1    The Perfect Storm:

2    Gene Tree Estimation Error, Incomplete Lineage Sorting, and Ancient Gene

3    Flow Explain the Most Recalcitrant Ancient Angiosperm Clade, Malpighiales

4

5    Liming Cai[1], Zhenxiang Xi[1,2], Emily Moriarty Lemmon[3], Alan R. Lemmon[4], Austin Mast[3],

6    Christopher E. Buddenhagen[3,5], Liang Liu[6], Charles C. Davis[1]

7

8    1 *Department of Organismic and Evolutionary Biology, Harvard University Herbaria,*

9    *Cambridge, MA 02138, USA;*

10   2 *Key Laboratory of Bio-Resource and Eco-Environment of Ministry of Education, College of*

11   *Life Sciences, Sichuan University, Chengdu 610065, China;*

12   3 *Department of Biological Sciences, 319 Stadium Dr., Florida State University, Tallahassee,*

13   *FL 32306, USA;*

14   4 *Department of Scientific Computing, Florida State University, Tallahassee, FL 32306, USA;*

15   5 *AgResearch, 10 Bisley Road, Hamilton 3214, New Zealand*

16   6 *Department of Statistics and Institute of Bioinformatics, University of Georgia, Athens, GA*

17   *30602, USA;*

18

19   Corresponding author:

20   Liming Cai, Department of Organismic and Evolutionary Biology, Harvard University

21   Herbaria, Cambridge, MA 02138, USA; E-mail: lcai@g.harvard.edu

22      Charles C. Davis, Department of Organismic and Evolutionary Biology, Harvard University

23      Herbaria, Cambridge, MA 02138, USA; E-mail: cdavis@oeb.harvard.edu

24

25      **ABSTRACT**

26              The genomic revolution offers renewed hope of resolving rapid radiations in the

27      Tree of Life. The development of the multispecies coalescent (MSC) model and  improved

28      gene tree estimation methods can better accommodate gene tree heterogeneity caused by

29      incomplete lineage sorting (ILS) and gene tree estimation error stemming from the short

30      internal branches. However, the relative influence of these factors in species tree inference

31      is not well understood. Using anchored hybrid enrichment, we generated a data set

32      including 423 single-copy loci from 64 taxa representing 39 families to infer the species

33      tree of the flowering plant order Malpighiales. This order alone includes nine of the top ten

34      most unstable nodes in angiosperms, and the recalcitrant relationships along the backbone

35      of the order have been hypothesized to arise from the rapid radiation during the

36      Cretaceous. Here, we show that coalescent-based methods do not resolve the backbone of

37      Malpighiales and concatenation methods yield inconsistent estimations, providing

38      evidence that gene tree heterogeneity is high in this clade. Despite high levels of ILS and

39      gene tree estimation error, our simulations demonstrate that these two factors alone are

40      insufficient to explain the lack of resolution in this order. To explore this further, we

41      examined triplet frequencies among empirical gene trees and discovered some of them

42      deviated significantly from those attributed to ILS and estimation error, suggesting gene

43      flow as an additional and previously unappreciated phenomenon promoting gene tree

44      variation in Malpighiales. Finally, we applied a novel method to quantify the relative

45    contribution of these three primary sources of gene tree heterogeneity and demonstrated

46    that ILS, gene tree estimation error, and gene flow contributed to 15%, 52%, and 32% of

47    the variation, respectively. Together, our results suggest that a perfect storm of factors

48    likely influence this lack of resolution, and further indicate that recalcitrant phylogenetic

49    relationships like the backbone of Malpighiales may be better represented as phylogenetic

50    networks. Thus, reducing such groups solely to existing models that adhere strictly to

51    bifurcating trees greatly oversimplifies reality, and obscures our ability to more clearly

52    discern the process of evolution.

53

54    Keywords: rapid radiation, triplet frequency, concatenation, coalescent, phylogenomics,

55    hybrid enrichment, flanking region

56

57    **INTRODUCTION**

58        One of the most difficult challenges in systematics is reconstructing evolutionary

59    history during periods of rapid radiation. During such intervals, few DNA substitutions

60    accrue, rendering little information for phylogenetic inference. The potentially large

61    population sizes and close evolutionary relationships create opportunities for widespread

62    incomplete lineage sorting (ILS) and gene flow, leading to excessive gene tree-species tree

63    conflict. The tremendous growth of genome-scale data sets, however, has greatly improved

64    researchers' ability to investigate rapid radiations by providing hundreds to thousands of

65    unlinked loci. Commonly applied approaches include not only whole genome sequencing

66    but also RNA-Seq, RAD-Seq, and anchored hybrid enrichment, which in general are cost

67    effective and efficient across broad taxonomic groups and yield data sets with dense locus

68    and taxon sampling (Lemmon and Lemmon 2013). These approaches are promising and

69    have been variously applied to successfully resolve a number of recalcitrant clades across

70    the Tree of Life, including in birds (Prum et al. 2015), mammals (Song et al. 2012), fish

71    (Wagner et al. 2013), and plants (Wickett et al. 2014).

72         Despite their promise, however, these enormous data sets also introduce new

73    methodological challenges and complexities. In particular, phylogenomic data sets may

74    yield strongly supported, yet conflicting or artifactual, results depending on the method of

75    inference or genomic regions sampled (Song et al. 2012; Jarvis et al. 2014; Xi et al. 2014;

76    Reddy et al. 2017; Shen et al. 2017). During rapid radiations, ILS can lead to extreme

77    conditions where the most probable gene tree differs from the topology of the true species

78    tree, which is referred to as the "anomaly zone" (Degnan and Rosenberg 2006; Rosenberg

79    and Tao 2008). Such pervasive genealogical discordance, in particular, can result in biased

80    species tree inference when applying concatenation methods, and produce inconsistent

81    and conflicting results with strong confidence (Song et al. 2012; Xi et al. 2014). The

82    multispecies coalescent (MSC) model, which explicitly accommodates gene tree

83    heterogeneity caused by ILS, in contrast, has been demonstrated to be more reliable under

84    these circumstances. Most recently, a class of "two-step" summary coalescent methods has

85    been the focus of substantial development and application (Nakhleh 2013). They are

86    demonstrated to be statistically consistent under the MSC model and can work efficiently

87    with genome-scale data (Liu et al. 2009; Liu et al. 2010; Chifman and Kubatko 2014;

88    Mirarab et al. 2014c). Their application has been successful in resolving mammalian, avian,

89    and seed plant relationships in cases where concatenation methods have been

90    demonstrated to be inconsistent (Song et al. 2012; Xi et al. 2013; Reddy et al. 2017).

91      In addition to ILS, gene tree estimation error has also been a major focus of work to

92      improve the accuracy of phylogenomic inference. This is especially relevant for summary

93      coalescent methods, which assume the input gene trees to be essentially error-free, e.g.,

94      Lanier et al. 2014; Mirarab et al. 2014c; Roch and Warnow 2015; Xu and Yang 2016; Blom

95      et al. 2017. Rapid radiations are particularly challenging in this regard. Here, short internal

96      branches may yield error-prone gene tree estimation when phylogenetically informative

97      characters are minimal (Xi et al. 2015). This may be further complicated if such radiations

98      are ancient and followed by long descendent branches, which may exacerbate long-branch

99      attraction artifacts (Whitfield and Kjer 2008). Though benchmark studies have

100     demonstrated the consistency of summary coalescent methods when substantial amounts

101     of such non-phylogenetic signal are included (Philippe et al. 2011; Roch and Warnow 2015;

102     Xi et al. 2015; Hahn and Nakhleh 2016), accurate gene tree inference remains of crucial

103     importance for reliable species tree estimation (Shen et al. 2017). A number of methods

104     have been developed to mitigate gene tree estimation error, including improving taxon

105     sampling, applying appropriate models of nucleotide evolution, reducing missing data,

106     subsampling informative genes, and locus binning (Zwickl and Hillis 2002; Lemmon et al.

107     2009; Salichos and Rokas 2013; Cox et al. 2014; Mirarab et al. 2014a; Hosner et al. 2015).

108     Beyond ILS and gene tree estimation error, gene flow between non-sister species

109     can similarly result in gene tree–species tree conflict and lead to incorrect species tree

110     estimation. Unlike the MSC model, gene flow from a non-sister species leads to an

111     overrepresentation of the parental allele in the descendants and therefore the frequencies

112     of the two minor topologies are asymmetrical (Durand et al. 2011). A number of species

113     network inference methods have been developed to detect and infer gene flow based on

114    such expectation. They either use counts of the shared derived alleles, such as the classic D-

115    statistic test (Green et al. 2010; Durand et al. 2011), or the gene tree topology as input (e.g.,

116    Huson et al. 2005; Meng and Kubatko 2009; Yu et al. 2011; Solís-Lemus et al. 2017). The

117    latter methods are often based on *a priori* evolutionary models and have been increasingly

118    applied to empirical data sets.

119            During periods of rapid radiation, all of the above phenomena—ILS, introgression,

120    and gene tree estimation error—may occur simultaneously to obscure phylogenetic signal

121    (Pease et al. 2016), culminating in a perfect storm confounding phylogenomic inference.

122    When a limited number of alternative species tree topologies are involved, these

123    phenomena can be distinguished from each other using methods discussed above (Zwickl

124    et al. 2014; Arcila et al. 2017; Meyer et al. 2017; Beckman et al. 2018; Glémin et al. 2019).

125    However, when the rapid radiation generates a cloud of alternative tree topologies, all of

126    which are weakly supported, such model-based methods become less practical because

127    priors necessary to test hypothesis of introgression are difficult to determine accurately.

128    Additional challenges arise from the excessive computational resources required to apply

129    such network inference methods to data sets involving hundreds of species. Moreover,

130    following the identification of ILS, introgression, and gene tree estimation error, a more

131    quantitative assessment characterizing their relative contribution to overall gene tree

132    variation has not been addressed in any empirical system to our knowledge.

133            Using anchored hybrid enrichment (Lemmon et al. 2012), we generated a large

134    phylogenomic data set including 423 single-copy nuclear loci with 64 taxa to infer

135    relationships of the flowering plant clade Malpighiales. The order Malpighiales comprise ca

136    7.8% of eudicot diversity (Magallon et al. 1999) and include more than 16,000 species in

137    ~36 families (Stevens and Davis 2001). Species in Malpighiales encompass astonishing

138    morphological and ecological diversity ranging from epiphytes (Clusiaceae), submerged

139    aquatics (Podostemaceae), to emergent rainforest canopy species (Callophyllaceae). The

140    order also includes numerous economically important crops with sequenced genomes, e.g.,

141    rubber (*Hevea*), cassava (*Manihot*), flax (*Linum*), and aspen (*Populus*). Despite their

142    ecological and economic importance, the evolutionary history of Malpighiales remains

143    poorly understood. While analyzing chloroplast genome sequences has greatly improved

144    the resolution of this clade, relationships among its major subclades remain uncertain (Xi

145    et al. 2012), and analyses using nuclear genes lack resolution along the spine of the clade

146    (Davis et al. 2005; Wurdack and Davis 2009). According to Smith et al. (2013), this region

147    of the Malpighiales phylogeny has been implicated in nine of the top ten most unstable

148    nodes across all angiosperms, including Pandaceae, Euphorbiaceae, Linaceae, the most

149    recent common ancestor (MRCA) of Salicaceae and Lacistemataceae, the MRCA of

150    Malpighiaceae and Elatinaceae, as well as the MRCA of putranjivoids, phyllanthoids,

151    chrysobalanoids, and rhizophoroids *sensu* Xi et al. (2012). In short, Malpighiales have been

152    coined one of the "thorniest nodes" in the angiosperm tree of life (Soltis et al. 2005). A long-

153    standing hypothesis for this lack of resolution has been attributed to the clade's rapid

154    radiation during the Albian and Cenomanian (112–94 million years ago [Ma]; Davis et al.

