1    Version dated: December 21, 2017

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    [1]*Department of Biological and Environmental Sciences, University of Gothenburg, SE-41319,*

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-41319, Göteborg, 41319, 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, Göteborg, SE-41319, Sweden*

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

18    Environmental Sciences, University of Gothenburg, Carl Skottsbergs Gata 22B, Göteborg,

19    41319, Sweden; E-mail: tobias.hofmann@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 prefer to 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, 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 show that analyzing allele sequences leads to more accurate estimates of tree topology and divergence times than the more common approach of using contig sequences. We conclude that allele phasing may be the most appropriate processing scheme for phylogenetic analyses of UCE data in particular, and sequence capture data, more generally.
(Keywords: UCE, SNP, heterozygous sites, Multispecies Coalescent, gene tree, species tree, Mitochondrial Genome, Trochilidae, Birds, Amazon)

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

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

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

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

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

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

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

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

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

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

56  output by these assemblers often ignore the variants at heterozygous positions that are

57  expected in diploid organisms. Typically, variable positions are treated as sequencing errors

58  and assembly algorithms output the contig containing the more probable (i.e., numerous)

59  variant while discarding the alternative (Iqbal et al. 2012). As a result, the contigs that are

60  produced contain only canonical nucleobases, losing the information for each alternative

61  allele present at each variable position (Fig. 4). Hereafter, we use "contigs" and "contig

62  sequences" to refer to the sequences that are output by *de novo* assemblers.

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

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

65  (also known as "calling" single nucleotide polymorphism (SNPs)) and subsequently sorting

66  (or "phasing") these SNPs into two allele sequences or "haplotypes" which represent alleles

67  on the same chromosome present at that locus. These approaches have been used to

68  estimate demographic parameters such as effective population size, rate of migration, and

69  the amount of gene flow between and within populations. However, it is rarely

70  acknowledged (*c.f.* Lischer et al. 2014; Potts et al. 2014; Schrempf et al. 2016; Eriksson

71  et al. 2017) that allelic sequences are useful for phylogenetic studies to improve the

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

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

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

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

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

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

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

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

80 into phylogenetic methods (Lischer et al. 2014; Potts et al. 2014; Schrempf et al. 2016).

81 One is to code heterozygous sites using IUPAC ambiguity codes and to include these as

82 additional characters in existing substitution models for gene tree and species tree inference

83 (Potts et al. 2014; Schrempf et al. 2016). While these studies demonstrate that integrating

84 additional allelic information in this manner increases accuracy in phylogenetic inference,

85 Lischer et al. (2014) found that coding heterozygous sites as IUPAC ambiguity codes in

86 phylogenetic models biases the results toward older divergence time estimates. Instead,

87 Lischer et al. (2014) introduced a method of repeated random haplotype sampling (RRHS)

88 in which allele sequences are repeatedly concatenated across many loci, using a random

89 haplotype for any given locus in each replicate. In their approach they then analyzed

90 thousands of concatenation replicates separately for phylogenetic tree estimation and

91 summarized the results between replicates, thereby integrating the allelic information in

92 form of uncertainty intervals. However there are two important shortcomings of this

93 approach: 1. concatenating unlinked loci (and in particular allele sequences from unlinked

94 loci) in a random manner is known to produce incorrect topologies (Degnan and Rosenberg

95 2009) often with false confidence (Edwards et al. 2007; Kolaczkowski and Thornton 2004;

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

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

estimation replicates based on extensive amounts of sequence data results in unfeasibly long computation times, particularly for Markov-Chain Monte Carlo (MCMC) based softwares such as MrBayes or BEAST. Hence there is need to find proper solutions to include heterozygous information in phylogenetic analyses, as concluded by Lischer et al. (2014).

Here, we introduce the bioinformatic assembly of allele sequences from UCE data and demonstrate a full integration of allele sequences to species tree estimation under the multispecies coalescent (MSC) model using empirical and simulated data. In our approach, we treat each allelic sequence of an individual at a given locus as an independent sample from the population, and we analyze these sequences using the species tree and delimitation software STACEY (Jones et al. 2014; Jones 2017), which does not require *a priori* clade- or species-assignments. We first demonstrate the empirical utility of our approach by resolving the shallow genetic structure (<1 Ma) within two recognized morphospecies of the South American hummingbird genus *Topaza* by analyzing a set of 2,386 ultraconserved elements (UCEs, see Faircloth et al. (2012)) collected using sequence capture of the 2.5k tetrapod baitset (see http://ultraconserved.org). We then validate this approach, using simulations, and show that allele sequences yield more accurate results in terms of species tree estimation and species delimitation than the contig sequence approach that ignores heterozygous information. Our simulation results further demonstrate that proper phasing of allele sequences outperforms alternative approaches of coding heterozygous information, such as analyzing sequences containing IUPAC ambiguity codes or SNPs. We conclude by demonstrating that phasing sequence capture data can be critical for correct species delimitation and phylogeny estimation, particularly in recently diverged groups, and that analyses using phased alleles should be considered as one "best practice" for analyzing sequence capture datasets in a phylogenetic context.

## Materials and Methods

## Study System

123
124     The genus *Topaza* and its sister genus *Florisuga* form the Topazes group, which
125 together with the Hermits represent the most ancient branch within the hummingbird
126 family (Trochilidae) (McGuire et al. 2014). Topazes are estimated to have diverged as a
127 separate lineage from all other hummingbirds around 21.5 Ma, whereas the most recent
128 common ancestor (MRCA) of *Topaza* and *Florisuga* lived approximately 19 Ma (McGuire
129 et al. 2014). At present, there are two morphospecies recognized within *Topaza*, namely the
130 Fiery Topaz, *T. pyra* (Gould, 1846), and the Crimson Topaz, *T. pella* (Linnaeus, 1758).
131 However, the species status of *T. pyra* has been challenged by some authors (Schuchmann
132 1999; Ornés-Schmitz and Schuchmann 2011), who consider this genus to be monotypic.
133 Topaz hummingbirds are endemic to the Amazonian rainforest and are some of the most
134 spectacular and largest hummingbirds worldwide, measuring up to 23 cm (adult males,
135 including tail feathers) and weighing up to 12 g (Schuchmann et al. 2016; del Hoyo et al.
136 2016a). These birds are usually found in the forest canopy along forest edges and clearings,
137 and are often seen close to river banks (Ornés-Schmitz and Schuchmann 2011). There is
138 morphological evidence for several subspecies within both currently recognized *Topaza*
139 species (Peters 1945; Schuchmann 1999; Hu et al. 2000; Ornés-Schmitz and Schuchmann
140 2011) that we investigate using genetic data.

## Sequence Data Generation

141
142     We extracted DNA from the muscle tissue of 10 vouchered hummingbirds (9
143 *Topaza*, one *Florisuga*, see Table 1) using the Qiagen DNeasy Blood and Tissue Kit
144 according to the manufacturer's instructions (Qiagen GmbH, Hilden, Germany). These
145 samples cover most of the genus' total geographic range (Fig. 1) and all morphologically
146 recognized intraspecific taxa (Schuchmann et al. 2016; del Hoyo et al. 2016a). All samples

Table 1: Specimens sequenced. 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 |

147 were sonicated with a Covaris S220 to a fragment length of 800 bp. Paired-end,

148 size-selected (range 600-800bp) DNA libraries were prepared for sequencing on the Illumina

149 MiSeq platform, using the magnetic-bead based NEXTflexTM Rapid DNA-Seq Kit (Bioo

150 Scientific Corporation, Austin, TX, USA), following the user's manual (v14.02).

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

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

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

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

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

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

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

158 sequence enrichment and sequencing were performed by Sahlgrenska Genomics Core

159 Facility in Gothenburg, Sweden.
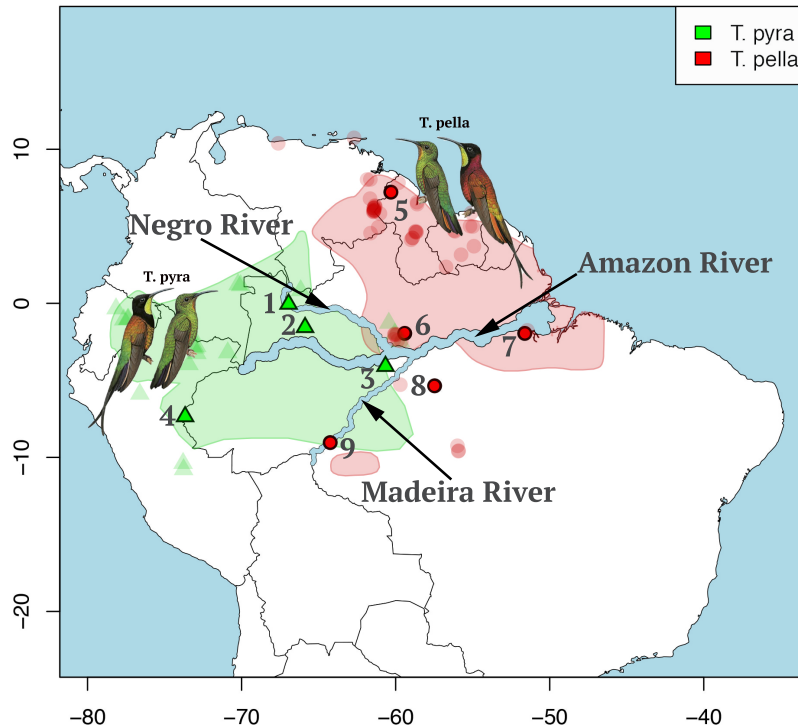
160 *Mitochondrial Genome*

Figure 1: Sample locations of *Topaza* specimens (numbered symbols) in northern South America. Numbers represent sample IDs (Table 1). The colored polygons show the distribution range of the two morphospecies (*T. pyra* = green, *T. pella* = red) 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 marked in blue (not in proportion). *Topaza* illustrations were provided by del Hoyo et al. (2016b).

