# Assessing the potential of RAD-sequencing to resolve phylogenetic relationships within species radiations: the fly genus *Chiastocheta* (Diptera: Anthomyiidae) as a case study

Tomasz Suchan[1,2], Anahí Espíndola[3*], Sereina Rutschmann[1,4*], Brent C. Emerson[5,6], Kevin Gori[7], Christophe Dessimoz[1,8,9,10], Nils Arrigo[1], Michał Ronikier[2], Nadir Alvarez[1]

[1] Department of Ecology and Evolution, University of Lausanne, Switzerland
[2] W. Szafer Institute of Botany, Polish Academy of Sciences, Kraków, Poland
[3] Department of Biological Sciences, University of Idaho, Moscow, ID, USA
[4] Department of Biochemistry, Genetics and Immunology, University of Vigo, Spain
[5] Island Ecology and Evolution Research Group, Instituto de Productos Naturales y Agrobiología (IPNA-CSIC), La Laguna, Tenerife, Canary Islands, Spain
[6] School of Biological Sciences, University of East Anglia, Norwich, UK
[7] Department of Veterinary Medicine, University of Cambridge, UK
[8] Center for Integrative Genomics, University of Lausanne, Switzerland
[9] Department of Genetics, Evolution & Environment and Department of Computer Science, University College London, UK
[10] Swiss Institute of Bioinformatics, Lausanne, Switzerland

[*] these authors are considered as second co-authors

corresponding authors:
Tomasz Suchan (t.suchan@botany.pl), W. Szafer Institute of Botany, Polish Academy of Sciences, Lubicz 46, 31-512 Kraków, Poland
Nadir Alvarez (nadir.alvarez@unil.ch), Department of Ecology and Evolution, University of Lausanne, Biophore Building, 1015 Lausanne, Switzerland; fax: +41 21 692 41 65

# Abstract

Determining phylogenetic relationships among recently diverged species has long been a challenge in evolutionary biology. Cytoplasmic markers, which have been widely used notably in the context of molecular barcoding, have not always proved successful in resolving such phylogenies, but phylogenies for closely related species have been resolved at a much higher detail in the last couple of years with the advent of next-generation-sequencing technologies and associated techniques of reduced genome representation. Here we examine the potential and limitations of one of such techniques — Restriction-site Associated DNA (RAD) sequencing, a method that produces thousands of (mostly) anonymous nuclear markers, in disentangling the phylogeny of the fly genus *Chiastocheta* (Diptera: Anthomyiidae). This genus encompasses seven described species of seed predators, which have been widely studied in the context of their ecological and evolutionary interactions with the plant *Trollius europaeus* (Ranunculaceae)*. So far, phylogenetic analyses using mitochondrial markers failed to resolve monophyly of most of the species from this recently diversified genus, suggesting that their taxonomy may need to be revised. However, relying on a single, non-recombining molecule and ignoring potential incongruences between mitochondrial and nuclear loci may provide incomplete account of a lineage history. In this study, we apply both classical Sanger sequencing of three mtDNA regions and RAD-sequencing, for reconstructing the phylogeny of the genus. Contrasting with results based on mitochondrial markers, RAD-sequencing analyses retrieved the monophyly of all seven species, in agreement with the morphological species assignment. We found robust nuclear-based species assignment of individual samples, and low levels of estimated contemporary gene flow among them. However, despite recovering species'

2

53     monophyly, interspecific relationships varied depending on the set of RAD loci

54     considered, producing contradictory topologies. Moreover, coalescence-based

55     phylogenetic analyses revealed low supports for most of the interspecific relationships.

56     Our results indicate that despite the higher performance of RAD-sequencing in terms of

57     species trees resolution compared to cytoplasmic markers, reconstructing inter-specific

58     relationships may lie beyond the possibilities offered by large sets of RAD-sequencing

59     markers in cases of strong gene tree incongruence.

60

61     Keywords: coalescent analysis; DNA barcoding; maximum likelihood; mito-nuclear

62     incongruence; single nucleotide polymorphisms; quartet inference

63

# 1. Introduction

Recently diverged lineages pose a problem for traditional phylogenetic approaches that typically rely on a small set of relatively slowly evolving loci (DeFilippis 2000), often lacking resolution at narrower evolutionary scales (Cariou et al. 2013). In addition, complex processes such as incomplete lineage sorting (Avise et al. 2008; Maddison & Knowles 2006; Pollard et al. 2006) and gene flow among species (Leaché et al. 2013) increase incongruences among gene trees and topological deviations from the species tree (Dengan & Rosenberg 2009; Maddison 1997). This is especially true for lineages that have undergone rapid radiations, in which ancestral polymorphisms sorted idiosyncratically into the descendant taxa through short evolutionary nodes (Avise et al. 2008), and in cases where subsequent evolutionary events may blur phylogenetic signal (Whitfield & Kjer 2008; Whitfield & Lockhart 2007). Sampling more loci has been shown to be a promising approach in such cases (Rokas & Carroll 2005; Townsend et al. 2011; Wielstra et al. 2014; Williams et al. 2013), but the spectrum of genetic markers developed for phylogeny estimation is still limited (Whitfield & Kjer 2008). Next-generation sequencing approaches, particularly reduced representation genome sequencing (Davey et al. 2011), offer the possibility to sample thousands of genomic markers from non-model species. Among them, Restriction site-Associated DNA (RAD; Baird et al. 2008) techniques rely on the sequencing of short DNA fragments flanking restriction sites, generating random anonymous genomic markers, homologous across the analyzed samples (Andrews et al. 2016; Davey & Blaxter 2010). From a phylogenetic perspective, an important aspect of RAD markers is the rise in the proportion of 'null alleles' as genome divergence across samples increases. This phenomenon is caused by random mutations occurring in the restriction sites that decrease the numbers of

88    shared RAD loci among taxa, resulting in data matrices containing large amounts of

89    missing data (Cariou et al. 2013; Chattopadhyay et al. 2014; Gautier et al. 2013).

90    However, using an in-silico approach Rubin et al. (2012) and Cariou et al. (2013) have

91    shown that RAD-seq data can be used successfully to resolve species relationships that

92    transcend timescales up to 60 Mya (million years ago). Experimentally sampled RAD

93    datasets have been applied to reconstruct phylogenetic relationships, mostly among

94    recently diverged taxa (e.g., Eaton & Ree 2013; Harvey et al. 2016; Jones et al. 2013;

95    Leaché et al. 2015; Nadeau et al. 2013; Wagner et al. 2013), with fewer studies involving

96    more distantly related ones, even up to even 80 Mya (e.g., Cruaud et al. 2014; Eaton et

97    al. 2016; Herrera and Shank 2016; Hipp et al. 2014; Pante et al. 2015). Although these

98    genomic datasets improved phylogenetic inferences for groups that were ambiguous

99    using classical markers (e.g., Escudero et al. 2014; Hipp et al. 2014), the potential utility

100    of RAD loci for resolving more complex phylogenetic histories, such as those where

101    historical introgression has occurred or those associated with incomplete lineage

102    sorting, remains still poorly explored (Combosch & Vollmer 2015; Eaton & Ree 2013).

103    Moreover, the use of RAD datasets as markers for evolutionary genetics has recently

104    been heavily discussed (Lowry et al. 2017; McKinney et al. 2017).

105    In this study, we test the utility of RAD-sequencing to recover phylogenetic

106    relationships in a genus of seed parasitic pollinators of *Trollius europaeus*

107    (Ranunculaceae) — flies from the genus *Chiastocheta* Pokorny, 1889 (Diptera:

108    Anthomyiidae). Here, sequencing of mitochondrial markers failed to reveal the

109    monophyly and phylogenetic relationships among previously morphologically

110    described species (Després et al. 2002; Espíndola et al. 2012). This discordance

111    between morphology and mitochondrial phylogeny has been interpreted as a call for a

taxonomic revision, and a possible reconsideration of conclusions from previous

ecological and evolutionary studies (Espíndola et al. 2012). However, several

mechanisms may cause mitochondria not to track species evolution (Funk & Omland

2003) and, indeed, there are many cases where mitochondrial and nuclear gene trees

have been shown to be incongruent (e.g., Govindarajulu et al. 2015; Phillips et al. 2013;

Seehausen et al. 2003). As relying on a single, non-recombining molecule may provide a

misleading account of a species history (Ballard & Whitlock 2004), utilizing a large set

of independent nuclear loci (sampled through RAD-sequencing) should allow us to test

the monophyly of the morphologically described species and resolve phylogenetic

relationships among them. Whether or not molecular markers are able to reveal

scenarios of rapid radiations is still an open question (Giarla & Esselstyn 2015). In

these, identifying a single species tree might lie beyond analytical possibilities due to

pervasive conflicts among the gene trees, particularly when population sizes are large

and speciation events happen at a higher rate than the mutation-drift equilibrium,

eventually producing conflicting topologies. In order to explore gene and species trees,

we applied both a concatenation-based phylogenetic approach (i.e., RAxML; Stamatakis

2014) and a coalescence-based inference method (i.e., SVDquartets; Chifman & Kubatko

2014) to a RAD-seq dataset encompassing specimens from 51 European populations,

representative of the seven recognized *Chiastocheta* morphospecies. In order to

examine the extent to which different combinations of RAD loci may produce distinct

species trees, we used a newly developed algorithm that performs loci binning, using

dissimilarity levels among phylogenetic patterns retrieved at single loci (treeCl; Gori et

al. 2016). We also applied population genetics clustering algorithms (i.e., STRUCTURE;

Pritchard et al. 2000) as a control. Eventually, we compared our results to those

136  obtained with classical phylogenetic inference based on concatenation of three

137  mitochondrial regions.