155    2005; Wurdack and Davis 2009; Xi et al. 2012). This radiation has produced a phylogeny

156    characterized by extremely short internal branches along the backbone of the phylogeny,

157    followed by long branches subtending most crown group families. This is particularly

158    problematic because, as we summarize above, short internal branches represent species

159    tree anomaly zones where ILS may be pervasive and gene tree estimation error is high (Liu

160    et al. 2015; Roch and Warnow 2015; Edwards et al. 2016). Incongruent phylogenetic

161    signals between organelle and nuclear genes also support introgression associated with the

162    origin of this order (Sun et al. 2015).

163         The development of next-generation sequencing, the MSC model that accommodate

164    ILS, and best practices to reduce gene tree estimation error offers a unique opportunity to

165    re-examine Malpighiales in the context of resolving rapid radiations. Here, we apply both

166    concatenation and coalescent-based methods for phylogenomic analyses and evaluate the

167    relationships and consistency of nodal resolution under a variety of conditions. We also

168    apply simulations to explore the impact of ILS and gene tree estimation error based on the

169    empirical parameters of our inferred species tree. We further apply a triplet analysis to

170    detect gene flow and identify hotspots of reticulate evolution in the species tree. And

171    finally, we develop a novel method to quantitatively assess the contribution of three

172    primary sources of gene tree variation in Malpighiales—ILS, gene tree estimation error,

173    and gene flow.

174

175    **MATERIALS AND METHODS**

176    *Taxon Sampling*

177         We sampled a total of 56 species in the order Malpighiales, representing 39 families

178    and all major clades *sensu* Wurdack and Davis (2009) and Xi et al. (2012) (Table S1).

179    Species were sampled to represent the breadth of Malpighiales diversity. Four species from

180    the order Celastrales and two species from the order Oxalidales were sampled as closely

181    related outgroups (Chase et al. 2016). Two species from the order Vitales were also

182    included as more distantly related outgroups (Chase et al. 2016, Table S1).

183

*Library Preparation, Enrichment, and Locus Assembly*

185       Data were collected at the Center for Anchored Phylogenomics at Florida State

186  University (http://www.anchoredphylogeny.com) using the anchored hybrid enrichment

187  method (Lemmon et al. 2012; Buddenhagen et al. 2016). This method targets universally

188  conserved single-copy regions of the genome that typically span 250 to 800 base pairs (bp),

189  thus mitigating the confounding effect of paralogy in gene tree estimates. Briefly, total

190  genomic DNA was sonicated to a fragment size of 300–800 bp using a Covaris E220

191  Focused-ultrasonicator. Library preparation and indexing was performed following the

192  protocol in Hamilton et al. (2016). A size-selection step was also applied after blunt-end

193  repair using SPRI select beads (Beckman-Coulter Inc). Indexed samples were then pooled

194  and enriched using the Angiosperm v1 kit (Agilent Technologies Custom SureSelect XT kit

195  ELID 623181; Buddenhagen et al. 2016). The resulting libraries were sequenced on an

196  Illumina HiSeq 2500 System using the PE150 protocol.

197       Quality-filtered sequencing reads were processed following the methods described

198  in Hamilton et al. (2016) to generate locus assemblies. Briefly, paired reads were merged

199  prior to assembly following Rokyta et al. (2012). Reads were then mapped to the probe

200  region sequences of the following reference genomes: *Arabidopsis thaliana* (Malvales,

201  Arabidopsis Genome Initiative 2000), *Populus trichocarpa* (Malpighiales, Tuskan et al.

202  2006), and *Billbergia nutans* (Poales, Buddenhagen et al. 2016). Finally, the assemblies

203  were extended into the flanking regions. Consensus sequences were generated from

204  assembly clusters with the most common base being called when polymorphisms could be

205  explained as sequencing error.

206

207 *Orthology Assignment*

208      Orthologous sequences were determined following Prum et al. (2015) and Hamilton

209 et al. (2016). The assembled sequences were grouped by locus and a pairwise distance was

210 calculated as the percent of shared 20-mers. Sequences were subsequently clustered based

211 on this distance matrix using the neighbor-joining algorithm (Saitou and Nei 1987). When

212 more than one cluster was detected for a target region, each cluster was treated as a

213 different locus in subsequent analyses. Clusters including less than 50% of the species were

214 discarded.

215

216 *Sequence Alignment, Masking, and Site-subsampling*

217      Each locus was first aligned using MAFFT v7.023b (Katoh and Standley 2013) with

218 "--genafpair --maxiterate 1000" flags imposed. Alignments were end trimmed and

219 internally masked to remove misassembled or misaligned regions (Buddenhagen et al.

220 2016). Firstly, conserved sites were identified in each alignment where >40% of the

221 nucleotides at that site were identical across species. For end trimming, sequences for each

222 gene accession were scanned from both ends towards the center until more than fourteen

223 nucleotides in a sliding window of 20 bp matched the conserved sites. Once the start and

224 end of each sequence was established, the internal masking then required that >50% of the

225 nucleotides in a sliding window of 30 bp matched the conserved sites. Regions that did not

226 meet this criterion were masked. Finally, we removed any gene sequence in the alignment

227 with >50% ambiguous nucleotide composition. We also required all locus alignments to

228 contain *Leea guineense* (Vitales) for rooting purposes.

229     To further explore the phylogenetic utility of the flanking regions of hybrid

230     enrichment data, we applied three increasingly stringent site-subsampling strategies using

231     trimAl v1.2 (Capella-Gutiérrez et al. 2009) following our masking steps described above. To

232     construct our "low-stringency data set", we set the gap threshold to be 0.8 (-gt 0.8) in

233     trimAl to remove sites containing >20% indels or missing data for each alignment.  This

234     data set includes the highest percentage of flanking regions and resulted in the longest

235     alignments. We then applied a site composition heterogeneity filter to this "low-stringency

236     data set" to create our "medium-" and "high-stringency data set" by setting the minimum

237     site similarity score to be 0.0002 and 0.001 (e.g., -st 0.001), respectively. This has the effect

238     of removing especially rapidly evolving sites within flanking regions for which we expect

239     higher composition heterogeneity. The resulting "medium-" and "high-stringency data set"

240     thus include lower percentage of flanking regions.

241

242     *Gene Tree Estimation*

243     To infer individual gene trees for coalescent-based analyses, we applied maximum

244     likelihood (ML) as well as Bayesian Inference (BI). To estimate ML trees, we used RAxML

245     v8.1.5 (Stamatakis 2014) under the GTR+$\Gamma$ model with 20 random starting points. We

246     chose the GTR+$\Gamma$ model because it accommodates rate heterogeneity among sites, while the

247     other available GTR model in RAxML, the GTRCAT model, is less appropriate due to our

248     small taxon sampling size (Stamatakis 2014). Statistical confidence of each gene tree was

249     assessed by performing 100 bootstrap (BP) replicates. We additionally inferred the

250     Bayesian posterior distribution of gene trees using MrBayes v3.2.1 (Ronquist and

251     Huelsenbeck 2003). We only applied BI to the low-stringency data set due to computational

252    cost and this data set yielded the best resolved gene trees (see Results below). We applied

253    the GTR+Γ model with two independent runs for each gene. Each run included four chains,

254    with the heated chain at temperature 0.20 and swapping attempts every 10 generations.

255    Initially, four million generations were used with 25% burn-in period, sampled every 1,000

256    generations. Runs that failed to reach the targeted standard deviation of split frequencies

257    ≤0.02 were rerun with the same settings but with 10 million generations, sampled every

258    5,000 generations until attaining a standard deviation of split frequencies ≤0.02. We

259    randomly sampled 100 trees in the posterior distribution of inferred gene trees to conduct

260    bootstrap replication in the coalescent analyses (Table S2). Trees sampled from the

261    posterior distribution are more similar to the optimum Bayesian tree than those sampled

262    from the non-parametric bootstrapping. Therefore, we also expect higher support values in

263    the species tree.

264

265    *Species Tree Inference Using Concatenation and Coalescent-based Methods*

266        Our trimmed gene matrices were concatenated and analyzed using both RAxML and

267    ExaML v3.0.18 (Kozlov et al. 2015). In our RAxML analyses, the species trees were inferred

268    under the GTR+Γ model with 100 rapid bootstrapping followed by a thorough search for

269    the ML tree. In ExaML analyses, species trees were inferred under the GTR+Γ model with

270    20 random starting points. We then conducted 100 bootstrap replicates to evaluate nodal

271    support. Partitions for both analyses were selected by PartitionFinder v2.1.1 based on AICc

272    (Akaike Information Criterion) criteria using the heuristic search algorithm "rcluster"

273    (Lanfear et al. 2012). We also conducted BI for species tree estimation as implemented in

274    PhyloBayes (Lartillot et al. 2013). For BI analysis, we applied the CAT-GTR model, which

275    accounts for across-site compositional heterogeneity using an infinite mixture model

276    (Lartillot and Philippe 2004). Two independent Markov chain Monte Carlo (MCMC)

277    analyses were conducted for each concatenated nucleotide matrix. Convergence and

278    stationarity from both MCMC analyses were determined using bpcomp and tracecomp

279    from PhyloBayes. We ran each MCMC analysis until the largest discrepancy observed

280    across all bipartitions was smaller than 0.1 and the minimum effective sampling size

281    exceeded 200 for all parameters in each chain.

282         To infer our species tree using coalescent-based models, we obtained ML gene trees

283    and BI consensus trees for each locus. MP-EST (Liu et al. 2010) and ASTRAL-II (Mirarab

284    and Warnow 2015) were subsequently used to perform species tree inference using

285    optimally estimated gene trees. Statistical confidence at each node was evaluated by

286    performing the same species tree inference analysis on 100 ML bootstrap gene trees or

287    trees sampled from our Bayesian posterior distributions. The resulting 100 species trees

288    estimated from bootstrapped samples were summarized onto the species tree inferred

289    from ML gene trees using the option "-f z" in RAxML.

290

291    *Simulation of gene alignments with realistic parameters of ILS and gene tree estimation error*

292         To investigate the impact of ILS and gene tree estimation error on the accuracy of

293    species tree inference we simulated sequences assuming a known species tree. Here, the

294    tree topology estimated by MP-EST with the low-stringency data set (analysis No. 15 in

295    Table S2) was invoked as the known species tree. We chose this best-supported MP-EST

296    topology because the branch lengths are estimated in coalescent unit, which is an essential

297    parameter for ILS simulation. We thus applied this species tree to all of the downstream
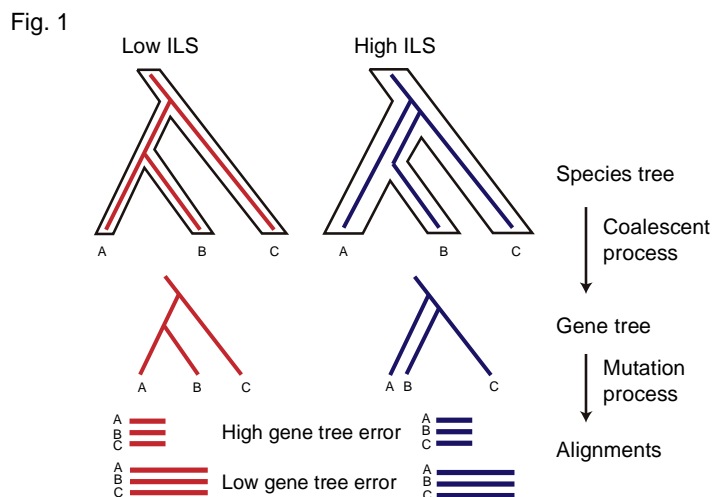
298    simulation-based analyses, including the triplet test for MSC model fitness and relative

299    importance analysis.

