

2 RH: Phasing improves utility of UCEs

3 **Allele Phasing Greatly Improves the Phylogenetic
4 Utility of Ultraconserved Elements**

5 TOBIAS ANDERMANN^{1,2}, ALEXANDRE M. FERNANDES³, URBAN OLSSON^{1,2}, MATS
6 TÖPEL^{2,4}, BERNARD PFEIL^{1,2}, BENGT OXELMAN^{1,2}, ALEXANDRE ALEIXO⁵, BRANT C.
7 FAIRCLOTH⁶ 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

20 *Abstract.*— In the era of high-throughput sequencing techniques, it is becoming
21 increasingly easy and affordable to sequence big portions of the genome, also for non-model
22 organisms. But instead of sequencing the complete genome, many phylogenetic studies
23 prefer to focus their sequencing efforts on a selected set of target loci. Sets of target loci
24 are commonly enriched through techniques such as sequence capture. The advantage of
25 this approach is that a higher sequencing depth is reached for these selected loci. Higher
26 sequencing depth does not only lead to more confidence in assembling the target sequences
27 but also reveals the allelic variation within the sequenced individual. Even though this
28 allelic variation constitutes important information of great phylogenetic utility, allele
29 sequences are rarely assembled in phylogenetic studies based on sequence capture data.
30 Instead many scientists use contig sequences, resulting from de novo assembly of
31 sequencing reads, for phylogenetic analyses. Here we show how phased allele sequences can
32 be retrieved from sequence capture data and we demonstrate for the first time a full
33 integration of these allele sequences under the Multispecies Coalescent (MSC) model. We
34 find, based on simulated data, that the common approach of analyzing contig sequences
35 produces biased divergence time estimates, while analyzing allele sequences leads to
36 accurate divergence time estimations. We benchmark our allele phasing workflow with an
37 empirical dataset of Ultraconserved Elements (UCEs), generated from the South American
38 hummingbird genus *Topaza* and demonstrate how UCE allele sequences can be used to
39 infer the genetic structure within a genus that diverged during the last three million years.
40 Our empirical results support the recognition of two species in *Topaza* whose taxonomy has
41 been a matter of recent debate. The results also suggest a high rate of gene flow across
42 large distances of rainforest habitats but rare admixture across the Amazon River. We
43 conclude based on comparisons of different data processing strategies that allele phasing is
44 the most proper processing scheme for UCE data and sequence capture data in general.
45 Therefore we recommend that analyzing allele sequences (and not contig sequences) should

⁴⁶ become the “best practice” for sequence capture datasets and thus make available this
⁴⁷ bioinformatic workflow for future phylogenetic studies.

⁴⁸ (Keywords: UCE, SNP, heterozygous sites, Multispecies Coalescent, gene tree, species tree,
⁴⁹ Mitochondrial Genome, Trochilidae, Birds, Amazon)

50 Massive Parallel Sequencing (MPS) techniques enable time- and cost-efficient
51 generation of DNA sequence datasets of large fractions or even the entirety of genomes.
52 But instead of sequencing complete genomes, researchers often choose to focus sequencing
53 efforts on a set of target loci in order to achieve higher coverage and thus more reliable
54 sequencing of these target regions (Faircloth et al. 2012, 2013; Mirarab et al. 2014; Smith
55 et al. 2014; Faircloth 2015; Harvey et al. 2016; Meiklejohn et al. 2016). Such multilocus
56 datasets, containing thousands of target loci, are commonly generated through enrichment
57 techniques such as sequence capture (synonym: target enrichment, Gnirke et al. (2009))
58 prior to sequencing. Sequence capture uses custom designed bait sequences (synonym:
59 probes) that bind to specific genetic regions of interest and which in turn enables specific
60 enrichment of these loci. This ensures that the targeted regions are present in higher copy
61 number and hence will be sequenced at higher coverage by MPS techniques.

62 Higher sequencing coverage at a given locus reveals important information about
63 allelic variation at that locus, which is particularly interesting for understanding
64 demographic processes such as migration, population sizes and the amount of gene flow
65 between and within populations. For this reason the information lying within allelic
66 variation is commonly applied in population genetic studies, often in form of single
67 nucleotide polymorphisms (SNPs). However, it is rarely acknowledged that allelic
68 sequences harbor huge potential for phylogenetic studies as well, in particular concerning
69 the estimation of gene trees and species trees (Potts et al. 2014), and the estimation of
70 divergence times (Lischer et al. 2014). Further, ignoring allelic variation can introduce
71 biases in downstream analyses (Lischer et al. 2014; Garrick et al. 2010).

72 Despite their excellent utility, allele sequence data are rarely compiled for
73 phylogenetic studies based on MPS data, with only few exceptions (Lischer et al. 2014;
74 Potts et al. 2014; Schrempf et al. 2016). This is surprising since most MPS generated
75 datasets, in particular those resulting from sequence capture enrichment, are very suitable

76 for compiling allelic sequences for any given locus in the dataset. Instead most phylogenetic
77 studies generate and analyze contig sequences for each locus, thereby masking the allelic
78 information that lies within the MPS data (see Results, Fig. 4). Part of the reason for the
79 common neglection of such useful information is that the computational tools that help
80 with the compilation of allelic information require bioinformatic knowledge that exceeds
81 the experience of the average biologist working with MPS sequence data. Up to this point
82 and to the authors best knowledge, none of the MPS processing pipelines contained user
83 friendly solutions for the compilation of allele sequences. We therefore integrated an allele
84 phasing function in the open-source PHYLUCE pipeline (Faircloth 2015).

85 Additionally to challenges with compiling allelic sequences, also the proper analysis
86 of allelic information in phylogenetic studies remains challenging and is an intensively
87 discussed topic (Garrick et al. 2010; Lischer et al. 2014; Potts et al. 2014; Schrempf et al.
88 2016). Various researchers have chosen different approaches to include this information into
89 phylogenetic methods. One such approach has been to code heterozygous sites using
90 IUPAC ambiguity codes and to include these as additional characters into existing
91 substitution models for gene tree and species tree inference (Potts et al. 2014; Schrempf
92 et al. 2016). While these studies demonstrated that integrating this additional information
93 leads to a higher accuracy in phylogenetic inference, Lischer et al. (2014) found that coding
94 heterozygous sites as IUPAC ambiguity codes in phylogenetic models biases the results
95 toward older divergence time estimates. Instead Lischer et al. (2014) introduce their
96 method of repeated random haplotype sampling (RRHS) in which they repeatedly
97 concatenate AFLP sequences (Amplified fragment length polymorphisms), choosing a
98 random haplotype for any given locus in each concatenation replicate. In their approach
99 they then analyze thousands of concatenation replicates separately for phylogenetic tree
100 estimation and summarize the results between replicates, thereby integrating the allelic
101 information in form of uncertainty intervals. However there are two important shortcomings

102 of the approach presented by Lischer et al. (2014): 1. concatenating unlinked loci and in
103 particular allele sequences from unlinked loci in a random manner is known to produce
104 incorrect topologies (Degnan and Rosenberg 2009) often with false confidence (Edwards
105 et al. 2007; Kolaczkowski and Thornton 2004; Kubatko and Degnan 2007; Mossel and
106 Vigoda 2005), which is not accounted for when doing so repeatedly and summarizing the
107 resulting trees and 2. running thousands of tree estimation replicates based on extensive
108 amounts of sequence data does result in unfeasibly long computation times, particularly for
109 Markov-Chain Monte Carlo (MCMC) based softwares such as MrBayes or BEAST.

110 Here we introduce the bioinformatic assembly of allele sequences for UCE data and
111 demonstrate a full integration of these sequences under the Multispecies Coalescent (MSC)
112 model for both empirical and simulated data. In our approach we treat each allelic
113 sequence of an individual at a given locus as an independent sample from the population.
114 We analyze these sequences using the species tree and delimitation software STACEY
115 (Jones et al. 2014; Jones 2017), which does not require a priori clade- or
116 species-assignments and therefore allows us to effectively treat each allele sequence as an
117 individual sample. We demonstrate the utility of our approach by resolving the shallow
118 genetic structure (<1 Ma) within the two recognized morphospecies of the South American
119 hummingbird genus *Topaza*. The underlying sequence data of 2,386 Ultraconserved
120 Elements (UCEs, see Faircloth et al. (2012)) was generated through sequence capture using
121 the bait set for Tetrapods (see <http://ultraconserved.org>) and subsequent Illumina
122 sequencing. We find, based on simulations, that allele sequences yield more accurate results
123 in terms of species tree estimation and species delimitation than the classic contig sequence
124 approach, which ignores heterozygous information. Our simulation results further show
125 that proper phasing of allele sequences outperforms alternative approaches of coding
126 heterozygous information, such as sequences containing IUPAC ambiguity codes or Single
127 Nucleotide Polymorphisms (SNPs). We conclude that phasing sequence capture data can

128 be critical for correct species delimitation and phylogeny estimation, particularly in
129 recently diverged groups, and that analyses using phased alleles should therefore become
130 the preferred choice for sequence capture datasets.

131 **MATERIALS AND METHODS**

132 *Study System*

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

150

Sequence Data Generation

151 We extracted DNA from 10 vouchered hummingbirds (9 *Topaza*, one *Florisuga*, see
152 Table 1) from muscle tissue using the Qiagen DNeasy Blood and Tissue Kit according to
153 the manufacturer's instructions (Qiagen GmbH, Hilden, Germany). These samples cover
154 most of the genus' total geographic range (Fig. 1) and all morphologically recognized
155 intraspecific taxa (Schuchmann et al. 2016; del Hoyo et al. 2016a). All samples were
156 sonicated with a Covaris S220 sonication device in order to break the genomic DNA into
157 shorter fragments with a targeted fragment length of 800 bp. Paired-end, size-selected
158 DNA libraries were prepared for sequencing on the Illumina MiSeq platform, using the
159 magnetic-bead based NEXTflexTM Rapid DNA-Seq Kit (Bioo Scientific Corporation,
160 Austin, TX, USA), following the user's manual (v14.02).

161 We used the "Tetrapods-UCE-2.5Kv1" bait set (`uce-2.5k-probes.fasta`),
162 consisting of 2,560 baits (each 120 bp), targeting 2,386 UCEs, as described by Faircloth
163 et al. (2012). The bait sequences were downloaded from <http://ultraconserved.org> and
164 synthesized by MYcroarray (Biodiscovery LLC, Ann Arbor, MI, USA). Sequence
165 enrichment was performed using a MYbaits kit according to the enclosed user manual
166 (v1.3.8). The enriched libraries were then sequenced using 250 bp, paired-end sequencing
167 on an Illumina MiSeq machine (Illumina Inc., San Diego, CA, USA). Library preparation,
168 sequence enrichment and sequencing were performed by Sahlgrenska Genomics Core
169 Facility in Gothenburg, Sweden.

