primarily *A. americanus* ancestry to samples with

primarily *A. terrestris* ancestry corresponding with the locations of streams and rivers. 1101  
In the Northern part of the sampling area, transitions occur at the Coosa River and at 1102  
Waxahatchee Creek (Fig. 3.4). In the Southern part, they occur at Sougahatchee Creek 1103  
(Fig. 3.4). Clearly these are not impassable boundaries as there has been introgression 1104  
beyond them. However, they likely reduce dispersal and as a result that the center of the 1105  
hybrid zone is caught in this location as described by Barton, 1979. 1106

### 3.4.2 Variability of introgression 1107

There are two primary parameters of interest in a genomic cline model that can be 1108  
interpreted in the evolutionary context of hybrid zones. The  $\alpha$  parameter specifies the 1109  
center of the cline and is dependent on the increase or decrease in the probability of 1110  
locus-specific ancestry from one of the parental populations. The  $\beta$  parameter specifies 1111  
the rate of change in probability of ancestry along the genome-wide admixture gradient. 1112  
Extreme estimates of these parameters may be associated with loci that cause reproduc- 1113  
tive incompatibility between hybridizing species. The Bayesian genomic cline analysis of 1114  
the genome-wide data in this study yielded extreme estimates for  $\alpha$  at some sites. Sites 1115  
were classified as having extreme values in two ways. First, sites could be classified as 1116  
having excess ancestry if the HDPI does not cover zero and is therefore extreme relative 1117  
to the genome-wide average of cline parameter estimates. Second, sites could be classi- 1118  
fied as being outliers if they are extreme relative to the genome-wide distribution of locus 1119  
specific effects under the cline model. A greater number of sites qualified as outliers for 1120  
estimates of  $\alpha$  than qualified as having excess ancestry. There were 116 loci classified as 1121  
outliers which make up 9.7% of the total number of sites. Of those, 16 were also classified 1122  
as having excess ancestry making up 1.3% of all sites. This difference is consistent with 1123  
other studies using both simulated and empirical data which typically find more outlier 1124  
loci than excess ancestry loci (Gompert & Buerkle, 2012). Both of these methods can 1125  
produce false positives as these extreme values can be produced solely by genetic drift 1126  
rather than by selection (Gompert & Buerkle, 2012). So not all sites with extreme 1127  
estimates will be associated with incompatibility loci. The false positive rate is exacer- 1128

bated when there are many loci with small effects on compatibility. However, these sites  
1129 should be enriched for loci associated with modest to strong reproductive incompatibility  
1130 and thus provide an upper estimate of the number of sites that are associated with these  
1131 modest to strong barriers to gene flow (Gompert & Buerkle, 2012).  
1132

None of the estimates for  $\beta$  were classified as either outliers or as having excess  
1133 ancestry. Simulations have demonstrated that the  $\alpha$  parameter is more impacted by  
1134 selection against hybrid genotypes than the  $\beta$  parameter (Gompert, Lucas, et al., 2012).  
1135 Other studies have also found no extreme estimates of  $\beta$  (Gompert, Lucas, et al., 2012;  
1136 Nikolakis et al., 2022)[CITATION NEEDED!]. One possible interpretation of the  
1137 absence of extreme values of  $\beta$  is that selection is only strong enough to have a significant  
1138 impact on  $\alpha$  but it is not strong enough to have a large impact on  $\beta$ . Unlike for  $\alpha$ ,  
1139 there is not a strong relationship between locally positive selection favoring introgressed  
1140 geneotypes and  $\beta$  (Gompert, Lucas, et al., 2012). Therefore, some of the extreme values  
1141 for  $\alpha$  could be due to adaptive introgression which does have much impact on estimates  
1142 of  $\beta$ . This is plausible given the large extent of introgression which is potentially due  
1143 to adaptive introgression. There is a negative relationship between  $\beta$  and dispersal rate  
1144 (Gompert, Lucas, et al., 2012). It is also plausible that high dispersal rates, rather than  
1145 selection is the cause of lower  $\beta$  values that do not reach the threshold to qualify as  
1146 extreme.  
1147

Of the 9.7% of sites that qualified as  $\alpha$  outliers, a substantially larger proportion had  
1148 positive values which represent greater *A. americanus* ancestry than expected at those  
1149 sites in admixed individuals. Negative  $\alpha$  estimates represent a greater probability of *A.*  
1150 *terrestris* ancestry at a site in within admixed individuals. Sites with positive outlier esti-  
1151 mates for  $\alpha$  made up 7.7% of all sites whereas those with negative outlier estimates made  
1152 up just 2%. This asymmetry suggests that introgression flows more in the direction of *A.*  
1153 *americanus* than it does in the direction of *A. terrestris*. This result is consistent with  
1154 a pattern evident upon visual inspection of the mapped *STRUCTURE* results. Samples  
1155 collected from sites adjacent to sites with admixed samples appear to have a greater  
1156 proportion of *A. americanus* ancestry than *A. terrestris* ancestry (Fig. 3.4). Taken to-  
1157

gether, these observations suggest that introgression at this hybrid zone is asymmetric  
1158  
(Yang et al., 2020). Asymmetries in introgression can arise for multiple reasons. There  
1159  
could differences in mate choice which make females of one species more selective than  
1160  
females of the other (Baldassarre et al., 2014). There can also be species differences in  
1161  
dispersal tendencies [CITATION NEEDED!]. Reciprocal-cross differences in repro-  
1162  
ductive isolation, termed Darwin’s Corollary, are very common (Turelli & Moyle, 2007).  
1163  
If one of the sexes is more prone to dispersal, introgression will flow more freely in one  
1164  
direction that it would in the other. It is possible that this observation is just an artifact  
1165  
of sampling. Particularly if this is a highly mosaic hybrid zone. However, many more  
1166  
samples with primarily *A. terrestris* ancestry were collected than samples with primarily  
1167  
*A. americanus* ancestry.  
1168

### 3.4.3 Relationship between introgression and differentiation

1169

Patterns of genetic differentiation and genomic introgression between *A. americanus*  
1170  
and *A. terrestris* are consistent with the hypothesis that regions of the genome ex-  
1171  
periencing divergent selection also affect hybrid fitness. As predicted, there is a positive  
1172  
association between locus specific estimates of  $F_{ST}$  and both the absolute value of the  $\alpha$   
1173  
and the  $\beta$  parameter estimates. Although this correlation supports the hypothesis that  
1174  
introgression outliers are linked to loci under selection, the association is only a mod-  
1175  
est one. Despite this, it is notable all of the outlier  $\alpha$  estimates as well as the highest  
1176  
 $\beta$  estimates have non-zero  $F_{ST}$  estimates. Whereas sites with lower  $\alpha$  and  $\beta$  estimates  
1177  
span the entire range from zero to one. This is consistent with expectations of secondary  
1178  
contact where not all loci that have undergone genomic divergence will necessarily result  
1179  
in reproductive isolation. A tighter coupling of divergence and resistance to gene flow  
1180  
would be expected under a scenario of divergence with gene flow.  
1181

<b>3.4.4 Conclusion</b>	1182
<b>References</b>	1183
Baldassarre, D. T., White, T. A., Karubian, J., & Webster, M. S. (2014). GENOMIC AND MORPHOLOGICAL ANALYSIS OF A SEMIPERMEABLE AVIAN HYBRID ZONE SUGGESTS ASYMMETRICAL INTROGRESSION OF A SEXUAL SIGNAL. <i>Evolution</i> , 68(9), 2644–2657. <a href="https://doi.org/10.1111/evo.12457">https://doi.org/10.1111/evo.12457</a>	1184
Barton, N. H. (1979). The dynamics of hybrid zones. <i>Heredity</i> , 43(3), 341–359. <a href="https://doi.org/10.1038/hdy.1979.87">https://doi.org/10.1038/hdy.1979.87</a>	1188
Barton, N. H., & Hewitt, G. M. (1985). Analysis of Hybrid Zones. <i>Annual Review</i> , 16, 113–148.	1190
Bayona-Vásquez, N. J., Glenn, T. C., Kieran, T. J., Pierson, T. W., Hoffberg, S. L., Scott, P. A., Bentley, K. E., Finger, J. W., Louha, S., Troendle, N., Diaz-Jaimes, P., Mauricio, R., & Faircloth, B. C. (2019). Adapterama III: Quadruple-indexed, double/triple-enzyme RADseq libraries (2RAD/3RAD). <i>PeerJ</i> , 7, e7724. <a href="https://doi.org/10.7717/peerj.7724">https://doi.org/10.7717/peerj.7724</a>	1192
Blackman, B. (2016). Speciation Genes. <i>Encyclopedia of Evolutionary Biology</i> (pp. 166–175). Elsevier. <a href="https://doi.org/10.1016/B978-0-12-800049-6.00066-4">https://doi.org/10.1016/B978-0-12-800049-6.00066-4</a>	1197
Blair, W. F. (1963). Intragroup genetic compatibility in the <i>Bufo americanus</i> species group of toads. <i>The Texas Journal of Science</i> , 13, 15–34.	1199
Blair, W. F. (1972). <i>Evolution in the genus Bufo</i> . University of Texas Press.	1201
Breden, F. (1987). The Effect of Post-Metamorphic Dispersal on the Population Genetic Structure of Fowler's Toad, <i>Bufo woodhousei fowleri</i> . <i>Copeia</i> , 1987(2), 386–395. <a href="https://doi.org/10.2307/1445775">https://doi.org/10.2307/1445775</a>	1202
Butlin, R., Debelle, A., Kerth, C., Snook, R. R., Beukeboom, L. W., RF, C. C., Diao, W., Maan, M. E., Paolucci, S., Weissing, F. J., et al. (2011). What do we need to know about speciation? <i>Trends in ecology &amp; evolution</i> , 27(1), 27–39.	1205

Cocroft, R. B., & Ryan, M. J. (1995). Patterns of advertisement call evolution in toads and chorus frogs. <i>Animal Behaviour</i> , 49(2), 283–303. <a href="https://doi.org/10.1006/anbe.1995.0043">https://doi.org/10.1006/anbe.1995.0043</a>	1208 1209 1210
Colliard, C., Sicilia, A., Turrisi, G. F., Arculeo, M., Perrin, N., & Stöck, M. (2010). Strong reproductive barriers in a narrow hybrid zone of West-Mediterranean green toads ( <i>Bufo viridissubgroup</i> ) with Plio-Pleistocene divergence. <i>BMC Evolutionary Biology</i> , 10(1), 232. <a href="https://doi.org/10.1186/1471-2148-10-232">https://doi.org/10.1186/1471-2148-10-232</a>	1211 1212 1213 1214
Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin, R., & 1000 Genomes Project Analysis Group. (2011). The variant call format and VCFtools. <i>Bioinformatics</i> , 27(15), 2156–2158. <a href="https://doi.org/10.1093/bioinformatics/btr330">https://doi.org/10.1093/bioinformatics/btr330</a>	1215 1216 1217 1218 1219
Đorđević, S., & Simović, A. (2014). STRANGE AFFECTION: MALE BUFO BUFO (ANURA: BUFONIDAE) PASSIONATELY EMBRACING A BULGE OF MUD. <i>Ecologica Montenegrina</i> , 1(1), 15–17. <a href="https://doi.org/10.37828/em.2014.1.4">https://doi.org/10.37828/em.2014.1.4</a>	1220 1221 1222
Evanno, G., Regnaut, S., & Goudet, J. (2005). Detecting the number of clusters of individuals using the software structure: A simulation study. <i>Molecular Ecology</i> , 14(8), 2611–2620. <a href="https://doi.org/10.1111/j.1365-294X.2005.02553.x">https://doi.org/10.1111/j.1365-294X.2005.02553.x</a>	1223 1224 1225
Forester, D. C., Snodgrass, J. W., Marsalek, K., & Lanham, Z. (2006). Post-Breeding Dispersal and Summer Home Range of Female American Toads ( <i>Bufo americanus</i> ). <i>Northeastern Naturalist</i> , 13(1), 59–72. <a href="https://doi.org/10.1656/1092-6194(2006)13[59:PDASHR]2.0.CO;2">https://doi.org/10.1656/1092-6194(2006)13[59:PDASHR]2.0.CO;2</a>	1226 1227 1228 1229
Francis, R. M. (2017). POPHELPER: An R package and web app to analyse and visualize population structure. <i>Molecular Ecology Resources</i> , 17(1), 27–32. <a href="https://doi.org/10.1111/1755-0998.12509">https://doi.org/10.1111/1755-0998.12509</a>	1230 1231 1232
Gompert, Z., & Buerkle, C. A. (2012). Bgc: Software for Bayesian estimation of genomic clines. <i>Molecular Ecology Resources</i> , 12(6), 1168–1176. <a href="https://doi.org/10.1111/1755-0998.12009.x">https://doi.org/10.1111/1755-0998.12009.x</a>	1233 1234 1235

Gompert, Z., & Buerkle, C. A. (2011). Bayesian estimation of genomic clines: BAYESIAN GENOMIC CLINES. <i>Molecular Ecology</i> , 20(10), 2111–2127. <a href="https://doi.org/10.1111/j.1365-294X.2011.05074.x">https://doi.org/10.1111/j.1365-294X.2011.05074.x</a>	1236
Gompert, Z., Lucas, L. K., Nice, C. C., Fordyce, J. A., Forister, M. L., & Buerkle, C. A. (2012). GENOMIC REGIONS WITH A HISTORY OF DIVERGENT SELECTION AFFECT FITNESS OF HYBRIDS BETWEEN TWO BUTTERFLY SPECIES: GENOMICS OF SPECIATION. <i>Evolution</i> , 66(7), 2167–2181. <a href="https://doi.org/10.1111/j.1558-5646.2012.01587.x">https://doi.org/10.1111/j.1558-5646.2012.01587.x</a>	1237
Gompert, Z., Mandeville, E. G., & Buerkle, C. A. (2017). Analysis of Population Genomic Data from Hybrid Zones. <i>Annual Review of Ecology, Evolution, and Systematics</i> , 48(1), 207–229. <a href="https://doi.org/10.1146/annurev-ecolsys-110316-022652">https://doi.org/10.1146/annurev-ecolsys-110316-022652</a>	1238
Gompert, Z., Parchman, T. L., & Buerkle, C. A. (2012). Genomics of isolation in hybrids. <i>Philosophical Transactions of the Royal Society B: Biological Sciences</i> , 367(1587), 439–450. <a href="https://doi.org/10.1098/rstb.2011.0196">https://doi.org/10.1098/rstb.2011.0196</a>	1239
Green, D. M. (1983). Allozyme Variation through a Clinal Hybrid Zone between the Toads <i>Bufo americanus</i> and <i>B. hemiophrys</i> in Southeastern Manitoba. <i>Herpetologica</i> , 39(1), 28–40.	1240
Green, D. M. (1996). The bounds of species: Hybridization in the <i>Bufo americanus</i> group of North American toads. <i>Israel Journal of Zoology</i> , 42, 95–109.	1241
Harrison, R. G. (1986). Pattern and process in a narrow hybrid zone. <i>Heredity</i> , 56(3), 337–349. <a href="https://doi.org/10.1038/hdy.1986.55">https://doi.org/10.1038/hdy.1986.55</a>	1242
Hedrick, P. W. (2013). Adaptive introgression in animals: Examples and comparison to new mutation and standing variation as sources of adaptive variation. <i>Molecular Ecology</i> , 22(18), 4606–4618. <a href="https://doi.org/10.1111/mec.12415">https://doi.org/10.1111/mec.12415</a>	1243
Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. <i>Computing in Science &amp; Engineering</i> , 9(3), 90–95. <a href="https://doi.org/10.1109/MCSE.2007.55">https://doi.org/10.1109/MCSE.2007.55</a>	1244
Jahner, J. P., Parchman, T. L., & Matocq, M. D. (2021). Multigenerational backcrossing and introgression between two woodrat species at an abrupt ecological transition. <i>Molecular Ecology</i> , 30(17), 4245–4258. <a href="https://doi.org/10.1111/mec.16056">https://doi.org/10.1111/mec.16056</a>	1245

Jombart, T. (2008). Adegenet : A R package for the multivariate analysis of genetic markers.	<i>Bioinformatics</i> , 24(11), 1403–1405. <a href="https://doi.org/10.1093/bioinformatics/btn129">https://doi.org/10.1093/bioinformatics/btn129</a>	1265 1266 1267
Larson, E. L., Andrés, J. A., Bogdanowicz, S. M., & Harrison, R. G. (2013). DIFFERENTIAL INTROGRESSION IN A MOSAIC HYBRID ZONE REVEALS CANDIDATE BARRIER GENES.	<i>Evolution</i> , 67(12), 3653–3661. <a href="https://doi.org/10.1111/evo.12205">https://doi.org/10.1111/evo.12205</a>	1268 1269 1270 1271
Mallet, J. (2005). Hybridization as an invasion of the genome.	<i>Trends in Ecology &amp; Evolution</i> , 20(5), 229–237. <a href="https://doi.org/10.1016/j.tree.2005.02.010">https://doi.org/10.1016/j.tree.2005.02.010</a>	1272 1273
Mallet, J. (2008). Hybridization, ecological races and the nature of species: Empirical evidence for the ease of speciation.	<i>Philosophical Transactions of the Royal Society B: Biological Sciences</i> , 363(1506), 2971–2986. <a href="https://doi.org/10.1098/rstb.2008.0081">https://doi.org/10.1098/rstb.2008.0081</a>	1274 1275 1276 1277
Malone, J. H., & Fontenot, B. E. (2008). Patterns of Reproductive Isolation in Toads (R. DeSalle, Ed.).	<i>PLoS ONE</i> , 3(12), e3900. <a href="https://doi.org/10.1371/journal.pone.0003900">https://doi.org/10.1371/journal.pone.0003900</a>	1278 1279 1280
Miller, C. J. J., & Matute, D. R. (2016). The Effect of Temperature on Drosophila Hybrid Fitness.	<i>G3: Genes/Genomes/Genetics</i> , 7(2), 377–385. <a href="https://doi.org/10.1534/g3.116.034926">https://doi.org/10.1534/g3.116.034926</a>	1281 1282 1283
Moran, B. M., Payne, C., Langdon, Q., Powell, D. L., Brandvain, Y., & Schumer, M. (2021). The genomic consequences of hybridization (P. J. Wittkopp, Ed.).	<i>eLife</i> , 10, e69016. <a href="https://doi.org/10.7554/eLife.69016">https://doi.org/10.7554/eLife.69016</a>	1284 1285 1286
Mount, R. H. (1975). <i>The Reptiles and Amphibians of Alabama</i> .	The University of Alabama Press.	1287 1288
Nikolakis, Z. L., Schield, D. R., Westfall, A. K., Perry, B. W., Ivey, K. N., Orton, R. W., Hales, N. R., Adams, R. H., Meik, J. M., Parker, J. M., Smith, C. F., Gompert, Z., Mackessy, S. P., & Castoe, T. A. (2022). Evidence that genomic incompatibilities and other multilocus processes impact hybrid fitness in a rattlesnake hybrid zone.	<i>Evolution</i> , 76(11), 2513–2530. <a href="https://doi.org/10.1111/evo.14612">https://doi.org/10.1111/evo.14612</a>	1289 1290 1291 1292 1293