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

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

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

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

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

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

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

168    (2011) and node-age priors within the genus *Topaza* that were estimated by McGuire et al.
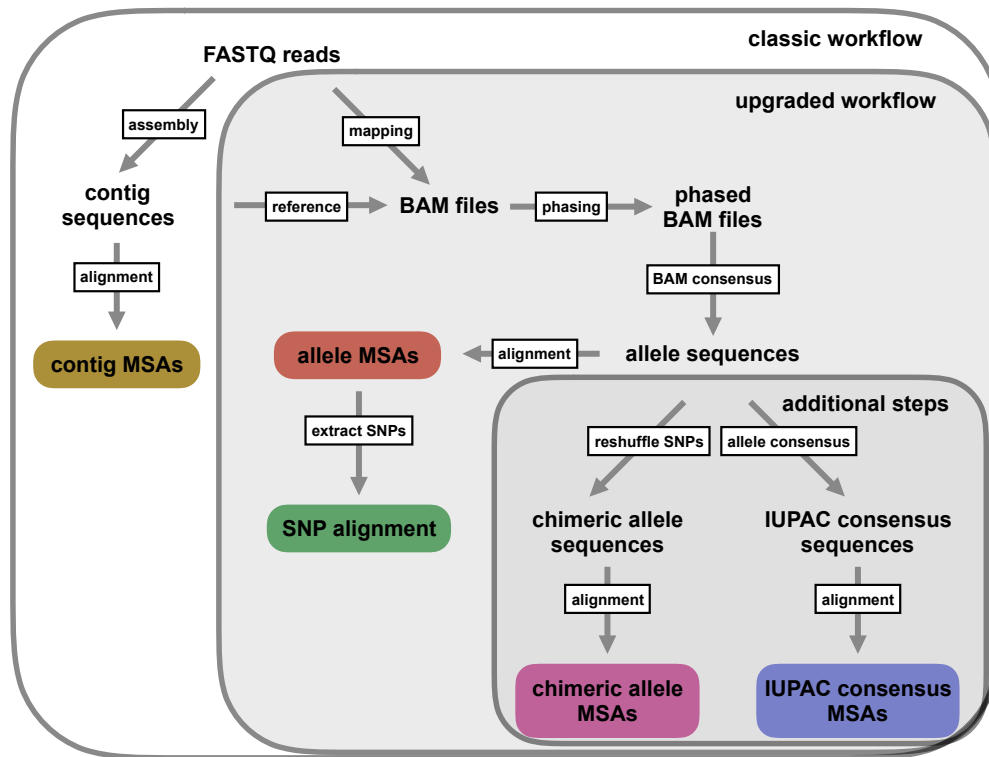
Figure 2: Depiction of the workflow developed here. Colored boxes represent different types of multiple sequence alignments (MSAs) used for phylogenetic inference in this study. In addition to the standard UCE workflow of generating contig MSAs (Faircloth et al. 2012; Smith et al. 2014; Faircloth 2015), we extended the bioinformatic processing in order to generate UCE allele MSAs, and to extract SNPs from these allele MSAs. We added these new functions to the PHYLUCE pipeline (Faircloth 2015). Additional data processing steps were executed in this study in order to test different codings of heterozygous positions.

169  (2014). A detailed description of the assembly and analysis of the mitochondrial genome

170  data can be found in online Appendix 1 (Supplemental Material available on Dryad).

171  *UCE Data Processing*

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

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

174  and are listed and described in the following sections.

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

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

177 (2015)), we generated multiple sequence alignments (MSAs) of contigs for all UCE loci in

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

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

180 the PHYLUCE documentation

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

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

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

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

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

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

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

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

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

190 samples (n=820). These UCE sequences were then aligned for each locus (Fig. 2) using

191 MAFFT (Katoh et al. 2009).

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

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

194 (Fig. 2). To create this new workflow, we extracted all UCE contigs for each sample

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

196 read mapping. We then mapped the cleaned reads against each reference library on a per

197 sample basis, using CLC-mapper from the CLC Workbench software. The mapped reads

198 were sorted and then phased with SAMtools v0.1.19 (Li et al. 2009), using the commands

199 `samtools sort` and `samtools phase`, respectively. This phasing function is based on a

200 dynamic programming algorithm that uses read connectivity across multiple variable sites

201 to determine the two phases of any given diploid locus (He et al. 2010). Further, this

202 algorithm uses paired-end read information to reach connectivity over longer distances and

203 it minimizes the problem of accidentally phasing a sequencing error, by applying the

204 minimum error correction function (He et al. 2010).

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

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

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

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

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

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

211 complete locus. The two phased output files (BAM format) were inspected for proper

212 variant separation for all loci using Tablet (Milne et al. 2013). We then collapsed each

213 BAM file into a single sequence and exported the two resulting allele sequences for each

214 sample in FASTA format. In the next, step we aligned the allele sequences between all

215 samples, separately for each UCE locus, using MAFFT (Fig. 2). We integrated this

216 complete workflow into the UCE processing software PHYLUCE (Faircloth 2015) with

217 slight alterations, one of which is the use of the open-source mapping program bwa (Li and

218 Durbin 2010) in place of CLC-mapper.

219 *3. UCE IUPAC consensus sequence alignments.*— We generated an additional set of

220 alignments by merging the two allele sequences for each individual into one consensus

221 sequence with heterozygous sites coded as IUPAC ambiguity codes

222 (`merge_allele_sequences_ambiguity_codes.py`, available from:

223 github.com/tobiashofmann88/UCE-data-management/). We used this dataset to test

224 whether our allele phasing approach improved phylogenetic inference when compared to

225  the IUPAC consensus approach applied in other studies (where heterozygous positions are

226  simply coded as IUPAC ambiguity codes in a consensus sequence for each locus and

227  individual (Potts et al. 2014; Schrempf et al. 2016)).

228  *4. UCE chimeric allele alignments.*— To investigate whether correct phasing of

229  heterozygous sites is essential or if similar results are achieved by randomly placing

230  variants in either allele sequence, we generated a dataset with chimeric allele sequence

231  alignments. We created these alignments by applying a custom python script

232  (`shuffle_snps_in_allele_alignments.py`, available from:

233  github.com/tobiashofmann88/UCE-data-management/) to the phased allele sequence

234  alignments and randomly shuffling the two variants at each polymorphic position between

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

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

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

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

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

240  individual).

241  *5. UCE SNP alignment.*— A common approach to analyze heterozygous information is to

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

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

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

245  data, due to extremely long computational times. Instead, alignments of unlinked SNPs

246  can be used to infer species trees and species demographics under the MSC model with the

247  BEAST2 package SNAPP (Bryant et al. 2012), a program specifically designed for such

248  data. However, extracting and filtering SNPs from BAM files with existing software (such

249  as the Genome Analysis Toolkit (GATK), McKenna et al. (2010)) and converting these

<sup>250</sup> into a SNAPP compatible format can be cumbersome, because SNAPP requires positions

<sup>251</sup> with exactly two different states, coded in the following manner: individual homozygous for

<sup>252</sup> the original state = "0", heterozygous = "1", and homozygous for the derived state = "2".

<sup>253</sup>      To alleviate this problem, we developed a python function that extracts biallelic

<sup>254</sup> SNPs directly from allele sequence MSAs (`snps_from_uce_alignments.py`, available from:

<sup>255</sup> github.com/tobiashofmann88/UCE-data-management/). Extracting SNPs from MSAs in

<sup>256</sup> this manner is a straightforward and simple way to generate a SNP dataset compatible

<sup>257</sup> with SNAPP, and does not require re-visiting the BAM files. Although a similar program

<sup>258</sup> already exists, which is implemented in the R-package `phrynomics` (Leaché et al. 2015), we

<sup>259</sup> integrated the SNP extraction from allele sequence MSAs into the PHYLUCE pipeline,

<sup>260</sup> and used this approach to extract one position per alignment (to ensure unlinked SNPs)

<sup>261</sup> that had exactly two different states among all *Topaza* samples, not allowing for positions

<sup>262</sup> with missing data or ambiguities. This produced a SNP dataset of 598 unlinked SNPs.

<sup>263</sup> <div align="center">*Generation of Simulated UCE Data*</div>

<sup>264</sup> To assess the accuracy of the phylogenetic inferences resulting from different data

<sup>265</sup> processing approaches, we simulated UCE data similar to those discussed in the five

<sup>266</sup> processing schemes we applied to the empirical *Topaza* data. However, because this

<sup>267</sup> approach required us to simulate allele alignments before generating contig alignments,

<sup>268</sup> steps one and two, below, are reversed from their order, above. For each of the five

<sup>269</sup> processing schemes, we generated and analyzed ten independent simulation replicates.