300         To simulate conditions of high and low levels of ILS, we modified the key population

301    mutation parameter "theta" when generating gene trees under the coalescent model using

302    the function "sim.coaltree.sp.mu" in the R package Phybase (Liu and Yu 2010). Theta was

303    set to be 0.01 and 0.1 to reflect low and high ILS, respectively. The range of theta was

304    determined based on our empirical data sets by following two steps. First, we inferred the

305    branch lengths of the species tree in mutation units in RAxML using the fixed topology of

306    the MP-EST species tree and the concatenated low-stringency data set. Second, theta for

307    each branch was calculated by dividing the branch lengths estimated from RAxML

308    (mutation units) by that estimated from MP-EST (coalescent units). The other input for

309    Phybase, the ultrametric species tree, was generated from this RAxML phylogeny using the

310    function "chronos" in the R package ape (Paradis et al. 2004). In addition, we set the

311    relative mutation rates to follow a Dirichlet distribution with alpha equal to 5.0. This alpha

312    reflected the large variance in gene mutation rates. Finally, 1,500 non-ultrametric gene

313    trees were simulated separately for each theta.

314         From these simulated gene trees, DNA alignments of different lengths were

315    subsequently generated to reflect various levels of gene tree estimation error since

316    alignment length is easy to manipulate and shorter alignments correspond to higher error

317    rates (Mirarab et al. 2014b). We used bppsuite (Guéguen et al. 2013) to simulate

318    alignments under the GTR+Γ model. Parameters of the model, including the substitution

319    matrix, base frequency, and the gamma rate distribution were extracted from the RAxML

320    phylogeny above inferred from the low-stringency data set. For each gene tree we

321    generated alignments of 300, 400, 500, 1,000, and 1,500 bp.

322         As a result, fifty data sets were generated by including 100, 200, 500, 1,000, and

323    1,500 simulated loci of five length categories and two theta categories (Table S3, Fig. 1).

324    Species trees were inferred using the concatenation and coalescent methods as described

325    above under these varying levels of ILS and gene tree estimation error. Finally, we

326    quantified gene tree–species tree discordance and species tree error by measuring the RF

327    distance between an estimated gene tree or species tree to the true species tree.



328

329    **Figure 1** Simulation of ILS and gene tree estimation error. ILS was simulated though the

330    coalescent process by setting low (0.01) and high (0.1) theta values. DNA alignments were

331    subsequently generated through the mutation process based on simulated gene trees. Five

332    alignments were generated for each gene tree with lengths of 300, 400, 500, 1000, and

333    1500 bp (only two are shown in the graph). Shorter alignment lengths increase in gene tree

334    estimation error.

335

336   In order to assess the sensitivity of our simulation results to the choice of input

337 species tree and theta values, we additionally examined gene tree–species tree discordance

338 among bootstrapped samples. We simulated 1,500 gene trees for each of the 100 MP-EST

339 bootstrapped species trees. Gene trees were simulated directly from each species tree

340 using the "sim.coal.mpest" function in the R package Phybase (Liu and Yu 2010). This

341 method does not require *a priori* theta parameters as was imposed in our simulation above

342 and so alleviates concerns of applying erroneous theta values. We subsequently quantified

343 the gene tree–species tree discordance for each bootstrap replicate as described above. We

344 did not use these gene trees to simulate alignments because these gene trees are

345 ultrametric (Liu and Yu 2010) and thus not suitable for such purpose.

346

347 *A Test of the MSC Model Using Triplet Frequencies*

348   To determine the fit of the MSC model to our empirical data we additionally

349 examined the triplet frequency for all 423 ML genes trees inferred from our low-stringency

350 data set using a custom R script available on Github

351 (http://github.com/lmcai/Coalescent_simulation_and_gene_flow_detection). We used the

352 asymmetrical triplet frequency as evidence for introgression (Fig. 2a). This metric has been

353 widely applied in parsimony, likelihood, and Bayesian based species network inference

354 methods to detect sources of gene flow (Nakhleh 2013) . Our method differs from these

355 methods in two aspects: first, the statistical significance of asymmetry in triplet frequency

356 is determined by a null distribution simulated from the empirical data. We took into

357 account variations from ILS and missing data, thus reducing the false positive rate. Second,

358 unlike other model-based species network inference methods, after identifying

359    significantly asymmetrical triplets, we used a novel method to summarize and visualize the

360    distribution of lineages involved in gene flow on a species tree without optimizing the

361    global network (Fig. 2b). As a result, our methods can be easily scaled to genomic data

362    involving hundreds of taxa.



363

364    **Figure 2** Identification of reticulate evolution using triplet frequency. (a) Theoretical

365    expectations of triplet frequency distribution under the multi-species coalescent (MSC)

366    model with and without introgression. In case of incomplete lineage sorting (ILS),

367    symmetrical distributions of the frequency of two minor topologies are excepted owing to

368    deep coalescence (left). In case of introgression, one of the minor topologies will occur with

369    higher frequency due to gene flow (right). (b) Mapped asymmetrical triplets to species tree

370    to identify reticulate nodes.

371

372          In order to identify a triplet with significantly asymmetrical frequencies, we

373    generated a null distribution of triplet frequencies for each triplet using simulated gene

374    trees under the MSC model. For each of the 100 MP-EST BP species trees, we simulated 423

375    gene trees using the "sim.coal.mpest" function in Phybase. For each set of simulated gene

376     trees, we then generated missing data for each species by pruning that species randomly

377     among all gene trees so that the number of sampled genes of that species was the same as

378     the empirical data. We subsequently counted triplet frequency for these gene trees in each

379     bootstrap replicate. This simulated distribution reflects the variation of triplet frequency

380     arising from ILS, estimation error, sampling error, and missing data. A triplet in the

381     empirical data was identified to be significantly asymmetrical if the difference between the

382     two less frequent triplets exceeded the maximum difference under simulated conditions.

383     Such triplets potentially violate the assumptions of the MSC model, and point towards gene

384     flow especially as an additional factor influencing gene tree heterogeneity, though ancestral

385     population structure (Slatkin and Pollack 2008) and biases in substitution or gene loss can

386     produce asymmetrical triplet as well (see Discussion below).

387

388     *Identifying hotspots of reticulate evolution using the Reticulation Index*

389         We developed a relative measurement statistic, the 'Reticulation Index', to quantify

390     the intensity of introgression at each node. First, for each asymmetrical triplet, we mapped

391     the two inferred introgression branches to the species tree (Fig. 2b). Second, for each node

392     on the species tree, we counted the number of introgression branches that were mapped to

393     it. These raw counts were then normalized by the total number of triplets associated with

394     that node. The resulting percentage is the Reticulation Index for each node. The R script for

395     calculating the Reticulation Index and visualizing the result on a species tree is available in

396     the above Github repository.

397

398    *A Novel Method to Quantify Gene Tree Variation Due to ILS, Gene Tree Estimation Error, and*

399    *Gene Flow*

400         Untangling the effects of ILS, gene tree estimation error, and gene flow is

401    challenging since they all lead to gene tree–species tree discordance. Here, based on a

402    multiple regression model (Grömping 2006), we assign shares of relative importance to ILS,

403    gene tree estimation error, and gene flow in generating gene tree variation by variance

404    decomposition.

405         For all 63 internal nodes in our species tree, we separately estimated the level of ILS,

406    gene tree estimation error, and gene flow for each node. ILS is represented by our

407    estimates of theta. Gene flow is represented by the Reticulation Index for each node. To

408    infer the level of gene tree estimation error at each node, we additionally simulated 423

409    gene alignments of 446 bp (median alignment length in low-stringency data set) from the

410    MP-EST species tree, but each with unique substitution model parameters estimated from

411    the empirical alignments. This simulation and phylogeny inference followed the same

412    strategy of alignment simulation described above (Fig. 1). We subsequently inferred

413    phylogenies for these alignments and summarized them on the species tree to obtain the

414    BP value at each node. Here, the BP values represent the gene tree variation generated by

415    estimation error.

416         The gene tree variation in the empirical data is obtained by summarizing bootstrap

417    trees from each of the 423 loci in our low-stringency data set onto the species tree. The

418    resulting BP values represented observed gene tree variation at each node. We then

419    inferred the relative contribution of ILS, estimation error, and gene flow in explaining gene

420    tree variation using linear regression methods implemented in the R package relaimpo

421    (Grömping 2006). We used four different methods, "lmg", "last", "first", and "Pratt", to

422    decompose the relative importance of the three regressors (Lindeman 1980; Pratt 1987).

423    All of these methods are capable of dealing with correlated regressors and "lmg" is the

424    most robust method among them (Grömping 2006). We applied the functions "boot.relimp"

425    and "booteval.relimp" to estimate the relative importance and their confidence interval by

426    bootstrapping 100 times.

427

428    *Testing the utility of the triplet-frequency-based method using a genomic data set from yeast*

429            To further validate our introgression detection method using the triplet frequency

430    distribution, we applied it to the benchmark multi-locus yeast data set from Salichos and

431    Rokas (2013). We obtained the 1,070 gene trees and inferred a species tree using MP-EST.

432    We also conducted 100 bootstrap replicates of species tree inference using the bootstrap

433    gene trees. We then applied our triplet method to identify asymmetrical triplets as

434    described above (*A Test of the MSC Model Using Triplet Frequencies*). Finally, all

435    asymmetrical triplets were mapped to the inferred species tree and the Reticulation Index

436    for each node was calculated and visualized as described above (*Identifying hotspots of*

437    *reticulate evolution using the Reticulation Index*).

438

439    **RESULTS**

440    *Hybrid Enrichment*

441            We successfully captured and sequenced 423 of our 491 targeted loci. The resulting

442    data matrix was densely sampled and included only 12% missing data. One hundred and

443    one loci included at least 61 taxa (>95% occupancy) and only four loci had more than 19

444    missing species (>30%). The locus sampling per taxon varied from 423 (*Leea guineense*) to

445    278 (*Ouratea sp.* and *Lophopyxis maingayi*, Table S4). After applying site subsampling, the

446    alignment lengths ranged from 190 to 885 bp (median 446 bp) for the low-stringency data

447    set, 157 to 791 bp (median 376 bp) for the medium-stringency data set, and 112 to 751 bp

448    (median 271 bp) for the high-stringency data set (Table S4). In all data sets, the number of

449    parsimony informative sites and the average nodal support was significantly positively

450    correlated with alignment length (*p*-value <1e-5, Fig. S1).

451

452    *Flanking Regions Increase Gene Tree and Species Tree Resolution*

453            We observed increasing mean BP support among gene trees as increasingly larger

454    percentages of the flanking region were included. The average gene tree nodal support

455    from our low-stringency data set (42 ML BP) was significantly higher than nodal support

456    estimates for the medium (39 ML BP, *p*-value = 1.2e-89 in paired t-test) and high-

457    stringency data sets (35 ML BP, *p*-value = 1.3e-77, Fig. S2a).

458    These increases in gene tree resolution also contributed to increased species tree

459    resolution as well as species tree inference congruency. For both concatenation and

460    coalescent analyses, species trees estimated from the low stringent data set with highest

461    amount of flanking regions, always resulted in the highest average BP support (Fig. S2b,c,

462    Table S2) and the lowest pairwise RF distances (Fig. S2d) indicating increased statistical

463    consistency when adding flanking regions.