# 138  2. Materials and methods

## 139  2.1 Study system

140  The center of origin and diversity of *Chiastocheta* has been inferred to be the Western

141  Palearctic, where seven fly species are involved in nursery pollination interactions with

142  *Trollius europaeus* L. (Espíndola et al. 2012; Pellmyr 1989, 1992; Suchan et al. 2015).

143  These seven morphologically delimited European *Chiastocheta* species, namely *C.*

144  *dentifera* Hennig 1953; *C. inermella* (Zetterstedt, 1838); *C. lophota* Karl, 1943; *C.*

145  *macropyga* Hennig, 1953; *C. rotundiventris* Hennig, 1953; *C. setifera* Hennig, 1953 and *C.*

146  *trollii* (Zetterstedt, 1838) are ecologically very similar and often sympatric (Collin 1954;

147  Hennig 1976; Michelsen 1985; Zetterstedt 1845; V. Michelsen pers. comm.). In his

148  monograph of this plant-pollinator interaction, Pellmyr (1992) discussed another

149  species, *C. abruptiventris* as a northern vicariant of *C. rotundiventris*, a taxon not

150  supported by previous molecular studies (Espíndola et al. 2012) and never formally

151  described. Although all *Chiastocheta* reproduce within the flowers of *T. europaeus*, with

152  potential cross-species mating possibilities, no putative hybrids have been observed

153  based on genital morphology (T. Suchan and A. Espíndola, pers. obs.).

154  Although the species are well defined morphologically, mitochondrial phylogenies

155  recovered only three monophyletic clades – *C. rotundiventris*, *C. dentifera*, and *C. lophota*

156  (Després et al. 2002; Espíndola et al. 2012), and suggested a polyphyletic origin for *C.*

157  *inermella* and *C. setifera* (Després et al. 2002; Espíndola et al. 2012), with *C. macropyga*

158  and *C. trollii* being paraphyletic (Espíndola et al. 2012). Molecular dating placed the

159    most recent common ancestor of all European species at the end of the Pliocene (2-3.4

160    Ma; Després et al. 2002; Espíndola et al. 2012), and indicated that most diversification

161    events occurred within the last 1.6 Ma.


162    ## 2.2 Sampling

163    *Chiastocheta* specimens were sampled from 51 European populations during spring and

164    summer 2006, 2007, and 2008 (Table 1; maps on Fig. S1). The flies were killed and

165    preserved in 70% ethanol and stored at room temperature until DNA extraction.

166    Collected specimens were identified to morphospecies following Hennig (1976) and

167    unpublished keys (V. Michelsen). All identifications were confirmed by an expert (V.

168    Michelsen, Natural History Museum of Denmark, Copenhagen), as the taxonomical

169    revision of the genus is not yet published.


170    ## 2.3 Sequencing mitochondrial regions and RAD markers

171    DNA was extracted from insect legs using a DNeasy Blood and Tissue Kit (Qiagen,

172    Hilden, Germany), following the manufacturer's instructions. We amplified three

173    mitochondrial regions: COI, COII, and the ultra-variable D-loop (control) region. We

174    followed Espíndola et al. (2012) for sequencing of the COI and COII regions. For D-loop

175    we used primers TM-N-193 and  SR-J-14612 as desxcribed in Simon et al. (1994) as

176    described by Espíndola et al. (2012) with the following modification of the PCR

177    program: 5 min at 95°C, followed by 35 cycles of 1 min at 95°C, 1 min of annealing

178    at 55°C and 2 min of elongation at 60°C, and 5 min of final elongation at 60°C. PCR

179    products were sequenced at Macrogen Inc. (South Korea) and Fasteris SA (Switzerland).

180    Chromatograms were visually corrected on ChromasPro 1.41 (Technelysium Pty. Ltd.).

181    Alignment was performed using MUSCLE algorithm (Edgar 2004) in Geneious 10.1.3

8

182 (Biomatters, Auckland, New Zeland) and gaps with more than 50% missing data in the

183 D-loop region were removed. Additionally, a dataset with D-loop removed was

184 analyzed. Double digest RAD (ddRAD) libraries were prepared according to Mastretta et

185 al. (2014), a modified protocol of Peterson et al. (2012), without performing the size-

186 selection of DNA fragments, and other minor modifications (see Supporting

187 Information). The enzymes used for DNA digestion were SbfI and MseI. Libraries were

188 sequenced at the Lausanne Genomic Technologies Facility (Switzerland) on three lanes

189 of the HiSeq2500 instrument (Illumina, San Diego, USA) using a 2x100 bp paired-end

190 reads protocol. For RAD-sequencing we introduced technical replicates for optimizing

191 *de novo* assembly and controls for the effects of sequencing errors and allele dropout on

192 the final results (Mastretta-Yanes et al. 2015; samples with "REPL" suffix in Table S1),

193 and DNA extraction replicates from the fly thoracic muscle (in order to control for flies'

194 body contamination with pollen; samples with "MUS" suffix in Table S1).


## 195  2.4 RAD-seq loci assembly

196 Two important considerations for *de novo* RAD loci assembly are the parameters for

197 clustering orthologous loci, while filtering out paralogs (Eaton 2014; Mastretta et al.

198 2015). If the sequence similarity required to consider sequence as orthologs is set too

199 high, real heterozygous alleles will be split into more than one cluster, therefore

200 creating false homozygous loci (Harvey et al. 2015). On the other hand, if the similarity

201 is set too low, this will result in paralogous sequences being clustered together. Several

202 methods were proposed for filtering out such sequences from the final dataset,

203 including ploidy filtering (removing clusters that have more than two sequences per

204 individual) and filtering out highly variable loci (Eaton 2014; Ilut et al. 2014). As there

9

205     are no general guidelines for fine-tuning the parameters mentioned above, we

206     empirically tested how the different clustering parameters affected the final dataset.

207     Finally, we chose the dataset with clustering parameter that maximized the loci overlap

208     between pairs of technical replicates (see below). Loci overlap among samples and pairs

209     of technical replicates were calculated using the RADami 1.0 library in R (Hipp et al.

210     2014).

211     Read demultiplexing and *de novo* assembly of RAD loci was performed using the pyRAD

212     3.0.1 package (Eaton 2014), based on an alignment-clustering algorithm. This approach

213     allows for indel variation among more diverged specimens. Moreover, it also allows for

214     lower similarity among the clustered reads, making it well-suited for phylogenetic-scale

215     analyses. First, the reads were demultiplexed according to the in-line 6-nucleotides

216     barcode present at the beginning of each sequenced fragment, while allowing for one

217     mismatch. Only reads with the restriction site present were retained for further

218     analyses. All nucleotides with Phred quality score lower than 20 were converted to

219     unknown bases and reads with more than four unknown sites were removed from the

220     dataset. Reads were then clustered within and between individuals, with a minimum

221     number of six reads to form a cluster and sequence similarity of 75%, 80%, 85%, 90%,

222     and 95%. Possible clustered paralogs or repetitive sequences were removed by filtering

223     out the loci that had more than five variable positions per locus or more than 10 shared

224     polymorphic sites in a locus among individuals, and the loci for which more than two

225     alleles were present per individual. Finally, datasets were produced by retaining the loci

226     present in a minimum of 10, 20, and 100 individuals and compared for the total number

227     of loci, proportion of missing data, loci overlap among replicates, and the mean number

228     of individuals per locus.

## 2.5 Phylogenetic analyses

229

230 We performed Maximum-likelihood (ML) analyses using RAxML (Stamatakis 2014)

231 with rapid bootstrap analyses and extended majority-rule consensus tree automatic

232 bootstrap stopping criterion, following search for the best-scoring ML tree. The

233 mitochondrial regions were partitioned using the PartitionFinder 2.1.1 software

234 (Lanefar et al. 2016). Analyses were performed in the RAxML 8.2.4, for the

235 mitochondrial dataset. For the dataset consisting of COI, COII, and D-loop regions the

236 GTR+G+I model with all three nucleotide positions on coding genes considered as

237 separate partitions and D-loop as a fourth partition. For the dataset consisting of COI

238 and COII the GTR+G model with the first two nucleotide positions considered as a first

239 and third nucleotide position considered as a second partition. Analyses for the RAD

240 dataset were performed using the GTRCAT model in the RAxML 8.2.10 version on the

241 CIPRES cluster (San Diego CA, USA; Miller et al. 2010). For the RAD-based dataset,

242 replicated samples were retained in the phylogenetic analyses and the concatenated

243 matrix was considered as a single partition. Additionally, ML analyses of the RAD

244 datasets with other clustering parameters were performed in order to evaluate how this

245 parameter affects the tree topology. The trees were rooted with *C. rotundiventris* as an

246 outgroup, as identified previously (Després et al. 2002; Espíndola et al. 2012).

