2  RH: PHASING IMPROVES UTILITY OF UCES

# Allele Phasing Greatly Improves the Phylogenetic Utility of Ultraconserved Elements

5  TOBIAS ANDERMANN[1,2], ALEXANDRE M. FERNANDES[3], URBAN OLSSON[1,2], MATS

6  TÖPEL[2,4], BERNARD PFEIL[1,2], BENGT OXELMAN[1,2], ALEXANDRE ALEIXO[5], BRANT C.

7  FAIRCLOTH[6] AND ALEXANDRE ANTONELLI[1,2,7,8]

8  [1]Department of Biological and Environmental Sciences, University of Gothenburg, SE-413 19,

9  Göteborg, Sweden;

10  [2]Gothenburg Global Biodiversity Centre, Box 461, SE-405 30, Göteborg, Sweden

11  [3]Universidade Federal Rural de Pernambuco, Serra Talhada, Brazil

12  [4]Department of Marine Sciences, University of Gothenburg, SE-413 19, Göteborg, Sweden;

13  [5]Museu Paraense Emílio Goeldi, Collection of Birds, Belém, Brazil

14  [6]Department of Biological Sciences and Museum of Natural Science, Louisiana State University,

15  Baton Rouge, LA, U.S.A.

16  [7]Gothenburg Botanical Garden, SE-413 19, Göteborg, Sweden

17  [8]Harvard University, Department of Organismic and Evolutionary Biology, Cambridge, MA,

18  U.S.A.

19  **Corresponding author:** Tobias Andermann, Department of Biological and

20  Environmental Sciences, University of Gothenburg, Carl Skottsbergs Gata 22B, SE-413 19,

21  Göteborg, Sweden; E-mail: tobias.andermann@bioenv.gu.se

*Abstract.*— Advances in high-throughput sequencing techniques now allow relatively easy and affordable sequencing of large portions of the genome, even for non-model organisms. Many phylogenetic studies reduce costs by focusing their sequencing efforts on a selected set of targeted loci, commonly enriched using sequence capture. The advantage of this approach is that it recovers a consistent set of loci, each with high sequencing depth, which leads to more confidence in the assembly of target sequences. High sequencing depth can also be used to identify phylogenetically informative allelic variation within sequenced individuals, but allele sequences are infrequently assembled in phylogenetic studies. Instead, many scientists perform their phylogenetic analyses using contig sequences which result from the *de novo* assembly of sequencing reads into contigs containing only canonical nucleobases, and this may reduce both statistical power and phylogenetic accuracy. Here, we develop an easy-to-use pipeline to recover allele sequences from sequence capture data, and we use simulated and empirical data to demonstrate the utility of integrating these allele sequences to analyses performed under the Multispecies Coalescent (MSC) model. Our empirical analyses of Ultraconserved Element (UCE) locus data collected from the South American hummingbird genus *Topaza* demonstrate that phased allele sequences carry sufficient phylogenetic information to infer the genetic structure, lineage divergence, and biogeographic history of a genus that diversified during the last three million years. The phylogenetic results support the recognition of two species, and suggest a high rate of gene flow across large distances of rainforest habitats but rare admixture across the Amazon River. Our simulations provide evidence that analyzing allele sequences leads to more accurate estimates of tree topology and divergence times than the more common approach of using contig sequences.

(Keywords: SNP, heterozygous sites, target enrichment, gene tree, species tree, mitochondrial genome, Trochilidae, Aves)

47        Massive Parallel Sequencing (MPS) techniques enable time- and cost-efficient

48 generation of DNA sequence data. Instead of using MPS to sequence complete genomes,

49 many researchers choose to focus their sequencing efforts on a set of target loci to lower

50 costs while achieving higher coverage and more reliable sequencing of these target regions

51 (Faircloth et al. 2012, 2013; Mirarab et al. 2014; Smith et al. 2014; Faircloth 2015; Harvey

52 et al. 2016; Meiklejohn et al. 2016). These multilocus datasets typically contain hundreds

53 or thousands of target loci, and most are generated through enrichment techniques such as

54 sequence capture (synonym: target enrichment, Gnirke et al. (2009)). After collecting

55 sequence data from these targeted loci, many researchers assemble their high coverage

56 sequence reads into "contigs" using *de novo* genome assembly software, and the "contig

57 sequence" output by these assemblers often ignore the variants at heterozygous positions

58 that are expected in diploid organisms. Typically, variable positions are treated as

59 sequencing errors and assembly algorithms output "contig sequences" containing the more

60 probable (i.e., numerous) variant while discarding the alternative (Iqbal et al. 2012). As a

61 result, the "contig sequences" that are produced contain only canonical nucleobases, losing

62 the information about read variability at variable positions. Hereafter, we use "contigs"

63 and "contig sequences" to refer to the sequences that are output by *de novo* assemblers.

64        One alternative approach to generating contig sequences uses the depth of

65 sequencing coverage to programatically identify variable positions within a targeted locus

66 (also known as "calling" single nucleotide polymorphisms (SNPs)) and subsequently

67 sorting (or "phasing") these SNPs into two allele sequences or "haplotypes" which

68 represent alleles on the same chromosome present at that locus. These approaches have

69 been used to estimate demographic parameters such as effective population size, rate of

70 migration, and the amount of gene flow between and within populations. However, it is

71 rarely acknowledged (*c.f.* Lischer et al. 2014; Potts et al. 2014; Schrempf et al. 2016;

72 Eriksson et al. 2017) that allelic sequences are useful for phylogenetic studies to improve

73 the estimation of gene trees, species trees, and divergence times (Garrick et al. 2010; Potts

74 et al. 2014; Lischer et al. 2014). The common practice of neglecting allelic information in

75 phylogenetic studies possibly results from historical inertia and a lack of computational

76 pipelines to prepare allele sequences for phylogenetic analysis using MPS data.

77 In addition to the problems of determining allelic sequences, the proper analysis of

78 allelic information in phylogenetic studies remains a challenging and intensively discussed

79 topic (Garrick et al. 2010; Lischer et al. 2014; Potts et al. 2014; Schrempf et al. 2016;

80 Leaché and Oaks 2017). Various approaches have been proposed to include this

81 information into phylogenetic methods (Lischer et al. 2014; Potts et al. 2014; Schrempf

82 et al. 2016). One is to code heterozygous sites using the International Union of Pure and

83 Applied Chemistry (IUPAC) ambiguity codes and to include these as additional characters

84 in existing substitution models for gene tree and species tree inference (Potts et al. 2014;

85 Schrempf et al. 2016). While these studies demonstrate that integrating additional allelic

86 information in this manner increases accuracy in phylogenetic inference, Lischer et al.

87 (2014) found that coding heterozygous sites as IUPAC ambiguity codes in phylogenetic

88 models biases the results toward older divergence time estimates. Instead, Lischer et al.

89 (2014) introduced a method of repeated random haplotype sampling (RRHS) in which

90 allele sequences are repeatedly concatenated across many loci, using a random haplotype

91 for any given locus in each replicate. In their approach, they then analyzed thousands of

92 concatenation replicates separately for phylogenetic tree estimation and summarized the

93 results between replicates, thereby integrating the allelic information in the form of

94 uncertainty intervals. However, there are two important shortcomings of this approach: 1.

95 concatenating unlinked loci (and in particular allele sequences from unlinked loci) in a

96 random manner is known to produce incorrect topologies (Degnan and Rosenberg 2009)

97 often with false confidence (Edwards et al. 2007; Kolaczkowski and Thornton 2004;

98 Kubatko and Degnan 2007; Mossel and Vigoda 2005), which is not accounted for when

99 doing so repeatedly and summarizing the resulting trees, and 2. running thousands of tree

100 estimation replicates based on extensive amounts of sequence data results in unfeasibly long

101 computation times, particularly for Markov-Chain Monte Carlo (MCMC) based softwares

102 such as MrBayes or BEAST. Hence, there is need to find proper solutions to include

103 heterozygous information in phylogenetic analyses, as concluded by Lischer et al. (2014).

104       Here, we introduce the bioinformatic assembly of allele sequences from UCE data

105 (Fig. 1) and demonstrate a full integration of allele sequences to species tree estimation

106 under the multispecies coalescent (MSC) model. In our approach, we treat each allelic

107 sequence of an individual at a given locus as an independent sample from the population,

108 and we analyze these sequences using the species tree and delimitation software STACEY

109 (Jones et al. 2014; Jones 2017), which allows for this approach by not requiring *a priori*

110 clade- or species-assignments. We first demonstrate the empirical utility of this approach

111 by resolving the shallow genetic structure ($<1$ Ma) within two recognized morphospecies of

112 the South American hummingbird genus *Topaza*, with a dataset of 2,386 ultraconserved

113 elements (UCEs, see Faircloth et al. (2012)). We then validate this approach, using

114 simulated data, and we find evidence that allele sequences yield more accurate results in

115 terms of species tree estimation and species delimitation than the contig sequence approach

116 that ignores heterozygous information. Further, our simulation results provide evidence

117 that compiling phased allele sequences and treating these as individual samples

118 outperforms alternative approaches of coding heterozygous information, such as analyzing

119 sequences containing IUPAC ambiguity codes or analyzing isolated SNPs. We conclude

120 that allele phasing for sequence capture data can be critical for correct species delimitation

121 and phylogeny estimation, particularly in recently diverged groups, and that analyses using

122 phased allele sequences should be considered as one, potential "best practice" for analyzing

123 sequence capture datasets in a phylogenetic context.

# Materials and Methods

[124]

## Study System

[125]

[126] The genus *Topaza* and its sister genus *Florisuga* form the Topazes group, which together
[127] with the Hermits represent the most ancient branch within the hummingbird family
[128] (Trochilidae) (McGuire et al. 2014). Topazes are estimated to have diverged as a separate
[129] lineage from all other hummingbirds around 21.5 Ma, whereas the most recent common
[130] ancestor (MRCA) of *Topaza* and *Florisuga* lived approximately 19 Ma (McGuire et al.
[131] 2014). At present, there are two morphospecies recognized within *Topaza*, namely the
[132] Fiery Topaz, *T. pyra* (Gould, 1846), and the Crimson Topaz, *T. pella* (Linnaeus, 1758).
[133] However, the species status of *T. pyra* has been challenged by some authors (Schuchmann
[134] 1999; Ornés-Schmitz and Schuchmann 2011), who consider this genus to be monotypic.
[135] Topaz hummingbirds are endemic to the Amazonian rainforest and are some of the most
[136] spectacular and largest hummingbirds worldwide, measuring up to 23 cm (adult males,
[137] including tail feathers) and weighing up to 12 g (Schuchmann et al. 2016; del Hoyo et al.
[138] 2016a). These birds are usually found in the forest canopy along forest edges and clearings,
[139] and are often seen close to river banks (Ornés-Schmitz and Schuchmann 2011). There is
[140] morphological evidence for several subspecies within both currently recognized *Topaza*
[141] species (Peters 1945; Schuchmann 1999; Hu et al. 2000; Ornés-Schmitz and Schuchmann
[142] 2011) that we investigate using genetic data.