<sup>270</sup> *1. Simulated allele alignments.*— From the empirical UCE allele alignments, we estimated

<sup>271</sup> species divergence times and population sizes under the MSC model (Rannala and Yang

<sup>272</sup> 2003) using the Bayesian MCMC program BPP v3.1 (Yang 2015). To do this, we used the

<sup>273</sup> A00 model with the species tree topology from the analysis of the allele sequence data in

274 STACEY, assigning the *Topaza* samples to five separate taxa (corresponding to colored

275 clades in Fig. 6b). An initial BPP analysis did not converge in reasonable computational

276 time, a problem that has previously been reported for UCE datasets containing several

277 hundred loci (Giarla and Esselstyn 2015). To avoid this issue, we split the 820 UCE

278 alignments randomly into 10 subsets of equal size (n=82) and analyzed these separately

279 with identical settings in BPP. The MCMC was set for 150,000 generations (burn-in

280 50,000), sampling every 10 generations. We summarized the estimates for population sizes

281 and divergence times across all 10 individual runs. We then applied the mean values of

282 these estimates to the species tree topology, by using the estimated divergence times as

283 branch lengths and estimated population sizes as node values, resulting in the species tree

284 in Fig. 6g. This tree was used to simulate sequence alignments with the MCcoal simulator,

285 which is integrated into BPP. Equivalent to the empirical data, we simulated sequence data

286 for five taxa (D, E, X, Y, and Z) and one outgroup taxon (F, not shown in Fig. 6g). In the

287 simulations, these taxa were simulated as true species under the MSC model. In order to

288 mimic the empirical allele data, we simulated four individuals for species 'D' (equivalent to

289 two allele sequences for 2 samples), four for species 'E', four for species 'X', two for species

290 'Y' (two allele sequences for one sample), four for species 'Z', and two for the outgroup

291 species 'F'. In this manner we simulated 820 UCE allele MSAs of 848 bp length (a value

292 equal to the average alignment length of the empirical allele alignments).

293 *2. Simulated contig alignments.*— To simulate UCE contig MSAs similar to those used in

294 previous studies (Faircloth et al. 2012; McCormack et al. 2012; Smith et al. 2014; Faircloth

295 2015) and output by assemblers like Velvet or Trinity which pick only one of the two

296 variants at a heterozygous site, we merged the sequences within each coalescent species in

297 pairs of two (equivalent to pairs of allele sequences). Each pair of allele sequences was

298 joined into one contig sequence by randomly picking one of the two variants at each

299   heterozygous site across all loci. As in the empirical contig assembly approach, our

300   simulation approach may generate chimeric contig sequences.

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

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

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

304   of randomly picking one of the two variants.

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

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

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

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

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

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

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

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

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

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

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

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

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

318                    *MSC Analyses of Empirical and Simulated UCE Data*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

337 the 150 loci with jModeltest (Supplementary Table S1 available on Dryad) using BIC. We

338 used BEAUTI v2.4.4 to create an input file for STACEY in which we unlinked substitution

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

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

341 fixed the average clock rate for one random locus to 1.0, while estimating all other clock

342 rates in relation to this locus. To ensure that all resulting species trees were scaled to an

343 average clock rate of 1.0, we rescaled every species tree from the posterior distribution

344 using the average clock rate of the respective MCMC step. We applied the

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

346 1e-5 for the collapseHeight parameter. Other settings were: bdcGrowthRate = log normal

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

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

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

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

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

352 every 20,000 generations. Convergence was assessed using Tracer v1.6 (Rambaut et al.

353 2013). We then summarized the posterior tree distribution into one maximum clade

354 credibility tree with TreeAnnotator v2.4.4, discarding the first 10% of trees as burn-in.

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

356 analysis with the program SpeciesDelimitationAnalyser (part of the STACEY

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

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

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

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

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

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

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

364 assigned as its own taxon), we set coalescent rate and mutation rates to be estimated based

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

366 0.00765). We ran the analysis for 10,000,000 generations, sampling trees and other

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

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

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

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

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

372 assignments (assignments as in Fig. 6g).

# RESULTS

## *Mitochondrial Tree (BEAST)*

The BEAST analysis of complete mitochondrial genomes (see online Appendix 1) produced a fully resolved topology (Fig. 3). All nodes were supported by 100% Bayesian posterior probability (PP). We inferred the divergence between the two lineages *T. pyra* and *T. pella* at 2.36 Ma, with 95% of the highest posterior density (HPD) ranging between 1.96 and 2.78 Ma. The tree also shows a separation of two distinct lineages within *T. pyra* at 0.68 Ma (95% HPD: 0.54 - 0.84 Ma), dividing the samples of this morphospecies into a northern and a southern clade, separated by the Amazon River (Fig. 1). A similar, yet slightly more recent split can be seen within *T. pella*. We inferred the age of this split to 0.39 Ma (95% HPD: 0.30 - 0.48 Ma), revealing the same pattern of one northern and one southern clade with the exception of sample 7; this sample from the southern bank of the Amazon River delta is placed together with the samples derived from localities north of the Amazon (samples 5 and 6). Below, we refer to those individuals sampled north of the Amazon River as "northern" and to those sampled south of the Amazon as "southern".

## *UCE Summary Statistics*

*Alignment statistics.*— 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* (for the following statistics we are excluding the outgroup). In this manner, we found that the empirical UCE contig sequence
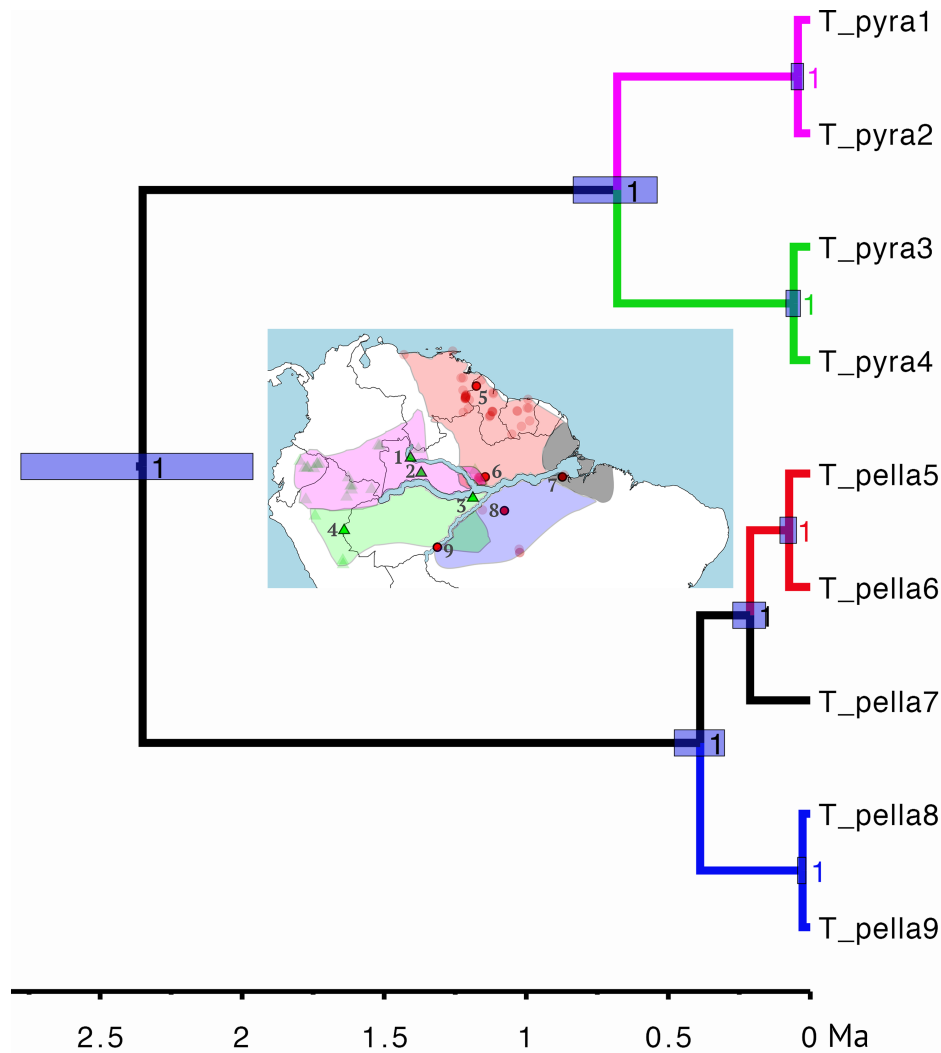
Figure 3: Phylogeny estimated from complete mitochondrial genomes in BEAST. Node support values represent PP. The blue bars at nodes represent the 95% HPD of divergence times. Scale axis shows time units in millions of years. The map in the center shows the potential ranges of the clades that are found in the mitochondrial tree (color-coded). The ranges are based on the BirdLife distribution ranges (Fig. 1) and have been expanded in order to accommodate all *Topaza* occurrence data.