247 To account for the effects of incongruence among nuclear loci on the inferred

248 phylogenies — for instance resulting from incomplete lineage sorting, we applied the

249 method of Gori et al. (2016) implemented in the treeCl package (http://git.io/treeCl).

250 Because the majority of RAD loci had sparse coverage over the individuals, we kept only

251 loci present in more than 100 individuals for this part of analysis. The ML

252 phylogenenies were first calculated for every locus using the GTR+G model as

11

253      implemented in RAxML 8.1.11 (Stamatakis 2014). Then, pairwise geodesic distances

254      between all the resulting single-locus phylogenies were measured, and the trees were

255      grouped based on the distance matrix using spectral clustering (a protocol hereafter

256      referred to as binning). The number of bins was estimated using the nonparametric

257      bootstrapping stopping criterion. Support for each branch in each topology was

258      calculated using aBayes in PhyML (Anisimova et al. 2011). We also analyzed the log-

259      likelihood improvement when analyzing the data with $n+1$ splits vs. $n$ splits, compared

260      to the null expectation (i.e. random loci clustering).

261      Additionally, we applied a coalescent-based inference method using SVDquartets

262      (Chifman and Kubatko 2014) as implemented in PAUP* v.4a150 and v.4a151 (Swofford

263      2002). This method infers the topology among randomly sampled quartets using a

264      coalescent model, and assembles the randomly sampled quartets using a quartet

265      amalgamation method. Breaking the sequence into quartets makes the analysis of large

266      numbers of loci feasible. We randomly sampled the maximum of all possible quartets

267      (i.e. 48,603,900 quartets = 200 taxa) with the multispecies coalescent option and 1,000

268      bootstrap replicates. The quartets were summarized with the QFM (Reaz et al. 2014)

269      quartet amalgamation program as implemented in PAUP*. Phylogenetic trees were

270      visualized using the ape 3.2 R package (Paradis et al. 2004).


271 ## 2.6 Population structure

272      We inferred population structure using the admixture model implemented in

273      STRUCTURE 2.3.4 (Pritchard et al. 2000), without prior population assignment and with

274      allele frequencies correlated among populations. The software uses a Bayesian

275      framework to estimate the likelihood of the data given a number of *a priori* defined *K*

276    population clusters, outputting the likelihood of each sample to belong to each possible

277    cluster. This analysis was performed after removing technical replicates from the

278    dataset, retaining only the loci present in a minimum of 20 individuals, and selecting

279    one random single SNP from each locus. Analyses were run for $K$ values ranging

280    between 1 and 8, with a burn-in of 200K cycles, followed by 1M cycles of sampling, with

281    3 replicates for each $K$ value. The optimal $K$ value was identified following Evanno et al.

282    (2005), as implemented in STRUCTURE HARVESTER (Earl & vonHoldt 2012). To

283    account for the phylogenetic component in the missing alleles, we ran STRUCTURE with

284    the recessive alleles model, with missing data coded as recessive.

# 3. Results

## 3.1 *Chiastocheta* sampling

287    We analyzed a total of 272 *Chiastocheta* specimens sampled from the entire European

288    range of the genus (see Table 1 and Fig. S1 for maps of the sampled specimens). Most

289    species displayed a broad spatial distribution. Up to six species could be found in one

290    single locality during a single visit (Table 1, mean = 2.7 species per locality, SD = 1.5),

291    confirming the sympatric nature of the species and the existing opportunities for

292    hybridization.

## 3.2 Sequencing and RAD loci assembly results

294    After initial screening, 21 samples were removed from the final dataset because of

295    insufficient coverage or technical errors. We successfully analyzed 260 specimens for

296    the mitochondrial dataset (255 for COI, 204 for COII, 141 for the first, and 152 for the

297    second fragment of the D-loop), and 263 for the RAD dataset, while 251 samples were

298  shared between the datasets (Table 1). For the RAD dataset, we also sequenced 22

299  technical replicates (samples with "REPL" suffix in Table S1), and 11 DNA extraction

300  replicates from the fly muscle vs. extractions from legs (samples with "MUS" suffix in

301  Table S1).

302  Sequencing of the mitochondrial regions yielded 1132 nucleotide positions for the

303  COI+COII dataset (of which 120 were variable) and 2003 for the COI+COII+D-loop

304  dataset (of which 334 were variable), after alignment and gap filtering. Three runs of

305  RAD sequencing output 552'425'482 of 2 x 100 bp reads, from which 340'598'636

306  (62%) passed the restriction site and barcode quality filters (Table S1).

307  After comparing the number of loci, coverage, and overlap of loci among replicates in

308  the obtained datasets (Fig. S2), we chose the dataset with a minimum of 75% sequence

309  similarity required for the sequences to cluster in a locus and a minimum of 20

310  individuals per locus for the main analyses. This dataset contained 1724 loci after

311  filtering and paralog removal, with 82'782 variable sites. The proportion of missing data

312  in the dataset was 0.84, with a strong phylogenetic component in the distribution of

313  missing loci (Fig. S3a). After sampling one SNP per locus for the STRUCTURE analysis,

314  we obtained 1669 SNPs, of which 159 were bi-allelic.

315  For the dataset used for assessing loci incongruence in the RAD-seq based phylogeny

316  (see below), we focused on loci present in at least 100 individuals. This resulted in a

317  matrix of 176 loci (among 1724 overall number of loci identified; i.e., 10.2%) with

318  missing data showing much less phylogenetic structuring (Fig. S3b).

## 3.3 Mitochondrial and nuclear-data phylogenies

319

320 The mitochondrial phylogeny on the COI+COII_D-loop dataset (Fig. 1a and Fig. S4a)

321 failed to resolve four of the clades identified based on the RAD-seq data (see below), but

322 retrieved well-supported monophyletic group for *C. rotundiventris* and, to the lesser

323 extent for *C. dentifera* and *C. inermella*, as both of the latter had two specimens placed

324 outside their clades. *C. inermella*, *C. setifera*, and *C. trollii* formed one clade with the

325 species extensively interdispersed and a clade containing mostly *C. macropyga* nested

326 within. As the analysis based based on the reduced COI+COII dataset recovered a similar

327 pattern, except placing *C. lophota* as sister to *C. macropyga*, we refer to the results of the

328 larger dataset in the rest of the paper. Most of the *C. lophota* samples also formed one

329 clade with lower support values. The relationships among samples from the remaining

330 four morphospecies remained unresolved, without clear support for the

331 morphologically described species.

332 In contrast to the ML mtDNA phylogeny, both ML and SVDquartets analysis of the RAD

333 analysis (Fig. 1b, c and Fig. S4b, c) confirmed monophyly of the seven morphologically

334 defined taxa. RAxML analysis revealed relatively high bootstrap supports (> 90%) for

335 all of the interspecific relationships, except the split between *C. setifera* and the clade (*C.*

336 *lophota*, (*C. macropyga*, (*C. dentifera*, *C. trollii*))) with bootstrap support > 80%. The split

337 of *C. rotundiventris* into two putative vicariant clades, informally proposed by Pellmyr

338 (1992) — northern *C. abruptriventris* and southern *C. rotundiventris*, was not recovered.

339 SVDquartets analysis also confirmed monophyly of the species, but only the split

340 between *C. dentifera* and *C. trollii* had a bootstrap support > 90%; the clade (*C.*

341 *macropyga*, (*C. dentifera*, *C. trollii*)) had bootstrap support > 80%; these two clades were

342 the only ones supported by both SVDquartets and the RAxML analyses (Fig. 1c).

15

343    Moreover, SVDquartets revealed two well-supported clades within *C. rotundiventris*.

344    These however do not show any pattern of vicariance and often occur together in a

345    single population, thus most likely do not correspond to the two vicariant species of *C.*

346    *abruptiventris* and *C. rotundiventris* as discussed by Pellmyr (1992). The technical

347    replicates were consistent in the placement of the sample within the proper clade, and

348    most replicates were placed as sister clades with both methods (Fig. S4b,c).