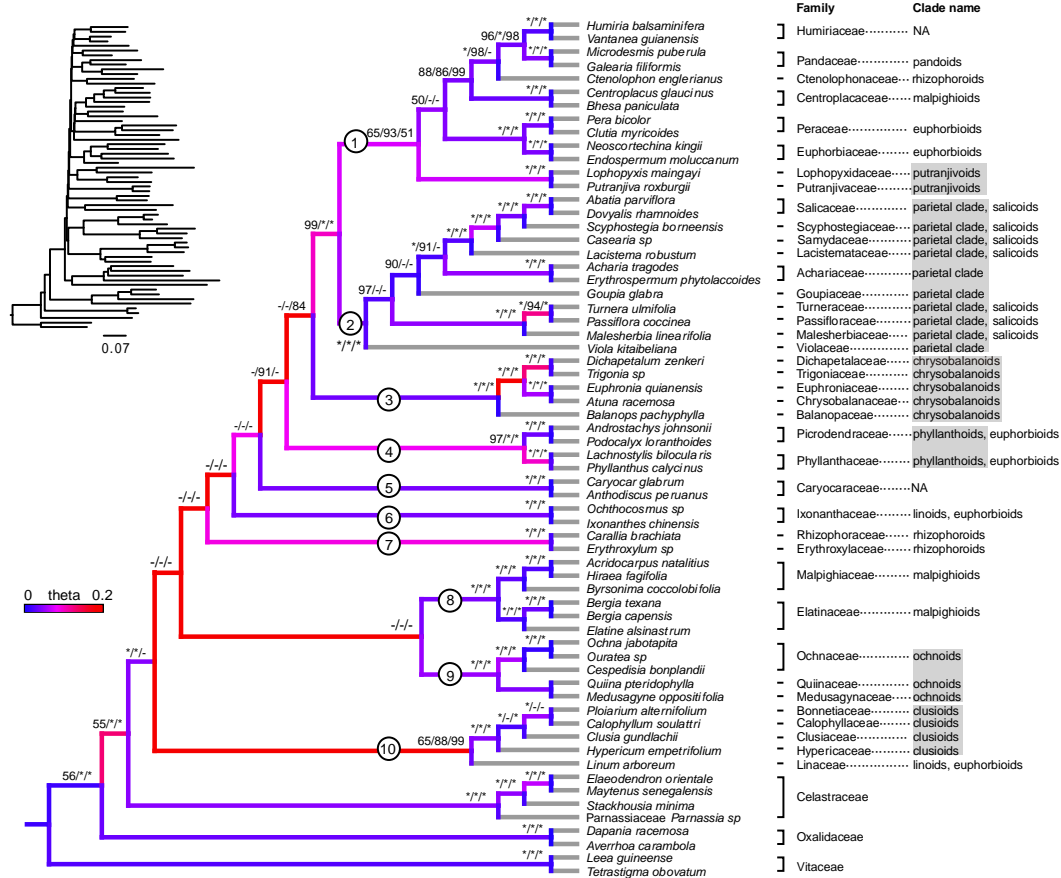464

465    *Malpighiales Species Tree Resolution*

466      We observed significantly higher average species tree nodal support in

467    concatenation compared to coalescent reconstructions (Table S2, *p*-value = 2.62e-28 in

468    paired *t*-test). However, our results also suggest statistical inconsistency across data sets

469    when applying concatenation (Fig. S3). The higher pairwise weighted Robinson–Foulds

470    distance (WRF) in concatenation indicate more well-supported conflicts among species

471    trees, which further supports mounting evidence that coalescent methods are more

472    consistent when reconstructing species tree relationships involving extensive ILS (i.e., the

473    anomaly zone, Degnan and Rosenberg, 2006, Rosenberg and Tao, 2008). In addition, we did

474    not find locus subsampling based on locus length, number of PI sites, or gene tree quality

475    help increase species tree resolution (see Supplementary Note 1, Table S2).

476      Our most well resolved species trees estimation inferred with ASTRAL and MP-EST

477    uncovered ten major subclades of Malpighiales (Clade 1 to 10 in Fig. 3). These relationships

478    corresponded to families or closely related clades of families, five of which have previously

479    been identified using plastid genome (Fig. 3, Xi et al. 2012). Five new clades were

480    supported with ≥50 BP, >0.90 PP. Three of these newly identified clades are in conflict (>70

481    BP) with the plastid phylogeny from Xi et al. (2012) and are discussed more extensively

482    below. Interrelationships among these ten major subclades, however, were not well

483    resolved (<50 BP).

484

485    *Simulated Levels of ILS and Gene Tree Estimation Error Reflects Empirical Data*

486      The 5% and 95% quantiles of theta were inferred to be 0.0254 and 0.176,

487    respectively, with a median of 0.0688. High theta was mostly found along the backbone of

488    the species tree, indicating the likelihood of extensive ILS within this region of the tree (Fig.

489



**Figure 3** Species phylogeny of Malpighiales derived from MP-EST with complete low-stringency data set (analysis No. 15 in Table S2). Gene trees are estimated using MrBayes. Branches are colored by the inferred population mutation parameter theta. Warmer colors indicate higher theta and thus higher level of ILS. Terminal branches are colored grey due to lack of data to infer theta. BP values from best-resolved MP-EST/ASTRAL/RAxML analyses (analysis No. 15, 17, and 11 in Table S2) are indicated above each branch; an asterisk indicates 100 BP support; a hyphen indicates less than 50 BP. Branch lengths estimated from RAxML by fixing the species tree topology are presented at the upper left corner. The eleven major clades highlighted in the discussion are identified with circled numbers along each relevant branch. The clade affiliation for each family based on the

500    plastid phylogeny (Xi et al. 2012) is indicated on the right. Clades identified by Xi et al.

501    (2012) that are also monophyletic in this study are highlighted using gray shades.

502

503    3). This is likely an overestimation of theta since all topological variations are attributed to

504    coalescent process including the ones originate from mutational variance (Huang and

505    Knowles 2009). We therefore set the theta parameter to be 0.01 and 0.1 in our coalescent

506    gene tree simulation, which reflected the left and right tails of low and high ILS estimated

507    from empirical data.

508         In our simulation, the average gene tree estimation error was 0.319 for alignments

509    of 300bp, 0.261 for 400bp, 0.221 for 500bp, 0.133 for 1000bp and 0.098 for 1500bp under

510    low ILS and 0.340 for alignments of 300bp, 0.286 for 400bp, 0.241 for 500bp, 0.161 for

511    1000bp and 0.120 for 1500bp under high ILS. Here, an RF distance of 0 signifies error-free

512    reconstruction versus 1 indicating that none of the true nodes are recovered. Gene tree

513    estimation error was therefore lower in low ILS ($p$-value=6.08e-16 in Student's $t$-test), but

514    was still significantly higher than that estimated from empirical data ($p$-value= 4.24e-65 in

515    Student's $t$-test; see Supplementary Note 2; Fig. S4).

516

517    *Simulation Yields Consistent and Accurate Species Tree Estimation*

518         In our empirical analyses, the low-stringency data set yielded the lowest average

519    gene tree–species tree conflict of 0.563 among the other data sets. In our simulations, the

520    highest average gene tree–species tree conflict observed was 0.507, by setting theta = 0.1

521    and alignment length = 300 bp. Therefore the lowest empirical gene tree–species tree

522    discordance was still significantly higher than the simulated conditions with extremely

523

**Figure 4** Extensive gene tree discordance in empirical versus simulated data. (a) Gene

tree–species tree (G-S) discordance in the empirical (Emp) and simulated (Sim) data

assuming fixed theta in simulation. Discordance is measured by RF distance between

inferred gene trees and the species tree. Under various simulated conditions of ILS (e.g.,

'Low ILS', theta = 0.01 and 'High ILS', theta = 0.1) and gene tree estimation error ('High ILS

+ High err.', theta = 0.1, alignment length=300bp), the simulated gene tree–species tree

discordance is significantly lower than that from empirical data. (b) Gene tree–species tree

discordance is higher in empirical versus simulated conditions without setting theta a

priori. For each BP data set, gene tree–species tree discordance is measured and compared

in both empirical and simulated data sets. Positive values indicate higher gene tree–species

tree discordance in our empirical data. (c) Species tree estimation discordance in empirical

data (left) and simulated data (right).

high level of ILS and gene tree estimation error ($p$-value = 2.2e-16, Student's $t$-test, Fig. 4a).

The same conclusion also applies when simulating gene trees directly from species tree

without setting theta *a priori* (Fig. 4b).
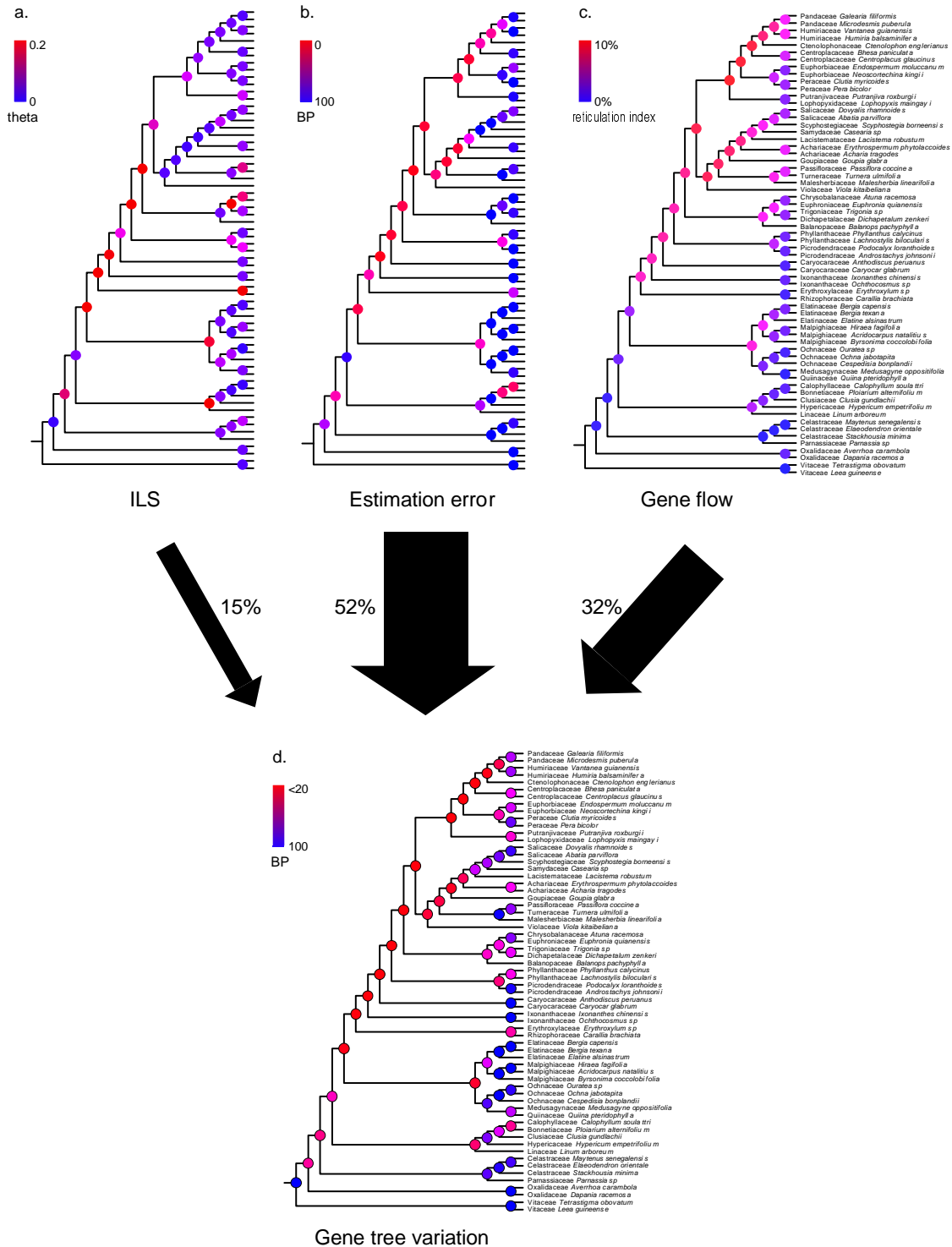
Moreover, even under such simulated conditions of extremely high ILS and gene

tree estimation error, both concatenation and coalescent-based methods yielded consistent

542    and accurate species tree estimation with no more than 12 nodes (< 0.10 RF distance, Fig.