## Sequence Data Generation

[143]

[144] We extracted DNA from the muscle tissue of 10 vouchered hummingbirds (9 *Topaza*, one
[145] *Florisuga*, see Table 1) using the Qiagen DNeasy Blood and Tissue Kit according to the

146 manufacturer's instructions (Qiagen GmbH, Hilden, Germany). These samples cover most

147 of the genus' total geographic range (Fig. 2) and all morphologically recognized

148 intraspecific taxa (Schuchmann et al. 2016; del Hoyo et al. 2016a). All samples were

149 sonicated with a Covaris S220 to a fragment length of 800 base pairs (bp). Paired-end,

150 size-selected (range 600-800bp) DNA libraries were prepared for sequencing, using the

151 magnetic-bead based NEXTflexTM Rapid DNA-Seq Kit (Bioo Scientific Corporation,

152 Austin, TX, USA), following the user's manual (v14.02).

153 We used the "Tetrapods-UCE-2.5Kv1" bait set (`uce-2.5k-probes.fasta`),

154 consisting of 2,560 baits (each 120 bp), targeting 2,386 UCEs, as described by Faircloth

155 et al. (2012). The bait sequences were downloaded from `http://ultraconserved.org` and

156 synthesized by MYcroarray (Biodiscovery LLC, Ann Arbor, MI, USA). Sequence

157 enrichment was performed using a MYbaits kit according to the enclosed user manual

158 (v1.3.8). The enriched libraries were then sequenced using 250 bp, paired-end sequencing

159 on an Illumina MiSeq machine (Illumina Inc., San Diego, CA, USA). Library preparation,

160 sequence enrichment and sequencing were performed by the Sahlgrenska Genomics Core

161 Facility in Gothenburg, Sweden.

## *Mitochondrial Genome*

163 To infer a dated mitochondrial phylogeny for the genus *Topaza* to compare with the

164 nuclear phylogeny, we used off-target mitochondrial reads to assemble the complete

165 mitochondrial genome for all samples. We found that as many as 4.5% of all sequence

166 reads were of mitochondrial origin, even though no baits targeting mitochondrial loci were

167 used during sequence capture. An alignment of the assembled mitochondrial genomes for

168 all samples was analyzed in BEAST (Drummond et al. 2012). Dating priors included

169 clock-rate priors for three mitochondrial genes, estimated for honeycreepers by Lerner et al.

170 (2011) and node-age priors within the genus *Topaza* that were estimated by McGuire et al.

171 (2014). The resulting phylogeny and estimated divergence times are shown in 2. A detailed

172 description of the assembly and phylogenetic analysis of the mitochondrial genome data

173 can be found in online Appendix 1 (Supplemental Material available on Dryad,

174 doi:10.5061/dryad.hq3vq).

## *UCE Data Processing*

176 For this study we generated five different types of datasets, which we analyzed under the

177 MSC. These five datasets represent different coding schemes for heterozygous information

178 and are listed and described in the following sections.

179 *1. UCE contig alignments.*— Because contig sequences are commonly used in phylogenetic

180 analyses of MPS datasets (e.g. Faircloth et al. (2012); Smith et al. (2014); Faircloth

181 (2015)), we generated multiple sequence alignments (MSAs) of contigs for all UCE loci in

182 order to test the accuracy of the phylogenetic estimation of this approach.

183 To create MSAs from UCE contig data, we followed the suggested workflow from

184 the PHYLUCE documentation

185 (`http://phyluce.readthedocs.io/en/latest/tutorial-one.html`). We applied the

186 PHYLUCE default settings unless otherwise stated. First we quality-filtered and cleaned

187 raw Illumina reads of adapter contamination with Trimmomatic (Bolger et al. 2014), which

188 is implemented in the PHYLUCE function `illumiprocessor`. The reads were then

189 assembled into contigs using the software ABYSS (Simpson et al. 2009) as implemented in

190 the PHYLUCE pipeline. In order to identify contigs representing UCE loci, all assembled

191 contigs were mapped against the UCE reference sequences from the bait sequence file

192 (`uce-2.5k-probes.fasta`), using the PHYLUCE function `match_contigs_to_probes.py`.

193 We extracted only those sequences that matched UCE loci and that were present in all

194 samples (n=820). These UCE sequences were then aligned for each locus (Fig. 1) using

195 MAFFT (Katoh et al. 2009).

196 *2. UCE allele alignments.*— We altered the typical UCE workflow in order to retrieve the

197 allelic information that is lost when collapsing multiple reads into a single contig sequence

198 (Fig. 1). To create this new workflow, we extracted all UCE contigs for each sample

199 separately and treated each resulting contig set as a sample-specific reference library for

200 read mapping (reference-based assembly). We then mapped the cleaned reads against each

201 reference library on a per sample basis, using CLC-mapper from the CLC Workbench

202 software. The mapped reads were sorted and then phased with SAMtools v0.1.19 (Li et al.

203 2009), using the commands `samtools sort` and `samtools phase`, respectively. This

204 phasing function is based on a dynamic programming algorithm that uses read connectivity

205 across multiple variable sites to determine the two phases of any given diploid locus (He

206 et al. 2010). Further, this algorithm uses paired-end read information to reach connectivity

207 over longer distances and it minimizes the problem of accidentally phasing a sequencing

208 error, by applying the minimum error correction function (He et al. 2010).

209 UCE data provide an excellent dataset for allele phasing based on read connectivity,

210 because the read coverage across any given UCE locus typically is highest in the center and

211 decreases toward the ends. This makes it possible to phase throughout the complete locus

212 without any breaks in the sequence. Even in cases where the only variable sites are found

213 on opposite ends of the locus, the insert size we targeted in this study (800 bp), in

214 combination with paired-end sequencing, enabled the phasing process to bridge the

215 complete locus (average length of compiled UCE-sequences in our study was 870 bp).

216 The two phased output files (BAM format) were inspected for proper variant

217 separation for all loci using Tablet (Milne et al. 2013). We then collapsed each allele BAM

218 file into a single consensus sequence per haplotype and exported the two resulting allele

219 sequences for each sample in FASTA format. In order to separate true heterozygous sites

220 from occasional variants introduced by sequencing errors, we only made a nucleotide call if

221 the respective nucleotide was supported by at least three reads. Ambiguous positions were

<sub>222</sub> coded with the IUPAC code 'N' in the allele consensus sequences. We explored the

<sub>223</sub> difference in the treatment of heterozygous positions between the contigs produced by the

<sub>224</sub> *de novo* assembler ABYSS and our phased allele sequences in detail (exemplary for one

<sub>225</sub> sample) in online Appendix 2 (Supplemental Material).

<sub>226</sub>      In the next, step we aligned the allele sequences between all samples, separately for

<sub>227</sub> each UCE locus, using MAFFT (Fig. 1). We integrated this complete workflow into the

<sub>228</sub> UCE processing software PHYLUCE (Faircloth 2015) with slight alterations, one of which

<sub>229</sub> is the use of the open-source mapping program bwa (Li and Durbin 2010) in place of

<sub>230</sub> CLC-mapper.

<sub>231</sub> *3. UCE IUPAC consensus sequence alignments.—* We generated an additional set of

<sub>232</sub> alignments by merging the two allele sequences for each individual into one consensus

<sub>233</sub> sequence with heterozygous sites coded as IUPAC ambiguity codes

<sub>234</sub> (`merge_allele_sequences_ambiguity_codes.py`, available from:

<sub>235</sub> github.com/tobiashofmann88/UCE-data-management/). We used this dataset to test

<sub>236</sub> whether our allele phasing approach improved phylogenetic inference when compared to

<sub>237</sub> the IUPAC consensus approach applied in other studies, where heterozygous positions are

<sub>238</sub> coded as IUPAC ambiguity codes in a consensus sequence for each locus and individual

<sub>239</sub> (Potts et al. 2014; Schrempf et al. 2016).

<sub>240</sub> *4. UCE chimeric allele alignments.—* To investigate whether correct phasing of

<sub>241</sub> heterozygous sites is essential or if similar results are achieved by randomly placing

<sub>242</sub> variants in either allele sequence, we generated a dataset with chimeric allele sequence

<sub>243</sub> alignments. We created these alignments by applying a custom python script

<sub>244</sub> (`shuffle_snps_in_allele_alignments.py`, available from:

<sub>245</sub> github.com/tobiashofmann88/UCE-data-management/) to the phased allele sequence

<sub>246</sub> alignments and randomly shuffling the two variants at each polymorphic position between

247 the two allele sequences for each individual. This process leads, in many cases, to an

248 incorrect combination of variants on each allele sequence, thereby creating chimeric allele

249 sequences. The resulting alignments contain the same number of sequences as the phased

250 allele alignments (two sequences per individual), whereas the contig alignments and the

251 IUPAC consensus alignments contain only half as many sequences (one sequence per

252 individual).

253 *5. UCE SNP alignment.*— A common approach to analyzing heterozygous information is

254 to reduce the sequence information to only a single variant SNP per locus. This

255 data-reduction approach is often chosen because multilocus datasets of the size generated

256 in this study can be incompatible with Bayesian MSC methods applied to the full sequence

257 data, due to extremely long computational times and convergence issues. Instead,

258 alignments of unlinked SNPs can be used to infer species trees and species demographics

259 under the MSC model with the BEAST2 package SNAPP (Bryant et al. 2012), a program

260 specifically designed for such data. However, extracting and filtering SNPs from BAM files

261 with existing software (such as the Genome Analysis Toolkit (GATK), McKenna et al.

262 (2010)) and converting these into a SNAPP compatible format can be cumbersome,

263 because SNAPP requires positions with exactly two different states, coded in the following

264 manner: individual homozygous for the original state = "0", heterozygous = "1", and

265 homozygous for the derived state = "2".

266 To alleviate this problem, we developed a python function that extracts biallelic

267 SNPs directly from allele sequence MSAs (`snps_from_uce_alignments.py`, available from:

268 github.com/tobiashofmann88/snp_extraction_from_alignments/). Extracting SNPs from

269 MSAs in this manner is a straightforward and simple way to generate a SNP dataset

270 compatible with SNAPP, and does not require re-visiting the BAM files. A similar

271 program is also available in the R-package `phrynomics` (Leaché et al. 2015). We used this

272 approach to extract one variable position per alignment (to ensure unlinked SNPs) that

273 had exactly two different states among all *Topaza* samples, not allowing for positions with

274 missing data or ambiguities. This produced a SNP dataset of 598 unlinked SNPs.