## 3.4 Incongruence among the RAD-sequencing loci

350    TreeCl analysis identified, in the most conservative interpretation, at least four clusters

351    of loci, as the largest likelihood improvement was obtained when increasing the number

352    of bins from three to four (Fig. S6). The bin sizes were of 29, 42, 47, and 58 loci,

353    therefore the identified groups were not simply consisting of a few outliers. The trees

354    inferred for the four bins confirmed the monophyly of the analysed species to a large

355    extent, although few individuals appeared outside their expected clades. The largest

356    departure from monophyly was observed for *C. lophota* in the smallest tree consisting of

357    29 loci (Fig. 2). The trees inferred for each cluster had branch supports for interspecific

358    nodes larger than 95%, and differed substantially in terms of topology and branch

359    lengths. Only one tree, with the largest number of loci (i.e., 58) supported the only clade

360    that was supported by both RAML and SVDquartets analysis (*C. macropyga*, (*C.*

361    *dentifera*, *C. trollii*)). Except that, the interspecific relationships retrieved with each of

362    the treeCl bins were different than with the concatenated RAxML analysis and

363    SVDquartets analysis (Fig. 1b, c).

16

## 3.5 Structure analysis

We found low levels of contemporary introgression, as shown by STRUCTURE analysis. The most likely *K* number of STRUCTURE groups was consistent with the number of morphological species (7), and all samples were assigned to their 'correct' morphospecies (Fig. 1d). Also for lower numbers of clusters, we did not observe signatures of introgression (Fig. S5).

# 4. Discussion

## 4.1 Utility and limits of RAD-sequencing for resolving phylogeny of a „difficult" genus

RAD-sequencing successfully discriminated all formally described European *Chiastocheta* species. The robust species delineation is striking when compared to mtDNA-based trees that failed to support monophyly of *C. inermella, C. macropyga*, *C. setifera,* and *C. trollii* (Fig. 1a and Fig. S4a; see also: Després et al. 2002; Espíndola et al. 2012). The ability to recover previously defined morphological species in our dataset, whatever analysis method used (i.e., maximum-likelihood phylogenetic reconstruction using a concatenated matrix with RAxML, coalescence-based phylogenetic inference with SVDquartets, or population-genetics clustering with STRUCTURE), supports the results of a previous simulation study by Hovmöller et al. (2013), that high amounts of missing data, typical for RAD-based datasets, should not interfere with clade (or cluster) identification. Recently, similar conclusions were drawn by Eaton et al. (2016) concerning the SVDquartets method.

385    In contrast, no consensus could be reached in retrieving inter-specific relationships.

386    Whereas RAxML identified relationships with high bootstrap support in four of the five

387    possible interspecific relationships, only two of them were also supported by the

388    SVDquartets analysis (Fig. 1b,c). Incongruence in the phylogenetic signals associated

389    with different sets of loci could explain the difficulty in resolving these interspecific

390    relationships. When performing loci binning using treeCL (Gori et al. 2016), we found

391    out that different subsets of loci (in our case, the optimal number of bins was equal to

392    four) produced different topologies, while still being largely congruent in the sample

393    assignment into species (Fig. 2).

394    Short interspecific branches in the resolved phylogenies confirm the conclusions of

395    Espíndola et al. (2012) that most of the species from the *Chiastocheta* genus underwent

396    a recent (less than 1.6 Mya), rapid radiation. These results highlight the fact that in such

397    cases it may be impossible to retrieve some of the phylogenetic relationships among the

398    taxa as fully bifurcating tree, because gene trees may depict different evolutionary

399    histories due to incomplete lineage sorting (Avise et al. 2008; Maddison 1997). This is a

400    limitation shared with classical markers (Walsh et al. 1999) and other NGS approaches

401    (see below), pointing to a possible constitutive limitation in resolving rapid radiations.

402    In rapidly diverging taxa, even the large number of nuclear markers, while being more

403    successful here in recovering species boundaries than mitochondrial markers may not

404    be informative-enough to retrieve all interspecific evolutionary relationships.

405    The extent to which the above limitation is the result of technical constraints of RAD

406    datasets or a true biological limitation remains to be investigated. RAD-seq targets

407    random, mostly neutral parts of the genome. This results in high number of lineage-

408    specific mutations that bear a strong signal to delineate species or populations – within

18

409  these fast-evolving parts of the genome, even varying allele copy-numbers (i.e. recent

410  paralogs) can appear as population-specific (Mastretta-Yanes et al. 2014). The

411  downside is however, missing data increases rapidly with evolutionary distance as a

412  result of the loss of restriction sites  (Cariou et al. 2013; Chattopadhyay et al. 2014;

413  DaCosta & Sorenson 2016; Gautier et al. 2013; Rubin et al. 2012; Wagner et al. 2013).

414  For instance, Leaché et al. (2015) found differences between phylogenies obtained

415  using RAD-seq vs. target enrichment techniques, whereas other studies have shown the

416  agreement among data types (Manthey et al. 2016). The latter techniques rely on

417  capture of a predefined (Faircloth et al. 2014; McCormack et al. 2012) or random

418  (Suchan et al. 2016, Schmid et al. 2017) subset of loci. By not relying on the presence of

419  restriction sites, and thus having less missing data, enrichment techniques may be

420  better suited for broader phylogenetic scales.

421  Nevertheless, it has been shown that even with hundreds of conserved loci, known

422  substitution models and several individuals per species, trees with short branches are

423  difficult to resolve, and ML analyses based on concatenated sequences may provide high

424  bootstrap values despite incorrectly resolved topologies (Giarla & Esselstyn 2015;

425  Kubatko & Degnan 2007; but see Gatesy & Springer 2013; Springer & Gatesy 2016; Roch

426  & Warnow 2015). This is exemplified by our study, in which using all RAD loci, we

427  obtained a ML phylogeny with highly supported interspecific nodes, whereas

428  coalescence-based phylogenetic inference did not show strong supports for most of the

429  interspecific relationships. Our exploration of explanations for such a discrepancy using

430  the loci binning approach showed support for at least four different underlying gene

431  tree topologies. In these analyses, we reduced the dataset to a non-random set of loci

432  when filtering for high loci coverage among samples. The retained loci, present in at

433 least 100 analyzed individuals, and with less phylogenetically-structured missing data

434 (see Fig. S3b), should be characterized by lower mutation rates or being under

435 stabilizing selection (Huang & Knowles 2014). Using binning, the best fit to the data was

436 not obtained with a single bin of loci but with four. We could therefore not identify one

437 single evolutionary history of the *Chiastocheta* genus, but rather equally-supported

438 gene trees topologies. Importantly, these different topologies cannot be attributed to a

439 few outlier loci, as their distribution was relatively even across the clusters (29, 58, 42

440 and 47 loci; Fig. S6), incongruence among these sets possibly impacting maximum-

441 likelihood phylogenetic reconstruction using a concatenated matrix and coalescence-

442 based phylogenetic inference. We have also confirmed that in such cases, ML methods

443 provide elevated bootstrap support values, and that lower bootstrap support values

444 resulting from coalescence-based methods may better reflect the biological uncertainty

445 of interspecific relationships.

## 4.2 Mitonuclear discordance in the phylogeny of *Chiastocheta*

446

447 While our RAD-sequence dataset delineated seven clades, with full agreement with the

448 morphological assignments, mitochondrial data failed to support species monophyly,

449 except for *C. rotundiventris* and, to a lesser extent, *C. lophota* and *C. dentifera*. The other

450 remaining species: *C. inermella*, *C. setifera* and *C. trollii* formed a large clade with the

451 species extensively interdispersed and with the clade consisting mostly of *C. macropyga*

452 nested within (Fig. 1a). Despite, on average, mitochondrial markers should be more

453 suited for capturing relationships among recently diverged lineages, due to an effective

454 population size four times less than that of nuclear genes (assuming neutral processes,

455 equal sex ratios, and unbiased mating systems), and thus shorter coalescence times

456    (Zink & Barrowclough 2008), analyzing a large dataset of nuclear markers provided

457    more power to discriminate the species in our case.

458    Mitonuclear discordance patterns can be explained either by the different biological

459    properties of mitochondrial DNA (vegetative segregation, uniparental inheritance,

460    intracellular selection, and reduced recombination; Birky 2001) or differences in the

461    evolutionary histories of nuclear and mitochondrial markers [e.g., direct selection on

462    the mitochondrial genes (Ballard et al. 2007; Ballard & Pichaud 2014; Boratyński et al.

463    2014; Dowling et al. 2008), incomplete lineage sorting, historical or ongoing gene flow

464    among species, or hybrid speciation].  Indeed, it has been shown before that relying on a

465    single, non-recombining mtDNA molecule may provide a misleading account of a

466    species history (e.g., Ballard & Whitlock 2004; Govindarajulu et al. 2015; Phillips et al.

467    2013; Seehausen et al. 2003; and reviews by Funk & Omland 2003; Rubinoff & Holland

468    2005).  While investigating the reasons for the mito-nuclear discordance was not within

469    the scope of this paper, we could reject the hypothesis of a contemporary gene flow or

470    hybrid origin of the taxa as responsible for this pattern. We did not detect signature of a

471    genetic mosaic in the—mostly—nuclear RAD data, which would be expected in the case

472    of hybrid origin (Ballard 2000; Brelsford et al. 2011; Mallet 2007; Pollard et al. 2006).