170

Mitochondrial Genome

171 Even though no baits targeting mitochondrial loci were used during sequence
172 capture, we found that as many as 4.5% of all sequence reads were of mitochondrial origin.
173 This was sufficient to assemble the complete mitochondrial genome for all samples, which

Table 1: Information of the specimens sequenced. Subspecies identifications based on morphological characters, which are hard to distinguish in this group. 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	<i>Topaza pyra</i>	<i>amaruni</i>	INPA A1106	-0.044167	-66.94944
2	<i>T. pyra</i>	<i>pyra</i>	MPEG 62475	-1.559444	-65.88006
3	<i>T. pyra</i>	<i>pyra</i>	MPEG 62474	-4.083889	-60.66050
4	<i>T. pyra</i>	<i>pyra</i>	MPEG 52721	-7.350000	-73.66667
5	<i>T. pella</i>	NA	USNM 586322	7.220000	-60.29000
6	<i>T. pella</i>	<i>pella</i>	INPA A3319	-1.927900	-59.41600
7	<i>T. pella</i>	<i>smaragdula</i>	MPEG 61688	-1.950000	-51.60000
8	<i>T. pella</i>	<i>microrhyncha</i>	MPEG 65603	-5.352417	-57.47500
9	<i>T. pella</i>	NA	INPA A6233	-9.028550	-64.24231
10	<i>Florisuga fusca</i>	NA	MPEG 70697	-15.15972	-39.04500

¹⁷⁴ we analyzed in BEAST (Drummond et al. 2012) in order to infer a dated mitochondrial
¹⁷⁵ genealogy. Dating priors included clock-rate priors for three mitochondrial genes, estimated
¹⁷⁶ for honeycreepers by Lerner et al. (2011) and node-age priors within the genus *Topaza* that
¹⁷⁷ were estimated by McGuire et al. (2014). A detailed description of the assembly and
¹⁷⁸ analysis of the mitochondrial genome data can be found in online Appendix 1
¹⁷⁹ (Supplemental Material available on Dryad).

¹⁸⁰ *UCE Data Processing*

¹⁸¹ For this study we generated five different types of datasets, which we analyzed under the
¹⁸² MSC. These five datasets, representing different coding schemes of heterozygous
¹⁸³ information, are listed and described in the following.

¹⁸⁴ 1. *UCE contig alignments*.— In this study we use the term contig in a simplified manner to
¹⁸⁵ refer to the consensus sequence derived from the actual contig (de novo assembly of
¹⁸⁶ overlapping sequence reads). Many de novo assemblers generate these contig consensus

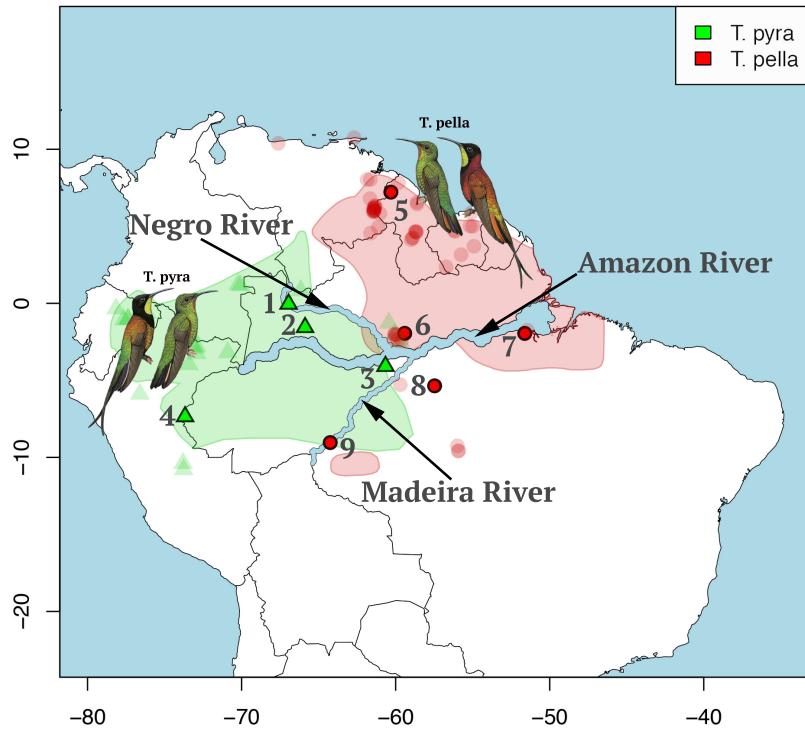


Figure 1: Map of northern South America, showing sampling locations (numbered symbols) of *Topaza* specimens. Numbers represent sample IDs (Table 1). The colored polygons show the distribution range of the two morphospecies (*T. pyra* = green, *T. pella* = red) as estimated by BirdLife International (<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).

187 sequences by ignoring variants at heterozygous positions and instead treat them like
 188 sequencing errors by choosing the more probable variant, while discarding the other (Iqbal
 189 et al. 2012). Since contig sequences generated in this manner are commonly used in
 190 phylogenetic analyses of MPS datasets (e.g. Faircloth et al. (2012); Smith et al. (2014);
 191 Faircloth (2015)), we generated contig MSAs for all UCE loci in order to test the accuracy
 192 of the phylogenetic estimation of this approach.

193 To create multiple sequence alignments (MSAs) from UCE contig data, we followed
 194 the suggested workflow from the PHYLUCE documentation

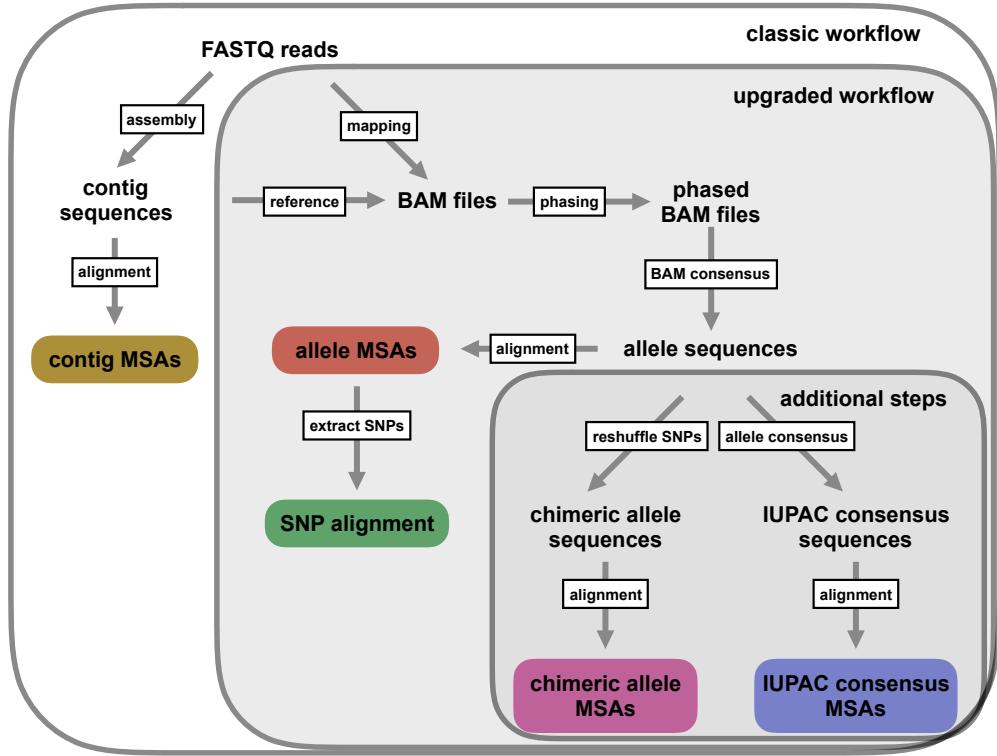


Figure 2: Depiction of the workflow developed here. Colored boxes represent different types of multiple sequence alignments (MSAs) used for phylogenetic inference in this study. In addition to the standard UCE workflow of generating contig MSAs (Faircloth et al. 2012; Smith et al. 2014; Faircloth 2015), 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 steps of data processing were executed in this study in order to test different codings of heterozygous positions.

195 (<http://phyluce.readthedocs.io/en/latest/>). We applied the PHYLUCE default
 196 settings unless otherwise stated. First we quality-filtered and cleaned raw Illumina
 197 reads of adapter contamination with Trimmomatic (Bolger et al. 2014), which is
 198 implemented in the PHYLUCE function `illumiprocessor`. The reads were then
 199 assembled into contigs using the software ABYSS (Simpson et al. 2009) as implemented in
 200 the PHYLUCE pipeline. In order to identify contigs representing UCE loci, all assembled
 201 contigs were mapped against the UCE reference sequences from the bait sequence file

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

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

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

205 MAFFT (Katoh et al. 2009).

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

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

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

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

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

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

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

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

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

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

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

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

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

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

220 since the read coverage across any given UCE locus typically is highest in the center and

221 decreases toward the ends, which makes it possible to phase throughout the complete locus

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

223 on opposite ends of the locus, the fragment length we selected for in this study (800 bp) in

224 combination with paired-end sequencing, enables the phasing process to bridge across the

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

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

227 BAM file into a single sequence and exported the two resulting allele sequences for each
228 sample in FASTA format. In the next step we aligned the allele sequences between all
229 samples, separately for each UCE locus, using MAFFT (Fig. 2). We integrated this
230 complete workflow into the UCE processing software PHYLUCE (Faircloth 2015) with
231 slight alterations, one of which is the use of the open-source mapping program bwa (Li and
232 Durbin 2010) in place of CLC-mapper.

233 *3. UCE IUPAC consensus sequence alignments.*— We generated an additional set of
234 alignments by merging the two allele sequences for each individual into one consensus
235 sequence with heterozygous sites coded as IUPAC ambiguity codes. We used this dataset
236 in order to test if our allele phasing approach leads to an improved phylogenetic inference
237 in comparison to the IUPAC consensus approach as applied in other studies, in which
238 heterozygous positions are coded as IUPAC ambiguity codes in a consensus sequence for
239 each locus and individual (Potts et al. 2014; Schrempf et al. 2016).

240 *4. UCE chimeric allele alignments.*— In this study we also investigate whether correct
241 phasing of heterozygous sites is essential or if equal results are achieved by randomly
242 placing variants on either one allele sequence. For this purpose we generated another
243 dataset with chimeric allele sequence alignments. These alignments were generated by
244 applying a custom python script to the phased allele sequence alignments, which randomly
245 shuffles the two variants at each polymorphic position between the two allele sequences for
246 each individual. This process leads in many cases to an incorrect combination of variants
247 on each allele sequence, thereby creating chimeric allele sequences. The resulting
248 alignments contain the same number of sequences as the phased allele alignments (two
249 sequences per individual), while the contig alignments and the IUPAC consensus
250 alignments contain only half as many sequences (one sequence per individual).

251 5. *UCE SNP alignment*.— A common approach of analyzing heterozygous information is
252 to reduce the sequence information to only a single variant SNP per locus. This
253 data-reduction approach is often chosen because multi locus datasets of the size generated
254 in this study can be incompatible for full sequence based phylogeny estimation with
255 Bayesian MSC methods, due to unfeasibly long computational times. Instead of full
256 sequence alignments, alignments of unlinked SNPs can be used to infer species trees and
257 species demographics under the MSC model with the BEAST2 package SNAPP (Bryant
258 et al. 2012), a program specifically designed for such data. However, extracting and
259 filtering SNPs from BAM files with existing software (such as the Genome Analysis Toolkit
260 (GATK), McKenna et al. (2010)) and converting these into a SNAPP compatible format
261 can be cumbersome, as SNAPP requires positions with exactly two different states, coded
262 in the following manner: individual homozygous for the original state = “0”, heterozygous
263 = “1”, and homozygous for the derived state = “2”.