## *Generation of Simulated UCE Data*

276 To assess the accuracy of the phylogenetic inferences resulting from different data

277 processing approaches, we simulated UCE data similar to those discussed in the five

278 processing schemes we applied to the empirical *Topaza* data. However, because this

279 approach required us to simulate allele alignments before generating contig alignments,

280 steps one and two, below, are reversed from their order, above. We repeated all steps

281 involving the generation and analyses of simulated data to produce 10 independent

282 simulation replicates.

283 *1. Simulated allele alignments.—* In order to simulate allele alignments similar to our

284 empirical data we first estimated species divergence times and population sizes from the

285 empirical UCE allele MSAs under the MSC model (Rannala and Yang 2003) using the

286 Bayesian MCMC program BPP v3.1 (Yang 2015). We applied the A00 model, which

287 estimates divergence times and population sizes from MSAs for a given species tree

288 topology. As input topology we used the species tree topology resulting from the analysis of

289 the empirical allele MSAs in STACEY, assigning the *Topaza* samples to five separate taxa

290 (corresponding to colored clades in Figure 3b). An initial BPP analysis did not converge in

291 reasonable computational time, a problem that has previously been reported for UCE

292 datasets containing several hundred loci (Giarla and Esselstyn 2015). To avoid this issue,

293 we split the 820 UCE alignments randomly into 10 subsets of equal size (n=82) and

294 analyzed these separately with identical settings in BPP. The MCMC was set for 150,000

295 generations (burn-in 50,000), sampling every 10 generations. We summarized the estimates

296 for population sizes and divergence times across all 10 individual runs. We then applied the

297 mean values of these estimates to the species tree topology, by using the estimated

298 divergence times as branch lengths and estimated population sizes as node values, resulting

299 in the species tree in Figure 4g. This tree was used to simulate sequence alignments with

300 the MCcoal simulator, which is integrated into BPP. Equivalent to the empirical data, we

301 simulated sequence data for five taxa (D, E, X, Y, and Z) and one outgroup taxon (F, not

302 shown in Figure 4g). In the simulations, these taxa were simulated as true species under

303 the MSC model. In order to mimic the empirical allele data, we simulated four individuals

304 for species 'D' (equivalent to two allele sequences for 2 samples), four for species 'E', four

305 for species 'X', two for species 'Y' (two allele sequences for one sample), four for species 'Z',

306 and two for the outgroup species 'F'. In this manner we simulated 820 UCE allele MSAs of

307 848 bp length (a value equal to the average alignment length of the empirical allele

308 alignments). The resulting simulated allele MSAs are equivalent to our empirical allele

309 MSAs, containing two phased allele sequences for every individual that differ only in true

310 heterozygous sites and which are not expected to contain read-errors.

311 *2. Simulated contig alignments.*— To simulate UCE contig MSAs that contain sequences

312 similar to contigs generated by assemblers like ABYSS, Velvet or Trinity, which pick only

313 one of the two variants at a heterozygous site, we merged the sequences within each

314 coalescent species in pairs of two (equivalent to pairs of allele sequences). Each pair of

315 allele sequences was joined into one contig sequence by randomly picking one of the two

316 variants at each heterozygous site across all loci. As in the empirical contig assembly

317 approach, our simulation approach may generate chimeric contig sequences.

318 *3. Simulated IUPAC consensus alignments.*— Next, we generated IUPAC consensus MSAs

319 in the same manner as we generated the simulated contig MSAs in the previous step, with

320 the exception that all heterozygous sites were coded with IUPAC ambiguity codes instead

321 of randomly picking one of the two variants.

322 *4. Simulated chimeric allele alignments.*— We generated chimeric allele sequence MSAs

323 from the simulated allele MSAs by randomly shuffling the heterozygous sites between each

324 pair of sequences using the same pairs as in the previous two steps.

325 *5. Simulated SNP alignment.*— Finally, we extracted two different SNP datasets from the

326 simulated phased allele MSAs. The first SNP dataset (SNPs complete) was extracted in

327 the same manner as described for the empirical data (one SNP per locus for all loci) which

328 resulted in a total alignment length of 820 SNPs for the simulated data. We extracted an

329 additional SNP dataset (SNPs reduced) from only the subset of the 150 simulated allele

330 alignments that were used for the sequence-based MSC analyses (see next section below).

331 The resulting dataset of 150 SNPs was used to compare the phylogenetic inference based

332 on SNP data versus that based on full sequence data, if the same number of loci is being

333 analyzed. This enabled us to evaluate the direct effect of reducing the full sequence

334 information in the MSAs to one single SNP for each of the selected 150 loci.


335                    *MSC Analyses of Empirical and Simulated UCE Data*


336 *Sequence-based tree estimation.*— To jointly infer gene trees and species trees, we analyzed

337 each of the generated sets of MSAs (processing schemes 1-4 for empirical and simulated)

338 under the MSC model, using the DISSECT method (Jones et al. 2014) implemented in

339 STACEY (Jones 2017), which is available as a BEAST2 (Bouckaert et al. 2014) package.

340 STACEY allows *BEAST analyses without prior taxonomic assignments, searching the tree

341 space while simultaneously collapsing very shallow clades in the species tree (controlled by

342 the parameter collapseHeight). This collapsing avoids a common violation of the MSC

343 model that occurs when samples belonging to the same coalescent species are assigned to

344 separate taxa in *BEAST. This feature makes STACEY suitable for analyzing allele

345 sequences, because they do not have to be constrained to belong to the same taxon and can

346 be treated as independent samples from a population. STACEY runs with the usual

347 *BEAST operators, but integrates out the population size parameter and has new MCMC

348 proposal distributions to more efficiently sample the species tree, which decreases the time

349 until convergence. In order to reach even faster convergence, we reduced the number of loci

350 for this analysis by selecting the 150 allele MSAs with the most parsimony informative

351 sites. This selection was made for both the empirical and the simulated allele MSAs. The

352 same 150 loci were selected for all other processing schemes.

353    Prior to analysis, we estimated the most appropriate substitution model for each of

354 the 150 loci with jModeltest (Supplementary Table S1) using BIC. We used BEAUTI

355 v2.4.4 to create an input file for STACEY in which we unlinked substitution models, clock

356 models and gene trees for all loci. We did not apply any taxon assignments, thereby

357 treating every sequence as a separate taxon. We chose a strict clock for all loci and fixed

358 the average clock rate for one random locus to 1.0, while estimating all other clock rates in

359 relation to this locus. To ensure that all resulting species trees were scaled to an average

360 clock rate of 1.0, we rescaled every species tree from the posterior distribution (post

361 analysis) using the average clock rate of the respective MCMC step. We applied the

362 STACEY-specific BirthDeathCollapse model as a species tree prior, choosing a value of

363 1e-5 for the collapseHeight parameter. Other settings were: bdcGrowthRate = log normal

364 (M=4.6, S=1.5); collapseWeight = beta (alpha=2, beta=2); popPriorScale = log normal

365 (M=-7, S=2); relativeDeathRate = beta (alpha=1.0, beta=1.0). For the IUPAC consensus

366 data, we enabled the processing of ambiguous sites by adding `useAmbiguities="true"` to

367 the gene tree likelihood priors for all loci in the STACEY XML file. All analyses were run

368 for 1,000,000,000 MCMC generations or until convergence (ESS values >200), logging every

369 20,000 generations. Convergence was assessed using Tracer v1.6 (Rambaut et al. 2013). We

370 then summarized the posterior tree distribution into one Maximum Clade Credibility tree

371 (i.e. tree in the posterior sample that has the maximum product of posterior clade

372 probabilities) with TreeAnnotator v2.4.4, discarding the first 10% of trees as burn-in.

373      For the simulated data, we analyzed the posterior species tree distributions of each

374 analysis with the program SpeciesDelimitationAnalyser (part of the STACEY

375 distribution). This program produces a similarity matrix that contains the posterior

376 probabilities of belonging to the same cluster for each pair of sequences. This analysis was

377 run with a collapseHeight value of 1e-5 (identical to the collapseHeight used in the

378 STACEY analysis), while discarding the first 10% of trees as burn-in.

379 *SNP-based tree estimation.—* To estimate the species tree phylogeny from the extracted

380 SNP data, we analyzed the empirical and simulated SNP data in SNAPP. We did not

381 apply prior clade assignments to the samples in the SNP alignment (each sample was

382 assigned as its own taxon). We set coalescent rate and mutation rates to be estimated

383 based on the input data, and we chose a Yule species tree model with default settings ($\lambda =$

384 0.00765). We ran the analysis for 10,000,000 generations, sampling trees and other

385 parameters from the posterior every 1,000 generations. Unlike STACEY, SNAPP assumes

386 correct assignments of all sequences to coalescent species. Using the simulated SNP data,

387 we therefore tested how our approach of assigning every individual as its own coalescent

388 species affects the resulting phylogenetic inference. We did so by running a separate

389 analysis for both simulated SNP datasets (complete and reduced) with correct species

390 assignments (assignments as in Figure 4g).

## *Additional Analyses*

392 We ran additional analyses of the contig and the phased allele MSAs for both the empirical

393 and simulated data using a summary coalescent approach as implemented in MP-EST (Yu

394  et al. 2007), which can be found in online Appendix 3 (Supplemental Material) and

395  Supplementary Figures S1-S3.

# Results

*UCE Summary Statistics*

*Alignment statistics.*— In the following we use the term "polymorphic sites" for those positions within a MSA alignment of a given locus where we find at least two different states at a particular position among the sequences for all samples. This does not require a particular individual being heterozygous for the given position, since we do not search for SNPs on a per sample basis but rather for SNPs within the genus *Topaza*. In this manner, we found that the empirical UCE contig sequence alignments had an average of 2.8 polymorphic sites per locus and an average alignment length of 870 bp. In contrast, phasing the empirical UCE data to create allele alignments led to 4.5 polymorphic sites per locus and an average alignment length of 848 bp, representing a 60% increase in polymorphic sites per locus. This increase of polymorphic sites was attributable to the fact that many variants get lost during contig assembly, because ABYSS and other tested contig assemblers, namely Trinity and Velvet, often eliminate one of the two variants at heterozygous positions (see below). The reduced length of the allele alignments in comparison to the contig alignments was due to conservative alignment clipping thresholds implemented in PHYLUCE, which clips alignment ends if less than 50% of sequences are present. Because the allele phasing algorithm divides the FASTQ reads into two allele bins and because a nucleotide is only called if it is supported by at least three high-quality FASTQ reads, we lost some of the nucleotide calls at areas of low read coverage (mostly at

416 the ends of a locus) when comparing the allele sequences to the contig sequences. More

417 information about the distribution of lengths and variable sites within the empirical UCE

418 data can be found in the Supplementary Figures S4 and S5. The simulated contig MSAs

419 had an average of 3.2 polymorphic sites per locus, after excluding the outgroup (average

420 calculated across all 10 simulation replicates). The simulated allele MSAs, on the other

421 hand, contained an average of 5.4 polymorphic sites (69% increase) across 10 independent

422 simulation replicates. An overview of parsimony informative sites, variable sites and length

423 of each alignment (simulated and empirical data) can be found in Supplementary Table S2.