543    4c, Fig. S5) failing to be recovered. The performance of coalescent-based methods is mainly

544    affected by gene tree estimation error (Fig. S5). Under the highest gene tree estimation

545    error (300bp), both ASTRAL and MP-EST require 1000 loci to recover the true species tree.

546    For concatenation methods, ML estimations are robust under low ILS levels, which is

547    consistent with previous findings (Mirarab et al. 2014b; Tonini et al. 2015). We were able

548    to recover the correct species tree with the smallest data set (100 loci with 300bp in length)

549    under low ILS (theta = 0.01). However, major challenges and inaccurate species trees are

550    generated under high ILS. Under such conditions, it requires the largest data set (1500 loci

551    with ≥400 bp length) to recover the true species tree (Fig. S5).

552

553    *MSC Model Fitness and the Relative Contribution of ILS, Gene Tree Estimation Error, and Gene*

554    *Flow to Gene Tree Variation*

555        Among all 41,664 triplets we examined, 553 (1.3%) have significant asymmetrical

556    minor frequencies. The node with the highest Reticulation Index is the MRCA of Clade 1 and

557    Clade 2 (the MRCA of Salicaceae and Euphorbiaceae; Fig. 5c). 10.3% of the triplets

558    associated with this node are significantly asymmetrical. According to our relative

559    importance decomposition analysis, ILS, gene tree estimation error, and gene flow explain

560    57.5% of the gene trees variation using the lmg algorithm ($R^2$ = 0.575). When scaling these

561    three factors to sum 100%, gene tree estimation error is the most dominant factor, which

562    explains 52% of the gene tree variation (Fig. S6). The second most significant factor is gene

563    flow, which explains 32% of the gene tree variation. And ILS explains the least variation

564

**Figure 5** Relative contributions of ILS, estimation error, and gene flow across Malpighiales.

566 (a) ILS. Nodes are colored by inferred population mutation parameter theta. (b) Gene tree

567    estimation error. Nodes are colored by BP values, which represent percentage of recovered

568    nodes from simulation (see Materials and Methods). (c) Gene flow. Nodes are colored by

569    Reticulation Index. (d) Gene tree variation. Nodal BPs reflect nodal recovery in gene trees.

570    Percentages of gene tree variation ascribed to ILS, estimation error, and gene flow are

571    indicated by black arrows.

572

573    (15%). The relative ranks of these three factors are consistent among regression methods

574    and bootstrap replicates (Fig. S6).

575            Further investigation revealed significant negative correlation ($p$-value 2.2e-16)

576    between the overall gene tree variation and species tree resolution (Fig. S7a). All of the

577    contributors to gene tree variation—ILS, tree estimation error, and introgression—are

578    strongly negatively correlated with species tree resolution ($p$-value <6.6e-4). We observed

579    the highest level of ILS, introgression and gene tree estimation error for the most

580    recalcitrant nodes along the backbone of the phylogeny using our methods (Fig. S7b–d).

581    This further corroborates our conclusion that a combination of all three factors contribute

582    to this low resolution. We did not find significant correlation between the estimated level of

583    introgression and ILS, suggesting that our triplet method can effectively distinguish these

584    two phenomena. However, both ILS and introgression are positively correlated to gene tree

585    estimation error ($p$-value < 6.8e-3).

586

587    *The triplet-frequency-based method identified three hotspots of introgression in yeasts*

588            Our species tree of yeast inferred using MP-EST is identical to the topology reported

589    in the original study by Salichos and Rokas (2013). We identified 116 asymmetrical triplets

590    among the 1,771 triplets in the yeast species tree. These triplets revealed three hotspots of

591    introgression that correspond to those identified by Yu and Nakhleh (2015): in the MRCA

592    of *Saccharomyces kluyveri* and *Kluyveromyces waltii*, the MRCA of *Zygosaccharomyces rouxii*

593    and *Saccharomyces castellii*, and the MRCA of *Candida guilliermondii* and *Debaryomyces*

594    *hansenii* (Fig. S8). The first two hotspots of reticulation (the MRCA of *S. kluyveri* and *K.*

595    *waltii,* the MRCA of *Z. rouxii* and *S. castellii*) reflect the donor and recipient lineage of one of

596    the two reticulation branches identified by Yu and Nakhleh (2015). The third introgression

597    hotspot involving the MRCA of *C. guilliermondii* and *D. hansenii* reflects the second

598    reticulation branch inferred in Yu and Nakhleh (2015).

599

600    **DISCUSSION**

601            Our results indicate that despite extensive phylogenomic data, the early branching

602    order of Malpighiales remains uncertain. We attribute this to a combination of factors—a

603    perfect storm—involving ILS, gene tree estimation error, and gene flow. Below we highlight

604    our findings in four subsections: the phylogenetic utility of flanking regions in sequence

605    capture data, novel phylogenetic relationships gleaned for Malpighiales, an efficient

606    method to investigate gene flow in large data sets, and a novel simulation-based method to

607    decompose gene tree variation into various contributing factors.

608

609    *Flanking Regions Greatly Enhance Phylogenetic Resolution*

610            Hybrid enrichment probes are designed to capture highly conserved anchor regions

611    as well as the more variable flanking regions adjacent to these anchors. Despite the

612    perceived utility of these flanking regions in mammals (McCormack et al. 2012) and more

613    recently in in plants (Fragoso-Martínez et al. 2017), assumptions of the enhanced

614    phylogenetic utility of these flanking regions have not been tested explicitly to our

615    knowledge. Here, we observed significantly higher average ML BP across gene trees,

616    increased species tree resolution, and most importantly, increased species tree estimation

617    congruency as flanking regions were increasingly added (Fig. S2). This suggests that longer

618    loci, favoring more phylogenetically informative flanking regions, should be prioritized in

619    future anchored hybrid enrichment kit designs. These flanking regions represent genomic

620    regions under nearly neutral selection where mutation rates are high, and thus appear to

621    be a rich source of phylogenetic utility. It has been demonstrated that the inclusion of genes

622    with higher mutation rates can greatly enhance phylogenetic resolution, even deep within

623    organismal phylogenies (Hilu et al. 2003; Lanier et al. 2014). Our site-subsampling strategy,

624    which includes increasingly larger proportions of these more rapidly evolving flanking

625    regions provides the first empirical evidence that these regions are particularly informative

626    for resolving phylogenetic relationships at shallow and deeper phylogenetic depths.

627

628    *Sequence Capture Data Confirms Malpighiales Relationships and Identifies Novel Clades*

629            We assessed the performance of hybrid enrichment markers by evaluating support

630    for major clades previously identified from plastome sequences (Xi et al., 2012; Fig. 3). The

631    majority of the well-supported (>90 BP) clades identified by Xi et al. (2012) are

632    corroborated in our analyses with high confidence (>97 BP). These include the parietal,

633    clusioid, phyllanthoid, ochnoid, chrysobalanoid, and putranjivoid subclades. With rare

634    exception, relationships within these clades were also identical to those by Xi et al. (2012).

635    In the case of the parietal and clusioid clades, internal resolutions were less well supported

636     owing to conflicting topologies recovered among coalescent and concatenation methods

637     (low nodal support indicated by '–' in Fig. 3). Within the parietal clade, for example, the

638     monophyly of the salicoids *sensu* Xi et al. (2012, Fig. 3) is supported by the RAxML

639     phylogeny with moderate support (69 BP) but is not supported in any of the coalescent

640     methods.

641             Additionally, we discovered several noteworthy clades that conflict with those

642     reported by Xi et al. (2012). The euphorbioids, malpighioids, and rhizophoroids were

643     paraphyletic in all of the best resolved MP-EST, ASTRAL, and RAxML analyses (Fig. 3). The

644     euphorbioids—including Euphorbiaceae, Peraceae, Lophopyxidaceae, Linaceae, and

645     Ixonanthaceae—were split into four polyphyletic groups. In particular, Linaceae was

646     placed as sister to the clusioid clade in all of the best resolved coalescent and concatenation

647     analyses (Fig. 3). The affiliation of Linaceae to the clusioids instead of to other members of

648     the euphorbioids is also supported in a recent transcriptomic study of this group with less

649     dense taxon sampling (Cai et al. 2019). Within malpighioids, Centroplacaceae is confidently

650     placed (>86 BP) with Humiriaceae, Pandaceae, and Ctenolophonaceae (Fig. 3) instead of

651     with Malpighiaceae and Elatinaceae. This relationship is partially supported by Wurdack et

652     al. (2004) in which Centroplacaceae was placed with Pandaceae, although with low support.

653     Within the rhizophoroids, Ctenolophonaceae was well nested (>98 BP for coalescent

654     methods) within a clade including Euphorbiaceae and Pandaceae (Clade1 in Fig. 3) rather

655     than with Rhizophoraceae and Erythroxylaceae.

656

657     *ILS and Gene Tree Estimation Error Alone Are Insufficient to Explain the Lack of Species Tree*

658     *Resolution in Malpighiales*

659     Our simulations to explore gene tree heterogeneity encompass the full

660     distributional range of ILS and gene tree estimation error inferred from the empirical data,

661     and clearly demonstrate that the data we have assembled should be sufficient to resolve

662     Malpighiales species tree relationships. Specifically, despite our inability to estimate a well-

663     resolved species tree from our empirical data, we were able to recover a species tree with

664     very high confidence in simulation (mean nodal support >91 BP). This is true even when

665     ILS (theta = 0.1) and gene tree estimation error (alignment length = 300bp) were set to the

666     highest levels inferred from our empirical data. Such extreme levels of theta, in particular,

667     are ten times higher than empirical estimations from *Arabidopsis* and *Drosophila* (0.01–

668     0.001 in both cases; Drost and Lee 1995; Fischer et al. 2017). Even when down sampling

669     our data set under these extreme conditions to include a mere 100 loci, both concatenation

670     and coalescent analyses recover the true species with no more than 10% error (Fig. S5). In

671     addition, we observed far fewer conflicts among species trees reconstructed from different

672     methods and data partitions in simulation versus from those estimated from the empirical

673     data (Fig. 4c). These results suggest that ILS and gene tree estimation error alone are

674     insufficient to explain the lack of resolution along the spine of Malpighiales, and suggest

675     that additional factors likely contribute to gene tree heterogeneity.

676

677     *Gene Flow Compromises Malpighiales Species Tree Resolution: A Novel Method for Assessing*

678     *Gene Tree Heterogeneity*

679     Beyond ILS and gene tree estimation error, gene tree heterogeneity is also

680     attributable to two other common biological factors: gene duplication and gene flow (Yang

681     2006). As we demonstrate above, the first two factors alone are insufficient to explain this

682     lack of resolution. Orthology assignment problems owing to gene duplications are also

683     highly unlikely for two reasons. First, our sequence capture data set was specifically

684     designed for single copy nuclear loci across land plants (Buddenhagen et al. 2016). Second,

685     large-scale genome duplication identified in Malpighiales all occurred *subsequent* to the

686     explosive radiation where discordance is localized (Cai et al. 2019). Thus, biased gene loss

687     arising from genome duplications are unlikely to hinder our ability to resolve backbone

688     relationships in the order. Additional analytical artifacts not reflected in our assessment

689     include homolog calls, alignment error, and most importantly, misspecification of DNA

690     substitution models, all of which can compromise species tree estimation. Though these

691     analytical errors may explain some discordance, it is quite possible that conflicts are

692     attributed to additional biological phenomena.