473    Using RAD-sequencing data, the assignment of samples into species was concordant

474    with morphology (Fig. 1b,c and S4b,c) and we did not detect significant levels of

475    contemporary gene flow using population genetics-based approaches (Fig. 1d), despite

476    apparent opportunities for hybridization. Most of *Chiastocheta* occur in sympatry (Fig.

477    S1), they also have very similar biologies, reproducing and spending most of their time

478    on or inside flowers of *Trollius europaeus* (Pellmyr 1989; Suchan et al. 2015). Although a

479    temporal sequence in oviposition has been observed (Després & Jaeger 1999;

480   Johannesen & Loeschcke 1996; Pellmyr 1989), most species co-occur temporally.

481   Despite these ecological similarities and the relatively young age of the genus (most of

482   the clades emerging less than 1.6 Ma; Espindola et al. 2012), a lack of nuclear evidence

483   for hybridization indicates strong contemporary reproductive barriers among the

484   species.

# 485   5. Conclusions

486   This study demonstrates how a combination of RAD-seq and mtDNA data can provide

487   insights into phylogenies of genera that are poorly resolved using mitochondrial

488   markers alone and reveal complex picture of mitonuclear discordance. It also

489   underlines the limits of RAD-seq-based phylogenies in case of rapid radiations. Our

490   results show that a scenario of rapid radiation can affect many loci across the genome,

491   leading to discordant gene trees, even when using methods controlling for incomplete

492   lineage sorting. This may point to an inherent limitation of using molecular markers to

493   resolve rapid radiations, at least at some of the inter-specific relationships, and suggests

494   that this limitation is not necessarily due to technical issues (e.g. low number of shared

495   markers).

496   Adding to the body of examples of mito-nuclear discordance (reviewed in Toews &

497   Brelsford 2012), our study warns against relying solely on mitochondrial markers (e.g.,

498   COI barcoding; Herbert et al. 2003) for species delimitation, especially when they show

499   incongruence with classical taxonomy. In the case presented here, mitochondrial

500   markers suggested poly- or paraphyly for most species, and proposed the need to

501   review the taxonomy of the genus (Espíndola et al. 2012). When tackled from the

502   genomic point of view, the genetic support of species status for these seven entities was

503    confirmed. Finally, we provide an example of how ML phylogenies based on large

504    concatenated datasets can provide erroneously high bootstrap supports for incorrect or

505    uncertain topologies (Giarla & Esselstyn 2015; Kubatko & Degnan 2007).

506    # Acknowledgments

# References

520 

521 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the

522 power of RADseq for ecological and evolutionary genomics. *Nature Reviews*

523 *Genetics*, 17, 81-92.

524 Anisimova M, Gil M, Dufayard J-F, Dessimoz C, Gascuel O (2011) Survey of Branch

525 Support Methods Demonstrates Accuracy, Power, and Robustness of Fast

526 Likelihood-based Approximation Schemes. *Systematic Biology*, 60, 685–699.

527 Avise JC, Robinson TJ, Kubatko L (2008) Hemiplasy: a new term in the lexicon of

528 phylogenetics. *Systematic Biology*, 57(3), 503-507.

529 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA,

530 Johnson EA (2008) Rapid SNP discovery and genetic mapping using sequenced RAD

531 markers. *PLoS ONE*, 3, e3376.

532 Ballard JWO (2000) When one is not enough: introgression of mitochondrial DNA in

533 *Drosophila*. *Molecular Biology and Evolution*, 17, 1126-1130.

534 Ballard JWO, Melvin RG, Katewa SD, Maas K (2007) Mitochondrial DNA variation is

535 associated with measurable differences in life-history traits and mitochondrial

536 metabolism in *Drosophila simulans. Evolution*, 61, 1735-1747.

537 Ballard JWO, Pichaud N (2014) Mitochondrial DNA: more than an evolutionary

538 bystander. *Functional Ecology*, 28, 218-231.

539 Ballard JWO, Whitlock MC (2004) The incomplete natural history of mitochondria.

540 *Molecular Ecology*, 13, 729-744.

24

541 Birky CW (2001) The inheritance of genes in mitochondria and chloroplasts: laws,

542     mechanisms, and models. *Annual Review of Genetics*, 35, 125-48.

543 Boratyński Z, Melo-Ferreira J, Alves PC, Berto S, Koskela E, Pentikäinen OT, Mappes T

544     (2014) Molecular and ecological signs of mitochondrial adaptation: consequences

545     for introgression. *Heredity*, 113, 277-286.

546 Brelsford A, Milá B, Irwin DE (2011) Hybrid origin of Audubon's warbler. *Molecular*

547     *Ecology*, 20, 2380-2389.

548 Cariou M, Duret L, Charlat S (2013) Is RAD-seq suitable for phylogenetic inference? An

549     *in silico* assessment and optimization. *Ecology and Evolution*, 3, 846-852.

550 Chattopadhyay B, Garg KM, Ramakrishnan U (2014) Effect of diversity and missing data

551     on genetic assignment with RAD-Seq markers. *BMC research notes*, 7, 841.

552 Chifman J, Kubatko L (2014) Quartet inference from SNP data under the coalescent

553     model. *Bioinformatics*, 30:3317-3324.

554 Collin JE (1954) The genus *Chiastocheta* Pokorny (Diptera: Anthomyiidae). *Proceedings*

555     *of the Royal Entomological Society London (B)*, 23, 95-102.

556 Combosch DJ, Vollmer SV (2015) Trans-Pacific RAD-Seq population genomics confirms

557     introgressive hybridization in Eastern Pacific *Pocillopora* corals. *Molecular*

558     *Phylogenetics and Evolution*, 88, 154-162.

559 Cruaud A, Gautier M, Galan M, Foucaud J, Sauné L, Genson G, Rasplus JY (2014)

560     Empirical assessment of RAD sequencing for interspecific phylogeny. *Molecular*

561     *Biology and Evolution*, 31, 1272-1274.

562 DaCosta JM, Sorenson MD (2016) ddRAD-seq phylogenetics based on nucleotide, indel,

563    and presence–absence polymorphisms: Analyses of two avian genera with

564    contrasting histories. *Molecular Phylogenetics and Evolution*, 94, 122-135.

565 Davey JL, Blaxter MW (2010) RADSeq: Next-generation population genetics. *Briefings in*

566    *Functional Genomics*, 9, 416-423.

567 Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML (2011) Genome-

568    wide genetic marker discovery and genotyping using next-generation sequencing.

569    *Nature Reviews Genetics*, 12, 499-510.

570 DeFilippis VR, Moore WS (2000) Resolution of phylogenetic relationships among

571    recently evolved species as a function of amount of DNA sequence: An empirical

572    study based on woodpeckers (Aves: Picidae). *Molecular Phylogenetics and*

573    *Evolution*, 16, 143-160.

574 Degnan, J.H., Rosenberg, N.A., 2009. Gene tree discordance, phylogenetic inference and

575    the multispecies coalescent. *Trends in Ecology and Evolution* 24, 332–340.

576 Després L, Jaeger N (1999) Evolution of oviposition strategies and speciation in the

577    globeflower flies *Chiastocheta* spp. (Anthomyiidae). *Journal of Evolutionary Biology*,

578    12, 822-831.

579 Després L, Pettex E, Plaisance V, Pompanon F (2002) Speciation in the globeflower fly

580    *Chiastocheta* spp. (Diptera: Anthomyiidae) in relation to host plant species,

581    biogeography, and morphology. *Molecular Phylogenetics and Evolution*, 22, 258-268.

582 Dowling DK, Friberg U, Lindell J (2008) Evolutionary implications of non-neutral

583    mitochondrial genetic variation. *Trends in Ecology and Evolution*, 23, 546-554.

584 Earl DA, vonHoldt BM (2012) STRUCTURE HARVESTER: a website and program for

585     visualizing STRUCTURE output and implementing the Evanno method. *Conservation*

586     *Genetics Resources*, 4, 359-361.

587 Eaton DAR (2014) PyRAD: assembly of de novo RADseq loci for phylogenetic analyses.

588     *Bioinformatics*, 30, 1844-1849.

589 Eaton DAR, Ree RH (2013) Inferring phylogeny and introgression using genomic

590     RADseq data: An example from flowering plants (*Pedicularis*: Orobanchaceae).

591     *Systematic Biology*, 62, 689-706.

592 Eaton, D. A., Spriggs, E. L., Park, B., & Donoghue, M. J. (2016). Misconceptions on Missing

593     Data in RAD-seq Phylogenetics with a Deep-scale Example from Flowering Plants.

594     *Systematic Biology*, syw092.

595 Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high

596     throughput. *Nucleic Acids Research* 32, 1792-1797.

597 Escudero M, Eaton DA, Hahn M, Hipp AL (2014) Genotyping-by-sequencing as a tool to

598     infer phylogeny and ancestral hybridization: A case study in *Carex* (Cyperaceae).

599     *Molecular Phylogenetics and Evolution*, 79, 359-367.

600 Espíndola A, Buerki S, Alvarez N (2012) Ecological and historical drivers of

601     diversification in the fly genus *Chiastocheta* Pokorny. *Molecular Phylogenetics and*

602     *Evolution*, 63, 466-474.