## *MSC Results of Empirical UCE Data*

425 The MSC species tree results for all tested processing schemes of the empirical UCE data

426 (contig sequences, allele sequences, IUPAC consensus sequences, chimeric allele sequences

427 and SNPs) strongly support the monophyly of both *T. pyra* and *T. pella* with 100%

428 Bayesian posterior probability (PP) (Fig. 3 and Supplementary Fig. S6). In all MSC

429 analyses, we also see strongly supported genetic structure within *T. pella* ($\geq$ 97% PP),

430 separating the northern samples (5 and 6, sampled north of the Amazon River) from the

431 southern ones (7, 8 and 9, sampled south of the Amazon River). Additionally, within the

432 shallow southern *T. pella* clade, all datasets, with exception of the IUPAC consensus data

433 (Fig. 3c), strongly support a genetic distinction ($\geq$ 99% PP) between sample 7 from the

434 Amazon River delta and the other southern *T. pella* samples (8 and 9). Further, the

435 analysis of the phased allele MSAs returns a phylogenetic signal, possibly also tracking a

436 genetic divergence between a northern and a southern clade within *T. pyra*, but their

437 monophyly is not very strongly supported (Fig. 3b). This pattern is further supported by

438 the mitochondrial phylogeny, which shows the same divergence within *T. pyra*, dated at

439 0.68 million years ago (Fig. 2 and online Appendix 1).

## MSC Results of Simulated Data

441  *Species tree topology.—* We analyzed six different datasets under the MSC model for each

442  of the ten simulation replicates: contig sequence MSAs (n=150, STACEY), allele sequence

443  MSAs (n=150, STACEY), IUPAC consensus MSAs (n=150, STACEY), chimeric allele

444  MSAs (n=150, STACEY), reduced SNP data (n=150, SNAPP), and the complete SNP

445  dataset (n=820, SNAPP). All resulting species trees (Fig. 4a-f) correctly return the

446  topology of the species tree that was used to simulate the data (Fig. 4g) across all ten

447  simulation replicates (Supplementary Fig. S7). All central nodes in the species trees are

448  supported by $\geq 90\%$ PP in all analyses, with the exception of the species tree resulting

449  from the reduced SNP dataset, which shows very weak support for two nodes and has a

450  large uncertainty interval around the root-height (Fig. 4e). However, these shortcomings

451  disappeared when we added more (unlinked) SNPs to the dataset (Fig. 4f). The full SNP

452  dataset (n=820) produced the correct species tree topology with high node support

453  consistently throughout all ten independently simulated datasets (Supplementary Fig. S8).

454  The SNAPP species tree topology appeared to be unaffected by the chosen clade

455  assignment model; while we allowed every sequence to be its own taxon in Figure 4e and f,

456  we also applied the correct species assignment (as in Fig. 4g) in two additional analyses for

457  one of the simulation replicates (reduced and complete SNP data) that returned the same

458  tree topology (Supplementary Figs. S9 and S10).

459  *Species delimitation.—* Although the inferred species tree topology was consistent among

460  all four sequence-based MSC analyses (Fig. 4a-d), the inferred node heights varied

461  considerably between the species trees resulting from the different data processing schemes.

462  For the contig sequence data (Fig. 4a) and the chimeric allele data (Fig. 4d), the node

463  heights within the five simulated species (D,E,X,Y,Z) were too high, which led to an

464  overestimation of the number of coalescent species in the dataset (see similarity matrices).

465  Conversely, the phased allele data (Fig. 4b) and the IUPAC consensus data (Fig. 4c)

466  correctly delimited the five coalescent species from the simulation input tree (Fig. 4g). The

467  STACEY results showed the same pattern in all ten simulation replicates (Fig. S7).

468  *Accuracy of divergence time estimation.—* For all four sequence-based analyses (Fig. 4a-d)

469  the average substitution rate across all loci was set to '1'. Under these settings, we

470  expected the absolute values of the sequence-based analyses to return the node height

471  values of the simulation input tree, which used substitution rates scaled in the same

472  manner. The phased allele MSAs produced the most accurate estimation of divergence

473  times out of all tested datasets (see proximity of estimates to simulation input value,

474  represented by green line in Figure 5). This was the case for all nodes in the species tree,

475  namely (D,E), (Y,Z), (X,(Y,Z)), and ((D,E)(X,(Y,Z))). The divergence time estimates

476  resulting from the phased allele data accurately recovered the true values and did not show

477  any bias throughout ten simulation replicates (Supplementary Fig. S11). This contrasts

478  with the contig MSAs and the chimeric allele MSAs that consistently overestimated the

479  height of all nodes and the IUPAC consensus MSAs which consistently underestimated the

480  height of all nodes (Figs. 5 and S11).

# DISCUSSION

481

## *Phased Allele Sequences Return The Most Accurate Phylogeny*

482

483       We tested whether phylogenetic inference improves by phasing sequence capture

484  data into allele sequences, in comparison to the standard workflow of analyzing contig

485  sequences (Faircloth et al. 2012; McCormack et al. 2012; Smith et al. 2014; Faircloth 2015).

486  The answer is yes. We find that phased allele data outperform contig sequences in terms of

species delimitation (Fig. 4) and divergence time estimation (Fig. 5). Contig sequence MSAs on the other hand lead to a consistent overestimation of divergence times (Fig. 5), which in turn lead to an overestimation of the number of coalescent species in our simulated data (Fig. 4a). These results support earlier work by Lischer et al. (2014), who concluded that consensus sequences introduce a bias towards older node heights. Because both our empirical and simulated data represent rather shallow phylogenetic relationships, future research is required to determine if these findings also apply to datasets representing divergence events occurring in deeper time.

Besides these practical advantages of using phased allele sequences for phylogenetic analyses, there are further theoretical arguments for compiling and analyzing allele sequence MSAs from sequence capture datasets.

First, allele sequences represent the smallest evolutionary unit on which selection and other evolutionary processes act. Therefore, the coalescent models that underlie our phylogenetic methods, including the MSC model Degnan and Rosenberg (2009), have been developed for allele sequences. Contig sequences, on the other hand, represent an artificial and possibly chimeric sequence construct that arises from merging all read variation at a given locus into a single sequence. This process masks information by eliminating one of the two variants at a heterozygous site (online Appendix 2). This shortcoming of the most common assemblers (e.g. ABYSS, Trinity and Velvet) is due to the fact that they were designed to assemble sequences of haploid genomes and they are not optimized for heterozygous sequences or genomes (Bodily et al. 2015).

Second, not only are allele sequences the more appropriate data type, but phasing sequence capture data also leads to a doubling of the effective sample size, since two sequences are compiled for a diploid individual, in contrast to the single sequence per individual that is recovered when taking the contig approach. Here, we demonstrate how these sequences can be properly applied as independent samples from a population by

513 using the assignment-free BirthDeathCollapse model as implemented in STACEY. Because

514 STACEY requires no *a priori* assignment of sequences to taxa, it avoids a violation of the

515 MSC that would occur when analyzing allele sequences as separate taxa in *BEAST, since

516 *BEAST assumes each taxon constitutes a separate coalescent species.

517 Third, sequence capture datasets such as UCEs are optimal for allele phasing

518 because they contain high read coverage collected across short genomic intervals that are

519 optimal for read-connectivity based phasing. The workflow developed in this study is now

520 fully integrated into the PHYLUCE pipeline, making allele phasing for sequence capture

521 data easily available to a broad user group.

## *Phasing of Heterozygous Sites Matters*

523 Several studies have accounted for heterozygosity by inserting IUPAC ambiguity codes into

524 their sequences at variable positions (Potts et al. 2014; Schrempf et al. 2016), rather than

525 phasing SNPs to produce separate allele sequences. Here, we directly compared these two

526 approaches, and found that the IUPAC consensus sequences performed equally well to the

527 phased allele sequences for estimating the species tree topology (Fig. 4). However, IUPAC

528 consensus sequence data led to a consistent underestimation of the divergence times of all

529 nodes in the species tree (Fig. 5). Our results contrast with those of (Lischer et al. 2014),

530 who reported an overestimation of divergence times for alignments containing IUPAC

531 ambiguity codes. The differences between our results may simply be caused by the different

532 tree inference programs used. Lischer et al. (2014) applied a Neighbour Joining (NJ) tree

533 algorithm as implemented in the software PHYLIP (Felsenstein 2005) that treats two

534 sequences containing the same ambiguity codes as identical. In effect, the approach used by

535 Lischer et al. (2014) did not directly investigate the effect of IUPAC ambiguity codes on

536 phylogenetic estimates but rather the effect of removing heterozygous sites. Our approach

537 of analyzing IUPAC consensus sequences under the MSC in STACEY, on the other hand,

538 properly integrates these IUPAC ambiguity codes into the calculation of the gene tree

539 likelihoods. Thus, we conclude that IUPAC ambiguity codes introduce a bias towards

540 younger divergence times, even when properly integrated into the phylogenetic model. The

541 underlying cause of this discrepancy should be further investigated in future studies.

542 We also tested whether the improved performance of phased allele sequences in

543 comparison to contig or IUPAC consensus sequence data may merely be an effect of

544 doubling the number of sequences in the MSAs, by analyzing a dataset of chimeric allele

545 sequences with randomly shuffled SNPs. As with the contig data, the chimeric allele data

546 led to an overestimation of the number of coalescent species (Fig. 4d) and to a biased

547 estimation towards older divergence times (Fig. 5). The fact that contig sequences and

548 chimeric allele sequences produce very similar results in our analyses is not surprising,

549 because contigs, themselves, represent chimeric consensus sequences of the variation found

550 at a locus within an individual. The similarity of the results between contig MSAs and

551 chimeric allele MSAs also shows that the number of sequences being analyzed does not

552 affect the estimated topology, species delimitation or divergence time estimates (Figs. 4

553 and 5).