693            Gene flow has yet to receive attention in phylogenomic studies, especially at deep-

694     time phylogenetic scales. It is estimated that at least 25% of plant species and 10% of

695     animal species hybridize (Mallet 2007) and various network inference methods have been

696     developed to assess gene flow in phylogenies (Nakhleh 2013). These methods have

697     provided valuable insights into reticulate evolution, including those associated with the

698     rapid radiations in wild tomatoes and heliconius butterflies (Pease et al. 2016; Edelman et

699     al. 2019). However, the performance of these methods often relies on accurate species tree

700     estimation and the generation of a handful of alternative species tree topologies to conduct

701     hypothesis testing. However, when alternative topologies are too numerous to evaluate,

702     such as along the backbone of Malpighiales, existing tools become quite limited. In

703     particular, these methods are computationally expensive and amenable only to small data

704     sets. For example, maximum likelihood can only be applied to networks involving fewer

705 than 10 taxa and three reticulations (Yu and Nakhleh 2015). We leveraged the theoretical

706 predictions of triplet frequencies to make inferences about gene flow by summarizing the

707 distribution of lineages involved in horizontal processes using our novel measurement

708 statistic, the Reticulation Index. Our method can effectively identify hotspots of reticulate

709 evolution, including both the donor and recipient lineage, in large clades and in deep time,

710 and provide valuable guidance to empirical studies. We further validated the application of

711 our Reticulation Index using the yeast data set from Salichos and Rokas (2013). The three

712 hotspots we identified in the yeast phylogeny (Fig. S8) correspond precisely to the two

713 reticulation branches previously inferred by Yu and Nakhleh (2015), thus demonstrating

714 the promise of our method for applications in larger phylogenies like Malpighiales.

715       In Malpighiales, the Reticulation Indices are especially high in deeper parts of the

716 phylogeny, suggesting that certain clades may contribute substantially to this phenomenon

717 (Fig. 5c). In particular, we hypothesize that the overabundance of asymmetrical triplets

718 observed within Clades 1 (MRCA of Euphorbiaceae and Putranjivaceae) and Clade 2 (MRCA

719 of Salicaceae and Violaceae) result from ancient and persistent gene flow between early

720 diverging members of these lineages. Specifically, Clade 1 contains six paralogous lineages

721 from the plastid phylogeny (Xi et al. 2012) and is a major hotspot for plastid-nuclear

722 conflict. Such conflict is widely recognized as an indicator of introgression (Soltis and

723 Kuzoff 1995; Baum et al. 1998). Moreover, members of two clades, the putranjivoids and

724 Pandaceae, have previously been implicated in the top three most unstable nodes of all

725 angiosperms (Smith et al. 2013). We hypothesize that this may be attributed to the

726 chimeric nature of their ancestral genealogy resulting from gene flow. The Reticulation

727 Index is also significantly negatively correlated with species tree resolution (Fig. S7d),

728     suggesting that introgression is an important barrier for robust species tree estimation in

729     Malpighiales. For the most recalcitrant nodes where almost no bootstrap replicates recover

730     the same topology, we also observed the highest values of inferred introgression. In the

731     meantime, no correlation is identified between the estimated level of ILS and introgression,

732     suggesting that our methods can effectively distinguish ILS and introgression. However,

733     nodes with strong introgression signals also have higher gene tree estimation error (Fig.

734     S7e). One possible explanation for such correlation is that the short branch lengths created

735     by introgression may lead to elevated estimation error at these nodes.

736            To better characterize gene tree variation attributable to ILS, gene tree estimation

737     error and gene flow, we devised a novel regression method to parse variation attributable

738     to these analytical and biological factors. Our method of decomposing gene tree variation

739     revealed that the majority of variation is due to estimation error (52%), while gene flow

740     and ILS account for 32% and 15%, respectively. This decomposition analysis is based on

741     estimations of ILS, gene tree error, and gene flow through simulation and is subject to

742     common limitations of regression analyses. As a result, errors from the simulation and

743     regression analysis can render the absolute values of these percentages less reliable.

744     Regardless, the relative influence of biological and analytical aspects of gene tree variation

745     as interpreted from these metrics can shed important light on empirical investigations and

746     the development of enhanced species tree inference methods. For example, though gene

747     tree error is to blame for the majority of gene tree variation in our test case, gene flow still

748     plays a significant role in gene tree variation. Therefore, a species network inference

749     method that accommodates gene flow is essential to better understand the early

750     evolutionary history of Malpighiales. Application of this method to other taxonomic groups

751  will also reveal the key factors contributing to recalcitrant relationships and provide

752  guidance for phylogenomic marker design targeting at specific questions.

753      Our results suggest that a confluence of factors—ILS, gene tree estimation error, and

754  gene flow—influence this lack of resolution and contribute to a perfect storm inhibiting our

755  ability to reconstruct branching order along the back of the Malpighiales phylogeny. Gene

756  flow, in particular, is a potentially potent, and overlooked factor accounting for this

757  phenomenon. Despite a relatively small percentage of asymmetrical triplets attributed to

758  gene flow (1.3% of all triplets), they appear to contribute substantially to gene tree

759  heterogeneity based on our relative importance decomposition (32%). Our approach of

760  interrogating this phenomenon using triplet frequencies and the relative importance

761  analyses can elucidate factors that give rise to gene tree variation. These approaches are

762  likely to be especially useful for investigating the causes of recalcitrant relationships,

763  especially at deeper phylogenetic nodes, and to highlight instances where relationships are

764  better modeled as a network rather than a bifurcating tree.

765

766  **ACKNOWLEDGEMENTS**

773

774     **FIGURE CAPTIONS**

775

776     **Figure 1** Simulation of ILS and gene tree estimation error. ILS was simulated though the

777     coalescent process by setting low (0.01) and high (0.1) theta values. DNA alignments were

778     subsequently generated through the mutation process based on simulated gene trees. Five

779     alignments were generated for each gene tree with lengths of 300, 400, 500, 1000, and

780     1500 bp (only two are shown in the graph). Shorter alignment lengths increase in gene tree

781     estimation error.

782

783     **Figure 2** Identification of reticulate evolution using triplet frequency. (a) Theoretical

784     expectations of triplet frequency distribution under the multi-species coalescent (MSC)

785     model with and without introgression. In case of incomplete lineage sorting (ILS),

786     symmetrical distributions of the frequency of two minor topologies are excepted owing to

787     deep coalescence (left). In case of introgression, one of the minor topologies will occur with

788     higher frequency due to gene flow (right). (b) Mapped asymmetrical triplets to species tree

789     to identify reticulate nodes.

790

791     **Figure 3** Species phylogeny of Malpighiales derived from MP-EST with complete low-

792     stringency data set (analysis No. 15 in Table S2). Gene trees are estimated using MrBayes.

793     Branches are colored by the inferred population mutation parameter theta. Warmer colors

794     indicate higher theta and thus higher level of ILS. Terminal branches are colored grey due

795     to lack of data to infer theta. BP values from best-resolved MP-EST/ASTRAL/RAxML

796     analyses (analysis No. 15, 17, and 11 in Table S2) are indicated above each branch; an

797    asterisk indicates 100 BP support; a hyphen indicates less than 50 BP. Branch lengths

798    estimated from RAxML by fixing the species tree topology are presented at the upper left

799    corner. The eleven major clades highlighted in the discussion are identified with circled

800    numbers along each relevant branch. The clade affiliation for each family based on the

801    plastid phylogeny (Xi et al. 2012) is indicated on the right. Clades identified by Xi et al.

802    (2012) that are also monophyletic in this study are highlighted using gray shades.

803

804    **Figure 4** Extensive gene tree discordance in empirical versus simulated data. (a) Gene

805    tree–species tree (G-S) discordance in the empirical (Emp) and simulated (Sim) data

806    assuming fixed theta in simulation. Discordance is measured by RF distance between

807    inferred gene trees and the species tree. Under various simulated conditions of ILS (e.g.,

808    'Low ILS', theta = 0.01 and 'High ILS', theta = 0.1) and gene tree estimation error ('High ILS

809    + High err.', theta = 0.1, alignment length=300bp), the simulated gene tree–species tree

810    discordance is significantly lower than that from empirical data. (b) Gene tree–species tree

811    discordance is higher in empirical versus simulated conditions without setting theta *a*

812    *priori*. For each BP data set, gene tree–species tree discordance is measured and compared

813    in both empirical and simulated data sets. Positive values indicate higher gene tree–species

814    tree discordance in our empirical data. (c) Species tree estimation discordance in empirical

815    data (left) and simulated data (right).

816

817    **Figure 5** Relative contributions of ILS, estimation error, and gene flow across Malpighiales.

818    (a) ILS. Nodes are colored by inferred population mutation parameter theta. (b) Gene tree

819    estimation error. Nodes are colored by BP values, which represent percentage of recovered

820 nodes from simulation (see Materials and Methods). (c) Gene flow. Nodes are colored by

821 Reticulation Index. (d) Gene tree variation. Nodal BPs reflect nodal recovery in gene trees.

822 Percentages of gene tree variation ascribed to ILS, estimation error, and gene flow are

823 indicated by black arrows.

824

825 **Figure S1** Number of PI sites and mean gene tree BP is positively correlated with

826 alignment length in high/medium/low-stringency data sets. (a,b) Correlation between

827 number of PI sites (a) or mean gene tree BP (b) with alignment lengths inferred from the

828 high-stringency data set. (c,d) Correlation between number of PI sites (c) or mean gene tree

829 BP (d) with alignment lengths inferred from the medium-stringency data set. (e,f)

830 Correlation between number of PI sites (e) or mean gene tree BP (f) with alignment lengths

831 inferred from low-stringency data set. Pearson's $R^2$ is presented at lower right corner of

832 each plot.

833

834 **Figure S2** Increased gene tree and species tree resolution as more flanking sites are

835 included in the analysis. (a) Distribution of mean gene tree BP in high/medium/low-

836 stringency data sets. (b,c) Increased species tree BP in concatenation (b) and coalescent

837 analysis (c). Analyses with same locus subsampling are connected by lines. (d) Increased

838 species tree inference consistency reflected by pairwise RF distance.

839

840 **Figure S3** Species tree discordance is more sensitive to site and locus subsampling in

841 coalescent (black) versus concatenation analyses (grey). Left, distribution of pairwise

842 species tree distances derived from all coalescent (black) and concatenation analyses (grey)

843 measured by RF distance. Right, distribution of pairwise species tree distances from

844 coalescent (black) and concatenation (grey) analyses measured by weighted RF (WRF)

845 distance (weighted by nodal support).

846

847 **Figure S4** Gene tree estimation error in empirical and simulated data. Gene tree estimation

848 error is measured by RF distance to the 'true gene tree' for both empirical and simulated

849 data sets. In both cases, gene tree estimation error is negatively correlated with alignment

850 length.

851

852 **Figure S5** Species tree estimation error in simulated data sets. Species tree estimation

853 error is measured by RF distance from inferred species in each analysis to the known

854 species tree. Results derived from alignments of varying lengths (300, 400, 500, 1000, 1500

855 bp) are marked by different color and shape. (a,b) Species tree estimation error of ExaML

856 under low (a) and high (b) ILS. (c,d) Species tree estimation error of MP-EST under low (c)

857 and high (d) ILS. (e,f) Species tree estimation error of ASTRAL-II under low (e) and high (f)

858 ILS.

859

860 **Figure S6** Relative importance of ILS, gene tree estimation error, and gene flow in

861 generating gene tree variation based on four regression methods. Percentages are

862 normalized to sum 100%. 95% confidence intervals are represented by bars.

