1     **Running Head: Phased Target Enrichment for Polyploids**

2     **Phasing Alleles Improves Network Inference with Allopolyploids**

3     George P. Tiley[1,†,*], Andrew A. Crowl[1,†], Paul S. Manos[1], Emily B. Sessa[2], Claudia Solís-

4     Lemus[3], Anne D. Yoder[1], J. Gordon Burleigh[2]

5     [1]Department of Biology, Duke University, Durham NC, 27708, USA

6     [2]Department of Biology, University of Florida, Gainesville FL, 32611, USA

7     [3]Wisconsin Institute for Discovery and Department of Plant Pathology, University of Wisconsin –

8     Madison, Madison WI, 53706, USA

9     [†]These authors contributed equally

10     [*]Author for correspondence: george.tiley@duke.edu

11

12

13    **Abstract**

14         Accurately reconstructing the reticulate histories of polyploids remains a central

15    challenge for understanding plant evolution. Although phylogenetic networks can provide

16    insights into relationships among polyploid lineages, inferring networks may be hampered by the

17    complexities of homology determination in polyploid taxa. We use simulations to show that

18    phasing alleles from allopolyploid individuals can improve inference of phylogenetic networks

19    under the multispecies coalescent. Phased allelic data can also improve divergence time

20    estimates for networks, which is helpful for evaluating allopolyploid speciation hypotheses and

21    proposing mechanisms of speciation. To achieve these outcomes, we present a novel pipeline

22    that leverages a recently developed phasing algorithm to reliably phase alleles from polyploids.

23    This pipeline is especially appropriate for target enrichment data, where depth of coverage is

24    typically high enough to phase entire loci. We provide an empirical example in the North

25    American *Dryopteris* fern complex that demonstrates how phasing can help reveal the mode of

26    polyploidization and improve network inference. We establish that our pipeline (PATÉ: Phased

27    Alleles from Target Enrichment data) is capable of recovering a high proportion of phased loci

28    from both diploids and polyploids, and that these data improve network estimates compared to

29    using haplotype consensus assemblies. This approach is shown to be especially effective in

30    reticulate complexes where there are multiple hybridization events. The pipeline is available at:

31    https://github.com/gtiley/Phasing.

32

33    **Key words:** Introgression; Hybridization; Reticulate Evolution; Multispecies Coalescent;

34    Divergence Time Estimation; Polyploidy; Target Enrichment; *Dryopteris*

35

36

37   **INTRODUCTION**

38          The phenomenon of polyploidy, or whole-genome duplication, occurs throughout the tree

39   of life. Nowhere, perhaps, is its evolutionary significance more evident than in plants, with recent

40   estimates suggesting up to 35% of vascular plant species are of recent polyploid origin (Wood

41   et al. 2009; Barker et al. 2016). Despite advances in genomic data generation and a long-term

42   interest in understanding the role of whole-genome duplication in driving plant speciation and

43   local adaptation (reviewed in Soltis et al. 2014), polyploids remain a central challenge for the

44   field of phylogenetics. One persistent problem when analyzing sequence data from polyploid

45   taxa, and especially allopolyploids, is identifying the alleles and divergent homeolog copies from

46   parental lineages. Most bioinformatic tools for processing next generation sequence data were

47   developed with diploids, or specifically humans, in mind. These approaches often collapse

48   variable homeolog sequences into a single consensus sequence for *de novo* assemblies or

49   assume the organism is diploid when performing genotyping and phasing for reference-based

50   assembly. For polyploids, this creates chimeric sequences that may interfere with phylogenetic

51   reconstruction and obscure signals of polyploidy and polyploid mode-of-origin. Using allelic data

52   that more accurately capture the complex genomic histories of polyploids should enable the

53   incorporation of divergent signals from polyploid loci into phylogenomic inference, distinguish

54   allopolyploidy from autopolyploidy, and identify parental taxa. However, few studies have

55   examined the potential benefits of using phased versus unphased data to reconstruct polyploid

56   histories, and there are few formal methods and little guidance for phasing alleles from polyploid

57   taxa. Here we explore the value of using phased data to reconstruct polyploid networks

58   leverage recent algorithmic advances in polyploid phasing (Xie et al. 2016) to develop a

59   bioinformatic pipeline that can phase alleles from polyploids using target enrichment sequence

60   data.

61          Previous studies of reticulate complexes have suggested phasing alleles is crucial for

62   accurate evolutionary reconstruction, at least when sampling relatively few loci (e.g., 4 to 10;

3

63    Rothfels et al. 2017; Eriksson et al. 2018). The applications of phased sequencing data for

64    phylogenomic studies of polyploid complexes are thus enticing; however, it remains challenging

65    to genotype and phase next generation sequencing data from polyploids. Methods exist to

66    genotype consensus loci from target enrichment data, but these have been either limited to

67    diploids (Kates et al. 2018; Andermann et al. 2019) or manual curation of variants with

68    polyploids where both parental populations are available (Eriksson et al. 2018). Otherwise,

69    obtaining phased sequence data for polyploids has largely depended on costly long-read

70    sequencing to recover complete haplotype sequences (e.g., Rothfels et al. 2017), or cloning

71    PCR products (e.g., Sessa et al. 2012; Oberprieler et al. 2017).

72        Target enrichment (or HybSeq), where specific regions of the genome are isolated and

73    sequenced (Faircloth et al. 2012; Lemmon et al. 2012), is an increasingly common method for

74    collecting large-scale phylogenomic datasets, and these data can provide insights into the

75    evolutionary history of reticulate complexes (e.g., Crowl et al. 2017; Karimi et al. 2020) and

76    sources of gene tree discordance (e.g., Morales-Briones et al. 2018; Stull et al. 2020). Probe

77    kits for target enrichment have been developed in many land plant lineages (Wolf et al. 2018;

78    Johnson et al. 2019; Liu et al. 2019; Breinholt et al. 2021), and there are bioinformatic pipelines

79    available for custom probe design (e.g., Jantzen et al. 2020). The most common approach to

80    assemble phylogenetic datasets from target enrichment data has been to use *de novo* assembly

81    pipelines (Faircloth 2016; Johnson et al. 2016; Andermann et al. 2018; Breinholt et al. 2018).

82    The assembly algorithms within these pipelines typically treat variable base calls as sequencing

83    errors and consider only the most frequent nucleotide sequence while discarding the

84    alternatives (Bankevich et al. 2012; Iqbal et al. 2012; Luo et al. 2012). This results in loci in

85    which variable positions are collapsed to a single base call (haplotype consensus loci), losing

86    information related to heterozygosity. While this may be appropriate, or at least benign, for

87    phylogenetic analyses of diploid taxa (Kates et al. 2018), it may pose substantial problems when

88    attempting to investigate the evolutionary history of polyploid taxa or reticulate lineages.

4

89   While phylogenetic studies in plants often infer strictly bifurcating trees, the complexities

90 of allopolyploid evolution can be represented more accurately using networks. Simulations and

91 empirical analyses have suggested that phylogenetic networks can recover the reticulate

92 histories of polyploid lineages with few loci, at least when gene tree discordance due to

93 incomplete lineage sorting (ILS; Hudson et al. 1983; Pamilo and Nei 1988) is low (Oberprieler et

94 al. 2017) using parsimony methods (Huber et al. 2006; Lott et al. 2009), or even with moderate

95 ILS when explicitly modeled (Jones et al. 2013). Contemporary phylogenetic network models

96 and software packages jointly consider gene tree variation due to allele sampling error as

97 described by the multispecies coalescent (MSC; Rannala and Yang 2003) and gene flow

98 modeled as episodic introgression events (Solis-Lemus and Ané 2016; Wen et al. 2016; Zhang,

99 Ogilvie et al. 2018, Flouri et al. 2020). Depending on the complexity and goals of the research

100 question, these methods can search for networks with a constrained number of reticulation

101 events using quartet-based maximum pseudolikelihood (Solis-Lemus and Ané 2016; Wen et al.

102 2018) or a full-likelihood Bayesian model where the number of reticulations is a parameter (Wen

103 et al. 2018; Zhang et al. 2018). Also, it is possible to estimate model parameters (divergence

104 times, population sizes, and the fraction of introgressed genes) on a fixed species network using

105 a full-likelihood Bayesian model that allows efficient computation with large numbers of loci

106 (Flouri et al. 2020). We refer to these network models from here on as the multispecies

107 coalescent with introgression (MSci), consistent with Flouri et al. (2020); although, other names

108 have been used, such as the network multispecies coalescent (NMSC; e.g Zhu and Degnan

109 2017) and multispecies network coalescent (MSNC; e.g., Wen et al. 2016). We emphasize that

110 network approaches to investigate polyploid complexes are not novel (e.g., Huber et al. 2006;

111 Lott et al. 2009; Jones et al. 2013), but the difficulty of collecting appropriate genomic data from

112 polyploids for such analyses has limited their use.

113   To address the issues outlined above, we have developed a pipeline, PATÉ (Phased

114 Alleles from Target Enrichment data), that can phase genotyping data for individuals of a known

115   ploidy without the need for sampling their parental lineages. PATÉ was designed with scalability

116   and population-level sampling in mind for target enrichment projects where deep coverage from

117   paired-end Illumina data allow calling of high-quality variants. In this study, we first use

118   simulations to explore the ability of network approaches to reconstruct the history of

119   allopolyploidy in the presence of ILS, and whether phasing the data affects the accuracy of the

120   reconstruction. We show that using phased allelic data can improve network estimation and

121   divergence time estimation compared to using haplotype consensus sequences, but also

122   highlight scenarios where phasing may not be necessary or beneficial. Next, we describe the

123   individual steps used by PATÉ to phase target enrichment data. For an empirical example of the

124   benefits of PATÉ, we compare phased and unphased (haplotype consensus) data to infer the

125   evolutionary history of the North American *Dryopteris* complex, a model system for reticulate

126   polyploid evolution (Sessa et al. 2012a; Sessa et al. 2012b), using new targeted enrichment

127   data. The system includes four diploid species, as well as one extinct diploid, that have formed

128   five allopolyploids in which there is high confidence in the parent-progeny relationships (Fig. 1),

129   although numerous sterile allopolyploid species have also been reported within the complex

130   (Montgomery and Paulton 1981). The allopolyploid species have relatively ancient origins, with

131   the best estimates placing hybridization events between six and 13 Ma (Sessa et al. 2012b).