Nosil, P., Funk, D. J., & Ortiz-Barrientos, D. (2009). Divergent selection and heterogeneous genomic divergence. <i>Molecular Ecology</i> , 18(3), 375–402. <a href="https://doi.org/10.1111/j.1365-294X.2008.03946.x">https://doi.org/10.1111/j.1365-294X.2008.03946.x</a>	1294 1295 1296
Nosil, P., & Schlüter, D. (2011). The genes underlying the process of speciation. <i>Trends in Ecology &amp; Evolution</i> , 26(4), 160–167. <a href="https://doi.org/10.1016/j.tree.2011.01.001">https://doi.org/10.1016/j.tree.2011.01.001</a>	1297 1298
Parchman, T. L., Gompert, Z., Braun, M. J., Brumfield, R. T., McDonald, D. B., Uy, J. a. C., Zhang, G., Jarvis, E. D., Schlinger, B. A., & Buerkle, C. A. (2013). The genomic consequences of adaptive divergence and reproductive isolation between species of manakins. <i>Molecular Ecology</i> , 22(12), 3304–3317. <a href="https://doi.org/10.1111/mec.12201">https://doi.org/10.1111/mec.12201</a>	1299 1300 1301 1302 1303
Powell, R., Conant, R., & Collins, J. T. (2016). <i>A Field Guide to Reptiles &amp; Amphibians: Eastern and Central North America</i> (4th ed.). Houghton Mifflin Harcourt.	1304 1305
Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of Population Structure Using Multilocus Genotype Data. <i>Genetics</i> , 155(2), 945–959. <a href="https://doi.org/10.1093/genetics/155.2.945">https://doi.org/10.1093/genetics/155.2.945</a>	1306 1307 1308
Rieseberg, L. H., Whitton, J., & Gardner, K. (1999). Hybrid zones and the genetic architecture of a barrier to gene flow between two sunflower species. <i>Genetics</i> , 152(2), 713–727.	1309 1310 1311
Rochette, N. C., Rivera-Colón, A. G., & Catchen, J. M. (2019). Stacks 2: Analytical methods for paired-end sequencing improve RADseq-based population genomics. <i>Molecular Ecology</i> , 28(21), 4737–4754. <a href="https://doi.org/10.1111/mec.15253">https://doi.org/10.1111/mec.15253</a>	1312 1313 1314
Servedio, M. R., & Noor, M. A. (2003). The Role of Reinforcement in Speciation: Theory and Data. <i>Annual Review of Ecology, Evolution, and Systematics</i> , 34(1), 339–364. <a href="https://doi.org/10.1146/annurev.ecolsys.34.011802.132412">https://doi.org/10.1146/annurev.ecolsys.34.011802.132412</a>	1315 1316 1317
Shankman, D., & Hart, J. L. (2007). The Fall Line: A Physiographic-Forest Vegetation Boundary. <i>Geographical Review</i> , 97(4), 502–519. <a href="https://doi.org/10.1111/j.1931-0846.2007.tb00409.x">https://doi.org/10.1111/j.1931-0846.2007.tb00409.x</a>	1318 1319 1320
Taylor, E. B., Boughman, J. W., Groenenboom, M., Sniatynski, M., Schlüter, D., & Gow, J. L. (2006). Speciation in reverse: Morphological and genetic evidence of	1321 1322

the collapse of a three-spined stickleback ( <i>Gasterosteus aculeatus</i> ) species pair.	1323
<i>Molecular Ecology</i> , 15(2), 343–355. <a href="https://doi.org/10.1111/j.1365-294X.2005.02794.x">https://doi.org/10.1111/j.1365-294X.2005.02794.x</a>	1324
Turelli, M., & Moyle, L. C. (2007). Asymmetric Postmating Isolation: Darwin's Corollary to Haldane's Rule. <i>Genetics</i> , 176(2), 1059–1088. <a href="https://doi.org/10.1534/genetics.106.065979">https://doi.org/10.1534/genetics.106.065979</a>	1325
Urban, M. C., Phillips, B. L., Skelly, D. K., & Shine, R. (2008). A Toad More Traveled: The Heterogeneous Invasion Dynamics of Cane Toads in Australia. <i>The American Naturalist</i> , 171(3), E134–E148. <a href="https://doi.org/10.1086/527494">https://doi.org/10.1086/527494</a>	1326
Van Riemsdijk, I., Arntzen, J. W., Bucciarelli, G. M., McCartney-Melstad, E., Rafajlović, M., Scott, P. A., Toffelmier, E., Shaffer, H. B., & Wielstra, B. (2023). Two transects reveal remarkable variation in gene flow on opposite ends of a European toad hybrid zone. <i>Heredity</i> . <a href="https://doi.org/10.1038/s41437-023-00617-6">https://doi.org/10.1038/s41437-023-00617-6</a>	1327
Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., ... SciPy 1.0 Contributors. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in python. <i>Nature Methods</i> , 17, 261–272. <a href="https://doi.org/10.1038/s41592-019-0686-2">https://doi.org/10.1038/s41592-019-0686-2</a>	1328
Weatherby, C. A. (1982). INTROGRESSION BETWEEN THE AMERICAN TOAD, <i>BUFO AMERICANUS</i> , AND THE SOUTHERN TOAD, <i>B. TERRESTRIS</i> , IN ALABAMA.	1329
Weir, B. S., & Cockerham, C. C. (1984). Estimating F-Statistics for the Analysis of Population Structure. <i>Evolution</i> , 38(6), 1358–1370. <a href="https://doi.org/10.2307/2408641">https://doi.org/10.2307/2408641</a>	1330
Wolf, J. B. W., & Ellegren, H. (2017). Making sense of genomic islands of differentiation in light of speciation. <i>Nature Reviews Genetics</i> , 18(2), 87–100. <a href="https://doi.org/10.1038/nrg.2016.133">https://doi.org/10.1038/nrg.2016.133</a>	1331

- Wu, C.-I. (2001). The genic view of the process of speciation. *Journal of Evolutionary Biology*, 14 (6), 851–865. <https://doi.org/10.1046/j.1420-9101.2001.00335.x> 1351
- Yang, W., Feiner, N., Laakkonen, H., Sacchi, R., Zuffi, M. A. L., Scali, S., While, G. M., & Uller, T. (2020). Spatial variation in gene flow across a hybrid zone reveals causes of reproductive isolation and asymmetric introgression in wall lizards\*. *Evolution*, 74 (7), 1289–1300. <https://doi.org/10.1111/evo.14001> 1352
- Yang, W., Feiner, N., Laakkonen, H., Sacchi, R., Zuffi, M. A. L., Scali, S., While, G. M., & Uller, T. (2020). Spatial variation in gene flow across a hybrid zone reveals causes of reproductive isolation and asymmetric introgression in wall lizards\*. *Evolution*, 74 (7), 1289–1300. <https://doi.org/10.1111/evo.14001> 1353
- Uller, T. (2020). Spatial variation in gene flow across a hybrid zone reveals causes of reproductive isolation and asymmetric introgression in wall lizards\*. *Evolution*, 74 (7), 1289–1300. <https://doi.org/10.1111/evo.14001> 1354
- Uller, T. (2020). Spatial variation in gene flow across a hybrid zone reveals causes of reproductive isolation and asymmetric introgression in wall lizards\*. *Evolution*, 74 (7), 1289–1300. <https://doi.org/10.1111/evo.14001> 1355
- Uller, T. (2020). Spatial variation in gene flow across a hybrid zone reveals causes of reproductive isolation and asymmetric introgression in wall lizards\*. *Evolution*, 74 (7), 1289–1300. <https://doi.org/10.1111/evo.14001> 1356

### 3.5 Figures

1357

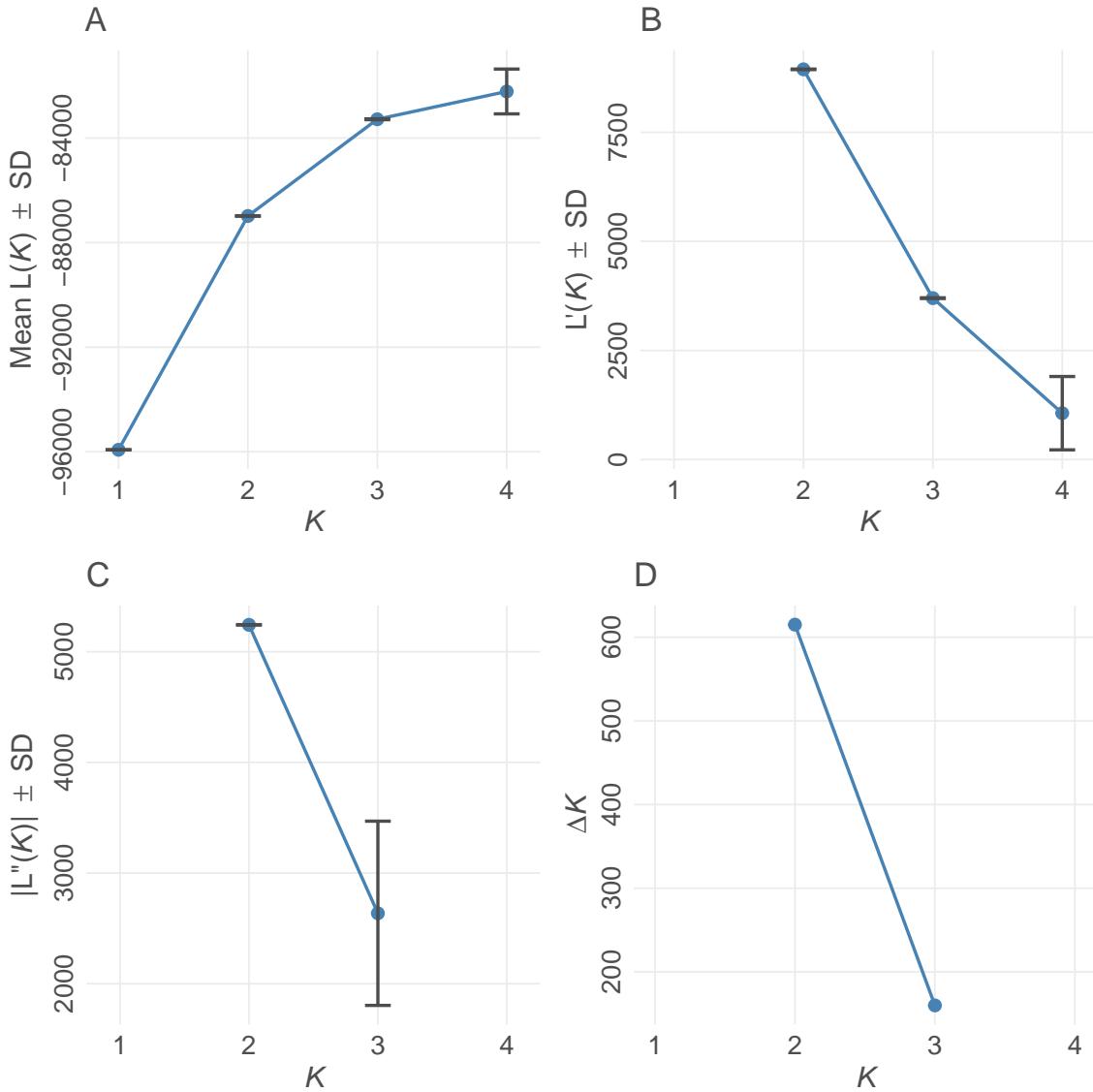


Figure 3.1. Evanno method for optimal value for  $K$  in *STRUCTURE* (Evanno et al., 2005).  $K$  refers to the number of populations for each of the different *STRUCTURE* models examined. (A) Mean estimated  $\ln$  probability of data over 10 iterations for each value of  $K \pm SD$ . (B) Rate of change of the likelihood distribution (mean  $\pm SD$ ) (C) Absolute values of the second order rate of change of the likelihood distribution (mean  $\pm SD$ ) (D)  $\Delta K$ . The modal value of this distribution is considered the true value of  $K$  for the data. Plot created using *POPHELPER* (Francis, 2017).

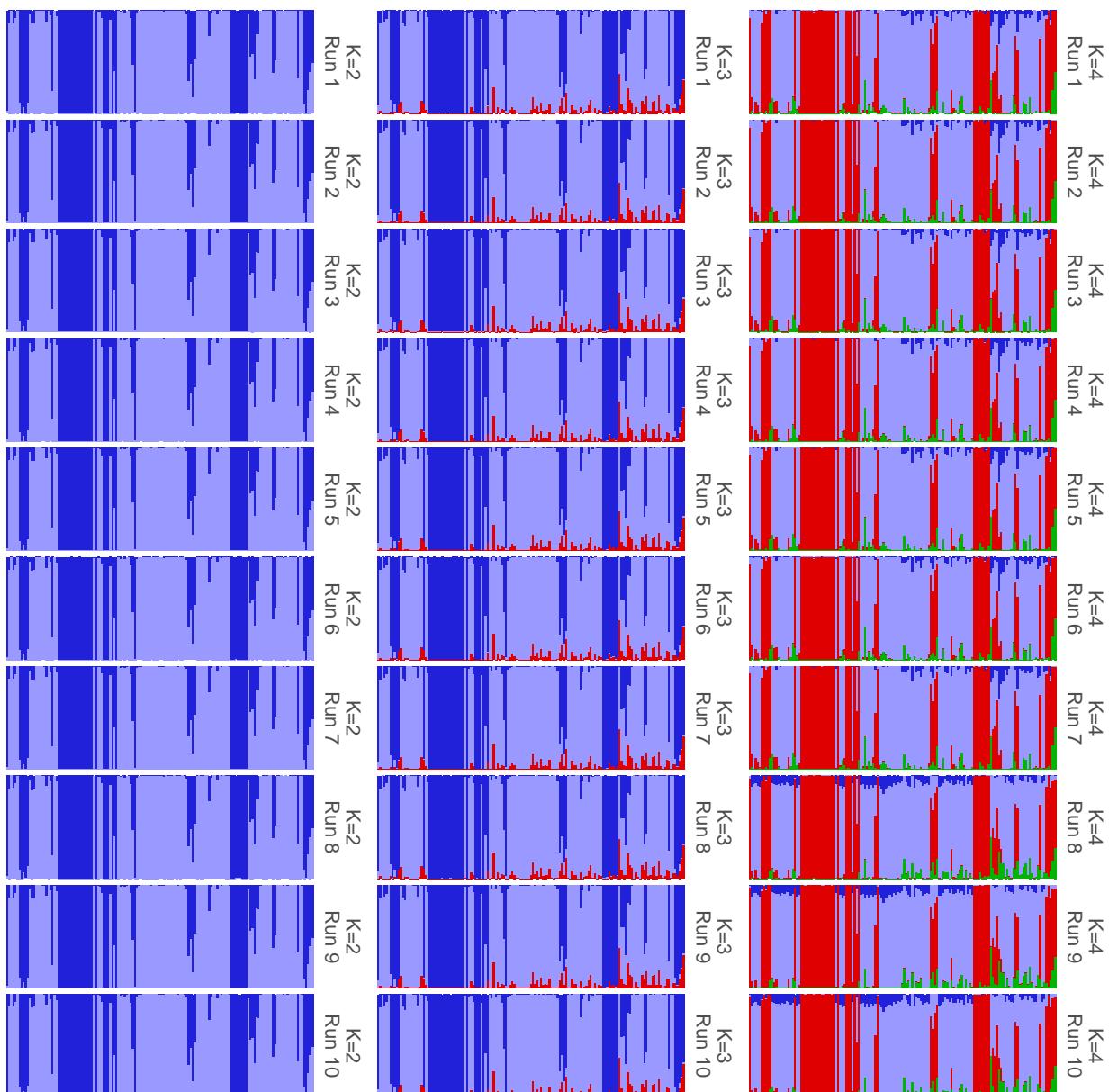


Figure 3.2. Results of each iteration of *STRUCTURE* showing convergence among iterations within runs having the same value for  $K$ . Plot was created with *POPHELPER* (Francis, 2017).