264 We therefore developed a python function, which extracts biallelic SNPs directly
265 from allele sequence MSAs. Extracting SNPs from MSAs in this manner is a
266 straightforward and simple way to generate a SNP dataset compatible with SNAPP, and
267 does not require re-visiting the BAM files. Our SNP extraction script outputs a SNP
268 alignment file that is formatted for input in SNAPP. We applied this function to the phased
269 UCE alignments and extracted one position per alignment (to ensure unlinked SNPs) that
270 had exactly two different states among all *Topaza* samples, not allowing for positions with
271 missing data or ambiguities. This resulted in a SNP dataset of 598 unlinked loci. We
272 integrated the SNP extraction from allele sequence MSAs into the PHYLUCE pipeline.

273 *Generation of Simulated UCE Data*

274 In order to assess the accuracy of the phylogenetic inferences resulting from different data
275 processing approaches, we generated simulated UCE data for all five processing schemes as

276 applied to the empirical *Topaza* data. For each of the five processing schemes we generated
277 and analyzed ten independent simulation replicates.

278 1. *Simulated allele alignments*.— We estimated parameters of species divergence times and
279 population sizes under the MSC model (Rannala and Yang 2003) from the empirical UCE
280 allele alignments using the Bayesian MCMC program BPP v3.1 (Yang 2015). We applied
281 the A00 model implemented in the software, which estimates divergence times and
282 population sizes when the species tree topology is provided. As input, we chose the species
283 tree topology that resulted from the analysis of the allele sequence data in STACEY,
284 assigning the *Topaza* samples to five separate taxa (corresponding to colored clades in Fig.
285 6b). An initial BPP analysis did not converge in reasonable computational time, a problem
286 that has been reported before for UCE datasets of several hundred loci (Giarla and
287 Esselstyn 2015). We therefore split the 820 UCE alignments randomly into 10 subsets of
288 equal size ($n=82$) and analyzed these separately with identical settings in BPP. The
289 MCMC was set for 150,000 generations (burn-in 50,000), sampling every 10 generations.
290 We summarized the estimates for population sizes and divergence times between all 10
291 individual runs. We then applied the mean values of these estimates to the species tree
292 topology, by using the estimated divergence times as branch lengths and estimated
293 population sizes as node values, resulting in the species tree in Fig. 6g. This tree was used
294 to simulate sequence alignments with the MCcoal simulator, which is integrated into BPP.
295 Equivalent to the empirical data, we simulated sequence data for five taxa (D, E, X, Y, and
296 Z) and one outgroup taxon (F, not shown in Fig. 6g). In the simulations these taxa were
297 simulated as true species under the MSC model. In order to mimic the empirical allele
298 data, we simulated four individuals for species ‘D’ (equivalent to two allele sequences for 2
299 samples), four for species ‘E’, four for species ‘X’, two for species ‘Y’ (two allele sequences
300 for one sample), four for species ‘Z’, and two for the outgroup species ‘F’. In this manner

301 we simulated 820 UCE allele MSAs of 848 bp length (= average alignment length of
302 empirical allele alignments).

303 *2. Simulated contig alignments.*— In order to simulated UCE contig MSAs, we merged the
304 sequences within each coalescent species in pairs of two (equivalent to pairs of allele
305 sequences). Each pair of sequences was joined into one sequence by randomly picking one
306 of the two variants in the case of a polymorphism (heterozygous site). This we used to
307 represent the UCE contig approach of previous studies (Faircloth et al. 2012; McCormack
308 et al. 2012; Smith et al. 2014; Faircloth 2015) that rely on UCE contig data as generated
309 by assemblers such as ABYSS, Velvet or Trinity, which pick only one of the two variants at
310 a heterozygous site. Just as in the empirical assembler approach, our simulation approach
311 may generate chimeric contig sequences.

312 *3. Simulated IUPAC consensus alignments.*— Next, we generated IUPAC consensus MSAs
313 in the same manner as we generated the simulated contig MSAs in the previous step, with
314 the exception that heterozygous sites were coded with IUPAC ambiguity codes (instead of
315 randomly picking one of the two variants).

316 *4. Simulated chimeric allele alignments.*— We generated chimeric allele sequence MSAs
317 from the simulated allele MSAs, by randomly shuffling the heterozygous sites between each
318 pair of sequences (same pairs as in the previous two steps).

319 *5. Simulated SNP alignment.*— Finally, we extracted two different SNP datasets from the
320 simulated phased allele MSAs. The first SNP dataset (SNPs complete) was extracted in
321 the same manner as described above for the empirical data (one SNP per locus) which
322 resulted in a total alignment length of 820 SNPs for the simulated data. We extracted an
323 additional SNP dataset (SNPs reduced) from only the subset of the 150 simulated allele
324 alignments that were used for the sequence-based MSC analyses (see next section below).

325 The resulting SNP dataset of 150 SNPs was used to compare the phylogenetic inference
326 resulting from SNP data versus that resulting from full sequence data, if the same number
327 of loci is being analyzed.

328 *MSC Analyses of Empirical and Simulated UCE Data*

329 *Sequence-based tree estimation.*— In order to jointly infer gene trees and species trees, we
330 analyzed each of the generated sets of MSAs (processing schemes 1-4 for empirical and
331 simulated) under the MSC model, using the DISSECT method (Jones et al. 2014)
332 implemented in STACEY (Jones 2017). STACEY is available as a BEAST2 (Bouckaert
333 et al. 2014) package. STACEY allows *BEAST analyses without prior taxonomic
334 assignments, searching the tree space while simultaneously collapsing very shallow clades in
335 the species tree (controlled by the parameter collapseHeight). This collapsing avoids a
336 common violation of the MSC model that occurs when samples belonging to the same
337 coalescent species are assigned to separate taxa in *BEAST. This feature makes STACEY
338 very suitable for analyzing allele sequences, as they don't have to be constrained to
339 belonging to the same taxon and hence can be treated as independent samples from a
340 population. STACEY runs with the usual *BEAST operators, but integrates out the
341 population size parameter and has more efficient species tree change operators that
342 decrease the time until convergence significantly. In order to reach even faster convergence,
343 we reduced the number of loci for this analysis by selecting the 150 allele MSAs with the
344 most parsimony informative sites. This selection was made for both the empirical and the
345 simulated allele MSAs. The same loci were selected for all other processing schemes.

346 We estimated the most appropriate substitution model for each of the 150 loci with
347 jModeltest (Supplementary Table S1 available on Dryad). We then applied the resulting
348 best substitution model for each locus selected by the BIC to the alignments in BEAUTI

349 v2.4.4. We did not apply any taxon assignments, thereby treating every sequence as a
350 separate taxon. We chose a strict clock and the STACEY-specific BirthDeathCollapse
351 model as species tree prior, choosing a value of 1e-5 for the collapseHeight parameter.
352 Other settings were: bdcGrowthRate = log normal (M=4.6, S=1.5); collapseWeight = beta
353 (alpha=2, beta=2); popPriorScale = log normal (M=-7, S=2); relativeDeathRate = beta
354 (alpha=1.0, beta=1.0). Additionally we set the clock rate for each locus to a log normal
355 distribution with M=0 and S=1. For the IUPAC consensus data, we enabled the
356 processing of ambiguous sites by adding `useAmbiguities="true"` to the gene tree
357 likelihood priors for all loci in the STACEY XML file. All analyses were run for
358 1,000,000,000 MCMC generations or until convergence (ESS values >200), logging every
359 20,000 generations. Convergence was assessed using Tracer v1.6 (Rambaut et al. 2013). We
360 then summarized the posterior tree distribution into one maximum clade credibility tree
361 with TreeAnnotator v2.4.4 discarding the first 10% of trees as burn-in. For the simulated
362 data, the posterior species tree distributions of each analysis were analyzed with the
363 program SpeciesDelimitationAnalyser (part of the STACEY distribution). This program
364 produces a similarity matrix that contains the posterior probabilities of belonging to the
365 same cluster for each pair of sequences. This analysis was run with a collapseHeight value
366 of 1e-5 (identical to the collapseHeight used in the STACEY analysis), while discarding the
367 first 10% of trees as burn-in.

368 *SNP-based tree estimation.*— In order to estimate the species tree phylogeny from the
369 extracted SNP data, we analyzed the empirical and simulated SNP data in SNAPP. We
370 did not apply any prior clade assignments to the samples in the SNP alignment (each
371 sample was assigned as its own taxon). Coalescent rate and mutation rates were set to be
372 estimated based on the input data. We chose a Yule species tree model with default
373 settings ($\lambda = 0.00765$). We ran the analysis for 10,000,000 generations, sampling trees and

374 other parameters from the posterior every 1,000 generations. Differently to STACEY,
375 SNAPP assumes correct assignments of all sequences to coalescent species. Using the
376 simulated SNP data we therefore tested how our approach of assigning every individual as
377 its own coalescent species affects the resulting phylogenetic inference. We did so by
378 running a separate analysis for both simulated SNP datasets (complete and reduced) with
379 correct species assignments (assignments as in Fig. 6g).

380 **RESULTS**

381 *Mitochondrial Tree (BEAST)*

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

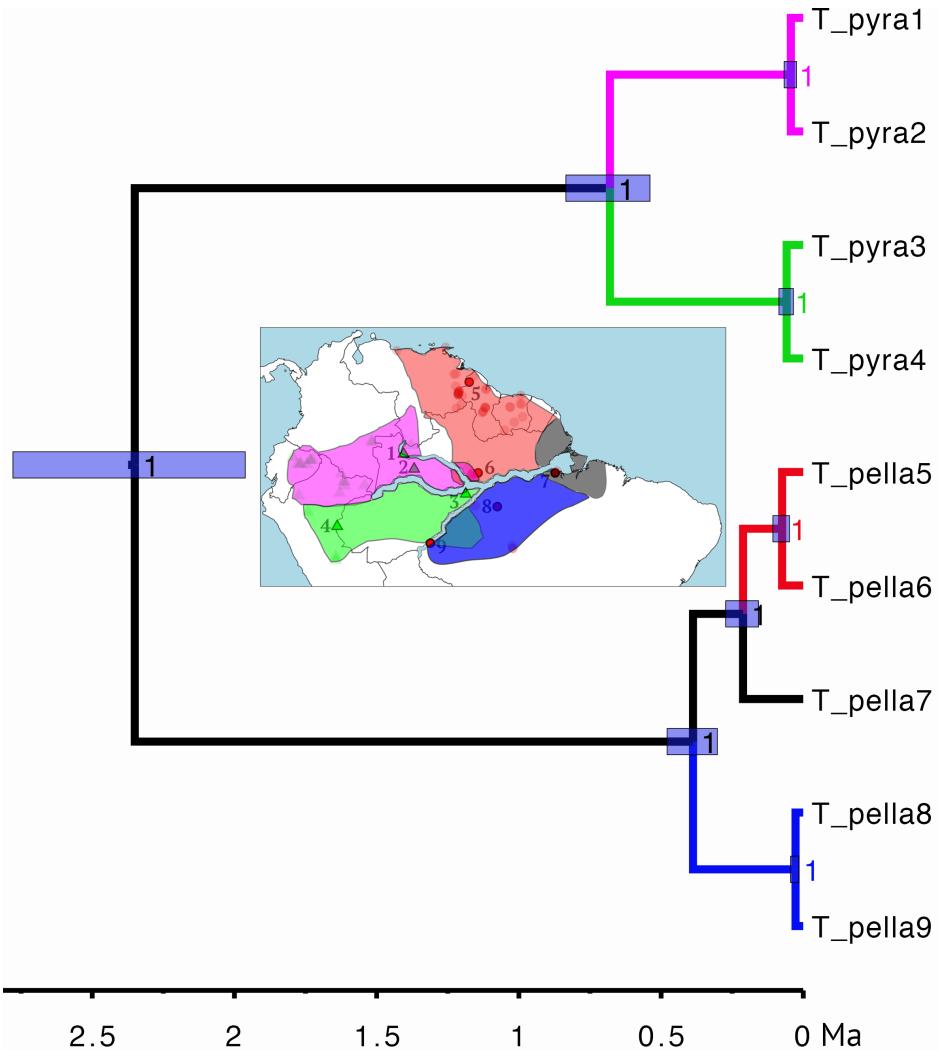


Figure 3: Mitochondrial phylogeny estimated from complete mitochondrial genomes in BEAST. Node support values represent PP. The blue bars at nodes represent the 95% HPD of divergence times. Scale axis shows time units in millions of years. The map in the center shows the potential ranges of the clades 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.