132   PATÉ is largely successful in recovering phased haploid sequences from polyploid individuals,

133   and networks inferred from phased data more accurately represent the North American

134   *Dryopteris* complex than those inferred from unphased data.

135

136   **MATERIALS & METHODS**

137   ***Testing the Effects of Phasing on Network Inference through Simulation***

138           *Simulating phased and unphased sequence data for an allopolyploid* — We simulated

139   gene trees and their nucleotide sequence data using the MSC model with BPP v1.4.1 (Flouri et

140   al. 2018) under a five-species network (Fig. 2). All simulations used the Jukes-Cantor (Jukes

6

141    and Cantor 1969) model of sequence evolution with no rate heterogeneity. The allopolyploid

142    species $E$ was treated as two lineages ($E$ sister to $B$ and $F$ sister to $C$) whose alleles were

143    pooled to form the hybrid species $E$ at time $\tau_h$. This makes species $E$ a tetraploid hybrid with the

144    parents $B$ and $C$. $\theta$ was constant among lineages and set at 0.01. Assuming a per-generation

145    mutation rate ($\mu$) of $1 \times 10^{-8}$ and one year per generation yields an effective population size

146    ($N_e$) of 250,000 and root age of 10 Ma for the simulation network (Fig. 2). We also simulated

147    data under a shallow divergence scenario, in which all node ages were divided by 10, and a

148    deep divergence scenario, in which all node ages were multiplied by 10. This changes the root

149    age ($\tau_r$) to 0.01 (1 Ma) and 1.0 (100 Ma), respectively.

150       Because the distance between speciation nodes is 0.025 and $\theta = 0.01$ for the simulation

151    network (Fig. 2), there are five coalescent units between nodes, $T = \frac{\tau}{\left(\frac{\theta}{2}\right)}$ (Yang 2006), which

152    implies a near-zero probability of gene tree discordance due to ILS (Hudson 1983). We

153    incorporated ILS into the simulation by reducing the node heights $\tau_u$ and $\tau_s$ to 0.0375 and 0.05

154    (for a low level of ILS) and 0.03 and 0.035 (for a moderate level of ILS). This increases the

155    probability of ILS at node $u$ and $s$ to 0.05 and 0.25, respectively. Thus, we used a total of nine

156    simulation conditions that combined three levels of evolutionary distance and three levels of ILS.

157    While our simulations do not explore extreme levels of gene tree discordance, they allow us to

158    learn about some general features of increasing ILS on network inference with different data

159    types. We simulated 1000 gene trees and their sequences, 500bp in length, under the MSC for

160    each of the nine conditions. We also explored the effects of sampling fewer genes on

161    downstream analyses. For each replicate of 1000 gene trees and their sequences, we randomly

162    sampled 400, 40, and four genes without replacement.

163       All simulations sampled haploid data. Two haploid sequences were sampled for each

164    diploid species and four haploid sequences were sampled for the allopolyploid species $E$, where

165    two sequences came from each parental lineage. We then investigated unphased data in three

7

166    ways. First, to generate unphased genotype sequences, the simulated haploid sequences for

167    each species were collapsed into a single sequence in which heterozygous sites were

168    represented by IUPAC ambiguity codes (genotype). The allopolyploid species was not restricted

169    to only biallelic sites. Second, we generated haploid consensus sequences, where for each

170    variable site, only one base was randomly retained (consensus). This could represent a case in

171    which read coverage across a locus is highly uneven such that a haploid sequence is actually a

172    chimera of two or more alleles. Finally, we simply picked one phased haploid sequence, which

173    is possible when one parental haplotype has a majority of reads for a locus (pick one). This

174    scenario where only one parent's sequence would be recovered in the offspring could be

175    anticipated in real data due to subgenome dominance (e.g., Buggs et al. 2014; Bird et al. 2018).

176    In practice, we expect most *de novo* assemblers to generate output in between the haploid

177    consensus data and pick one data.

178

179        *Inferring species networks with phased and unphased simulated data* — We estimated

180    species networks with PhyloNetworks v0.12.0 (Solis-Lemus et al. 2017) using Julia v.1.4.1

181    (Bezanson et al. 2015) from either the true gene trees used to simulate the data, or gene trees

182    estimated from the phased or unphased sequence data. For estimated trees, we used IQTREE

183    v1.6.10 (Nguyen et al. 2015) with the same model used for simulation. Each PhyloNetworks

184    analysis used the species tree (A,((B,E),(C,D))) as the starting tree and allowed zero, one, or

185    two reticulation events. Each analysis included ten independent optimizations of the

186    pseudolikelihood score. We considered larger numbers of reticulations an improvement if they

187    were two or more pseudolikelihood units lower than the best model. We compared the

188    estimated networks with one reticulation to the true network with the *hardwiredClusterDistance*

189    function (Huson et al. 2010) in PhyloNetworks. This allowed us to score the number of

190    replicates that 1) recovered the correct number of reticulations and 2) matched the true network

191    when the number of reticulations was set to one. We estimated networks for samples of four,

8

192    40, 400, and 1000 gene trees for each of the 30 replicates for each of the nine simulation

193    divergence and ILS conditions.

194

195        *Effect of phasing on divergence time estimation* — We also used our simulated phased

196    and unphased data to estimate divergence times under the MSci model (Flouri et al. 2020)

197    using BPP v4.2.9. Here, we estimate divergence times ($\tau s$), population sizes ($\theta s$), and the

198    proportion of introgressed loci ($\varphi$) on the correct fixed species network. MSCi analyses used

199    diffuse priors on $\tau_r$ and $\theta$ with a mean on their simulated values with $\varphi \sim \beta(1,1)$. Phased,

200    consensus, and pick one sequences were treated as haploids while genotype sequences were

201    treated as unphased diploids and used the analytical phasing (Gronau et al. 2011) implemented

202    in BPP. Although this is not correct for the tetraploid, it is arguably more appropriate than

203    treating all of the genotype sequences as haploid. Each Markov chain Monte Carlo (MCMC)

204    analysis collected 10,000 posterior samples, saving every 100 generations, while discarding the

205    first 100,000 generations (i.e., 10% of the total run) as burnin. All scripts for simulation and

206    subsequent analyses of simulated data are available in Dryad (X).

207

208    **A Phasing Pipeline for Polyploids**

209        *Target enrichment data* — We were motivated by the general premise of using phased

210    data to infer reticulate evolutionary histories of polyploids based on the success of empirical

211    studies where phasing was informative about hybridization or introgression events (e.g.,

212    Oberprieler et al. 2017; Eriksson et al. 2018). We were aware of few instances of phasing

213    genomic or phylogenomic data in polyploids, except in cases where chromosome-level whole-

214    genome assemblies have characterized subgenomes in allopolyploid crops (Yang et al. 2017;

215    Colle et al. 2019) or emerging results that are dependent on the sampling of parental lineages

216    (Freyman et al. 2020; Nauheimer et al. 2020). We designed PATÉ for target enrichment data

9

217    because of the availability of such data for many of taxa, but it is applicable to other types of

218    data with paired-end Illumina reads.

219         The end product of a *de novo* target enrichment assembly pipeline (such as HybPiper;

220    Johnson et al. 2016) generally is a single consensus sequence for each locus for each

221    individual. Allelic variation may be represented by ambiguous nucleotide codes within the single

222    consensus sequence or lost when the pipeline outputs the haplotype consensus sequence

223    where the majority vote from a collection of reads is used. We use these existing *de novo*

224    assembly pipelines as a starting point to provide the reference sequence for each locus for each

225    individual and leverage a recent phasing algorithm with high-quality variant calls to recover

226    phased haplotype sequences for taxa with known ploidy levels. Ploidy levels were well-

227    characterized for individuals in our *Dryopteris* analyses, but for unknown systems there are

228    existing methods to estimate ploidy directly from target enrichment data (Weiss et al. 2018;

229    Viruel et al. 2019) in the absence of other sources, such as flow cytometry (Farhat et al. 2019).

230

231         *Phasing alleles within loci* — PATÉ (Fig. 3) starts with assembled target enrichment loci,

232    such as the supercontig files output from HybPiper (Johnson et al. 2016) that contain a single

233    haplotype consensus sequence from each individual per locus. Reads for each individual are

234    then realigned to their consensus locus using BWA v0.7.17 (Li and Durbin 2009). PCR

235    duplicates are flagged with MarkDuplicates in Picard v2.9.2

236    (http://broadinstitute.github.io/picard), and variant calls are computed with HaplotypeCaller in

237    GATK v.4.1.4 (McKenna et al. 2010). We applied the following hard filters with VariantFiltration

238    in GATK: (1) QD < 2.0, (2) FS > 60.0, (3) MQ < 40.0, (4) ReadPosRankSum < −8.0, (5) AF <

239    0.05 || AF > 0.95. These loosely follow community recommendations on filters for germline

240    variant discovery (DePristo et al. 2011). Notably, we do not perform quality score recalibration

241    or filter on the mapping quality rank sum, as we anticipate allopolyploids could have a lower

242    mapping quality associated with an alternate allele due to sequence divergence or structural

10

243    variation among homeologous chromosomes. We also consider a very narrow window for

244    filtering on allele frequency. Because increasing ploidy levels will generate smaller anticipated

245    ratios of alternate to reference alleles, coupled with sequencing error and read stochasticity, we

246    only aim to remove the most extreme cases. For example, if almost all reads support the

247    alternate allele at a site, it is difficult to diagnose if the error lies in the consensus assembly or

248    the read alignment. In these cases, only the reference site is retained, and the variant does not

249    pass the allele frequency filter. However, the allele frequency filter could be removed if

250    investigators are focused on organisms with extremely high ploidy levels.