863

864 **Figure S7** ILS, gene tree estimation error, and introgression contribute to low species tree

865 resolution in Malpighiales. Species tree resolution is represented by nodal support from the

866    MP-EST phylogeny in Figure 3 from the main text. The other statistics reflect the variables

867    presented in Figure 5. The *p*-value of the Pearson's correlation test is indicated in the upper

868    right corner in each panel. (a) Significant negative correlation between gene tree variation

869    and species tree resolution. (b) Significant negative correlation between ILS and species

870    tree resolution. (c) Significant negative correlation between gene tree estimation error and

871    species tree resolution. (d) Significant negative correlation between introgression and

872    species tree resolution. (e) Significant positive correlation between gene tree estimation

873    error and introgression. (f) No significant correlation between gene tree estimation error

874    and ILS.

875

876    **Figure S8** Hotspots of reticulate evolution in baker's yeast. Species phylogeny is inferred

877    from MP-EST with the 1,070 genes trees in Salichos and Rokas (2013). Nodes are colored

878    by Reticulation Index. Black thick arrows indicate inferred reticulation by Yu and Nakhleh

879    (2015) for comparative purpose.

880

881    **Table S1** Voucher and GenBank information for 64 species in Malpighiales, Celastrales,

882    Oxalidales, and Vitales used for anchored hybrid enrichment.

883

884    **Table S2** Species tree estimation strategies using various phylogenetic estimation methods

885    and phylogenetic subsampling methods (see Supplementary Note 1) with high-, medium-,

886    and low-stringency data sets.

887

888    **Table S3** Coalescent and mutational parameters of simulated data sets.

889

890    **Table S4** Summary statistics of 423 loci in high/medium/low-stringency data sets,

891    including number of captured taxa, alignment length, number of PI sites, and mean gene

892    tree BP.

893

894

895

896    **LITERATURE CITED**

897    Arcila, D., Ortí, G., Vari, R., Armbruster, J.W., Stiassny, M.L., Ko, K.D., Sabaj, M.H., Lundberg, J.,
898        Revell, L.J., Betancur-R, R. 2017. Genome-wide interrogation advances resolution of
899        recalcitrant groups in the tree of life. Nat Ecol Evol, 1:1-10.
900    Baum, D.A., Small, R.L., Wendel, J.F. 1998. Biogeography and floral evolution of baobabs
901        (*Adansonia*, Bombacaceae) as inferred from multiple data sets. Syst Biol, 47:181-
902        207.
903    Beckman, E.J., Benham, P.M., Cheviron, Z.A., Witt, C.C. 2018. Detecting introgression despite
904        phylogenetic uncertainty: The case of the South American siskins. Mol Ecol,
905        27:4350-4367.
906    Blom, M.P., Bragg, J.G., Potter, S., Moritz, C. 2017. Accounting for uncertainty in gene tree
907        estimation: summary-coalescent species tree inference in a challenging radiation of
908        Australian lizards. Syst Biol, 66:352-366.
909    Buddenhagen, C., Lemmon, A.R., Lemmon, E.M., Bruhl, J., Cappa, J., Clement, W.L., Donoghue,
910        M., Edwards, E.J., Hipp, A.L., Kortyna, M. 2016. Anchored phylogenomics of
911        angiosperms I: Assessing the robustness of phylogenetic estimates. bioRxiv:086298.
912    Cai, L., Xi, Z., Amorim, A.M., Sugumaran, M., Rest, J.S., Liu, L., Davis, C.C. 2019. Widespread
913        ancient whole-genome duplications in Malpighiales coincide with Eocene global
914        climatic upheaval. New Phytol, 221:565-576.
915    Capella-Gutiérrez, S., Silla-Martínez, J.M., Gabaldón, T. 2009. trimAl: a tool for automated
916        alignment trimming in large-scale phylogenetic analyses. Bioinformatics, 25:1972-
917        1973.
918    Chase, M.W., Christenhusz, M., Fay, M., Byng, J., Judd, W.S., Soltis, D., Mabberley, D.,
919        Sennikov, A., Soltis, P.S., Stevens, P.F. 2016. An update of the Angiosperm Phylogeny
920        Group classification for the orders and families of flowering plants: APG IV. Bot J
921        Linn Soc, 181:1-20.
922    Chifman, J., Kubatko, L. 2014. Quartet inference from SNP data under the coalescent model.
923        Bioinformatics, 30:3317-3324.

924 Cox, C.J., Li, B., Foster, P.G., Embley, T.M., Civáň, P. 2014. Conflicting phylogenies for early
925     land plants are caused by composition biases among synonymous substitutions. Syst
926     Biol, 63:272-279.
927 Davis, C.C., Webb, C.O., Wurdack, K.J., Jaramillo, C.A., Donoghue, M.J. 2005. Explosive
928     radiation of Malpighiales supports a mid-Cretaceous origin of modern tropical rain
929     forests. Am Nat, 165:E36-E65.
930 Degnan, J.H., Rosenberg, N.A. 2006. Discordance of species trees with their most likely gene
931     trees. PLOS Genet, 2:e68.
932 Drost, J.B., Lee, W.R. 1995. Biological basis of germline mutation: comparisons of
933     spontaneous germline mutation rates among drosophila, mouse, and human.
934     Environ Mol Mutagen, 25:48-64.
935 Durand, E.Y., Patterson, N., Reich, D., Slatkin, M. 2011. Testing for ancient admixture
936     between closely related populations. Mol Biol Evol, 28:2239-2252.
937 Edelman, N.B., Frandsen, P.B., Miyagi, M., Clavijo, B., Davey, J., Dikow, R.B., García-Accinelli,
938     G., Van Belleghem, S.M., Patterson, N., Neafsey, D.E. 2019. Genomic architecture and
939     introgression shape a butterfly radiation. Science, 366:594-599.
940 Edwards, S.V., Xi, Z., Janke, A., Faircloth, B.C., McCormack, J.E., Glenn, T.C., Zhong, B., Wu, S.,
941     Lemmon, E.M., Lemmon, A.R. 2016. Implementing and testing the multispecies
942     coalescent model: a valuable paradigm for phylogenomics. Mol Phylogenet Evol,
943     94:447-462.
944 Fischer, M.C., Rellstab, C., Leuzinger, M., Roumet, M., Gugerli, F., Shimizu, K.K., Holderegger,
945     R., Widmer, A. 2017. Estimating genomic diversity and population differentiation–an
946     empirical comparison of microsatellite and SNP variation in *Arabidopsis halleri*. BMC
947     Genomics, 18:69.
948 Fragoso-Martínez, I., Salazar, G.A., Martínez-Gordillo, M., Magallón, S., Sánchez-Reyes, L.,
949     Lemmon, E.M., Lemmon, A.R., Sazatornil, F., Mendoza, C.G. 2017. A pilot study
950     applying the plant Anchored Hybrid Enrichment method to New World sages (*Salvia*
951     subgenus *Calosphace*; Lamiaceae). Mol Phylogenet Evol, 117:124-134.
952 Glémin, S., Scornavacca, C., Dainat, J., Burgarella, C., Viader, V., Ardisson, M., Sarah, G.,
953     Santoni, S., David, J., Ranwez, V. 2019. Pervasive hybridizations in the history of
954     wheat relatives. Sci Adv, 5:eaav9188.
955 Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li, H.,
956     Zhai, W., Fritz, M.H.-Y. 2010. A draft sequence of the Neandertal genome. Science,
957     328:710-722.
958 Grömping, U. 2006. Relative importance for linear regression in R: the package relaimpo. J
959     Stat Softw, 17:1-27.
960 Guéguen, L., Gaillard, S., Boussau, B., Gouy, M., Groussin, M., Rochette, N.C., Bigot, T.,
961     Fournier, D., Pouyet, F., Cahais, V. 2013. Bio++: efficient extensible libraries and tools
962     for computational molecular evolution. Mol Biol Evol, 30:1745-1750.
963 Hahn, M.W., Nakhleh, L. 2016. Irrational exuberance for resolved species trees. Evolution
964     (N Y), 70:7-17.
965 Hamilton, C.A., Lemmon, A.R., Lemmon, E.M., Bond, J.E. 2016. Expanding anchored hybrid
966     enrichment to resolve both deep and shallow relationships within the spider tree of
967     life. BMC Evol Biol, 16:212.

968    Hilu, K.W., Borsch, T., Müller, K., Soltis, D.E., Soltis, P.S., Savolainen, V., Chase, M.W., Powell,
969         M.P., Alice, L.A., Evans, R. 2003. Angiosperm phylogeny based on matK sequence
970         information. Am J Bot, 90:1758-1776.
971    Hosner, P.A., Faircloth, B.C., Glenn, T.C., Braun, E.L., Kimball, R.T. 2015. Avoiding missing
972         data biases in phylogenomic inference: an empirical study in the landfowl (Aves:
973         Galliformes). Mol Biol Evol:msv347.
974    Huang, H., Knowles, L.L. 2009. What is the danger of the anomaly zone for empirical
975         phylogenetics? Syst Biol, 58:527-536.
976    Huson, D.H., Klöpper, T., Lockhart, P.J., Steel, M.A. 2005. Reconstruction of reticulate
977         networks from gene trees. Annual International Conference on Research in
978         Computational Molecular Biology, Springer, p. 233-249.
979    Arabidopsis Genome Initiative. 2000. Analysis of the genome sequence of the flowering
980         plant Arabidopsis thaliana. Nature, 408:796.
981    Jarvis, E.D., Mirarab, S., Aberer, A.J., Li, B., Houde, P., Li, C., Ho, S.Y., Faircloth, B.C., Nabholz,
982         B., Howard, J.T. 2014. Whole-genome analyses resolve early branches in the tree of
983         life of modern birds. Science, 346:1320-1331.
984    Katoh, K., Standley, D.M. 2013. MAFFT multiple sequence alignment software version 7:
985         improvements in performance and usability. Mol Biol Evol, 30:772-780.
986    Kozlov, A.M., Aberer, A.J., Stamatakis, A. 2015. ExaML version 3: a tool for phylogenomic
987         analyses on supercomputers. Bioinformatics, 31:2577-2579.
988    Lanfear, R., Calcott, B., Ho, S.Y., Guindon, S. 2012. PartitionFinder: combined selection of
989         partitioning schemes and substitution models for phylogenetic analyses. Mol Biol
990         Evol, 29:1695-1701.
991    Lanier, H.C., Huang, H., Knowles, L.L. 2014. How low can you go? The effects of mutation
992         rate on the accuracy of species-tree estimation. Mol Phylogenet Evol, 70:112-119.
993    Lartillot, N., Philippe, H. 2004. A Bayesian mixture model for across-site heterogeneities in
994         the amino-acid replacement process. Mol Biol Evol, 21:1095-1109.
995    Lartillot, N., Rodrigue, N., Stubbs, D., Richer, J. 2013. PhyloBayes MPI: phylogenetic
996         reconstruction with infinite mixtures of profiles in a parallel environment. Syst Biol,
997         62:611-615.
998    Lemmon, A.R., Brown, J.M., Stanger-Hall, K., Lemmon, E.M. 2009. The effect of ambiguous
999         data on phylogenetic estimates obtained by maximum likelihood and Bayesian
1000        inference. Syst Biol, 58:130-145.
1001   Lemmon, A.R., Emme, S.A., Lemmon, E.M. 2012. Anchored hybrid enrichment for massively
1002        high-throughput phylogenomics. Syst Biol:sys049.
1003   Lemmon, E.M., Lemmon, A.R. 2013. High-throughput genomic data in systematics and
1004        phylogenetics. Annu Rev Ecol Evol Syst, 44:99-121.
1005   Lindeman, R.H. 1980. Introduction to bivariate and multivariate analysis.
1006   Liu, L., Xi, Z., Wu, S., Davis, C.C., Edwards, S.V. 2015. Estimating phylogenetic trees from
1007        genome-scale data. Ann N Y Acad Sci, 1360:36-53.
1008   Liu, L., Yu, L. 2010. Phybase: an R package for species tree analysis. Bioinformatics, 26:962-
1009        963.
1010   Liu, L., Yu, L., Edwards, S.V. 2010. A maximum pseudo-likelihood approach for estimating
1011        species trees under the coalescent model. BMC Evol Biol, 10:302.
1012   Liu, L., Yu, L., Kubatko, L., Pearl, D.K., Edwards, S.V. 2009. Coalescent methods for estimating
1013        phylogenetic trees. Mol Phylogenet Evol, 53:320-328.