603 Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals

604     using the software STRUCTURE: a simulation study. *Molecular Ecology*, 14, 2611-

605     2620.

606 Faircloth BC, Branstetter MG, White ND, Brady SG (2014) Target enrichment of

607    ultraconserved elements from arthropods provides a genomic perspective on

608    relationships among Hymenoptera. *Molecular Ecology Resources*, 15, 489-501

609 Funk DJ, Omland KE (2003) Species-level paraphyly and polyphyly: frequency, causes,

610    and consequences, with insights from animal mitochondrial DNA. *Annual Review of*

611    *Ecology, Evolution and Systematics*, 34, 397-423.

612 Gatesy, J, Springer MS (2013) Concatenation versus coalescence versus

613    "concatalescence". *Proceedings of the National Academy of Sciences*, 110, E1179-

614    E1179.

615 Gautier M, Gharbi K, Cezard T, Foucaud J, Kerdelhué C, Pudlo P, Cornuet J-M, Estoup A

616    (2013) The effect of RAD allele dropout on the estimation of genetic variation

617    within and between populations. *Molecular Ecology*, *22*, 3165-3178.

618 Giarla TC, Esselstyn JA (2015) The challenges of resolving a rapid, recent radiation:

619    Empirical and simulated phylogenomics of philippine shrews. *Systematic Biology*,

620    64, 727-740.

621 Gori K, Suchan T, Alvarez N, Goldman N, Dessimoz C (2016) Clustering genes of common

622    evolutionary history. *Molecular Biology and Evolution*, 33,1590–1605

623 Govindarajulu R, Parks M, Tennessen JA, Liston A, Ashman TL (2015) Comparison of

624    nuclear, plastid, and mitochondrial phylogenies and the origin of wild octoploid

625    strawberry species. *American Journal of Botany*, 102, 544-554.

626 Harvey MG, Judy CD, Seeholzer GF, Maley JM, Graves GR, Brumfield RT (2015) Similarity

627      thresholds used in DNA sequence assembly from short reads can reduce the

628      comparability of population histories across species. *PeerJ*, 3, e895.

629 Harvey MG, Smith BT, Glenn TC, Faircloth BC, Brumfield RT (2016) Sequence capture

630      versus restriction site associated DNA sequencing for shallow systematics.

631      *Systematic Biology*, 65, 910-924.

632 Hebert PDN, Cywinska A, Ball SL, deWaard JR (2003) Biological identifications through

633      DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences*, 270, 313-322.

634 Hennig W (Ed.) (1976) Anthomyiidae. Die Fliegen der Palaearktischen Region. Stuttgart,

635      E. Schweizerbart.

636 Herrera S, Shank TM (2016) RAD sequencing enables unprecedented phylogenetic

637      resolution and objective species delimitation in recalcitrant divergent taxa.

638      *Molecular Phylogenetics and Evolution*, 100, 70-79.

639 Hipp AL, Eaton DAR, Cavender-Bares J, Fitzek E, Nipper R, Manos PS (2014) A

640      framework phylogeny of the american oak clade based on sequenced RAD data.

641      *PLoS ONE*, 9, e93975.

642 Hovmöller R, Knowles LL, Kubatko LS (2013) Effects of missing data on species tree

643      estimation under the coalescent. *Molecular Phylogenetics and Evolution*, 69, 1057-

644      1062.

645 Huang H, Knowles LL (2014) Unforeseen consequences of excluding missing data from

646      next-generation sequences: simulation study of RAD sequences. *Systematic Biology*,

647      doi:10.1093/sysbio/syu046.

648    Ilut DC, Nydam ML, Hare MP (2014) Defining loci in restriction-based reduced

649        representation genomic data from nonmodel species: Sources of bias and

650        diagnostics for optimal clustering. *BioMed Research International*, 2014, 675158.

651    Johannesen J, Loeschcke V (1996) Distribution, abundance and oviposition patterns of

652        four coexisting *Chiastocheta* species (Diptera: Anthomyiidae). *Journal of Animal*

653        *Ecology*, 65, 567-576.

654    Jones JC, Fan S, Franchini P, Schartl M, Meyer A (2013) The evolutionary history of

655        *Xiphophorus* fish and their sexually selected sword: a genome-wide approach using

656        restriction site-associated DNA sequencing. *Molecular Ecology*, 22, 2986-3001.

657    Kubatko LS, Degnan JH (2007) Inconsistency of phylogenetic estimates from

658        concatenated data under coalescence. *Systematic Biology*, 56, 17-24.

659    Lanfear R, Frandsen PB, Wright AM, Senfeld T, Calcott B (2017) PartitionFinder 2: new

660        methods for selecting partitioned models of evolution for molecular and

661        morphological phylogenetic analyses. *Molecular Biology and Evolution.* 34, 772-773.

662    Leaché AD, Chavez AS, Jones LN, Grummer JA, Gottscho AD, Linkem CW (2015)

663        Phylogenomics of phrynosomatid lizards: conflicting signals from sequence capture

664        versus restriction site associated DNA sequencing. *Genome Biology and Evolution*, 7,

665        706-719.

666    Leaché AD, Harris RB, Rannala B, Yang Z (2013) The influence of gene flow on species

667        tree estimation: a simulation study. *Systematic Biology*, 63, 17-30.

668    Lowry DB, Hoban S, Kelley JL, Lotterhos KE, Reed LK, Antolin MF, Storfer A (2017)

669        Breaking RAD: an evaluation of the utility of restriction site-associated DNA

670 sequencing for genome scans of adaptation. *Molecular Ecology Resources*, 17, 142–

671 152.

672 Maddison WP (1997) Gene trees in species trees. *Systematic Biology*, 46, 523-536.

673 Maddison WP, Knowles LL (2006) Inferring phylogeny despite incomplete lineage

674 sorting. *Systematic Biology*, 55, 21-30.

675 Mallet J (2007) Hybrid speciation. *Nature*, 446, 279-283.

676 Manthey JD, Campillo LC, Burns KJ, Moyle RG (2016) Comparison of target-capture and

677 restriction-site associated DNA sequencing for phylogenomics: a test in cardinalid

678 tanagers (Aves, Genus: *Piranga*). *Systematic biology*, syw005.

679 Mastretta-Yanes A, Zamudio S, Jorgensen TH, Arrigo N, Alvarez N, Piñero D, Emerson BC

680 (2014) Gene duplication, population genomics, and species-level differentiation

681 within a tropical mountain shrub. *Genome biology and evolution*, 6, 2611-2624.

682 Mastretta-Yanes A, Arrigo N, Alvarez N, Jorgensen TH, Piñero D, Emerson BC (2015)

683 Restriction site-associated DNA sequencing, genotyping error estimation and de

684 novo assembly optimization for population genetic inference. *Molecular Ecology*

685 *Resources*, 15, 28-41.

686 McCormack JE, Faircloth BC, Crawford NG, Gowaty PA, Brumfield RT, Glenn TC (2012)

687 Ultraconserved elements are novel phylogenomic markers that resolve placental

688 mammal phylogeny when combined with species tree analysis. *Genome Research*,

689 22, 746-754.

690 McKinney GJ, Larson WA, Seeb LW, Seeb JE (2017) RADseq provides unprecedented

691  insights into molecular ecology and evolutionary genetics: comment on Breaking

692  RAD by Lowry *et al.* (2016). *Molecular Ecology Resources*, doi:10.1111/1755-

693  0998.12649.

694 Michelsen V (1985) A revision of the Anthomyiidae (Diptera) described by J.W.

695  Zetterstedt. *Steenstrupia*, 11, 37-65

696 Miller MA, Pfeiffer W, Schwartz T (2010) Creating the CIPRES Science Gateway for

697  inference of large phylogenetic trees. In: Proceedings of the Gateway Computing

698  Environments Workshop (GCE), 14 Nov. 2010, New Orleans, LA pp. 1-8.

699 Nadeau NJ, Martin SH, Kozak KM, Salazar C, Dasmahapatra K, Davey JW, Baxter SW,

700  Blaxter ML, Mallet J, Jiggins CD (2013) Genome-wide patterns of divergence and

701  gene flow across a butterfly radiation. *Molecular Ecology*, 22, 814–826.

702 Pante E, Abdelkrim J, Viricel A, Gey D, France SC, Boisselier MC, Samadi S (2015) Use of

703  RAD sequencing for delimiting species. *Heredity*, 114, 450-459.

704 Paradis E, Claude J, Strimmer K (2004) APE: analyses of phylogenetics and evolution in

705  R language. *Bioinformatics*, 20, 289-290.

706 Pellmyr O (1989) The cost of mutualism – interactions between *Trollius europaeus* and

707  its pollinating parasites. *Oecologia*, 78, 53-59.

708 Pellmyr O (1992) The phylogeny of a mutualism: evolution and coadaptation between

709  *Trollius* and its seed-parasitic pollinators. *Biological Journal of the Linnean Society*,

710  47, 337-365.