251        Biallelic SNPs that pass filters are then phased with H-PoPG v.0.2.0 (Xie et al. 2016). H-

252    PopG solves a heuristic phasing problem efficiently using dynamic programming. Although not

253    guaranteed to be an optimal solution, H-PoPG has been shown to have high accuracy while

254    also being fast (Xie et al. 2016; He et al. 2018; Moeinzadeh et al. 2020). Phasing variants in

255    polyploids is difficult because for $n$ variants and $k$ ploidy, there are $2^{n-1}(k-1)^n$ possible ways

256    to link the sites together. H-PoPG evaluates possible solutions efficiently by grouping reads into

257    $k$ groups in such a way that differences within the groups are minimized. Focusing on target

258    enrichment data also constrains the complexity of the phasing problem compared to whole-

259    genome alignments (i.e., haplotype blocks are constrained to about 1000 bp). We then used the

260    phased variants to create individual allele sequences, where invariable sites are filled in based

261    on the reference sequence. In cases where there is no linkage information to phase across the

262    entire locus, we retain the phasing only for the longest block. The variants for the shorter

263    haplotype blocks can be collapsed into IUPAC ambiguity codes or treated as missing data

264    based on the investigator's preferences. PATÉ outputs analysis-ready fasta files with multiple

265    alleles per species. Variants are only phased within loci; we do not attempt to assign loci to

266    parental subgenomes. While this may complicate analyses of concatenated multi-locus

267    datasets, it is ideal for the MSC that assumes free recombination between loci and can leverage

268    multiple alleles per species for estimating $\theta s$. Those interested in concatenated analyses can

11

269    use other recent approaches that assign gene copies to parental subgenomes (Freyman et al.

270    2020; Nauheimer et al. 2020).

271

272    ***Analyses of a Species Complex with Allopolyploidy***

273          *The North American wood fern complex (Dryopteris)* — We tested PATÉ using new

274    target enrichment data from nine North American *Dryopteris* species, including both

275    allotetraploid and allohexaploid taxa, with well-studied reticulate relationships (Sessa et al.

276    2012a, b) as well as two outgroup taxa from the sister genus *Polystichum*. All putative parental

277    lineages are represented in our dataset, with the exception of a hypothesized extinct lineage (*D.*

278    *semicristata*; Sessa et al. 2012b). We sampled two or three individuals for each *Dryopteris*

279    taxon (Table S1). The target enrichment data were generated from the GoFlag 408 flagellate

280    land plant probe set (Breinholt et al. 2021) at RAPID Genomics (Gainesville, FL). The target

281    regions for this probe set are 408 exons found in 229 single or low-copy nuclear genes. We

282    generated haplotype consensus assemblies for each with HybPiper (Johnson et al. 2016). The

283    resulting supercontig sequences became our reference sequences for genotyping and phasing

284    with PATÉ. We aligned both phased and unphased (i.e., the reference supercontig sequences)

285    with MUSCLE with default settings (Edgar 2004).

286

287          *Three species tests* — We first explored the value of phasing data when estimating

288    reticulate relationships among three species, including two diploid parental lineages (*D.*

289    *expansa* and *D. intermedia*) and their putative allotetraploid descendent (*D. campyloptera*). The

290    two diploid parents last shared a common ancestor during the Late Eocene and Early Miocene,

291    approximately 23 Ma (Sessa et al. 2012b). We used both a full-likelihood Bayesian approach

292    and a topology-based pseudolikelihood approach to estimate the correct species relationships

293    from phased and unphased data. First, using BPP v.4.1.4 (Flouri et al. 2020), we estimated log-

294    marginal likelihoods (ln *mL*) with stepping-stone sampling (Xie et al. 2011) for the three possible

295    rooted three taxon trees and twelve possible network models that imply differences for the

296    timing and direction of allopolyploidy and the presence of unsampled ancestral lineages

297    (Supplementary Fig. S1). Each ln *mL* estimate used 24 steps, and each step had a posterior

298    sample of 10,000, saving every 100 generations after a 100,000 generation burnin (10% of the

299    total run). The ln *mL* values were then used to calculate the fifteen model probabilities following

300    equation 1 (e.g., Beerli et al. 2019).

301    $$P(model) = \frac{exp(ln\ mL_{model} - ln\ mL_{max})}{\Sigma_i[exp(ln\ mL_i - ln\ mL_{max})]} \qquad Equation\ 1$$

302         We repeated analyses for 30 random subsets of 40 and then four loci to explore the

303    effects of the number of loci on inferring allopolyploidy. Next, we used PhyloNetworks to test the

304    presence and placement of gene flow between the three species. For that analysis, we also

305    included the sequences from the two *Polystichum* outgroups. IQ-TREE v1.6.10 (Nguyen et al.

306    2015) and model selection by ModelFinder (Kalyaanamoorthy et al. 2017) was used to estimate

307    gene trees for the phased and unphased data. The starting tree was obtained with ASTRAL III

308    v5.6.3 (Zhang, Rabiee et al. 2018). We tested the presence of zero, one, or two reticulations

309    with slope heuristics (Solis-Lemus and Ané 2016). Each analysis used ten independent

310    optimizations. In addition to the dataset of all loci, we analyzed the same 30 random subsets of

311    40 and four loci from the marginal likelihood analyses.

312

313         *Nine species tests* — We also investigated the differences in results from phased versus

314    consensus sequences when estimating a network for the nine-species complex, which involves

315    multiple reticulation events on an edge and thus should be difficult for network estimation (Solis-

316    Lemus and Ané 2016). Because the increased number of species and complexity of reticulation

317    in *Dryopteris,* evaluation with marginal likelihoods was not computationally feasible. Instead, we

318    performed analyses of the nine-species complex and two outgroups with phased and unphased

319    haplotype consensus data with PhyloNetworks as described above, but we allowed up to six

13

320    reticulation events. We used ASTRAL III v5.6.3 (Zhang, Rabiee et al. 2018) to generate the

321    starting species tree for network estimation using gene trees inferred from IQ-TREE v1.6.10

322    (Nguyen et al. 2015) with the best model selected by ModelFinder (Kalyaanamoorthy et al.

323    2017).

324

325    **RESULTS**

326    ***Simulation Shows Benefits of Phased Data***

327        *Network inference* — In our simulation results, both phased (i.e., the haplotypic allele

328    sequences) and unphased data (i.e., genotype, consensus, and pick one) performed well when

329    the goal is only to detect the correct number of reticulations (Supplementary Fig. S2). The only

330    case where analyses did not converge to the true number of reticulations was in the presence of

331    moderate ILS and high nucleotide divergence; however, this appears largely due to gene tree

332    estimation error, as analyses using the true simulated gene trees greatly outperformed those

333    using gene trees estimated from the simulated sequence data. However, using phased data

334    provides more accurate estimates of the placement of the reticulation edge in comparison to

335    using genotype data, and to a lesser extent, consensus sequences (Fig. 4). When the true gene

336    trees were used, which have information about the allopolyploid's hybridization event (i.e., the

337    allele sequences are sister to their respective parents in every tree), the correct network can be

338    inferred with 40 or fewer loci when nucleotide divergence is moderate. The gene trees

339    estimated from phased data perform equally well, although they require a few more gene trees

340    when nucleotide divergence is low and ILS moderate. Analyses using the genotype data almost

341    never recover the true network for these medium and low divergence scenarios, regardless of

342    the amount of ILS. Analyses based on pick one data perform similarly well to the phased and

343    true data when sampling 400 loci, but they are less accurate for four or 40 loci at low and

344    medium divergences. Analyses using the consensus data perform poorly for low numbers of loci

345    under a low divergence and no ILS scenario, but they are capable of recovering the true

14

346    reticulation when sampling 400 or more loci for the other low and medium divergence cases. In

347    the high divergence simulations, all data types could infer the true network if enough genes

348    were sampled, but analyses with the phased data required fewer genes than others.

349

350        *Divergence times* — Using phased data also improves divergence time estimates when

351    nucleotide divergence is low, but not when divergence is moderate or high (Fig. 5). When

352    divergence was low, the timing of reticulate events was accurately estimated when using

353    phased haplotypic data, but was overestimated when using genotype and especially consensus

354    data. For analyses with genotype and consensus data, as the number of loci increased and

355    uncertainty in the posterior was reduced, the posterior mean did not converge to the true

356    estimate and the simulated value was not within the highest posterior density (HPD) interval.

357    Additionally, all other nodes in the species network were overestimated with genotype or

358    consensus data, while the phased data were capable of recovering the simulated divergence

359    times (Supplementary Figs. S3-S5). Aside from some cases with the consensus sequences, all

360    four data types performed similarly with four loci; however, this is likely due to the posterior

361    being dominated by the prior in the absence of enough data. There was little improvement in

362    divergence time estimates when going from 40 to 400 loci, aside from further reduction in the

363    HPD intervals.

364        For medium sequence divergence, there was little difference between the phased and

365    unphased data. Phased sequences slightly underestimated the timing of hybridization while

366    unphased data slightly overestimated the timing of these events. However, phased data

367    accurately estimated the age of older speciation nodes that were again systematically

368    overestimated with unphased data (Supplementary Figs. S6-S8). At a high level of nucleotide

369    divergence, genotype data were capable of accurately estimating all divergence time

370    parameters while the phased data underestimated the timing of hybrid events (Fig. 5;

371    Supplementary Figs. S9-S11). Notably, the pick one data performed very well for all divergence

15

372    time estimation scenarios. Divergence times were not strongly affected by increasing levels of

373    ILS, likely because estimates were performed with the MSci model and we did not explore very

374    high ILS scenarios, but age estimates improved slightly for the genotypes, consensus, and pick

375    one data with increasing ILS.

376

377    ***Analyses of target enrichment data from Dryopteris***

378        *Recovery of Phased Loci* — On average, 62% of loci sequenced for an individual were

379    phased with eight variants passing filters (Table S2). The ploidy level appears to be strongly

380    associated with the number of phased loci. Among diploids, only 30% of loci were phased; the

381    other 70% of diploid loci were either homozygous or had too few linked variants for phasing. For

382    tetraploids and hexaploids, 87% and 94% of loci, respectively, were phased such that two or

383    more phased haplotype sequences could be recovered. Among loci where phasing was

384    possible, variants were almost always resolved as a single haplotype block, as opposed to

385    being split into two or more blocks because not enough reads were available to physically link

386    variants. For polyploids, the number of phased haplotype sequences most frequently matched

387    the ploidy level except in the case of a single *D. campyloptera* individual (B087-D08), which also

388    had few recovered loci. Phasing data only extended sequence alignment length by about four

389    base pairs, but it more than doubled the number of parsimony informative sites (Table S3).