554 Based on the findings discussed above, we conclude that proper phasing of

555 heterozygous positions is preferable to the alternative of coding heterozygous sites as

556 IUPAC ambiguity codes, particularly when the estimation of divergence times is of interest.

557 Further, allele sequences are theoretically more appropriate input for coalescent models and

558 should be the preferred data type input to these models. The scalability of this approach

559 to larger sample sizes and the applicability of our results to studies of older divergences are

560 questions that should be investigated in future studies.

561 One additional issue that we do not address in this study are the effects of

562 sequencing errors. While sequencing errors can potentially be a serious issue particularly for

563 datasets affected by low read coverage, we do not expect sequencing errors to be assembled

564 into our final allele sequences, due to our relatively high read coverage per exported variant

565 (>three reads each). The effects of sequencing errors and incorrectly inferred read

566 variability on downstream analyses are subjects that need to be explored in future studies.

567 *Practicality of Using Phased Allele Data in Multilocus Phylogenetics*

568 In this study, we analyze MSAs resulting from the different processing schemes in a MSC

569 framework using the STACEY BirthDeathCollapse tree model. However, due to the size

570 (number of samples and loci) of many sequence capture datasets, it is often unfeasible to

571 analyze all MSAs jointly in one MSC analysis because of computational limitations (Smith

572 et al. 2014; Manthey et al. 2016). This problem is exacerbated when working with allele

573 MSAs compared to the contig or IUPAC consensus approach, because each alignment

574 contains twice the number of sequences, leading to a doubling of tips in all estimated gene

575 trees. Here we outline three different strategies of addressing this problem:

576 1. One reasonable approach to data reduction is to use a subset of the allele MSAs

577 for phylogeny estimation. We chose this approach here and reduced the UCE dataset from

578 820 MSAs to 150 MSAs in order to reach convergence of the MCMC (BirthDeathCollapse

579 without taxon-assignments) within a reasonable time frame (three to four days, single core

580 on a Mac Pro, Late 2013, 3.5 GHz 6-Core Intel Xeon E5 processor). This approach has the

581 advantage that we can fully integrate the allelic sequence information and avoid *a priori*

582 assignments of allele sequences to taxa. However this approach discards the majority of the

583 multilocus information by excluding most MSAs from the analysis.

584 2. An alternative approach to data reduction, while keeping the multilocus

585 information of all loci, is to analyze only a single polymorphic position (SNP) per MSA

586 using SNAPP (Bryant et al. 2012). We find that phased allele MSAs provide an excellent

587 template for SNP extraction; since all polymorphisms present in the allele sequences have

588 already undergone quality and coverage filters, it is very straightforward to extract SNPs

589 directly from the allele MSAs. We provide an open-source script for this purpose which

590 also converts the extracted SNPs into a SNAPP compatible format. In our study, this

591 approach produced the correct species tree topology and also estimated the relative

592 node-heights correctly (Fig. 4f). However, SNAPP can only estimate relative and not

593 absolute values for divergence times (Bryant et al. 2012), in contrast to sequence-based

594 analyses (Fig. 4a-d) that deliver absolute divergence time estimates. A more thorough

595 discussion about extracting SNPs from sequence capture data can be found in online

596 Appendix 4 (Supplemental Material).

597        3. Another common approach is to abdicate the more appropriate but

598 computationally heavy co-estimation of gene trees and species trees of the MCMC-based

599 MSC methods and chose species tree methods that separate gene tree and species tree

600 estimation into two consecutive steps. This family of methods is often referred to as

601 summary coalescent methods. In this approach gene trees are estimated separately for each

602 MSA. In a subsequent step, the estimated gene trees are used to infer the most likely

603 species tree. The advantage of this approach is that the number of independent loci being

604 analyzed does not constitute a serious computational limitation, because every gene tree is

605 estimated independently, which allows for efficient computational parallelization. On the

606 other hand, summary coalescent methods are sensitive to the number of informative sites

607 per individual locus (Gatesy and Springer 2014; Springer and Gatesy 2014). Given that

608 our phased allele MSAs contained on average 60% more polymorphic sites than the contig

609 MSAs (69% for the simulated data), we argue that phased allele MSAs may lead to more

610 precise phylogenetic estimates under the summary coalescent approach in comparison to

611 contig MSAs. In our case, the summary coalescent approach was not very suitable, due to

612 rather conserved alignments with limited number of informative sites for individual gene

613 tree inference, which obscured the inference of branch lengths in the species tree (online

614 Appendix 3). However, in the case of our simulated data, we observed a more precise

615 estimate of the species tree topology based on phased allele MSAs when compared to those

616 based on contig MSAs (online Appendix 3). In conclusion the summary coalescent

617 approach can be suitable if the individual alignments contain a sufficient number of

618 parsimony informative sites for gene tree inference, and for this reason it is likely that

619 phased allele MSAs might return more precise phylogenetic estimates than contig MSAs.

620 However, further simulation studies are required to properly test this hypothesis.


621                    *Phylogenetic relationships in Topaza*


622 *One or two species?.*— Our results show a separation of two lineages within the genus

623 *Topaza* that is dated at ca. 2.4 Ma in the mitochondrial tree (Fig. 2 and online Appendix

624 1). These lineages are consistent with the previously described morphospecies *T. pyra*

625 (Gould, 1846) and *T. pella* (Linnaeus, 1758) that are generally accepted in the

626 ornithological community (Hu et al. 2000; del Hoyo et al. 2016a). However, the species

627 status of *T. pyra* has been challenged by some authors (Ornés-Schmitz and Schuchmann

628 2011; Schuchmann 1999). These authors concluded that *Topaza* is a monotypic genus with

629 *T. pyra* being a subspecies of *T. pella*, which they refer to as *T. pella pyra*. Our results

630 consistently support *T. pyra* as a separate lineage across all analyses, lending no support

631 for the conspecificity of these two taxa (Fig. 3).

632 *Genetic divergence within morphospecies.*— One aim of this study was to evaluate the

633 genetic structure within the two morphospecies, *T. pyra* and *T. pella*. The mitochondrial

634 tree shows two divergent clades within *T. pyra* (Fig. 2 and online Appendix 1), but these

635 clades are not strongly supported by the UCE data (Fig. 3), even though the allele

636 sequence data are picking up a signal that possibly indicates two clades are in the process

637 of diversifying (Fig. 3b). For *T. pella*, on the other hand, we consistently find the same

638 clades throughout all multilocus MSC analyses (Fig. 3), leading us to distinguish between

639 the following populations that are congruent with previous morphological subspecies

640 descriptions: a northern *T. pella* population (*T. pella pella*), a southern *T. pella*

641 population (*T. pella microrhyncha*) and a separate population occupying the estuary

642 region of Amazon River (*T. pella smaragdula*). We discuss these phylogenetic conclusions

643 in more detail in online Appendix 5 (Supplemental Material).

644 *Summarizing biogeographic remarks.*— The presence of genetically similar individuals

645 sampled at great geographic distances (e.g. samples 5 and 6) suggests that *Topaza*

646 hummingbirds maintain high levels of gene flow across vast distances of rainforest habitat.

647 At the same time, we find indicators of phylogenetic structure within species,

648 distinguishing samples that are separated by only a small geographic distance (see e.g.

649 samples 6 and 8). These samples are however separated by the Amazon River, which has

650 been found to constitute a dispersal barrier for various species of birds and many other

651 animals (Remsen and Parker 1983; Clair 2003; Hayes and Sewlal 2004; Moore et al. 2008;

652 Fernandes et al. 2012; Ribas et al. 2012; Thom and Aleixo 2015). Even though some

653 hummingbird species are known to disperse across large distances (Wyman et al. 2004;

654 Russell et al. 1994), the Amazon River and its associated habitats (such as seasonally

655 flooded forests) may be part of a complex network of factors that inhibit gene flow among

656 populations of *Topaza* hummingbirds.

# Conclusions

658 This study provides evidence that the assembly of phased allele sequence MSAs improves

659 phylogenetic inference under the MSC model. We find that contig sequences, on the other

660 hand, which are commonly used for phylogenetic inference, lead to biases in the estimation

661 of divergence times. Additionally, phased allele sequence MSAs provide a useful template

662 for the extraction of SNP data, and SNP data can be applied as an alternative dataset for

663 phylogenetic inference, circumventing some computational limitations when analyzing

664 multilocus full-sequence data with MCMC-based MSC methods. Our empirical results

665 suggest the separation of two species within the genus *Topaza*, and we further find genetic

666 structure within one of these species, justifying the definition of separate subspecies. Based

667 on our empirical and simulated results, we conclude that allele phasing should be

668 considered as one "best practice" for processing sequence capture data, although the

669 sample-size, phylogenetic scale, and analytical limitations of this approach have not yet

670 been well-established.

# Supplementary Material

672 Supplementary material, including Supplemental Figs. S1-S11, Supplemental Tables S1

673 and S2, online Appendices 1-5 as well as data files, can be found in the Dryad data

674 repository at https://doi.org/10.5061/dryad.hq3vq.

# Availability

676 The documentation for the allele phasing workflow, which we included into the PHYLUCE

677 pipeline, can be found here:

678 `http://phyluce.readthedocs.io/en/latest/tutorial-two.html`. The script for

679 extracting SNPs from MSAs is available here:

680 `https://github.com/tobiashofmann88/snp_extraction_from_alignments`. All

681 processing and analyses steps executed on the data are stored in bash-scripts on our

682 project GitHub page at `https://github.com/tobiashofmann88/topaza_uce`. The raw

683 sequencing reads are stored in the NCBI Short Read Archive (SRA) at

684 `https://www.ncbi.nlm.nih.gov/sra/SRP135707`.

# ACKNOWLEDGMENTS

# FUNDING

∗

References

708 Bodily, P. M., M. Fujimoto, C. Ortega, N. Okuda, J. C. Price, M. J. Clement, and Q. Snell.
709  2015. Heterozygous genome assembly via binary classification of homologous sequence.
710  BMC Bioinformatics 16:S5.

711 Bolger, A. M., M. Lohse, and B. Usadel. 2014. Trimmomatic: a flexible trimmer for
712  Illumina sequence data. Bioinformatics 30:2114–20.

713 Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard,
714  A. Rambaut, and A. J. Drummond. 2014. BEAST 2: a software platform for Bayesian
715  evolutionary analysis. PLoS Computational Biology 10:e1003537.

716 Bryant, D., R. Bouckaert, J. Felsenstein, N. A. Rosenberg, and A. RoyChoudhury. 2012.
717  Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a
718  full coalescent analysis. Molecular Biology and Evolution 29:1917–32.

719 Clair, C. C. S. 2003. Comparative permeability of roads, rivers, and meadows to songbirds
720  in Banff national park. Conservation Biology 17:1151–1160.