395  alignments had an average of 2.8 polymorphic sites per locus and an average alignment

396  length of 870 bp. In contrast, phasing the empirical UCE data to create allele alignments

397  led to 4.5 polymorphic sites per locus and an average alignment length of 848 bp,

398  representing a 60% increase in polymorphic sites per locus. This increase of polymorphic

399  sites was attributable to the fact that many variants get lost during contig assembly,

400 because ABYSS and other tested contig assemblers, namely Trinity and Velvet, often

401 eliminate one of the two variants at heterozygous positions (see below). The reduced

402 length of the allele alignments in comparison to the contig alignments was due to

403 conservative alignment clipping thresholds implemented in PHYLUCE, which clip

404 alignment ends if less than 50% of sequences are present. Because the allele phasing

405 algorithm divides the FASTQ reads into two allele bins and because a nucleotide is only

406 called if it is supported by at least three high-quality FASTQ reads, we lost some of the

407 nucleotide calls at areas of low read coverage (mostly at the ends of a locus) when

408 comparing the allele sequences to the contig sequences. More information about the

409 distribution of lengths and variable sites within the empirical UCE data can be found in

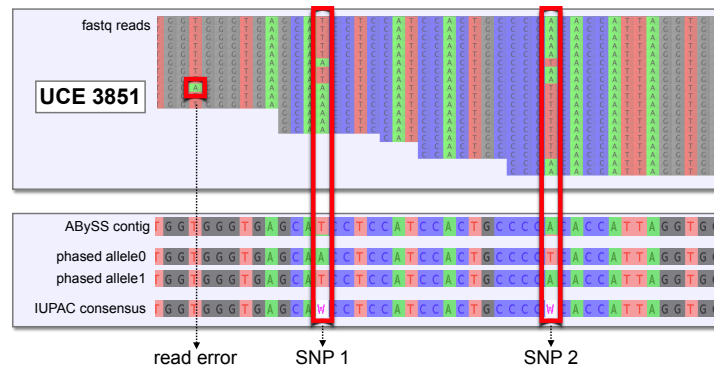410 the Supplementary Figs. S1 and S2 available on Dryad.

411 The simulated contig MSAs had an average of 3.2 polymorphic sites per locus, after

412 excluding the outgroup. The simulated allele MSAs, on the other hand, contained an

413 average of 5.4 polymorphic sites (69% increase). An overview of parsimony informative

414 sites, variable sites and length of each alignment (simulated and empirical data) can be

415 found in Supplementary Table S2 available on Dryad.

416 *ABYSS does not detect heterozygous sites.*— ABYSS occasionally produces contig

417 sequences containing IUPAC ambiguity codes, which suggests that these sites may

418 accurately represent heterozygosity in the read data and that assembly with ABYSS may

419 be preferred to using other assembly algorithms because the resulting contigs contain more

420 information. To validate this assumption, we checked one randomly selected sample

421 (sample 5, *T. pella*) to see if degenerate sites in the contig sequences produced by ABYSS

422 were heterozygous in the phased allele sequences. The results are striking, because there

423 are zero heterozygous sites within the allele sequences for sample 5 that were correctly

424 coded as IUPAC ambiguity codes in the ABYSS contigs (e.g. Fig 4a). Moreover, our
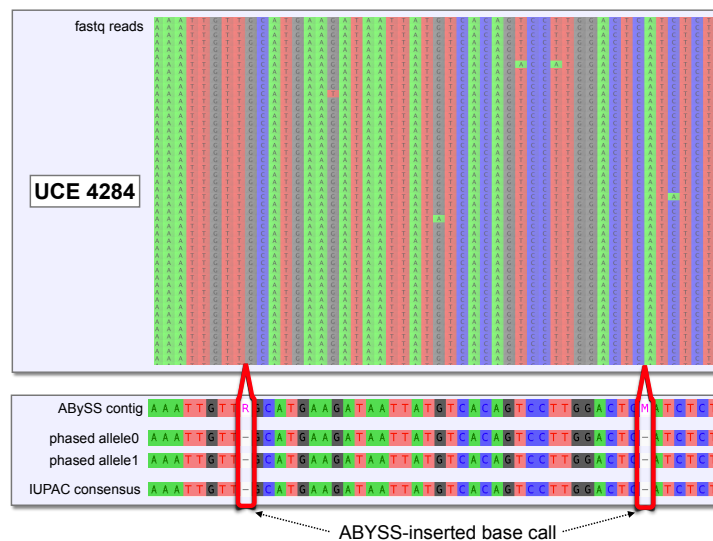
425  phasing approach revealed 343 heterozygous UCE loci with a total of 728 SNPs in sample 5

426  while contig sequences from the ABYSS *de novo* assembly only contained IUPAC

427  ambiguity codes (degenerate bases) at 26 UCE loci. For all other loci, ABYSS output

428  homozygous contig sequences, indicated by the fact that all are free of ambiguity codes.

429  Within the 26 loci containing IUPAC ambiguity codes, ABYSS introduced 473 degenerate

430  bases, most of which constitute blocks of N's. Effectively, all of these ambiguous positions

431  are in places of extremely low FASTQ read coverage (<2 reads per haplotype), with the

432  exception of six positions that are covered by greater than two reads per haplotype.

433  However, even those six positions do not represent true heterozygous sites within sample 5,

434  which becomes apparent when comparing aligned FASTQ reads at those loci with the

435  phased allele sequences with the contig sequences produced by ABYSS (e.g. Fig. 4b).

## *MSC Results of Empirical UCE Data*

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

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

439  and SNPs) converge on similar topologies for relationships among *T. pella*, yet the

440  relationships inferred among *T. pyra* are less clear (Fig. 5 and Supplementary Fig. S3

441  available on Dryad). All analyses strongly support the monophyly of both *T. pyra* and *T.*

442  *pella* with 100% PP. In all MSC analyses, we also see strongly supported genetic structure

443  within *T. pella* ($\geq$ 97% PP), separating the northern samples (5 and 6) from the southern

444  ones (7, 8 and 9). Additionally, within the shallow southern *T. pella* clade, all datasets,

445  with exception of the IUPAC consensus data (Fig. 5c), strongly support a genetic

446  distinction ($\geq$ 99% PP) between sample 7 from the Amazon River delta and the other

447  southern *T. pella* samples (8 and 9). The deep split between northern and southern

448  samples within *T. pyra* on the other hand, which we find in the mitochondrial tree (Fig.

449  3), is not well-supported by the multilocus MSC analyses. However, the analysis of the

(a) Heterozygous position picked up by allele phasing



(b) Erroneous insertion of IUPAC ambiguity by ABYSS

Figure 4: Detection of heterozygous sites in FASTQ reads. The figure shows two UCE loci for sample 5 (*T. pella*). Displayed in both cases are the FASTQ reads, the ABYSS contig sequence, the two phased allele sequences and the correct IUPAC consensus sequence generated from our phased allele sequences. (a) An example of true heterozygous sites, which are correctly represented in the phased allele sequences sequences but are not coded as IUPAC ambiguities in resulting ABYSS contig. Instead ABYSS makes a majority call for this position, thereby masking the heterozygous site by eliminating one of the two variants. This is the case for all heterozygous sites that were picked up by the allele sequences in our data. (b) An example of a UCE locus that contains IUPAC ambiguity codes in the ABYSS contig sequence. Contrary to expectations, the ambiguity calls at these positions are not supported by the FASTQ reads and appear to be inserted by ABYSS at random positions. Our phased allele sequences, on the other hand, represent the FASTQ reads correctly and do not call this position as heterozygous. We observed this same patters across all 26 loci in our data with ABYSS-inserted IUPAC ambiguity codes.

450 allele dataset returns a phylogenetic signal, possibly tracking a genetic divergence between

451 these two clades, but their monophyly is not very strongly supported (Fig. 5b).

## *MSC Results of Simulated Data*

453 *Species tree topology.*— For the simulated data, we analyzed six different datasets under

454 the MSC model: contig sequence MSAs (n=150, STACEY), allele sequence MSAs (n=150,

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

456 STACEY), reduced SNP data (n=150, SNAPP), and the complete SNP dataset (n=820,

457 SNAPP). All resulting species trees (Figs. 6a to 6f) correctly return the topology of the

458 species tree that was used to simulate the data (Fig. 6g). All central nodes in the species

459 trees are supported by $\geq 90\%$ PP in all analyses, with the exception of the species tree

460 resulting from the reduced SNP dataset, which shows very weak support for two nodes and

461 has a large uncertainty interval around the root-height (Fig. 6e). However, these

462 shortcomings disappeared when we added more (unlinked) SNPs to the dataset (Fig. 6f).

463 The full SNP dataset (n=820) produced the correct species tree topology with high node

464 support consistently throughout ten independently simulated datasets (Supplementary Fig.

465 S4 available on Dryad). The SNAPP species tree topology appeared to be unaffected by

466 the chosen clade assignment model; while we allowed every sequence to be its own taxon in

467 Figs. 6e and 6f, we also applied the correct species assignment (Fig. 6g) in two additional

468 analyses (reduced and complete SNP data) that returned the same tree topology

469 (Supplementary Figs. S5 and S6 available on Dryad).