390

391        *Placing a single reticulation event* — For our three-species full-likelihood analyses with

392    the MSci, both phased and unphased data recovered the anticipated reticulation hypothesis,

393    identifying *D. campyloptera* as an allotetraploid with the two diploid parental lineages *D.*

394    *expansa* and *D. intermedia*, when using all loci (Table 1). Model probabilities show decisive

395    support for a scenario where there are two unsampled ancestral populations that were the

396    progenitors of *D. campyloptera*, as opposed to *D. campyloptera* being a hybrid species with

397    extant *D. expansa* and *D. intermedia* as parents. All model parameters converged for both

16

398     phased (Supplementary Figs. S12 and S13) and unphased data (Supplementary Figs. S14 and

399     S15). Divergence time estimates were older in the analysis of phased data, although the relative

400     order of divergence events was the same using phased and unphased data (Fig. 6). There was

401     more uncertainty in the $\theta$ estimates, especially for the alloplyploid species *D. campyloptera* and

402     the two ancestral populations of the parental lineages that formed the polyploid ancestor

403     (Supplementary Figure S16). Repeating the analyses with fewer loci did not always produce the

404     same result, but there was either decisive support for a model or multiple plausible models that

405     all had the correct species relationships and direction of introgression for both phased and

406     unphased data when using 40 loci (Fig. 7). The only difference between models was the

407     presence or absence of ancestral $\theta s$ for unsampled lineages. Analyses with four loci produced

408     less reliable results for both phased and unphased data. The four-locus analyses of phased and

409     unphased data found some non-negligible model probability for trees without hybridization or

410     networks where hybridization was incorrect in nine and twelve out of 30 replicates, respectively

411     (Fig. 7).

412         When performing a network search based on gene tree distributions, both the phased

413     and unphased data were able to recover the hypothesized allopolyploidy event (Fig. 8;

414     Supplementary Table S5). Both analyses inferred the major branch to be *D. intermedia* with the

415     minor branch from *D. expansa*. Phased data estimated a slightly smaller inheritance probability

416     compared to the unphased data. These findings from PhyloNetworks are consistent with

417     parameter estimation under the MSci model, where phased data estimated a slightly smaller

418     mean $\varphi_h$ compared to unphased data (Table S4). When sampling 100 replicates of 40 loci, 98%

419     of replicates for phased data and 100% for unphased data were capable of detecting a single

420     hybrid event in the data. When sampling four loci, this drops to 66% and 85%, respectively (Fig.

421     8). Phased data more frequently recovered the correct network (58%) compared to unphased

422     data (38%) with 40 loci; however, the converse is true for four loci, with 34% correct for phased

423     and 41% correct for unphased (Fig. 8). Unphased data also got the network wrong more

17

424    frequently than phased data, such that phased data had the direction of introgression incorrect

425    in 36% and 12% of replicates while unphased data were incorrect in 60% and 38% replicates for

426    40 and four loci, respectively (Fig. 8).

427

428        *Inferring relationships among a complex with multiple reticulation events* — The network

429    recovered for phased data identified three out of five hypothesized reticulation events among

430    the nine *Dryopteris* species (Fig. 9; Supplementary Table S5). The allotetraploid *D. celsa* was

431    correctly identified with diploid *D. ludoviciana* and *D. goldiana* as parents. Analysis of the

432    phased data detected a low level of gene flow from the common ancestor of this clade into

433    tetraploid *D. cristata,* which has *D. ludoviciana* as one hypothesized parent while the other

434    parent lineage (*D. semicristata*) is assumed to be extinct (Sessa et al. 2012b). *Dryopteris*

435    *carthusiana* is another tetraploid that is assumed to share the extinct common ancestor with *D.*

436    *cristata*, but has experienced introgression from *D. intermedia*, with a high inheritance

437    probability of 0.43 (Fig. 9). However, the phased data missed the putative allotetraploid case of

438    *D. campyloptera*, despite the strong evidence for this reticulation event in our earlier three-taxon

439    analyses. Our network with phased data also failed to identify the putative reticulate

440    evolutionary history of *D. clintoniana*, an allohexaploid where allotetraploid *D. cristata* and

441    diploid *D. goldiana* are assumed to be the parents (Sessa et al. 2012b). The spectra of quartet

442    concordance factors were overall similar between the phased and unphased data

443    (Supplementary Fig. S17), but the phased data were arguably more accurate.

444        Although the phased data were not successful in recovering all hypothesized reticulate

445    relationships, they performed better than the unphased data. Analyses of unphased data were

446    capable of finding the allotetraploid history of *D. celsa* with an inheritance probability similar to

447    the phased data (Fig. 9). Hybridization between *D. intermedia* and *D. carthusiana* was also

448    detected; however, the directionality was reversed, with gene flow going from the allotetraploid

449    into the diploid. A third reticulation edge was found in the unphased analysis, from the common

18

450    ancestor of *Dryopteris* into the common ancestor of *D. clintoniana* and its sister clade. This

451    hybrid edge is difficult to reconcile because of the hypothesized extinct common ancestor that

452    contributed to both *D. cristata* and *D. carthusiana*. *Dryopteris clintoniana* is correctly placed in

453    the major species tree topology, as a grade between *D. cristata* and *D. goldiana*. This

454    reticulation edge from the *Dryopteris* common ancestor may reflect the extinct lineage and a

455    high degree of gene tree variation.

456

457    **DISCUSSION**

458        New phylogenetic network methods offer the promise of elucidating the often complex

459    reticulate histories of polyploid lineages, even in the presence of ILS. Our results demonstrate

460    that phasing polyploid target enrichment data can improve the accuracy of such network

461    inferences as well as divergence time estimates for the networks, and we describe a novel

462    pipeline (PATÉ) to address the difficult problem of phasing polyploid data. Although PATÉ could

463    handle different types of genomic data, such as transcriptomes and whole genomes, target

464    enrichment data are ideal for investigation because they often yield high and even coverage

465    across loci. Because MSC methods assume treat loci independently and assume free

466    recombination between loci, it is not necessary to assign individual loci to parental subgenomes.

467    However, the allele sequences output by PATÉ can also be used as input for the recently

468    developed *Homologizer* (Freyman et al. 2020) or *HybPhaser* (Nauheimer et al. 2020), which

469    attempt to phase across loci and recover parental subgenomes. Our methods also enable

470    population genomic studies of polyploids where accurate estimation of site frequency spectra

471    can be used for demographic analyses otherwise complicated by polyploidy (e.g., Excoffier et

472    al. 2013; Liu and Fu 2020) or SNP-based network inference in the absence of variation suitable

473    for gene tree estimation (Blischak et al 2018; Olave and Meyer 2020).

474

475    ***Promises and Pitfalls of Phasing***

476         The prospect of using alleles from phased genomic data presents an exciting step

477    towards revealing the evolutionary history of polyploids, which remains a critical impediment

478    within the plant evolution community (McKain et al. 2018). Strategies for explicitly addressing

479    this challenge are only now emerging (Freyman et al. 2020; Nauheimer et al. 2020), and PATÉ

480    can be a useful tool by phasing variants for many individuals while leveraging genotyping

481    information. Our simulations demonstrate that phasing can improve estimates of reticulate

482    evolutionary relationships using network methods. Phased data can more accurately recover

483    the placement and directionality of hybrid edges than various types of unphased data in

484    simulations (Fig. 4) and empirical analyses with limited numbers of loci (Fig. 7). The advantages

485    of using phased versus unphased data for network estimation decrease when a large number of

486    loci are sampled (Table 1; Fig. 8).

487         Perhaps an underappreciated aspect of phased data is their ability to improve

488    divergence times estimates (Fig. 5; Supplementary Figs. S3-S11; Anderman et al. 2019). Our

489    empirical analyses also demonstrated how the timing of introgression $\tau_h$ can be greatly affected

490    by whether phased or consensus data are used. In our *Dryopteris* analyses, the estimate from

491    phased data was nearly four times older than the estimate from unphased data (Supplementary

492    Table S4). The direction of this difference was unanticipated, because our simulations

493    suggested the consensus data should overestimate age compared to phased data. This

494    highlights the difficulties of simulating data that capture real complexities and makes deciding

495    which estimate is more reliable somewhat difficult; however, the uncertainty of $\tau_h$ for haplotype

496    consensus data reflected in its posteriors (Supplementary Fig. S14) gives us more confidence in

497    the phased estimates (Supplementary Fig. S12). Because we used the MSci model for

498    divergence time estimation, we did not observe any effect of ILS on age estimates in our

499    simulations, but if we were using concatenation methods to date the divergence times for the

20

500    allopolyploids two subgenomes, nodes affected by ILS should be overestimated (Tiley et al.

501    2020).

502          In most cases phasing appears to be beneficial, but it may be problematic when the

503    parental lineages are deeply diverged. Although the phased data were able to accurately

504    estimate the age of older speciation nodes (Supplementary Figs. S3-S11), as phylogenetic

505    information is lost from multiple hits, the influence of the prior becomes more substantial. These

506    deep divergence simulation scenarios may border on being biologically unrealistic, as identifying

507    an allopolyploid and its two parental lineages becomes more difficult over time due to extinction

508    and population genetic processes, but it provides some expectations for the performance of

509    phased and unphased data in the presence of high nucleotide divergence. Our simulations

510    showed that when alleles are phased but only one is sampled, as in our pick one simulations,

511    network and divergence time estimation is similar to having all phased alleles present. We also

512    showed where consensus data can perform poorly through simulation (Figs. 3 and 4) as well as

513    empirical analyses where the direction of introgression was more frequently reversed in a

514    simple example of two parental lineages and an allopolyploid (Fig. 8). When enough reads are

515    available to call high-confidence variants, we suggest that phasing with PATÉ can improve

516    network and divergence time estimation for species complexes with low to moderate sequence

517    divergence. However, when investigating very ancient hybrid events, unphased genotype data

518    may be preferential, and using analytical approaches that integrate over phases (Gronau et al.

519    2011; Flouri et al. 2020) may outperform analyses with phased sequences because allele-

520    specific information no longer captures shared ancestry with parental lineages and haplotype

521    blocks become smaller due to recombination.