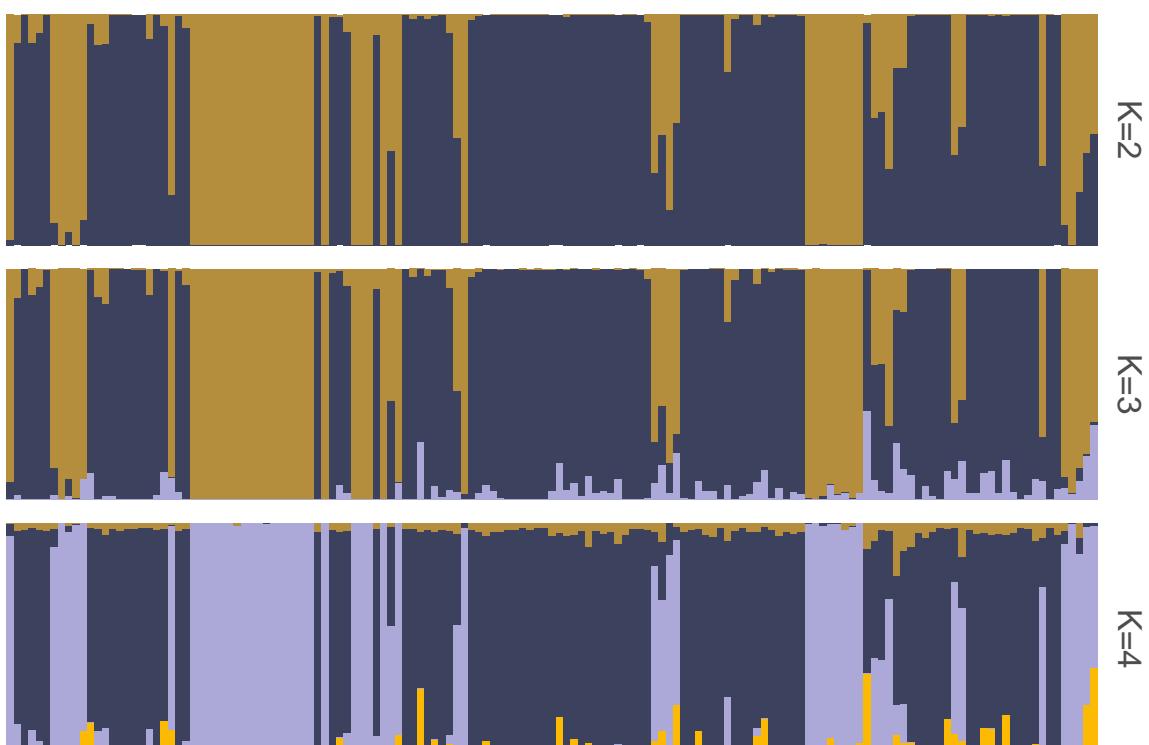


Figure 3.3. Summarized *STRUCTURE* results for each value of  $K$ . Ancestry proportions shown are the mean of ancestry proportions across all iterations. Summarization and plotting done using *POPHELPER* (Francis, 2017).

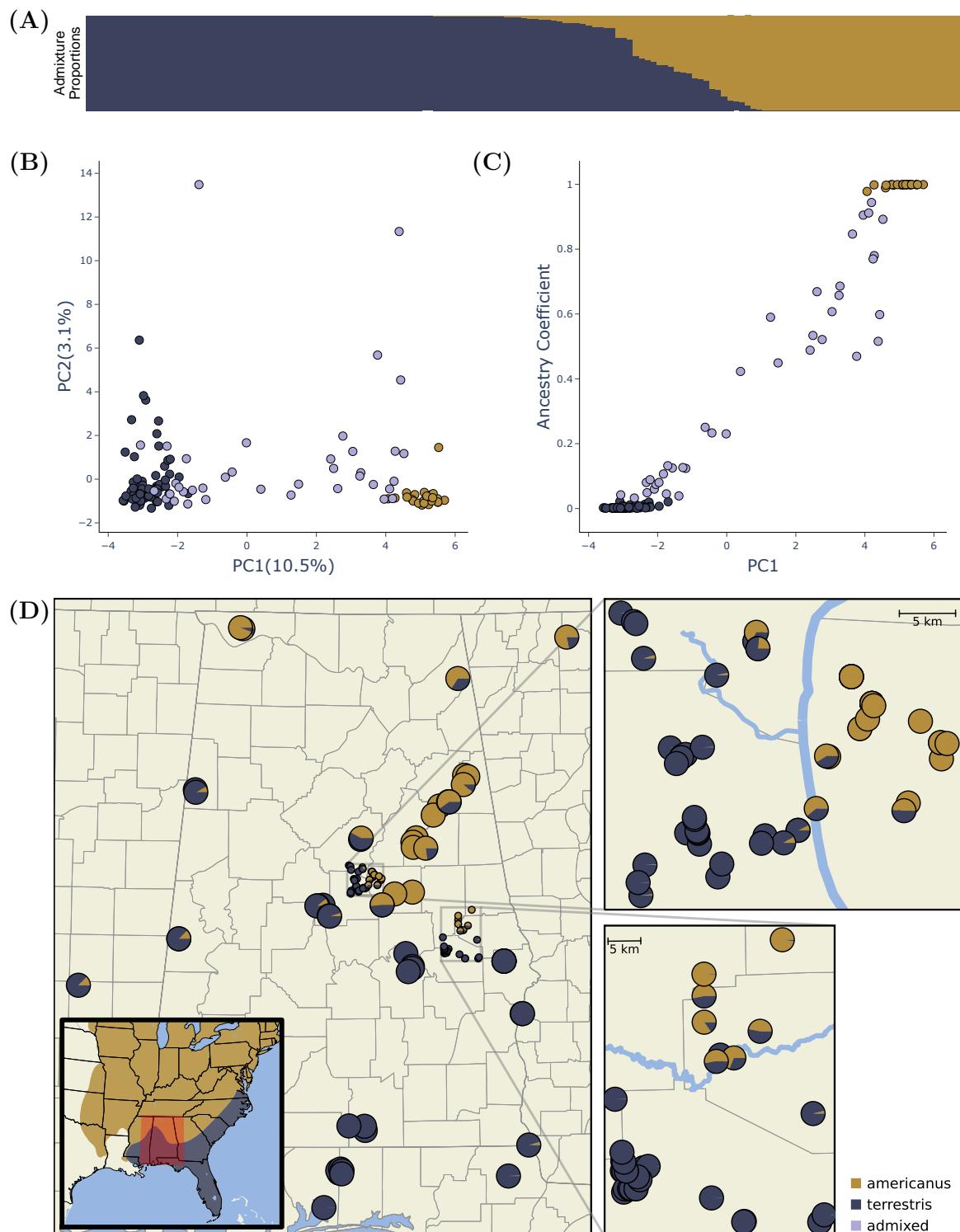


Figure 3.4. Genetic evidence of hybridization between *A. americanus* and *A. terrestris*. (A) *STRUCTURE* plot showing estimated ancestry proportions. (B) Summary of population genetic structure based on the principle component axes one (PC1) and two (PC2). These axes explain 10.5% (PC1) and 3.1% (PC2) of the genetic variation among individuals. (C) Relationship between the first principal component axis and the admixture proportions estimated with *STRUCTURE*. (D) Sample map showing the sampling location and estimated ancestry proportion of each sample. The inset map shows the approximate ranges of each species and the study area highlighted in red. Figure created using *POPHelper* (Francis, 2017) and *Matplotlib* (Hunter, 2007)

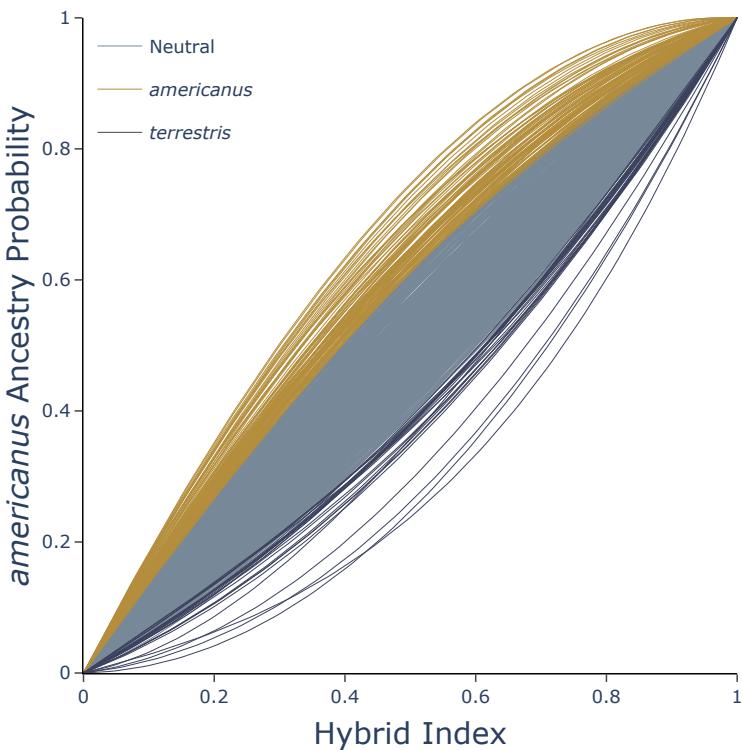


Figure 3.5. Shape of genomic clines estimated for each locus with BGC. Outliers are highlighted with XX.

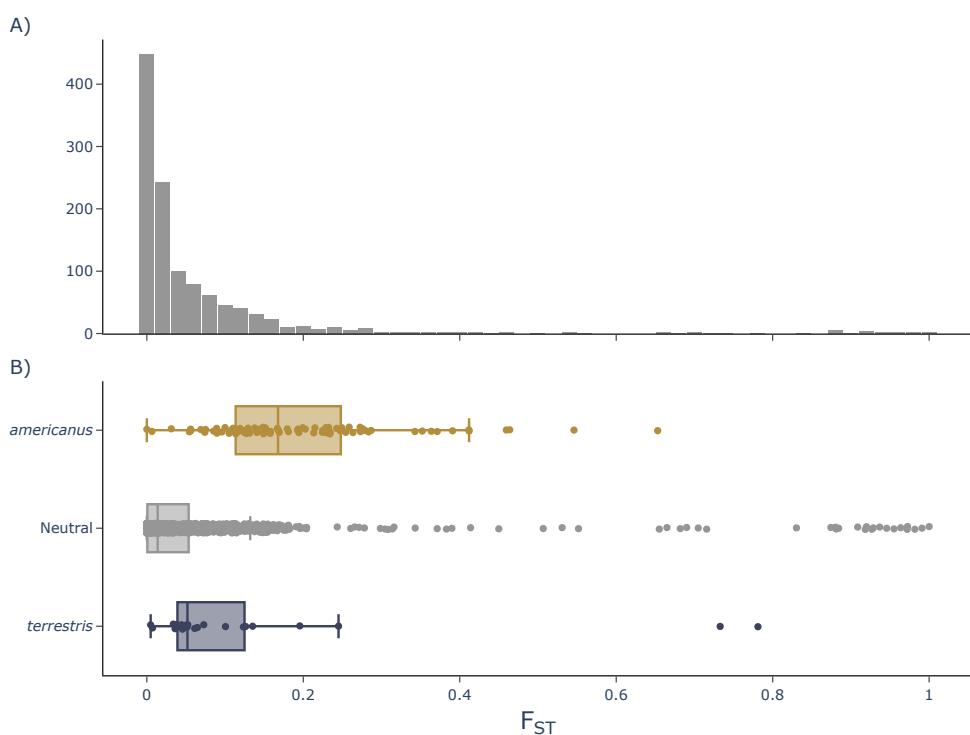


Figure 3.6. A) Distribution of per site  $F_{ST}$  estimates. B) Box plots showing the distribution and mean of  $F_{ST}$  for three categories of  $\alpha$  estimates, outliers with greater than expected *A. americanus* ancestry (gold), outliers with greater than expected *A. terrestris* ancestry (violet), and non-outliers (gray).

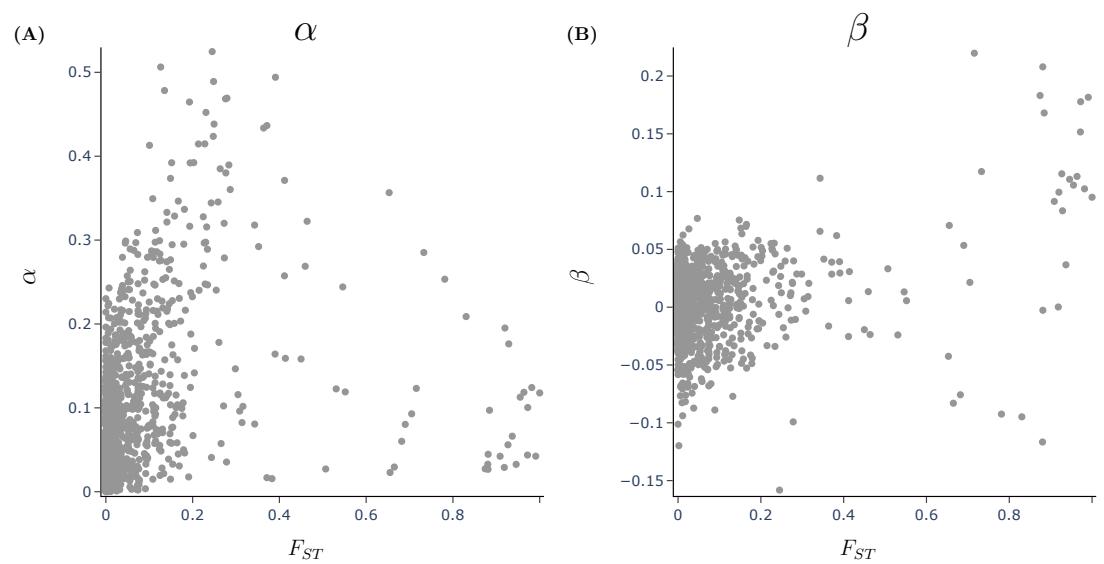


Figure 3.7. Relationship between genetic divergence measured with Weir and Cockerham, 1984  $F_{ST}$  and BGC cline parameters A)  $\alpha$  and B)  $\beta$ .