396

UCE Summary Statistics

397 *Alignment statistics.*— In the following we use the term polymorphic sites for those
 398 positions within a MSA alignment of a given locus, for which we find at least two different
 399 states among the sequences for all samples. This does not require any of the samples being

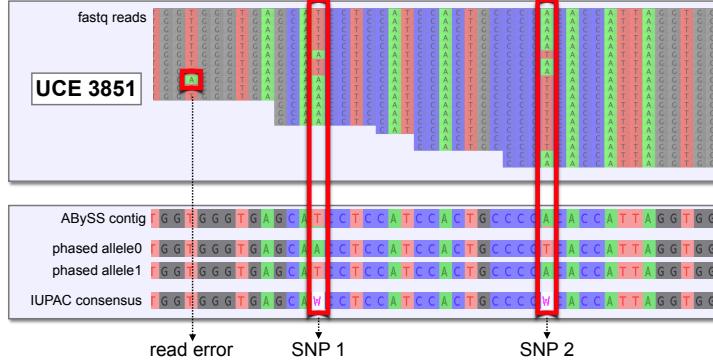
400 heterozygous for the given position, since we do not search for SNPs on a per sample basis
401 but rather for SNPs within the genus *Topaza* (for the following statistics we are excluding
402 the outgroup). In this manner we find that the empirical UCE contig alignments have an
403 average of 2.8 polymorphic sites per locus (average alignment length = 870 bp). In
404 contrast, phasing the empirical UCE data leads to 4.5 polymorphic sites per locus and an
405 average alignment length of 848 bp for the UCE allele alignments, representing a 60%
406 increase in polymorphic sites per locus. This increase of polymorphic sites is attributable
407 to the fact that many variants get lost during contig assembly, since ABYSS (and other
408 tested contig assemblers, namely Trinity and Velvet) eliminate one of the two variants at a
409 heterozygous locus. The reduced alignment length of the allele alignments in comparison to
410 the contig alignments is due to conservative alignment clipping thresholds, which lead to a
411 clipping of the alignment ends, if less than 50% of sequences are present. Since the allele
412 phasing algorithm divides the FASTQ reads into two allele bins and since a nucleotide is
413 only called if it is supported by at least three high-quality FASTQ reads, we loose some of
414 the nucleotide calls at areas of low read coverage (mostly at the ends of a locus) in the
415 allele sequences in comparison to the contig sequences. Hence the alignment clipping
416 mechanism reduces the alignment length of allele sequence MSAs in average more than for
417 contig sequence MSA.

418 More information about the distribution of lengths and variable sites within the
419 empirical UCE data can be found in the Supplementary Figs. S1 and S2 available on
420 Dryad. The simulated contig MSAs have an average of 3.2 polymorphic sites per locus,
421 after excluding the outgroup, resulting to a clade, which is the equivalent to the *Topaza*
422 clade in the empirical data. The simulated allele MSAs, on the other hand, contain an
423 average of 5.4 polymorphic sites (69% increase). An overview of parsimony informative
424 sites, variable sites and length of each alignment (simulated and empirical data) can be
425 found in Supplementary Table S2 available on Dryad.

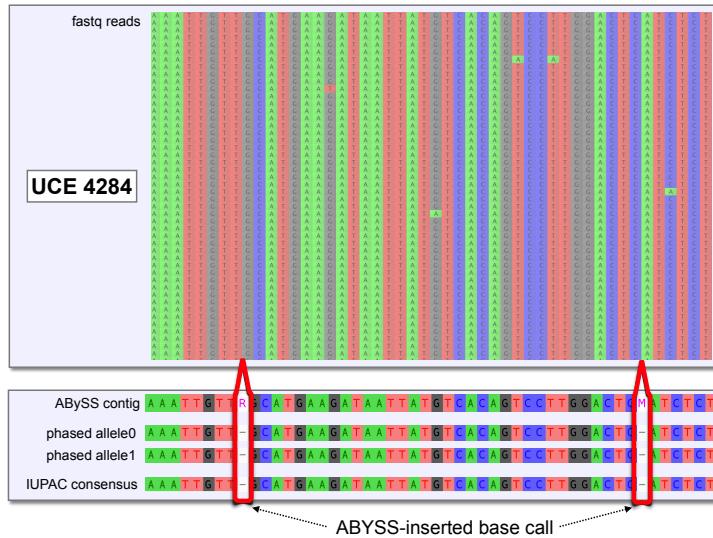
426 *ABYSS does not detect heterozygous sites.*— Since ABYSS occasionally produces sequences
427 containing IUPAC ambiguity codes, possibly representing heterozygous sites, we checked
428 exemplary for one sample (sample 5, *T. pella*) whether we find the same sites to be
429 heterozygous in the phased allele sequences. The results are striking, since not a single
430 heterozygous site within sample 5 that was detected in our allele sequence was coded as a
431 IUPAC ambiguity code in the ABYSS contigs (e.g. Fig 4a). While our phasing approach
432 revealed 343 heterozygous UCE loci with a total of 728 SNPs in sample 5, only 26 UCE
433 loci contained IUPAC ambiguity codes (degenerate bases) in the contigs resulting from the
434 ABYSS de novo assembly. For all other loci, ABYSS assumes complete homozygosity, as
435 the resulting contig sequences are free of ambiguity codes. Within the 26 loci with IUPAC
436 ambiguity codes, ABYSS introduced 473 degenerate bases, most of which constitute blocks
437 of N's. Effectively all of these ambiguous positions are in places of extremely low FASTQ
438 read coverage (<2 reads per haplotype), with the exception of 6 positions that are covered
439 by >2 reads per haplotype. However, even those 6 positions don't represent true
440 heterozygous sites within sample 5, which becomes apparent when looking at the phased
441 allele sequences or the FASTQ reads at those loci (e.g. Fig. 4b).

442 *MSC Results of Empirical UCE Data*

443 The MSC species tree results for all tested processing schemes of the empirical UCE data
444 (contig sequences, allele sequences, IUPAC consensus sequences, chimeric allele sequences
445 and SNPs) converge on a similar species tree topology, yet with no consistent topology
446 being inferred within *T. pyra* (Fig. 5 and Supplementary Fig. S3 available on Dryad). All
447 analyses strongly support the monophyly of both *T. pyra* and *T. pella* with 100% PP. In
448 all MSC analyses we also see strongly supported genetic structure within *T. pella* ($\geq 97\%$
449 PP), separating the northern samples (5 and 6) from the southern ones (7, 8 and 9).
450 Additionally, within the shallow southern *T. pella* clade, all datasets, with exception 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 but not coded as IUPAC ambiguity in the 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 carries IUPAC ambiguity codes in the ABYSS contig sequence, as was the case for 26 UCE loci in our data. The ambiguity calls at these positions are not supported by the FASTQ reads but are inserted by ABYSS at random positions and are thus not correct. Our phased allele sequences on the other hand represent the FASTQ reads correctly and do not call this position as heterozygous. This was the case for all 26 loci in our data with ABYSS-inserted IUPAC ambiguity codes.

451 IUPAC consensus data (Fig. 5c), strongly support a genetic distinction ($\geq 99\%$ PP)
452 between sample 7 from the Amazon River delta and the other southern *T. pella* samples (8
453 and 9). The deep split between northern and southern samples within *T. pyra* on the other
454 hand, which we find in the mitochondrial tree (Fig. 3), is not very well supported by the
455 multilocus MSC analyses. However, the analysis of the allele dataset returns a phylogenetic
456 signal, possibly tracking a genetic divergence between these two clades, but their
457 monophyly is not very strongly supported (Fig. 5b).

458 *MSC Results of Simulated Data*

459 *Species tree topology.*— For the simulated data we analyzed six different datasets under the
460 MSC model: contig sequence MSAs (n=150, STACEY), allele sequence MSAs (n=150,
461 STACEY), IUPAC consensus MSAs (n=150, STACEY), chimeric allele MSAs (n=150,
462 STACEY), reduced SNP data (n=150, SNAPP), and the complete SNP dataset (n=820,
463 SNAPP). All resulting species trees (Figs. 6a to 6f) correctly return the topology of the
464 species tree that was used to simulate the data (Fig. 6g). All central nodes in the species
465 trees are supported by $\geq 90\%$ PP in all analyses, with the exception of the species tree
466 resulting from the reduced SNP dataset, which shows very weak support for two nodes and
467 has a large uncertainty interval around the root-height (Fig. 6e). However, these
468 shortcomings disappear when adding more (unlinked) SNPs to the dataset (Fig. 6f). The
469 full SNP dataset (n=820) produces the correct species tree topology with high node
470 support consistently throughout ten independently simulated datasets (Supplementary Fig.
471 S4 available on Dryad). The SNAPP species tree topology appears in our case to be
472 unaffected by the chosen clade assignments model; while we allowed every sequence to be
473 its own taxon in Figs. 6e and 6f we ran additional analyses in which we applied the correct
474 species assignment (Fig. 6g) to both SNP datasets (reduced and complete SNP data),