522

523    ***Challenges of Network Estimation***

524          Our analyses highlight the difficulty of estimating the evolutionary histories of reticulate

525    complexes using phylogenetic methods, regardless of data type. The full-likelihood

21

526    implementation of the MSci model appears to be useful for limited cases, but these methods are

527    computationally demanding. Thus, they may not be practical for generating hypotheses and

528    exploring unknown relationships for large numbers of taxa (e.g., Zhang et al. 2018). Quartet-

529    based methods are fast and accurate when there are limited numbers of reticulations (Solis-

530    Lemus and Ané 2016) and show similar accuracy to full-likelihood methods for estimating

531    introgression probabilities (Flouri et al. 2020). However, there are scenarios where true network

532    topologies become non-identifiable (Solis-Lemus and Ané 2016). For example, when multiple

533    introgression events affect the same lineage, the expected quartet distribution under the MSci

534    model becomes a poor fit for the empirical data (Cai and Ané 2020) and the network estimated

535    may be incorrect. These effects were evident in our empirical analyses where the relationship

536    between *D. campyloptera*, *D. expansa*, and *D. intermedia* was missing in the nine-species

537    analysis (Fig. 9), which we expect is due to a reticulation edge present between *D. intermedia*

538    and *D. carthusiana*. Phasing sequence data adds information that can improve estimates

539    (Supplementary Table S3), but unsampled or extinct lineages, such as the hypothesized *D.*

540    *semicristata,* can create significant barriers to recovering the true evolutionary history of

541    reticulate complexes, regardless of how many loci or individuals are available.

542

543    ***Insights into Dryopteris Evolution***

544        The North American *Dryopteris* complex has been well-characterized through the study

545    of multi-locus nuclear and chloroplast phylogenies, morphology, cytological observations,

546    chromatography, and isozyme analyses (reviewed in Sessa et al. 2012b). This makes it a useful

547    system for testing our phasing pipeline, and our analyses add nuance to our understanding of

548    some of the relationships among *Dryopteris* species. For example, our results indicate *D.*

549    *campyloptera* has received slightly more loci from *D. intermedia* than *D. expansa* (Fig. 8), and

550    *D. intermedia* is thought to be the likely maternal progenitor (Sessa et al. 2012b). Although our

551    analysis of marginal likelihoods for all target enrichment loci suggested the presence of

22

552     unsampled ancestral populations (Table 1), the age of introgression and divergence of the *D.*

553     *campyloptera* ancestral population from the *D. intermedia* and *D. expansa* parental lineages is

554     consistent with hybrid speciation rather than a lineage that was isolated from *D. intermedia* and

555     received more recent gene flow from *D. expansa* (Fig. 6). Following the allopolyploidy event, the

556     *D. intermedia* genome was likely dominant, providing some selective advantage for *D.*

557     *campyloptera* in its distribution at the time (Bird et al. 2018). Similar insights can be gained from

558     the nine-taxon analyses, which suggests *D. ludoviciana* is the dominant genome in *D. celsa*.

559     *Dryopteris ludoviciana* is the maternal parent of *D. celsa*, again suggesting some bias in

560     retaining homeologous alleles from the maternal lineage, which provides the chloroplast

561     genome in ferns (Sessa et al. 2012b).

562             The nine-taxon analyses also support the hypothesis of the unsampled diploid lineage *D.*

563     *semicristata*, based on the placement of *D. carthusiana* as sister to the rest of *Dryopteris* in the

564     phased nine-taxon analyses, but with *D. carthusiana* having received over 40% of its genes

565     from *D. intermedia* more recently. Our analyses suggest that *D. cristata* did not have *D.*

566     *ludoviciana* as a progenitor, but rather an unsampled common ancestor of *D. goldiana* and *D.*

567     *ludoviciana*. In the case of *D. cristata*, both parental diploid lineages, including *D. semicristata,*

568     may have gone extinct.

569

570     ***Conclusions***

571             Combining phased data with recent network methods holds much promise for

572     confronting a major challenge of plant phylogenetics: resolving the complex histories of

573     polyploids. The PATÉ pipeline can enhance systematic, speciation genomic, and population

574     genomic analyses of groups containing polyploids. While haplotype consensus sequences may

575     be adequate for resolving single reticulation events where both parents are sampled, using

576     phased sequences can improve inferences of more complicated allopolyploid events,

577     demonstrating how allelic variation can be leveraged for MSC methods that account for

23

578     reticulation. Still, some reticulate complexes can be difficult to disentangle with any data when

579     there are multiple hybrid events involving the same branch. PATÉ is available through GitHub

580     (https://github.com/gtiley/Phasing) and can be run on any UNIX environment after installing

581     basic genotyping software and H-PoPG.

582

583     **DATA AVAILABILITY**

584          PATÉ is freely available through GitHub at https://github.com/gtiley/Phasing. Simulated

585     and empirical data supporting findings and files for replicating analyses are available from the

586     Dryad Digital Repository: (X). Raw Fastq reads for *Dryopteris* individuals are available through

587     the NCBI SRA database and are associated with BioProject PRJNA725004. Individual SRA

588     Identifiers are available in Supplementary Table S1.

589

590     **SUPPLEMENTARY MATERIAL**

591          Data available from the Dryad Digital Repository: (X).

592

593     **ACKNOWLEDGEMENTS**

594          The authors thank A.M. Duffy, K. Imwattana, M. Nieto-Lugilde, B.T. Piatkowski, and A.J.

595     Shaw for helpful discussions and providing feedback on the manuscript. We also thank M.G

596     Johnson, J. Mendez Reneau, and L. Nauheimer for discussions and sharing their strategies for

597     phasing sequence data. We are grateful to the people who made data collection possible; we

598     used *Dryopteris* samples collected by C.J. Rothfels, M.A. Sundue, and W. Testo, and DNA

599     extractions were performed by S.B. Carey and E. Lockwood.

600

601     **FUNDING**

602          Funding was provided by National Science Foundation awards DEB-2038213 to AAC

603     and PSM, and from DEB-1541506 to JGB and EBS. This work was also partially supported by

24

606

607

**LITERATURE CITED**

608

609 Andermann T., Cano A., Zizka A., Bacon C., Antonelli A. 2018. SECAPR-a bioinformatics
610 pipeline for the rapid and user-friendly processing of targeted enriched Illumina sequences, from
611 raw reads to alignments. *PeerJ* 6:e5175.
612

613 Andermann T., Fernandes A.M., Olsson U., Topel M., Pfeil B., Oxelman B., Aleixo A., Faircloth
614 B.C., Antonelli A. 2019. Allele Phasing Greatly Improves the Phylogenetic Utility of
615 Ultraconserved Elements. *Syst. Biol.* 68:32-46.
616

617 Baaijens J.A., Schonhuth A. 2019. Overlap graph-based generation of haplotigs for diploids and
618 polyploids. *Bioinformatics* 35:4281-4289.
619

620 Bankevich A., Nurk S., Antipov D., Gurevich A.A., Dvorkin M., Kulikov A.S., Lesin V.M.,
621 Nikolenko S.I., Pham S., Prjibelski A.D., Pyshkin A.V., Sirotkin A.V., Vyahhi N., Tesler G.,
622 Alekseyev M.A., Pevzner P.A. 2012. SPAdes: a new genome assembly algorithm and its
623 applications to single-cell sequencing. *J. Comput. Biol.* 19:455-477.
624

625 Barker M.S., Arrigo N., Baniaga A.E., Li Z., Levin D.A. 2016. On the relative abundance of
626 autopolyploids and allopolyploids. *New Phytol.* 210:391-398.
627

628 Beerli P., Mashayekhi S., Sadeghi M., Khodaei M., Shaw K. 2019. Population Genetic Inference
629 With MIGRATE. *Curr. Protoc. Bioinformatics* 68:e87.
630

631 Berger E., Yorukoglu D., Peng J., Berger B. 2014. HapTree: a novel Bayesian framework for
632 single individual polyplotyping using NGS data. *PLoS Comput. Biol.* 10:e1003502.
633

634 Bezanson J., Edelman A., Karpinski S., Shah V.B. 2015. Julia: A Fresh Approach to Numerical
635 Computing. *arXiv*1411.1607v4.
636

637 Bird K.A., VanBuren R., Puzey J.R., Edger P.P. 2018. The causes and consequences of
638 subgenome dominance in hybrids and recent polyploids. *New Phytol.* 220:87-93.
639

640 Blischak P.D., Chifman J., Wolfe A.D., Kubatko, L.S. 2018. HyDe: a Python package for
641 genome-scale hybridization detection. *Syst. Biol.*, 67:821-829.
642

643 Breinholt J.W., Carey S.B., Tiley G.P., Davis E.C., Endara L., McDaniel S.F., Neves L.G., Sessa
644 E.B., von Konrat M., Chantanaorrapint S., Fawcett S., Ickert-Bond S.M., Labiak P.H., Larraín J.,
645 Lehnert M., Lewis L.R., Nagalingum N.S., Patel N., Rensing S.A., Testo W., Vasco A., Villarreal
646 J.C., Williams E.W., Burleigh J.G. 2021. A target enrichment probe set for resolving the
647 flagellate plant tree of life. *Appl. Plant Sci.* 9:e11406.
648

649 Breinholt J.W., Earl C., Lemmon A.R., Lemmon E.M., Xiao L., Kawahara A.Y. 2018. Resolving
650 relationships among the megadiverse butterflies and moths with a novel pipeline for anchored
651 phylogenomics. *Syst. Biol.* 67:78-93.
652

653 Buggs J.A., Wendel J.F., Doyle J.J., Soltis D.E., Soltis P.S., Coate J.E. 2014. The legacy of
654 diploid progenitors in allopolyploid gene expression patterns. *Philos. Trans. R. Soc. Lond. B*
655 *Biol. Sci.* 368:20130354.
656

26

657    Cai R., Ané C. 2020. Assessing the fit of the multi-species network coalescent to multi-locus
658    data. *Bioinformatics* btaa863. doi: https://doi.org/10.1093/bioinformatics/btaa863
659

660    Colle M., Leisner C.P., Wai C.M., Ou S., Bird K.A., Wang J., Wisecaver J.H., Yocca A.E., Alger
661    E.I., Tang H., Xiong Z., Callow P., Ben-Zvi G., Brodt A., Baruch K., Swale T., Shiue L., Song G.-
662    Q., Childs K.L., Schilmiller A., Vorsa N., Buell C.R., VanBuren R., Jiang N., Edger P.P. 2019.
663    Haplotype-phased genome and evolution of phytonutrient pathways of tetraploid blueberry.
664    *GigaScience* 8:giz012.
665

666    Crowl A.A., Myers C., Cellinese N. 2017. Embracing discordance: Phylogenomic analyses
667    provide evidence for allopolyploidy leading to cryptic diversity in a Mediterranean Campanula
668    (Campanulaceae) clade. *Evolution* 71:913-922.
669