711    Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double digest RADseq:

712        An inexpensive method for de novo SNP discovery and genotyping in model and

713        non-model species. *PLoS ONE*, 7, e37135.

714    Phillips MJ, Haouchar D, Pratt RC, Gibb GC, Bunce M (2013) Inferring Kangaroo

715        Phylogeny from Incongruent Nuclear and Mitochondrial Genes. *PLoS ONE*, 8,

716        e57745.

717    Pollard DA, Iyer VN, Moses AM, Eisen MB (2006) Widespread discordance of gene trees

718        with species tree in *Drosophila*: Evidence for incomplete lineage sorting. *PLoS*

719        *Genetics*, 2, e173.

720    Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using

721        multilocus genotype data. *Genetics*, 155, 945-959.

722    Reaz R, Bayzid MS, Rahman MS (2014) Accurate phylogenetic tree reconstruction from

723        quartets: A heuristic approach. *PloS One*, 9:e104008.

724    Roch S, Warnow T (2015) On the robustness to gene tree estimation error (or lack

725        thereof) of coalescent-based species tree methods. *Systematic Biology*, syv016.

726    Rokas A, Carroll SB (2005) More genes or more taxa? The relative contribution of gene

727        number and taxon number to phylogenetic accuracy. *Molecular Biology and*

728        *Evolution*, 22, 1337-1344.

729    Rubin BER, Ree RH, Moreau CS (2012) Inferring phylogenies from RAD sequence data.

730        *PLoS ONE*, 7, e33394.

731  Rubinoff D, Holland BS (2005) Between two extremes: mitochondrial DNA is neither the

732      panacea nor the nemesis of phylogenetic and taxonomic inference. *Systematic*

733      *Biology*, 54, 952-961.

734  Schmid S, Genevest R, Gobet E, Suchan T, Sperisen C, Tinner W, Alvarez N (2017)

735      HyRAD-X, a versatile method combining exome capture and RAD sequencing to

736      extract genomic information from ancient DNA. *Methods in Ecology and Evolution*,

737      doi:10.1111/2041-210X.12785

738  Seehausen O, Koetsier E, Schneider MV, Chapman LJ, Chapman CA, Knight ME, Turner

739      GF, van Alphen JJM, Bills R (2003) Nuclear markers reveal unexpected genetic

740      variation and a Congolese-Nilotic origin of the Lake Victoria cichlid species flock.

741      *Proceedings of the Royal Society of London Series B*, 270, 129-137.

742  Springer MS, Gatesy J (2016) The gene tree delusion. *Molecular Phylogenetics and*

743      *Evolution*, 94, 1-33.

744  Stamatakis A (2014) RAxML Version 8: A tool for phylogenetic analysis and post-

745      analysis of large phylogenies. *Bioinformatics*, 30, 1312-1313.

746  Suchan T, Beauverd M, Trim N, Alvarez N (2015) Asymmetrical nature of the *Trollius-*

747      *Chiastocheta* interaction: insights into the evolution of nursery pollination systems.

748      *Ecology and Evolution*, 5, 4766-4777.

749  Suchan T, Pitteloud C, Gerasimova NS, Kostikova A, Schmid S, Arrigo N, Pajkovic M,

750      Ronikier M, Alvarez N (2016) Hybridization Capture Using RAD Probes (hyRAD), a

751      New Tool for Performing Genomic Analyses on Collection Specimens. *PLoS ONE*, 11,

752      e0151651.

753 Suchan T, Espíndola A, Rutschmann S, Emerson BC, Gori K, Dessimoz C, Arrigo N, Ronikier

754    M, Alvarez N (2017) Data from: Assessing the potential of RAD-sequencing to resolve phy-

755    logenetic relationships within species radiations: the fly genus *Chiastocheta* (Diptera: An-

756    thomyiidae) as a case study. *Mendeley Data*, https://doi.org/XXX

757 Swofford DL (2002) PAUP*: Phylogenetic analysis using parsimony (*and other

758    methods). Version 4. Sinauer, Sunderland, Massachusetts, USA.

759 Toews DP, Brelsford A (2012) The biogeography of mitochondrial and nuclear

760    discordance in animals. *Molecular Ecology*, 21, 3907-3930.

761 Townsend TM, Mulcahy DG, Noonan BP, Sites JW, Kuczynski CA, Wiens JJ, Reeder TW

762    (2011) Phylogeny of iguanian lizards inferred from 29 nuclear loci, and a

763    comparison of concatenated and species-tree approaches for an ancient, rapid

764    radiation. *Molecular Phylogenetics and Evolution*, 61, 363-380.

765 Wagner CE, Keller I, Wittwer S, Selz OM, Mwaiko S, Greuter L, Sivasundar A, Seehausen

766    O (2013) Genome-wide RAD sequence data provide unprecedented resolution of

767    species boundaries and relationships in the Lake Victoria cichlid adaptive radiation.

768    *Molecular Ecology*, 22, 787–798.

769 Walsh HE, Kidd MG, Moum T, Friesen VL (1999) Polytomies and the power of

770    phylogenetic inference. *Evolution*, 53, 932-937.

771 Whitfield JB, Kjer KM (2008) Ancient rapid radiations of insects: challenges for

772    phylogenetic analysis. *Annual Review of Entomology*, 53, 449-472.

773 Whitfield JB, Lockhart PJ (2007) Deciphering ancient rapid radiations. Trends in

774    Ecology and Evolution, 22, 258-265.

775 Wielstra B, Arntzen JW, van der Gaag KJ, Pabijan M, Babik W (2014) Data Concatenation,

776     Bayesian Concordance and Coalescent-Based Analyses of the Species Tree for the

777     Rapid Radiation of *Triturus* Newts. *PLoS ONE*, 9, e111011.

778 Williams JS, Niedzwiecki JH, Weisrock DW (2013) Species tree reconstruction of a

779     poorly resolved clade of salamanders (Ambystomatidae) using multiple nuclear

780     loci. *Molecular Phylogenetics and Evolution*, 68, 671-682.

781 Zetterstedt JW (1845) Diptera scandinaviae disposita et descripta. IV, 1281-1738. Lund,

782     Sweden.

783 Zink RM, Barrowclough GF (2008) Mitochondrial DNA under siege in avian

784     phylogeography. *Molecular Ecology*, 17, 2107-2121.

## 785 Data Accessibility

786 DNA sequences are available on Genbank under accessions no XXX-XXX. Nexus and

787 STRUCTURE datasets used for the analyses, ML phylogeny inferred for the main RAD

788 dataset, SVDquartet analyses, and ML phylogeny inferred for the mtDNA dataset are

789 available on Mendeley Data https://doi .org/XXX (Suchan et al. 2017).

## 790 Author Contributions

791 TS, NA, AE, KG and CD designed research, TS, AE, KG, and SR performed research and

792 analyzed data. All authors wrote the paper.

# Figures



Figure 1. Phylogenies obtained for a) ML analysis of the mtDNA dataset; b) ML analysis of the RAD dataset; c) SVDquartets analysis of RAD dataset; bootstrap node supports > 80 are shown denoted by gray points, bootstrap node supports > 90 are shown denoted by black points. d) Population clustering of the sampled *Chiastocheta* specimens, estimated with STRUCTURE using $K = 7$ value.

37

801

802

803    Figure 2. Phylogenetic trees on the four bins, as identified by the treeCl analysis,

804    considering only the loci present in at least 100 specimens. Bootstrap node supports >

805    80 are shown denoted by gray points, bootstrap node supports > 90 are shown denoted

806    by black points.

# Tables

Table 1. Populations included in the study, with geographical coordinates and the number of specimens used in the final analyses. Letter codes denote *Chiastocheta* species: *C. dentifera* (D), *C. inermella* (I), *C. lophota* (L), *C. macropyga* (M), *C. rotundiventris* (R), *C. setifera* (S), and *C. trollii* (T).