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

471 all four sequence-based MSC analyses (Figs. 6a to 6d), the inferred node heights varied

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

473 For the contig sequence data (Fig. 6a) and the chimeric allele data (Fig. 6d), the node

(a) Contig sequence alignments (n=150)



(b) Phased allele alignments (n=150)



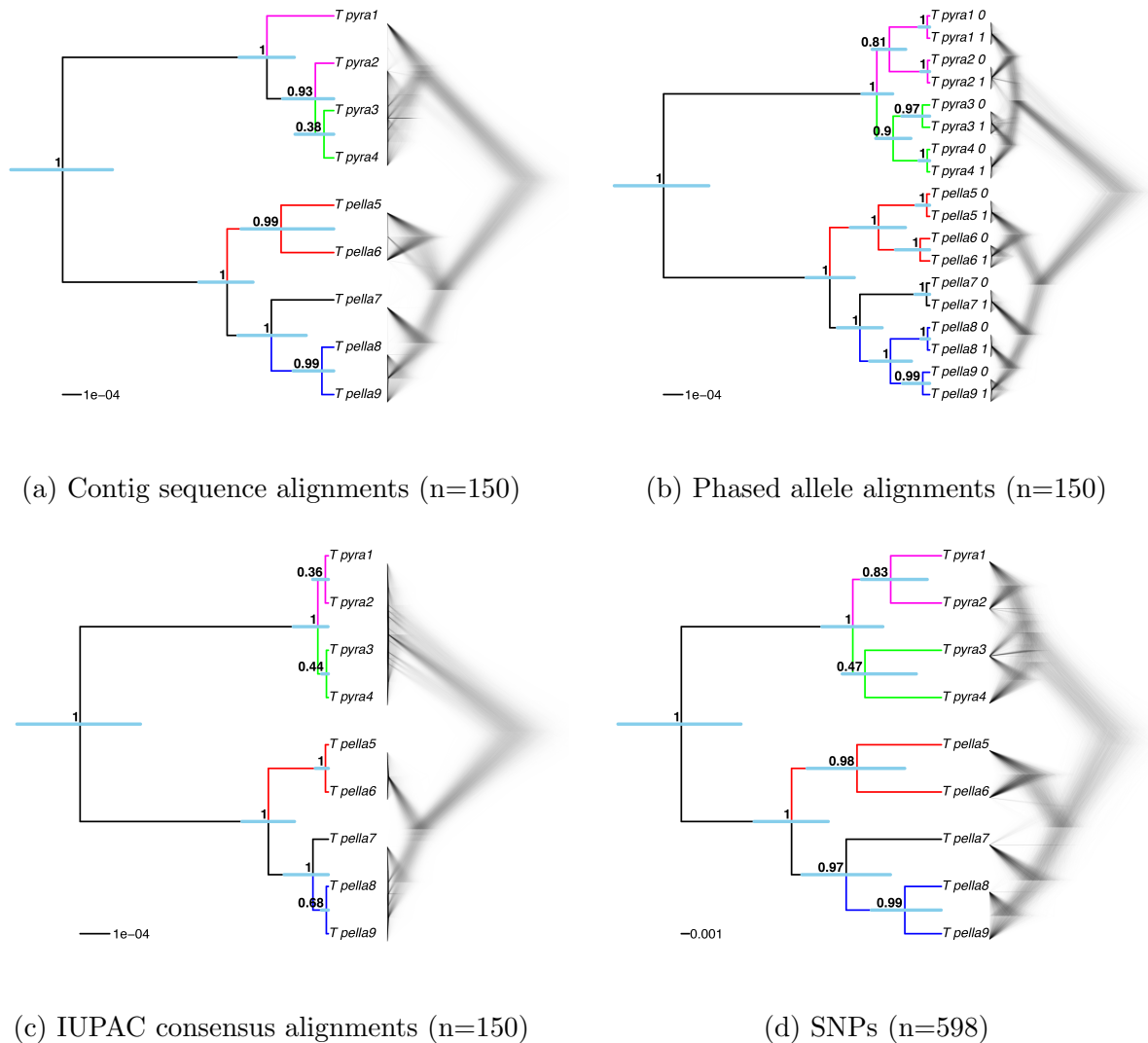(c) IUPAC consensus alignments (n=150)



(d) SNPs (n=598)

Figure 5: 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 SNP data (1 SNP per locus if present, n=598). Shown are the maximum clade credibility tree (node values = PP, error-bars = 95% HPD of divergence times) and a plot of the complete posterior species tree distribution (excluding burn-in).

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

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

476 Conversely, the phased allele data (Fig. 6b) and the IUPAC consensus data (Fig. 6c)

477 correctly delimited the five coalescent species from the simulation input tree (Fig. 6g). The

478 STACEY results showed the same pattern in all ten simulation replicates (Supplementary

479 Fig. S7 available on Dryad).

480 *Accuracy of divergence time estimation.*— For all four sequence-based analyses (Figs. 6a

481 to 6d) the average substitution rate across all loci was set to '1'. Under these settings, we

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

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

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

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

486 represented by green line in Fig. 7). This was the case for all nodes in the species tree,

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

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

489 any bias throughout ten simulation replicates (Supplementary Fig. S8 available on Dryad).

490 This contrasts with the contig MSAs and the chimeric allele MSAs that consistently

491 overestimated the height of all nodes and the IUPAC consensus MSAs which consistently

492 underestimated the height of all nodes (Fig. 7, Supplementary Fig. S8).

## *Additional Analyses*

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

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

496 et al. 2007), which can be found in online Appendix 2 and Supplementary Figs. S9 to S11

497 (available on Dryad).

(a) Contig sequence alignments (n=150)

(b) Phased allele alignments (n=150)

(c) IUPAC consensus alignments (n=150)

(d) Chimeric allele alignments (n=150)

(e) SNPs reduced (n=150)  (f) SNPs complete (n=820)
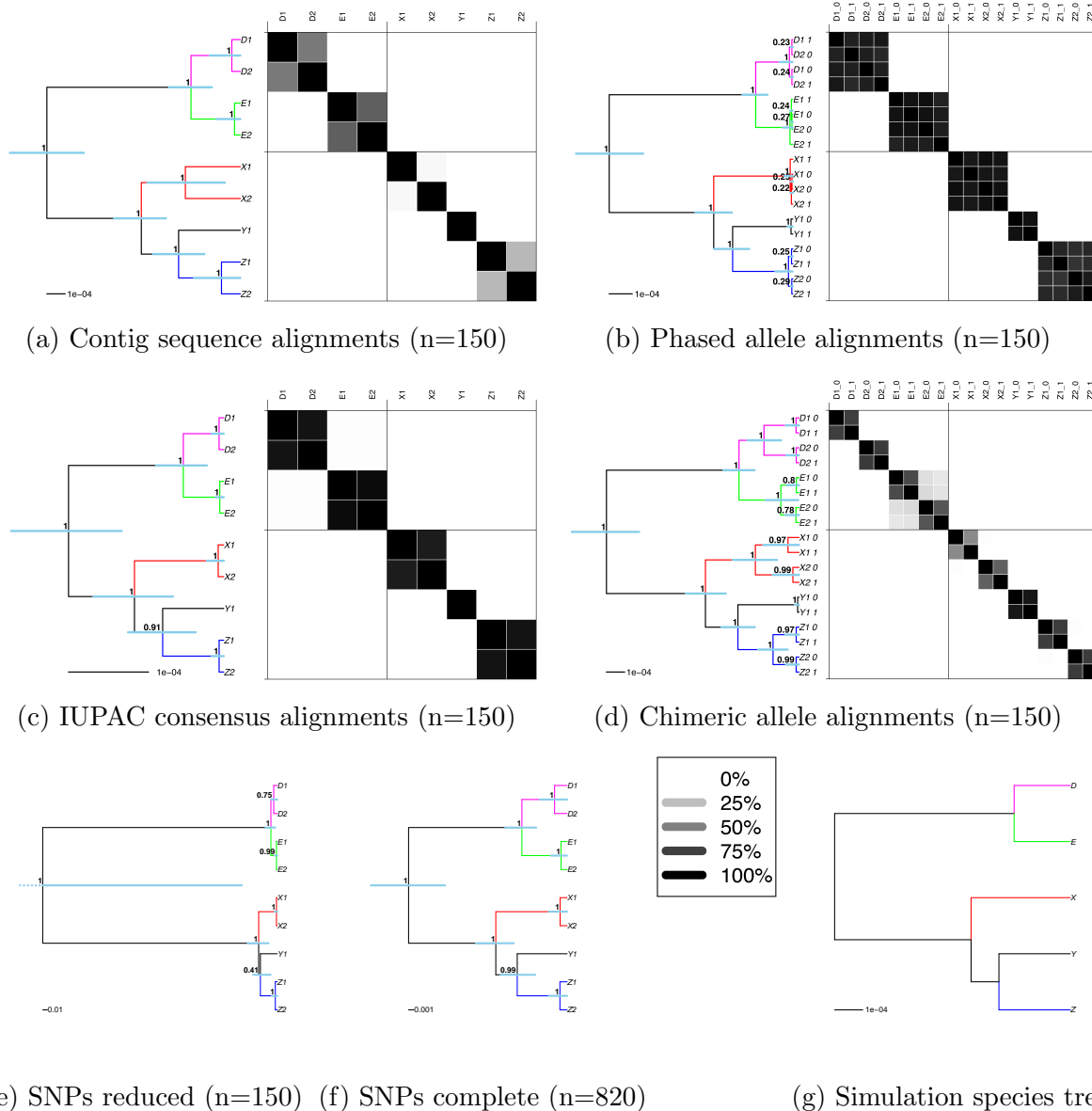
(g) Simulation species tree