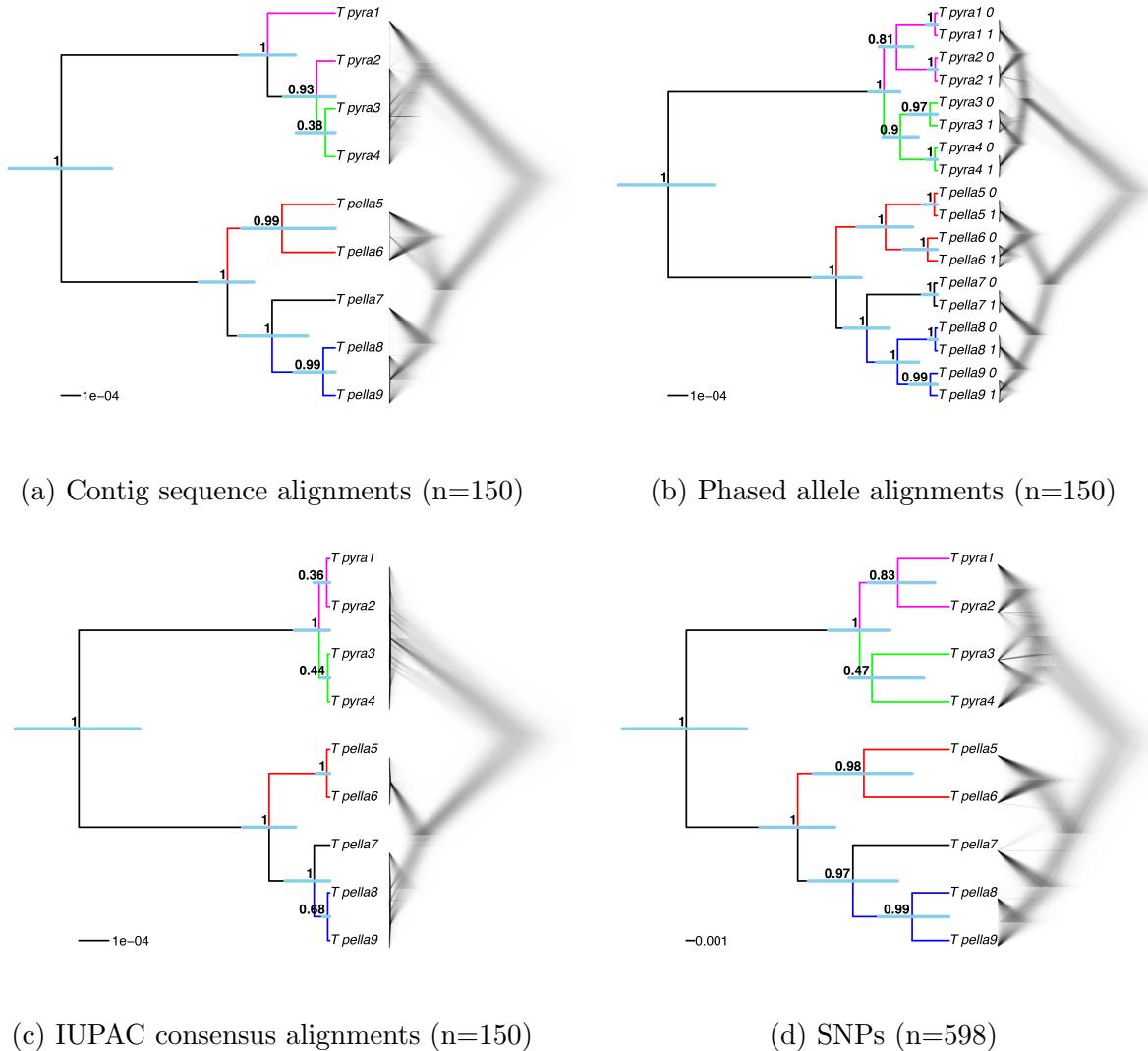


Figure 5: MSC species trees for the empirical *Topaza* data, based on four different data types used in this study: contig sequence MSAs, phased allele sequence MSAs, IUPAC consensus sequence MSAs and SNP data. (a) STACEY species tree from UCE contig alignments ($n=150$), (b) STACEY species tree from 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).

475 which returned the same tree topology (Supplementary Figs. S5 and S6 available on
476 Dryad).

477 *Species delimitation*.— While the inferred species tree topology is consistent between all
478 four sequence-based MSC analyses (Figs. 6a to 6d), the inferred node heights on the other
479 hand vary greatly between the species trees resulting from the different data processing
480 schemes. For both, the contig sequence data (Fig. 6a) and the chimeric allele data (Fig.
481 6d), the node heights within the five simulated species (D,E,X,Y,Z) are estimated too high,
482 which leads to an overestimation of the number of coalescent species in the dataset (see
483 similarity matrices). The phased allele data (Fig. 6b) and the IUPAC consensus data (Fig.
484 6c) on the other hand correctly estimate and delimit the five coalescent species from the
485 simulation input tree (Fig. 6g). The STACEY results show the same pattern in all ten
486 simulation replicates (Supplementary Fig. S7 available on Dryad)

487 *Accuracy of divergence time estimation*.— For all four sequence-based analyses (Figs. 6a
488 to 6d) the substitution rate was set to ‘1’ for one random locus, and substitution rates for
489 all other loci were estimated relative to this value. Under these settings, we expect the
490 absolute values of the sequence-based analyses to return the node height values of the
491 simulation input tree, which used substitution rates scaled in the same manner. The
492 phased allele MSAs produce the most accurate estimation of divergence times out of all
493 tested datasets (see proximity of estimates to simulation input value, represented by green
494 line in Fig. 7). This is the case for all nodes in the species tree, namely (D,E), (Y,Z),
495 (X,(Y,Z)), and ((D,E)(X,(Y,Z))). While the phased allele data lead to a slight
496 underestimation of divergence times for the very young clades (D,E) and (Y,Z), the heights
497 of the deeper nodes on the other hand are very accurately estimated by these data (Fig. 7).
498 This is in contrast to the contig MSAs and the chimeric allele MSAs, which consistently
499 overestimate the height of all nodes, while the IUPAC consensus MSAs lead to a consistent

500 underestimation of the height of all nodes (Fig. 7. These biases are present across all
501 simulation replicates while the phased allele data produces the best estimation of
502 divergence times in 88% of the cases (Supplementary Fig. S8 available on Dryad).

503 *Additional Analyses*

504 We ran additional analyses of the contig and the phased allele MSAs for both the empirical
505 and simulated data using a summary coalescent approach as implemented in MP-EST (Yu
506 et al. 2007), which can be found in online Appendix 2 and Supplementary Figs. S9 to S11
507 (available on Dryad).

508 DISCUSSION

509 *Allele Phasing is the Preferred Data Processing Scheme*

510 In this study we tested whether phylogenetic inference improves by phasing sequence
511 capture data into allele sequences, in comparison to the standard workflow of analyzing
512 contig sequences (Faircloth et al. 2012; McCormack et al. 2012; Smith et al. 2014; Faircloth
513 2015). The answer is yes. We find that phased allele data outperform contig sequences in
514 terms of estimation of divergence times (Fig. 7) and species delimitation (Fig. 6). Contig
515 sequence MSAs lead to a consistent overestimation of divergence times (Fig. 7), which in
516 turn lead to an overestimation of the number of coalescent species in our simulated data
517 (Fig. 6a. This constitutes further support to the finding of Lischer et al. (2014), who
518 conclude that consensus sequences introduce a bias toward older node heights. The phased
519 allele MSAs on the other hand accurately estimate divergence times (Fig. 7) and the
520 number of coalescent species in our simulated data (Fig. 6b).

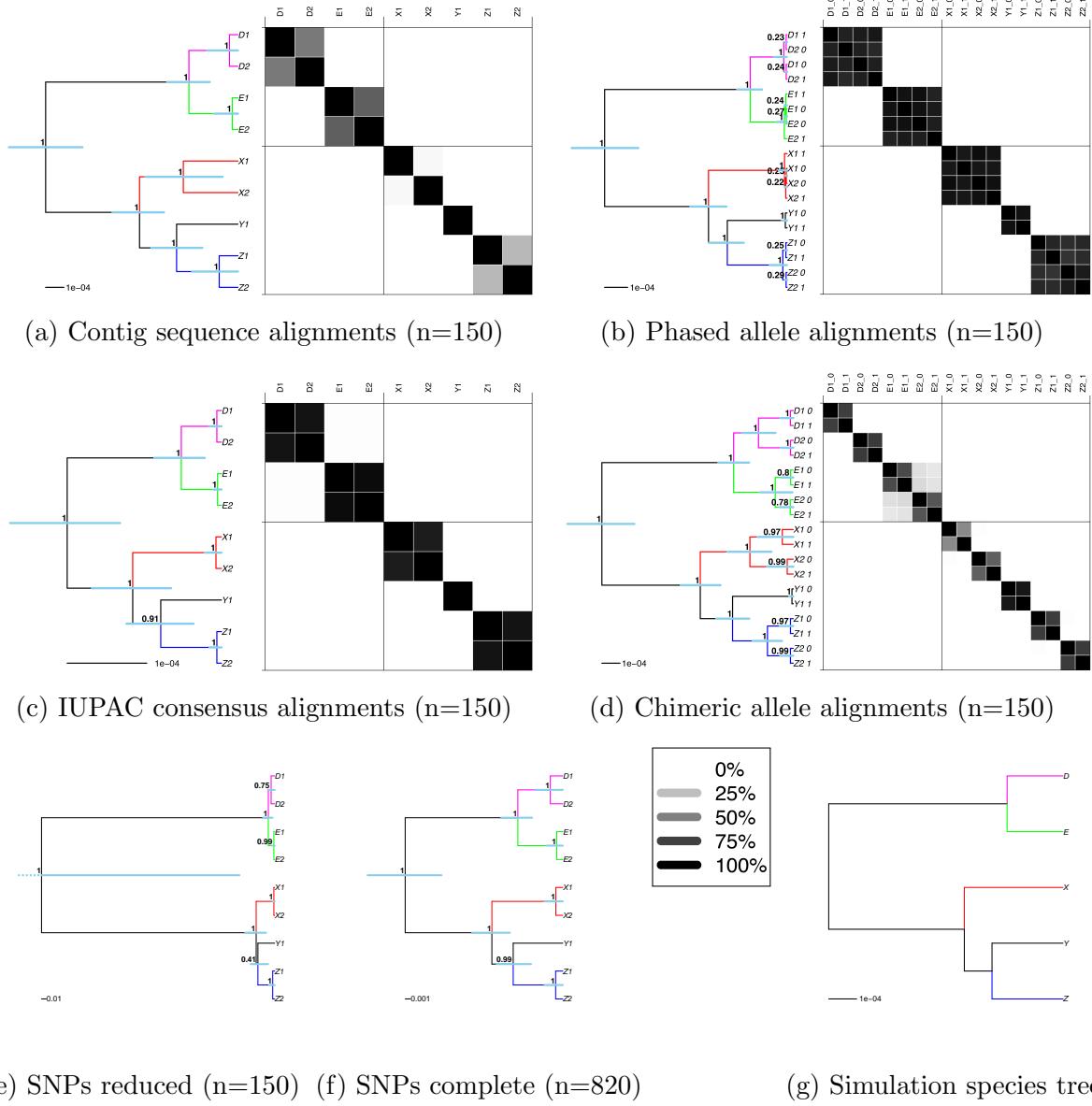


Figure 6: MSC species tree results for different data processing schemes of simulated data. (a) to (d) show the STACEY results of the four different types of MSAs analyzed in this study (see sub-figure captions). Displayed in these panels are the maximum clade credibility trees and the similarity matrices depicting the posterior probability of two samples belonging to the same clade, as calculated with SpeciesDelimitationAnalyser. Dark panels depict a high pairwise similarity, while 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 was simulated in this study. Node support values in PP, blue bars representing 95% HPD confidence intervals.