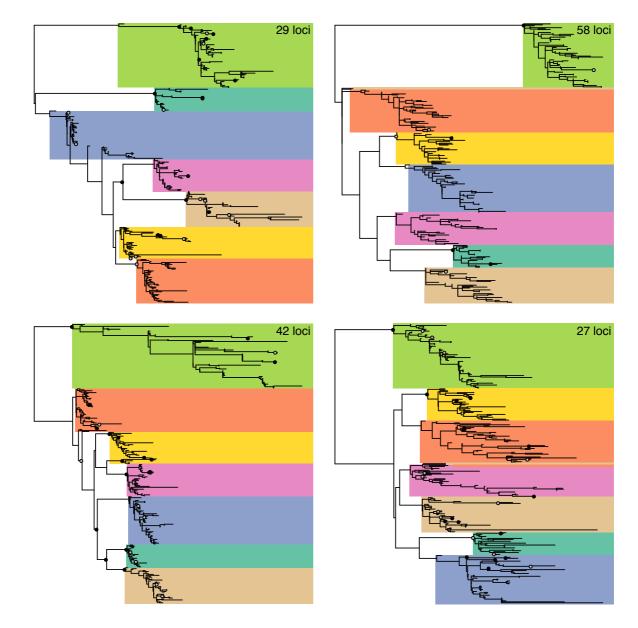| code | site | latitude | longitude | year | D | I | L | M | R | S | T | sum |
|------|------|----------|-----------|------|---|---|---|---|---|---|---|-----|
| AMB | Ambri | 46.50680 | 8.70292 | 2008 | | | 2 | 4 | 2 | | 1 | 9 |
| AMO | Amot | 59.62199 | 8.42346 | 2007 | | 4 | | | | | | 4 |
| BAY | Bayasse | 44.30814 | 6.74067 | 2007 | | 1 | 2 | 2 | 2 | | 3 | 10 |
| BEI | Beistohlen | 61.20761 | 8.95473 | 2007 | | 5 | | | | | | 5 |
| BID | Bidjovagge | 69.29778 | 22.47808 | 2008 | | 1 | | | 1 | | | 2 |
| BON | Col de Bonnecombe | 44.57557 | 3.11410 | 2007 | | | 3 | 1 | 1 | | 2 | 7 |
| BRA | Braas | 57.09309 | 15.06817 | 2007 | | 1 | | | | | | 1 |
| CCO | Col de la Co-lombière | 45.98722 | 6.46972 | 2006 | | | 2 | 1 | 2 | | 1 | 6 |
| CDV | Creux du Van | 46.93526 | 6.74119 | 2006 | | 1 | 2 | | 2 | 2 | | 7 |
| CHA | Chasseral | 47.12569 | 7.02130 | 2006 | | | 5 | | 3 | | 1 | 9 |
| CHE | Chemin | 46.08993 | 7.08978 | 2006 | 3 | 1 | | 3 | 2 | 1 | 2 | 12 |
| CRA | Crans-Montana | 46.34650 | 7.53890 | 2006 | | 1 | 1 | | 3 | 1 | 1 | 7 |
| CRE | Cressbrook Dale | 53.26724 | -1.74041 | 2008 | | | | | | 2 | | 2 |
| CTP | Colt Park | 54.19365 | -2.35247 | 2008 | | 4 | | | 1 | | 1 | 6 |
| DON | Donovaly | 48.88922 | 19.23068 | 2008 | | | 3 | 1 | 2 | | | 6 |
| EID | Eidda Pastures | 53.03720 | -3.74190 | 2008 | | | | | | 2 | | 2 |
| ELL | Ellingsrudelva | 59.91771 | 10.91844 | 2007 | | 1 | | | | | | 1 |
| EPO | Esposouille | 42.62341 | 2.09450 | 2008 | 1 | 1 | | | 2 | 2 | 3 | 9 |
| FRO | Froson | 63.18205 | 14.60268 | 2007 | | 2 | | | | | | 2 |

| Code | Name | Lat | Lon | Year | | | | | | | | Total |
|------|------|-----|-----|------|---|---|---|---|---|---|---|-------|
| GAL | Col du Galibier | 45.08528 | 6.43861 | 2006 | | 2 | 2 | | 3 | | 3 | 10 |
| GLE | Glen Fender | 56.78138 | -3.79485 | 2008 | | | | | 2 | 2 | | 4 |
| HT1 | Haute Tinee 1 | 44.29617 | 6.81871 | 2007 | | | 3 | | | | | 3 |
| HT2 | Haute Tinee 2 | 44.28426 | 6.85581 | 2007 | | | | 3 | | 1 | | 4 |
| KRA | Krasno Polje | 44.80869 | 14.97271 | 2008 | | | | | 1 | | | 1 |
| LAK | Laktatjakka | 68.42931 | 18.40674 | 2007 | | 1 | | 4 | 2 | | | 7 |
| LFE | Lough Fern | 55.06569 | -7.71130 | 2008 | | | | | 1 | | | 1 |
| LOS | Loser | 47.66052 | 13.78485 | 2007 | | | 4 | | 3 | 1 | | 8 |
| MOE | Moerlimatt | 47.90597 | 8.07760 | 2007 | | | 1 | | 1 | 1 | 1 | 4 |
| MTP | Monte Pizi | 41.91524 | 14.16714 | 2008 | | | | | | 2 | | 2 |
| NAV | Naverdal | 62.70417 | 10.13002 | 2007 | | 3 | | | | | 1 | 4 |
| PAJ | Pajino Preslo | 43.27799 | 20.81970 | 2008 | | | | | | 2 | | 2 |
| PAN | Puerto de Panderrueda | 43.12743 | -4.97223 | 2008 | | | 1 | 1 | 3 | 2 | 1 | 8 |
| PIL | Pila | 48.90017 | 20.29449 | 2008 | 3 | | | 2 | 1 | | | 6 |
| POD | Podlesok | 48.94962 | 20.35190 | 2008 | | | | | 1 | 1 | | 2 |
| PPN | Petit Papa Noel | 66.51647 | 25.79386 | 2007 | 1 | 3 | | | 2 | 3 | | 9 |
| PYD | Puy de Dome | 45.77222 | 2.96333 | 2006 | | | 2 | 2 | 2 | | | 6 |
| PYM | Puy Mary | 45.11139 | 2.68083 | 2006 | | | 1 | | | | | 1 |
| PYS | Puy de Sancy | 45.53500 | 2.80972 | 2006 | | | 3 | 2 | 1 | | | 6 |
| RAD | Radkow | 50.46866 | 16.35321 | 2008 | 2 | 4 | | | 2 | | | 8 |
| RIS | Risnjak - Snjeznik | 45.43871 | 14.58494 | 2008 | | | | | 1 | 2 | | 3 |
| SAL | Salla | 66.83020 | 28.65427 | 2007 | 1 | | | | 2 | 1 | 4 | 8 |
| SED | Sede de Pan | 43.03949 | -0.48651 | 2008 | | | | 2 | | | | 2 |
| SET | Seterasen | 65.53432 | 13.67744 | 2007 | 1 | 4 | | | 1 | 1 | 1 | 8 |
| SOL | Solberga | 57.95194 | 13.56116 | 2007 | 4 | 2 | | | 2 | 1 | | 9 |
| STE | Steingaden | 47.59529 | 11.01296 | 2007 | | | | | 1 | 1 | 31 | 5 |

| STR | Straumen | 67.38440 | 15.64921 | 2007 | | | | | 2 | | | 2 |
|-----|----------|----------|----------|------|---|---|---|---|---|---|---|---|
| SUS | Susch | 46.74728 | 10.07473 | 2006 | 2 | | 2 | | 2 | 1 | 3 | 10 |
| SVA | Svartla | 65.99583 | 21.22062 | 2007 | 3 | 3 | | | | | | 6 |
| TAR | Tarasp | 46.77730 | 10.25056 | 2006 | | | 1 | | 2 | 1 | 3 | 7 |
| VIT | Vitosha | 42.59032 | 23.29342 | 2008 | | | | | | 2 | | 2 |
| ZAL | Zali Log | 46.20342 | 14.11080 | 2008 | | | | 3 | 2 | 1 | | 6 |
| | | | | total: | 21 | 44 | 42 | 33 | 59 | 34 | 38 | 271 |

812

# Appendix A. Supplementary material

Table S1. Summary statistics for the RAD-sequenced samples: number of RAD fragments clusters and mean coverage after retaining clusters with a coverage >5, estimated heterozygosities, number of consensus loci after paralog filtering, and the numbers of loci retained for each dataset after filtering for coverage among the samples.

Fig. S1. Map of the sampled *Chiastocheta* specimens used in the study.

Fig. S2. The effect of different clustering thresholds (*X*-axis) and minimum loci coverages (indicated by colors: red – 10, blue – 20, green – 100 individuals) on the total number of assembled loci, proportion of missing data, loci overlap among the technical replicates, and mean number of individuals per locus.

Fig. S3. Pattern of RAD-seq loci sharing among the sequenced individuals for datasets: a) the main dataset using clustering similarity of 75% and minimum loci coverage among individuals of 20; b) using clustering similarity of 75% and minimum loci coverage among individuals of 100.

Fig. S4. a) ML phylogeny inferred for the mtDNA dataset; b) ML phylogeny inferred for the RAD-seq dataset; c) SVDquartets phylogeny inferred for the RAD-seq dataset; bootstrap node supports > 80 are shown denoted by gray points, bootstrap node supports > 90 are shown denoted by black points.

Fig. S5. STRUCTURE runs for *K*=2 to 7, plotted against the RAD-seq based phylogeny.

837    Fig. S6. Phylogenetic trees on the loci partitioned into the sets of 2 to 6 clusters,

838    considering only the loci present in at least 100 specimens. Bootstrap node supports >

839    80 are shown denoted by gray points, bootstrap node supports > 90 are shown denoted

840    by black points. Numbers below the trees denote the number of clusters into which the

841    dataset was divided. Plot of log-likelihood improvement versus the number of clusters

842    is presented in the first box.

843

844    Appendix S1. RAD-sequencing protocol.