## 3.6 Tables

1358

Table 3.1. Samples collected for this study

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 016	<i>terrestris</i>	30.54819	-86.93067	X
KAC 038	<i>terrestris</i>	32.81470	-86.93968	X
KAC 039	<i>terrestris</i>	32.81094	-86.98967	X
KAC 040	<i>terrestris</i>	32.80985	-86.99795	X
KAC 042	<i>terrestris</i>	32.82406	-86.99314	
KAC 043	<i>terrestris</i>	32.82406	-86.99314	
KAC 044	<i>terrestris</i>	32.80450	-87.03078	
KAC 045	<i>terrestris</i>	32.76703	-87.07073	
KAC 046	<i>terrestris</i>	32.76592	-87.07184	
KAC 047	<i>terrestris</i>	32.78932	-86.90850	
KAC 048	<i>terrestris</i>	32.73575	-86.88149	X
KAC 049	<i>terrestris</i>	32.73291	-86.87707	X
KAC 050	<i>terrestris</i>	32.74822	-86.79806	
KAC 051	<i>terrestris</i>	32.78742	-86.75847	
KAC 052	<i>terrestris</i>	32.78044	-86.73877	
KAC 070	<i>americanus</i>	34.79963	-84.57678	X
KAC 071	<i>terrestris</i>	32.43478	-85.64630	
KAC 074	<i>terrestris</i>	30.77430	-85.22690	X
KAC 075	<i>terrestris</i>	32.94778	-86.63224	X
KAC 076	<i>terrestris</i>	32.94970	-86.52687	
KAC 077	<i>terrestris</i>	32.94970	-86.52687	
KAC 078	<i>americanus</i>	33.00267	-86.38960	X
KAC 079	<i>americanus</i>	33.01205	-86.47872	
KAC 080	<i>americanus</i>	33.04456	-86.45547	
KAC 081	<i>americanus</i>	33.04456	-86.45547	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 082	<i>americanus</i>	33.04456	-86.45547	X
KAC 083	<i>americanus</i>	33.04456	-86.45547	X
KAC 084	<i>americanus</i>	33.04456	-86.45547	X
KAC 085	<i>americanus</i>	33.04456	-86.45547	
KAC 086	<i>americanus</i>	33.04456	-86.45547	X
KAC 087	<i>americanus</i>	33.01484	-86.39040	X
KAC 089	<i>americanus</i>	33.01484	-86.39040	X
KAC 090	<i>americanus</i>	33.06472	-86.47496	X
KAC 091	<i>americanus</i>	33.06472	-86.47496	X
KAC 092	<i>americanus</i>	33.06472	-86.47496	
KAC 093	<i>americanus</i>	33.06472	-86.47496	X
KAC 094	<i>americanus</i>	33.06472	-86.47496	X
KAC 095	<i>americanus</i>	33.06472	-86.47496	X
KAC 096	<i>americanus</i>	33.06472	-86.47496	X
KAC 097	<i>americanus</i>	33.06472	-86.47496	X
KAC 098	<i>americanus</i>	33.02572	-86.46711	X
KAC 099	<i>americanus</i>	33.02572	-86.46711	X
KAC 100	<i>terrestris</i>	32.92374	-86.67199	X
KAC 101	<i>americanus</i>	33.03283	-86.45975	X
KAC 102	<i>terrestris</i>	32.94544	-86.55777	X
KAC 103	<i>terrestris</i>	32.94947	-86.52630	X
KAC 104	<i>terrestris</i>	32.94947	-86.52630	X
KAC 105	<i>americanus</i>	33.04278	-86.45377	X
KAC 106	<i>americanus</i>	33.00464	-86.49692	X
KAC 107	<i>americanus</i>	33.01416	-86.38417	X
KAC 108	<i>terrestris</i>	32.94013	-86.54004	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 109	<i>terrestris</i>	32.94173	-86.55787	
KAC 110	<i>americanus</i>	33.03099	-86.40941	X
KAC 111	<i>americanus</i>	33.00518	-86.49895	X
KAC 112	<i>terrestris</i>	32.95011	-86.53723	
KAC 113	<i>americanus</i>	33.00528	-86.38897	
KAC 114	<i>americanus</i>	33.01617	-86.40318	
KAC 115	<i>americanus</i>	32.98218	-86.40488	
KAC 116	<i>americanus</i>	32.96964	-86.42137	X
KAC 117	<i>terrestris</i>	32.97146	-86.52901	
KAC 121	<i>terrestris</i>	32.44120	-85.65386	X
KAC 122	<i>terrestris</i>	32.85411	-86.76619	
KAC 123	<i>terrestris</i>	32.90084	-86.67587	X
KAC 124	<i>terrestris</i>	32.91060	-86.67850	X
KAC 125	<i>terrestris</i>	32.91715	-86.68208	
KAC 126	<i>terrestris</i>	32.92717	-86.67407	
KAC 127	<i>terrestris</i>	32.97159	-86.62516	
KAC 128	<i>terrestris</i>	33.00585	-86.63703	
KAC 129	<i>terrestris</i>	33.00797	-86.64210	
KAC 130	<i>terrestris</i>	33.00818	-86.64333	
KAC 131	<i>terrestris</i>	33.01508	-86.64937	
KAC 132	<i>terrestris</i>	33.02034	-86.66651	
KAC 133	<i>terrestris</i>	33.01163	-86.64759	X
KAC 134	<i>terrestris</i>	33.00537	-86.63652	X
KAC 135	<i>terrestris</i>	33.00644	-86.63368	X
KAC 136	<i>terrestris</i>	33.00673	-86.63316	X
KAC 138	<i>americanus</i>	32.70224	-85.66196	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 139	<i>americanus</i>	32.73042	-85.66173	X
KAC 140	<i>terrestris</i>	32.62553	-85.63684	X
KAC 141	<i>terrestris</i>	32.41032	-85.60107	X
KAC 142	<i>terrestris</i>	32.57011	-85.80888	X
KAC 143	<i>terrestris</i>	32.47773	-85.79824	X
KAC 144	<i>terrestris</i>	32.47707	-85.79577	X
KAC 145	<i>terrestris</i>	32.48128	-85.76354	X
KAC 146	<i>terrestris</i>	32.48291	-85.75622	X
KAC 147	<i>terrestris</i>	32.45001	-85.79652	X
KAC 148	<i>terrestris</i>	32.45420	-85.79408	X
KAC 149	<i>terrestris</i>	32.45449	-85.78664	X
KAC 150	<i>terrestris</i>	32.45449	-85.78664	X
KAC 151	<i>terrestris</i>	32.45451	-85.78416	X
KAC 152	<i>terrestris</i>	32.45423	-85.77634	X
KAC 153	<i>terrestris</i>	32.45423	-85.77634	X
KAC 154	<i>terrestris</i>	32.46574	-85.76977	X
KAC 155	<i>terrestris</i>	32.46961	-85.77369	X
KAC 156	<i>terrestris</i>	32.47709	-85.79175	X
KAC 158	<i>terrestris</i>	32.47709	-85.79175	X
KAC 159	<i>terrestris</i>	32.49000	-85.79741	X
KAC 160	<i>terrestris</i>	32.40809	-85.47857	X
KAC 161	<i>terrestris</i>	32.41744	-85.47117	X
KAC 162	<i>terrestris</i>	32.35417	-86.09838	X
KAC 163	<i>terrestris</i>	32.33994	-86.09946	X
KAC 164	<i>terrestris</i>	32.31562	-86.13789	X
KAC 167	<i>terrestris</i>	33.06620	-86.60328	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 172	<i>americanus</i>	32.62171	-85.61467	X
KAC 173	<i>americanus</i>	32.61751	-85.64335	X
KAC 176	<i>americanus</i>	32.66836	-85.66233	X
KAC 177	<i>americanus</i>	32.65571	-85.57134	X
KAC 181	<i>terrestris</i>	32.38644	-85.23561	X
KAC 182	<i>terrestris</i>	32.38579	-85.23565	X
KAC 183	<i>terrestris</i>	32.38579	-85.23565	X
KAC 184	<i>terrestris</i>	32.38579	-85.23565	X
KAC 185	<i>terrestris</i>	32.38579	-85.23565	X
KAC 187	<i>americanus</i>	32.64548	-85.55135	
KAC 188	<i>terrestris</i>	32.40976	-85.60208	X
KAC 189	<i>terrestris</i>	33.09152	-86.56686	X
KAC 190	<i>terrestris</i>	33.11298	-86.69434	X
KAC 191	<i>terrestris</i>	33.10659	-86.68228	X
KAC 192	<i>terrestris</i>	33.10509	-86.68014	X
KAC 193	<i>terrestris</i>	33.07896	-86.67286	X
KAC 194	<i>terrestris</i>	32.93933	-86.62008	X
KAC 195	<i>terrestris</i>	32.94745	-86.62146	X
KAC 196	<i>terrestris</i>	32.94829	-86.62190	X
KAC 197	<i>terrestris</i>	32.94929	-86.62241	X
KAC 198	<i>terrestris</i>	32.95077	-86.62306	
KAC 199	<i>terrestris</i>	32.95794	-86.62477	X
KAC 200	<i>terrestris</i>	32.95940	-86.62489	X
KAC 205	<i>terrestris</i>	32.54852	-85.48692	X
KAC 206	<i>americanus</i>	33.30759	-86.58201	X
KAC 207	<i>americanus</i>	33.31685	-86.57596	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC 208	<i>americanus</i>	33.09829	-86.56529	X
KAC 209	<i>terrestris</i>	33.08600	-86.56394	X
KAC 210	<i>terrestris</i>	33.08600	-86.56394	X
KAC 211	<i>terrestris</i>	33.01464	-86.60995	
KAC 212	<i>terrestris</i>	33.01208	-86.61707	X
KAC 213	<i>terrestris</i>	33.00435	-86.63710	X
KAC 214	<i>terrestris</i>	32.99991	-86.64181	X
KAC 215	<i>terrestris</i>	32.99605	-86.64526	
KAC 216	<i>terrestris</i>	33.01346	-86.60960	
KAC 217	<i>terrestris</i>	32.91470	-86.60270	X
KAC 218	<i>terrestris</i>	32.92432	-86.59895	X
KAC 219	<i>terrestris</i>	32.93987	-86.56113	X
KAC 220	<i>americanus</i>	32.96579	-86.50892	X
KAC 221	<i>americanus</i>	32.96389	-86.42549	X
KAC 223	<i>terrestris</i>	32.53362	-85.79839	
KAC 224	<i>terrestris</i>	32.48869	-85.79555	X
KAC 225	<i>terrestris</i>	32.50159	-85.79860	X
KAC 230	<i>terrestris</i>	30.80933	-86.77686	X
KAC 232	<i>terrestris</i>	30.80922	-86.78994	X
KAC 233	<i>terrestris</i>	30.80922	-86.78994	X
KAC 234	<i>terrestris</i>	30.80922	-86.78994	X
KAC 236	<i>terrestris</i>	30.82632	-86.80258	X
KAC 237	<i>terrestris</i>	30.83733	-86.77630	X
KAC 238	<i>terrestris</i>	30.82433	-86.76284	X
KAC 239	<i>terrestris</i>	30.80162	-86.76659	X
KAC 242	<i>americanus</i>	34.50446	-85.63768	X

Continued on next page

Table 3.1 – continued from previous page

Sample ID	Species	Latitude	Longitude	Passed Filtering
KAC t1020	<i>terrestris</i>	31.10783	-86.62247	
KAC t1030	<i>terrestris</i>	31.99042	-85.07423	X
KAC t1040	<i>terrestris</i>	31.99016	-85.07046	X
KAC t2004	<i>americanus</i>	33.58295	-85.73524	X
KAC t2015	<i>americanus</i>	33.58435	-85.74064	X
KAC t2018-02-17-01	<i>americanus</i>	33.55274	-85.82913	X
KAC t2018-02-17-04	<i>americanus</i>	33.48548	-85.88857	X
KAC t2018-02-17-05	<i>americanus</i>	33.31649	-86.05293	X
KAC t2018-02-17-06	<i>americanus</i>	33.28443	-86.08443	X
KAC t2018-02-17-07	<i>americanus</i>	33.24576	-86.08168	X
KAC t2018-03-10-1	<i>americanus</i>	32.91057	-86.09272	X
KAC t2018-03-10-3	<i>americanus</i>	32.95104	-86.14539	
KAC t2018-03-10-4	<i>americanus</i>	32.89787	-86.26061	X
KAC t2018-03-10-5	<i>americanus</i>	32.81642	-86.38018	X
KAC t2019-08-25-1	<i>americanus</i>	34.21852	-87.36662	
KAC t2020	<i>americanus</i>	33.23853	-85.96270	X
KAC t2040	<i>americanus</i>	33.58295	-85.73539	X
KAC t2043	<i>americanus</i>	32.81642	-86.38018	X

Table 3.2. Samples loaned from museums

Sample ID	Species	Latitude	Longitude	Passed Filtering
AHT 1975	<i>americanus</i>	32.77356	-85.53325	X
AHT 2456	<i>terrestris</i>	32.19494	-89.23629	X
AHT 2885	<i>terrestris</i>	32.45090	-86.15934	X
AHT 3419	<i>terrestris</i>	33.67290	-88.16068	X
AHT 3421	<i>terrestris</i>	33.65420	-88.15580	X
AHT 3428	<i>terrestris</i>	31.12679	-86.54755	X
AHT 3459	<i>americanus</i>	34.88028	-87.71849	X
AHT 3460	<i>americanus</i>	33.78013	-85.58421	X
AHT 3461	<i>americanus</i>	34.88779	-87.74103	X
AHT 3462	<i>americanus</i>	33.77001	-85.55434	X
AHT 3463	<i>americanus</i>	33.71125	-85.59762	X
AHT 3813	<i>terrestris</i>	31.13854	-86.53906	
AHT 3833	<i>terrestris</i>	31.00422	-85.03427	X
AHT 3997	<i>terrestris</i>	32.55607	-88.29975	X
AHT 3998	<i>terrestris</i>	32.55607	-88.29975	X
AHT 5276	<i>terrestris</i>	31.55613	-86.82514	
AHT 5277	<i>terrestris</i>	31.15830	-86.55430	X
AHT 5278	<i>terrestris</i>	31.16105	-86.69868	X
UTEP 19947	<i>terrestris</i>	31.22432	-88.77548	

# Chapter 4

1359

## Comparison of Linked versus Unlinked

1360

## Character Models for Species Tree

1361

## Inference

1362

### 4.1 Introduction

1363

Current model-based methods of species tree inference require biologists to make difficult decisions about their genomic data. They must decide whether to assume (1) sites in their alignments are each inherited independently (“unlinked”), or (2) groups of sites are inherited together (“linked”). If assuming the former, they must then decide whether to analyze all of their data or only putatively unlinked variable sites. Our goal in this chapter is to use simulated data to help guide these choices by comparing the robustness of different approaches to errors that are likely common in high-throughput genetic datasets.

1364

1365

1366

1367

1368

1369

1370

1371

Reduced-representation genomic data sets acquired from high-throughput instruments are becoming commonplace in phylogenetics (Leaché & Oaks, 2017), and usually comprise hundreds to thousands of loci from 50 to several thousand nucleotides long. Full likelihood approaches for inferring species trees from such datasets can be classified into two groups based on how they model the evolution of orthologous DNA sites along gene trees within

1372

1373

1374

1375

1376

the species tree—those that assume (1) each site evolved along its own gene tree (i.e., each 1377  
site is “unlinked”) (Bryant et al., 2012; De Maio et al., 2015), or (2) contiguous, linked 1378  
sites evolved along a shared gene tree (Heled & Drummond, 2010; Liu & Pearl, 2007; 1379  
Ogilvie et al., 2017; Yang, 2015). We will refer to these as unlinked and linked-character 1380  
models, respectively. For both models, the gene tree of each locus (whether each locus 1381  
is a single site or a segment of linked sites) is assumed to be independent of the gene 1382  
trees of all other loci, conditional on the species tree. Methods using linked character 1383  
models become computationally expensive as the number of loci grows large, due to the 1384  
estimation or numerical integration of all of the gene trees (Ogilvie et al., 2017; Yang, 1385  
2015). Unlinked-character models on the other hand are more tractable for a large number 1386  
of loci, because estimating individual gene trees is avoided by analytically integrating over 1387  
all possible gene trees (Bryant et al., 2012; De Maio et al., 2015). Whereas unlinked- 1388  
character models can accommodate a larger number of loci than linked-character models, 1389  
most genetic data sets comprise linked sites and unlinked-character models are unable to 1390  
utilize the aggregate information about ancestry contained in such linked sites. 1391

Investigators are thus faced with decisions about how best to use their data to infer 1392  
a species tree. Should they use a linked-character method that assumes the sites 1393  
within each locus evolved along a shared gene tree? Ideally, the answer would be “yes,” 1394  
however this is not always computationally feasible and the model could be violated by 1395  
intralocus recombination. Alternatively, should investigators remove all but one single- 1396  
nucleotide polymorphism (SNP) from each locus and use an unlinked-character model? 1397  
Or, perhaps they should apply the unlinked-character method to all of their sites, even if 1398  
this violates the assumption that each site evolved along an independent gene tree. Im- 1399  
portant considerations in such decisions include the sources of error and bias that result 1400  
from reduced-representation protocols, high-throughput sequencing technologies, and the 1401  
processing of these data. 1402

Most reduced-representation sequencing workflows employ amplification of DNA using 1403  
polymerase chain reaction (PCR) which can introduce mutational error at a rate of 1404  
up to  $1.5 \times 10^{-5}$  substitutions per base (Potapov & Ong, 2017). Furthermore, current 1405

high-throughput sequencing technologies have non-negligible rates of error. For example, Illumina sequencing platforms have been shown to have error rates as high as 0.25% per base (Pfeiffer et al., 2018). In hope of removing such errors, it is common for biologists to filter out variants that are not found above some minimum frequency threshold (Linck & Battey, 2019; Rochette et al., 2019). The effect of this filtering will be more pronounced in data sets with low or highly variable coverage. Also, to avoid aligning paralogous sequences, it is common to remove loci that exceed an upper threshold on the number of variable sites (Harvey et al., 2015). These processing steps can introduce errors and acquisition biases, which have been shown to affect estimates derived from the assembled alignments (Harvey et al., 2015; Huang & Knowles, 2016; Linck & Battey, 2019). Given these issues are likely common in high-throughput genomic data, downstream decisions about what methods to use and what data to include in analyses should consider how sensitive the results might be to errors and biases introduced during data collection and processing.

Our goal is to determine whether linked and unlinked character models differ in their robustness to errors in reduced-representation genomic data, and whether it is better to use all sites or only SNPs for unlinked character methods. Linked-character models can leverage shared information among linked sites about each underlying gene tree. Thus, these models might be able to correctly infer the general shape and depth of a gene tree, even if the haplotypes at some of the tips have errors. Unlinked character models have very little information about each gene tree, and rely on the frequency of allele counts across many characters to inform the model about the relative probabilities of all possible gene trees. Given this reliance on accurate allele count frequencies, we predict that unlinked character models will be more sensitive to errors and acquisition biases in genomic data. To test this prediction that linked character models are more robust to the types of errors contained in reduced-representation data, we simulated data sets with varying degrees of errors related to miscalling rare alleles and heterozygous sites. Our results support this prediction, but also show that with only two species, the region of parameter space where there are differences between linked and unlinked character

models is quite limited. Further work is needed to determine whether this difference in robustness between linked and unlinked character models will increase for larger species trees. 1435  
1436  
1437

## 4.2 Methods 1438

### 4.2.1 Simulations of error-free data sets 1439

For our simulations, we assumed a simple two-tipped species tree with one ancestral population with a constant effective size of  $N_e^R$  that diverged at time  $\tau$  into two descendant populations (terminal branches) with constant effective sizes of  $N_e^{D1}$  and  $N_e^{D2}$  (Fig. 4.1). For two diploid individuals sampled from each of the terminal populations (4 sampled gene copies per population), we simulated 100,000 orthologous biallelic characters under a finite-sites, continuous-time Markov chain (CTMC) model of evolution. We simulated 100 data sets comprised of loci of four different lengths—1000, 500, 250, and 1 characters. We assume each locus is effectively unlinked and has no intra-locus recombination; i.e., each locus evolved along a single gene tree that is independent of the other loci, conditional on the species tree. We chose this simple species tree model for our simulations to help ensure any differences in estimation accuracy or precision were due to differences in the underlying linked and unlinked character models, and *not* due to differences in numerical algorithms for searching species and gene tree space. Furthermore, we simulated biallelic characters, because unlinked-character multi-species coalescent models (Bryant et al., 2012; Oaks, 2019) that are most comparable to linked-character models (Heled & Drummond, 2010; Ogilvie et al., 2017) are limited to characters with (at most) two states. 1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461

We simulated the two-tipped species trees under a pure-birth process (Yule, 1925) with a birth rate of 10 using the *Python* package *DendroPy* (Version 4.40, Commit eb69003; Sukumaran & Holder, 2010). This is equivalent to the time of divergence between the two species being Exponentially distributed with a mean of 0.05 substitutions per site. We drew population sizes for each branch of the species tree from a Gamma distribution with 1458  
1459  
1460  
1461

a shape of 5.0 and mean of 0.002. We simulated 100, 200, 400, and 100,000 gene trees  
1462  
for data sets with loci of length 1000, 500, 250, and 1, respectively, using the contained  
1463  
coalescent implemented in *DendroPy*. We simulated linked biallelic character alignments  
1464  
using *Seq-Gen* (Version 1.3.4) (Rambaut & Grass, 1997) with a GTR model with base  
1465  
frequencies of A and C equal to 0 and base frequencies of G and T equal to 0.5. The  
1466  
transition rate for all base changes was 0, except for the rate between G and T which  
1467  
was 1.0.  
1468

#### 4.2.2 Introducing Site-pattern Errors 1469

From each simulated dataset containing linked characters described above, we created  
1470  
four datasets by introducing two types of errors at two levels of frequency. The first  
1471  
type of error we introduced was changing singleton character patterns (i.e., characters  
1472  
for which one gene copy was different from the other seven gene copies) to invariant  
1473  
patterns by changing the singleton character state to match the other gene copies. We  
1474  
introduced this change to all singleton site patterns with a probability of 0.2 and 0.4 to  
1475  
create two datasets from each simulated dataset. The second type of error we introduced  
1476  
was missing heterozygous gene copies. To do this, we randomly paired gene copies from  
1477  
within each species to create two diploid genotypes for each locus, and with a probability  
1478  
of 0.2 or 0.4 we randomly replaced one allele of each genotype with the other. For  
1479  
the unlinked character dataset comprised of a single site per locus, we only simulated  
1480  
singleton character pattern error at a probability of 0.4.  
1481

#### 4.2.3 Assessing Sensitivity to Errors 1482

For each simulated data set with loci of 250, 500, and 1000 characters, we approx-  
1483  
imated the posterior distribution of the divergence time ( $\tau$ ) and effective population  
1484  
sizes ( $N_e^R$ ,  $N_e^{D1}$ , and  $N_e^{D2}$ ) under an unlinked-character model using *ecoevolity* (Version  
1485  
0.3.2, Commit a7e9bf2; Oaks, 2019) and a linked-character model using the *StarBEAST2*  
1486  
package (Version 0.15.1; Ogilvie et al., 2017) in *BEAST2* (Version 2.5.2; Bouckaert et al.,  
1487  
2014). For both methods, we specified a CTMC model of character evolution and prior  
1488

distributions that matched the model and distributions from which the data were generated. The prior on the effective size of the root population in the original implementation of *ecoevolity* was parameterized to be relative to the mean effective size of the descendant populations. We added an option to *ecoevolity* to compile a version where the prior is specified as the absolute effective size of the root population, which matches the model in *StarBEAST2* and the model we used to generate the data. The linkage of sites within loci of our simulated data violates the unlinked-character model of *ecoevolity* (Bryant et al., 2012; Oaks, 2019). Therefore, we also analyzed each data set with *ecoevolity* after selecting, at most, one variable character from each locus; loci without variable sites were excluded.