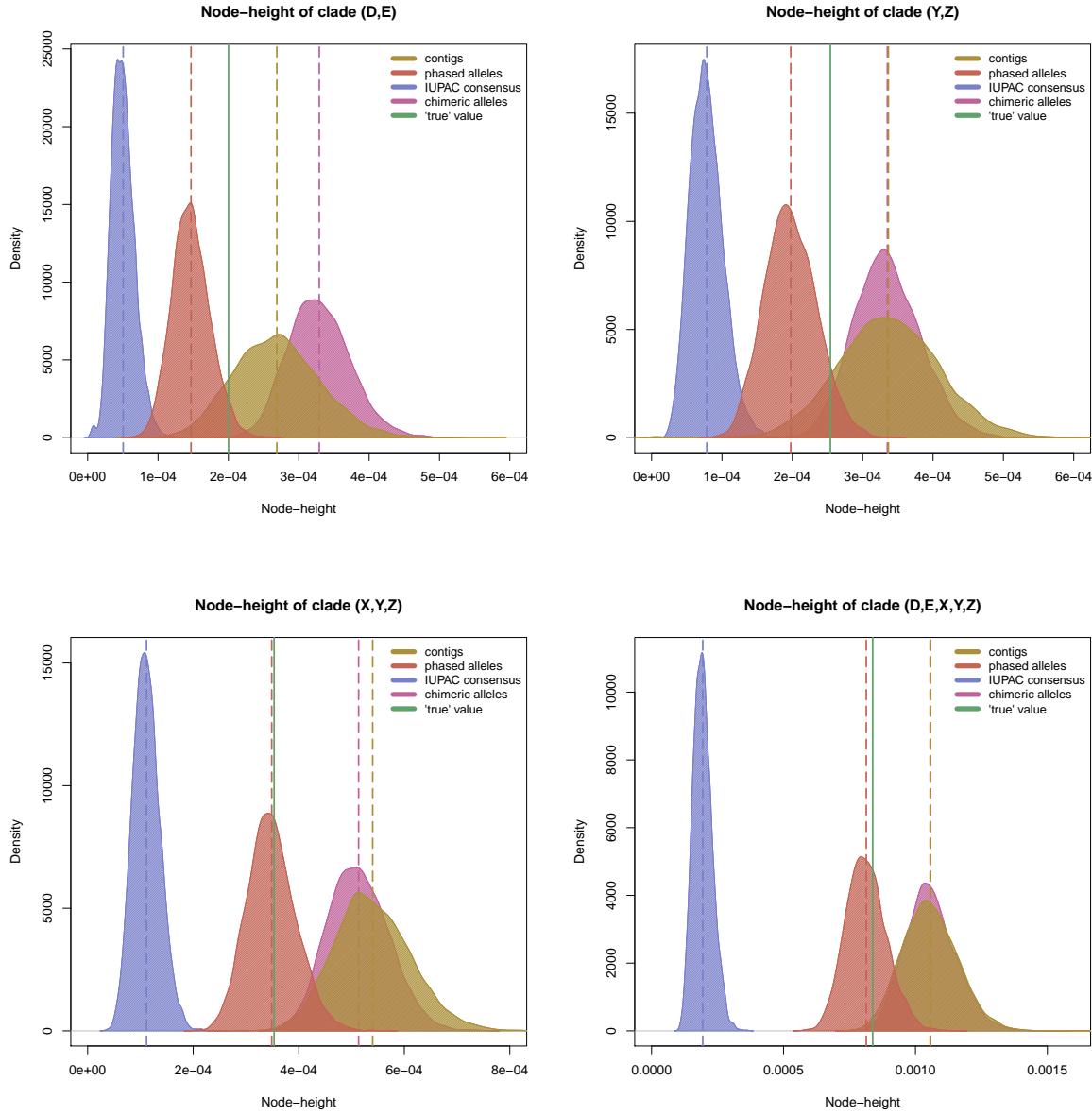


Figure 7: Allele sequences return highest accuracy of node-height estimates. Each sub-figure shows several density plots of node-height estimates for a different node in the STACEY species tree (see sub-figure titles). The four density plots in each sub-figure represent the posterior distribution of node heights (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.

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

524 First and foremost allele sequences represent the smallest evolutionary unit on
525 which selection and other evolutionary processes act. Therefore the coalescent models that
526 underly our phylogenetic methods, including the MSC model Degnan and Rosenberg
527 (2009), have been developed for allele sequences. Contig sequences on the other hand
528 represent an artificial and possibly chimeric sequence construct which arises from merging
529 all read variation at a given locus into a single sequence. This process masks heterozygous
530 information by eliminating one of the two variants at a heterozygous site (Fig. 4). This
531 shortcoming of the most common assemblers (e.g. ABYSS, Trinity and Velvet) is due to
532 the fact that they were designed to assemble haploid sequences and are not optimized for
533 heterozygous sequences or genomes (Bodily et al. 2015). This makes it obvious why it is
534 preferable to apply allele sequences and not contig sequences for inferring phylogenies.

535 Second, not only are allele sequences the more appropriate data type but phasing
536 sequence capture data also leads to a doubling of the effective sample size, since two
537 sequences are compiled for a diploid individual, in contrast to only a single sequence per
538 individual in the contig approach. We demonstrate in this study how these sequences can
539 effectively be treated as independent samples from a population by using the
540 assignment-free BirthDeathCollapse model as implemented in STACEY. Since STACEY
541 requires no a priori assignment of sequences to taxa, this avoids a violation of the MSC,
542 which would occur when analyzing allele sequences as separate taxa in the regular
543 *BEAST, since *BEAST assumes each taxon to constitute a separate coalescent species.

544 Third, sequence capture datasets such as UCEs are perfectly suitable for allele
545 phasing due to their increased read coverage of specific regions which are typically in a
546 suitable size range for read-connectivity based phasing (several hundred to few thousand

547 base pairs). The workflow developed in this study is now fully integrated into the
548 PHYLUCE pipeline, making allele phasing and SNP extraction for sequence capture data
549 easily available to a broad user group.

550 Given the mentioned advantages of allele sequences over contig sequences and given
551 the easy availability of the processing workflow, we recommend that allele phasing should
552 become the standard practice for future sequence capture studies.

553 *Phasing of Heterozygous Sites Matters*

554 Several studies have been accounting for heterozygosity by inserting IUPAC ambiguity
555 codes into their sequences at heterozygous sites Potts et al. (2014); Schrempf et al. (2016),
556 rather than phasing heterozygous positions and producing separate allele sequences. Here
557 we directly compare these two approaches and find that the IUPAC consensus sequences
558 perform equally well to the phased allele sequences in terms of species tree topology and
559 species delimitation (Fig. 6). However, the IUPAC consensus sequence data led to a
560 consistent underestimation of the divergence times of all node in the species tree (Fig. 7).
561 This is the opposite pattern found by (Lischer et al. 2014), who reported an overestimation
562 of divergence times for alignments containing IUPAC ambiguity codes. These opposite
563 findings may be caused by the different tree inference programs that were applied between
564 our study and (Lischer et al. 2014). In their study (Lischer et al. 2014) applied a
565 Neighbour Joining (NJ) tree algorithm as implemented in the software PHYLIP Felsenstein
566 (2005), which treats two sequences containing the same ambiguity codes as identical. In
567 effect, (Lischer et al. 2014) did not truly investigate the effect of IUPAC ambiguity codes
568 on phylogenetic estimates but rather the effect of removing heterozygous sites. Our
569 approach of analyzing IUPAC consensus sequences under the MSC in STACEY on the
570 other hand properly integrates these IUPAC ambiguity codes into the calculation of the
571 gene tree likelihoods. Thus we conclude that IUPAC ambiguity codes introduce a bias

572 toward younger divergence times, even when properly integrating IUPAC ambiguities into
573 the phylogenetic model. COAUTHORS: ANY HYPOTHESES WHY??

574 We also tested if the better performance of phased allele sequence data may be a
575 mere effect of doubling the number of sequences in the MSAs, since we are producing two
576 allele sequences for each individual rather than one contig sequence. Therefore we
577 generated a dataset of chimeric allele sequences, which contains the same number of
578 sequences as the phased allele data, but all heterozygous positions within an individual are
579 randomly shuffled between the two allele sequences. Just as the contig data, the chimeric
580 allele data led to an overestimation of the number of coalescent species (Fig. 6d) and to a
581 biased estimation toward older divergence times (Fig. 7). The fact that contig sequences
582 and chimeric allele sequences produce very similar results in our analyses is not surprising,
583 since contigs themselves represent chimeric consensus sequences of the variation found at a
584 locus within an individual. However, this result does demonstrate that the number of
585 sequences being analyzed does not affect the phylogenetic estimation in terms of topology,
586 species delimitation and divergence time estimation (Figs. 6 and 7).

587 Based on these findings we conclude that proper phasing of heterozygous positions
588 is crucial and is clearly preferable to the alternative of coding heterozygous sites as IUPAC
589 ambiguity codes, particularly when the estimation of divergence times is of interest.
590 Further, allele sequences are the theoretically more appropriate input for coalescent models
591 and should thus always be the preferred data type for phylogenetic studies.

592 *UCEs as source for SNP data*

593 Due to the size (number of loci) of many sequence capture datasets, it is often unfeasible to
594 analyze all MSAs jointly in one MSC analysis (Smith et al. 2014; Manthey et al. 2016) due
595 to computational limitations. For all sequence based MSC analyses in this study, we
596 reduced the UCE dataset from 820 loci to 150 loci in order to reach convergence of the

597 MCMC within a reasonable time frame (three to four days, single core on a Mac Pro, Late
598 2013, 3.5 GHz 6-Core Intel Xeon E5 processor). However, a viable approach of data
599 reduction while keeping the multilocus information of all loci, is to analyze only a single
600 polymorphic position per MSA using SNAPP Bryant et al. (2012). In our study, this
601 approach produces the correct species tree topology and also estimates the relative
602 node-heights correctly (Fig. 6f). However, SNAPP can only estimate relative and not
603 absolute values for divergence times Bryant et al. (2012), in contrast to the sequence based
604 analyses Figs. 6a to 6d, which deliver absolute divergence time estimates.

605 Sequence capture datasets such as UCEs provide a suitable data source to extract
606 both, full sequence alignments and SNP datasets of sufficient size for a robust species tree
607 estimation. Even though sequence capture data are not commonly thought about as a
608 source of SNP data, they can in many cases even be preferable to other sequencing
609 techniques, such as RAD sequencing, for producing SNP data. This is because sequence
610 capture data yields a sizable complete SNP matrix (SNPs recovered for all individuals),
611 due to the targeted sequence enrichment prior to sequencing. In our case the complete
612 matrix of unlinked SNPs in the empirical data consisted of 598 positions, which were
613 present and sufficiently supported (>three high-quality reads per haplotype) in all taxa.
614 Particularly when evolutionary distances between individuals are large, RAD sequencing
615 and other restriction site based sequencing techniques are not expected to yield many loci
616 shared between all individuals, while UCE data on the other hand are not sensitive to
617 great evolutionary distances Harvey et al. (2016). In those cases the size of the complete
618 SNP matrix resulting from UCE data can exceed that resulting from RAD sequencing.
619 Additionally, UCE data provide hundreds to thousands of full sequence MSAs as well as
620 the complete mitochondrial genome as a byproduct of the sequence enrichment. The
621 mitochondrial genome provides an excellent marker for estimating absolute divergence
622 times (Fig. 3), based on substitution rates of mitochondrial markers which are known for

623 birds Lerner et al. (2011) and thus is a valuable phylogenetic data source.