670    DePristo M., Banks E., Poplin R., Garimella K., Maguire J., Hartl C., Philippakis A., del Angel G.,
671    Rivas M.A., Hanna M., McKenna A., Fennell T., Kernytsky A., Sivachenko A., Cibulskis K.,
672    Gabriel S., Altschuler S., Daly M. 2011. A framework for variation discovery and genotyping
673    using next-generation DNA sequencing data. *Nat. Genet.* 43:491-498.
674

675    Edgar R.C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
676    throughput. *Nucleic Acids Res.* 32:1792-1797.
677

678    Eriksson J.S., de Sousa F., Bertrand Y.J.K., Antonelli A., Oxelman B., Pfeil B.E. 2018. Allele
679    phasing is critical to revealing a shared allopolyploid origin of Medicago arborea and M.
680    strasseri (Fabaceae). *BMC Evol. Biol.* 18:9.
681

682    Excoffier L., Dupanloup I., Huerta-Sanchez E., Sousa V.C., Foll M. 2013. Robust demographic
683    inference from genomic and SNP data. *PLoS Genet.* 9:e1003905.
684

685    Faircloth B.C. 2016. PHYLUCE is a software package for the analysis of conserved genomic
686    loci. *Bioinformatics* 32:786-788.
687

688    Faircloth B.C., McCormack J.E., Crawford N.G., Harvey M.G., Brumfield R.T., Glenn T.C. 2012.
689    Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary
690    timescales. *Syst. Biol.* 61:717-726.
691

692    Farhat P., Hidalgo O., Robert T., Siljak-Yakovlev S., Leitch I.J., Adams R.P., Bou Dagher-
693    Kharrat M. 2019. Polyploidy in the Conifer Genus *Juniperus*: An Unexpectedly High Rate. *Front.*
694    *Plant Sci.* 10:676.
695

696    Flouri T., Jiao X., Rannala B., Yang Z. 2020. A Bayesian Implementation of the Multispecies
697    Coalescent Model with Introgression for Phylogenomic Analysis. *Mol. Biol. Evol.* 37:1211-1223.
698

699    Flouri T., Jiao X., Rannala B., Yang Z.. 2018. Species Tree Inference with BPP Using Genomic
700    Sequences and the Multispecies Coalescent. *Mol. Biol. Evol.* 35:2585-2593.
701

702    Freyman W.A,. Johnson M.G., Rothfels C.J. 2020. Homologizer: Phylogenetic phasing of gene
703    copies into polyploid subgenomes. *bioRxiv* doi: https://doi.org/10.1101/2020.10.22.351486.
704

705    Gronau I., Hubisz M.J., Gulko B., Danko C.G., Siepel A. 2011. Bayesian inference of ancient
706    human demography from individual genome sequences. *Nat. Genet.* 43:1031-1034.
707

27

708    He D., Saha S., Finkers R., Parida L. 2018. Efficient algorithms for polyploid haplotype phasing.
709    *BMC Genomics* 19:110.
710
711    Huang J., Flouri T., Yang Z. 2020. A simulation study to examine the information content in
712    phylogenomic datasets under the multispecies coalescent model. *Mol. Biol. Evol.* 37:3211-3224.
713
714    Huber K.T., Oxelman B., Lott M., Moulton V. 2006. Reconstructing the evolutionary history of
715    polyploids from multi-labelled trees. *Mol. Biol. Evol.* 23:1784-1791.
716
717    Hudson R.R. 1983. Testing the Constant-Rate Neutral Allele Model with Protein Sequence
718    Data. *Evolution* 37:203-217.
719
720    Huson D.H., Rupp R., Scornavacca C. 2010. Phylogenetic networks: concepts, algorithms and
721    applications. Cambridge University Press.
722
723    Iqbal Z., Caccamo M., Turner I., Flicek P., McVean G. 2012. De novo assembly and genotyping
724    of variants using colored de Bruijn graphs. *Nat. Genet.* 44:226-232.
725
726    Jantzen J.R., Amarasinghe P., Folk R.A., Reginato M., Michelangeli F.A., Soltis D.E., Cellinese
727    N., Soltis P.S. 2020. A two-tier bioinformatic pipeline to develop probes for target capture of
728    nuclear loci with applications in Melastomataceae. *Appl. Plant Sci.* 8:e11345.
729
730    Johnson M.G., Gardner E.M., Liu Y., Medina R., Goffinet B., Shaw A.J., Zerega N.J., Wickett
731    N.J. 2016. HybPiper: Extracting coding sequence and introns for phylogenetics from high-
732    throughput sequencing reads using target enrichment. *Appl. Plant Sci.* 4:1600016.
733
734    Johnson M.G., Pokorny L., Dodsworth S., Botigue L.R., Cowan R.S., Devault A., Eiserhardt
735    W.L., Epitawalage N., Forest F., Kim J.T., Leebens-Mack J.H., Leitch I.J., Maurin O., Soltis
736    D.E., Soltis P.E., Wong G.K.-S., Baker W.J., Wickett N.J. 2019. A Universal Probe Set for
737    Targeted Sequencing of 353 Nuclear Genes from Any Flowering Plant Designed Using k-
738    Medoids Clustering. *Syst. Biol.* 68:594-606.
739
740    Jones G., Sagitov S., Oxelman B. 2013. Statistical inference of allopolyploid species networks in
741    the presence of incomplete lineage sorting. *Syst. Biol.* 62:467-478.
742
743    Jukes T.H., Cantor C.R. 1969. Evolution of protein molecules. In: Munro H.N., editor.
744    Mammalian Protein Metabolism. New York, NY: Acedemic Press. p. 21-132.
745
746    Kalyaanamoorthy S., Minh B.Q., Wong T.K.F., von Haeseler A., Jermiin L.S. 2017.
747    ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* 14:587-
748    589.
749
750    Karimi N., Grover C.E., Gallagher J.P., Wendel J.F., Ané C., Baum D.A. 2020. Reticulate
751    Evolution Helps Explain Apparent Homoplasy in Floral Biology and Pollination in Baobabs
752    (Adansonia; Bombacoideae; Malvaceae). *Syst. Biol.* 69:462-478.
753
754    Kates H.R., Johnson M.G., Gardner E.M., Zerega N.J.C., Wickett N.J. 2018. Allele phasing has
755    minimal impact on phylogenetic reconstruction from targeted nuclear gene sequences in a case
756    study of Artocarpus. *Am. J. Bot.* 105:404-416.
757

758 Lemmon A.R., Emme S.A., Lemmon E.M. 2012. Anchored hybrid enrichment for massively
759 high-throughput phylogenomics. *Syst. Biol.* 61:727-744.
760
761 Li H., Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
762 *Bioinformatics* 25:1754-1760.
763
764 Liu X., Fu Y.X. 2020. Stairway Plot 2: demographic history inference with folded SNP frequency
765 spectra. *Genome Biol.* 21:280.
766
767 Liu Y., Johnson M.G., Cox C.J., Medina R., Devos N., Vanderpoorten A., Hedenas L., Bell N.E.,
768 Shevock J.R., Aguero B., Quandt D., Wickett N.J., Shaw A.J., Goffinet B. 2019. Resolution of
769 the ordinal phylogeny of mosses using targeted exons from organellar and nuclear genomes.
770 *Nat. Commun.* 10:1485.
771
772 Lott M., Spillner A., Huber K.T., Moulton V. 2009. PADRE: a package for analyzing and
773 displaying reticulate evolution. *Bioinformatics* 25:1199-1200.
774
775 Luo R., Liu B., Xie Y., Li Z., Huang W., Yuan J., He G., Chen Y., Pan Q., Liu Y., Tang J., Wu G.,
776 Zhang H., Shi Y., Liu Y., Yu C., Wang B., Lu Y., Han C., Cheung D.W., Yiu S.-M., Peng S.,
777 Xiaoqian Z., Liu G., Liao X., Li Y., Yang H., Wang J., Lam T-W., Wang J. 2012. SOAPdenovo2:
778 an empirically improved memory-efficient short-read de novo assembler. *GigaScience* 1:18.
779
780 McKain M.R., Johnson M.G., Uribe-Convers S., Eaton D., Yang Y. 2018. Practical
781 considerations for plant phylogenomics. *Appl. Plant Sci.* 6:e1038.
782
783 McKenna A., Hanna M., Banks E., Sivachenko A., Cibulskis K., Kernytsky A., Garimella K.,
784 Altshuler D., Gabriel S., Daly M., DePristo M.A. 2010. The Genome Analysis Toolkit: a
785 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.*
786 20:1297-1303.
787
788 Moeinzadeh M.H., Yang J., Muzychenko E., Gallone G., Heller D., Reinert K., Haas S., Vingron
789 M. 2020. Ranbow: A fast and accurate method for polyploid haplotype reconstruction. *PLoS*
790 *Comput. Biol.* 16:e1007843.
791
792 Montgomery J.D., Paulton E.M. 1981. *Dryopteris* in North America. *Fiddlehead Forum* 8:25-31.
793
794 Morales-Briones D.F., Liston A., Tank D.C. 2018. Phylogenomic analyses reveal a deep history
795 of hybridization and polyploidy in the Neotropical genus Lachemilla (Rosaceae). *New Phytol.*
796 218:1668-1684.
797
798 Nauheimer L., Weigner N., Joyce E., Crayn D., Clarke C., Nargar K. 2020. HybPhaser: a
799 workflow for the detection and phasing of hybrids in target capture datasets. bioRxiv. doi:
800 https://doi.org/10.1101/2020.10.27.354589.
801
802 Nguyen L.-T., Schmidt H.A., von Haeseler A., Minh B.Q. 2015. IQ-TREE: a fast and effective
803 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32:268-
804 274.
805
806 Oberprieler C., Wagner F., Tomasello S., Konowalik K. 2017. A permutation approach for
807 inferring species networks from gene trees in polyploid complexes by minimising deep
808 coalescences. *Methods in Ecology and Evolution* 8:835-849.