We analyzed the data sets simulated with 1-character per locus (i.e., unlinked data) with *ecoevolity*. Our goal with these analyses was to verify that the generative model of our simulation pipeline matched the underlying model of *ecoevolity*, and to confirm that any behavior of the method with the other simulated data sets was not being caused by the linkage violation.

For *ecoevolity*, we ran four independent Markov chain Monte Carlo (MCMC) analyses with 75,000 steps and a sample frequency of 50 steps. For *StarBEAST2*, we ran two independent MCMC analyses with 20 million steps and a sample frequency of 5000 steps. To assess convergence and mixing of the *ecoevolity* and *StarBEAST2* MCMC chains, we computed the effective sample size (ESS; Gong & Flegal, 2016) and potential scale reduction factor (PSRF; the square root of Equation 1.1 in Brooks & Gelman, 1998) from the samples of each parameter, and considered an ESS value greater than 200 and PSRF less than 1.2 (Brooks & Gelman, 1998) to indicate adequate convergence and mixing of the chains. Based on preliminary analyses of simulated data sets without errors, we chose to discard the first 501 and 201 samples from the MCMC chains of *ecoevolity* and *StarBEAST2*, leaving 4000 and 7600 posterior samples for each data set, respectively.

#### 4.2.4 Project repository

1515

The full history of this project has been version-controlled and is available at <https://github.com/kerrycobb/align-error-sp-tree-sim>, and includes all of the data and scripts necessary to produce our results.

1516

1517

1518

### 4.3 Results

1519

#### 4.3.1 Behavior of linked (*StarBEAST2*) versus unlinked (*eco-evolity*) character models

1520

1521

The divergence times estimated by the linked-character method, *StarBEAST2*, were very accurate and precise for all alignment lengths and types and degrees errors, despite poor MCMC mixing (i.e., low ESS values) for shorter loci (Figs. 4.2–4.4). For data sets without error, the unlinked-character method, *ecoevolity*, estimated divergence times with similar accuracy and precision as *StarBEAST2* when all characters are analyzed (Figs. 4.2–4.4). However when alignments contained errors, *ecoevolity* underestimated very recent divergence times with increasing severity as the frequency of errors increased (Figs. 4.2–4.4); estimates of older divergence times were unaffected.

1522

1523

1524

1525

1526

1527

1528

1529

The biased underestimation of divergence times by *ecoevolity* in the face of errors was coupled with overestimation of the ancestral effective population sizes (Figs. 4.5–4.7). When analyzing the alignments without errors, *ecoevolity* essentially returned the prior distribution on the effective size of the ancestral population (Figs. 4.5–4.7). Despite poor MCMC mixing, *StarBEAST2* consistently estimated the effective size of the ancestral population better than *ecoevolity* and was unaffected by errors in the data (Figs. 4.5–4.7), and the precision of *StarBEAST2*'s estimates of  $N_e^R$  increased with locus length.

1530

1531

1532

1533

1534

1535

1536

Estimates of the effective size of the descendant populations are largely similar between *StarBEAST2* and *ecoevolity*; both methods underestimate the descendant population sizes when the data sets contain errors, and this downward bias is generally worse for *ecoevolity* (Figs. 4.8–4.10). The degree of underestimation increases with the rate of

1537

1538

1539

1540

errors in the data sets for both *StarBEAST2* and *ecoevolity*, and the results were largely  
1541  
consistent across different locus lengths. (Figs. 4.8–4.10).  
1542

When we apply *ecoevolity* to data sets simulated with unlinked characters (i.e., data  
1543  
sets simulated with 1-character per locus), we see the same patterns of biased parameter  
1544  
estimates in response to errors (Fig. 4.11) as we did with the linked loci (Figs. 4.2–4.4).  
1545  
These results rule out the possibility that the greater sensitivity of *ecoevolity* to the  
1546  
errors we simulated is due to violation of the method’s assumption that all characters are  
1547  
unlinked.  
1548

### 4.3.2 Analyzing all sites versus SNPs with *ecoevolity* 1549

The unlinked character model implemented in *ecoevolity* assumes that orthologous  
1550  
nucleotide sites evolve independently along separate gene trees. The data however, were  
1551  
simulated under a model assuming that contiguous linked sites evolve along a shared  
1552  
gene tree. It would thus be a violation of the *ecoevolity* model to include all sites in  
1553  
the analysis. However, avoiding this violation by removing all but one variable site per  
1554  
locus drastically reduces the amount of data. When analyzing the simulated data sets  
1555  
without errors, the precision and accuracy of parameter estimates by *ecoevolity* was much  
1556  
greater when all sites of the alignment were used relative to when a single SNP per locus  
1557  
was used despite violating the model (Figs. 4.2–4.10). This was generally true across the  
1558  
different lengths of loci, however, the coverage of credible intervals is lower with longer  
1559  
loci. Analyzing only SNPs does make *ecoevolity* more robust to the errors we introduced.  
1560  
However, this robustness is due to the lack of information in the SNP data leading to wide  
1561  
credible intervals, and in the case of population size parameters, the marginal posteriors  
1562  
essentially match the prior distribution (Figs. 4.8–4.10).  
1563

### 4.3.3 Coverage of credible intervals 1564

The 95% credible intervals for divergence times and effective population sizes esti-  
1565  
mated from alignments without error in *StarBEAST2* had the expected coverage fre-  
1566  
quency in that the true value was within approximately 95% of the estimated credible  
1567

intervals. This was also true for *ecoevolity* when analyzing data sets simulated with un-linked characters (i.e., no linked sites). This coverage behavior is expected, and helps to confirm confirm that our simulation pipeline generated data under the same model used for inference by *StarBEAST2* and *ecoevolity*. As seen previously (Oaks, 2019), analyzing longer linked loci causes the coverage of *ecoevolity* to be lower, due to the violation of the model’s assumption that the sites are unlinked.

#### 4.3.4 MCMC convergence and mixing

Most sets of *StarBEAST2* and *ecoevolity* MCMC chains yielded samples of parameters with a PSRF less than 1.2, indicative of convergence. However, we do see poor mixing (ESS < 200) of the *StarBEAST2* chains as the length of loci decreases (Figs. 4.2–4.10; yellow indicates ESS < 200, red indicates PSRF > 1.2, green indicates both) We only see evidence of poor mixing and convergence for *ecoevolity* when applied to data sets with errors. This is in contrast to *StarBEAST2*, for which the frequency and degree of poor MCMC behavior is largely unaffected by the type or frequency of errors. The proportion of simulation replicates where *StarBEAST2* had an ESS of the ancestral population size less than 200 was high across all analyses (Figs. 4.5–4.7). For the descendant population size, *StarBEAST2* had better ESS values across all analyses, with the exception of rare estimates of essentially zero when analyzing 250 bp loci (Figs. 4.8–4.10).

## 4.4 Discussion

Phylogeneticists seeking to infer species trees from large, multi-locus data sets are faced with difficult decisions regarding assumptions about linkage across sites and, if assuming all sites are unlinked, what data to include in their analysis. With the caveat that we only explored trees with two species, the results of our simulations provide some guidance for these decisions. As we predicted, the linked-character method we tested, *StarBEAST2*, was more robust to the sequencing errors we simulated than the unlinked character method, *ecoevolity*. However, even with only two species in our simulations, the

current computational limitations of linked-character models was apparent from the poor 1594  
sampling efficiency of the MCMC chains, especially with shorter loci. For data sets with 1595  
more species and many short loci, linked character models are theoretically appealing, but 1596  
current implementations may not be computationally feasible. The unlinked character 1597  
method, *ecoevolity*, was more sensitive to sequence errors, but was still quite robust to 1598  
realistic levels of errors and is more computationally feasible thanks to the analytical 1599  
integration over gene trees. 1600

Overall, for data sets with relatively long loci, as is common with sequence-capture 1601  
approaches, it might be worth trying a linked-character method. If computationally 1602  
practical, you stand to benefit from the aggregate information about each gene tree 1603  
contained in the linked sites of each locus. However, if your loci are shorter, as in 1604  
restriction-site-associated DNA (RAD) markers, you are likely better off applying an 1605  
unlinked-character model to all of your data, even though this violates an assumption of 1606  
the model. Below we discuss why performance differs between methods, locus lengths, 1607  
and degree of error in the data, and what this means for the analyses of empirical data. 1608

#### 4.4.1 Robustness to character-pattern errors 1609

As predicted, the linked-character model of *StarBEAST2* was more robust to erro- 1610  
neous character patterns in the alignments than the unlinked-character model of *eco- 1611  
evolity*. This is most evident in the estimates of divergence times, for which the two 1612  
methods perform very similarly when there are no errors in the data (Row 1 of Figs. 4.2– 1613  
4.4). When errors are introduced, the divergence time estimates of *StarBEAST2* are 1614  
unaffected, but *ecoevolity* underestimates recent divergence times as both singleton and 1615  
heterozygosity errors become more frequent (Rows 2–5 of Figs. 4.2–4.4). However, *eco- 1616  
evolity* divergence-time estimates are only biased at very recent divergence times, and 1617  
the effect disappears when the time of divergence is larger than about  $8N_e\mu$ . 1618

These patterns make sense given that both types of errors we simulated reduce varia- 1619  
tion *within* each species. Thus, it is not too surprising that the unlinked-character model 1620  
in *ecoevolity* struggles when there is shared variation between the two populations (i.e., 1621

most gene trees have more than two lineages that coalesce in the ancestral population). 1622  
The erroneous character patterns mislead both models that the effective size of the de- 1623  
scendant branches is smaller than they really are (Figs. 4.8–4.10). To explain the shared 1624  
variation between the species (i.e. deep coalescences) when underestimating the descen- 1625  
dant population sizes, the unlinked-character model of *ecoevolity* simultaneously reduces 1626  
the divergence time and increases the effective size of the ancestral population. De- 1627  
spite also being misled about the size of the descendant populations (Figs. 4.8–4.10), the 1628  
linked-character model of *StarBEAST2* seems to benefit from more information about the 1629  
general shape of each gene tree across the linked sites and can still maintain an accurate 1630  
estimate of the divergence time (Figs. 4.2–4.4) and ancestral population size (Figs. 4.5– 1631  
4.7). 1632

This downward biased variation within each species becomes less of a problem for the 1633  
unlinked-character model as the divergence time gets larger, likely because the average 1634  
gene tree only has a single lineage from each species that coalesces in the ancestral 1635  
population. As the coalesced lineage within each species leading back to the ancestral 1636  
population becomes a large proportion of the overall length of the average gene tree, 1637  
the proportion of characters that either show fixed differences between the species or 1638  
are invariant likely provides enough information to the unlinked character model about 1639  
the time of divergence to overcome the downward biased estimates of the descendant 1640  
population sizes. 1641

From the *ecoevolity* results, we also see that when faced with heterozygosity errors, 1642  
accuracy decreases as locus length increases. In contrast, accuracy of *ecoevolity* is not 1643  
affected by locus length when analyzing data sets with singleton errors. This pattern 1644  
makes sense in light of how we generated these errors. We introduced singleton errors per- 1645  
site and heterozygosity errors per-locus. Thus, the same per-locus rate of heterozygosity 1646  
errors affects many more sites of a dataset with 1000bp loci compared to dataset with 1647  
250bp loci. 1648

Unsurprisingly, the MCMC sampling performance of *StarBEAST2* declines with de- 1649  
creasing locus length. There is less information in the shorter loci about ancestry, and 1650

thus more posterior uncertainty about the gene trees. This forces *StarBEAST2* to traverse a much broader distribution of gene trees during MCMC sampling, which is difficult due to the constraints imposed by the species tree. This decline in MCMC performance in *StarBEAST2* does not appear to correlate with poor parameter estimates and the distribution of estimates is generally as good or better than those from *ecoevolity*. However, this might be due to fact that there is no uncertainty in the species tree in any of our analyses, because there are only two species. As the number of species increases, it seems likely that the MCMC performance will further decline and start to affect parameter and topology estimates.

#### 4.4.2 Relevance to empirical data sets

It is reassuring to see the effect of sequence errors on the unlinked-character model is limited to a small region of parameter space, and is only severe when the frequency of errors in the data is large. Our simulated error rate of 40% is likely higher than the rate that these types of errors occur during most sample preparation, high-throughput sequencing, and bioinformatic processing. However, empirical alignments likely contain a mix of different sources of errors and biases from various steps in the data collection process. Also, real data are not be generated under a known model with no prior misspecification. Violations of the model might make these methods of species-tree inference more sensitive to lower rates of error.

The degree to which a dataset will be affected by errors from missing heterozygote haplotypes and missing singletons will be highly dependent on the method used to reduce representation of the genome, depth of sequencing coverage (i.e., the number of overlapping sequence reads at a locus), and how the data are processed. To filter out sequencing errors, most pipelines for processing sequence reads set a minimum coverage threshold for variants or a minimum minor allele frequency. This can result in the miscalling or removal of true variation, especially if coverage is low due to random chance or biases in PCR amplification and sequencing. Processing the data in this way can result in biased estimates of parameters that are sensitive to the frequencies of rare alleles (Huang &

Knowles, 2016; Linck & Battey, 2019). If the thresholds for such processing steps are 1679  
stringent, it could introduce levels of error greater than our simulations. 1680

#### 4.4.3 Recommendations for using unlinked-character models 1681

When erroneous character patterns cause *ecoevolity* to underestimate the divergence 1682  
time it also inflates the effective population size of the ancestral population. We are 1683  
seeing values of  $N_e^R \mu$  consistent with an average sequence divergence between individuals 1684  
*within* the ancestral population of 3%, which is almost an order of magnitude larger than 1685  
our prior mean expectation (0.4%). Thus, looking for unrealistically large population 1686  
sizes estimated for internal branches of the phylogeny might provide an indication that 1687  
the unlinked-character model is not explaining the data well. However, there is little 1688  
information in the data about the effective population sizes along ancestral branches, so 1689  
the parameter that might indicate a problem is going to have very large credible intervals. 1690  
Nonetheless, many of the posterior estimates of the ancestral population size from our 1691  
data sets simulated with character-pattern errors are well beyond the prior distribution. 1692

Whether using linked or unlinked-character models with empirical high-throughput 1693  
data sets, it is good practice to perform analyses on different versions of the aligned data 1694  
that are assembled under different coverage thresholds for variants or alleles. Variation 1695  
of estimates derived from different assemblies of the data might indicate that the model 1696  
is sensitive to the errors or acquisition biases in the alignments. This is especially true 1697  
for data where sequence coverage is low for samples and/or loci. Given our findings, it 1698  
might be helpful to compare the estimates of the effective population sizes along internal 1699  
branches of the tree. Seeing unrealistically large estimates for some assemblies of the 1700  
data might indicate that the model is being biased by errors or acquisition biases present 1701  
in the character patterns. 1702