624 In this study we present and make available a new SNP calling approach for
625 sequence capture data. In contrast to other SNP calling softwares such as GATK
626 (McKenna et al. 2010), which use BAM files as input, our approach uses full sequence
627 MSAs as input (see Fig. 2), in order to identify and extract sites in the alignments that
628 show variation between any user-defined group of sequences. Our SNP calling algorithm
629 can technically be applied to any type of sequence alignments (i.e. allele or contig sequence
630 alignments), even though we recommend to work with SNPs extracted from phased allele
631 alignments, as they represent the true heterozygous information. The user can choose
632 whether or not to allow missing data or ambiguities in the extracted positions, whether to
633 extract them in binary format (as e.g. required by SNAPP) or as nucleotides, and if only
634 one SNP per locus or all SNPs should be extracted. Overall we find that extracting SNPs
635 in this manner (from MSAs) is more efficient and straightforward for sequence capture data
636 than alternative methods such as implemented in GATK. This is mainly because, in our
637 approach, the user does not have to revisit the BAM files of all samples for cumbersome and
638 superfluous re-cleaning and re-filtering of the data, but can instead use the already cleaned
639 and filtered information in the MSAs. Particularly after phasing the FASTQ reads into
640 allele sequences in a previous step, there is no need to go back to the BAM files to extract
641 bi-allelic SNPs. Another advantage of our SNP calling script is that the user can choose a
642 binary SNP output, a feature particularly interesting to users planning downstream
643 analyses with SNAPP. Thus our SNP calling mechanism is an easy, open-source and
644 straightforward tool to derive SNP data from any set of multiple sequence alignments.

645 *Phylogenetic relationships in Topaza*

646 In contrast to previous findings that UCE sequences may be too conserved to resolve
647 phylogenetic questions at shallow taxonomic scales (Giarla and Esselstyn 2015), we find

648 that UCE sequences are well suited for resolving the intraspecific genetic structure in
649 *Topaza*.

650 *One or two species?*— Our results show a separation of two lineages within the genus
651 *Topaza* that is dated at ca. 2.4 Ma in the mitochondrial tree (Fig. 3). These lineages are
652 consistent with the previously described morphospecies *T. pyra* (Gould, 1846) and *T. pella*
653 (Linnaeus, 1758), which are generally accepted in the ornithological community (Hu et al.
654 2000; del Hoyo et al. 2016a). However, the species status of *T. pyra* has been challenged by
655 some authors (Ornés-Schmitz and Schuchmann 2011; Schuchmann 1999). These authors
656 concluded that *Topaza* is a monotypic genus with *T. pyra* being a subspecies of *T. pella*,
657 which they refer to as *T. pella pyra*. Their findings are based on the analyses of plumage
658 coloration, in which they found an “east-west clinal trend of characters” (Ornés-Schmitz
659 and Schuchmann 2011). In contrast, we do not find such an east-west clinal trend in the
660 genetic data. Instead, *T. pyra* is consistently supported as a separate lineage across all
661 analyses, lending no support for the conspecificity of these two taxa (Figs. 3 and 5).

662 One aim of this study was to evaluate the genetic structure within these two
663 morphospecies *T. pyra* and *T. pella*. The mitochondrial tree shows two divergent clades
664 within *T. pyra* (Fig. 3), but these clades are not supported by the UCE data (Fig. 5), even
665 though the allele sequence data are picking up a signal that possibly indicates that two
666 such clades are in the process of diversifying (Fig. 5b). For *T. pella* on the other hand, we
667 consistently find the same genetically distinct clades throughout all multilocus MSC
668 analyses (Fig. 5), leading us to distinguish between the following populations that are
669 congruent with previous morphological subspecies descriptions:

670 *Northern T. pella population: T. pella pella*.— For the mitochondrial tree and all MSC
671 species trees, we find the northern *T. pella* samples 5 and 6 to be sister taxa with high
672 support values (98-100% PP, Figs. 3 and 5). Particularly in the mitochondrial tree (Fig.

673 3), these two samples appear as close sister taxa, separated by only very short terminal
674 branches. Their close position in the mitochondrial tree shows that, even though
675 geographically far apart, samples 5 and 6 share a relatively recent MRCA in the
676 mitochondrial genealogy, indicating some rather recent gene flow. The sampling locality of
677 sample 5 is within the range of the subspecies *T. pella pella*, which mainly extends across
678 the Guiana shield (Peters 1945; Schuchmann 1999; Hu et al. 2000; Ornés-Schmitz and
679 Schuchmann 2011). Given the sampling location of genetically linked sample 6, which also
680 has been morphologically identified as *T. pella pella* (Table 1), we propose to extend the
681 distribution range of *T. pella pella* from the Guiana shield all the way south to the
682 northern Amazon River bank (see map in Fig. 3).

683 *Southern T. pella population: T. pella microrhyncha*.—In the same manner as for the
684 northern population *T. pella pella*, we also consistently find the southern *T. pella* samples
685 8 and 9 to be sister taxa (99-100% PP, Figs. 3 and 5). The sampling locations of these two
686 samples are included in the distribution range of the previously recognized subspecies *T.*
687 *pella microrhyncha*, extending from the southern bank of the Amazon River as far South as
688 Porto Velho (Brazil) at the Madeira River, close to the border to Bolivia (Peters 1945;
689 Schuchmann 1999; Ornés-Schmitz and Schuchmann 2011). This southernmost boundary of
690 *T. pella microrhyncha* is not accepted by Hu et al. (2000), who instead conclude that this
691 southernmost population belongs to *T. pella pella*. In contrast to the findings by Hu et al.,
692 our genetic data clearly support the southernmost sample 9 belonging to the same
693 population as sample 8, which was morphologically identified as *T. pella microrhyncha*.
694 This leads us to propose a distribution range of *T. pella microrhyncha* as shown in Fig. 3,
695 in agreement with the findings by Peters (1945), Schuchmann (1999), and Ornés-Schmitz
696 and Schuchmann (2011).

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

698 concerning the phylogenetic placement of sample 7, which was collected from the southern
699 estuary region of the Amazon River and morphologically identified as *T. pella smaragdula*.
700 The sampling locality also falls into the range of the subspecies *T. pella smaragdula* (Peters
701 1945; Hu et al. 2000; Ornés-Schmitz and Schuchmann 2011), with a distribution including
702 the Amazon River estuary and extending north along the coast to French Guiana. All
703 MSC analyses of the UCE sequence and SNP data place sample 7 with high confidence
704 (97-100% PP) as sister to the southern clade *T. pella microrhyncha* (Fig. 5), whereas in
705 the mitochondrial phylogeny this sample is placed as sister to *T. pella pella* in the North.

706 The discordance between a gene tree and the species tree in a scenario such as this
707 could be the effect of incomplete lineage sorting, which is most likely if the species or clades
708 in question have diverged rather recently and if population sizes are large. Given that the
709 divergence between *T. pella pella* and *T. pella microrhyncha* appears to be considerably
710 deep based on the multilocus data (crown height of *T. pella* see Fig. 5) and given that
711 mitochondria are generally considered to have only 25% of the population size of nuclear
712 loci, it is rather unlikely that the position of sample 7 in the mitochondrial tree is a result
713 of incomplete lineage sorting in this case. It seems more likely that the separate position of
714 sample 7 in the mitochondrial tree is a result of introgression of the mitochondrial genome
715 from *T. pella pella* into the gene pool of *T. pella smaragdula*. However, a denser taxon
716 sampling would be necessary to further evaluate the evolutionary history of this particular
717 population. The case of sample 7 highlights that the mitochondrial tree presents a single
718 gene tree phylogeny that only shows one of many genealogies and therefore must not be
719 equaled to a species tree phylogeny. Hence it is important to generate multilocus data for
720 an informed inference of the species tree phylogeny.

721 *Summarizing biogeographic remarks.*— The presence of genetically similar individuals
722 sampled at great geographic distances (e.g. samples 5 and 6) suggests that *Topaza*

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

734 SUPPLEMENTARY MATERIAL

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

738 AVAILABILITY

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

ACKNOWLEDGMENTS

747 We wish to thank all those ornithologists who have dedicated their time to collecting
748 samples in Amazonia; museum curators for providing us with samples for this study;
749 Brazilian authorities for issuing the permits needed for this work; our lab engineer Anna
750 Ansebo for laboratory assistance; Alexander Zizka for assistance in creating the sampling
751 and range maps; HBW Alive for providing the *Topaza* illustrations; and colleagues at our
752 labs for discussions and feedback. Computational analyses were performed on the
753 bioinformatics computer cluster Albiorix at the Department of Biological and
754 Environmental Sciences, University of Gothenburg.

FUNDING

756 This work was funded by the Swedish Research Council to A. Antonelli (B0569601) and B.
757 Oxelman (2012-3917); the CNPq (grants 310593/2009-3; ‘INCT em Biodiversidade e Uso
758 da Terra da Amazônia’ 574008/2008-0; 563236/2010-8; and 471342/2011-4), FAPESPA
759 (ICAAF 023/2011), and NSF-FAPESP (grant 1241066 - Dimensions US-BIOTA-São Paulo:
760 Assembly and evolution of the Amazonian biota and its environment: an integrated
761 approach) to A. Aleixo; the European Research Council under the European Union’s
762 Seventh Framework Programme (FP/2007-2013, ERC Grant Agreement n. 331024), the
763 Swedish Foundation for Strategic Research and a Wallenberg Academy Fellowship to A.
764 Antonelli.

765 *

766 References

- 767 Bodily, P. M., M. Fujimoto, C. Ortega, N. Okuda, J. C. Price, M. J. Clement, and Q. Snell.
768 2015. Heterozygous genome assembly via binary classification of homologous sequence.
769 BMC Bioinformatics 16:S5.
- 770 Bolger, A. M., M. Lohse, and B. Usadel. 2014. Trimmomatic: a flexible trimmer for
771 Illumina sequence data. Bioinformatics 30:2114–20.
- 772 Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard,
773 A. Rambaut, and A. J. Drummond. 2014. BEAST 2: a software platform for Bayesian
774 evolutionary analysis. PLoS Computational Biology 10:e1003537.
- 775 Bryant, D., R. Bouckaert, J. Felsenstein, N. A. Rosenberg, and A. RoyChoudhury. 2012.
776 Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a
777 full coalescent analysis. Molecular Biology and Evolution 29:1917–32.
- 778 Clair, C. C. S. 2003. Comparative permeability of roads, rivers, and meadows to songbirds
779 in Banff national park. Conservation Biology 17:1151–1160.
- 780 Degnan, J. H. and N. a. Rosenberg. 2009. Gene tree discordance, phylogenetic inference
781 and the multispecies coalescent. Trends in Ecology and Evolution 24:332–340.
- 782 del Hoyo, J., N. Collar, G. Kirwan, and P. Boesman. 2016a. Fiery Topaz (*Topaza pyra*). in
783 Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,
784 D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.
- 785 del Hoyo, J., A. Elliott, J. Sargatal, D. Christie, and E. de Juana. 2016b. Handbook of the
786 Birds of the World Alive. Lynx Edicions, Barcelona, Spain.
- 787 Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics
788 with BEAUTi and the BEAST 1.7. Molecular Biology and Evolution 29:1969–73.