809
810  Olave M., Meyer A. 2020. Implementing Large Genomic Single Nucleotide Polymorphism Data
811  Sets in Phylogenetic Network Reconstructions: A Case Study of Particularly Rapid Radiations of
812  Cichlid Fish. *Syst. Biol.* 69:848-862.
813
814  Pamilo P., Nei M. 1988. Relationships between gene trees and species trees. *Mol. Biol. Evol.*
815  5:568-583.
816
817  Rothfels C.J., Pryer K.M., Li F.-W. 2017. Next-generation polyploid phylogenetics: rapid
818  resolution of hybrid polyploid complexes using PacBio single-molecule sequencing. *New Phytol.*
819  213:413-429.
820
821  Saada O.A., Tsouris A., Freidrich A., Schachrer J. 2020. nPhase: An accurate and contiguous
822  phasing method for polyploids. *bioRxiv* doi:https://doi.org/10.1101/2020.07.24.219105.
823
824  Sessa E.B., Zimmer E.A., Givnish T.J. 2012a. Reticulate evolution on a global scale: a nuclear
825  phylogeny for New World Dryopteris (Dryopteridaceae). *Mol. Phylogenet. Evol.* 64:563-581.
826
827  Sessa E.B., Zimmer E.A., Givnish T.J. 2012b. Unraveling reticulate evolution in North American
828  Dryopteris (Dryopteridaceae). *BMC. Evol. Biol.* 12:104.
829
830  Solis-Lemus C., Ané C. 2016. Inferring Phylogenetic Networks with Maximum Pseudolikelihood
831  under Incomplete Lineage Sorting. *PLoS Genet.* 12:e1005896.
832
833  Solis-Lemus C., Bastide P., Ané C. 2017. PhyloNetworks: A Package for Phylogenetic
834  Networks. *Mol. Biol. Evol.* 34:3292-3298.
835
836  Soltis, D.E., Visger C.J., Soltis P.S. 2014. The polyploidy revolution then...and now: Stebbins
837  revisited. *American Journal of Botany* 101: 1057–1078.
838
839  Stull G.W., Soltis P.S., Soltis D.E., Gitzendanner M.A., Smith S.A. 2020. Nuclear phylogenomic
840  analyses of asterids conflict with plastome trees and support novel relationships among major
841  lineages. *Am. J. Bot.* 107:790-805.
842
843  Tiley G.P., Poelstra J.W., dos Reis M., Yang Z., Yoder A.D. 2020. Molecular Clocks without
844  Rocks: New Solutions for Old Problems. *Trends Genet.* 36:845-856.
845
846  Viruel J., Conejero M., Hidalgo O., Pokorny L., Powell R.F., Forest F., Kantar M.B., Soto Gomez
847  M., Graham S.W., Gravendeel B., Wilkin P., Leitch I.J. 2019. A Target Capture-Based Method
848  to Estimate Ploidy from Herbarium Specimens. *Front. Plant Sci.* 10:937.
849
850  Weiss C.L., Pais M., Cano L.M., Kamoun S., Burbano H.A. 2018. nQuire: a statistical framework
851  for ploidy estimation using next generation sequencing. *BMC Bioinformatics* 19:122.
852
853  Wen D., Yu Y., Nakhleh L. 2016. Bayesian Inference of Reticulate Phylogenies under the
854  Multispecies Network Coalescent. *PLoS Genet.* 12:e1006006.
855
856  Wen D., Yu Y., Zhu J., Nakhleh L. 2018. Inferring Phylogenetic Networks Using PhyloNet. *Syst.*
857  *Biol.* 67:735-740.
858

859    Wolf P.G., Robison T.A., Johnson M.G., Sundue M.A., Testo W.L., Rothfels C.J. 2018. Target
860    sequence capture of nuclear-encoded genes for phylogenetic analysis in ferns. *Appl. Plant Sci.*
861    6:e01148.
862
863    Wood T.E., Takebayashi N., Barker M.S., Mayrose I., Greenspoon P.B., Rieseberg L.H. 2009.
864    The frequency of polyplid speciation in vascular plants. *Proc. Natl. Acad. Sci. U. S. A.*
865    106:13875-13879.
866
867    Xie M., Wu Q., Wang J., Jiang T. 2016. H-PoP and H-PoPG: heuristic partitioning algorithms for
868    single individual haplotyping of polyploids. *Bioinformatics* 32:3735-3744.
869
870    Xie W., Lewis P.O., Fan Y., Kuo L., Chen M.H. 2011. Improving marginal likelihood estimation
871    for Bayesian phylogenetic model selection. *Syst. Biol.* 60:150-160.
872
873    Yang J., Moeinzadeh M-H., Kuhl H., Helmuth J., Xiao P., Haas S., Liu G., Zheng J., Sun Z., Fan
874    W., Deng G., Wang H., Hu F., Zhao S., Fernie A.R., Boerno S., Timmermann B., Zhang P.,
875    Vingron M. 2017. Haplotype-resolved sweet potato genome traces back its hexaploidization
876    history. *Nat Plants* 3:696-703.
877
878    Yang Z. 2006. Computational molecular evolution: Oxford University Press.
879
880    Zhang C., Ogilvie H.A., Drummond A.J., Stadler T. 2018. Bayesian Inference of Species
881    Networks from Multilocus Sequence Data. *Mol. Biol. Evol.* 35:504-517.
882
883    Zhang C., Rabiee M., Sayyari E., Mirarab S. 2018. ASTRAL-III: polynomial time species tree
884    reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19:153.
885
886    Zhu S., Degnan J.H. 2017. Displayed trees do not determine distinguishability under the
887    network multispecies coalescent. *Syst. Biol.* 66:283-298.
888

889  **TABLES**

| Topology | Phased Marginal lnL | Unphased Marginal lnL | Phased Model Probability | Unphased Model Probability |
|---|---|---|---|---|
| 1 | -613405.3154 | -552720.6776 | 1.006E-278 | 3.2574E-214 |
| 2 | -613223.0755 | -552534.6540 | 1.4073E-199 | 2.0039E-133 |
| 3 | -613366.4128 | -552644.7603 | 7.9026E-262 | 3.0431E-181 |
| 4 | -612866.6368 | -552245.2218 | 8.86633E-45 | 1.00164E-07 |
| 5 | -612828.7348 | -552243.7091 | 2.56082E-28 | 4.54611E-07 |
| 6 | -612822.4509 | -552238.6462 | 1.37227E-25 | 7.18535E-05 |
| 7† | -612765.2027 | -552229.1054 | 1 | 0.999927592 |
| 8 | -613211.8583 | -552484.8320 | 1.047E-194 | 8.6956E-112 |
| 9 | -613168.7562 | -552479.3732 | 5.4825E-176 | 2.0419E-109 |
| 10 | -613102.1509 | -552454.2002 | 4.6266E-147 | 1.74796E-98 |
| 11 | -613079.9316 | -552426.8867 | 2.0654E-137 | 1.27244E-86 |
| 12 | -613379.6350 | -552538.6937 | 1.4303E-267 | 3.5276E-135 |
| 13 | -613390.9585 | -552533.6798 | 1.7287E-272 | 5.3085E-133 |
| 14 | -613272.7561 | -552561.3165 | 3.7359E-221 | 5.2787E-145 |
| 15 | -613247.7503 | -552566.6741 | 2.7054E-210 | 2.4873E-147 |

890

891  **Table 1 — Marginal Likelihoods for Possible Topological Hypotheses.** Topologies of the 15

892  models are displayed in Figure S1.

893  †*a priori* allopolyploid hypothesis

894

32

895    **FIGURES**



896

897    **Figure 1 — Hypothesized Relationships among North American *Dryopteris*.** Synthesis of

898    results from Sessa et al. 2012a and Sessa et al. 2012b. A) Links between shapes show the

899    putative parents and their allopolyploid derivatives. Black circles are diploids, squares are

900    tetraploids, and the hexagon is the one hexaploid species in the group. *Dryopteris semicristata*

901    is presumed extinct. B) Placement of allopolyploids in the context of the backbone relationships

902    among diploids. Tetraploids are indicated with solid orange lines and the hexaploid with dotted

903    lines. The grey line denoting a sister relationship between *D. semicristata* and *D. expansa*

904    reflects one possible placement for the extinct taxon based on previous analyses (Sessa et al.,

905    2012b).

33

906



907

908 **Figure 2 — Species network used for simulation.** The divergence times in expected

909 substitutions per site are given for each node, and *h* is the hybrid node where two alleles enter

910 from both *t* and *v*. E is an allotetraploid while other species are diploid. Nucleotide divergence

911 was reduced by dividing all $\tau$ by 10 or increased by multiplying all $\tau$ by 10. ILS was increased by

912 halving the distance between $\tau_h$ and $\tau_u$ and between $\tau_u$ and $\tau_s$ either once (for the medium ILS

913 condition) or twice (for the high ILS condition).

914

34

915



916 **Figure 3 — PATÉ Phasing Pipeline.** Overview of data input, output, and steps taken to phase

917 alleles. Input data are assumed to be paired-end Illumina reads and reference sequences for

918 each individual are required (consensus loci from HybPiper can be used). The ploidy of each

919 individual must be specified. Only biallelic sites are used for phasing.

920

**Figure 4 — Proportion of simulations that correctly identify the allopolyploid lineage.** The x-axes are the number of loci sampled for each simulation. Th y-axes are the proportions of correct networks. Results are based on networks estimated with a single reticulation, even if that network was considered less optimal than networks with zero or two reticulations. We saved the true gene trees from the simulations, while we estimated gene trees with the phased and unphased (Genotype, Consensus, and Pick One) data.