Consistent with what has been shown in previous work (Oaks, 2019; Oaks et al., 2019), 1703  
*ecoevolity* performed better when all sites were utilized despite violating the assumption 1704  
that all sites are unlinked. This suggests that investigators might obtain better estimates 1705  
by analyzing all their data under unlinked-character models, rather than discarding much 1706

of it to avoid violating an assumption of the model. Given that the model of unlinked  
1707  
characters implemented in *ecoevolity* does not use information about linkage among sites  
1708  
(Bryant et al., 2012; Oaks, 2019), it is not surprising that this model violation does not  
1709  
introduce a bias. Linkage among sites does not change the gene trees and site patterns  
1710  
that are expected under the model, but it does reduce the variance of those patterns  
1711  
due to them evolving along fewer gene trees. As a result, the accuracy of the parameter  
1712  
estimates is not affected by the linkage among sites within loci, but the credible intervals  
1713  
become too narrow as the length of loci increase (Oaks, 2019; Oaks et al., 2019). However,  
1714  
it remains to be seen whether the robustness of the model’s accuracy to linked sites holds  
1715  
true for larger species trees.  
1716

#### 4.4.4 Other complexities of empirical data in need of exploration

1717

Our goal was to compare the theoretical performance of linked and unlinked character  
1718  
models, not their current software implementations. Accordingly, to minimize differences  
1719  
in performance that are due to differences in algorithms for exploring the space of gene  
1720  
and species trees, we restricted our simulations to two species model and a small number  
1721  
of individuals. Nonetheless, exploring how character-pattern errors and biases affect the  
1722  
inference of larger species trees would be informative. The species tree topology is usually  
1723  
a parameter of great interest to biologists, so it would be interesting to know whether  
1724  
the linked model continues to be more robust to errors than the unlinked model as the  
1725  
number of species increases. We saw the MCMC performance of *StarBEAST2* decline  
1726  
concomitantly with locus length in our simulations due to greater uncertainty in gene  
1727  
trees. Given that data sets frequently contain loci shorter than 250 bp, it is important  
1728  
to know whether good sampling of the posterior of linked-character models becomes  
1729  
prohibitive for larger trees. Also, *ecoevolity* greatly overestimated the effective size of  
1730  
the ancestral population in the face of high rates of errors in the data. Exploring larger  
1731  
trees will also determine whether this behavior is limited to the root population or is a  
1732  
potential problem for all internal branches of the species tree.  
1733

Exploring other types of errors and biases would also be informative. To generate  
1734

alignments of orthologous loci from high-throughput data, sequences are matched to a 1735  
similar portion of a reference sequence or clustered together based on similarity. To avoid 1736  
aligning paralogous sequences it is necessary to establish a minimum level of similarity for 1737  
establishing orthology between sequences. This can lead to an acquisition bias due to the 1738  
exclusion of more variable loci or alleles from the alignment (Huang & Knowles, 2016). 1739  
Furthermore, when a reference sequence is used, this data filtering will not be random 1740  
with respect to the species, but rather there will be a bias towards filtering loci and alleles 1741  
with greater sequence divergence from the reference. Simulations exploring the affect of 1742  
these types of data acquisition biases would complement the errors we explored here. 1743

In our analyses, there was no model misspecification other than the introduced er- 1744  
rors (except for the linked sites violating the unlinked-character model). With empirical 1745  
data, there are likely many model violations, and our prior distributions will never match 1746  
the distributions that generated the data. Introducing other model violations and mis- 1747  
specified prior distributions would thus help to better understand how species-tree models 1748  
behave on real data sets. Of particular concern is whether misspecified priors will amplify 1749  
the effect of character-pattern errors or biases. 1750

We found that character-pattern errors that remove variation from within species 1751  
can cause unlinked-character models to underestimate divergence times and overestimate 1752  
ancestral population sizes in order to explain shared variation among species. This raises 1753  
the question of whether we can explicitly model and correct for these types of data 1754  
collection errors in order to avoid biased parameter estimates. An approach that could 1755  
integrate over uncertainty in the frequency of these types of missing-allele errors would 1756  
be particularly appealing. 1757

## References

- Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.-H., Xie, D., Suchard, M. A., 1759  
Rambaut, A., & Drummond, A. J. (2014). BEAST 2: A Software Platform for 1760

Bayesian Evolutionary Analysis (A. Prlic, Ed.). <i>PLoS Computational Biology</i> , 10(4), e1003537. <a href="https://doi.org/10.1371/journal.pcbi.1003537">https://doi.org/10.1371/journal.pcbi.1003537</a>	1761 1762
Brooks, S. P., & Gelman, A. (1998). General methods for monitoring convergence of iterative simulations. <i>Journal of computational and graphical statistics</i> , 7(4), 434–455.	1763 1764 1765
Bryant, D., Bouckaert, R., Felsenstein, J., Rosenberg, N. A., & RoyChoudhury, A. (2012). Inferring Species Trees Directly from Biallelic Genetic Markers: Bypassing Gene Trees in a Full Coalescent Analysis. <i>Molecular Biology and Evolution</i> , 29(8), 1917–1932. <a href="https://doi.org/10.1093/molbev/mss086">https://doi.org/10.1093/molbev/mss086</a>	1766 1767 1768 1769
De Maio, N., Schrempf, D., & Kosiol, C. (2015). PoMo: An Allele Frequency-Based Approach for Species Tree Estimation. <i>Systematic Biology</i> , 64(6), 1018–1031. <a href="https://doi.org/10.1093/sysbio/syv048">https://doi.org/10.1093/sysbio/syv048</a>	1770 1771 1772
Gong, L., & Flegal, J. M. (2016). A Practical Sequential Stopping Rule for High-Dimensional Markov Chain Monte Carlo. <i>Journal of Computational and Graphical Statistics</i> , 25(3), 684–700. <a href="https://doi.org/10.1080/10618600.2015.1044092">https://doi.org/10.1080/10618600.2015.1044092</a>	1773 1774 1775
Harvey, M. G., Judy, C. D., Seeholzer, G. F., Maley, J. M., Graves, G. R., & Brumfield, R. T. (2015). Similarity thresholds used in DNA sequence assembly from short reads can reduce the comparability of population histories across species. <i>PeerJ</i> , 3, e895. <a href="https://doi.org/10.7717/peerj.895">https://doi.org/10.7717/peerj.895</a>	1776 1777 1778 1779
Heled, J., & Drummond, A. J. (2010). Bayesian Inference of Species Trees from Multilocus Data. <i>Molecular Biology and Evolution</i> , 27(3), 570–580. <a href="https://doi.org/10.1093/molbev/msp274">https://doi.org/10.1093/molbev/msp274</a>	1780 1781 1782
Huang, H., & Knowles, L. L. (2016). Unforeseen Consequences of Excluding Missing Data from Next-Generation Sequences: Simulation Study of RAD Sequences. <i>Systematic Biology</i> , 65(3), 357–365. <a href="https://doi.org/10.1093/sysbio/syu046">https://doi.org/10.1093/sysbio/syu046</a>	1783 1784 1785
Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. <i>Computing in Science &amp; Engineering</i> , 9(3), 90–95. <a href="https://doi.org/10.1109/MCSE.2007.55">https://doi.org/10.1109/MCSE.2007.55</a>	1786 1787

Leaché, A. D., & Oaks, J. R. (2017). The Utility of Single Nucleotide Polymorphism (SNP) Data in Phylogenetics. <i>Annual Review of Ecology, Evolution, and Systematics</i> , 48(1), 69–84. <a href="https://doi.org/10.1146/annurev-ecolsys-110316-022645">https://doi.org/10.1146/annurev-ecolsys-110316-022645</a>	1788
	1789
Linck, E., & Battey, C. J. (2019). Minor allele frequency thresholds strongly affect population structure inference with genomic data sets. <i>Molecular Ecology Resources</i> , 19(3), 639–647. <a href="https://doi.org/10.1111/1755-0998.12995">https://doi.org/10.1111/1755-0998.12995</a>	1790
	1791
Liu, L., & Pearl, D. K. (2007). Species Trees from Gene Trees: Reconstructing Bayesian Posterior Distributions of a Species Phylogeny Using Estimated Gene Tree Distributions (T. Buckley, Ed.). <i>Systematic Biology</i> , 56(3), 504–514. <a href="https://doi.org/10.1080/10635150701429982">https://doi.org/10.1080/10635150701429982</a>	1794
	1795
Oaks, J. R. (2019). Full Bayesian Comparative Phylogeography from Genomic Data (L. Kubatko, Ed.). <i>Systematic Biology</i> , 68(3), 371–395. <a href="https://doi.org/10.1093/sysbio/syy063">https://doi.org/10.1093/sysbio/syy063</a>	1798
	1799
Oaks, J. R., Siler, C. D., & Brown, R. M. (2019). The comparative biogeography of Philippine geckos challenges predictions from a paradigm of climate-driven vicariant diversification across an island archipelago. <i>Evolution</i> , 73(6), 1151–1167. <a href="https://doi.org/10.1111/evo.13754">https://doi.org/10.1111/evo.13754</a>	1800
	1801
Ogilvie, H. A., Bouckaert, R. R., & Drummond, A. J. (2017). StarBEAST2 Brings Faster Species Tree Inference and Accurate Estimates of Substitution Rates. <i>Molecular Biology and Evolution</i> , 34(8), 2101–2114. <a href="https://doi.org/10.1093/molbev/msx126">https://doi.org/10.1093/molbev/msx126</a>	1802
	1803
Pfeiffer, F., Gröber, C., Blank, M., Händler, K., Beyer, M., Schultze, J. L., & Mayer, G. (2018). Systematic evaluation of error rates and causes in short samples in next-generation sequencing. <i>Scientific Reports</i> , 8(1), 10950. <a href="https://doi.org/10.1038/s41598-018-29325-6">https://doi.org/10.1038/s41598-018-29325-6</a>	1804
	1805
Potapov, V., & Ong, J. L. (2017). Examining Sources of Error in PCR by Single-Molecule Sequencing (R. Kalendar, Ed.). <i>PLOS ONE</i> , 12(1), e0169774. <a href="https://doi.org/10.1371/journal.pone.0169774">https://doi.org/10.1371/journal.pone.0169774</a>	1806
	1807
	1808
	1809
	1810
	1811
	1812
	1813
	1814
	1815

- Rambaut, A., & Grass, N. C. (1997). Seq-Gen: An application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics*, 13(3), 235–238. <https://doi.org/10.1093/bioinformatics/13.3.235>
- Rochette, N. C., Rivera-Colón, A. G., & Catchen, J. M. (2019). Stacks 2: Analytical methods for paired-end sequencing improve RADseq-based population genomics. *Molecular Ecology*, 28(21), 4737–4754. <https://doi.org/10.1111/mec.15253>
- Sukumaran, J., & Holder, M. T. (2010). DendroPy: A Python library for phylogenetic computing. *Bioinformatics*, 26(12), 1569–1571. <https://doi.org/10.1093/bioinformatics/btq228>
- Yang, Z. (2015). The BPP program for species tree estimation and species delimitation. *Current Zoology*, 61(5), 854–865. <https://doi.org/10.1093/czoolo/61.5.854>
- Yule, G. U. (1925). A mathematical theory of evolution, based on the conclusions of Dr. JC Willis, FR S. *Philosophical transactions of the Royal Society of London. Series B, containing papers of a biological character*, 213(402-410), 21–87.

## 4.5 Figures

1830

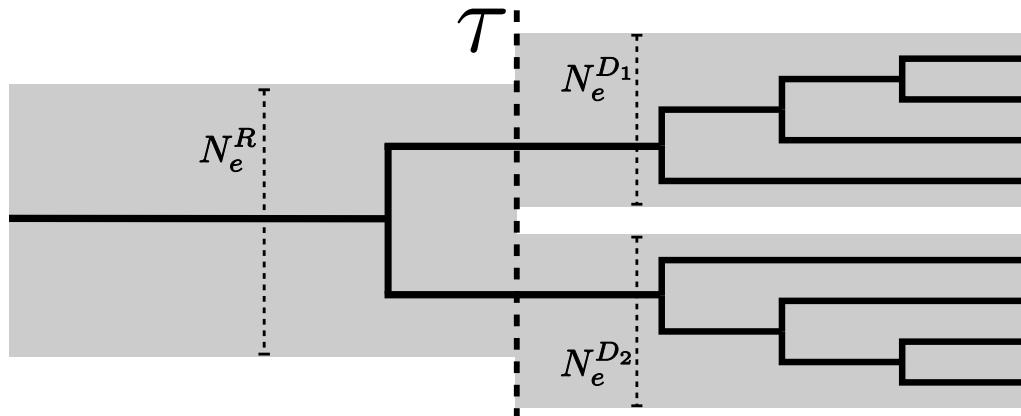
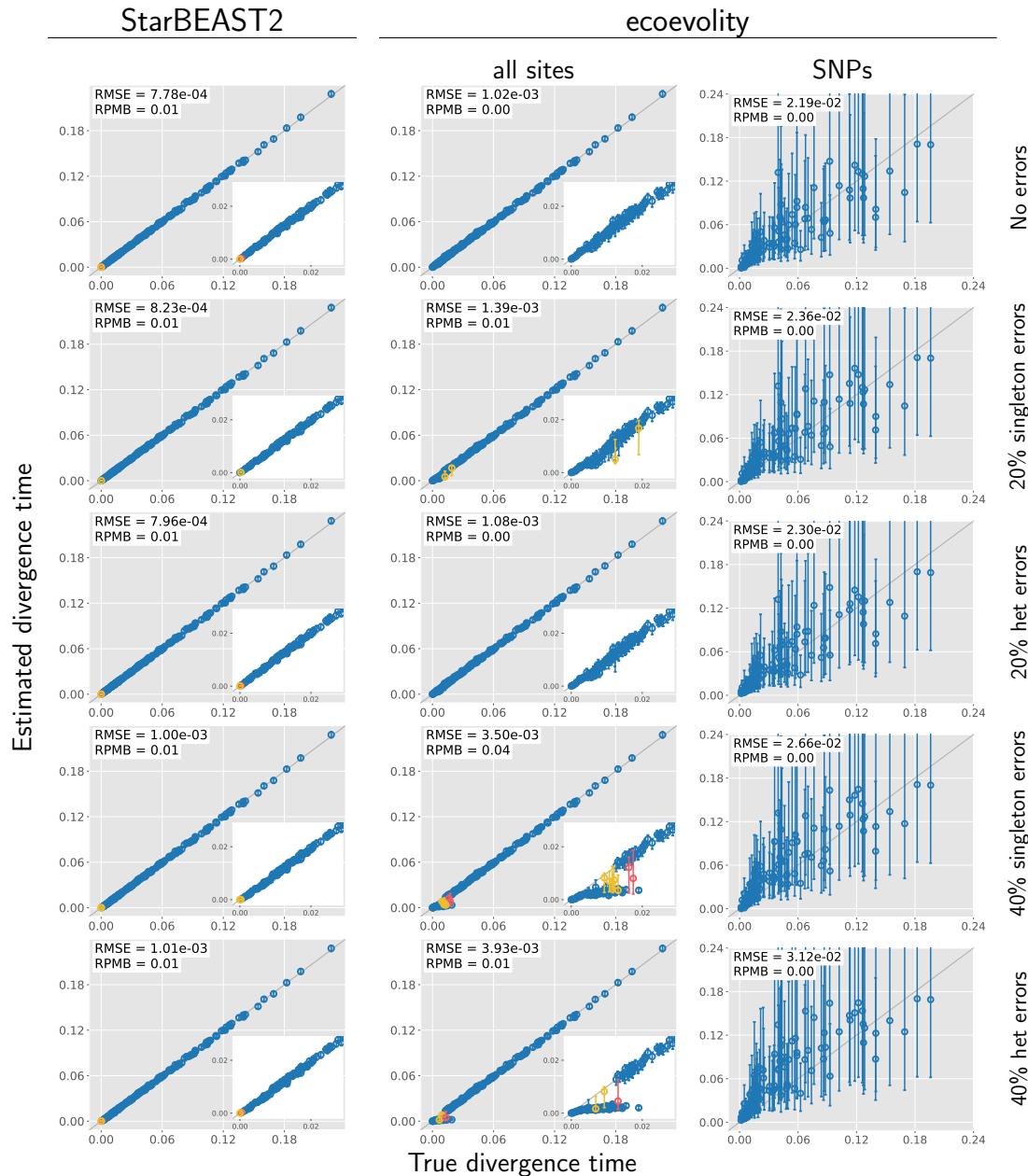


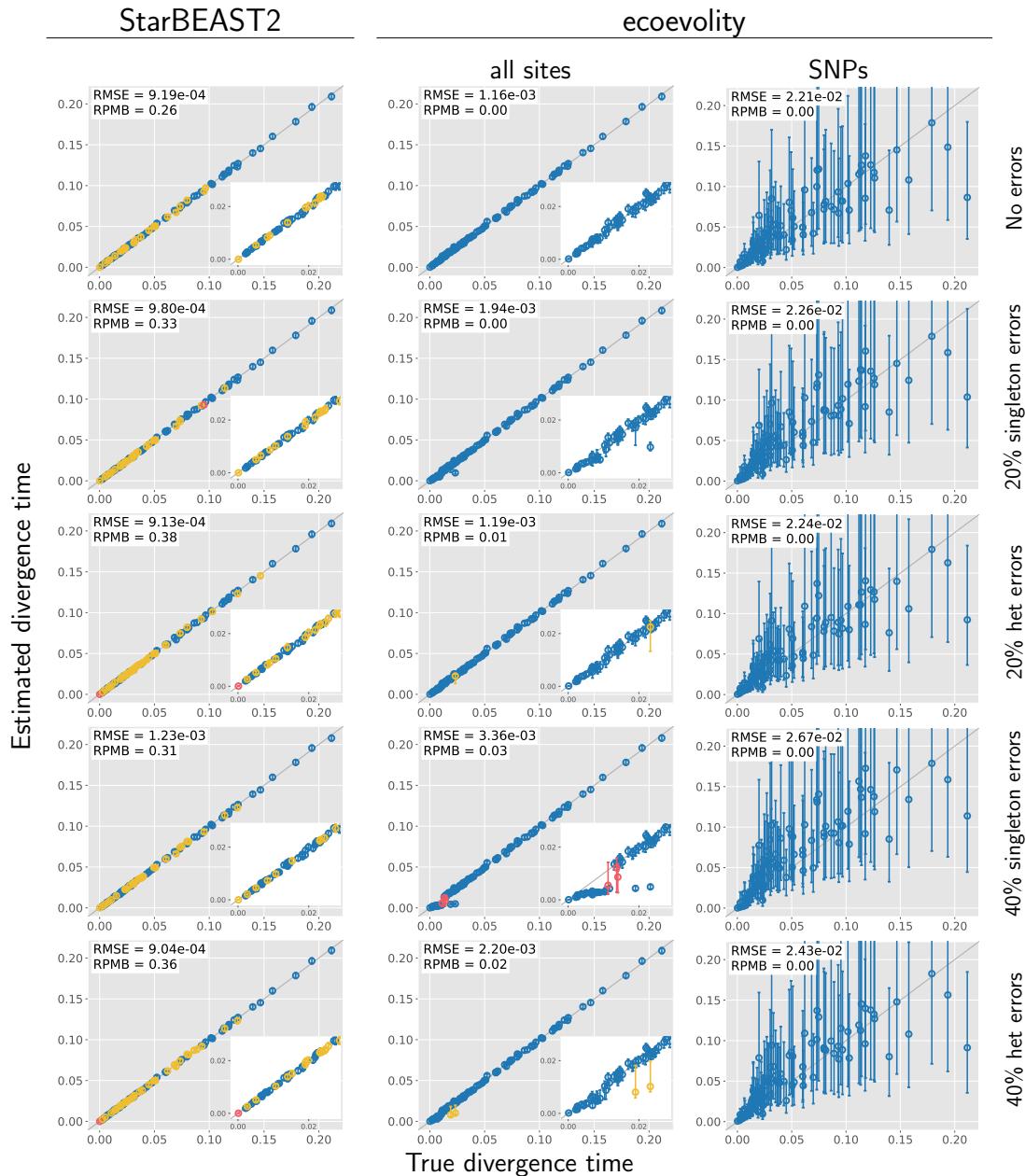
Figure 4.1. An illustration of the species-tree model we used to simulate data.  $N_e^R$ ,  $N_e^{D1}$ , and  $N_e^{D2}$  represent the constant effective population sizes of the root, and each of the two terminal populations.  $\tau$  represents the instantaneous separation of the ancestral population into two descendant populations. One hypothetical gene tree is shown to illustrate the gene trees simulated under a contained coalescent process for 4 haploid gene copies sampled from each of the terminal branches of the species tree.