Figure 6: 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 (see sub-figure captions). 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.
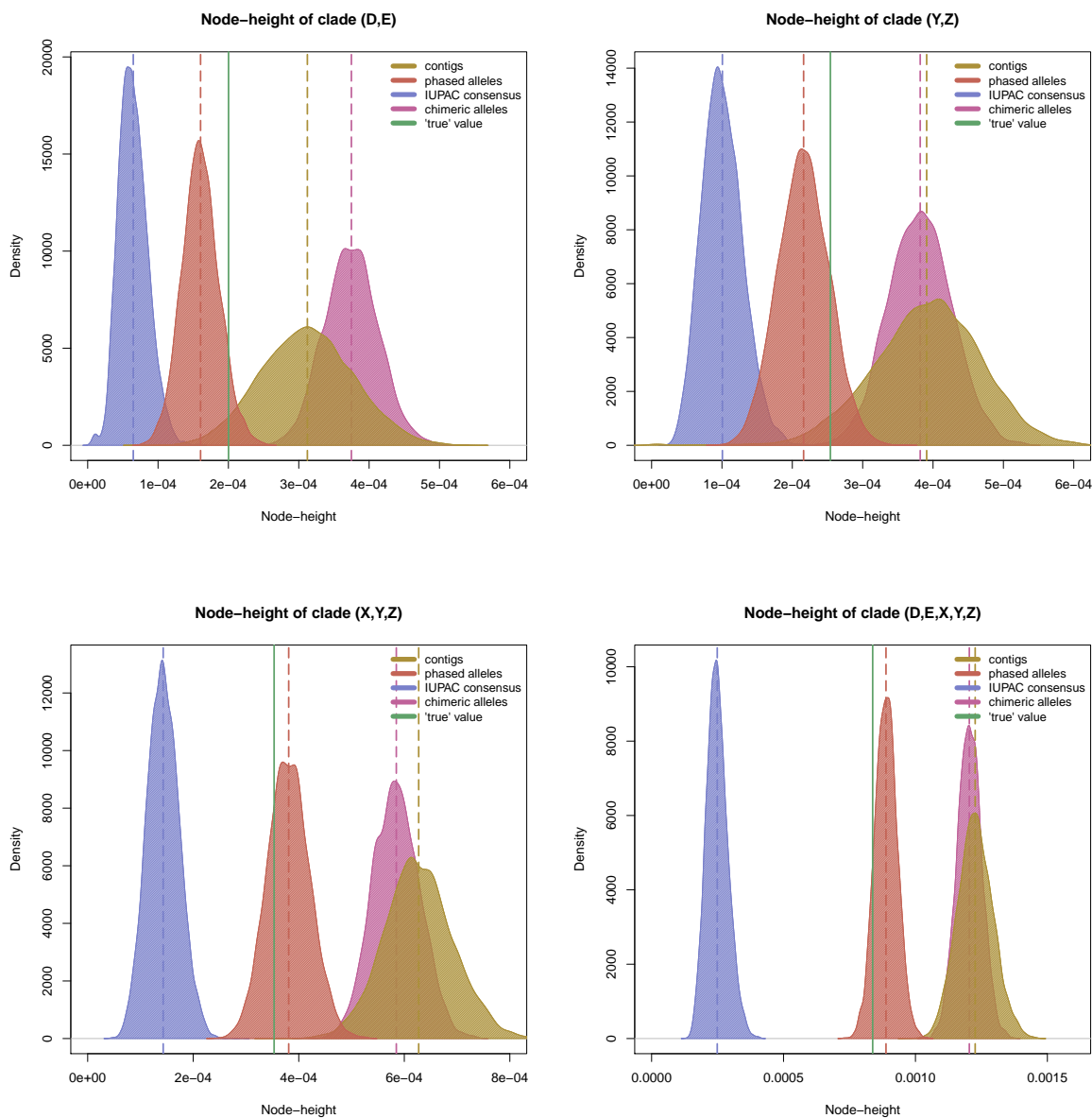
Figure 7: Posterior distributions of divergence times estimated with STACEY. Each sub-figure shows several density plots of node-height estimates for a different node in the STACEY species tree (see sub-figure titles). The four density plots in each sub-figure are approximated from all node height values in the posterior sample (excl. 10% burnin), as estimated by STACEY for the four different data processing schemes tested in this study: contig sequences (yellow), phased allele sequences (red), IUPAC consensus sequences (blue) and chimeric allele sequences (pink). The dotted lines show the means of these posterior distributions. The solid green 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.

# DISCUSSION

## *Allele Phasing is the Preferred Data Processing Scheme*

We tested whether phylogenetic inference improves by phasing sequence capture data into allele sequences, in comparison to the standard workflow of analyzing contig sequences (Faircloth et al. 2012; McCormack et al. 2012; Smith et al. 2014; Faircloth 2015). The answer is yes. We find that phased allele data outperform contig sequences in terms of species delimitation (Fig. 6) and the estimation of divergence times (Fig. 7). Contig sequence MSAs lead to a consistent overestimation of divergence times (Fig. 7), which in turn lead to an overestimation of the number of coalescent species in our simulated data (Fig. 6a). These results support earlier work by Lischer et al. (2014), who concluded that consensus sequences introduce a bias towards older node heights.

Besides the qualitative 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 underly 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 (Fig. 4). This shortcoming of the most common assemblers (e.g. ABYSS, Trinity and Velvet) is due to the fact that they were designed to assemble haploid sequences and 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

522  data also leads to a doubling of the effective sample size, since two sequences are compiled

523  for a diploid individual, in contrast to the single sequence per individual that is recovered

524  when taking the contig approach. Our results demonstrate how these sequences can be

525  treated as independent samples from a population by using the assignment-free

526  BirthDeathCollapse model as implemented in STACEY. Because STACEY requires no *a*

527  *priori* assignment of sequences to taxon, this avoids a violation of the MSC that would

528  occur when analyzing allele sequences as separate taxa in *BEAST, because *BEAST

529  assumes each taxon constitutes a separate coalescent species. Third, sequence capture

530  datasets such as UCEs are optimal for allele phasing because they contain high read

531  coverage collected across short genomic intervals that are optimal for read-connectivity

532  based phasing. The workflow developed in this study is now fully integrated into the

533  PHYLUCE pipeline, making allele phasing and SNP extraction for sequence capture data

534  easily available to a broad user group. Given these advantages of allele sequences over

535  contig sequences and given the easy availability of the processing workflow, we recommend

536  that allele phasing be considered as a standard practice for future sequence capture studies.

## *Phasing of Heterozygous Sites Matters*

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

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

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

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

542  phased allele sequences for estimating the species tree topology (Fig. 6). However, IUPAC

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

544  nodes in the species tree (Fig. 7). Our results contrast with those of (Lischer et al. 2014),

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

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

547 tree inference programs that we used. Lischer et al. (2014) applied a Neighbour Joining

548 (NJ) tree algorithm as implemented in the software PHYLIP (Felsenstein 2005) that treats

549 two sequences containing the same ambiguity codes as identical. In effect, the approach

550 used by Lischer et al. (2014) did not truly investigate the effect of IUPAC ambiguity codes

551 on phylogenetic estimates but rather the effect of removing heterozygous sites. Our

552 approach of analyzing IUPAC consensus sequences under the MSC in STACEY, on the

553 other hand, properly integrates these IUPAC ambiguity codes into the calculation of the

554 gene tree likelihoods. Thus, we conclude that IUPAC ambiguity codes introduce a bias

555 towards younger divergence times, even when properly integrating IUPAC ambiguities into

556 the phylogenetic model. The underlying cause of this discrepancy should be further

557 investigated in future studies.

558 We also tested whether the improved performance of phased allele sequence data

559 may merely be an effect of doubling the number of sequences in the MSAs, since we are

560 producing two allele sequences for each individual rather than one contig sequence.

561 Therefore, we generated a dataset of chimeric allele sequences that contains the same

562 number of sequences as the phased allele data, but we randomly shuffled all heterozygous

563 positions within an individual between the two allele sequences. As with the contig data,

564 the chimeric allele data led to an overestimation of the number of coalescent species (Fig.

565 6d) and to a biased estimation towards older divergence times (Fig. 7). The fact that

566 contig sequences and chimeric allele sequences produce very similar results in our analyses

567 is not surprising, because contigs, themselves, represent chimeric consensus sequences of

568 the variation found at a locus within an individual. The similarity of the results between

569 contig MSAs and chimeric allele MSAs also shows that the number of sequences being

570 analyzed does not affect our topology, species delimitation and divergence time estimates

571 (Figs. 6 and 7).

572 Based on these findings, we conclude that proper phasing of heterozygous positions

573  is clearly preferable to the alternative of coding heterozygous sites as IUPAC ambiguity

574  codes, particularly when the estimation of divergence times is of interest. Further, allele

575  sequences are theoretically more appropriate input for coalescent models and should be the

576  preferred data type input to these models. The scalability of this approach to larger sample

577  sizes and the applicability of our results to studies of older divergences are questions that

578  should be investigated in future studies.