721 Degnan, J. H. and N. a. Rosenberg. 2009. Gene tree discordance, phylogenetic inference
722  and the multispecies coalescent. Trends in Ecology and Evolution 24:332–340.

723 del Hoyo, J., N. Collar, G. Kirwan, and P. Boesman. 2016a. Fiery Topaz (*Topaza pyra*). *in*
724  Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,
725  D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.

726 del Hoyo, J., A. Elliott, J. Sargatal, D. Christie, and E. de Juana. 2016b. Handbook of the
727  Birds of the World Alive. Lynx Edicions, Barcelona, Spain.

728 Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics
729  with BEAUti and the BEAST 1.7. Molecular Biology and Evolution 29:1969–73.

731  Edwards, S. V., L. Liu, and D. K. Pearl. 2007. High-resolution species trees without
732      concatenation. Proceedings of the National Academy of Sciences 104:5936–5941.

733  Eriksson, J. S., J. L. Blanco-Pastor, F. Sousa, Y. J. Bertrand, and B. E. Pfeil. 2017. A
734      cryptic species produced by autopolyploidy and subsequent introgression involving
735      Medicago prostrata (Fabaceae). Molecular Phylogenetics and Evolution 107:367–381.

736  Faircloth, B. C. 2015. PHYLUCE is a software package for the analysis of conserved
737      genomic loci. Bioinformatics 32:786–788.

738  Faircloth, B. C., J. E. McCormack, N. G. Crawford, M. G. Harvey, R. T. Brumfield, and
739      T. C. Glenn. 2012. Ultraconserved elements anchor thousands of genetic markers
740      spanning multiple evolutionary timescales. Systematic Biology 61:717–26.

741  Faircloth, B. C., L. Sorenson, F. Santini, and M. E. Alfaro. 2013. A phylogenomic
742      perspective on the radiation of ray-finned fishes based upon targeted sequencing of
743      ultraconserved elements (UCEs). PLoS ONE 8:e65923.

744  Felsenstein, J. 2005. Phylip (phylogeny inference package) version 3.6. distributed by the
745      author. dep genome sci univ washington, seattle.

746  Fernandes, A. M., M. Wink, and A. Aleixo. 2012. Phylogeography of the chestnut-tailed
747      antbird (*Myrmeciza hemimelaena*) clarifies the role of rivers in Amazonian biogeography.
748      Journal of Biogeography 39:1524–1535.

749  Garrick, R. C., P. Sunnucks, and R. J. Dyer. 2010. Nuclear gene phylogeography using
750      PHASE: dealing with unresolved genotypes, lost alleles, and systematic bias in
751      parameter estimation. BMC Evolutionary Biology 10:118.

752  Gatesy, J. and M. S. Springer. 2014. Phylogenetic analysis at deep timescales: unreliable

753  gene trees, bypassed hidden support, and the coalescence/concatalescence conundrum.

754  Molecular Phylogenetics and Evolution 80:231–266.

755  Giarla, T. C. and J. A. Esselstyn. 2015. The challenges of resolving a rapid, recent

756  radiation: empirical and simulated phylogenomics of philippine shrews. Systematic

757  Biology 64:727–740.

758  Gnirke, A., A. Melnikov, J. Maguire, P. Rogov, E. M. LeProust, W. Brockman, T. Fennell,

759  G. Giannoukos, S. Fisher, C. Russ, S. Gabriel, D. B. Jaffe, E. S. Lander, and

760  C. Nusbaum. 2009. Solution hybrid selection with ultra-long oligonucleotides for

761  massively parallel targeted sequencing. Nature Biotechnology 27:182–189.

762  Harvey, M. G., B. T. Smith, T. C. Glenn, B. C. Faircloth, and R. T. Brumfield. 2016.

763  Sequence capture versus restriction site associated DNA sequencing for shallow

764  systematics. Systematic Biology Advance Access syw036.

765  Hayes, F. E. and J. A. N. Sewlal. 2004. The Amazon River as a dispersal barrier to

766  passerine birds: effects of river width, habitat and taxonomy. Journal of Biogeography

767  31:1809–1818.

768  He, D., A. Choi, K. Pipatsrisawat, A. Darwiche, and E. Eskin. 2010. Optimal algorithms

769  for haplotype assembly from whole-genome sequence data. Bioinformatics 26:i183–i190.

770  Hu, D.-S., L. Joseph, and D. J. Agro. 2000. Distribution, variation, and taxonomy of

771  *Topaza* Hummingbirds (Aves: Trochilidae). Ornitologia Neotropical 11:123–142.

772  Iqbal, Z., M. Caccamo, I. Turner, P. Flicek, and G. McVean. 2012. De novo assembly and

773  genotyping of variants using colored de Bruijn graphs. Nature Genetics 44:226–232.

774  Jones, G. 2017. Algorithmic improvements to species delimitation and phylogeny estimation

775  under the multispecies coalescent. Journal of Mathematical Biology 74:447–467.

776  Jones, G., Z. Aydin, and B. Oxelman. 2014. DISSECT: an assignment-free Bayesian

777    discovery method for species delimitation under the multispecies coalescent.

778    Bioinformatics 31:991–998.

779  Katoh, K., G. Asimenos, and H. Toh. 2009. Multiple alignment of DNA sequences with

780    MAFFT. Methods in Molecular Biology 537:39–64.

781  Kolaczkowski, B. and J. W. Thornton. 2004. Performance of maximum parsimony and

782    likelihood phylogenetics when evolution is heterogeneous. Nature 431:980–984.

783  Kubatko, L. S. and J. H. Degnan. 2007. Inconsistency of Phylogenetic Estimates from

784    Concatenated Data under Coalescence. Systematic Biology 56:17–24.

785  Leaché, A. D., B. L. Banbury, J. Felsenstein, A. N. M. De Oca, and A. Stamatakis. 2015.

786    Short tree, long tree, right tree, wrong tree: New acquisition bias corrections for inferring

787    SNP phylogenies. Systematic Biology 64:1032–1047.

788  Leaché, A. D. and J. R. Oaks. 2017. The Utility of Single Nucleotide Polymorphism (SNP)

789    Data in Phylogenetics. Annual Review of Ecology, Evolution, and Systematics 48:69–84.

790  Lerner, H. R., M. Meyer, H. F. James, M. Hofreiter, and R. C. Fleischer. 2011. Multilocus

791    resolution of phylogeny and timescale in the extant adaptive radiation of Hawaiian

792    honeycreepers. Current Biology 21:1838–1844.

793  Li, H. and R. Durbin. 2010. Fast and accurate long-read alignment with Burrows-Wheeler

794    transform. Bioinformatics 26:589–595.

795  Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis,

796    and R. Durbin. 2009. The Sequence Alignment/Map format and SAMtools.

797    Bioinformatics 25:2078–9.

798  Lischer, H. E., L. Excoffier, and G. Heckel. 2014. Ignoring heterozygous sites biases

799    phylogenomic estimates of divergence times: Implications for the evolutionary history of

800    microtus voles. Molecular Biology and Evolution 31:817–831.

801  Manthey, J. D., L. C. Campillo, K. J. Burns, and R. G. Moyle. 2016. Comparison of

802    target-capture and restriction-site associated DNA sequencing for phylogenomics: a test

803    in cardinalid tanagers (Aves, Genus: *Piranga*). Systematic Biology Advance Access

804    syw005.

805  McCormack, J. E., B. C. Faircloth, N. G. Crawford, P. A. Gowaty, R. T. Brumfield, and

806    T. C. Glenn. 2012. Ultraconserved elements are novel phylogenomic markers that resolve

807    placental mammal phylogeny when combined with species-tree analysis. Genome

808    Research 22:746–754.

809  McGuire, J., C. C. Witt, J. V. Remsen, A. Corl, D. L. Rabosky, D. L. Altshuler, and

810    R. Dudley. 2014. Molecular phylogenetics and the diversification of hummingbirds.

811    Current Biology 24:910–916.

812  McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky,

813    K. Garimella, D. Altshuler, S. Gabriel, M. Daly, and M. A. DePristo. 2010. The Genome

814    Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA

815    sequencing data. Genome research 20:1297–303.

816  Meiklejohn, K. A., B. C. Faircloth, T. C. Glenn, R. T. Kimball, and E. L. Braun. 2016.

817    Analysis of a rapid evolutionary radiation using ultraconserved elements (UCEs):

818    Evidence for a bias in some multispecies coalescent methods. Systematic Biology

819    Advance Access syw014.

820  Milne, I., G. Stephen, M. Bayer, P. J. A. Cock, L. Pritchard, L. Cardle, P. D. Shaw, and

821    D. Marshall. 2013. Using Tablet for visual exploration of second-generation sequencing

822    data. Briefings in Bioinformatics 14:193–202.

823  Mirarab, S., R. Reaz, M. S. Bayzid, T. Zimmermann, M. S. Swenson, and T. Warnow.

824    2014. ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics

825    30:541–548.

826  Moore, R. P., W. D. Robinson, I. J. Lovette, and T. R. Robinson. 2008. Experimental

827    evidence for extreme dispersal limitation in tropical forest birds. Ecology Letters

828    11:960–968.

829  Mossel, E. and E. Vigoda. 2005. Phylogenetic MCMC algorithms are misleading on

830    mixtures of trees. Science 309:2207–9.

831  Ornés-Schmitz, A. and K. L. Schuchmann. 2011. Taxonomic review and phylogeny of the

832    hummingbird genus *Topaza* (Gray, 1840) using plumage color spectral information.

833    Ornitologia Neotropical Pages 25–38.

834  Peters, J. L. 1945. Check-list of birds of the world. Volume 5 ed. Harvard Univ. Press,

835    Cambridge, Massachusetts.

836  Potts, A. J., T. A. Hedderson, and G. W. Grimm. 2014. Constructing Phylogenies in the

837    Presence Of Intra-Individual Site Polymorphisms (2ISPs) with a Focus on the Nuclear

838    Ribosomal Cistron. Systematic Biology 63:1–16.

839  Rambaut, A., M. A. Suchard, W. Xie, and A. Drummond. 2013. Tracer v1.6.

840  Rannala, B. and Z. Yang. 2003. Bayes estimation of species divergence times and ancestral

841    population sizes using DNA sequences from multiple loci. Genetics 164:1645–1656.

842  Remsen, J. V. and T. A. Parker. 1983. Contribution of river-created habitats to bird

843    species richness in Amazonia. Biotropica 15:223–231.

844    Ribas, C. C., a. Aleixo, a. C. R. Nogueira, C. Y. Miyaki, and J. Cracraft. 2012. A
845        palaeobiogeographic model for biotic diversification within Amazonia over the past three
846        million years. Proceedings of the Royal Society B: Biological Sciences 279:681–689.