## Divergence Time — 1000bp loci



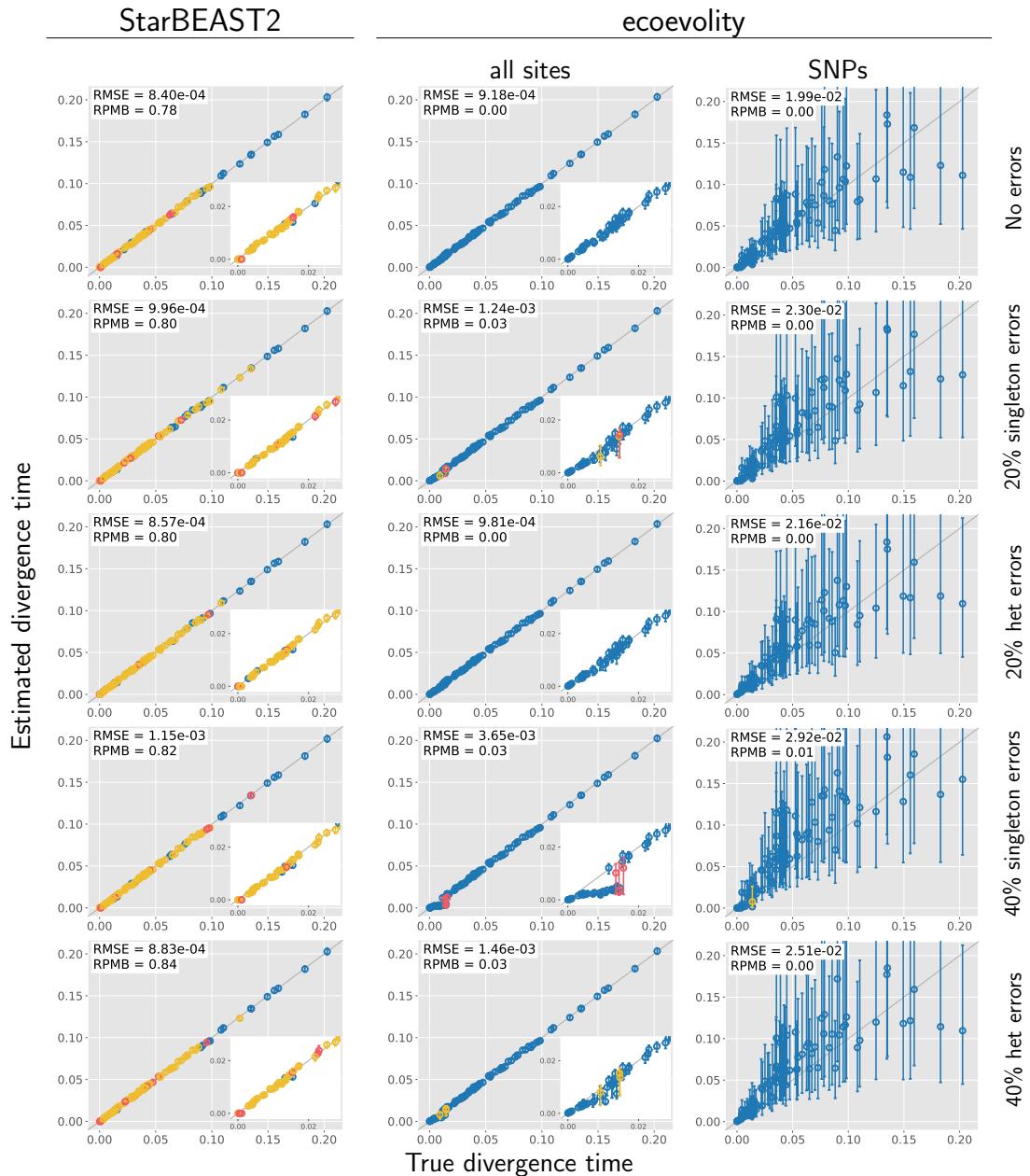
**Figure 4.2. Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 1000 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

### Divergence Time — 500bp loci

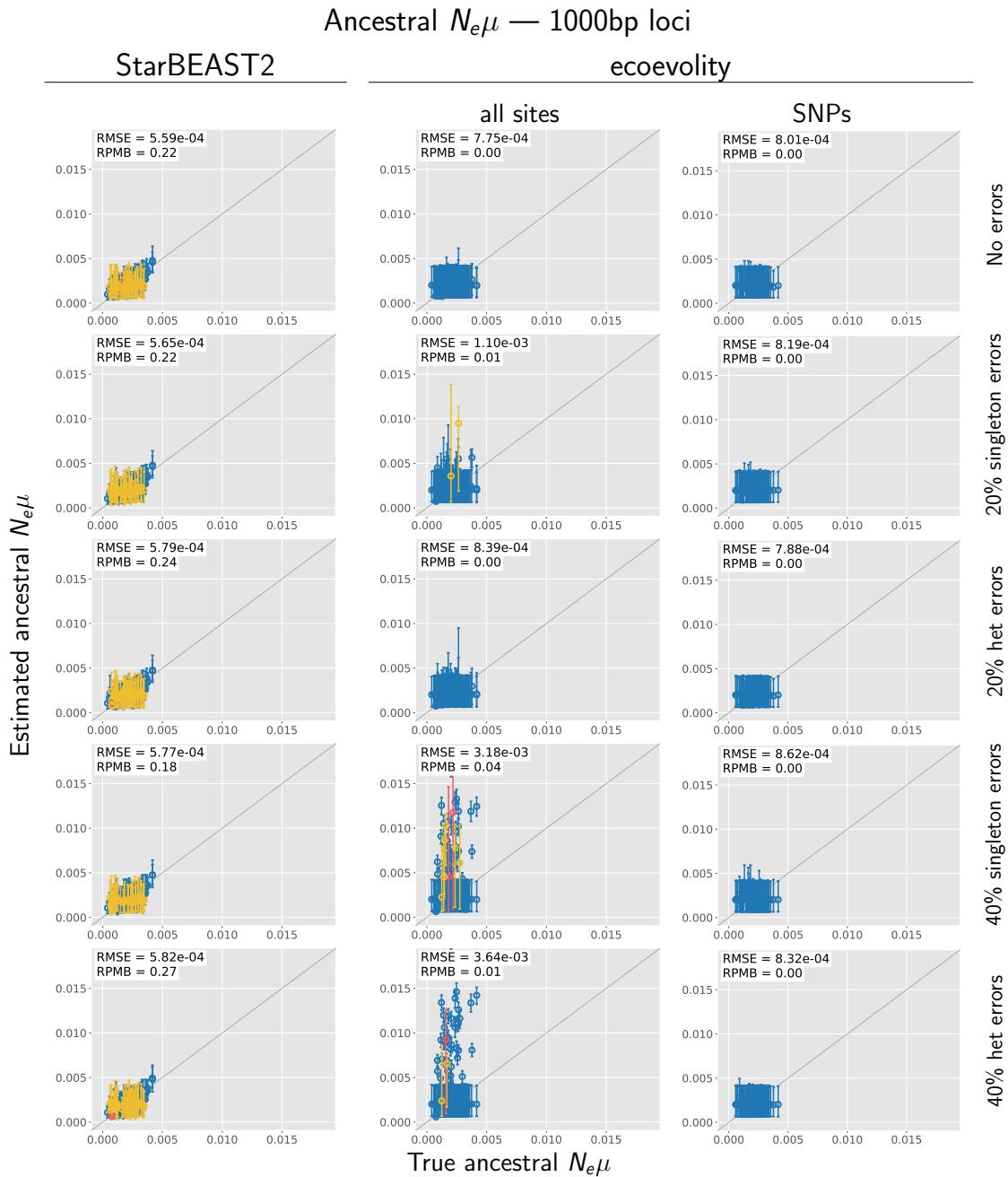


**Figure 4.3. Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 500 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

## Divergence Time — 250bp loci

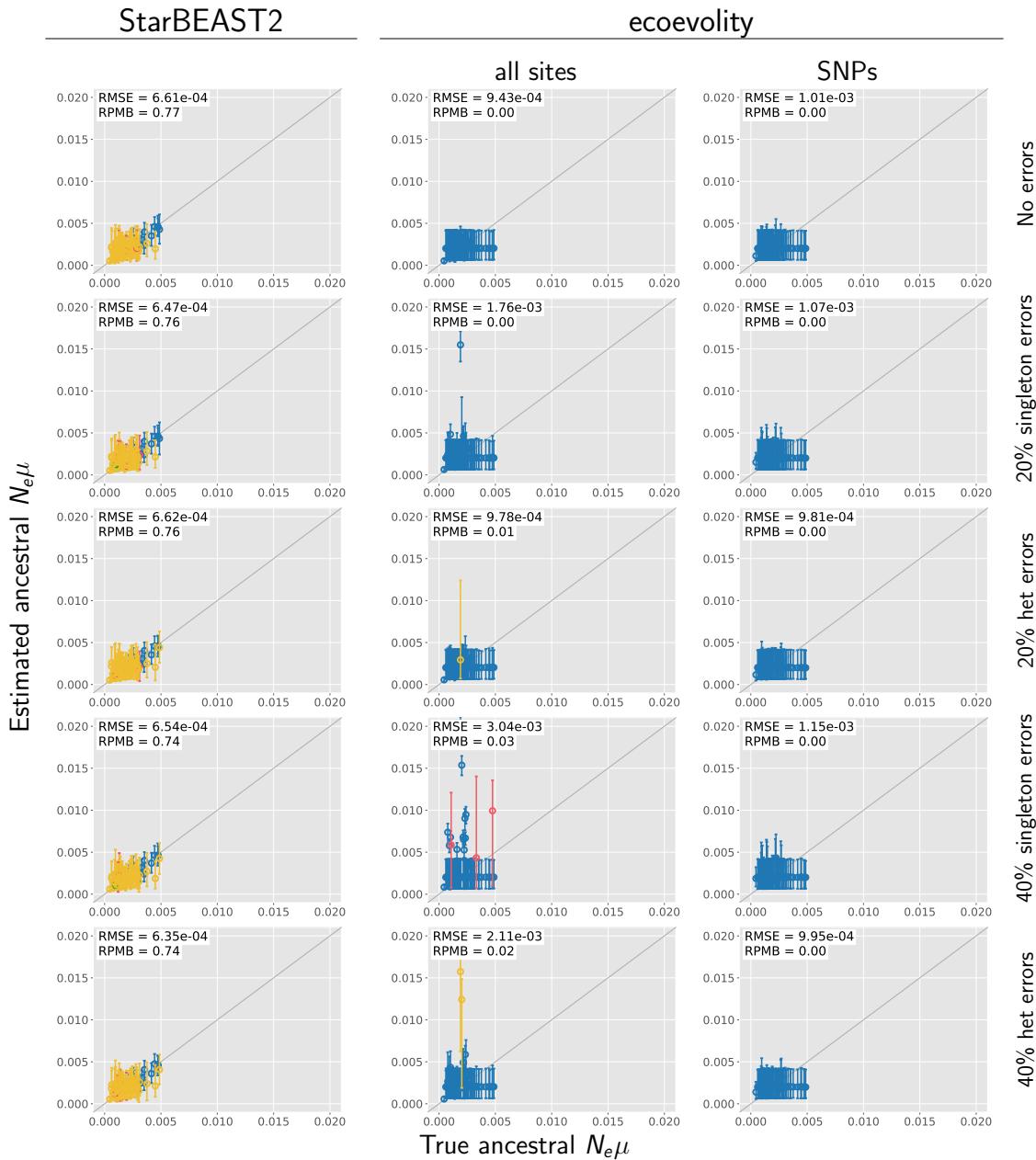


**Figure 4.4. Accuracy and precision of divergence-time estimates (in units of expected substitutions per site) with 250 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).



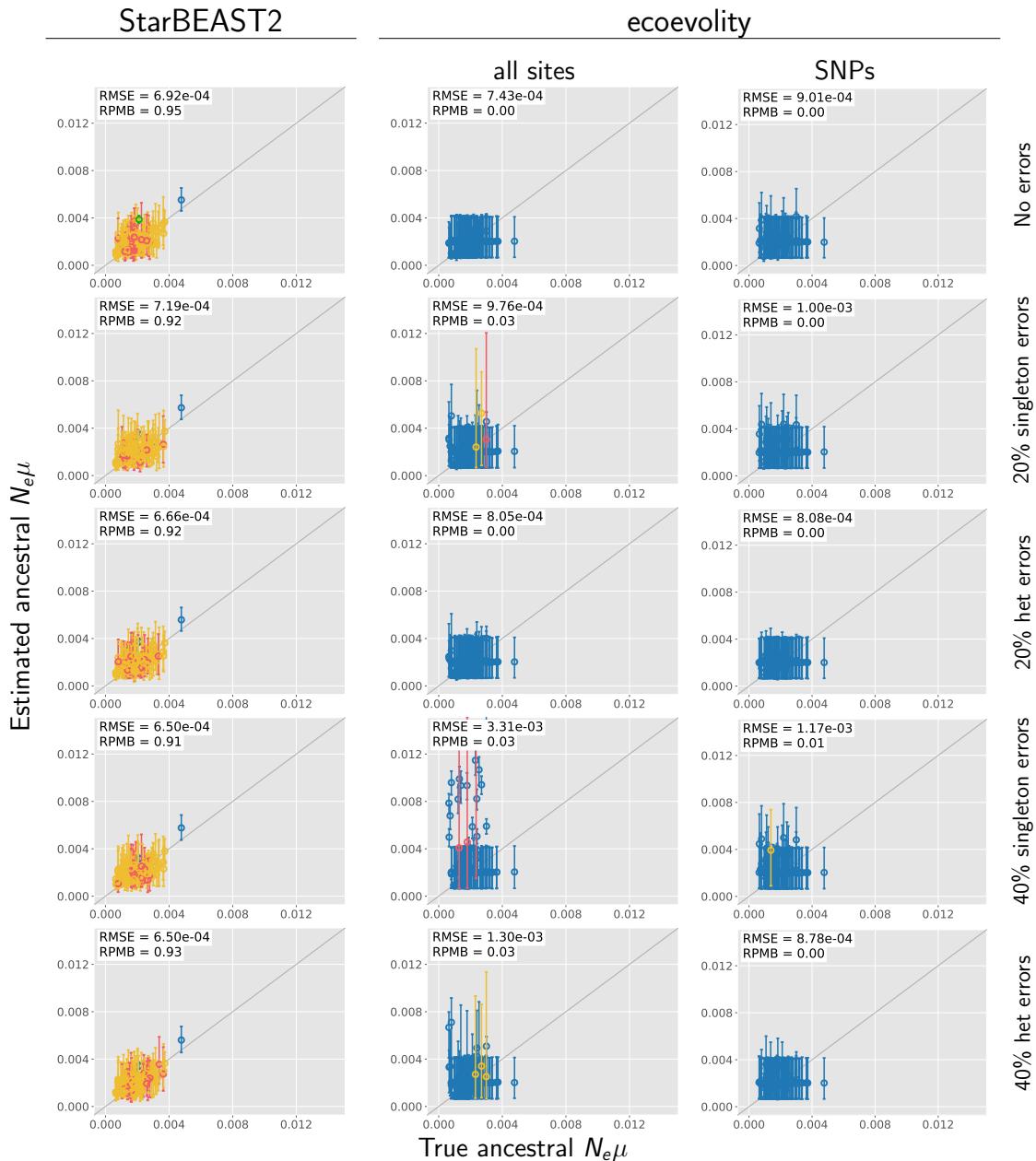
**Figure 4.5. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R \mu$ ) with 1000 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