## *UCEs as source for SNP data*

580  Due to the size (number of loci) of many sequence capture datasets, it is often unfeasible to

581  analyze all MSAs jointly in one MSC analysis (Smith et al. 2014; Manthey et al. 2016)

582  because of computational limitations. For all sequence-based MSC analyses in this study,

583  we reduced the UCE dataset from 820 loci to 150 loci in order to reach convergence of the

584  MCMC within a reasonable time frame (three to four days, single core on a Mac Pro, Late

585  2013, 3.5 GHz 6-Core Intel Xeon E5 processor). However, a viable approach to data

586  reduction, while keeping the multilocus information of all loci, is to analyze only a single

587  polymorphic position per MSA using SNAPP (Bryant et al. 2012). In our study, this

588  approach produces the correct species tree topology and also estimated the relative

589  node-heights correctly (Fig. 6f). However, SNAPP can only estimate relative and not

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

591  analyses Figs. 6a to 6d that deliver absolute divergence time estimates.

592       Sequence capture datasets such as UCEs provide a suitable data source to extract

593  both full sequence alignments and SNP datasets of sufficient size for robust species tree

594  estimation. Even though sequence capture data are not commonly thought of as a source

595  of SNPs, they can, in many cases, be preferable to other sequencing techniques, such as

596  RAD sequencing, for producing SNP data. This is because sequence capture data yields a

597  sizable, complete SNP matrix (SNPs recovered for all individuals), due to targeted

598 sequence enrichment. In this study, the complete matrix of unlinked SNPs in the empirical

599 data consisted of 598 positions, which were present and sufficiently supported (>three

600 high-quality reads per haplotype) in all taxa. Particularly when evolutionary distances

601 between individuals are large, RAD sequencing and other restriction-site based sequencing

602 techniques are not expected to yield many loci shared by all individuals, whereas UCE

603 data are less sensitive to large evolutionary distances (Harvey et al. 2016). In these cases,

604 the size of the complete SNP matrix resulting from UCE data can exceed that resulting

605 from RAD sequencing. Additionally, UCE data provide hundreds to thousands of full

606 sequence MSAs as well as the complete mitochondrial genome as a byproduct of the

607 sequence enrichment. The mitochondrial genome provides an excellent marker for

608 estimating absolute divergence times (Fig. 3), based on substitution rates of mitochondrial

609 markers which are known for birds (Lerner et al. 2011), and thus remains a valuable source

610 of phylogenetic information.

611      In this study, we present and make available a new SNP calling pipeline for

612 sequence capture data. In contrast to other SNP calling software such as GATK (McKenna

613 et al. 2010) that uses BAM files, our approach uses full sequence MSAs as input (see Fig.

614 2), in order to identify and extract sites in the alignments that show variation between any

615 user-defined group of sequences. Although our SNP calling script can be applied to any

616 type of sequence alignments (i.e. allele or contig sequence alignments), we recommend

617 using SNPs extracted from phased allele alignments for phylogenetic analyses, because they

618 represent the true heterozygous information. The user can choose whether or not to allow

619 missing data or ambiguities in the extracted positions, whether to extract them in binary

620 format (as e.g. required by SNAPP) or as nucleotides, and if only a single SNP per locus

621 or all SNPs should be extracted. Thus our SNP calling mechanism is an easy, open-source

622 and straightforward tool to derive SNP data from any set of multiple sequence alignments.

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

625    *Topaza* that is dated at ca. 2.4 Ma in the mitochondrial tree (Fig. 3). These lineages are

626    consistent with the previously described morphospecies *T. pyra* (Gould, 1846) and *T. pella*

627    (Linnaeus, 1758) that are generally accepted in the ornithological community (Hu et al.

628    2000; del Hoyo et al. 2016a). However, the species status of *T. pyra* has been challenged by

629    some authors (Ornés-Schmitz and Schuchmann 2011; Schuchmann 1999). These authors

630    concluded that *Topaza* is a monotypic genus with *T. pyra* being a subspecies of *T. pella*,

631    which they refer to as *T. pella pyra*. Their findings are based on the analyses of plumage

632    coloration, in which they found an "east-west clinal trend of characters" (Ornés-Schmitz

633    and Schuchmann 2011). In contrast, we do not find such an east-west clinal trend in the

634    genetic data. Instead, *T. pyra* is consistently supported as a separate lineage across all

635    analyses, lending no support for the conspecificity of these two taxa (Figs. 3 and 5).

636        One aim of this study was to evaluate the genetic structure within these two

637    morphospecies, *T. pyra* and *T. pella*. The mitochondrial tree shows two divergent clades

638    within *T. pyra* (Fig. 3), but these clades are not strongly supported by the UCE data (Fig.

639    5), even though the allele sequence data are picking up a signal that possibly indicates two

640    clades are in the process of diversifying (Fig. 5b). For *T. pella*, on the other hand, we

641    consistently find the same clades throughout all multilocus MSC analyses (Fig. 5), leading

642    us to distinguish between the following populations that are congruent with previous

643    morphological subspecies descriptions:

644    *Northern T. pella population: T. pella pella.*— For the mitochondrial tree and all MSC

645    species trees, we find the northern *T. pella* samples 5 and 6 to be sister taxa with high

646    support values (98-100% PP, Figs. 3 and 5). Particularly in the mitochondrial tree (Fig.

647    3), these two samples appear as close sister taxa, separated by only very short terminal

648 branches. Their close position in the mitochondrial tree shows that, even though

649 geographically far apart, samples 5 and 6 share a relatively recent MRCA in the

650 mitochondrial genealogy, indicating some rather recent gene flow. The sampling locality of

651 sample 5 is within the range of the subspecies *T. pella pella*, which extends mainly across

652 the Guiana shield (Peters 1945; Schuchmann 1999; Hu et al. 2000; Ornés-Schmitz and

653 Schuchmann 2011). Given the sampling location of genetically related sample 6, which also

654 has been morphologically identified as *T. pella pella* (Table 1), we propose that the

655 distribution range of *T. pella pella* extends from the Guiana shield all the way south to the

656 northern Amazon River bank (see map in Fig. 3).

657 *Southern T. pella population: T. pella microrhyncha.*— In the same manner as for the

658 northern population *T. pella pella*, we also consistently find the southern *T. pella* samples

659 8 and 9 to be sister taxa (99-100% PP, Figs. 3 and 5). The sampling locations of these two

660 samples are included in the distribution range of the previously recognized subspecies *T.*

661 *pella microrhyncha*, extending from the southern bank of the Amazon River as far South as

662 Porto Velho (Brazil) at the Madeira River, close to the border to Bolivia (Peters 1945;

663 Schuchmann 1999; Ornés-Schmitz and Schuchmann 2011). This southernmost boundary of

664 *T. pella microrhyncha* is not accepted by Hu et al. (2000), who instead conclude that this

665 southernmost population belongs to *T. pella pella*. In contrast to the findings by Hu et al.,

666 our genetic data clearly support the southernmost sample 9 belonging to the same

667 population as sample 8, which was morphologically identified as *T. pella microrhyncha*.

668 This leads us to propose that the distribution range of *T. pella microrhyncha* is in fact as

669 shown in Fig. 3, in agreement with the findings by Peters (1945),Schuchmann (1999), and

670 Ornés-Schmitz and Schuchmann (2011).

671 *Estuary region of Amazon River: T. pella smaragdula.*— Our results show a mixed signal

672 concerning the phylogenetic placement of sample 7, which was collected from the southern

673 estuary region of the Amazon River and morphologically identified as *T. pella smaragdula.*

674 The sampling locality also falls into the range of the subspecies *T. pella smaragdula* (Peters

675 1945; Hu et al. 2000; Ornés-Schmitz and Schuchmann 2011), with a distribution including

676 the Amazon River estuary and extending north along the coast to French Guiana. All

677 MSC analyses of the UCE sequence and SNP data place sample 7 with high confidence

678 (97-100% PP) as sister to the southern clade *T. pella microrhyncha* (Fig. 5), whereas in

679 the mitochondrial phylogeny this sample is placed as sister to *T. pella pella* in the North.

680    The discordance between a gene tree and the species tree in a scenario such as this

681 could be the effect of incomplete lineage sorting, which is most likely if the species or clades

682 in question have diverged rather recently and if population sizes are large. Given that the

683 divergence between *T. pella pella* and *T. pella microrhyncha* appears to be considerably

684 deep based on the multilocus data (crown height of *T. pella*) see Fig. 5) and given that

685 mitochondria are generally considered to have only 25% of the population size of nuclear

686 loci, it is rather unlikely that the position of sample 7 in the mitochondrial tree is the result

687 of incomplete lineage sorting in this case. It seems more likely that the separate position of

688 sample 7 in the mitochondrial tree is the result of introgression of the mitochondrial

689 genome from *T. pella pella* into the gene pool of *T. pella smaragdula.* However, a denser

690 taxon sampling would be necessary to further evaluate the evolutionary history of this

691 particular population. The case of sample 7 highlights that the mitochondrial tree presents

692 a single gene tree phylogeny that only shows one of many genealogies and therefore must

693 not be equated with a species tree phylogeny. Hence it is important to generate multilocus

694 data for an informed inference of the species tree phylogeny.

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

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

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

At the same time, we find indicators of phylogenetic structure within species, distinguishing samples that are separated by only a small geographic distance (see e.g. samples 6 and 8). These samples are however separated by the Amazon River, which has been found to constitute a dispersal barrier for various species of birds and many other animals (Remsen and Parker 1983; Clair 2003; Hayes and Sewlal 2004; Moore et al. 2008; Fernandes et al. 2012; Ribas et al. 2012; Thom and Aleixo 2015). Even though some hummingbird species are known to disperse across large distances (Wyman et al. 2004; Russell et al. 1994), the Amazon River and its associated habitats (such as seasonally flooded forests) may be part of a complex network of factors that inhibit gene flow among populations of *Topaza* hummingbirds.