847    Russell, R. W., F. L. Carpenter, M. A. Hixon, and D. C. Paton. 1994. The impact of
848        variation in stopover habitat quality on migrant rufous hummingbirds. Conservation
849        Biology 8:483–490.

850    Schrempf, D., B. Q. Minh, N. De Maio, A. von Haeseler, and C. Kosiol. 2016. Reversible
851        polymorphism-aware phylogenetic models and their application to tree inference. Journal
852        of Theoretical Biology 407:362–370.

853    Schuchmann, K., G. Kirwan, and P. Boesman. 2016. Crimson Topaz (*Topaza pella*). *in*
854        Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,
855        D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.

856    Schuchmann, K. L. 1999. Family Trochilidae (hummingbirds). Pages 468–680 *in* Handbook
857        of the Birds of the World Alive (J. del Hoyo, A. Elliott, and J. Sargatal, eds.) volume 5
858        ed. Lynx Edicions, Barcelona, Spain.

859    Simpson, J. T., K. Wong, S. D. Jackman, J. E. Schein, S. J. M. Jones, and I. Birol. 2009.
860        ABySS: a parallel assembler for short read sequence data. Genome Research 19:1117–23.

861    Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014.
862        Target capture and massively parallel sequencing of ultraconserved elements for
863        comparative studies at shallow evolutionary time scales. Systematic Biology 63:83–95.

864    Springer, M. S. and J. Gatesy. 2014. Land plant origins and coalescence confusion. Trends
865        in Plant Science 19:267–9.

866 Sullivan, B. L., C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling. 2009. eBird:

867     A citizen-based bird observation network in the biological sciences. Biological

868     Conservation 142:2282–2292.

869 Thom, G. and A. Aleixo. 2015. Cryptic speciation in the white-shouldered antshrike

870     (*Thamnophilus aethiops*, Aves - Thamnophilidae): The tale of a transcontinental

871     radiation across rivers in lowland Amazonia and the northeastern Atlantic Forest.

872     Molecular Phylogenetics and Evolution 82:95–110.

873 Wyman, S. K., R. K. Jansen, and J. L. Boore. 2004. Automatic annotation of organellar

874     genomes with DOGMA. Bioinformatics 20:3252–5.

875 Yang, Z. 2015. The BPP program for species tree estimation and species delimitation.

876     Current Zoology 61:854–865.

877 Yu, L., Y.-W. Li, O. a. Ryder, and Y.-P. Zhang. 2007. Analysis of complete mitochondrial

878     genome sequences increases phylogenetic resolution of bears (Ursidae), a mammalian

879     family that experienced rapid speciation. BMC Evolutionary Biology 7:198.

Table 1: Sequenced specimens and coordinates of their sampling locations, subspecies identifications based on morphological characters, abbreviation for sample providers: INPA = Instituto Nacional de Pesquisas da Amazônia, MPEG = Museum Paraense Emílio Goeldi, USNM = NMNH, Smithsonian Institution, Washington DC, USA.

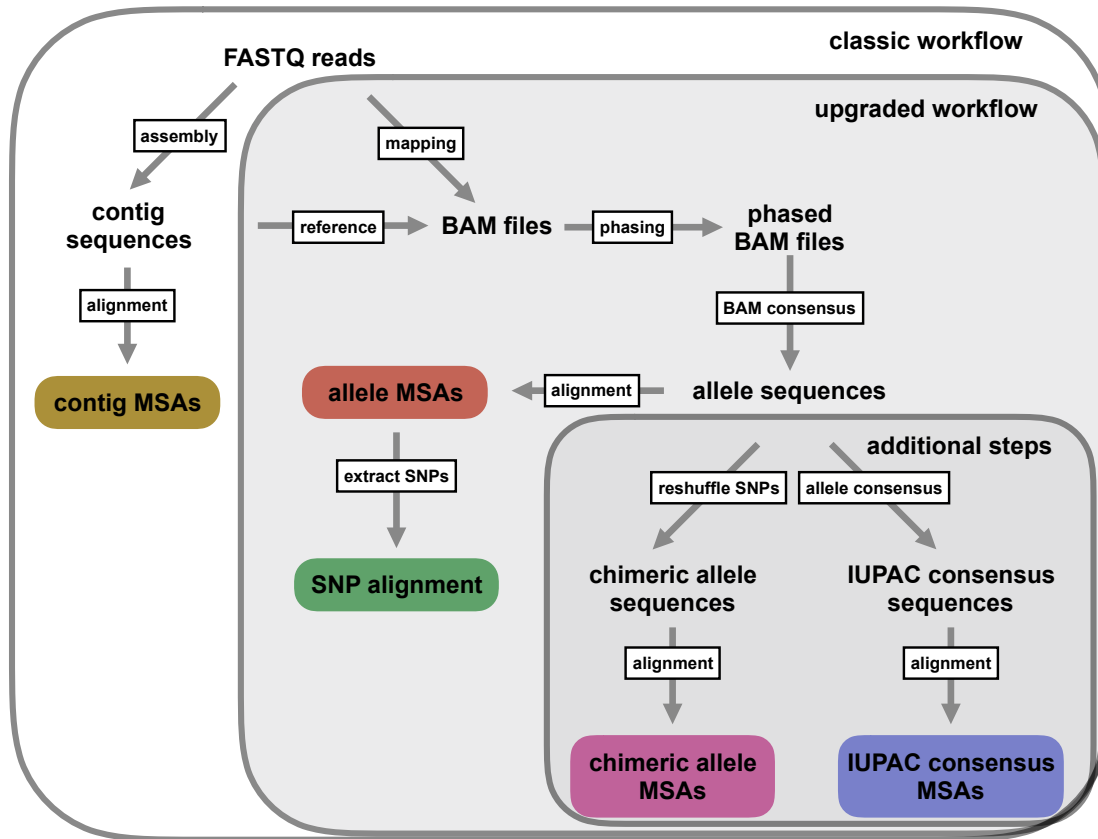| ID | Taxon | Subspecies | Voucher number | Latitude | Longitude |
|---|---|---|---|---|---|
| 1 | *Topaza pyra* | *amaruni* | INPA A1106 | -0.044167 | -66.94944 |
| 2 | *T. pyra* | *pyra* | MPEG 62475 | -1.559444 | -65.88006 |
| 3 | *T. pyra* | *pyra* | MPEG 62474 | -4.083889 | -60.66050 |
| 4 | *T. pyra* | *pyra* | MPEG 52721 | -7.350000 | -73.66667 |
| 5 | *T. pella* | NA | USNM 586322 | 7.220000 | -60.29000 |
| 6 | *T. pella* | *pella* | INPA A3319 | -1.927900 | -59.41600 |
| 7 | *T. pella* | *smaragdula* | MPEG 61688 | -1.950000 | -51.60000 |
| 8 | *T. pella* | *microrhyncha* | MPEG 65603 | -5.352417 | -57.47500 |
| 9 | *T. pella* | NA | INPA A6233 | -9.028550 | -64.24231 |
| 10 | *Florisuga fusca* | NA | MPEG 70697 | -15.15972 | -39.04500 |