### Ancestral $N_e\mu$ — 500bp loci

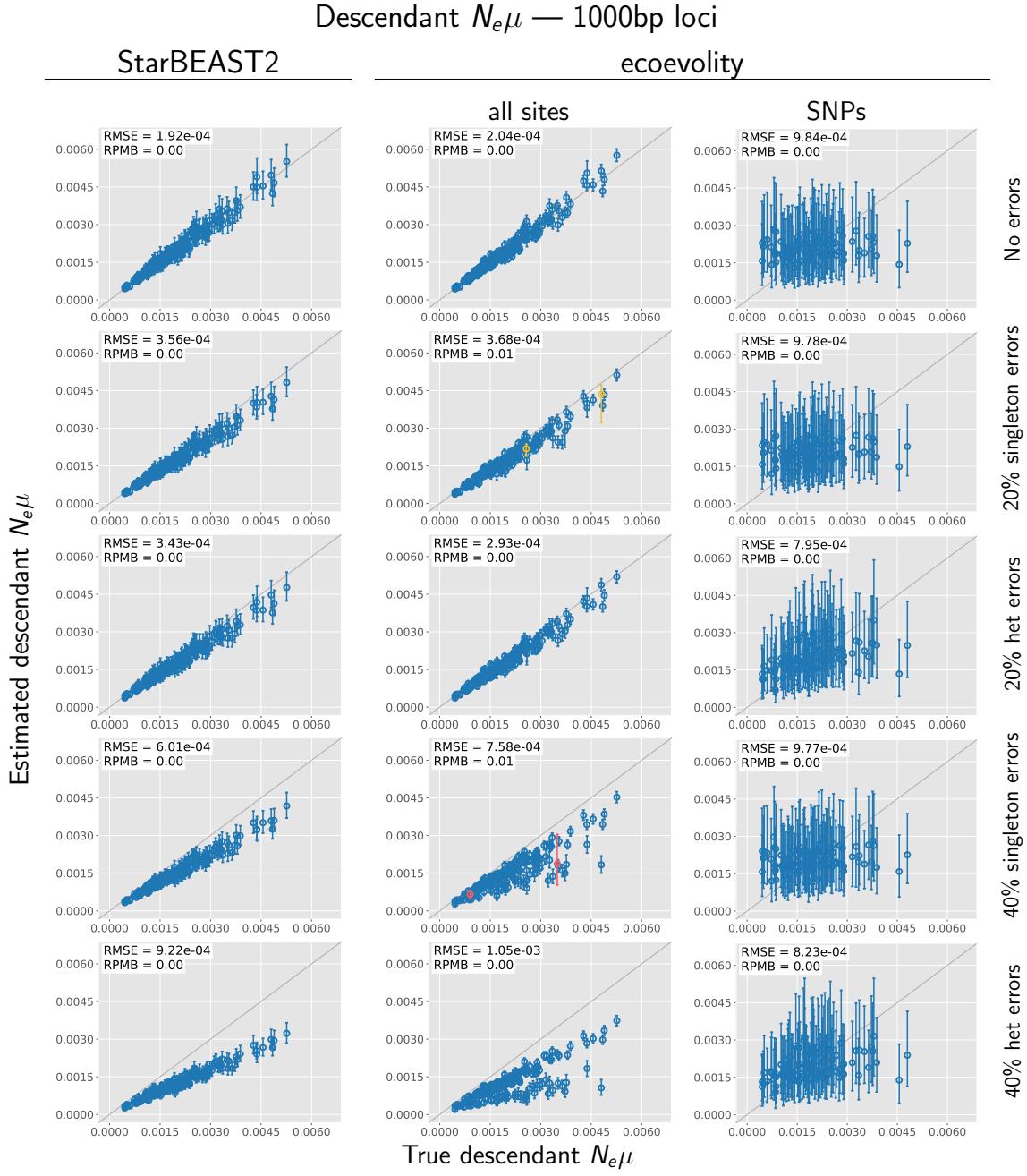


**Figure 4.6. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R\mu$ ) with 500 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

### Ancestral $N_e\mu$ — 250bp loci



**Figure 4.7. Accuracy and precision of estimates of root effective population size scaled by the mutation rate ( $N_e^R \mu$ ) with 250 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

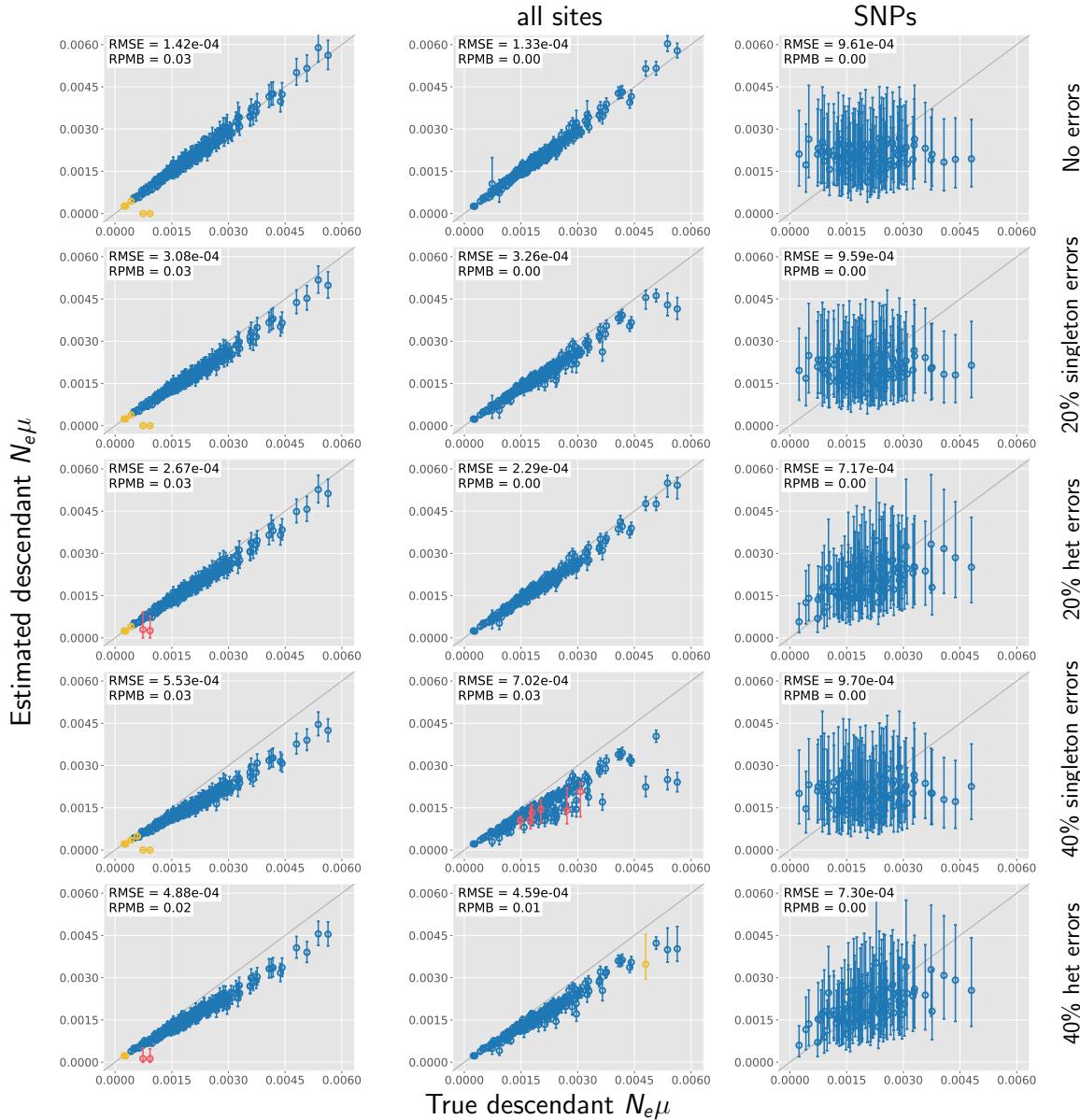


**Figure 4.8. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D\mu$ ) with 1000 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

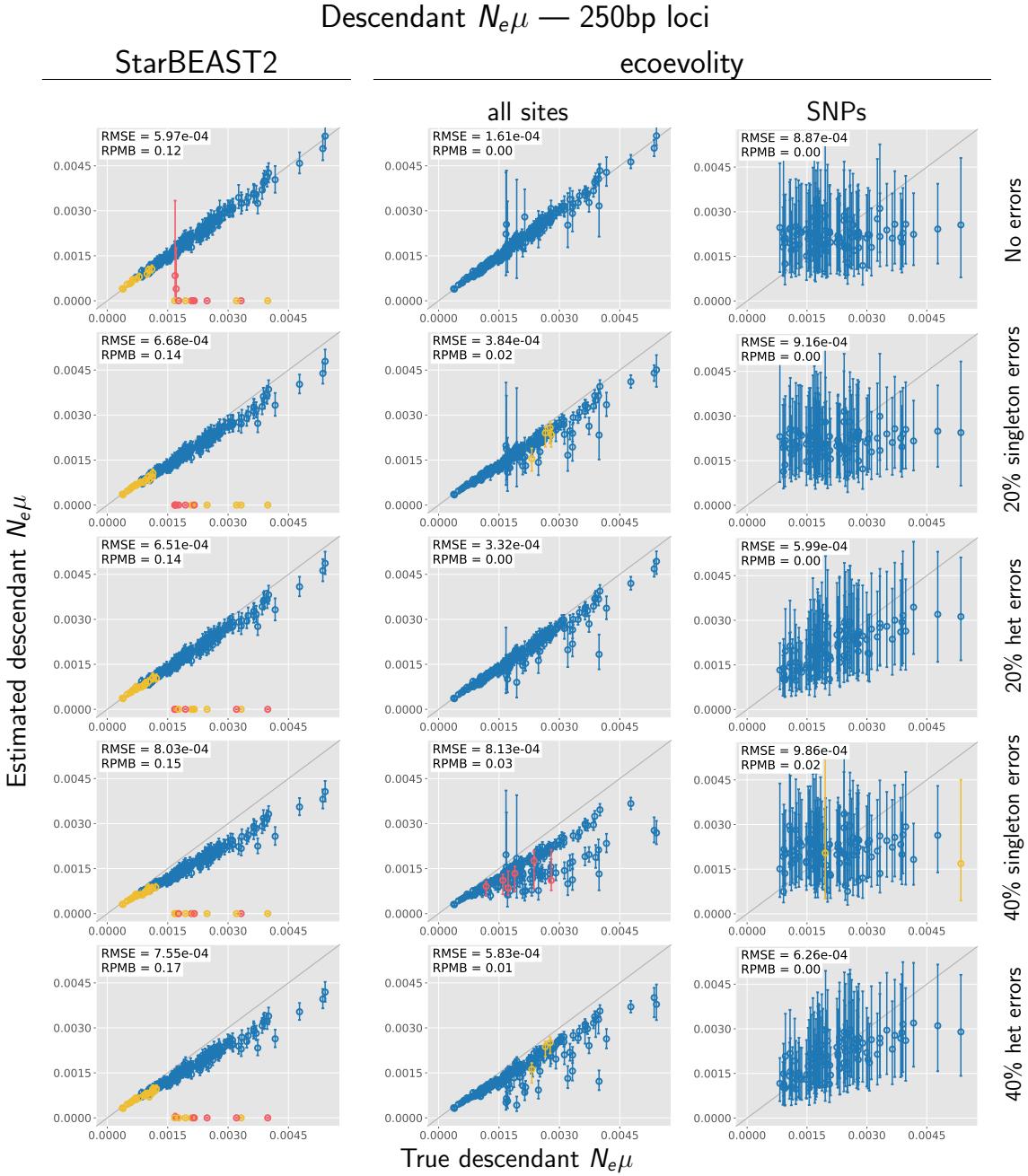
### Descendant $N_e\mu$ — 500bp loci

StarBEAST2

ecoevolity



**Figure 4.9. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D\mu$ ) with 500 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).



**Figure 4.10. Accuracy and precision of estimates of effective population sizes of the descendant branches of the tree scaled by the mutation rate ( $N_e^D\mu$ ) with 250 base pair loci.** The left column shows estimates from *StarBEAST2*, and the center and right column shows estimates from *ecoevolity* using all sites and (at most) one SNP per locus. The top row shows estimates from 200 data sets simulated without character-pattern errors. Rows labelled 20% and 40% singleton errors show estimates from the same alignments after singleton site patterns were changed to invariant sites with probabilities 0.2 and 0.4, respectively. Rows labelled 20% and 40% het errors show estimates from the same (error-free) alignments after we randomly paired gene copies within each species into two diploid genotypes, and for each genotype we randomly replaced one allele with the other with probability 0.2 and 0.4, respectively. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with ESS < 200 and/or PSRF > 1.2. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).

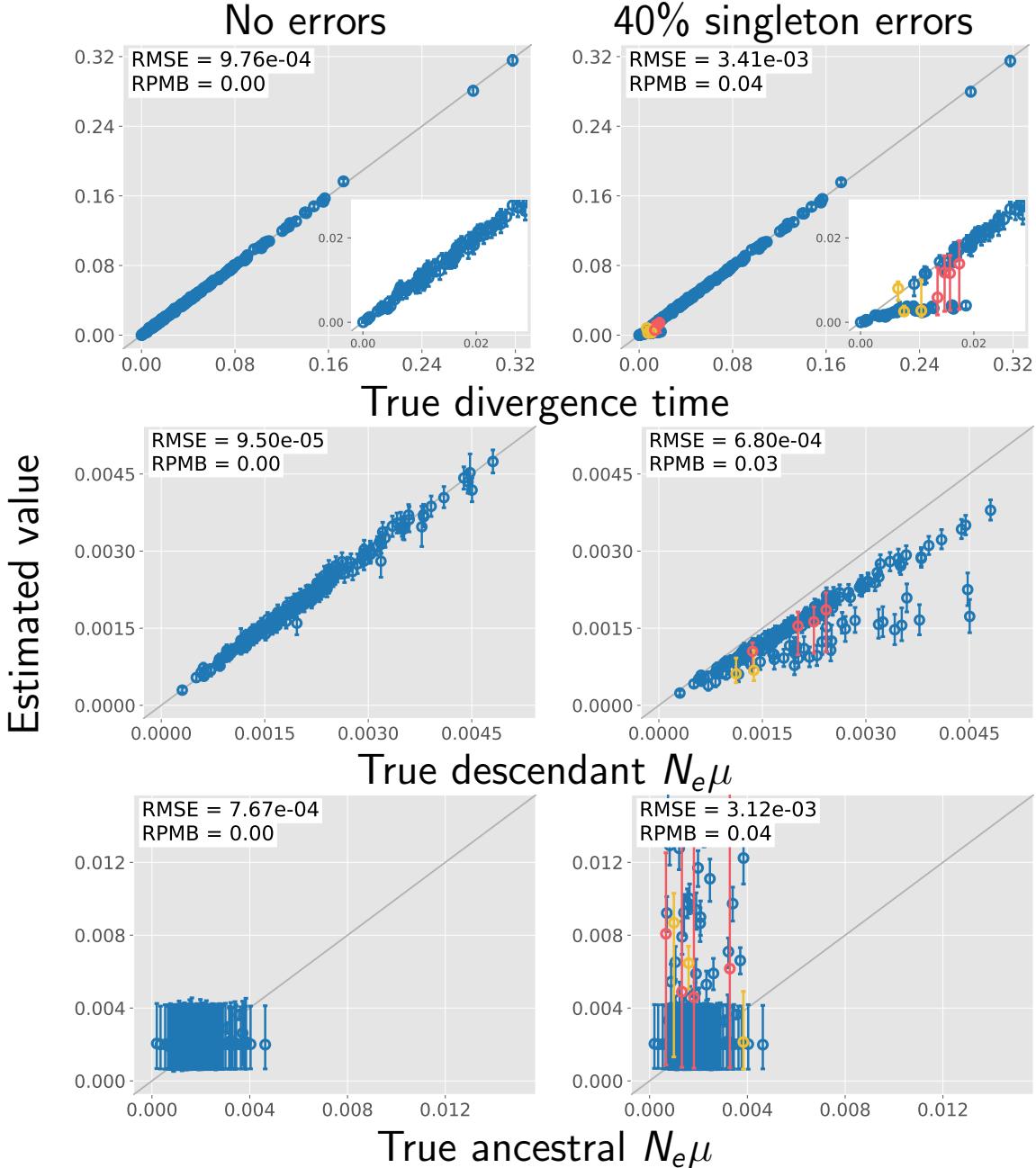


Figure 4.11. The performance of *ecoevolity* with data sets simulated with unlinked characters. Each plotted circle and associated error bars represent the posterior mean and 95% credible interval. Circles and error bars are colored yellow if the effective sample size (ESS) of the estimate was less than 200, red if the potential scale reduction factor (PSRF) was greater than 1.2, and green if both conditions were true. The root mean square error (RMSE) and rate of poor MCMC behavior (RPMB) is given for each plot, the latter of which is the proportion of estimates with  $\text{ESS} < 200$  and/or  $\text{PSRF} > 1.2$ . Inset plots magnify estimates of most recent divergence times. We generated the plots using matplotlib Version 3.1.1 (Hunter, 2007).