# CONCLUSIONS

In this study, we demonstrate that properly phasing allele sequences produces the most suitable dataset for phylogenetic analyses, particularly when these allele sequences are treated as independent sequences under the MSC. Contig sequences, on the other hand, which are commonly used for phylogenetic inference, lead to biases in the estimation of divergence times and may cause problems for certain types of phylogenetic analyses. Additionally, phased allele sequences provide a useful template for the extraction of SNPs, and we argue that sequence capture data can provide sizable SNP datasets that can be also used for phylogenetic analyses. Our empirical results suggest the separation of two species within the genus *Topaza*, and we further find genetic structure within both of these species, justifying the definition of separate subspecies. Based on our empirical and simulated results, we conclude that allele phasing should be considered as one "best practice" for processing sequence capture data, although the sample-size, time, and analytical limitations of this approach have not been well-established.

# Supplementary Material

Supplemental Figs. S1-S11, Supplemental Tables S1 and S2, online Appendices 1 and 2 and all scripts, data and setup-files relevant to analyses and figures in the manuscript are available from the Dryad Digital Repository:

# Availability

We integrated all scripts and documentation necessary for phasing and SNP extraction as open-source into the PHYLUCE pipeline (`http://https://github.com/faircloth-lab/phyluce/blob/working/bin/snps/`). All data processing and analyses steps executed on the data are stored in bash-scripts on our project GitHub page at `https://github.com/tobiashofmann88/topaza_uce`. We further provide a documented workflow of processing the raw reads into UCE contig alignments at `https://github.com/tobiashofmann88/UCE-data-management/wiki`.

# Acknowledgments

# Funding

\*

References

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

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

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

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

766 Clair, C. C. S. 2003. Comparative permeability of roads, rivers, and meadows to songbirds

767     in Banff national park. Conservation Biology 17:1151–1160.

768 Degnan, J. H. and N. a. Rosenberg. 2009. Gene tree discordance, phylogenetic inference

769     and the multispecies coalescent. Trends in Ecology and Evolution 24:332–340.

770 del Hoyo, J., N. Collar, G. Kirwan, and P. Boesman. 2016a. Fiery Topaz (*Topaza pyra*). *in*

771     Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,

772     D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.

773 del Hoyo, J., A. Elliott, J. Sargatal, D. Christie, and E. de Juana. 2016b. Handbook of the

774     Birds of the World Alive. Lynx Edicions, Barcelona, Spain.

775 Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics

776     with BEAUti and the BEAST 1.7. Molecular Biology and Evolution 29:1969–73.

777 Edwards, S. V., L. Liu, and D. K. Pearl. 2007. High-resolution species trees without

778     concatenation. Proceedings of the National Academy of Sciences 104:5936–5941.

779 Eriksson, J. S., J. L. Blanco-Pastor, F. Sousa, Y. J. Bertrand, and B. E. Pfeil. 2017. A

780     cryptic species produced by autopolyploidy and subsequent introgression involving

781     Medicago prostrata (Fabaceae). Molecular Phylogenetics and Evolution 107:367–381.

782 Faircloth, B. C. 2015. PHYLUCE is a software package for the analysis of conserved

783     genomic loci. Bioinformatics 32:786–788.

784 Faircloth, B. C., J. E. McCormack, N. G. Crawford, M. G. Harvey, R. T. Brumfield, and

785     T. C. Glenn. 2012. Ultraconserved elements anchor thousands of genetic markers

786     spanning multiple evolutionary timescales. Systematic Biology 61:717–26.

787  Faircloth, B. C., L. Sorenson, F. Santini, and M. E. Alfaro. 2013. A phylogenomic

788     perspective on the radiation of ray-finned fishes based upon targeted sequencing of

789     ultraconserved elements (UCEs). PLoS ONE 8:e65923.

790  Felsenstein, J. 2005. Phylip (phylogeny inference package) version 3.6. distributed by the

791     author. dep genome sci univ washington, seattle.

792  Fernandes, A. M., M. Wink, and A. Aleixo. 2012. Phylogeography of the chestnut-tailed

793     antbird (*Myrmeciza hemimelaena*) clarifies the role of rivers in Amazonian biogeography.

794     Journal of Biogeography 39:1524–1535.

795  Garrick, R. C., P. Sunnucks, and R. J. Dyer. 2010. Nuclear gene phylogeography using

796     PHASE: dealing with unresolved genotypes, lost alleles, and systematic bias in

797     parameter estimation. BMC Evolutionary Biology 10:118.

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

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

800     Biology 64:727–740.

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

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

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

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

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

806     Sequence capture versus restriction site associated DNA sequencing for shallow

807     systematics. Systematic Biology Advance Access syw036.

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

passerine birds: effects of river width, habitat and taxonomy. Journal of Biogeography 31:1809–1818.

He, D., A. Choi, K. Pipatsrisawat, A. Darwiche, and E. Eskin. 2010. Optimal algorithms for haplotype assembly from whole-genome sequence data. Bioinformatics 26:i183–i190.

Hu, D.-S., L. Joseph, and D. J. Agro. 2000. Distribution, variation, and taxonomy of *Topaza* Hummingbirds (Aves: Trochilidae). Ornitologia Neotropical 11:123–142.

Iqbal, Z., M. Caccamo, I. Turner, P. Flicek, and G. McVean. 2012. De novo assembly and genotyping of variants using colored de Bruijn graphs. Nature Genetics 44:226–232.

Jones, G. 2017. Algorithmic improvements to species delimitation and phylogeny estimation under the multispecies coalescent. Journal of Mathematical Biology 74:447–467.

Jones, G., Z. Aydin, and B. Oxelman. 2014. DISSECT: an assignment-free Bayesian discovery method for species delimitation under the multispecies coalescent. Bioinformatics 31:991–998.

Katoh, K., G. Asimenos, and H. Toh. 2009. Multiple alignment of DNA sequences with MAFFT. Methods in Molecular Biology 537:39–64.

Kolaczkowski, B. and J. W. Thornton. 2004. Performance of maximum parsimony and likelihood phylogenetics when evolution is heterogeneous. Nature 431:980–984.

Kubatko, L. S. and J. H. Degnan. 2007. Inconsistency of Phylogenetic Estimates from Concatenated Data under Coalescence. Systematic Biology 56:17–24.

Leaché, A. D., B. L. Banbury, J. Felsenstein, A. N. M. De Oca, and A. Stamatakis. 2015. Short tree, long tree, right tree, wrong tree: New acquisition bias corrections for inferring SNP phylogenies. Systematic Biology 64:1032–1047.

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

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

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

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

835      honeycreepers. Current Biology 21:1838–1844.

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

837      transform. Bioinformatics 26:589–595.

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

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

840      Bioinformatics 25:2078–9.

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

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

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

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

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

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

847      syw005.

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

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

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

851      Research 22:746–754.

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

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

854   Current Biology 24:910–916.

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

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

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

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

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

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

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

862   Advance Access syw014.

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

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

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

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

867   2014. ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics

868   30:541–548.

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

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

871   11:960–968.

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

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

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

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

876     Ornitologia Neotropical Pages 25–38.

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

878     Cambridge, Massachusetts.

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

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

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

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

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

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

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

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

887 Ribas, C. C., a. Aleixo, a. C. R. Nogueira, C. Y. Miyaki, and J. Cracraft. 2012. A

888     palaeobiogeographic model for biotic diversification within Amazonia over the past three

889     million years. Proceedings of the Royal Society B: Biological Sciences 279:681–689.

890 Russell, R. W., F. L. Carpenter, M. A. Hixon, and D. C. Paton. 1994. The impact of

891     variation in stopover habitat quality on migrant rufous hummingbirds. Conservation

892     Biology 8:483–490.

893 Schrempf, D., B. Q. Minh, N. De Maio, A. von Haeseler, and C. Kosiol. 2016. Reversible

894     polymorphism-aware phylogenetic models and their application to tree inference. Journal

895     of Theoretical Biology 407:362–370.

896  Schuchmann, K., G. Kirwan, and P. Boesman. 2016. Crimson Topaz (*Topaza pella*). *in*

897  Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,

898  D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.

899  Schuchmann, K. L. 1999. Family Trochilidae (hummingbirds). Pages 468–680 *in* Handbook

900  of the Birds of the World Alive (J. del Hoyo, A. Elliott, and J. Sargatal, eds.) volume 5

901  ed. Lynx Edicions, Barcelona, Spain.

902  Simpson, J. T., K. Wong, S. D. Jackman, J. E. Schein, S. J. M. Jones, and I. Birol. 2009.

903  ABySS: a parallel assembler for short read sequence data. Genome Research 19:1117–23.

904  Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014.

905  Target capture and massively parallel sequencing of ultraconserved elements for

906  comparative studies at shallow evolutionary time scales. Systematic Biology 63:83–95.

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

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

909  Conservation 142:2282–2292.

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

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

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

913  Molecular Phylogenetics and Evolution 82:95–110.

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

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

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

917  Current Zoology 61:854–865.

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

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

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