1014 Magallon, S., Crane, P.R., Herendeen, P.S. 1999. Phylogenetic pattern, diversity, and
1015     diversification of eudicots. Annals of the Missouri Botanical Garden:297-372.
1016 Mallet, J. 2007. Hybrid speciation. Nature, 446:279-283.
1017 McCormack, J.E., Faircloth, B.C., Crawford, N.G., Gowaty, P.A., Brumfield, R.T., Glenn, T.C.
1018     2012. Ultraconserved elements are novel phylogenomic markers that resolve
1019     placental mammal phylogeny when combined with species-tree analysis. Genome
1020     Res, 22:746-754.
1021 Meng, C., Kubatko, L.S. 2009. Detecting hybrid speciation in the presence of incomplete
1022     lineage sorting using gene tree incongruence: a model. Theor Popul Biol, 75:35-45.
1023 Meyer, B.S., Matschiner, M., Salzburger, W. 2017. Disentangling incomplete lineage sorting
1024     and introgression to refine species-tree estimates for Lake Tanganyika cichlid fishes.
1025     Syst Biol, 66:531-550.
1026 Mirarab, S., Bayzid, M.S., Boussau, B., Warnow, T. 2014a. Statistical binning enables an
1027     accurate coalescent-based estimation of the avian tree. Science, 346:1250463.
1028 Mirarab, S., Bayzid, M.S., Warnow, T. 2014b. Evaluating summary methods for multilocus
1029     species tree estimation in the presence of incomplete lineage sorting. Syst Biol,
1030     65:366-380.
1031 Mirarab, S., Reaz, R., Bayzid, M.S., Zimmermann, T., Swenson, M.S., Warnow, T. 2014c.
1032     ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics,
1033     30:i541-i548.
1034 Mirarab, S., Warnow, T. 2015. ASTRAL-II: coalescent-based species tree estimation with
1035     many hundreds of taxa and thousands of genes. Bioinformatics, 31:i44-i52.
1036 Nakhleh, L. 2013. Computational approaches to species phylogeny inference and gene tree
1037     reconciliation. Trends Ecol Evol, 28:719-728.
1038 Paradis, E., Claude, J., Strimmer, K. 2004. APE: analyses of phylogenetics and evolution in R
1039     language. Bioinformatics, 20:289-290.
1040 Pease, J.B., Haak, D.C., Hahn, M.W., Moyle, L.C. 2016. Phylogenomics reveals three sources of
1041     adaptive variation during a rapid radiation. PLOS Biol, 14.
1042 Philippe, H., Brinkmann, H., Lavrov, D.V., Littlewood, D.T.J., Manuel, M., Wörheide, G.,
1043     Baurain, D. 2011. Resolving difficult phylogenetic questions: why more sequences
1044     are not enough. PLOS Biol, 9:e1000602.
1045 Pratt, J.W. 1987. Dividing the indivisible: Using simple symmetry to partition variance
1046     explained. Proceedings of the second international Tampere conference in statistics,
1047     1987, Department of Mathematical Sciences, University of Tampere, p. 245-260.
1048 Prum, R.O., Berv, J.S., Dornburg, A., Field, D.J., Townsend, J.P., Lemmon, E.M., Lemmon, A.R.
1049     2015. A comprehensive phylogeny of birds (Aves) using targeted next-generation
1050     DNA sequencing. Nature, 526:569.
1051 Reddy, S., Kimball, R.T., Pandey, A., Hosner, P.A., Braun, M.J., Hackett, S.J., Han, K.-L.,
1052     Harshman, J., Huddleston, C.J., Kingston, S. 2017. Why do phylogenomic data sets
1053     yield conflicting trees? Data type influences the avian tree of life more than taxon
1054     sampling. Syst Biol, 66:857-879.
1055 Roch, S., Warnow, T. 2015. On the robustness to gene tree estimation error (or lack thereof)
1056     of coalescent-based species tree methods. Syst Biol, 64:663-676.
1057 Rokas, A., Ladoukakis, E., Zouros, E. 2003. Animal mitochondrial DNA recombination
1058     revisited. Trends Ecol Evol, 18:411-417.

1059   Rokyta, D.R., Lemmon, A.R., Margres, M.J., Aronow, K. 2012. The venom-gland
1060         transcriptome of the eastern diamondback rattlesnake (*Crotalus adamanteus*). BMC
1061         Genomics, 13:312.
1062   Ronquist, F., Huelsenbeck, J.P. 2003. MrBayes 3: Bayesian phylogenetic inference under
1063         mixed models. Bioinformatics, 19:1572-1574.
1064   Rosenberg, N.A., Tao, R. 2008. Discordance of species trees with their most likely gene
1065         trees: the case of five taxa. Syst Biol, 57:131-140.
1066   Saitou, N., Nei, M. 1987. The neighbor-joining method: a new method for reconstructing
1067         phylogenetic trees. Mol Biol Evol, 4:406-425.
1068   Salichos, L., Rokas, A. 2013. Inferring ancient divergences requires genes with strong
1069         phylogenetic signals. Nature, 497:327-331.
1070   Shen, X. X., Hittinger, C.T., Rokas, A. 2017. Contentious relationships in phylogenomic
1071         studies can be driven by a handful of genes. Nat Ecol Evol, 1:0126.
1072   Slatkin, M., Pollack, J.L. 2008. Subdivision in an ancestral species creates asymmetry in gene
1073         trees. Mol Biol Evol, 25:2241-2246.
1074   Smith, S.A., Brown, J.W., Hinchliff, C.E. 2013. Analyzing and synthesizing phylogenies using
1075         tree alignment graphs. PLOS Comput Biol, 9:e1003223.
1076   Solís-Lemus, C., Bastide, P., Ané, C. 2017. PhyloNetworks: a package for phylogenetic
1077         networks. Mol Biol Evol, 34:3292-3298.
1078   Soltis, D.E., Kuzoff, R.K. 1995. Discordance between nuclear and chloroplast phylogenies in
1079         the *Heuchera* group (Saxifragaceae). Evolution (N Y), 49:727-742.
1080   Soltis, P., Soltis, D., Edwards, C. 2005. Angiosperms, Flowering Plants. The Tree of Life Web
1081         Project, http://tolweb. org/Version, 3.
1082   Song, S., Liu, L., Edwards, S.V., Wu, S. 2012. Resolving conflict in eutherian mammal
1083         phylogeny using phylogenomics and the multispecies coalescent model. Proc Natl
1084         Acad Sci USA, 109:14942-14947.
1085   Stamatakis, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
1086         large phylogenies. Bioinformatics, 30:1312-1313.
1087   Stevens, P.F., Davis, H. 2001. Angiosperm phylogeny website.
1088   Sun, M., Soltis, D.E., Soltis, P.S., Zhu, X., Burleigh, J.G., Chen, Z. 2015. Deep phylogenetic
1089         incongruence in the angiosperm clade Rosidae. Mol Phylogenet Evol, 83:156-166.
1090   Tonini, J., Moore, A., Stern, D., Shcheglovitova, M., Ortí, G. 2015. Concatenation and species
1091         tree methods exhibit statistically indistinguishable accuracy under a range of
1092         simulated conditions. PLOS Curr, 7.
1093   Tuskan, G.A., Difazio, S., Jansson, S., Bohlmann, J., Grigoriev, I., Hellsten, U., Putnam, N.,
1094         Ralph, S., Rombauts, S., Salamov, A. 2006. The genome of black cottonwood, *Populus
1095         trichocarpa* (Torr. & Gray). Science, 313:1596-1604.
1096   Wagner, C.E., Keller, I., Wittwer, S., Selz, O.M., Mwaiko, S., Greuter, L., Sivasundar, A.,
1097         Seehausen, O. 2013. Genome-wide RAD sequence data provide unprecedented
1098         resolution of species boundaries and relationships in the Lake Victoria cichlid
1099         adaptive radiation. Mol Ecol, 22:787-798.
1100   Whitfield, J.B., Kjer, K.M. 2008. Ancient rapid radiations of insects: challenges for
1101         phylogenetic analysis. Annu Rev Entomol, 53:449-472.
1102   Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam,
1103         S., Barker, M.S., Burleigh, J.G., Gitzendanner, M.A. 2014. Phylotranscriptomic analysis

1104    of the origin and early diversification of land plants. Proc Natl Acad Sci USA,
1105    111:E4859-E4868.
1106 Wurdack, K.J., Davis, C.C. 2009. Malpighiales phylogenetics: gaining ground on one of the
1107    most recalcitrant clades in the angiosperm tree of life. Am J Bot, 96:1551-1570.
1108 Xi, Z., Liu, L., Davis, C.C. 2015. Genes with minimal phylogenetic information are
1109    problematic for coalescent analyses when gene tree estimation is biased. Mol
1110    Phylogenet Evol, 92:63-71.
1111 Xi, Z., Liu, L., Rest, J.S., Davis, C.C. 2014. Coalescent versus concatenation methods and the
1112    placement of Amborella as sister to water lilies. Syst Biol, 63:919-932.
1113 Xi, Z., Rest, J.S., Davis, C.C. 2013. Phylogenomics and coalescent analyses resolve extant seed
1114    plant relationships. PLOS One, 8:e80870.
1115 Xi, Z., Ruhfel, B.R., Schaefer, H., Amorim, A.M., Sugumaran, M., Wurdack, K.J., Endress, P.K.,
1116    Matthews, M.L., Stevens, P.F., Mathews, S. 2012. Phylogenomics and a posteriori data
1117    partitioning resolve the Cretaceous angiosperm radiation Malpighiales. Proc Natl
1118    Acad Sci USA, 109:17519-17524.
1119 Xu, B., Yang, Z. 2016. Challenges in species tree estimation under the multispecies
1120    coalescent model. Genetics, 204:1353-1368.
1121 Yang, Z. 2006. Computational molecular evolution. Oxford University Press.
1122 Yu, Y., Nakhleh, L. 2015. A maximum pseudo-likelihood approach for phylogenetic
1123    networks. BMC Genomics, 16:S10.
1124 Yu, Y., Than, C., Degnan, J.H., Nakhleh, L. 2011. Coalescent histories on phylogenetic
1125    networks and detection of hybridization despite incomplete lineage sorting. Syst
1126    Biol, 60:138-149.
1127 Zwickl, D.J., Hillis, D.M. 2002. Increased taxon sampling greatly reduces phylogenetic error.
1128    Syst Biol, 51:588-598.
1129 Zwickl, D.J., Stein, J.C., Wing, R.A., Ware, D., Sanderson, M.J. 2014. Disentangling
1130    methodological and biological sources of gene tree discordance on *Oryza* (Poaceae)
1131    chromosome 3. Syst Biol, 63:645-659.
1132
1133