Figure 1: Depiction of the workflow used in this manuscript. Colored boxes represent different types of multiple sequence alignments (MSAs) used for phylogenetic inference in this study. In addition to the standard UCE workflow (boxlabel: classic workflow) of generating contig MSAs (Faircloth et al. 2012; Smith et al. 2014; Faircloth 2015), we extended the bioinformatic processing in order to generate UCE allele MSAs, and to extract single nucleotide polymorphism (SNPs) from these allele MSAs (boxlabel: upgraded workflow). We added these new functions to the PHYLUCE pipeline (Faircloth 2015). Additional data processing steps (boxlabel: additional steps) were executed in this study in order to test different codings of heterozygous positions.
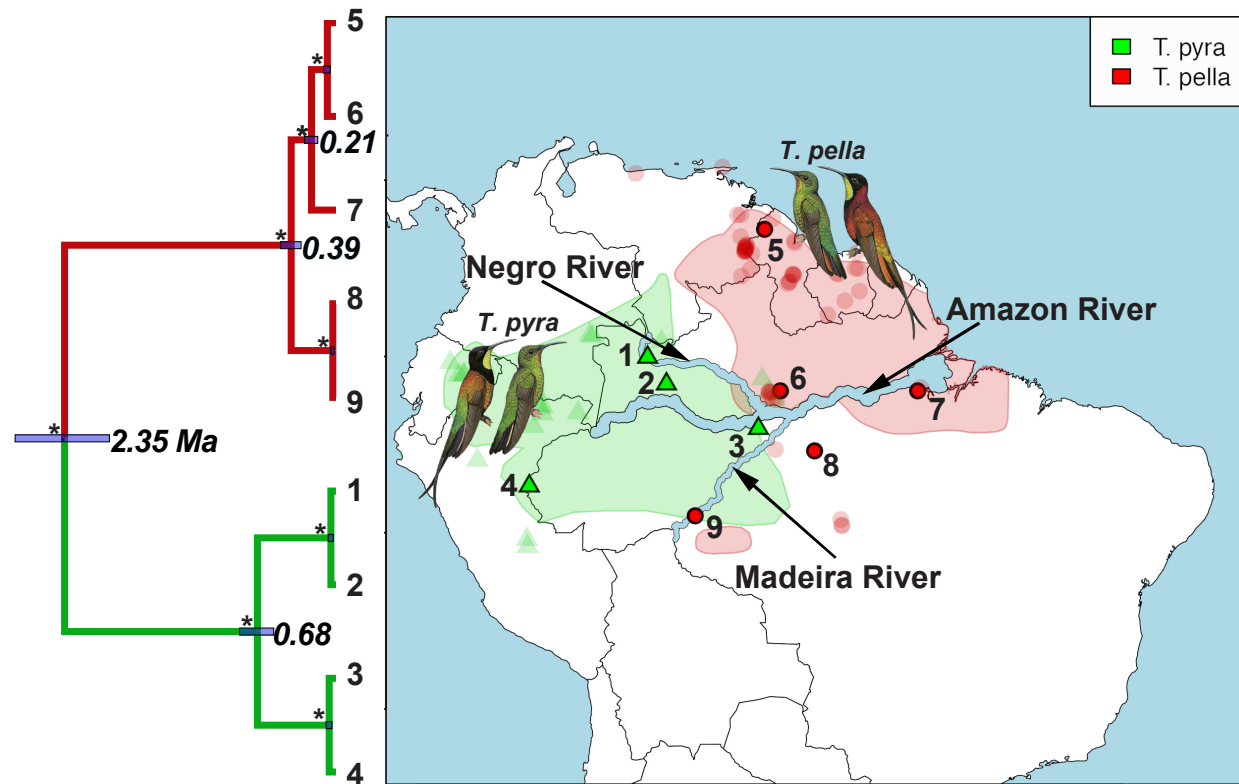
Figure 2: Distribution ranges and mitochondrial phylogeny of the South American hummingbird genus *Topaza*. Tip labels of phylogeny and numbers on map represent sample IDs (Table 1) of sequenced *Topaza* specimens. Node labels in phylogeny show mean divergence time estimates for mitochondrial lineages, with node bars representing the surrounding uncertainty (95% highest posterior density (HPD)). All nodes are supported with 100% posterior probability (PP), as indicated by asterisks. Polygons on map represent distribution ranges of the two morphospecies (*T. pyra* and *T. pella*) as estimated by BirdLife International (http://www.birdlife.org). Transparent symbols (triangles and circles) represent *Topaza* sightings, which were downloaded from the eBird database (Sullivan et al. 2009). The major river systems in the Amazon drainage basin are labeled and emphasized in size for better visibility. *Topaza* illustrations were provided by del Hoyo et al. (2016b).
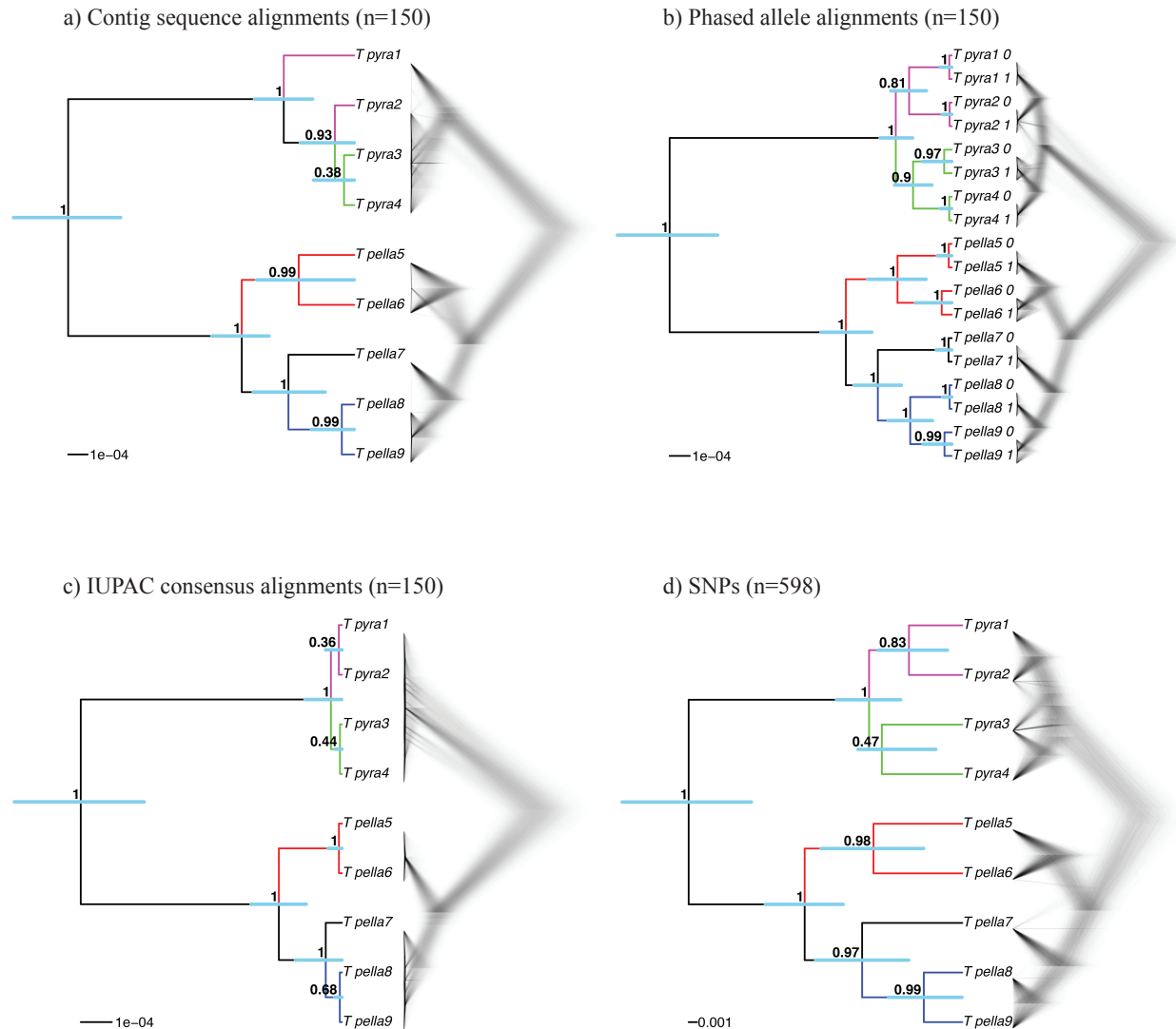
Figure 3: Multispecies Coalescent (MSC) species trees for the empirical *Topaza* data, based on four different data types used in this study: contig sequence MSAs, phased allele sequence MSAs, IUPAC consensus sequence MSAs and SNP data. a) STACEY species tree from UCE contig alignments (n=150), b) STACEY species tree from UCE allele alignments (n=150), c) STACEY species tree from UCE IUPAC consensus alignments (n=150) and d) SNAPP species tree from UCE SNP data (1 SNP per locus if present, n=598). Shown are the Maximum Clade Credibility trees (node values = PP, error-bars = 95% HPD of divergence times) and a plot of the complete posterior species tree distribution (excluding burn-in).
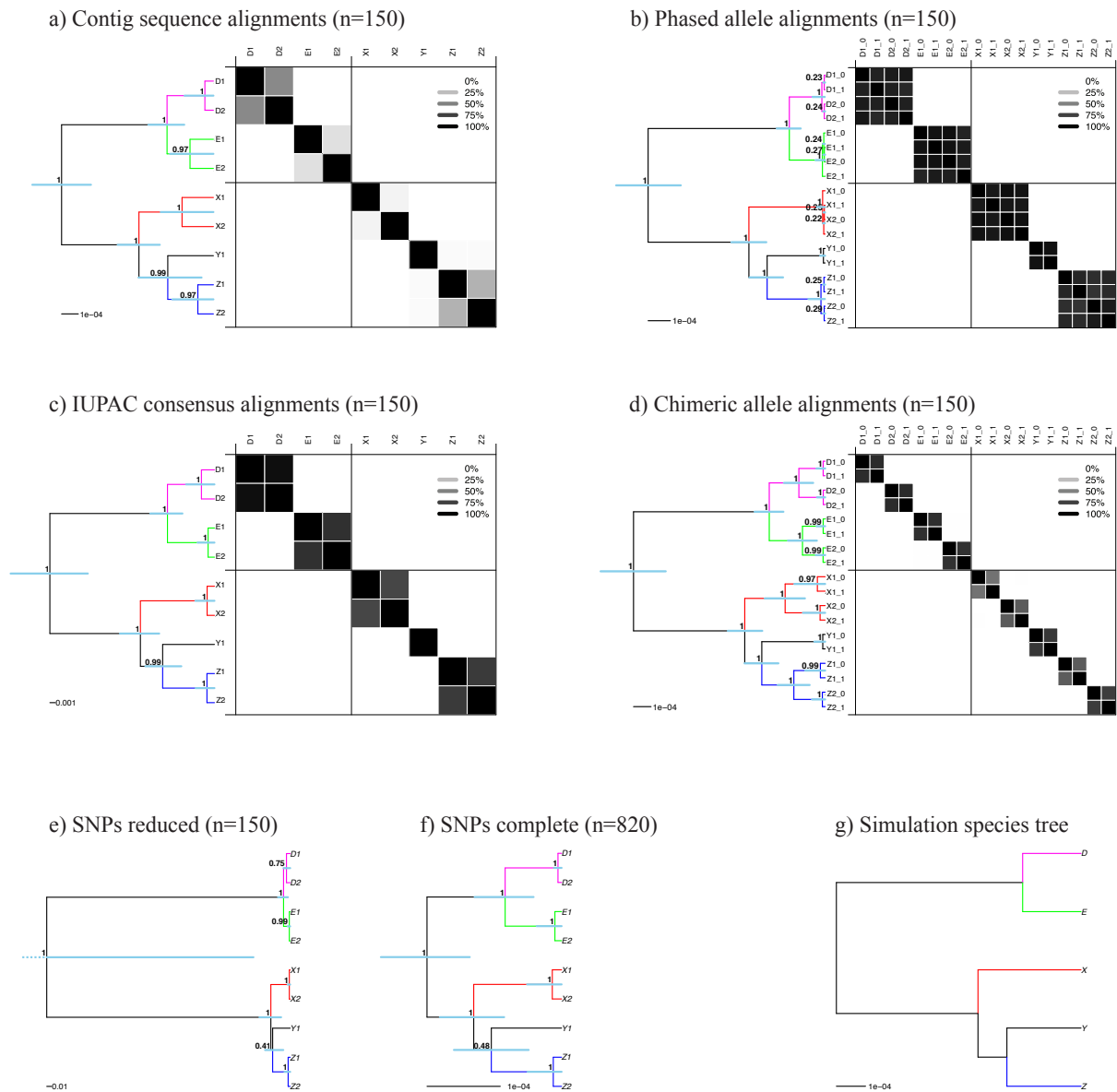
11

Figure 4: MSC species tree results for different data processing schemes of simulated data. a) to d) show the STACEY results of the four different types of MSAs analyzed in this study. Displayed in these panels are the Maximum Clade Credibility trees and the similarity matrices depicting the posterior probability of two samples belonging to the same clade, as calculated with SpeciesDelimitationAnalyser. Dark panels depict a high pairwise similarity, whereas light panels depict low similarity scores (see legend). e) and f) show the Maximum Clade Credibility trees resulting from SNAPP for our two SNP datasets, (reduced and complete). g) shows the species tree under which the sequence data were simulated in this study. Node support values in PP, blue bars representing 95% HPD confidence intervals.

Figure 5: Posterior distributions of divergence times, estimated with STACEY. Each panel represents a different node in the STACEY species tree (see panel titles) and shows density plots of the posterior node-height distribution (excl. 10% burnin) for each of the 4 sequence-based processing schemes: contig sequences, phased allele sequences, IUPAC consensus sequences and chimeric allele sequences (see legend for color-codes). The dotted vertical lines show the means of these posterior distributions. The solid vertical line shows the true node height value, which is the node height for the respective clade in the input species tree, under which the sequence alignments were simulated.