- 789 Edwards, S. V., L. Liu, and D. K. Pearl. 2007. High-resolution species trees without
790 concatenation. *Proceedings of the National Academy of Sciences* 104:5936–5941.
- 791 Faircloth, B. C. 2015. PHYLUCE is a software package for the analysis of conserved
792 genomic loci. *Bioinformatics* 32:786–788.
- 793 Faircloth, B. C., J. E. McCormack, N. G. Crawford, M. G. Harvey, R. T. Brumfield, and
794 T. C. Glenn. 2012. Ultraconserved elements anchor thousands of genetic markers
795 spanning multiple evolutionary timescales. *Systematic Biology* 61:717–26.
- 796 Faircloth, B. C., L. Sorenson, F. Santini, and M. E. Alfaro. 2013. A phylogenomic
797 perspective on the radiation of ray-finned fishes based upon targeted sequencing of
798 ultraconserved elements (UCEs). *PLoS ONE* 8:e65923.
- 799 Felsenstein, J. 2005. Phylip (phylogeny inference package) version 3.6. distributed by the
800 author. dep genome sci univ washington, seattle.
- 801 Fernandes, A. M., M. Wink, and A. Aleixo. 2012. Phylogeography of the chestnut-tailed
802 antbird (*Myrmeciza hemimelaena*) clarifies the role of rivers in Amazonian biogeography.
803 *Journal of Biogeography* 39:1524–1535.
- 804 Garrick, R. C., P. Sunnucks, and R. J. Dyer. 2010. Nuclear gene phylogeography using
805 PHASE: dealing with unresolved genotypes, lost alleles, and systematic bias in
806 parameter estimation. *BMC Evolutionary Biology* 10:118.
- 807 Giarla, T. C. and J. A. Esselstyn. 2015. The challenges of resolving a rapid, recent
808 radiation: empirical and simulated phylogenomics of philippine shrews. *Systematic
809 Biology* 64:727–740.
- 810 Gnirke, A., A. Melnikov, J. Maguire, P. Rogov, E. M. LeProust, W. Brockman, T. Fennell,
811 G. Giannoukos, S. Fisher, C. Russ, S. Gabriel, D. B. Jaffe, E. S. Lander, and

- 812 C. Nusbaum. 2009. Solution hybrid selection with ultra-long oligonucleotides for
813 massively parallel targeted sequencing. *Nature Biotechnology* 27:182–189.
- 814 Harvey, M. G., B. T. Smith, T. C. Glenn, B. C. Faircloth, and R. T. Brumfield. 2016.
815 Sequence capture versus restriction site associated DNA sequencing for shallow
816 systematics. *Systematic Biology Advance Access* syw036.
- 817 Hayes, F. E. and J. A. N. Sewlal. 2004. The Amazon River as a dispersal barrier to
818 passerine birds: effects of river width, habitat and taxonomy. *Journal of Biogeography*
819 31:1809–1818.
- 820 He, D., A. Choi, K. Pipatsrisawat, A. Darwiche, and E. Eskin. 2010. Optimal algorithms
821 for haplotype assembly from whole-genome sequence data. *Bioinformatics* 26:i183–i190.
- 822 Hu, D.-S., L. Joseph, and D. J. Agro. 2000. Distribution, variation, and taxonomy of
823 *Topaza* Hummingbirds (Aves: Trochilidae). *Ornitologia Neotropical* 11:123–142.
- 824 Iqbal, Z., M. Caccamo, I. Turner, P. Flicek, and G. McVean. 2012. De novo assembly and
825 genotyping of variants using colored de Bruijn graphs. *Nature Genetics* 44:226–232.
- 826 Jones, G. 2017. Algorithmic improvements to species delimitation and phylogeny estimation
827 under the multispecies coalescent. *Journal of Mathematical Biology* 74:447–467.
- 828 Jones, G., Z. Aydin, and B. Oxelman. 2014. DISSECT: an assignment-free Bayesian
829 discovery method for species delimitation under the multispecies coalescent.
830 *Bioinformatics* 31:991–998.
- 831 Katoh, K., G. Asimenos, and H. Toh. 2009. Multiple alignment of DNA sequences with
832 MAFFT. *Methods in Molecular Biology* 537:39–64.
- 833 Kolaczkowski, B. and J. W. Thornton. 2004. Performance of maximum parsimony and
834 likelihood phylogenetics when evolution is heterogeneous. *Nature* 431:980–984.

- 835 Kubatko, L. S. and J. H. Degnan. 2007. Inconsistency of Phylogenetic Estimates from
836 Concatenated Data under Coalescence. *Systematic Biology* 56:17–24.
- 837 Lerner, H. R., M. Meyer, H. F. James, M. Hofreiter, and R. C. Fleischer. 2011. Multilocus
838 resolution of phylogeny and timescale in the extant adaptive radiation of Hawaiian
839 honeycreepers. *Current Biology* 21:1838–1844.
- 840 Li, H. and R. Durbin. 2010. Fast and accurate long-read alignment with Burrows-Wheeler
841 transform. *Bioinformatics* 26:589–595.
- 842 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis,
843 and R. Durbin. 2009. The Sequence Alignment/Map format and SAMtools.
844 *Bioinformatics* 25:2078–9.
- 845 Lischer, H. E., L. Excoffier, and G. Heckel. 2014. Ignoring heterozygous sites biases
846 phylogenomic estimates of divergence times: Implications for the evolutionary history of
847 *microtus* voles. *Molecular Biology and Evolution* 31:817–831.
- 848 Manthey, J. D., L. C. Campillo, K. J. Burns, and R. G. Moyle. 2016. Comparison of
849 target-capture and restriction-site associated DNA sequencing for phylogenomics: a test
850 in cardinalid tanagers (Aves, Genus: *Piranga*). *Systematic Biology Advance Access*
851 syw005.
- 852 McCormack, J. E., B. C. Faircloth, N. G. Crawford, P. A. Gowaty, R. T. Brumfield, and
853 T. C. Glenn. 2012. Ultraconserved elements are novel phylogenomic markers that resolve
854 placental mammal phylogeny when combined with species-tree analysis. *Genome*
855 *Research* 22:746–754.
- 856 McGuire, J., C. C. Witt, J. V. Remsen, A. Corl, D. L. Rabosky, D. L. Altshuler, and

- 857 R. Dudley. 2014. Molecular phylogenetics and the diversification of hummingbirds.
- 858 Current Biology 24:910–916.
- 859 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky,
- 860 K. Garimella, D. Altshuler, S. Gabriel, M. Daly, and M. A. DePristo. 2010. The Genome
- 861 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
- 862 sequencing data. Genome research 20:1297–303.
- 863 Meiklejohn, K. A., B. C. Faircloth, T. C. Glenn, R. T. Kimball, and E. L. Braun. 2016.
- 864 Analysis of a rapid evolutionary radiation using ultraconserved elements (UCEs):
- 865 Evidence for a bias in some multispecies coalescent methods. Systematic Biology
- 866 Advance Access syw014.
- 867 Milne, I., G. Stephen, M. Bayer, P. J. A. Cock, L. Pritchard, L. Cardle, P. D. Shaw, and
- 868 D. Marshall. 2013. Using Tablet for visual exploration of second-generation sequencing
- 869 data. Briefings in Bioinformatics 14:193–202.
- 870 Mirarab, S., R. Reaz, M. S. Bayzid, T. Zimmermann, M. S. Swenson, and T. Warnow.
- 871 2014. ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics
- 872 30:541–548.
- 873 Moore, R. P., W. D. Robinson, I. J. Lovette, and T. R. Robinson. 2008. Experimental
- 874 evidence for extreme dispersal limitation in tropical forest birds. Ecology Letters
- 875 11:960–968.
- 876 Mossel, E. and E. Vigoda. 2005. Phylogenetic MCMC algorithms are misleading on
- 877 mixtures of trees. Science 309:2207–9.
- 878 Ornés-Schmitz, A. and K. L. Schuchmann. 2011. Taxonomic review and phylogeny of the

- 879 hummingbird genus *Topaza* (Gray, 1840) using plumage color spectral information.
- 880 *Ornitologia Neotropical* Pages 25–38.
- 881 Peters, J. L. 1945. Check-list of birds of the world. Volume 5 ed. Harvard Univ. Press,
- 882 Cambridge, Massachusetts.
- 883 Potts, A. J., T. A. Hedderson, and G. W. Grimm. 2014. Constructing Phylogenies in the
- 884 Presence Of Intra-Individual Site Polymorphisms (2ISPs) with a Focus on the Nuclear
- 885 Ribosomal Cistron. *Systematic Biology* 63:1–16.
- 886 Rambaut, A., M. A. Suchard, W. Xie, and A. Drummond. 2013. Tracer v1.6.
- 887 Rannala, B. and Z. Yang. 2003. Bayes estimation of species divergence times and ancestral
- 888 population sizes using DNA sequences from multiple loci. *Genetics* 164:1645–1656.
- 889 Remsen, J. V. and T. A. Parker. 1983. Contribution of river-created habitats to bird
- 890 species richness in Amazonia. *Biotropica* 15:223–231.
- 891 Ribas, C. C., a. Aleixo, a. C. R. Nogueira, C. Y. Miyaki, and J. Cracraft. 2012. A
- 892 palaeobiogeographic model for biotic diversification within Amazonia over the past three
- 893 million years. *Proceedings of the Royal Society B: Biological Sciences* 279:681–689.
- 894 Russell, R. W., F. L. Carpenter, M. A. Hixon, and D. C. Paton. 1994. The impact of
- 895 variation in stopover habitat quality on migrant rufous hummingbirds. *Conservation*
- 896 *Biology* 8:483–490.
- 897 Schrempf, D., B. Q. Minh, N. De Maio, A. von Haeseler, and C. Kosiol. 2016. Reversible
- 898 polymorphism-aware phylogenetic models and their application to tree inference. *Journal*
- 899 *of Theoretical Biology* 407:362–370.

- 900 Schuchmann, K., G. Kirwan, and P. Boesman. 2016. Crimson Topaz (*Topaza pella*). *in*
901 Handbook of the Birds of the World Alive (J. del Hoyo, A. Elliott, J. Sargatal,
902 D. Christie, and E. de Juana, eds.). Lynx Edicions, Barcelona, Spain.
- 903 Schuchmann, K. L. 1999. Family Trochilidae (hummingbirds). Pages 468–680 *in* Handbook
904 of the Birds of the World Alive (J. del Hoyo, A. Elliott, and J. Sargatal, eds.) volume 5
905 ed. Lynx Edicions, Barcelona, Spain.
- 906 Simpson, J. T., K. Wong, S. D. Jackman, J. E. Schein, S. J. M. Jones, and I. Birol. 2009.
907 ABySS: a parallel assembler for short read sequence data. *Genome Research* 19:1117–23.
- 908 Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014.
909 Target capture and massively parallel sequencing of ultraconserved elements for
910 comparative studies at shallow evolutionary time scales. *Systematic Biology* 63:83–95.
- 911 Sullivan, B. L., C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling. 2009. eBird:
912 A citizen-based bird observation network in the biological sciences. *Biological
913 Conservation* 142:2282–2292.
- 914 Thom, G. and A. Aleixo. 2015. Cryptic speciation in the white-shouldered antshrike
915 (*Thamnophilus aethiops*, Aves - Thamnophilidae): The tale of a transcontinental
916 radiation across rivers in lowland Amazonia and the northeastern Atlantic Forest.
917 *Molecular Phylogenetics and Evolution* 82:95–110.
- 918 Wyman, S. K., R. K. Jansen, and J. L. Boore. 2004. Automatic annotation of organellar
919 genomes with DOGMA. *Bioinformatics* 20:3252–5.
- 920 Yang, Z. 2015. The BPP program for species tree estimation and species delimitation.
921 *Current Zoology* 61:854–865.

922 Yu, L., Y.-W. Li, O. a. Ryder, and Y.-P. Zhang. 2007. Analysis of complete mitochondrial
923 genome sequences increases phylogenetic resolution of bears (Ursidae), a mammalian
924 family that experienced rapid speciation. *BMC Evolutionary Biology* 7:198.