**Figure 5 — Estimating the Timing of Introgression.** Divergence times are for node *h* in Figure 2. The low divergence simulation corresponds to the y-axis units of $\tau_h \times 10^{-2}$ while the high divergence case is represented by $\tau_h$. Divergence times are measured as the expected number of substitutions per site. The dashed line represents the true simulated values. Points are posterior means, and error bars are 95% HPD intervals, averaged across 30 replicates.
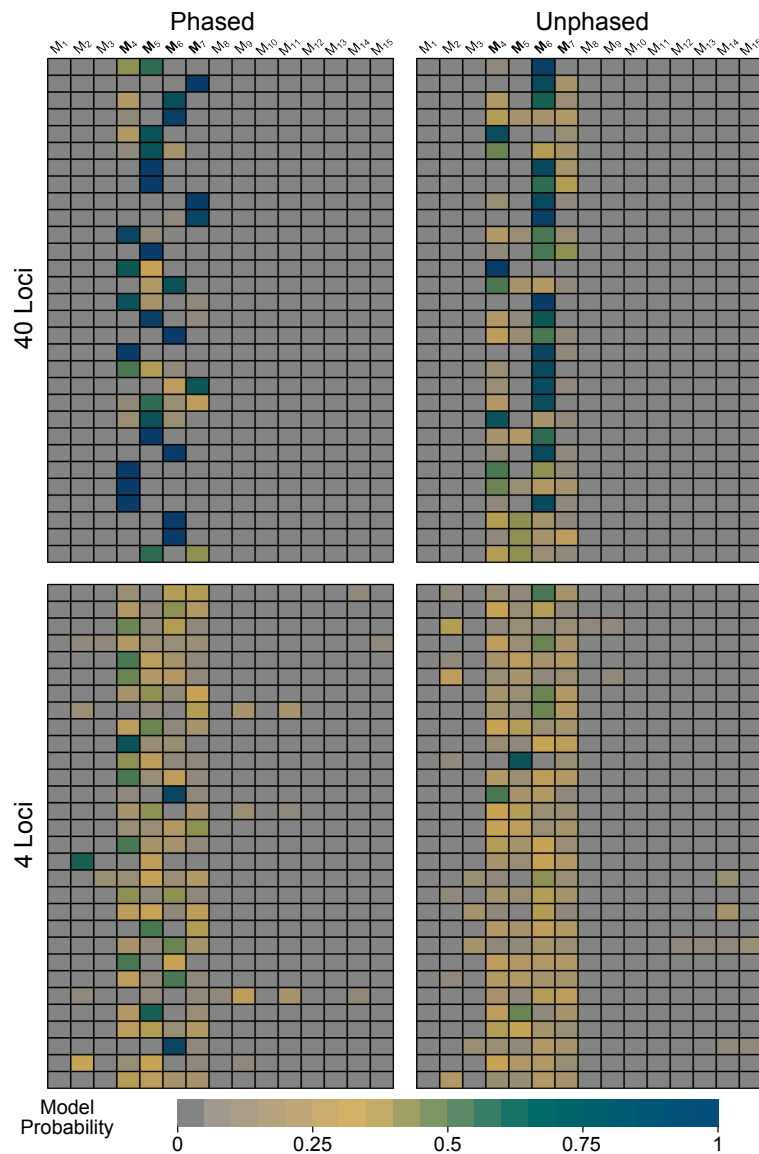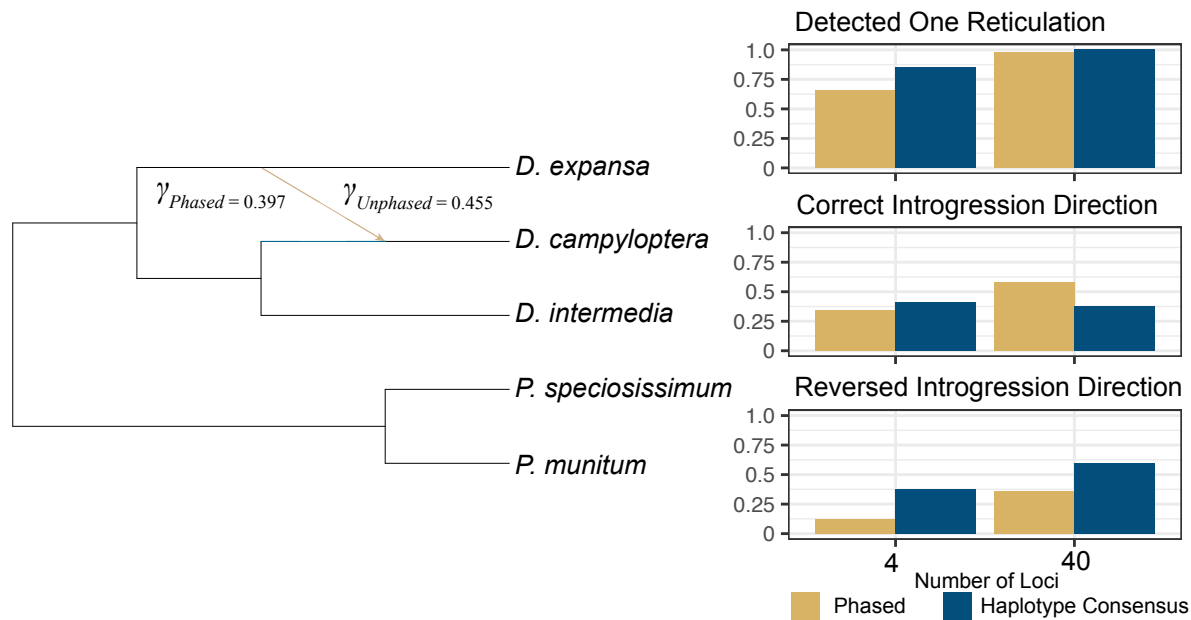
**Figure 6— *Dryopteris* Divergence Times under the MSci Model.** Divergence times are measured as the expected number of substitutions per site. The x-axis shows estimates from phased data and the y-axis shows estimates from unphased (haplotype consensus) data for the inset network (equivalent to model 7 from Supplementary Fig. S1). The dashed one-to-one line shows where older age estimates are consistently obtained from phased data. Error bars on points show the 95% HPD intervals for phased and unphased data.
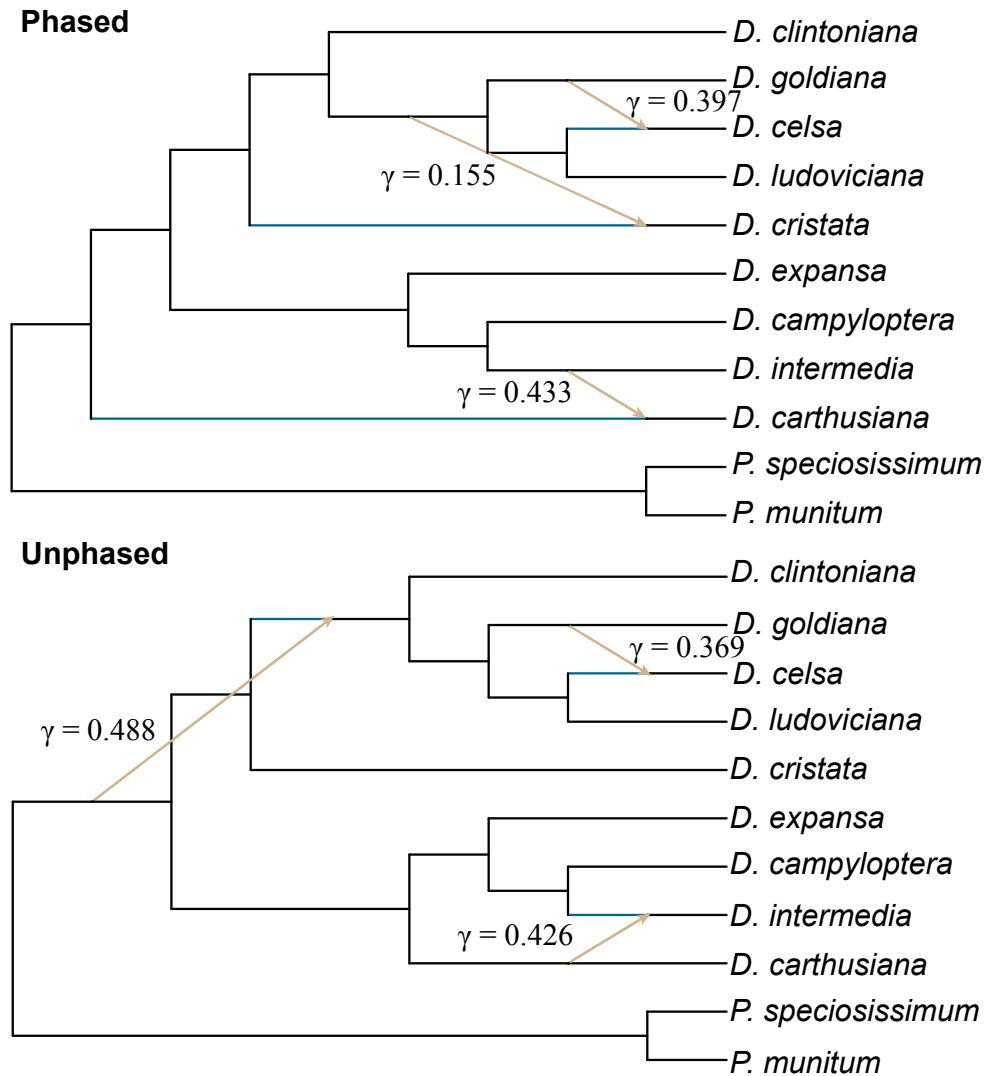
**Figure 7— Probabilities for three-taxon *Dryopteris* MSci Models with fewer Loci.** Marginal

likelihoods were estimated for the 15 MSci models (Supplementary Fig. S1). Model weights

were used to obtain probabilities. We consider model probability greater than 0.95 as decisive

evidence in favor of a model, and a probability less than 0.05 is evidence against a model. We

considered probabilities between 0.05 and 0.95 to be ambiguous. Models four through seven (in

bold) all have the correct reticulate relationships between the parental diploid lineages and the

allopolyploid. Models one through three do not have introgression, and models eight through 15

have incorrect introgression events. Each row represents one of 30 sampling replicates.

953

**Figure 8 — Network search results for three-taxon *Dryopteris* example.** Both phased and unphased haplotype consensus data recovered the same network topology with introgression occurring in the expected direction. The major topology is indicated by the blue edge and the minor edge (direction of introgression) is shown in tan. The inheritance probability $\gamma$ was slightly higher in the haplotype consensus data ($\gamma_{unphased} = 0.455$). Bar plots show the proportion of 100 replicates when sampling four or 40 loci that correctly detect one reticulation based on the pseudolikelihood scores (top), correctly estimate the network with introgression going from one of the diploids into *D. campyloptera* (middle), and estimate a network where the direction of introgression is from *D. campyloptera* into one of the diploids (bottom).

963

964

**Figure 9 — Networks for Nine-Taxon *Dryopteris* example.** Both data types recovered

optimal networks with three reticulation events. The major topology edge is blue and the minor

(reticulation) edges shown in tan, with the direction of introgression flowing into the major edge.

The position of *D. carthusiana* changes in the major topology between phased and unphased

data but the relationships are otherwise the same. All three reticulation events in the phased

data are plausible, but in the unphased data, the direction of introgression from *D. carthusiana*

into *D. intermedia* is incorrect and the reticulate edge from the common ancestor of *Dryopteris* is

difficult to reconcile. Inheritance probabilities for each introgression event are shown next to

reticulation edges.